



ERNIE 5.0 Technical Report

ERNIE Team, Baidu
ernie@baidu.com

Abstract

In this report, we introduce *ERNIE 5.0*, a natively autoregressive foundation model designed for unified multimodal understanding and generation across text, image, video, and audio. All modalities are trained from scratch under a unified next-group-of-tokens prediction objective, based on an ultra-sparse mixture-of-experts (MoE) architecture with modality-agnostic expert routing. To address practical challenges in large-scale deployment under diverse resource constraints, *ERNIE 5.0* adopts a novel elastic training paradigm. Within a single pre-training run, the model learns a family of sub-models with varying depths, expert capacities, and routing sparsity, enabling flexible trade-offs among performance, model size, and inference latency in memory- or time-constrained scenarios. Moreover, we systematically address the challenges of scaling reinforcement learning to unified foundation models, thereby guaranteeing efficient and stable post-training under ultra-sparse MoE architectures and diverse multimodal settings. Extensive experiments demonstrate that *ERNIE 5.0* achieves strong and balanced performance across multiple modalities. To the best of our knowledge, among publicly disclosed models, *ERNIE 5.0* represents the first production-scale realization of a trillion-parameter unified autoregressive model that supports both multimodal understanding and generation. To facilitate further research, we present detailed visualizations of modality-agnostic expert routing in the unified model, alongside comprehensive empirical analysis of elastic training, aiming to offer profound insights to the community.

Feb 4, 2026

Contents

1	Introduction	3
2	Architecture	4
2.1	Unified Autoregressive Backbone with Ultra-Sparse Mixture-of-Experts	4
2.2	Visual Modeling	5
2.2.1	Vision Tokenization	5
2.2.2	Visual Understanding with Dual-Path Hybrid Representation	6
2.2.3	Visual Generation with Next-Frame-and-Scale Prediction	7
2.3	Audio Modeling	7
2.3.1	Audio Tokenization	8
2.3.2	Audio Understanding and Generation with Next-Codec Prediction	8
3	Pre-Training	9
3.1	Pre-Training Data	9
3.2	Training Recipe	10
3.3	Once-For-All with Elastic Training	10
4	Post-Training	11
4.1	Enhancing Rollout Efficiency with Unbiased Replay Buffer	12
4.2	Stabilizing Training with Mitigated Entropy Collapse	13
4.3	Boosting Sample Efficiency with Hint-based Learning	14
5	Infrastructures	15
5.1	Hybrid Parallelism for Training at Scale	16
5.2	Disaggregation Architecture for Multimodal Training	16
5.3	FlashMask for Flexible Multimodal Attention	16
5.4	Scalable and Disaggregated RL Infrastructure	17
6	Evaluations	17
6.1	Evaluation on Language Benchmarks	17
6.2	Evaluation on Vision Benchmarks	20
6.3	Evaluation on Audio Benchmarks	22
6.4	Discussion	23
6.4.1	Modality-Agnostic Expert Routing	23
6.4.2	Elastic Training	25
7	Conclusion	27
8	Contributors	27

1 Introduction

Recent advances in large language and vision-language models, including ERNIE (ERNIE Team, 2025), Gemini (DeepMind, 2025a;b), GPT (OpenAI, 2024; 2025), Claude (Anthropic, 2025), DeepSeek (Liu et al., 2024; 2025a), and Qwen (Yang et al., 2025; Bai et al., 2025), demonstrate that large-scale autoregressive sequence modeling provides a powerful foundation for language and multimodal understanding. By modeling diverse inputs as sequences of tokens, these models exhibit strong reasoning and alignment capabilities across modalities. However, in most existing systems, autoregressive modeling serves multimodal understanding while the output is still text-centric, which restricts the model’s ability to engage in various multimodal interactions. To overcome this limitation, recent approaches augment pre-trained language models with modality-specific decoders or generators, which are connected to the language backbone through late-fusion designs (Xu et al., 2025a; Seedream et al., 2025). Although effective for individual modalities, these designs decouple multimodal generation from understanding and rely on modality-specific, non-autoregressive objectives, which hinder deep cross-modal integration and often force a trade-off between multimodal integration and core language performance. Against this backdrop, designing a unified autoregressive paradigm remains a prominent open challenge. Such a framework must natively support both multimodal understanding and generation, preserve strong unimodal capabilities, and scale effectively as model and data sizes continue to grow.

In this report, we introduce *ERNIE 5.0*, a next-generation foundation model natively designed to integrate text, image, audio, and video capabilities *under a unified autoregressive framework* for both multimodal understanding and generation. Rather than augmenting a pre-trained language model with modality-specific components, *ERNIE 5.0* trains all modalities simultaneously *from scratch*, which alleviates the “ability seesaw” problem observed in later-fusion approaches and ensures that all modalities evolve collectively without sacrificing performance. Specifically, heterogeneous inputs are mapped into a shared token space, and modeling across all modalities is formulated under a unified *Next-Group-of-Tokens Prediction* objective, which avoids explicit modality boundaries and inconsistent optimization trajectories. To support scalability, *ERNIE 5.0* leverages an ultra-sparse Mixture-of-Experts (MoE) backbone with *modality-agnostic expert routing*. Routing decisions are conditioned on unified token representations rather than modality identifiers, allowing tokens from various modalities to be dispatched to a shared pool of experts. This ultra-sparse, modality-agnostic architecture eliminates the need for heuristic modality-specific expert allocation, offering sufficient capacity for both differentiation and collaboration among modality-specialized behaviors.

During pre-training of *ERNIE 5.0*, we propose a novel elastic training paradigm that enables a single pre-training run to produce a family of models with varying capacity–efficiency trade-offs. Instead of optimizing a static architecture, our elastic training approach dynamically samples sub-models with varying depth, width, and routing sparsity for each training instance, guided by a pre-defined schedule. Both the sampled sub-models and the full-model are optimized in one backpropagation process under the same autoregressive objective. The sub-model sampling strategy improves the functional integrity of parameters and maintains competitive performance even when only a subset of parameters is available. Elasticity in depth and width facilitates the production of models with smaller sizes, whereas the sparsity elasticity reduces the number of activated experts during inference, leading to higher throughput and improved computational efficiency. Meanwhile, elastic training enables sub-models to inherit knowledge from the full model and provides flexible instantiation of smaller models in subsequent post-training stages, thereby eliminating the need to pre-train multiple models of various sizes or rely on customized compression, and making *ERNIE 5.0* well suited for deployment under diverse hardware, memory, and latency constraints.

Following unified pre-training, we conduct a multi-stage post-training pipeline that combines supervised fine-tuning (SFT) with unified multimodal reinforcement learning (UMRL). The coexistence of heterogeneous multimodal inputs and ultra-sparse MoE architecture introduces substantial optimization challenges, which increases the sensitivity of UMRL to sampling bias, sparse reward signals, and entropy collapse. To cope with these issues, we build a unified verifier system and develop a suite of scalable techniques to improve the stability and efficiency of RL training. An unbiased replay buffer is employed to improve rollout efficiency while preserving a balanced data distribution. Multi-granularity importance sampling, together with positive sample masking, stabilizes policy optimization and effectively mitigates entropy collapse. For difficult tasks with sparse rewards, adaptive hint-based RL is introduced to provide auxiliary guidance when needed. By ensuring stable and efficient post-training, these designs support the excellent multimodal reasoning ability in *ERNIE 5.0*.

For infrastructure, we utilize hybrid parallelism with fine-grained memory control to support effective training of a trillion-parameter ultra-sparse MoE model. For unified multimodal training, tokenizers are decoupled from the MoE backbone and deployed on separate GPU nodes, so that each component can adopt its most suitable parallelization strategy. To accommodate local bidirectional attention in vision, we

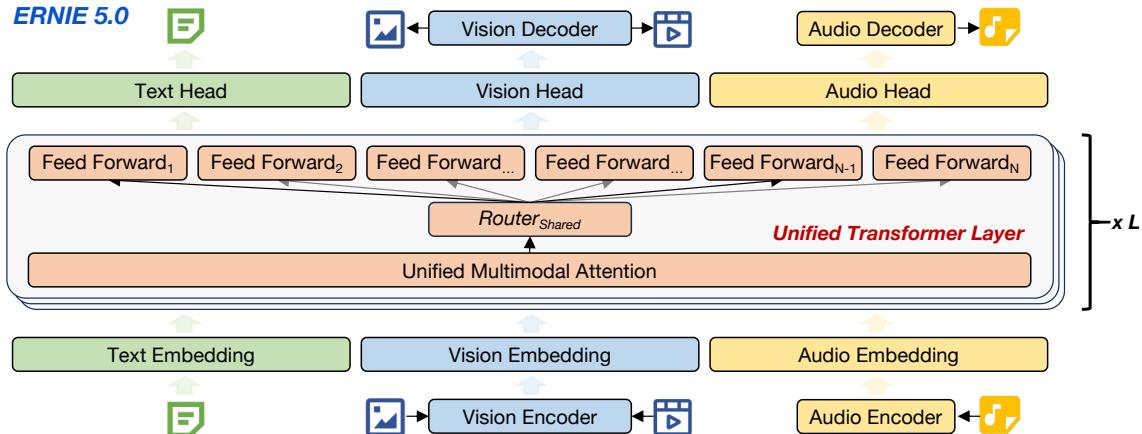


Figure 1: *ERNIE 5.0* architecture. It is trained from scratch under a unified autoregressive paradigm that integrates multimodal understanding and generation. Text, vision, and audio are encoded and serialized, then processed by a unified backbone. An ultra-sparse MoE architecture with modality-agnostic routing is employed, in which tokens from different modalities are dispatched to a shared expert pool.

employ FlashMask (Wang et al., 2024a) to efficiently handle per-sample heterogeneous attention masks. Finally, we design a scalable and disaggregated RL infrastructure that coordinates training, rollout, and environment interaction, to ensure high-throughput and numerically consistent post-training.

We evaluate *ERNIE 5.0* on a diverse set of text and multimodal benchmarks spanning perception, reasoning, understanding, and generation. Across these tasks, *ERNIE 5.0* consistently matches or outperforms specialized baselines, indicating that unified training retains strong modality-specific performance without architectural fragmentation. Ablation results further highlight the effectiveness of modality-agnostic expert routing and elastic training. Despite employing a single shared routing mechanism across modalities, experts exhibit clear specialization patterns that are primarily shaped by task requirements rather than modality boundaries. Reducing routing top- k to 25% during inference yields over 15% decoding speedup with only minor accuracy loss, while elastic training across depth, width, and sparsity preserves near-full performance using only 53.7% activated parameters and 35.8% total parameters, suggesting a scalable and efficient foundation for next-generation unified multimodal models.

In the following sections, we systematically present the design of the model architecture and its core technical components, followed by a detailed description of the training and optimization pipeline. We then evaluate *ERNIE 5.0* on a comprehensive set of benchmarks to validate the effectiveness of the proposed unified framework. Finally, we share some key technical insights during the model training process, hoping to be helpful for future research on scalable and general-purpose foundation models.

2 Architecture

As shown in Figure 1, *ERNIE 5.0* adopts an ultra-sparse mixture-of-experts architecture that integrates language, image, video, and audio within a single autoregressive framework for both multimodal understanding and generation. The model consists of a shared backbone for unified sequence modeling, together with visual and audio tokenizers that convert multimodal inputs into a unified token sequence. All modalities are trained under a shared *Next-Group-of-Tokens Prediction* objective, enabling deep cross-modal interactions with end-to-end optimization. In this section, we first describe the autoregressive backbone in Sec. 2.1, which serves as the core of *ERNIE 5.0*, followed by the visual and audio processing pipelines in Secs. 2.2 and 2.3.

2.1 Unified Autoregressive Backbone with Ultra-Sparse Mixture-of-Experts

Heterogeneous modalities differ substantially in token semantics, temporal structures, and optimization dynamics, making naive cross-modal parameter sharing prone to unstable optimization and performance degradation, especially when modeling understanding and generation in a single model.

To address these challenges, *ERNIE 5.0* is designed to be trained from scratch under a unified autoregressive framework. Multimodal inputs, including text, image, video, and audio, are projected into a shared token space, serialized into a unified sequence, and optimized under the *Next-Group-of-Tokens Prediction* objective. Specifically, text generation adheres to the standard *Next-Token Prediction* (NTP) paradigm,

augmented by the Multi-Token Prediction (MTP) mechanism (Gloeckle et al., 2024; Liu et al., 2024) to enhance both output quality and inference efficiency. For vision and audio modalities, generation is formulated as a group-of-tokens prediction task, so as to align their generative processes with the text autoregressive objective. Vision generation employs *Next-Frame-and-Scale Prediction* (NFSP) (Ji et al., 2026), and audio generation utilizes *Next-Codec Prediction* (NCP) to capture temporal and spectral structure. By unifying heterogeneous modalities under a single optimization objective, all tokens are trained within a consistent sequence prediction paradigm, enabling training from scratch and deep token-level multimodal interactions throughout the unified backbone.

Although unifying the learning paradigm helps reduce modality discrepancies, it cannot fully eliminate the intrinsic differences across modalities, necessitating substantial model capacity to capture diverse multimodal knowledge. To this end, *ERNIE 5.0* adopts a sparse Mixture-of-Experts (MoE) architecture to scale model capacity efficiently while controlling both training and inference costs. At the core of this architecture is *modality-agnostic expert routing*, where routing decisions are conditioned on unified token representations rather than explicit modality identifiers. The router dispatches tokens from different modalities to a shared pool of experts, enabling effective cross-modal parameter sharing. In contrast to modality-isolated routing strategies in our previous models (ERNIE Team, 2025), such a unified routing mechanism promotes cross-modal knowledge generalization and improves single-modality performance through the emergent specialization of shared experts. Moreover, it obviates the requirement of heuristic modality-specific expert allocation, which is often non-trivial in practice, especially when more than two modalities are involved. By employing an ultra-sparse and fine-grained MoE architecture, *ERNIE 5.0* achieves an activation rate below 3%, allowing the model to substantially expand its effective capacity without incurring a proportional increase in computational overhead. The training is further stabilized by an auxiliary-loss-free load balancing (Wang et al., 2024c), ensuring robust expert utilization at a trillion-parameter scale.

Based on the unified optimization objective and the shared parameter space, *ERNIE 5.0* formally integrates multimodal understanding and generation within a single autoregressive backbone. Despite such formal unification, some challenges remain in learning representations that can convincingly support both tasks. Typically, multimodal understanding focuses on abstract and semantic-level concepts, whereas generation requires accurate modeling of fine-grained perceptual details. *ERNIE 5.0* is therefore designed to learn a unified representation that captures high-level semantics while preserving fine-grained details, enabling both comprehension and synthesis in a unified manner. In the unified framework, semantic-level signals guide generative modeling toward global consistency, while generative training, in turn, strengthens fine-grained perception and detail-sensitive reasoning. This mutual reinforcement allows a single backbone to robustly support perception, reasoning, and creative generation. Based on this design philosophy, we further introduce unified visual and audio input-output interfaces and their corresponding processing pipelines in the following sections, which constitute a distinctive feature of *ERNIE 5.0* compared to previous models.

2.2 Visual Modeling

In *ERNIE 5.0*, image is treated as a special case of video (e.g., a single-frame video), sharing the unified design philosophy of visual understanding and generation. Visual understanding is built upon a hybrid representation that encodes both global semantic information and local perceptual details, allowing high-level reasoning while preserving fine-grained visual sensitivity. For image and video generation, a visual autoregressive paradigm is proposed to support coherent modeling across both spatial and temporal dimensions in discrete token space, as shown in Figure 2. In the following, we first introduce vision tokenization, followed by our tailored designs for visual understanding and generation.

2.2.1 Vision Tokenization

To support autoregressive visual modeling across both spatial and temporal dimensions, *ERNIE 5.0* propose *Next-Frame-and-Scale Prediction* (NFSP), where image generation is formulated as a *Next-scale Prediction* problem, and video generation further extends this formulation with *Next-Frame Prediction* (Ji et al., 2026). To this end, we first train a causal 2D multi-scale tokenizer for images, which provides strong spatial representations through large-scale image pre-training. Building upon this image tokenizer, we inflate it into a causal 3D convolutional tokenizer, thereby unifying image and video tokenization within a single model. The progressive design preserves the spatial modeling capabilities learned from images while introducing temporal perception for videos, leading to faithful reconstruction of high-level visual elements such as scene text and human faces.

During tokenizer training, we incorporate auxiliary supervision signals to enhance representation quality and training stability. Specifically, we utilize the adversarial loss (Karras et al., 2019) from GAN-based discriminators to improve distributional fidelity. Meanwhile, we incorporate a semantic branch and apply

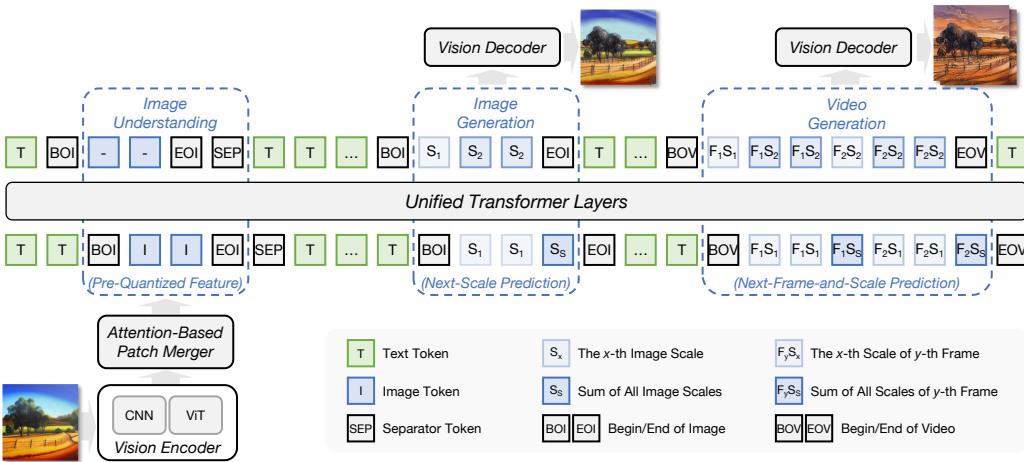


Figure 2: Overview of the unified vision understanding and generation architecture. For understanding, visual features are extracted by a hybrid CNN–ViT representation and then compressed via an Attention-Based Patch Merger. For generation, we introduce the Next-Frame-and-Scale Prediction (NFSP) paradigm, where image generation is formulated as Next-Scale Prediction, and video generation further extends this process with Next-Frame prediction along the temporal dimension.

a semantic regularization loss derived from large-scale vision foundation models to preserve high-level semantic consistency. These complementary objectives improve the learnability and stability of visual tokens, as well as facilitate effective autoregressive modeling in the unified backbone.

