

Mixture of Experts in Large Language Models

Danyang Zhang

*Department of Research
ByteDance Inc*

San Jose, California, United States
joseph.zhang@bytedance.com

Junhao Song[†]

*Department of Computing
Imperial College London
London, United Kingdom*
junhao.song23@imperial.ac.uk

Ziqian Bi

*Department of Computer Science
Purdue University
West Lafayette, Indiana, United States*
bi32@purdue.edu

Yingfang Yuan

*Department of Computer Science
Heriot-Watt University
Edinburgh, United Kingdom*
y.yuan@hw.ac.uk

Tianyang Wang

*AI Agent Lab
Vokram Group
London, United Kingdom*
tianyang.wg35@gmail.com

Joe Yeong

*Department of Anatomical Pathology
Singapore General Hospital
Singapore*
yeongps@imcb.a-star.edu.sg

Junfeng Hao

*AI Agent Lab
Vokram Group
London, United Kingdom*
ygzjhj85@gmail.com

Abstract—This paper presents a comprehensive review of the Mixture-of-Experts (MoE) architecture in large language models, highlighting its ability to significantly enhance model performance while maintaining minimal computational overhead. Through a systematic analysis spanning theoretical foundations, core architectural designs, and large language model (LLM) applications, we examine expert gating and routing mechanisms, hierarchical and sparse MoE configurations, meta-learning approaches, multimodal and multitask learning scenarios, real-world deployment cases, and recent advances and challenges in deep learning. Our analysis identifies key advantages of MoE, including superior model capacity compared to equivalent Bayesian approaches, improved task-specific performance, and the ability to scale model capacity efficiently. We also underscore the importance of ensuring expert diversity, accurate calibration, and reliable inference aggregation, as these are essential for maximizing the effectiveness of MoE architectures. Finally, this review outlines current research limitations, open challenges, and promising future directions, providing a foundation for continued innovation in MoE architecture and its applications.

Index Terms—Large language models, mixture of experts, expert routing, meta learning, knowledge transfer, sparse activation, large language models architecture, natural language processing

I. INTRODUCTION AND FUNDAMENTALS

Over the past decade, deep learning has witnessed an explosive increase in model capacity, particularly with the emergence of large-scale transformer architectures. These models, while powerful, are often computationally intensive and memory hungry. Scaling such systems linearly in parameter count leads to exponential increases in FLOPs and energy consumption. This trend is rapidly becoming unsustainable for deployment in real-world scenarios. In response, researchers have begun exploring sparse and modular computation as viable alternatives to dense, monolithic models. Among these, Mixture-of-Experts (MoE) architectures stand out as a promising solution. By conditionally activating only a subset of a model’s parameters, typically a small number of specialized

expert modules, MoEs decouple inference cost from total model size, enabling efficient scaling of model capacity [61].

The conceptual foundations of MoE trace back to early work in adaptive learning systems, where models were organized as ensembles of specialized experts, each handling a specific sub-region of the input space [26, 77, 69]. These early frameworks emphasized modularity and competitive learning, often using a gating function to route each input to the most relevant expert. However, computational constraints and the lack of scalable training mechanisms limited their practical impact. MoEs did not become viable for high-performance applications until sparse routing was integrated with modern deep networks, made possible by large-scale distributed computation. The breakthrough came with the development of sparsely gated networks [61], which showed that it is possible to maintain model accuracy while activating only a small fraction of the parameters during both training and inference.

More than just a tool for parameter efficiency, MoE represents a broader architectural shift toward modular design in deep learning. Traditional networks treat all parts of the model as uniformly important across all tasks and inputs. In contrast, MoEs introduce conditional computation: each expert module can specialize in a certain domain, linguistic pattern, or modality. During training, a learnable gating function, often implemented as a lightweight neural network, determines which experts are best suited to handle each input [14, 4]. This routing process encourages functional diversity among the experts, allowing them to focus on distinct aspects of the input distribution, which improves generalization and robustness across tasks.