Following the bit-wise quantization strategy, we quantize the unified visual latent representation into a group of bit-codes, where the number of bits directly corresponding to the size of discrete vocabulary (Han et al., 2025). Based on this mechanism, we pre-train a series of tokenizers with progressively increasing bit numbers. During the training of *ERNIE 5.0*, we adopt a progressive tokenizer switching strategy, starting with a low-bit tokenizer (i.e., small vocabulary) and gradually transitioning to higher-bit variants (i.e., larger vocabularies). By first learning coarse-grained, low-bit representations with small vocabularies and progressively introducing finer-grained, higher-bit tokenizers, the backbone follows a smoother and stable optimization trajectory, effectively alleviating early-stage training instability and leading to improved visual generation quality.

2.2.2 Visual Understanding with Dual-Path Hybrid Representation

While the above design mainly targets visual tokenization for unified autoregressive modeling and generation, we observe that, prior to quantization, visual features are typically compressed by a downsampling module into low-dimensional representations, whose dimensionality is aligned with the tokenizer bit-width. The compression inevitably leads to the loss of fine-grained semantic information, which has also been widely observed to limit the performance of visual understanding tasks (Ma et al., 2025a).

To address this issue, *ERNIE 5.0* directly leverages the dual-path visual features prior to quantization. We integrate perceptual features extracted by Convolutional Neural Networks (CNNs) with semantic features encoded by a Vision Transformer (ViT). However, the representations produced by these two paths are misaligned in spatial structure, our empirical study shows that roughly fusing CNN and ViT features through MLP-based adapters often fails to fully exploit their complementary strengths and introduces representational interference, resulting in degraded understanding performance. This observation motivates the following Attention-based Patch Merger.

Formally, given a spatial token in an image or a spatio-temporal token in a video, we extract two sets of features, $\mathbf{F}_{cnn} \in \mathbb{R}^{N \times K \times D_{cnn}}$ and $\mathbf{F}_{vit} \in \mathbb{R}^{N \times K \times D_{vit}}$, where N is the number of visual understanding tokens and K is the number of local patches grouped for each token, D_{cnn} and D_{vit} are the feature dimensions. In image understanding tasks, we group $K = 4$ spatially adjacent patches, while in video understanding tasks, we group $K = 16$ patches spanning 4 neighboring frames. Before feature fusion, we project CNN features to match the ViT feature space, and then concatenate the aligned CNN and ViT patch features along the patch dimension to obtain $\mathbf{F}_{mrg} \in \mathbb{R}^{N \times 2K \times D_{vit}}$. Next, multi-head self-attention is applied to the concatenated patch tokens, $\mathbf{Z} = \text{Attn}(\mathbf{F}_{mrg})$, where the attention mechanism jointly models correlations between CNN and ViT features, as well as spatial and temporal dependencies among the group of patches. The output preserves the same shape, $\mathbf{Z} \in \mathbb{R}^{N \times 2K \times D_{vit}}$. Finally, we perform mean pooling over

the patch dimension to obtain a compact representation $\mathbf{F}_{out} \in \mathbb{R}^{N \times D_{vit}}$, which is then projected to align with the embedding dimension of the unified backbone.

The choice of feature fusion has a significant impact on model performance, and the naive MLP-based fusion is proven inadequate to effectively integrate CNN and ViT features. In contrast, the proposed attention-based aggregation module consistently outperforms both CNN-only and ViT-only baselines on a wide range of benchmarks without introducing noticeable computational overhead, with particularly pronounced gains in document and chart understanding as well as general visual understanding tasks. Importantly, *ERNIE 5.0* is designed as a unified framework for visual understanding and generation, where visual representations must support not only discriminative tasks but also fine-grained generative tasks such as pixel-level image and video editing. By utilizing attention, the module adaptively aggregates local patches together with high-level semantic information, capturing critical visual features while reducing the number of visual tokens. It results in a compact yet expressive visual representation that provides a strong and stable foundation for various vision-language tasks.

2.2.3 Visual Generation with Next-Frame-and-Scale Prediction

Recall that the Next-Frame-and-Scale Prediction (NFSP) paradigm introduced in vision tokenization formulates visual generation in an autoregressive manner, where image generation is viewed as a special case of single-frame video generation. Under this formulation, the model predicts visual tokens across multiple spatial scales within each image (or each frame) for image generation, while performing frame-wise prediction along the temporal dimension for video generation. When predicting tokens at a certain scale, the model takes the previous generated scales as input, and a *scale-wise causal attention mask* is applied, where tokens within the current scale are bidirectionally visible and are predicted in parallel, while tokens from all previous scales and historical frames are visible in a causal (uni-directional) manner. The NFSP paradigm disentangles spatial and temporal modeling, that is, intra-frame prediction from low-resolution to high-resolution captures fine-grained spatial structures, whereas next-frame prediction models inter-frame temporal dependencies.

To support positional modeling of heterogeneous tokens across spatial and temporal dimensions, we introduce a *Unified Spatiotemporal Rotary Positional Embedding (Uni-RoPE)* and apply it to all tokens in *ERNIE 5.0*. For a unified sequence of length N , the positional encoding of the i -th token is defined as $\text{Uni-RoPE}_i = (t_i, h_i, w_i), i \in \{1, \dots, N\}$. For text and audio tokens, we set $t_i = h_i = w_i$, where the shared value follows token index in the sequence. For visual tokens, t_i is used for frame indexing, which increases monotonically to preserve temporal ordering, and (h_i, w_i) corresponds to spatial locations within each frame. To ensure spatial consistency across multi-scales, we adopt a center-aligned coordinate strategy, where tokens at different scales are aligned based on the geometric centers of the scale.

Empirically, auto-regressive visual generation is susceptible to error accumulation over extremely long token sequences. To mitigate this issue, we corrupt historical tokens during training by *randomly flipping* their bits, while supervising the model to self-correct toward the ground-truth tokens of the current scale. The corruption-based training strategy improves robustness against compounding errors in long-horizon generation. Meanwhile, we apply a *loss reweighting* strategy to emphasize early-stage predictions and alleviate the token imbalance introduced by multi-scale tokenization. For video generation in particular, we further introduce *windowed temporal attention* and *random historical frame masking* to encourage the model to focus on relevant temporal context and improve robustness (Ji et al., 2026).

Within the token-based modeling paradigm, training visual generation abilities under a fixed modality token budget poses a fundamental challenge for high-resolution images and videos. Increasing visual resolution enlarges the token sequence length, which in turn reduces the effective training batch size and degrades optimization stability. To address this challenge, we adopt a cascaded diffusion refiner on top of the autoregressive backbone. The backbone generates low-resolution samples with precise semantics and structural layout, while the refiner focuses on enhancing fine-grained visual details at higher resolution. The diffusion refiner is trained separately from the backbone, using paired low-resolution samples with controlled degradation, together with their corresponding high-resolution images or videos. The decoupled training scheme enables high-fidelity refinement while preserving the complete semantic and structure produced by the autoregressive model, and avoids optimization conflicts caused by introducing autoregressive and diffusion losses within a shared backbone.

2.3 Audio Modeling

Similar to the vision modality, audio modeling in *ERNIE 5.0* is also formulated under a unified autoregressive, token-based framework that supports understanding and high-fidelity generation. Inspired by the success of neural audio codecs (Kumar et al., 2023; Zhang et al., 2024b), audio signals are represented as hierarchical discrete codec tokens that capture higher-level semantics and fine-grained acoustic details. To

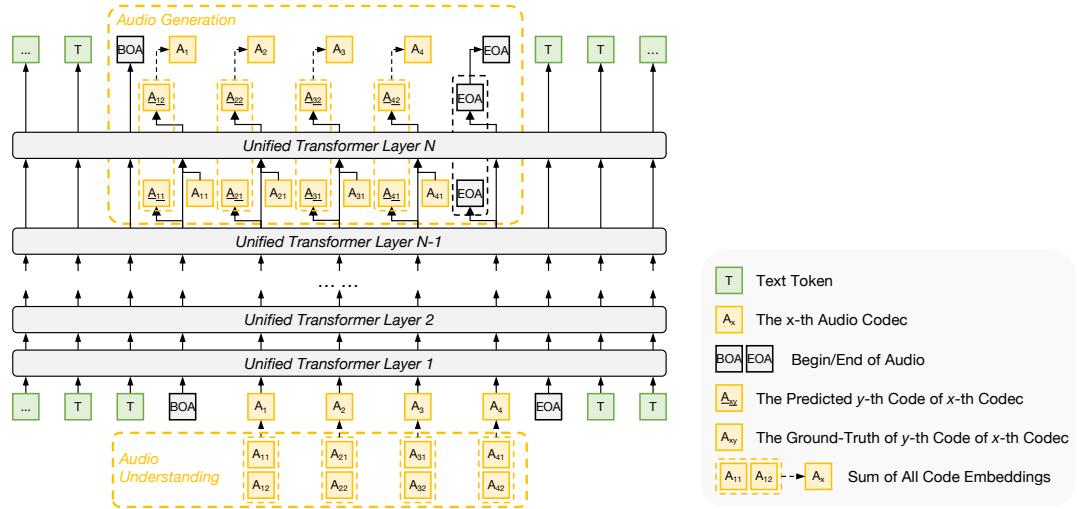


Figure 3: Overview of the depth-wise architecture for audio understanding and generation. For understanding, embeddings from multiple residual levels are additively combined to form audio token representation. For generation, *ERNIE 5.0* introduces Next-Codec Prediction (NCP) to achieve hierarchical prediction across transformer layers, where the ground-truth code embedding (or the predicted one during inference) is fed back to condition subsequent predictions.

avoid the prohibitive sequence length caused by flattening multi-codebook tokens into a single sequence, we introduce a depth-wise autoregressive architecture that performs structured prediction across codec dimensions, as illustrated in Figure 3.

2.3.1 Audio Tokenization

Given an input waveform, we first map the continuous audio signal into a sequence of discrete tokens using a codec-style tokenizer with a token rate of 12.5 Hz. The audio quantization module follows a Residual Vector Quantization (RVQ) design, decomposing the signal into multiple tokens at different levels of granularity. In that case, the first token is explicitly assigned to encode high-level audio semantics, while the remaining tokens encode residual acoustic information with progressively finer details. To ensure that the first token captures rich semantic information, including linguistic and phonetic cues, required for audio-text modeling, we distill knowledge from a pretrained Whisper model (Radford et al., 2023). Specifically, we align the representation of the first audio token with the encoder outputs of Whisper. Average pooling is applied to the Whisper representations to match the token rate of our audio tokenizer, which resolves the temporal mismatch between the teacher and student models. Complementary to the first semantic token, residual acoustic tokens preserve fine-grained characteristics of the audio signal, such as timbre and prosody. Together, the hierarchical tokenization process disentangles semantic content from acoustic realization and provides a compact audio representation that is integrated naturally into *ERNIE 5.0*’s unified autoregressive backbone.

2.3.2 Audio Understanding and Generation with Next-Codec Prediction

Based on the audio tokenization described above, we utilize a depth-wise autoregression architecture to model audio tokens for both understanding and generation, drawing inspiration from coarse-to-fine prediction paradigms developed in visual generation (Chen et al., 2024a). Instead of flattening all residual audio tokens into a single long sequence, *ERNIE 5.0* distributes the prediction of residual codes across transformer layers. Each layer models audio information at a specific level of granularity, which allows multi-level audio representations to be efficiently handled within a unified framework.

For audio understanding tasks, text tokens are embedded using standard text embedding layers, while audio tokens are represented through a depth-wise additive embedding mechanism. Each audio token consists of multiple discrete codes that correspond to different residual levels. At each level, the code is mapped to an embedding through a level-specific embedding matrices, and embeddings from all levels are summed to form the final codec representation. The additive aggregation reflects the residual nature of audio representation, where each depth contributes complementary information at different granularities from coarse to fine. The resulting audio token representations are placed at the corresponding positions in the input sequence and processed uniformly with text tokens by the autoregressive backbone, enabling seamless multimodal understanding.

For audio generation tasks, *ERNIE 5.0* introduces Next-Codec Prediction (NCP) to generate hierarchical audio tokens in a coarse-to-fine manner. Multiple audio heads are inserted into the top transformer layers to support depth-wise prediction. Conditioned on the multimodal context, the model first predicts the first semantic code and then sequentially generates codes for subsequent residual levels. After each prediction, the generated code is mapped to its corresponding embedding and added back to the hidden state, which then conditions the prediction at the next level. During training, teacher forcing is applied, and the feedback embedding is derived from the ground-truth code. Such iterative process continues until all levels are predicted, allowing high-level semantic information to guide the synthesis of increasingly fine-grained acoustic details. Once the complete set of hierarchical audio codes is obtained, the audio decoder converts them into waveforms. For speech synthesis, a speaker embedding is inserted as part of the conditioning context to enable controllable voice timbre, guiding acoustic realization without altering deep semantic content or depth-wise prediction structure. Overall, the NCP formulation is compatible with the residual design used during audio understanding, which ensures structural alignment between audio input and output.

3 Pre-Training

The pre-training phase forms the foundation of *ERNIE 5.0*, where the model learns generalizable representations across multiple modalities. In this section, we first describe the composition of pre-training data (Sec. 3.1), followed by a part of training recipe (Sec. 3.2). Finally, we focus on the *Elastic Training* technique (Sec. 3.3) introduced in *ERNIE 5.0*, which enables the production of multiple models of different sizes within a single pre-training run, significantly reducing the computational costs associated with training a series of models while maintaining both efficiency and performance.

3.1 Pre-Training Data

ERNIE 5.0 is trained on a large, high-quality multimodal dataset that reflects its natively omni design and dual capabilities for both understanding and generation. Unlike conventional late-fusion approaches, *ERNIE 5.0* is simultaneously exposed to text, images, videos, and audios from the very beginning of training. The unified training paradigm enables the model to learn representations that integrate semantic information across all modalities, while also requiring a large amount of pre-training data to support learning from scratch for each modality. To manage such diverse data, we build a standardized platform and organize all data according to their input and output modalities. Based on this organization, the pre-training data are broadly categorized into two groups: text data and multimodal data.

Text Data The textual component spans a vast collection of multilingual web crawls, curated corpora, books, scientific publications, code repositories, and structured knowledge sources selected for breadth, diversity, and linguistic richness. We retrain the text tokenizer to better support large-scale multilingual modeling. Specifically, we encode text in UTF-16BE to provide stable byte-level fallback and a more compact representation for many non-Latin symbols, improving data throughput in multilingual training, and we use BPE dropout (Provilkov et al., 2020) to reduce overfitting to frequent patterns. It is worth noting that, for languages without explicit whitespace word boundaries (e.g., Chinese), we filter out long unspaced phrases that can be decomposed by standard word-segmentation tools, which helps reduce vocabulary sparsity, improve training efficiency, and enhance model generalization.

Multimodal Data For visual and audio modalities, we curate a dataset comprising paired image–text, video–text, audio–text, as well as diverse interleaved multimodal sequences where text is integrated with images, videos, and audio, all accompanied by metadata and captions. Such data composition connects textual concepts with visual and audio contexts across both spatial and temporal dimensions. By explicitly modeling cross-modal alignment, the model is able to learn semantic relationships not only within individual modalities and across disparate ones, which supports a wide range of tasks from multimodal understanding to creative multimodal generation.

Rigorous preprocessing and quality controls are applied at scale to maintain both signal integrity and diversity. Heuristic and model-based filters remove low-quality and unsafe content, extensive deduplication prevents memorization artifacts, and decontamination safeguards keep benchmarks out of the training data. The finalized pre-training corpus comprises trillions of text tokens and multimodal instances that balance scale with high-fidelity semantic content. The large, diverse and well-filtered dataset is fundamental to *ERNIE 5.0*’s strong performance on text and multimodal understanding, reasoning, and generation benchmarks.

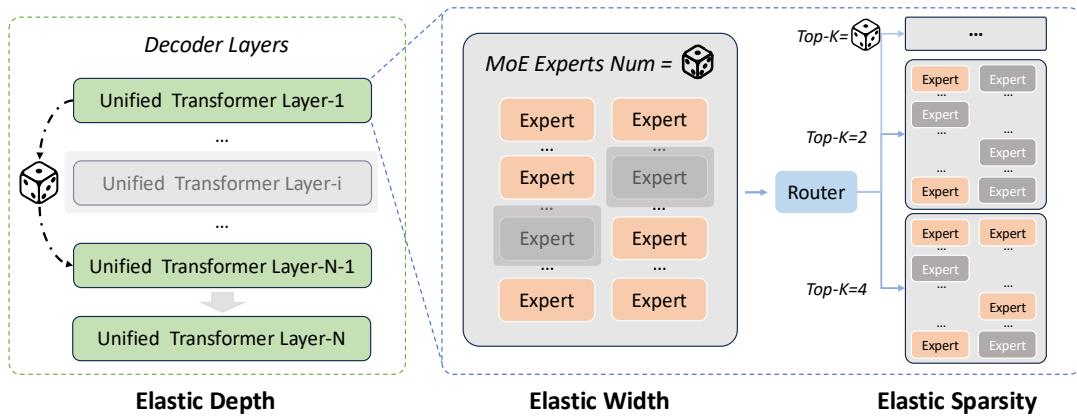


Figure 4: Overview of the elastic training framework in *ERNIE 5.0*. The framework supports elastic depth, width, and sparsity in a unified MoE architecture. *Elastic Depth* randomly adapts the number of active layers, *Elastic Width* varies the total number of experts in each MoE layer, and *Elastic Sparsity* changes the top- k routing per token. These mechanisms collectively enable flexible deployment under different compute, memory, and latency constraints without retraining.

3.2 Training Recipe

ERNIE 5.0 is trained with a carefully designed recipe to ensure training stability, scalability, and efficient utilization of compute resources. The training follows a multi-stage pre-training strategy to progressively extend context length while maintaining stable optimization dynamics.

Stage 1: 8K Pre-Training The initial stage uses a maximum context length of 8K tokens. We adopt a Warmup-Stable-Decay (WSD) learning rate schedule (Hu et al., 2024) in this stage. The learning rate is linearly warmed up for 2,000 steps from zero to a peak value of 1×10^{-4} , and then kept constant for the remainder of the 8K training stage. To improve training efficiency and stability at scale, we employ a batch size scheduling strategy, where the global batch size is gradually increased from 14M tokens to 56M tokens during early training. To facilitate seamless long-context extension, the RoPE base is set to 1,000,000 starting from the 8K stage. This design choice avoids the need for reparameterization or interpolation during subsequent context length expansion, ensuring lossless and stable long-context training.

Stage 2: 32K&128K Mid-Training In the mid-training stage, we progressively extend the context length to 32K and 128K tokens while keeping the global batch size unchanged. During this phase, we switch to a cosine learning rate schedule and anneal the learning rate from 1×10^{-4} to 1×10^{-5} .

For MoE-specific optimization, the bias update speed for auxiliary-loss-free load balancing (Wang et al., 2024c) is set to 1×10^{-4} in the 8K pre-training stage and reduced to 1×10^{-5} during mid-training, which effectively suppresses iteration-level oscillations observed in large-scale MoE training. The MTP loss weight (Liu et al., 2024) is decreased from 0.3 in the 8K stage to 0.1 during mid-training, ensuring stable adaptation as the model scales to longer contexts. Besides, we introduce a posterior-based loss weighting strategy that rescales the autoregressive losses of different modalities to the same interval, thereby improving training stability and preventing imbalance across modalities.

3.3 Once-For-All with Elastic Training

The scaling of modern models presents a fundamental challenge. Although models with trillions of parameters achieve remarkable performance across a wide range of tasks, their high computational and deployment costs limit the applicability in scenarios that demand flexibility and efficiency.

Traditional approaches typically follow a “train-then-compress” pipeline, employing techniques such as pruning (Sajjad et al., 2023; Xia et al., 2023; Men et al., 2025), knowledge distillation (Gu et al., 2023; Xu et al., 2024) or more efficient fusion-based variants (Wang et al., 2023; Chen et al., 2024d;e) to produce smaller models. However, this paradigm still suffers from notable limitations. Model compression requires a dedicated pruning or distillation stage, which demands specialized infrastructure and incurs substantial computational overhead. Moreover, once a model is compressed, its architecture becomes fixed. Consequently, creating models of other sizes necessitates repeating the full compression process, thereby constraining deployment flexibility.

To address these issues, we propose a novel elastic training strategy and apply it to *ERNIE 5.0* for the first time. Rather than compressing a pre-trained model post hoc, elastic training simultaneously optimizes a family of sub-networks during pre-training, so that a single large model to efficiently produce smaller, deployable variants on demand. It extends the design philosophy of Once-For-All (Devrit et al., 2023; Cai et al., 2024; Gu et al., 2025) to pre-training, in which sub-network configurations of varying depth, width, and sparsity are trained together with the full-scale model. As a result, *ERNIE 5.0* can flexibly selects subsets of parameters to construct models at different scales, which reduces the computational and engineering overhead compared with traditional pruning or distillation methods.

The elastic training is shown in Figure 4, which introduces structural flexibility along three orthogonal dimensions:

Elastic Depth To support elastic depth, *ERNIE 5.0* randomly varies the number of active transformer layers during training, enabling the extraction of sub-networks with different depths. Most of the time, the full-depth network is used to ensure all layers are well-optimized, while shallower sub-networks are occasionally sampled to foster resilience against layer removal. Specifically, the full model is trained with a probability of 75%, while a reduced-depth sub-network is activated with a probability of 25%. Through this training scheme, intermediate representations are encouraged to remain informative even when some layers are bypassed, and the model supports flexible deployment across different depth configurations without requiring separate training.

Elastic Width Complementary to elastic depth, *ERNIE 5.0* also supports elastic width by varying the total number of experts in each Mixture-of-Experts (MoE) layer. Instead of always activating the full expert pool, the training process alternates between two modes. With a probability of 80%, all experts participate in routing, preserving the full-width configuration. In the remaining 20% of cases, routing is restricted to a randomly sampled subset of experts, leading to a narrower effective model width. By exposing the model to both full and reduced expert configurations, the resulting model supports different capacity budgets with only partial experts, making it suitable for deployment in memory-constrained environments where hosting all experts is impractical.

Elastic Sparsity To improve inference efficiency without changing the deployed model size, elastic sparsity is introduced by varying the number of activated experts per token. Similar to the elasticity of total number of experts, the default routing configuration is applied with a probability of 80% during training. With a probability of 20%, the routing top- k is randomly sampled from a predefined range, where k is smaller than the standard configuration. In other words, the number of activated experts for each token is decreased. Finally, the model is compatible with different compute budgets and exhibits improved robustness during latency-constrained inference.

By training an elastic super-network, *ERNIE 5.0* is able to produce smaller models of varying configurations by selecting subsets of parameters along the layer number, total expert number, and activated expert number. Elasticity along the representation dimension (i.e., hidden size), as explored in recent work (Chen et al., 2024c), is orthogonal to our design and can be naturally incorporated as a future extension. These sub-networks can be instantiated on demand to meet different latency and memory constraints, serving as effective starting points for mid-training or fine-tuning. Compared with training separate models from scratch or relying on post-hoc compression techniques, our elastic training strategy substantially lowers overall computational overhead and engineering complexity. Detailed experimental results and ablation studies are provided in Sec 6.4.2.

4 Post-Training

After unified pre-training, we follow the same post-training pipeline as *ERNIE 4.5* (ERNIE Team, 2025) to obtain the final *ERNIE 5.0*, which includes two stages, supervised fine-tuning (SFT) and unified multimodal reinforcement learning (UM-RL). With curating comprehensive set of high-quality instruction pairs, SFT endows the model with fundamental instruction-following capability and strengthens its ability to think through long chains-of-thought. During UM-RL phase, we merge the training of various tasks such as reasoning, agent, and instruction following into a multi-stage RL pipeline, enabling balanced performance across diverse tasks and modalities. We further extended the unified verifier system to generate accurate and consistent reward signals for model responses in a wide range of multimodal scenarios, providing reliable supervision for unified multimodal RL training. In this section, we will discuss the key challenges of RL training and describe the solutions we propose.

The RL training of *ERNIE 5.0* faces several challenges. Firstly, RL training is computationally expensive, which is further amplified by the large scale of *ERNIE 5.0*. Secondly, the ultra-sparse MoE architecture exacerbates the training–inference discrepancy and undermines stability. Finally, compared to standalone

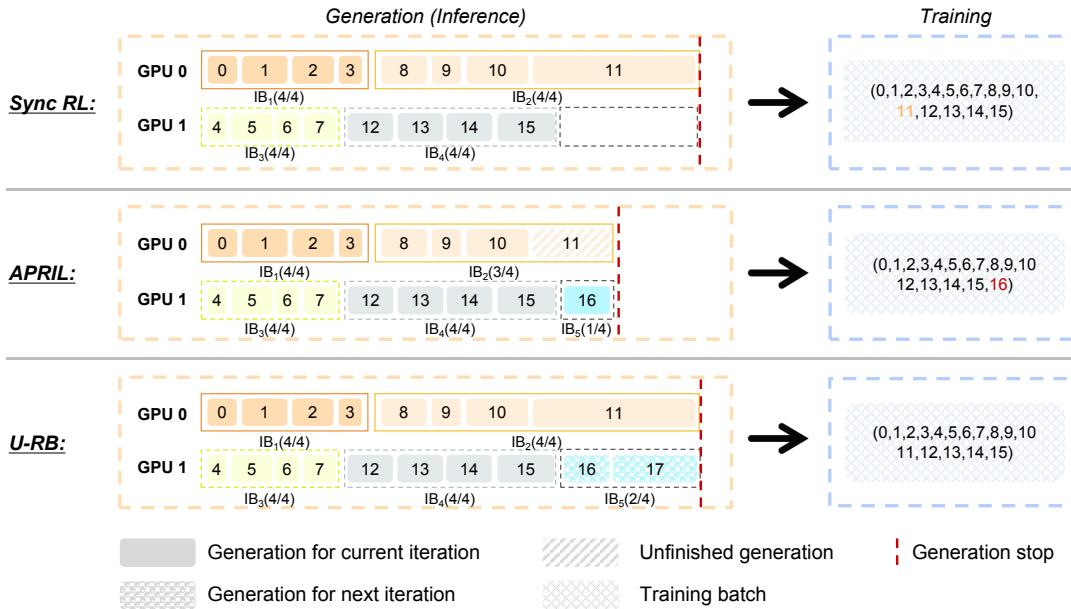


Figure 5: Visualization of the Unbiased Replay Buffer (U-RB) in *ERNIE 5.0*, in comparison with existing methods, where each query is assigned a unique index, and IB denotes the inference batch. In Sync RL, a long-tail query (e.g., index 11) blocks the entire batch, leaving GPUs idle and poorly utilized. APRIL stops generation once the target number of responses (16) is reached, which leads to a non-stationary data difficulty distribution. U-RB extends APRIL with a data-ordering constraint that prepares future batches while waiting for long-tail queries, preserving query order and mitigating inefficiency.

RLVR tasks such as mathematical reasoning or code generation, training a model that simultaneously supports multiple scenarios and modalities introduces substantially higher complexity. To address these bottlenecks, we implement a suite of synergistic engineering and algorithmic optimizations that enable stable RL training for large-scale, ultra-sparse MoE models. In the sections that follow, we delineate the formidable challenges encountered in this endeavor and present our corresponding solutions.

4.1 Enhancing Rollout Efficiency with Unbiased Replay Buffer

Rollout generation accounts for more than 90% of the total training time in RL, and efficiency is often limited by the long-tail distribution of rollout response lengths. In that case, a small number of unusually long responses stall entire batches, leaving GPUs idle and underutilized. Recent work, such as APRIL (Zhou et al., 2025), seeks to mitigate long-tail inefficiency by over-provisioning rollout requests. Generation is terminated once a target number of responses is collected, and incomplete responses are recycled for continuation in subsequent steps. However, APRIL tends to update model parameters using trajectories with shorter reasoning steps, which usually correspond to easier queries. In contrast, longer-horizon samples are deferred, leading to a non-stationary distribution of data difficulty. Consequently, periodic shifts in data difficulty may hinder convergence and ultimately degrade model performance.

U-RB: Unbiased Replay Buffer Generation We introduce U-RB, an unbiased extension of APRIL that accelerates rollout generation in RL. As illustrated in Figure 5, U-RB introduces a data-ordering constraint, under which only the data group assigned to the current iteration at initialization is allowed to participate in subsequent training process. Specifically, U-RB builds two modules. The first is a high-throughput inference pool, $\mathcal{P}_{inference}$, with capacity $\Omega_{RBS} = \Omega_{BS} * N$, where Ω_{BS} is the training batch size and N is the buffer size. The second component is a training pool \mathcal{P}_{train} with capacity Ω_{BS} , which collects completed trajectories for RL training. At iteration t , the inference engine $\pi_{inference, \theta_t}$ populates the inference pool by generating rollouts in parallel. Inference proceeds until the terminal state (i.e., [EOS]) is reached for the longest rollout belonging to the data group \mathcal{D}_t assigned to iteration t . At this point, rollouts associated with \mathcal{D}_t are moved from $\mathcal{P}_{inference}$ to the training pool \mathcal{P}_{train} , enabling the training engine π_{train, θ_t} to update model parameters. These rollouts may include trajectories resumed from earlier inference runs. By dynamically partitioning rollout generation, U-RB prevents computational idleness caused by individual long rollouts, while maintaining an unbiased data distribution.

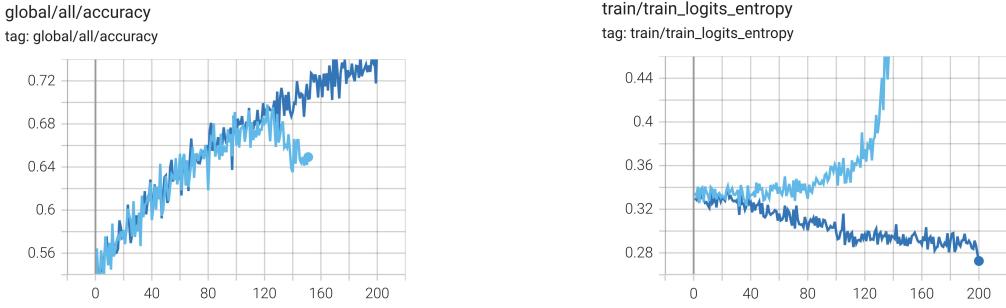


Figure 6: Training dynamics of applying $\mathfrak{J}_{IcePop}^{Mixed}$ (dark-blue) and $\mathfrak{J}_{IcePop}^{GSPO}$ (light-blue) to conduct RL training on ERNIE 5.0. By using Multi-granularity Importance Sampling Clipping (MISC), we avoid the entropy collapse during early stage and achieve stable RL training.

4.2 Stabilizing Training with Mitigated Entropy Collapse

The phenomenon of rapid entropy collapse in a multimodal model is manifested as a sharp increase or decrease in policy entropy during the early stages of RL. In multimodal decision-making tasks that integrate text, vision and audio information, such collapse gradually erodes the model’s ability to fuse information across modalities for flexible reasoning and reveals a pronounced modality bias.

Recent studies (Cui et al., 2025; Wang et al., 2025) attribute entropy collapse mainly to two factors. First, most contemporary RL frameworks rely on separate engines for training and inference, which introduces inconsistencies in numerical computation, and ultimately destabilizes policy optimization. The problem becomes more severe for MoE models, where dynamic routing further amplifies the numerical mismatch problem. Second, the policy model often overfits easy queries in the early stage of training. Such behavior accelerates entropy collapse and limits the model’s ability to discover alternative reasoning paths. To address these issues, we introduce Multi-granularity Importance Sampling Clipping (MISC) and Well-learned Positive Sample Mask (WPSM) to stabilize RL training at scale.

MISC: Multi-granularity Importance Sampling Clipping IcePop (Ling-Team et al., 2025) suppresses training-inference mismatch through double-sided masking calibration on GRPO (Guo et al., 2025):

$$\begin{aligned} \mathfrak{J}_{IcePop}^{GRPO}(\theta) &= \mathbb{E}_{x \sim D, \{y_i\}_{i=1}^G \sim \pi_{infer}(\cdot | x; \theta_{old})} \\ &\left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|y_i|} \sum_{j=1}^{|y_i|} [\mathfrak{M}\left(\frac{\pi_{train}(y_{i,j}|x, y_{i,<j}; \theta_{old})}{\pi_{infer}(y_{i,j}|x, y_{i,<j}; \theta_{old})}; \alpha, \beta\right) \cdot \min(r_{i,j} \hat{A}_{i,j}, \text{clip}(r_{i,j}, 1 - \epsilon, 1 + \epsilon) \hat{A}_{i,j}) \right] \\ r_{i,j} &= \frac{\pi_{train}(y_{i,j}|x, y_{i,<j}; \theta)}{\pi_{train}(y_{i,j}|x, y_{i,<j}; \theta_{old})} \\ \mathfrak{M}(k) &= \begin{cases} k & \text{if } k \in [\alpha, \beta], \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (1)$$

where α, β controls the lower and upper limits. We apply this technology to GSPO (Zheng et al., 2025):

$$\begin{aligned} \mathfrak{J}_{IcePop}^{GSPO}(\theta) &= \mathbb{E}_{x \sim D, \{y_i\}_{i=1}^G \sim \pi_{infer}(\cdot | x; \theta_{old})} \\ &\left[\frac{1}{G} \sum_{i=1}^G [\mathfrak{M}\left(\left(\frac{\pi_{train}(y_i|x; \theta_{old})}{\pi_{train}(y_i|x; \theta_{old})}\right)^{\frac{1}{|y_i|}}; \alpha, \beta\right) \cdot \min(s_i(\theta) \hat{A}_i, \text{clip}(s_i(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_i) \right] \\ s_i(\theta) &= \left(\frac{\pi_{train}(y_i|x; \theta)}{\pi_{train}(y_i|x; \theta_{old})}\right)^{\frac{1}{|y_i|}} = \exp\left(\frac{1}{|y_i|} \sum_{j=1}^{|y_i|} \log \frac{\pi_{train}(y_{i,j}|x, y_{i,<j}; \theta)}{\pi_{train}(y_{i,j}|x, y_{i,<j}; \theta_{old})}\right) \\ \mathfrak{M}(k) &= \begin{cases} k & \text{if } k \in [\alpha, \beta], \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (2)$$

However, our experiments show that directly applying $\mathfrak{J}_{IcePop}^{GSPO}$ to RL training on ERNIE 5.0 leads to rapid entropy collapse, as illustrated by the light-blue line in Figure 6. The phenomenon is caused by sequence-level truncated importance sampling, which prunes a large number of low-entropy responses due to the training–inference mismatch.

Original Prompt

Let S be the set of vertices of a regular 24-gon. Find the number of ways to draw 12 segments of equal lengths so that each vertex in S is an endpoint of exactly one of the 12 segments.

**Augmented Prompt**

Let S be the set of vertices of a regular 24-gon. Find the number of ways to draw 12 segments of equal lengths so that each vertex in S is an endpoint of exactly one of the 12 segments.

Please reason step by step, and put your final answer within `\boxed{}`.

<think> The problem states: Let S be the set of vertices of a regular 24-gon. I need to find the number of ways to draw 12 segments of equal lengths such that each vertex is an endpoint of exactly one segment. Since it's a regular 24-gon, all vertices are equally spaced on a circle.

Figure 7: Overview of the Adaptive Hint-based Reinforcement Learning (AHRL) in *ERNIE 5.0*, which introduces think skeletons to guide hard queries and mitigate sparse rewards.

To address this issue, we revise $\mathfrak{J}_{IcePop}^{GSPo}$ to $\mathfrak{J}_{IcePop}^{Mixed}$:

$$\begin{aligned} \mathfrak{J}_{IcePop}^{Mixed}(\theta) &= \mathbb{E}_{x \sim D, \{y_i\}_{i=1}^G \sim \pi_{infer}(\cdot | x; \theta_{old})} \\ &\quad \left[\frac{1}{G} \sum_{i=1}^G [\mathfrak{M}_{j \in [1, |y_i|]} \left(\frac{\pi_{train}(y_{i,j} | x, y_{i,<j}; \theta_{old})}{\pi_{infer}(y_{i,j} | x, y_{i,<j}; \theta_{old})}; \alpha, \beta \right) \cdot \min(s_i(\theta) \hat{A}_i, \text{clip}(s_i(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_i) \right] \\ s_i(\theta) &= \left(\frac{\pi_{train}(y_i | x; \theta)}{\pi_{train}(y_i | x; \theta_{old})} \right)^{\frac{1}{|y_i|}} = \exp \left(\frac{1}{|y_i|} \sum_{j=1}^{|y_i|} \log \frac{\pi_{train}(y_{i,j} | x, y_{i,<j}; \theta)}{\pi_{train}(y_{i,j} | x, y_{i,<j}; \theta_{old})} \right) \\ \mathfrak{M}(k) &= \begin{cases} k & \text{if } k \in [\alpha, \beta], \\ 0 & \text{otherwise} \end{cases} \end{aligned} \tag{3}$$

By modulating the trust region according to the modality sensitivity, we achieve a more balanced exploration-exploitation trade-off. The mechanism avoids premature convergence to a “safe” yet suboptimal strategy in complex multi-scenario settings, and preserves flexibility across various inputs.

WPSM: Well-learned Positive Sample Mask We introduce a sample mask strategy to prevent the model from over-optimizing on already mastered queries, in which the proficiency is tracked by maintaining a success-rate for each query. For a given query x with a rollout group $\mathcal{Y}^x = \{y_1^x, y_2^x, \dots, y_G^x\}$, where G is the group size, if the average accuracy acc_t^x in iteration t exceeds a threshold τ , we flag the rollout y_i^x in \mathcal{Y}^x as a “well-learned” response when its policy entropy $\mathcal{H}_{y_i^x}(\pi_\theta)$ falls below a stability bound η . During training, the “well-learned” responses are masked as follows:

$$\begin{aligned} \mathfrak{J}(\theta) &= \mathbb{E}_{x \sim D, \{y_i\}_{i=1}^G \sim \pi_{\theta_{old}}(\cdot | x)} \left[\frac{1}{G} \sum_{i=1}^G [1 - \mathbb{M}_{mask}^i] \min(s_i(\theta)) \hat{A}_i, \text{clip}(s_i(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_i \right] \\ \mathbb{M}_{mask}^i &= \begin{cases} \alpha & \mathcal{H}_{y_i^x}(\pi_\theta) < \eta \text{ and } acc_t^x > \tau \\ 0 & \text{otherwise} \end{cases} \end{aligned} \tag{4}$$

where $s_i(\theta)$ denotes the importance ratio and $\alpha \in [0, 1]$ controls the degree of supplementary learning applied to “well-learned” responses. Under this design, the gradient budget is shifted toward harder samples, such as those with sparse rewards or diverse reasoning paths. By masking redundant positive signals, WPSM alleviates the entropy collapse problem caused by over-fitting to easy queries, and encourages the model to improve the performances of challenging, low-performing tasks.

4.3 Boosting Sample Efficiency with Hint-based Learning

Recent studies (Yue et al., 2025b; Liu et al., 2025b; Zhao et al., 2025b) indicate that although state-of-the-art RL methods, such as GRPO (Guo et al., 2025) and DAPO (Yu et al., 2025), enhance the pass@1 metric by

reinforcing high-reward completions, they exhibit clear limitations on challenging tasks where the base model performs poorly. Specifically, when all rollouts receive zero reward, the GRPO framework fails to provide effective gradient signals for policy optimization. In such cases, RL training on hard queries tends to progress much more slowly because sparse rewards and limited sample efficiency impede learning. To address this challenge, we propose Adaptive Hint-based Reinforcement Learning (AHRL), a method that mitigates the issues of sparse rewards on hard queries. As can be seen in Figure 7, AHRL introduces partial hints that decompose the complex problem into intermediate steps and gradually increase the performance of trained models. The mechanism is described in detail below.

AHRL: Adaptive Hint-based Reinforcement Learning Unlike approaches that modify reward functions or optimization algorithms, AHRL is designed to inject partial think sketches into queries during RL training. By decomposing problems into intermediate steps, AHRL increases the propensity of the base model to generate correct responses and improves sample efficiency. As a result, the model is driven to master the hardest problems, which accelerates the RL training process.

For a given query x with a response consisting of a thinking trajectory and a final solution, denoted as $y = (\text{think}, \text{solution})$, AHRL augments x into $\tilde{x}^{(p)}$ by attaching the first p_{hint} tokens of the think to the original query. p_{hint} denotes the fraction of revealed thinking, allowing fine-grained control over query difficulty. Specifically, The probability p_{hint} follows an annealing schedule:

$$p_{\text{hint}}(x^t) = p_{\text{initial}} \cdot \exp(-\gamma \cdot t \cdot \text{pass}_{\text{initial}}^x) \quad (5)$$

where t is the training iteration, γ is the decay rate, and $\text{pass}_{\text{initial}}^x$ is the pass@k score of query x evaluated on the SFT model. As training progresses and model performance improves, the fraction of revealed hints is gradually reduced, transitioning the model to full self-exploration. The mechanism provides necessary “scaffolding” to bridge the gap between initial exploration and successful task completion. It prevents training from stalling in complex tasks where valid reasoning paths are statistically rare, and ensures consistent performance improvements across modalities.

5 Infrastructures

The training of *ERNIE 5.0* is built upon PaddlePaddle (Ma et al., 2019). Based on the infrastructure of *ERNIE 4.5* (ERNIE Team, 2025), we further address the unique challenges introduced by native multimodal training, ultra-sparse MoE models, and large-scale RL pipelines.

Tremendous Memory Pressure and Communication Overhead The ultra-sparse MoE architecture of *ERNIE 5.0* poses two major challenges for efficient training: ultra-sparse expert activation leads to heavy inter-node communication, and large-scale expert parameters impose significant memory pressure. To mitigate these challenges, we propose a hybrid parallel strategy tailored for the MoE architecture, integrated with fine-grained memory control.

Multimodal Training with Multiple Tokenizers As a unified model supporting text, vision, and audio, *ERNIE 5.0* relies on multiple modality-specific tokenizer models. Their computational characteristics differ significantly from those of the MoE backbone, making it challenging to apply traditional end-to-end optimization strategy across all parts. To resolve this mismatch, we adopt a decoupled architecture that physically separates tokenizers from the MoE backbone and deploys them on different GPU nodes, allowing each component to use its most suitable parallelization strategy.

Flexible Attention Patterns across Modalities *ERNIE 5.0* processes inputs from multiple modalities with heterogeneous attention patterns. Visual inputs typically require bidirectional attention, whereas text and audio inputs rely on unidirectional attention. Although existing solutions such as FlexAttention (Dong et al., 2024) support flexible attention, they are less efficient when attention mask patterns vary across samples within the same batch. We therefore adopt the self-developed FlashMask (Wang et al., 2024a) to accelerate attention mask computation here.

Scalable RL under Throughput and Consistent Constraints Reinforcement learning for a trillion-parameter model requires coordinated execution across training, inference, and environment interaction. Such heterogeneous workloads pose challenges in maintaining numerical consistency between training and rollout, mitigating data distribution bias in asynchronous pipelines, and maximizing the utilization of diverse hardware resources. We therefore design a scalable and disaggregated RL infrastructure to resolve these bottlenecks while preserving strict computational determinism.

5.1 Hybrid Parallelism for Training at Scale

To address the severe memory pressure and communication overhead introduced by the ultra-sparse MoE architecture, we develop a distributed parallel strategy that adapts to varying training resources. The final configuration combines 4-way tensor parallelism (Shoeybi et al., 2019), 12-way pipeline parallelism (Huang et al., 2019) with virtual stages, 64-way expert parallelism (Lepikhin et al., 2020), ZeRO-1 data parallelism (Rajbhandari et al., 2020), and context parallelism (Liu et al., 2023a) for long context training. We also use DeepEP (Zhao et al., 2025a) to enable efficient inter-node communication, while virtual pipeline parallelism is employed to minimize pipeline bubbles. To ensure final performance, *ERNIE 5.0* utilize a no-token-dropping strategy throughout training, and out-of-memory (OOM) issues sometimes occur due to unbalanced expert routing, especially in the initial stage of training. To mitigate these issues and stabilize large-scale MoE training, we develop a set of techniques as follows.

To ensure memory sufficiency, we implement the following strategies:

- **FP8 Mixed-Precision Training.** Following the practice of *ERNIE 4.5* (ERNIE Team, 2025), we adopt FP8 mixed-precision training and store activation tensors in FP8 format, which effectively reduces peak memory consumption during training.
- **Dynamic Adaptive Offloading of Activation Memory.** During forward propagation, all activation tensors retained for backward computation are tracked. We extend the memory allocator to enable adaptive offloading of selected activation tensors when an OOM event is encountered. No offloading is triggered when sufficient memory is available. This technology ensures sufficient total memory for training with minimal performance overhead.

To reduce memory fragmentation, we implement the following strategies:

- **Sub-batch Computations.** We decompose large memory allocation requests into a series of smaller requests with sub-batch computation, which reduces the probability of OOM caused by memory fragmentation.
- **Automatic Memory Defragmentation.** Based on CUDA virtual memory management (VMM), we develop a memory allocation allocator capable of automatic memory defragmentation, ensuring successful memory allocation even under extreme conditions.

Through the aforementioned techniques, we ensure the feasibility and reliability of pre-training *ERNIE 5.0* in memory-constrained scenarios.

5.2 Disaggregation Architecture for Multimodal Training

ERNIE 5.0 integrates heterogeneous multimodal inputs, where each modality is handled by a dedicated tokenizer that converts raw signals into token sequences. Both sequence lengths and computational costs vary substantially across modalities and even among samples within the same modality. Deploying all tokenizers together with backbone on homogeneous hardware would induce significant load imbalance, which in turn degrades overall training efficiency.

To solve this problem, we design a tokenizer-backbone disaggregation architecture. We decouple the tokenizers from the backbone by deploying them as independent, horizontally scalable services on dedicated compute nodes, under a data-parallel configuration. During training, the backbone interacts with these tokenizer services via remote calls to retrieve encoded representations. The architectural separation allows each component to adopt parallelization strategies suited to its own workload, improving scalability and efficiency in distributed multimodal training.

5.3 FlashMask for Flexible Multimodal Attention

In *ERNIE 5.0*, attention patterns vary across modalities and may even differ across samples. Text modeling typically relies on causal attention, while visual features often employ bidirectional attention. That is, for visual samples, attention follows a globally causal structure while allowing local bidirectional interactions to capture spatial dependencies. Efficiently supporting such heterogeneous masking patterns is therefore essential for unified multimodal training.

To meet these requirements, we employ FlashMask (Wang et al., 2024a), which not only meets the flexible and diverse attention masking needs of *ERNIE 5.0*, but also significantly accelerates the computational efficiency of attention mask operations. In practice, FlashMask achieves up to a 200% speedup over FlexAttention (Dong et al., 2024) at the operator-level, and delivers more than 20% end-to-end training acceleration. In addition, we integrate FlashMask with context parallelism (Liu et al., 2023a) at the kernel level, achieving an 80% performance improvement compared to the Megatron-LM solution.

5.4 Scalable and Disaggregated RL Infrastructure

Scaling Reinforcement Learning (RL) to unified multimodal models with trillions of parameters presents unique challenges regarding computational consistency, data distribution bias, and heterogeneous resource utilization. To address these bottlenecks, we introduce the *ERNIE 5.0* RL Infrastructure, a disaggregated system designed to orchestrate large-scale asynchronous training. By prioritizing high-throughput execution and computation determinism, our system ensures stable and efficient RL training.

The key architectural components are summarized as follows:

- **Disaggregated Control Plane for Asynchronous RL.** We introduce a fully disaggregated control plane built around a centralized RL controller to maximize system throughput, which coordinates training, inference, environment interaction, and reward evaluation in an asynchronous manner. Logical decoupling across these subsystems enables flexible scaling and efficient pipeline management, forming the foundation for large-scale asynchronous multimodal RL.
- **Unified FP8 Stack for Consistent Training and Inference.** Precision divergence is a common issue in low-bit RL training. To alleviate this situation, we build a unified FP8 execution engine. By employing identical high-performance operators across both training and inference (rollout) stages and integrating the Rollout Router Replay (Ma et al., 2025b) strategy, the engine minimizes numerical mismatch and ensures stable convergence under low-precision settings.
- **Replay Buffer for Sequence-Length Bias Mitigation.** Asynchronous rollout in RL may introduce sequence-length bias, where shorter responses enter training earlier and distort the data distribution. We design an unbiased replay buffer, in collaboration with the algorithm described in Section 4.1, that preserves the original data order, ensuring consistent data arrival and mitigating bias caused by asynchronous completion.
- **Heterogeneous Resource Optimization with Elastic CPU Pooling.** To address the under utilization of CPU resources commonly observed in GPU-dominated AI clusters, we implement an elastic CPU pooling strategy. The elastic mechanism isolates and virtualizes idle CPU capacity from the cluster to power logic-intensive tasks such as intensive RL environment interactions and result verification. It effectively amplifies the computational resources available for environment rollouts, enabling massive-scale parallel simulation. Consequently, it reduces the wall-clock time of training iterations while significantly improving the total cost of ownership (TCO) efficiency of the underlying hardware.

6 Evaluations

We conduct systematic evaluations of *ERNIE 5.0* against state-of-the-art models across a wide range of text (Sec. 6.1), vision (Sec. 6.2), and audio benchmarks (Sec. 6.3) with internal evaluation framework, *ERNIE-Eval*¹. Next, we further provide an in-depth analysis of our two key design choices, namely *modality-agnostic expert routing* (Sec. 6.4.1) and *elastic training* (Sec. 6.4.2).

6.1 Evaluation on Language Benchmarks

To comprehensively assess the text-centric capabilities learned during large-scale pre-training and post-training, we evaluate *ERNIE 5.0* on a diverse set of benchmarks covering factual knowledge, reasoning, mathematical problem solving, coding, multilingual understanding, instruction following and agent-oriented tasks. These benchmarks are selected to encompass both core language modeling abilities and advanced reasoning and decision-making skills. This enables a systematic examination of how the unified architecture and training strategies translate into downstream text performance.

Benchmarks

- **Knowledge:** PreciseWikiQA (Bang et al., 2025), PopQA (Mallen et al., 2023), HotPotQA (Yang et al., 2018), ChineseSimpleQA (He et al., 2025), SimpleQA (Wei et al., 2024).
- **General:** MMLU-Pro (Wang et al., 2024d), MMCU (Zeng, 2023), AGIEval (Zhong et al., 2024), MMLU (Hendrycks et al., 2020), MMLU-Redux (Gema et al., 2025), C-Eval (Huang et al., 2023), CMMLU (Li et al., 2024), BBH (Suzgun et al., 2023), WinoGrande (Sakaguchi et al., 2021), Humanity’s Last Exam (HLE) (Phan et al., 2025).

¹For code-related tasks, evaluations are conducted with LiveCodeBench (Jain et al., 2024) and SandboxFusion (Cheng et al., 2024) to ensure reliable and execution-based assessment.

Category	Benchmark	DS V3.2-Exp-Base	Kimi K2-Base	ERNIE 5.0-Base
Knowledge	PreciseWikiQA10-shot	52.60	61.66	74.48
	PopQA10-shot	48.73	51.74	65.24
	HotPotQA10-shot	53.11	57.07	67.08
	ChineseSimpleQA10-shot	74.19	78.29	90.09
General	MMLU-Pro5-shot	68.27	67.19	75.58
	MMCUShot	91.16	90.63	93.92
	AGIEval5-shot	78.20	76.43	80.22
	MMLU5-shot	88.60	88.40	90.58
	MMLU-Redux5-shot	90.50	89.88	92.19
	C-Eval5-shot	90.70	92.49	91.98
	CMMLU5-shot	88.43	90.61	89.69
	BBH5-shot, w/o.-cot (Direct)	73.50	70.07	75.69
STEM	WinoGrande5-shot	88.87	87.61	92.66
	MATH (CoT)8-shot	65.70	65.90	73.89
Coding	GPQA-Diamond5-shot	53.01	48.10	57.30
	LiveCodeBench v61-shot	24.90	26.30	31.94
	HumanEval+0-shot	53.99	70.73	80.86
	MBPP+0-shot	78.77	79.36	79.10
	CRUXEval-I1-shot	59.54	73.64	79.75
Multilingual	CRUXEval-O1-shot	71.28	81.61	84.01
	MMMLU5-shot	70.99	61.50	78.94
	INCLUDE5-shot	77.45	72.29	77.81

Table 1: Comparison of pre-trained models on various language benchmarks. The best performance in each row is highlighted in **bold**.

- **STEM:** MATH (CoT) ([Hendrycks et al., 2021](#)), GPQA-Diamond ([Rein et al., 2024](#)), AIME 2025 ([AIME, 2025](#)), HMMT 2025 ([HMMT Organization, 2025](#)).
- **Coding:** LiveCodeBench v6 (2408 to 2505) ([Jain et al., 2024](#)), HumanEval+ ([Chen, 2021](#)), MBPP+ ([Austin et al., 2021](#)), CRUXEval ([Gu et al., 2024](#)).
- **Multilingual:** MMMLU ([Hendrycks et al., 2020](#)), INCLUDE ([Romanou et al., 2024](#)).
- **Reasoning:** ZebraLogic ([Lin et al., 2025](#)), BBEH ([Kazemi et al., 2025](#)).
- **Instruct Following:** IFEval ([Zhou et al., 2023](#)), Multichallenge ([Deshpande et al., 2025](#)), Multi-IF ([He et al., 2024](#)).
- **Agent:** TAU2-Bench ([Barres et al., 2025](#)), ACEBench ([Chen et al., 2025](#)), BFCL v4 ([Patil et al., 2025](#)), BrowseComp ([Wei et al., 2025](#)), SpreadSheetBench ([Ma et al., 2024](#)).

Evaluation of Pre-trained Models. Table 1 summarizes the pre-training results of *ERNIE 5.0* in comparison with strong open-source baselines on a diverse set of text benchmarks. Across these benchmarks, *ERNIE 5.0* exhibits consistently strong and well-balanced performance in knowledge, reasoning, mathematics, coding, and multilingual tasks:

On knowledge-intensive benchmarks, *ERNIE 5.0-Base* demonstrates clear and substantial advantages over DeepSeek V3.2-Exp-Base (DS V3.2-Exp-Base) and Kimi K2-Base, particularly on both English and Chinese question answering tasks. The large margins observed on these datasets indicate that large-scale unified pre-training effectively consolidates factual knowledge and supports robust retrieval-style reasoning in multilingual settings.

On general reasoning and exam-style benchmarks, *ERNIE 5.0-Base* achieves the best results on a wide range of challenging evaluations, including MMLU-Pro, MMLU, MMCU, AGIEval, MMLU-Redux, and BBH. Notably, the pronounced gains on harder benchmarks such as MMLU-Pro suggest improved reasoning depth and robustness, which can be attributed to the shared backbone and modality-agnostic MoE routing that encourage effective expert specialization.

On STEM tasks, *ERNIE 5.0-Base* consistently outperforms strong baselines on MATH (CoT) and GPQA-Diamond, demonstrating robust multi-step reasoning and solution consistency. These improvements reflect enhanced long-horizon dependency modeling enabled by the deep unified architecture, together with stable optimization under elastic pre-training.

In the coding domain, *ERNIE 5.0-Base* attains state-of-the-art performance on LiveCodeBench v6 and CRUXEval, while remaining competitive on MBPP+, indicating strong generalization of algorithmic and

Category	Benchmark	DS V3.2-Thinking	Gemini 2.5-Pro	GPT-5 (High)	Gemini 3-Pro	ERNIE 5.0
Knowledge	SimpleQA	28.02	54.00	51.30	69.33	74.01
	ChineseSimpleQA	72.37	76.50	75.10	84.08	86.03
General	MMLU-Pro	85.00	86.20	87.10	86.88	83.80
	HLE	25.10	21.60	24.80	37.50	25.81
STEM	GPQA-Diamond	82.40	86.40	85.70	91.90	86.36
	AIME 2025	93.10	88.00	94.60	95.00	89.06
	HMMT 2025	86.67	81.20	93.30	93.33	79.58
Coding	LiveCodeBench v6	81.06	72.90	81.70	86.34	76.21
	HumanEval+	90.80	94.50	92.70	95.12	94.48
	MBPP+	81.48	73.80	83.10	86.21	82.54
Reasoning	ZebraLogic	97.60	92.90	98.80	95.50	96.50
	BBEH	67.04	68.80	69.00	78.80	66.63
Instruction Following	IFEval	91.87	89.50	94.10	92.24	93.35
	MultiChallenge	42.43	51.50	58.30	62.50	65.98
	Multi-IF	71.17	76.10	70.00	81.15	85.56
Agent	TAU2-Bench	80.30	56.20	80.10	85.40	78.79
	ACEBench-en	81.40	80.90	79.30	80.90	87.70
	ACEBench-zh	83.40	87.50	83.60	85.00	89.60
	BFCL v4	61.18	52.30	61.60	68.14	66.47
	BrowseComp-zh	65.00	28.70	61.90	63.67	64.71
	SpreadSheetBench	35.29	27.70	34.00	55.36	40.08

Table 2: Evaluation of post-trained models across a wide range of language benchmarks.

procedural reasoning learned during large-scale pre-training.

On multilingual benchmarks, *ERNIE 5.0-Base* significantly surpasses baselines on MMMLU and IN-CLUE, validating the effectiveness of unified tokenization and shared expert routing in learning robust multilingual representations.

The results demonstrate that the proposed unified architecture and elastic pre-training strategy jointly yield strong generalization across diverse text-centric tasks, providing a solid and versatile foundation for subsequent post-training and deployment.

Evaluation of Post-trained Models. As shown in Table 2, the post-trained *ERNIE 5.0* achieves competitive or leading performance across a broad set of text-centric benchmarks, matching strong open-source and proprietary models on knowledge, instruction-following, coding, and agent-oriented tasks, while maintaining strong general reasoning ability, despite being a natively unified omni model:

On knowledge-intensive benchmarks, such as SimpleQA and ChineseSimpleQA, *ERNIE 5.0* further improves upon the already strong pre-training results. This indicates that post-training effectively enhances factual recall and answer calibration, while preserving the underlying knowledge representations learned during large-scale pre-training.

On general reasoning, mathematical, and coding benchmarks, *ERNIE 5.0* demonstrates stable and competitive performance against strong post-trained baselines. Although Gemini 3-Pro achieves leading results on several particularly challenging benchmarks, including GPQA-Diamond, AIME 2025, HMMT 2025, and LiveCodeBench, *ERNIE 5.0* consistently matches or outperforms DeepSeek-V3.2-Thinking (DS-V3.2-Thinking), Gemini 2.5-Pro, and GPT-5 (High) on most of the benchmarks. This behavior aligns well with the pre-training observations, suggesting that *ERNIE 5.0* emphasizes robust and balanced capability development rather than aggressive optimization toward extreme reasoning or competition-style tasks.

On instruction-following benchmarks, *ERNIE 5.0* shows clear advantages on multi-instruction evaluations, achieving the best performance on MultiChallenge and Multi-IF, as well as near-top results on IFEval. These results suggest that post-training effectively strengthens instruction compliance and compositional instruction understanding, complementing the strong pre-trained foundation.

On agent-oriented benchmarks, *ERNIE 5.0* demonstrates competitive and often leading performance, particularly on ACEBench (in both English and Chinese) and BrowseComp-zh. While Gemini 3-Pro excels on certain tool-intensive tasks, *ERNIE 5.0* exhibits strong generalization across diverse agent scenarios, indicating practical usability in complex interactive environments.

Category	Benchmark	Qwen3-VL Thinking	Gemini 2.5-Pro	GPT-5 (High)	Gemini 3-Pro	ERNIE 5.0
STEM & Reasoning	MMMU-Pro	68.28	68.80	78.40	81.00	68.63
	MathVista	86.80	82.70	82.10	89.20	84.80
	MathVerse	83.96	86.27	84.19	91.62	85.13
	MathVision	71.84	73.30	78.06	87.27	74.34
	VisualPuzzle	57.01	61.51	57.75	71.48	64.82
	VisuLogic	31.93	32.80	29.80	37.60	32.00
Document Understanding	VLMAreBlind	75.13	76.54	69.60	80.83	91.38
	ChartQA	84.60	84.08	78.24	89.44	87.80
	AI2D	96.96	97.09	95.63	97.70	96.89
	DocVQA _{val}	95.44	91.43	94.16	90.70	95.45
	OCRBench	863	866	804	909	878
	ChartXiv-RQ	63.00	67.80	81.10	81.40	67.10
General VQA	ChartXiv-DQ	92.92	93.38	91.17	95.95	89.05
	SimpleVQA	62.37	68.19	55.84	74.06	67.64
	HallusionBench	64.01	63.70	66.58	73.48	63.87
	MMStar	76.88	77.50	82.10	82.96	75.54
	BLINK	66.60	70.60	70.39	77.49	70.02
	CV-Bench	87.57	84.87	84.99	90.07	87.19
Video Understanding	CountBench	92.67	91.00	88.88	97.35	96.54
	VideoMME _(w sub)	80.97	86.90	87.36	88.40	81.35
	Video-MMMU	80.00	83.60	84.60	87.60	81.11
	MMVU	71.10	76.10	87.34	76.30	72.24

Table 3: Evaluation of post-trained models across a wide range of vision benchmarks.

The comparison between pre-training and post-training results highlights a consistent capability trajectory of *ERNIE 5.0*. Strengths in factual knowledge and general robustness established during pre-training are almost retained after post-training, while instruction following and agent capabilities are significantly enhanced. Although *ERNIE 5.0* remains competitive across general reasoning, math, and coding tasks, a moderate gap persists on the most challenging reasoning benchmarks compared to models such as Gemini 3-Pro. Our future work will further leverage architectural design and advanced training strategies to better support complex, long-horizon reasoning.

6.2 Evaluation on Vision Benchmarks

To evaluate the visual understanding and generation capabilities enabled by native multimodal training, we assess *ERNIE 5.0* on various image-centric and video-centric benchmarks. These evaluations cover visual reasoning, document understanding, visual question answering, video understanding, as well as image and video generation, providing a holistic assessment of how the unified multi-modal architecture manifests as performances across perception, reasoning, and generation tasks.

Visual Understanding & Generation Benchmarks

- **STEM and Reasoning:** MMMU-Pro (Yue et al., 2025a), MathVista (Lu et al., 2023), MathVerse (Zhang et al., 2024a), MathVision (Wang et al., 2024b), VisualPuzzle (Song et al., 2025), VisuLogic (Xu et al., 2025b), VLMAreBlind (Rahmanzadehgervi et al., 2024), MMMU (Yue et al., 2024).
- **Document Understanding:** ChartQA (Masry et al., 2022), AI2D (Kembhavi et al., 2016), DocVQA(val) (Mathew et al., 2021), OCRBench (Liu et al., 2023b), ChartXiv-RQ, ChartXiv-DQ (Wang et al., 2024e).
- **General VQA:** SimpleVQA (Cheng et al., 2025), HallusionBench (Guan et al., 2024), MMStar (Chen et al., 2024b), BLINK (Fu et al., 2024), CV-Bench (Zhu et al., 2025), CountBench (Paiss et al., 2023).
- **Video Understanding:** VideoMME (Fu et al., 2025), Video-MMMU (Hu et al., 2025), MMVU (Zhao et al., 2025c).
- **Image and Video Generation:** GenEval (Ghosh et al., 2023), VBench (Huang et al., 2024).

Model	MathVista	MathVerse	MathVision	VisualPuzzle	VisuLogic	ChartQA	AI2D
ERNIE 5.0-Base	84.40	81.45	68.75	54.24	28.50	87.68	96.02
	OCRBench	ChartXiv-RQ	ChartXiv-DQ	SimpleVQA	MMStar	BLINK	CV-Bench
	875	62.70	90.35	61.54	74.07	64.12	86.96

Table 4: Evaluation of our pre-trained model across a wide range of vision benchmarks. There are few publicly available visual base models and their results, here we only report the results of our model.

Benchmark	Nano Banana Pro	Seedream 4.0	GPT-Image	Qwen-Image	ERNIE 5.0-Base	ERNIE 5.0
GenEval	89.0	85.4	84.0	91.0	88.4	90.1

Table 5: Comparison of the image generation ability on GenEval against specialized models.

Benchmark	Metric	HunyuanVideo-1	Wan2.1-14B-0725	Veo3	ERNIE 5.0-Base	ERNIE 5.0
VBench	Quality	85.07	85.59	85.70	84.14	84.40
	Semantic	76.88	76.11	82.49	82.31	83.40
	Overall	83.43	83.69	85.06	83.78	84.20

Table 6: Comparison of the video generation ability on VBench against specialized models.

Evaluation of Visual Understanding Table 3 and Table 4 illustrate the multimodal capabilities of *ERNIE 5.0* under the instruction-tuned setting (*ERNIE 5.0*) and its native multimodal foundation (*ERNIE 5.0-Base*), respectively. This comparison enables a clearer understanding of how instruction tuning interacts with a unified multimodal pre-training paradigm.

ERNIE 5.0-Base exhibits strong and well-rounded performances across visual reasoning, document understanding, and general VQA benchmarks, despite the absence of instruction tuning. It suggests that core multimodal perception and cross-modal reasoning abilities are largely acquired during pre-training, benefiting from early fusion of visual and textual tokens and modality-agnostic representation learning. This behavior differs from modular or late-fusion designs, where comparable capabilities typically emerge only after extensive task-specific tuning. Based on this foundation, *ERNIE 5.0* consistently improves performance on most benchmarks, particularly on tasks requiring explicit reasoning, compositional understanding, and robust visual-language alignment.

Compared with other strong multimodal baselines, *ERNIE 5.0* achieves competitive or superior results on a wide range of reasoning, document understanding, and general VQA tasks, while maintaining balanced performance rather than optimizing for individual benchmarks. It indicates that instruction tuning in *ERNIE 5.0* primarily refines reasoning and alignment behaviors, instead of compensating for perceptual capacity.

In document understanding scenarios, the relatively strong results of *ERNIE 5.0-Base* already demonstrate effective layout-aware reasoning and text extraction, while *ERNIE 5.0* further enhances question understanding and structured reasoning over complex visual documents. A similar trend is observed in general VQA and video understanding benchmarks, where instruction tuning improves robustness and temporal reasoning without altering the underlying architecture.

Evaluation of Visual Generation We evaluate *ERNIE 5.0* on widely used benchmarks for both image and video generation, and compare it with leading commercial and open-source models.

For image generation, as shown in Table 5, *ERNIE 5.0* achieves competitive performance on the GenEval benchmark. In particular, *ERNIE 5.0* performs on par with state-of-the-art commercial systems such as Nano-Banana Pro (DeepMind, 2025b) and Qwen-Image (Wu et al., 2025), and is comparable to GPT-Image (OpenAI, 2024) and Seedream 4.0 (Seedream et al., 2025). It demonstrates that *ERNIE 5.0* is capable of producing high-aesthetic images with strong semantic alignment and fine-grained visual details.

For video generation, Table 6 summarizes the results on the VBench benchmark. *ERNIE 5.0* achieves the best performance on VBench-Semantic, surpassing strong commercial models such as Veo3 (DeepMind, 2025c), which indicates its superior semantic alignment in video generation. This advantage aligns with the unified multimodal architecture, where high-level semantic representations are effectively transferred to generative tasks. Meanwhile, *ERNIE 5.0* remains competitive on the overall and quality metrics, demonstrating solid visual fidelity and temporal consistency. In addition, *ERNIE 5.0* performs on par with leading open-source models such as HunyuanVideo-1 (Kong et al., 2024) and Wan2.1-14B-0725 (Wan

Benchmark	Kimi Audio	GPT-4o -Audio	Qwen3-Omni -Instruct	LongCat-Flash -Omni	Gemini-3 -Pro	ERNIE 5.0-Base	ERNIE 5.0
<i>Automatic Speech Recognition (↓)</i>							
AISHELL-1	0.60	3.52	0.84	0.63	3.04	0.75	0.31
AISHELL-2	2.56	4.26	2.34	2.78	4.98	2.90	2.64
WenetSpeech net meeting	5.37 6.28	15.30 32.27	4.69 5.89	6.09 6.69	10.94 12.08	11.57 22.85	7.27 7.36
LibriSpeech clean other	1.28 2.42	1.39 3.75	1.22 2.48	1.57 4.01	2.74 4.40	1.47 3.73	1.16 2.61
Fleurs-en	4.44	3.32	2.72	5.02	3.63	4.39	3.14
Fleurs-zh	2.69	2.44	2.20	3.99	4.98	1.58	0.83
<i>VoiceBench (↑)</i>							
AlpacaEval	4.46	4.73	4.74	4.94	4.80	4.65	4.62
CommonEval	3.97	4.37	4.54	4.32	4.68	4.39	3.74
SD-QA	63.12	90.10	76.90	82.46	94.39	86.44	77.58
MMSU	62.17	78.90	69.00	81.95	92.16	76.61	84.68
OpenBookQA	83.52	87.90	89.70	93.41	96.26	88.35	92.97
IFEval	61.10	66.81	77.80	77.99	87.45	71.17	72.67
AdvBench	100.00	99.23	99.30	100.00	98.46	98.65	99.23
<i>Audio Understanding (↑)</i>							
MMAU	65.20	68.40	77.50	75.90	80.80	80.80	80.40
TUT2017	65.25	20.74	40.74	65.43	61.42	57.65	68.09
CochlScene	80.42	34.94	43.03	70.02	74.60	75.24	82.77
ClothoAQA	72.21	61.87	75.16	72.83	74.41	65.70	73.68
VocalSound	94.85	82.37	91.60	92.76	92.01	91.48	90.73

Table 7: Comparison of the automatic speech recognition, speech-to-text dialogue, and audio understanding abilities against specialized models. Automatic speech recognition uses word error rate (WER) as the metric (lower is better), while other tasks use accuracy or score (higher is better).

et al., 2025), indicating that the proposed architecture provides a strong video generation foundation even before post-training. Overall, these results validate the effectiveness of *ERNIE 5.0* in producing semantically accurate and visually coherent videos.

6.3 Evaluation on Audio Benchmarks

To evaluate the audio understanding and generation capabilities enabled by native multimodal training, we assess *ERNIE 5.0* on a diverse suite of speech- and audio-centric benchmarks. These evaluations cover automatic speech recognition, speech-based dialogue, general audio understanding, as well as speech generation, providing a comprehensive view of how the unified multimodal architecture translates into robust audio perception, semantic understanding, and generation performance across across diverse linguistic and acoustic environments.

Audio Understanding & Generation Benchmarks

- **Automatic Speech Recognition:** AISHELL-1 (Bu et al., 2017), AISHELL-2 (Du et al., 2018), WenetSpeech (Zhang et al., 2022), LibriSpeech (Panayotov et al., 2015), Fleurs (Conneau et al., 2023).
- **Voice Chatting:** VoiceBench (Chen et al., 2024f).
- **Audio Understanding:** MMAU (Sakshi et al., 2024), CochlsScene (Jeong & Park, 2022), TUT2017 (Mesaros et al., 2016), ClothoAQA (Lipping et al., 2022), VocalSound (Gong et al., 2022).
- **Speech Generation:** Seed-TTS (Anastassiou et al., 2024).

Evaluation of Audio Understanding We compare *ERNIE 5.0* with other leading specialist and generalist models on ASR, voice-chatting, and other audio understanding benchmarks. As shown in Table 7 and Table 8, *ERNIE 5.0* exhibits state-of-the-art or competitive performance across speech recognition, speech-to-text dialogue, audio understanding and speech generation benchmarks.

On automatic speech recognition (ASR) tasks, *ERNIE 5.0* achieves low word error rate (WER) on both Chinese and English benchmarks, including AISHELL and LibriSpeech, indicating robust performance

Benchmark	Model	Performance	
SEED-TTS (\downarrow) <i>test-zh test-en</i>	Seed-TTS _{ICL}	1.11	2.24
	Seed-TTS _{RL}	1.00	1.94
	MaskGCT	2.27	2.62
	E2 TTS	1.97	2.19
	F5-TTS	1.56	1.83
	Spark TTS	1.20	1.98
	CosyVoice 2	1.45	2.57
	CosyVoice 3	0.71	1.45
	Qwen2.5-Omni	1.42	2.33
	Qwen3-Omni	1.07	1.39
<i>ERNIE 5.0-Base</i>		3.41	2.44
<i>ERNIE 5.0</i>		1.35	1.54

Table 8: Comparison of the text-to-speech ability on SEED-TTS against specialized models. Word Error Ratio (WER) is used to evaluate content consistency, the lower the better.

across languages and acoustic conditions. While some models achieve stronger results on specific datasets such as WenetSpeech, *ERNIE 5.0* maintains consistently stable performance across a wide range of ASR benchmarks, reflecting good generalization.

On VoiceBench, *ERNIE 5.0* remains competitive on tasks such as MMSU and OpenBookQA, reflecting its ability to handle speech-based interaction and knowledge-grounded reasoning.

For audio understanding tasks, *ERNIE 5.0* performs favorably on acoustic scene and environmental sound benchmarks such as TUT2017 and CochlScene, while achieving comparable results on more diverse benchmarks including MMAU and ClothoAQA. Overall, the model shows consistent performance across a wide range of non-speech and general audio understanding tasks, indicating its ability to capture audio semantics beyond speech content.

Evaluation of Audio Generation On SEED-TTS, *ERNIE 5.0* achieves competitive content consistency on both Chinese and English test sets (*test-zh*, *test-en*), performing comparably to recent audio-language models such as Qwen3-Omni. While specialist TTS systems (e.g., CosyVoice-3) obtain stronger results, *ERNIE 5.0* demonstrates reliable content preservation without task-specific TTS optimization.

Taken together, the results in Table 7 indicate that *ERNIE-5.0* is a competitive unified audio-language model that balances speech recognition accuracy, speech interaction, general audio understanding and speech generation within a single framework.

6.4 Discussion

To better understand the behavior of *ERNIE 5.0* in large-scale multimodal training, we dissect two pivotal architectural designs: modality-agnostic expert routing and elastic training. The following subsections analyze these two components and their effects on efficiency and scalability.

6.4.1 Modality-Agnostic Expert Routing

ERNIE 5.0 employs a modality-agnostic expert routing mechanism, without introducing modality-specific parameters or routing rules. All inputs, including text, image, video, and audio, are processed by the same routing network and share a common pool of experts. In this discussion, we examine the expert routing behavior of the MoE model, with particular attention to modality-specific utilization patterns, cross-modality expert overlap, and routing balance across layers.

Perspective of Expert Utilization and Specialization Figure 8 reveals a clearly non-uniform expert activation distribution, suggesting that different experts still play distinct functional roles under the modality-agnostic expert routing setting. A subset of experts is repeatedly activated across text, image, video, and audio inputs, whereas the remaining experts exhibit strong modality-specific activation patterns. Expert activations for image, video, and audio inputs are noticeably more concentrated than those for text-only inputs. From a task-level perspective, visual generation and audio-related tasks lead to more concentrated expert activations than text and visual understanding tasks.

Perspective of Cross-Modality Expert Collaboration Figure 9 further analyzes cross-modality expert collaboration by reporting the Intersection over Union (IoU) of the top 25% of experts by activation frequency for each modality across different layers. From the overlap patterns, text shows a higher degree

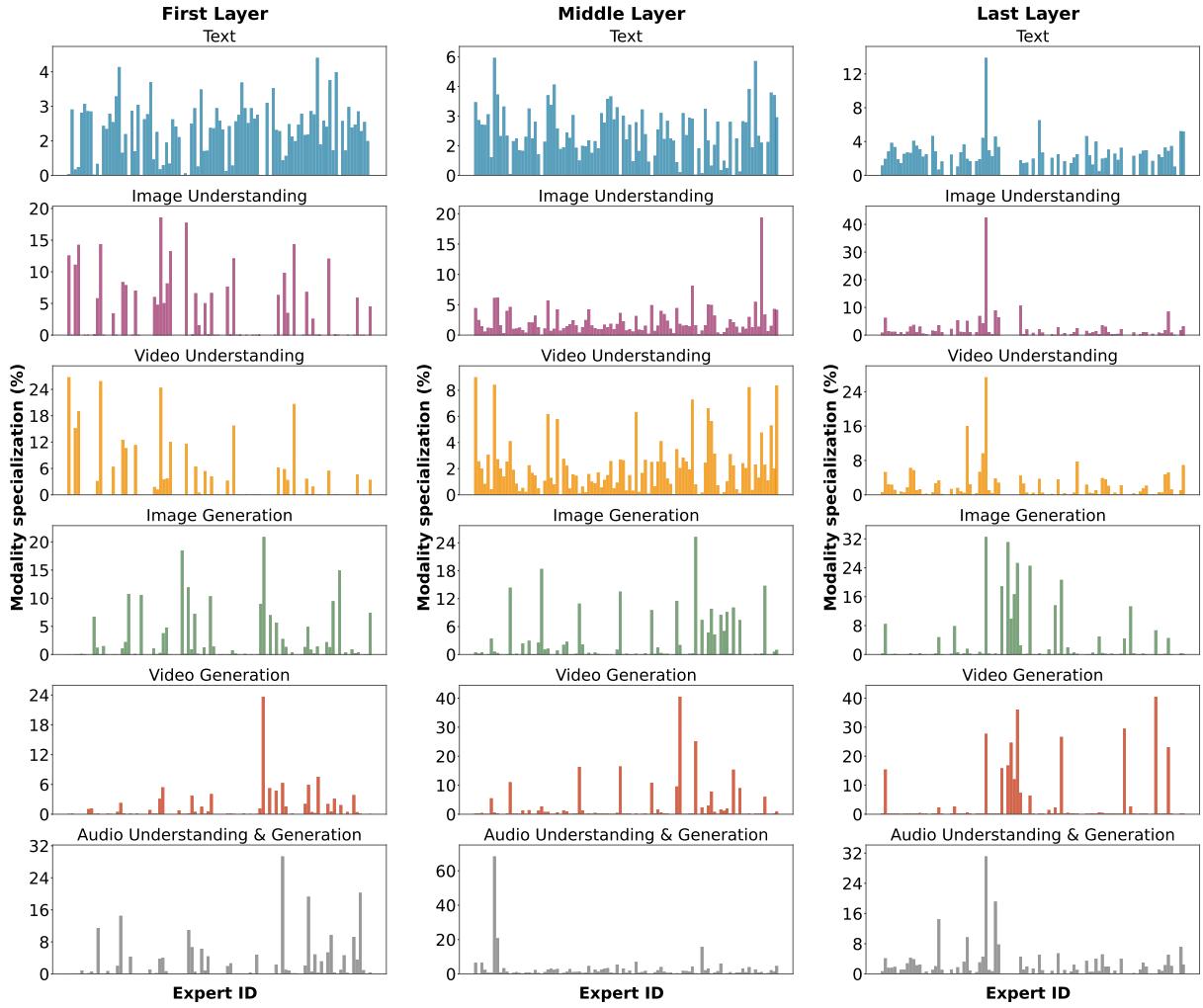


Figure 8: Visualization of expert utilization patterns across modalities and tasks for the representative first, middle and last layers. The y-axis indicates the frequency of expert activation.

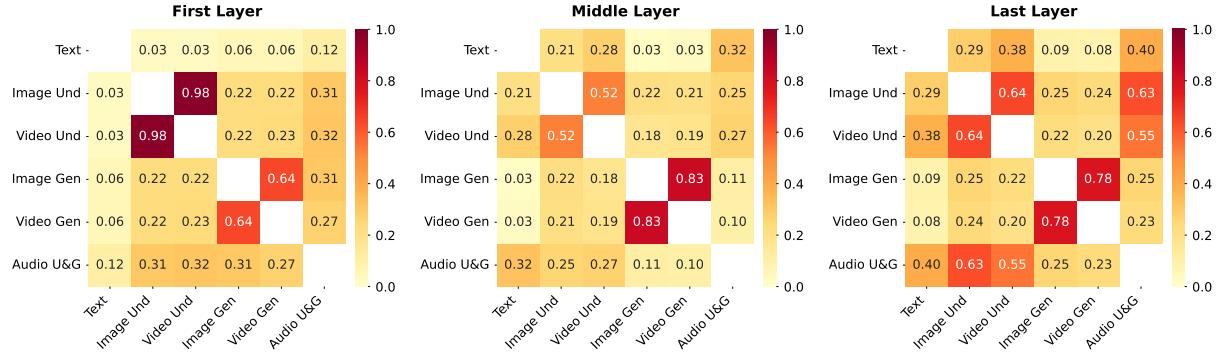


Figure 9: Visualization of expert collaboration across modalities and tasks for the representative first, middle and last layers, using the Intersection over Union (IoU) of the top 25% most frequently activated experts for each modality.

of co-activation with audio than with other modalities, and this overlap with image and video becomes increasingly pronounced in deeper layers. This trend suggests that multimodal representations gradually shift from low-level modality-specific features toward higher-level unified semantic features, enabling stronger collaboration with text at deeper layers. For visual modalities, image and video understanding tasks exhibit high expert overlap, and a similar pattern is observed for image and video generation tasks, which aligns with the architectural design in which an image is regarded as a single-frame video. In contrast, the overlap between visual understanding and generation remains relatively low, with no clear

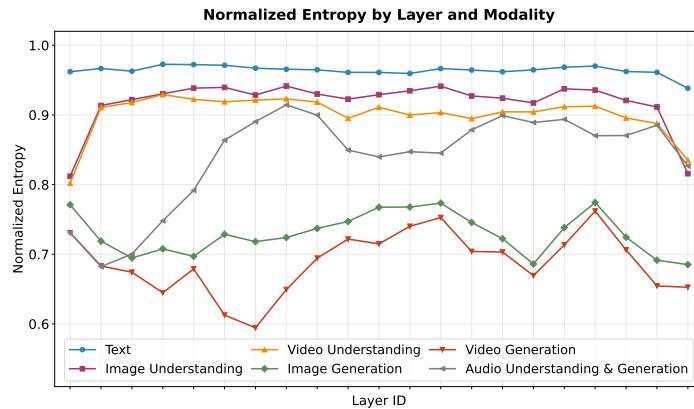


Figure 10: Visualization of the load balancing across layers and modalities, using the normalized entropy (NE) metric of expert routing.

Training Configuration	Inference Configuration (# of Layers)	Validation Loss
Baseline Layers=16	Layers=16	1.945
Elastic Depth Layers $\in [1, 16]$	Layers=16 Layers=12	1.941 2.137

Table 9: Validation loss under elastic depth training and inference configurations for small-scale MoE model. Elastic depth slightly improves full-depth performance and yields reduced-depth sub-networks with predictable degradation, enabling flexible deployment without retraining.

preference toward either direction.

Perspective of Load Balancing Figure 10 reports the normalized entropy (NE) of expert routing across layers and modalities, defined as $NE = -\frac{\sum_{i=1}^N p_i \log(p_i)}{\log N}$, which provides a quantitative measure of expert load balance. Here, N is the number of experts and p_i is the fraction of tokens routed to i -th expert. The value of NE lies in $[0, 1]$, where larger values correspond to more uniform expert utilization.

From a layer-wise perspective, the text modality exhibits consistently high and stable normalized entropy across almost all layers, with only a mild drop observed at the final layer. Notably, the first layer does not show severe imbalance, contradicting the common assumption that early MoE layers require dense designs to maintain routing stability (Liu et al., 2024).

For visual understanding, expert utilization becomes less balanced at both the lowest and highest layers, while the middle layers maintain relatively uniform routing, suggesting that visual perception benefits from balanced expert sharing during intermediate feature abstraction stages. In contrast, visual generation and audio-related tasks display a different trend: the first layer shows a moderate level of balance, followed by a decrease in entropy in lower layers, a partial recovery in lower-to-mid layers, and a fluctuating drop at the higher layers. This pattern indicates alternating phases of expert specialization and re-integration along depth for generative and audio-centric tasks.

Taken together, these analyses reveal clear expert utilization patterns across layers, modalities, and tasks in *ERNIE 5.0*. Although the routing mechanism is modality-agnostic during training, pronounced expert specialization across modalities still emerges, which suggests that the router could capture modality structure and allocate expert capacity in a self-driven manner. From a system perspective, the unified routing mechanism simplifies the overall design and avoids manual expert partitioning across modalities. Beyond empirical observations, the routing behaviors and load-balancing characteristics also offer practical guidance for future model design. For example, layer-aware expert allocation, adaptive balancing strategies, and modality-shared expert configurations may have the potential to become effective principles for building the next generation of native multimodal architectures.

6.4.2 Elastic Training

Elastic training plays a central role in enabling *ERNIE 5.0* to adapt to different compute, memory, and latency constraints. In this section, we analyze elastic training strategies through controlled ablation studies on a small-scale MoE model and report key findings on *ERNIE 5.0*. Unless otherwise specified, all

Training Configuration	Inference Configuration (# of Experts)	Validation Loss
Baseline Experts=64	Experts=64	1.957
Elastic Width Experts $\in \{64, 32\}$	Experts=64 Experts=32	1.964 2.218

Table 10: Validation loss under elastic width configurations for the small-scale MoE model. Elastic width introduces only minor degradation at full capacity, while reduced-width models remain usable, supporting deployment under constrained parameter budgets.

Training Configuration	Inference Configuration (MoE Routing Top- k)	Validation Loss
Baseline Top- k = 8	Top- k = 8	1.945
Elastic Sparsity Top- k $\in [1, 8]$	Top- k = 8	1.969
	Top- k = 4	1.971
	Top- k = 2	2.003
	Top- k = 1	2.175

Table 11: Validation loss under different routing sparsity levels for the small-scale MoE model. Elastic sparsity allows stable inference across varying routing budgets, with graceful performance degradation under aggressive sparsification.

elastic configurations are derived from the same pretrained checkpoint.

Controlled-Scale Experiments We conduct ablation studies using a small-scale MoE model with 64 experts, 454M activated parameters, and 3.2B total parameters. The model is trained on 250B tokens with a default routing configuration of $top-k = 8$. Validation loss on held-out data is reported throughout. Based on this experimental setup, we analyze elasticity along the three dimensions separately.

- **Elastic Depth.** We first study elastic depth by varying the number of active transformer layers during training. Consistent with the configuration of *ERNIE 5.0*, 80% of training samples retain the full-depth model, while the remaining 20% reduced-depth sub-networks. As reported in Table 9, elastic depth slightly improves full-depth performance, which indicates a regularization effect introduced by occasional layer dropping. Moreover, reduced-depth sub-networks exhibit smooth and predictable performance degradation, suggesting that intermediate representations remain robust under layer removal. Overall, elastic depth offers a low-risk and cost-efficient mechanism for deriving multiple deployable sub-models from a single pretrained checkpoint.
- **Elastic Width.** We next evaluate elasticity over MoE width. During training, 80% of samples retain all 64 experts, while the remaining 20% samples configurations with 32 experts. Table 10 shows that elastic width introduces only a slight degradation at full capacity. Although reduced-width configurations incur minor performance fluctuations of larger models, they remain functional without retraining, enabling deployment under strict memory constraints.
- **Elastic Sparsity.** We also analyze elastic sparsity by varying routing top- k at inference. During training, most instances are equipped with default configuration ($k = 8$), while the remaining instances are trained with reduced activated experts. As shown in Table 11, elastic top- k training incurs a modest degradation under the full-activation configuration, while enabling stable and effective inference under substantially reduced routing budgets.

Scaling to *ERNIE 5.0* Beyond controlled-scale experiments, we further evaluate the effectiveness of elastic training on the pre-trained experimental model *ERNIE 5.0-Exp-Base* and its post-trained model *ERNIE 5.0-Exp*. Starting from *ERNIE 5.0-Exp*, we first analyze elastic sparsity by reducing the routing top- k to 25% during inference. As summarized in Table 12, *ERNIE 5.0-Exp-ES_{25.0%}* retains comparable performance across a wide range of text and visual benchmarks. The resulting accuracy drop remains minor, while the reduction in routing sparsity brings substantial efficiency gains, providing a *more than 15% improvement in decoding speed*. Moreover, we jointly activate elasticity along depth, width, and sparsity, deriving a compact pre-trained model from *ERNIE 5.0-Exp-Base*, which operates with only 53.7% of the activated parameters and 35.8% of the total parameters. Using the same data and training strategy for mid-training and post-training, we obtain the post-trained model, *ERNIE 5.0-Exp-EA_{35.8%}*. Despite its substantially reduced computational footprint, the elastic variant achieves competitive performance across benchmarks, attaining an average score of 75.17 compared to 75.55 for the full *ERNIE 5.0-Exp*. Moreover, strong robustness on challenging reasoning and perception tasks, such as ZebraLogic and VisualPuzzle, indicates that elastic training effectively mitigates the performance degradation typically

Model	AVG.	ZebraLogic	LiveCodeBench v6	TAU2	MMMU	MathVista	VisualPuzzle	SimpleVQA
ERNIE 5.0-Exp	75.55	95.00	73.35	79.35	74.11	83.70	59.93	63.40
ERNIE 5.0-Exp-ES _{25.0%}	74.43	94.10	70.70	77.34	73.78	84.90	57.98	62.19
ERNIE 5.0-Exp-EA _{35.8%}	75.17	95.20	70.93	77.23	75.11	84.50	60.39	62.86

Table 12: Performance comparison between the full *ERNIE 5.0*-Exp and its elastic variants under reduced routing sparsity (*ERNIE 5.0*-Exp-ES) and all elastic configurations (*ERNIE 5.0*-Exp-EA). Elastic sparsity preserves comparable accuracy with improved decoding efficiency, while the fully elastic model achieves competitive performance despite substantially reduced activated computation and parameter usage.

associated with aggressive reductions in computation and parameters.

Overall, these results highlight that elasticity is not merely as a post-hoc compression technique, but a principled training paradigm. By jointly optimizing depth, width, and sparsity during pre-training and, the elastic model learns to redistribute representational capacity across layers and modalities, leading to favorable performance-efficiency trade-offs.

7 Conclusion

We introduce *ERNIE 5.0*, a natively unified foundation model that integrates multimodal understanding and generation of text, image, video, and audio with a shared next-group-of-tokens prediction objective. To the best of our knowledge, *ERNIE 5.0* represents the first realization of multimodal understanding and generation within a unified, trillion-level autoregressive framework. An ultra-sparse mixture-of-experts architecture with modality-agnostic expert routing enables scalable cross-modal modeling without relying on modality-specific designs. We explore a novel elastic training paradigm that supports flexible deployment configurations within a single pre-trained model, which is proved successful to maintain both training efficiency and performances. We also address key challenges in reinforcement learning for large-scale multimodal models and ensure pos-training stability and efficiency. Extensive experiments demonstrate competitive and balanced performance across modalities. Overall, the results indicate that autoregressive unified multimodal and elastic pre-training provides a scalable pathway toward the next generation of foundational models.

8 Contributors

Haifeng Wang, Hua Wu, Tian Wu, Yu Sun, Jing Liu, Dianhai Yu, Yanjun Ma, Jingzhou He, Zhongjun He, Dou Hong, Qiwen Liu, Shuohuan Wang, Junyuan Shang, Zhenyu Zhang, Yuchen Ding, Jinle Zeng, Jiabin Yang, Liang Shen, Ruibiao Chen, Weichong Yin, Siyu Ding, Dai Dai, Shikun Feng, Siqi Bao, Bolei He, Yan Chen, Zhenyu Jiao, Ruiqing Zhang, Zeyu Chen, Qingqing Dang, Kaipeng Deng, Jiajun Jiang, Enlei Gong, Guoxia Wang, Yanlin Sha, Yi Liu, Yehan Zheng, Weijian Xu, Jiaxiang Liu, Zengfeng Zeng, Yingqi Qu, Zhongli Li, Zhengkun Zhang, Xiyang Wang, Zixiang Xu, Xinshao Xu, Zhengjie Huang, Dong Wang, Bingjin Chen, Yue Chang, Xing Yuan, Shiwei Huang, Qiao Zhao, Xinzhe Ding, Shuangshuang Qiao, Baoshan Yang, Bihong Tang, Bin Li, Bingquan Wang, Binhan Tang, Binxiong Zheng, Bo Cui, Bo Ke, Bo Zhang, Bowen Zhang, Boyan Zhang, Boyang Liu, Caiji Zhang, Can Li, Chang Xu, Chao Pang, Chao Zhang, Chaoyi Yuan, Chen Chen, Cheng Cui, Chenlin Yin, Chun Gan, Chunguang Chai, Chuyu Fang, Cuiyun Han, Dan Zhang, Danlei Feng, Danxiang Zhu, Dong Sun, Dongbo Li, Dongdong Li, Dongdong Liu, Dongxue Liu, Fan Ding, Fan Hu, Fan Li, Fan Mo, Feisheng Wu, Fengwei Liu, Gangqiang Hu, Gaofeng Lu, Gaopeng Yong, Gexiao Tian, Guan Wang, Guangchen Ni, Guangshuo Wu, Guanzhong Wang, Guihua Liu, Guishun Li, Haibin Li, Haijian Liang, Haipeng Ming, Haisu Wang, Haiyang Lu, Haiye Lin, Han Zhou, Hangting Lou, Hanwen Du, Hanzhi Zhang, Hao Chen, Hao Du, Hao Liu, Hao Zhou, Haochen Jiang, Haodong Tian, Haoshuang Wang, Haozhe Geng, Heju Yin, Hong Chen, Hongchen Xue, Hongen Liu, Honggeng Zhang, Hongji Xu, Hongwei Chen, Hongyang Zhang, Hongyuan Zhang, Hua Lu, Huan Chen, Huan Wang, Huang He, Hui Liu, Hui Zhong, Huibin Ruan, Jiafeng Lu, Jiage Liang, Jiahao Hu, Jiahao Hu, Jiajie Yang, Jialin Li, Jian Chen, Jian Wu, Jianfeng Yang, Jianguang Jiang, Jianhua Wang, Jianye Chen, Jiaodi Liu, Jiarui Zhou, Jiawei Lv, Jiaxin Zhou, Jiaxuan Liu, Jie Han, Jie Sun, Jiefan Fang, Jihan Liu, Jihua Liu, Jing Hu, Jing Qian, Jing Yan, Jingdong Du, Jingdong Wang, Jingjing Wu, Jingyong Li, Jinheng Wang, Jinjin Li, Jinliang Lu, Jinlin Yu, Jinnan Liu, Jixiang Feng, Jiyi Huang, Jiyuan Zhang, Jun Liang, Jun Xia, Jun Yu, Junda Chen, Junhao Feng, Junhong Xiang, Junliang Li, Kai Liu, Kailun Chen, Kairan Su, Kang Hu, Kangkang Zhou, Ke Chen, Ke Wei, Kui Huang, Kun Wu, Kunbin Chen, Lei Han, Lei Sun, Lei Wen, Linghui Meng, Linhao Yu, Liping Ouyang, Liwen Zhang, Longbin Ji, Longzhi Wang, Meng Sun, Meng Tian, Mengfei Li, Mengqi Zeng, Mengyu Zhang, Ming Hong, Mingcheng Zhou, Mingming Huang, Mingxin Chen, Mingzhu Cai, Naibin Gu, Nemin Qiu, Nian Wang, Peng Qiu, Peng

Zhao, Pengyu Zou, Qi Wang, Qi Xin, Qian Wang, Qiang Zhu, Qianhui Luo, Qianwei Yang, Qianyue He, Qifei Wu, Qinrui Li, Qiwen Bao, Quan Zhang, Quanxiang Liu, Qunyi Xie, Rongrui Zhan, Rufeng Dai, Rui Peng, Ruian Liu, Ruihao Xu, Ruijie Wang, Ruixi Zhang, Ruixuan Liu, Runsheng Shi, Ruting Wang, Senbo Kang, Shan Lu, Shaofei Yu, Shaotian Gong, Shenwei Hu, Shifeng Zheng, Shihao Guo, Shilong Fan, Shiqin Liu, Shiwei Gu, Shixi Zhang, Shuai Yao, Shuang Zhang, Shuangqiao Liu, Shuhao Liang, Shuwei He, Shuwen Yang, Sijun He, Siming Dai, Siming Wu, Siyi Long, Songhe Deng, Suhui Dong, Suyin Liang, Teng Hu, Tianchan Xu, Tianliang Lv, Tianmeng Yang, Tianyi Wei, Tiezhu Gao, Ting Sun, Ting Zhang, Tingdan Luo, Wei He, Wei Luan, Wei Yin, Wei Zhang, Wei Zhou, Weibao Gong, Weibin Li, Weicheng Huang, Weichong Dang, Weiguo Zhu, Weilong Zhang, Weiqi Tan, Wen Huang, Wenbin Chang, Wenjing Du, Wenlong Miao, Wenpei Luo, Wenquan Wu, Xi Shi, Xi Zhao, Xiang Gao, Xiangguo Zhang, Xiangrui Yu, Xiangsen Wang, Xiangzhe Wang, Xianlong Luo, Xianying Ma, Xiao Tan, Xiaocong Lin, Xiaofei Wang, Xiaofeng Peng, Xiaofeng Wu, Xiaojian Xu, Xiaolan Yuan, Xiaopeng Cui, Xiaotian Han, Xiaoxiong Liu, Xiaoxu Fei, Xiaoxuan Wu, Xiaoyu Wang, Xiaoyu Zhang, Xin Sun, Xin Wang, Xinhui Huang, Ximming Zhu, Xintong Yu, Xinyi Xu, Xinyu Wang, Xiuxian Li, XuanShi Zhu, Xue Xu, Xueying Lv, Xuhong Li, Xulong Wei, Xuyi Chen, Yabing Shi, Yafeng Wang, Yamei Li, Yan Liu, Yanfu Cheng, Yang Gao, Yang Liang, Yang Wang, Yang Wang, Yang Yang, Yanlong Liu, Yannian Fu, Yanpeng Wang, Yanzheng Lin, Yao Chen, Yaozong Shen, Yaqian Han, Yehua Yang, Yekun Chai, Yesong Wang, Yi Song, Yichen Zhang, Yifei Wang, Yifeng Guo, Yifeng Kou, Yilong Chen, Yilong Guo, Yiming Wang, Ying Chen, Ying Wang, Yingsheng Wu, Yingzhan Lin, Yinqi Yang, Yiran Xing, Yishu Lei, Yixiang Tu, Yiyan Chen, Yong Zhang, Yonghua Li, Yongqiang Ma, Yongxing Dai, Yongyue Zhang, Yu Ran, Yu Sun, Yu-Wen Michael Zhang, Yuang Liu, Yuanle Liu, Yuanyuan Zhou, Yubo Zhang, Yuchen Han, Yucheng Wang, Yude Gao, Yuedong Luo, Yuehu Dong, Yufeng Hu, Yuhui Cao, Yuhui Yun, Yukun Chen, Yukun Gao, Yukun Li, Yumeng Zhang, Yun Fan, Yun Ma, Yunfei Zhang, Yunshan Xie, Yuping Xu, Yuqin Zhang, Yuqing Liu, Yurui Li, Yuwen Wang, Yuxiang Lu, Zefeng Cai, Zelin Zhao, Zelun Zhang, Zenan Lin, Zehao Dong, Zhaowu Pan, Zhaoyu Liu, Zhe Dong, Zhe Zhang, Zhen Zhang, Zhengfan Wu, Zhengrui Wei, Zhengsheng Ning, Zhenxing Li, Zhenyu Li, Zhenyu Qian, Zhenyun Li, Zhi Li, Zhichao Chen, Zhicheng Dong, Zhida Feng, Zhifan Feng, Zhihao Deng, Zhijin Yu, Zhiyang Chen, Zhonghui Zheng, Zhuangzhuang Guo, Zhujun Zhang, Zhuo Sun, Zichang Liu, Zihan Lin, Zihao Huang, Zihe Zhu, Ziheng Zhao, Ziping Chen, Zixuan Zhu, Ziyang Xu, Ziyi Liang, Ziyuan Gao

References

AIME. Aime problems and solutions, 2025. URL https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions.

Philip Anastassiou, Jiawei Chen, Jitong Chen, Yuanzhe Chen, Zhuo Chen, Ziyi Chen, Jian Cong, Lelai Deng, Chuang Ding, Lu Gao, et al. Seed-tts: A family of high-quality versatile speech generation models. *arXiv preprint arXiv:2406.02430*, 2024.

Anthropic. Claude Opus 4.5, 2025. URL <https://www-cdn.anthropic.com/bf10f64990cfda0ba858290be7b8cc6317685f47.pdf>.

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.

Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, et al. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*, 2025.

Yejin Bang, Ziwei Ji, Alan Schelten, Anthony Hartshorn, Tara Fowler, Cheng Zhang, Nicola Cancedda, and Pascale Fung. Hallulens: Llm hallucination benchmark. *arXiv preprint arXiv:2504.17550*, 2025.

Victor Barres, Honghua Dong, Soham Ray, Xujie Si, and Karthik Narasimhan. τ^2 -Bench: Evaluating Conversational Agents in a Dual-Control Environment. *arXiv preprint arXiv:2506.07982*, 2025.

Hui Bu, Jiayu Du, Xingyu Na, Bengu Wu, and Hao Zheng. Aishell-1: An open-source mandarin speech corpus and a speech recognition baseline. In *2017 20th conference of the oriental chapter of the international coordinating committee on speech databases and speech I/O systems and assessment (O-COCOSDA)*, pp. 1–5. IEEE, 2017.

Ruisi Cai, Saurav Muralidharan, Greg Heinrich, Hongxu Yin, Zhangyang Wang, Jan Kautz, and Pavlo Molchanov. Flextron: Many-in-one flexible large language model. *arXiv preprint arXiv:2406.10260*, 2024.

Chen Chen, Xinlong Hao, Weiwen Liu, Xu Huang, Xingshan Zeng, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Yuefeng Huang, et al. Acebench: Who wins the match point in tool usage? *arXiv preprint arXiv:2501.12851*, 2025.

Liang Chen, Sinan Tan, Zefan Cai, Weichu Xie, Haozhe Zhao, Yichi Zhang, Junyang Lin, Jinze Bai, Tianyu Liu, and Baobao Chang. A spark of vision-language intelligence: 2-dimensional autoregressive transformer for efficient finegrained image generation. *arXiv preprint arXiv:2410.01912*, 2024a.

Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language models? *Advances in Neural Information Processing Systems*, 37:27056–27087, 2024b.

Mark Chen. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.

Yilong Chen, Junyuan Shang, Zhengyu Zhang, Jiawei Sheng, Tingwen Liu, Shuohuan Wang, Yu Sun, Hua Wu, and Haifeng Wang. Mixture of hidden-dimensions transformer. *arXiv preprint arXiv:2412.05644*, 2024c.

Yilong Chen, Junyuan Shang, Zhenyu Zhang, Shiyao Cui, Tingwen Liu, Shuohuan Wang, Yu Sun, and Hua Wu. Lemon: Reviving stronger and smaller lms from larger lms with linear parameter fusion. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8005–8019, 2024d.

Yilong Chen, Linhao Zhang, Junyuan Shang, Zhenyu Zhang, Tingwen Liu, Shuohuan Wang, and Yu Sun. Dha: Learning decoupled-head attention from transformer checkpoints via adaptive heads fusion. *Advances in Neural Information Processing Systems*, 37:45879–45913, 2024e.

Yiming Chen, Xianghu Yue, Chen Zhang, Xiaoxue Gao, Robby T Tan, and Haizhou Li. Voicebench: Benchmarking llm-based voice assistants. *arXiv preprint arXiv:2410.17196*, 2024f.

Xianfu Cheng, Wei Zhang, Shiwei Zhang, Jian Yang, Xiangyuan Guan, Xianjie Wu, Xiang Li, Ge Zhang, Jiaheng Liu, Yuying Mai, et al. Simplevqa: Multimodal factuality evaluation for multimodal large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4637–4646, 2025.

Yao Cheng, Jianfeng Chen, Jie Chen, Li Chen, Liyu Chen, Wentao Chen, Zhengyu Chen, Shijie Geng, Aoyan Li, Bo Li, et al. Fullstack bench: Evaluating llms as full stack coders. *arXiv preprint arXiv:2412.00535*, 2024.

Alexis Conneau, Min Ma, Simran Khanuja, Yu Zhang, Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara Rivera, and Ankur Bapna. Fleurs: Few-shot learning evaluation of universal representations of speech. In *2022 IEEE Spoken Language Technology Workshop (SLT)*, pp. 798–805. IEEE, 2023.

Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen Fan, Huayu Chen, Weize Chen, et al. The entropy mechanism of reinforcement learning for reasoning language models. *arXiv preprint arXiv:2505.22617*, 2025.

Google DeepMind. Gemini 2.5, 2025a. URL <https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/>.

Google DeepMind. Gemini 3.0, 2025b. URL <https://blog.google/products-and-platforms/products/gemini/gemini-3/>.

Google DeepMind. Introducing Veo 3, 2025c. URL <https://deepmind.google/models/veo/>.

Kaustubh Deshpande, Ved Sirdeshmukh, Johannes Baptist Mols, Lifeng Jin, Ed-Yeremai Hernandez-Cardona, Dean Lee, Jeremy Kritz, Willow E Primack, Summer Yue, and Chen Xing. Multichallenge: A realistic multi-turn conversation evaluation benchmark challenging to frontier llms. In *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 18632–18702, 2025.

Fnu Devvrit, Sneha Kudugunta, Aditya Kusupati, Tim Dettmers, Kaifeng Chen, Inderjit S Dhillon, Yulia Tsvetkov, Hannaneh Hajishirzi, Sham M Kakade, Ali Farhadi, et al. Matformer: Nested transformer for elastic inference. In *Workshop on Advancing Neural Network Training: Computational Efficiency, Scalability, and Resource Optimization (WANT@ NeurIPS 2023)*, 2023.

Juechu Dong, Boyuan Feng, Driss Guessous, Yanbo Liang, and Horace He. Flex attention: A programming model for generating optimized attention kernels. *arXiv preprint arXiv:2412.05496*, 2024.

Jiayu Du, Xingyu Na, Xuechen Liu, and Hui Bu. Aishell-2: Transforming mandarin asr research into industrial scale. *arXiv preprint arXiv:1808.10583*, 2018.

Baidu ERNIE Team. Ernie 4.5 technical report, 2025. URL https://yiyuan.baidu.com/blog/publication/ERNIE_Technical_Report.pdf.

Chaoyou Fu, Yuhai Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 24108–24118, 2025.

Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not perceive. In *European Conference on Computer Vision*, pp. 148–166. Springer, 2024.

Aryo Pradipta Gema, Joshua Ong Jun Leang, Giwon Hong, Alessio Devoto, Alberto Carlo Maria Mancino, Rohit Saxena, Xuanli He, Yu Zhao, Xiaotang Du, Mohammad Reza Ghasemi Madani, et al. Are we done with mmlu? In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 5069–5096, 2025.

Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36:52132–52152, 2023.

Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière, David Lopez-Paz, and Gabriel Synnaeve. Better & faster large language models via multi-token prediction. *arXiv preprint arXiv:2404.19737*, 2024.

Yuan Gong, Jin Yu, and James Glass. Vocalsound: A dataset for improving human vocal sounds recognition. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 151–155. IEEE, 2022.

Alex Gu, Baptiste Rozière, Hugh Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida I Wang. Cruxeval: A benchmark for code reasoning, understanding and execution. *arXiv preprint arXiv:2401.03065*, 2024.

Naibin Gu, Zhenyu Zhang, Yuchen Feng, Yilong Chen, Peng Fu, Zheng Lin, Shuohuan Wang, Yu Sun, Hua Wu, Weiping Wang, et al. Elastic moe: Unlocking the inference-time scalability of mixture-of-experts. *arXiv preprint arXiv:2509.21892*, 2025.

Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Minillm: Knowledge distillation of large language models. *arXiv preprint arXiv:2306.08543*, 2023.

Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14375–14385, 2024.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

Jian Han, Jinlai Liu, Yi Jiang, Bin Yan, Yuqi Zhang, Zehuan Yuan, Bingyue Peng, and Xiaobing Liu. Infinity: Scaling bitwise autoregressive modeling for high-resolution image synthesis. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 15733–15744, 2025.

Yancheng He, Shilong Li, Jiaheng Liu, Yingshui Tan, Weixun Wang, Hui Huang, Xingyuan Bu, Hangyu Guo, Chengwei Hu, Boren Zheng, et al. Chinese simpleqa: A chinese factuality evaluation for large language models. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 19182–19208, 2025.

Yun He, Di Jin, Chaoqi Wang, Chloe Bi, Karishma Mandyam, Hejia Zhang, Chen Zhu, Ning Li, Tengyu Xu, Hongjiang Lv, et al. Multi-if: Benchmarking llms on multi-turn and multilingual instructions following. *arXiv preprint arXiv:2410.15553*, 2024.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.

HMMT Organization. Hmmt, 2025. URL <https://www.hmmt.org/>.

Kairui Hu, Penghao Wu, Fanyi Pu, Wang Xiao, Yuanhan Zhang, Xiang Yue, Bo Li, and Ziwei Liu. Videomm: Evaluating knowledge acquisition from multi-discipline professional videos. *arXiv preprint arXiv:2501.13826*, 2025.

Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, et al. Minicpm: Unveiling the potential of small language models with scalable training strategies. *arXiv preprint arXiv:2404.06395*, 2024.

Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat, Dehao Chen, Mia Chen, HyoukJoong Lee, Jiquan Ngiam, Quoc V Le, Yonghui Wu, et al. Gpipe: Efficient training of giant neural networks using pipeline parallelism. *Advances in neural information processing systems*, 32, 2019.

Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Yao Fu, et al. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. *Advances in Neural Information Processing Systems*, 36:62991–63010, 2023.

Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21807–21818, 2024.

Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.

Il-Young Jeong and Jeongsoo Park. Cochlsene: Acquisition of acoustic scene data using crowdsourcing. In *2022 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pp. 17–21. IEEE, 2022.

Longbin Ji, Xiaoxiong Liu, Junyuan Shang, Shuohuan Wang, Yu Sun, Hua Wu, and Haifeng Wang. Videoar: Autoregressive video generation via next-frame & scale prediction. *arXiv preprint arXiv:2601.05966*, 2026.

Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4401–4410, 2019.

Mehran Kazemi, Bahare Fatemi, Hritik Bansal, John Palowitch, Chrysovalantis Anastasiou, Sanket Vaibhav Mehta, Lalit K Jain, Virginia Aglietti, Disha Jindal, Yuanzhu Peter Chen, et al. Big-bench extra hard. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 26473–26501, 2025.

Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images. In *European conference on computer vision*, pp. 235–251. Springer, 2016.

Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuancode: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024.

Rithesh Kumar, Prem Seetharaman, Alejandro Luebs, Ishaan Kumar, and Kundan Kumar. High-fidelity audio compression with improved rvqgan. *Advances in Neural Information Processing Systems*, 36: 27980–27993, 2023.

Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional computation and automatic sharding. *arXiv preprint arXiv:2006.16668*, 2020.

Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. Cmmlu: Measuring massive multitask language understanding in chinese. In *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 11260–11285, 2024.

Bill Yuchen Lin, Ronan Le Bras, Kyle Richardson, Ashish Sabharwal, Radha Poovendran, Peter Clark, and Yejin Choi. Zebralogic: On the scaling limits of llms for logical reasoning. *arXiv preprint arXiv:2502.01100*, 2025.

Ling-Team, Anqi Shen, Baihui Li, Bin Hu, Bin Jing, Cai Chen, Chao Huang, Chao Zhang, Chaokun Yang, Cheng Lin, et al. Every step evolves: Scaling reinforcement learning for trillion-scale thinking model. *arXiv preprint arXiv:2510.18855*, 2025.

Samuel Lipping, Parthasarathy Sudarsanam, Konstantinos Drossos, and Tuomas Virtanen. Clotho-aqa: A crowdsourced dataset for audio question answering. In *2022 30th European Signal Processing Conference (EUSIPCO)*, pp. 1140–1144. IEEE, 2022.

Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.

Aixin Liu, Aoxue Mei, Bangcai Lin, Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao Wu, Bowei Zhang, Chaofan Lin, Chen Dong, et al. Deepseek-v3.2: Pushing the frontier of open large language models. *arXiv preprint arXiv:2512.02556*, 2025a.

Hao Liu, Matei Zaharia, and Pieter Abbeel. Ring attention with blockwise transformers for near-infinite context. *arXiv preprint arXiv:2310.01889*, 2023a.

Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin Dong, Yejin Choi, Jan Kautz, and Yi Dong. Prorl: Prolonged reinforcement learning expands reasoning boundaries in large language models. *arXiv preprint arXiv:2505.24864*, 2025b.

Yuliang Liu, Zhang Li, Hongliang Li, Wenwen Yu, Mingxin Huang, Dezhi Peng, Mingyu Liu, Mingrui Chen, Chunyuan Li, Lianwen Jin, et al. On the hidden mystery of ocr in large multimodal models. *arXiv preprint arXiv:2305.07895*, 2(5):6, 2023b.

Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating math reasoning in visual contexts with gpt-4v, bard, and other large multimodal models. *arXiv preprint arXiv:2310.02255*, 2023.

Chuofan Ma, Yi Jiang, Junfeng Wu, Jihan Yang, Xin Yu, Zehuan Yuan, Bingyue Peng, and Xiaojuan Qi. Unitok: A unified tokenizer for visual generation and understanding. *arXiv preprint arXiv:2502.20321*, 2025a.

Wenhan Ma, Hailin Zhang, Liang Zhao, Yifan Song, Yudong Wang, Zhifang Sui, and Fuli Luo. Stabilizing moe reinforcement learning by aligning training and inference routers. *arXiv preprint arXiv:2510.11370*, 2025b.

Yanjun Ma, Dianhai Yu, Tian Wu, and Haifeng Wang. Paddlepaddle: An open-source deep learning platform from industrial practice. *Frontiers of Data and Computing*, 1(1):105–115, 2019.

Zeyao Ma, Bohan Zhang, Jing Zhang, Jifan Yu, Xiaokang Zhang, Xiaohan Zhang, Sijia Luo, Xi Wang, and Jie Tang. Spreadsheetbench: Towards challenging real world spreadsheet manipulation. *Advances in Neural Information Processing Systems*, 37:94871–94908, 2024.

Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9802–9822, 2023.

Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the association for computational linguistics: ACL 2022*, pp. 2263–2279, 2022.

Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 2200–2209, 2021.

Xin Men, Mingyu Xu, Qingyu Zhang, Qianhao Yuan, Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng Chen. Shortgpt: Layers in large language models are more redundant than you expect. In *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 20192–20204, 2025.

Annamaria Mesaros, Toni Heittola, and Tuomas Virtanen. Tut database for acoustic scene classification and sound event detection. In *2016 24th European signal processing conference (EUSIPCO)*, pp. 1128–1132. IEEE, 2016.

OpenAI. Hello GPT-4o, 2024. URL <https://openai.com/index/hello-gpt-4o/>.

OpenAI. Introducing gpt-5, 2025. URL <https://openai.com/index/introducing-gpt-5/>.

Roni Paiss, Ariel Ephrat, Omer Tov, Shiran Zada, Inbar Mosseri, Michal Irani, and Tali Dekel. Teaching clip to count to ten. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3170–3180, 2023.

Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: an asr corpus based on public domain audio books. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pp. 5206–5210. IEEE, 2015.

Shishir G. Patil, Huanzhi Mao, Charlie Cheng-Jie Ji, Fanjia Yan, Vishnu Suresh, Ion Stoica, and Joseph E. Gonzalez. The berkeley function calling leaderboard (bfcl): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*, 2025.

Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, et al. Humanity’s last exam. *arXiv preprint arXiv:2501.14249*, 2025.

Ivan Prosvirkov, Dmitrii Emelianenko, and Elena Voita. Bpe-dropout: Simple and effective subword regularization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 1882–1892, 2020.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pp. 28492–28518. PMLR, 2023.

Pooyan Rahmazadehgervi, Logan Bolton, Mohammad Reza Taesiri, and Anh Totti Nguyen. Vision language models are blind: Failing to translate detailed visual features into words. *arXiv preprint arXiv:2407.06581*, 2024.

Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations toward training trillion parameter models. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pp. 1–16. IEEE, 2020.

David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.

Angelika Romanou, Negar Foroutan, Anna Sotnikova, Zeming Chen, Sree Harsha Nelaturu, Shivalika Singh, Rishabh Maheshwary, Micol Altomare, Mohamed A Haggag, Alfonso Amayuelas, et al. Include: Evaluating multilingual language understanding with regional knowledge. *arXiv preprint arXiv:2411.19799*, 2024.

Hassan Sajjad, Fahim Dalvi, Nadir Durrani, and Preslav Nakov. On the effect of dropping layers of pre-trained transformer models. *Computer Speech & Language*, 77:101429, 2023.

Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.

S Sakshi, Utkarsh Tyagi, Sonal Kumar, Ashish Seth, Ramaneswaran Selvakumar, Oriol Nieto, Ramani Duraiswami, Sreyan Ghosh, and Dinesh Manocha. Mmau: A massive multi-task audio understanding and reasoning benchmark. *arXiv preprint arXiv:2410.19168*, 2024.

Team Seedream, Yunpeng Chen, Yu Gao, Lixue Gong, Meng Guo, Qiushan Guo, Zhiyao Guo, Xiaoxia Hou, Weilin Huang, Yixuan Huang, et al. Seedream 4.0: Toward next-generation multimodal image generation. *arXiv preprint arXiv:2509.20427*, 2025.

Mohammad Shoeybi, Mostafa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*, 2019.

Yueqi Song, Tianyue Ou, Yibo Kong, Zecheng Li, Graham Neubig, and Xiang Yue. Visualpuzzles: Decoupling multimodal reasoning evaluation from domain knowledge. *arXiv preprint arXiv:2504.10342*, 2025.

Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, et al. Challenging big-bench tasks and whether chain-of-thought can solve them. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 13003–13051, 2023.

Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025.

Guoxia Wang, Jinle Zeng, Xiyuan Xiao, Siming Wu, Jiabin Yang, Lujing Zheng, Zeyu Chen, Jiang Bian, Dianhai Yu, and Haifeng Wang. Flashmask: Efficient and rich mask extension of flashattention. *arXiv preprint arXiv:2410.01359*, 2024a.

Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie Zhan, and Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset. *Advances in Neural Information Processing Systems*, 37:95095–95169, 2024b.

Lean Wang, Huazuo Gao, Chenggang Zhao, Xu Sun, and Damai Dai. Auxiliary-loss-free load balancing strategy for mixture-of-experts. *arXiv preprint arXiv:2408.15664*, 2024c.

Peihao Wang, Rameswar Panda, Lucas Torroba Hennigen, Philip Greengard, Leonid Karlinsky, Rogerio Feris, David Daniel Cox, Zhangyang Wang, and Yoon Kim. Learning to grow pretrained models for efficient transformer training. *arXiv preprint arXiv:2303.00980*, 2023.

Shenzi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen, Jianxin Yang, Zhenru Zhang, et al. Beyond the 80/20 rule: High-entropy minority tokens drive effective reinforcement learning for llm reasoning. *arXiv preprint arXiv:2506.01939*, 2025.

Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyuan Jiang, et al. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *Advances in Neural Information Processing Systems*, 37:95266–95290, 2024d.

Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, et al. Charxiv: Charting gaps in realistic chart understanding in multimodal llms. *Advances in Neural Information Processing Systems*, 37:113569–113697, 2024e.

Jason Wei, Nguyen Karina, Hyung Won Chung, Yunxin Joy Jiao, Spencer Papay, Amelia Glaese, John Schulman, and William Fedus. Measuring short-form factuality in large language models. *arXiv preprint arXiv:2411.04368*, 2024.

Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won Chung, Alex Tachard Passos, William Fedus, and Amelia Glaese. Browsecomp: A simple yet challenging benchmark for browsing agents. *arXiv preprint arXiv:2504.12516*, 2025.

Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng-ming Yin, Shuai Bai, Xiao Xu, Yilei Chen, et al. Qwen-image technical report. *arXiv preprint arXiv:2508.02324*, 2025.

Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. Sheared llama: Accelerating language model pre-training via structured pruning. *arXiv preprint arXiv:2310.06694*, 2023.

Jin Xu, Zhifang Guo, Hangrui Hu, Yunfei Chu, Xiong Wang, Jinzheng He, Yuxuan Wang, Xian Shi, Ting He, Xinfu Zhu, et al. Qwen3-omni technical report. *arXiv preprint arXiv:2509.17765*, 2025a.

Weiye Xu, Jiahao Wang, Weiyun Wang, Zhe Chen, Wengang Zhou, Aijun Yang, Lewei Lu, Houqiang Li, Xiaohua Wang, Xizhou Zhu, et al. Visulogic: A benchmark for evaluating visual reasoning in multi-modal large language models. *arXiv preprint arXiv:2504.15279*, 2025b.

Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen, Reynold Cheng, Jinyang Li, Can Xu, Dacheng Tao, and Tianyi Zhou. A survey on knowledge distillation of large language models. *arXiv preprint arXiv:2402.13116*, 2024.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, pp. 2369–2380, 2018.

Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.

Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9556–9567, 2024.

Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, et al. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15134–15186, 2025a.

Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Shiji Song, and Gao Huang. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model? *arXiv preprint arXiv:2504.13837*, 2025b.

Hui Zeng. Measuring massive multitask chinese understanding. *arXiv preprint arXiv:2304.12986*, 2023.

Binbin Zhang, Hang Lv, Pengcheng Guo, Qijie Shao, Chao Yang, Lei Xie, Xin Xu, Hui Bu, Xiaoyu Chen, Chenchen Zeng, et al. Wenetspeech: A 10000+ hours multi-domain mandarin corpus for speech recognition. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6182–6186. IEEE, 2022.

Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Yu Qiao, et al. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? In *European Conference on Computer Vision*, pp. 169–186. Springer, 2024a.

Xin Zhang, Dong Zhang, Shimin Li, Yaqian Zhou, and Xipeng Qiu. Speechtokenerizer: Unified speech tokenizer for speech language models. In *The Twelfth International Conference on Learning Representations*, 2024b.

Chenggang Zhao, Shangyan Zhou, Liyue Zhang, Chengqi Deng, Zhean Xu, Yuxuan Liu, Kuai Yu, Jiashi Li, and Liang Zhao. Deepep: an efficient expert-parallel communication library. <https://github.com/deepseek-ai/DeepEP>, 2025a.

Rosie Zhao, Alexandru Meterez, Sham Kakade, Cengiz Pehlevan, Samy Jelassi, and Eran Malach. Echo chamber: RL post-training amplifies behaviors learned in pretraining. *arXiv preprint arXiv:2504.07912*, 2025b.

Yilun Zhao, Haowei Zhang, Lujing Xie, Tongyan Hu, Guo Gan, Yitao Long, Zhiyuan Hu, Weiyuan Chen, Chuhan Li, Zhijian Xu, et al. Mmvu: Measuring expert-level multi-discipline video understanding. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 8475–8489, 2025c.

Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, et al. Group sequence policy optimization. *arXiv preprint arXiv:2507.18071*, 2025.

Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. Agieval: A human-centric benchmark for evaluating foundation models. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pp. 2299–2314, 2024.

Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*, 2023.

Yuzhen Zhou, Jiajun Li, Yusheng Su, Gowtham Ramesh, Zilin Zhu, Xiang Long, Chenyang Zhao, Jin Pan, Xiaodong Yu, Ze Wang, et al. April: Active partial rollouts in reinforcement learning to tame long-tail generation. *arXiv preprint arXiv:2509.18521*, 2025.

Nannan Zhu, Yonghao Dong, Teng Wang, Xueqian Li, Shengjun Deng, Yijia Wang, Zheng Hong, Tiantian Geng, Guo Niu, Hanyan Huang, et al. Cvbench: Benchmarking cross-video synergies for complex multimodal reasoning. *arXiv preprint arXiv:2508.19542*, 2025.