The renewed interest in MoEs aligns with a growing recognition that monolithic deep learning models, while powerful, are not necessarily optimal from a computational or representational standpoint. MoE architectures offer a more flexible and scalable alternative by embracing heterogeneity in both structure and computation. Their success across diverse application domains, including language modeling, machine translation, and vision-language reasoning, demonstrates the

[†]: Corresponding author: Junhao Song (junhao.song23@imperial.ac.uk)

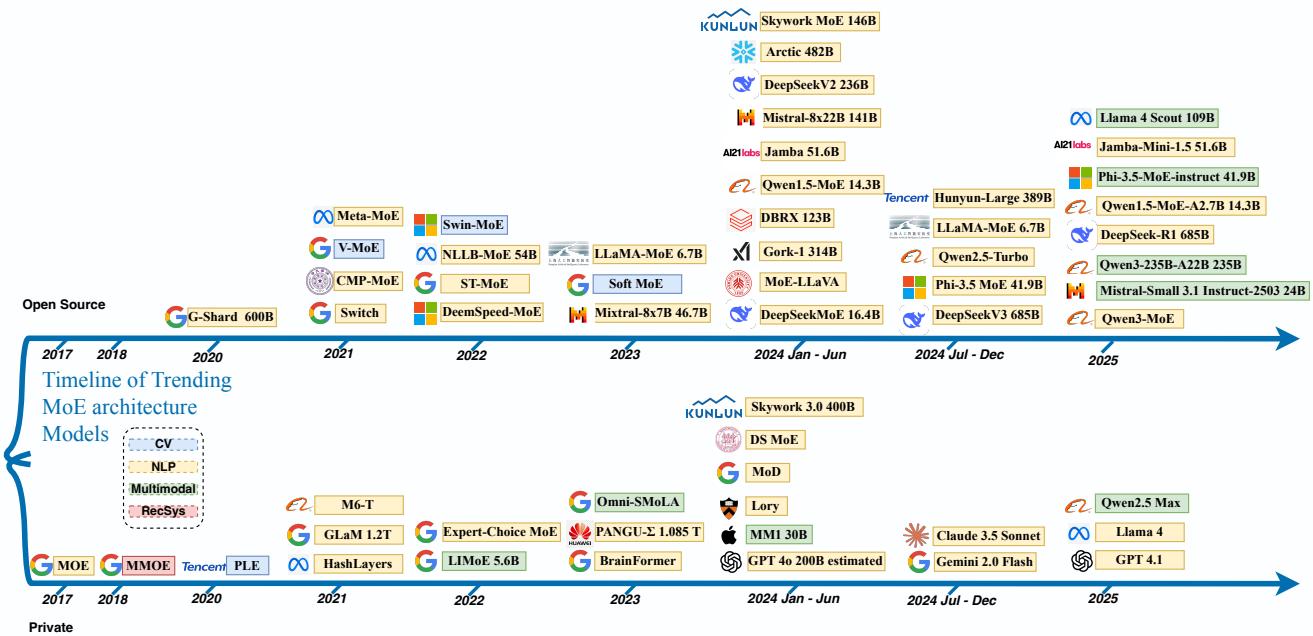


Fig. 1. **Timeline of mixture of experts (MoE) models development.** The timeline shows key milestones in MoE architecture evolution from foundational concepts to modern large-scale implementations.

generalizability of the approach. Crucially, MoEs challenge the prevailing assumption that scale must come at the cost of efficiency, providing a viable path toward building more intelligent, adaptable, and resource-conscious AI systems.

From sparse gating to billion-scale deployment. The development of MoE architectures underwent a pivotal transformation starting in 2020, when research shifted from conceptual formulations to practical, high-scale implementations. As shown in **Figure 1**, this turning point was marked by the release of GShard [33], a large-scale multilingual model with 600 billion parameters. GShard pioneered auto-sharding and token-level expert routing, making it one of the first successful demonstrations of sparse MoE architectures at the trillion-parameter scale. Shortly after, Switch Transformer [16] and GLaM [13] extended this paradigm to language modeling. These models employed token-choice gating with only 1–2 experts activated per input, significantly reducing compute overhead and establishing MoE as a credible scaling alternative to dense transformers.

By 2021–2022, MoE architectures had matured from proof-of-concept models into a broadly adopted computational framework. Open-source releases such as Meta-MoE, CMP-MoE, and V-MoE demonstrated the community-wide interest in modular computation. At the same time, commercial labs began investing in domain-specific variants: NLLB-MoE for multilingual translation, Swin-MoE for vision tasks, and LIMoE [47] for multimodal learning. This diversification marked the beginning of MoE’s adaptation beyond NLP, driven by the recognition that sparse expert activation could generalize to vision, audio, and cross-modal scenarios.

The post-2023 landscape, as shown in the upper half of **Figure 1**, marks a phase of rapid industrial scaling and architectural diversification. Models such as DeepSeekV3 (685B), Skywork 3.0 (400B), and Arctic (482B) demonstrate that MoE has become a core component of modern foundation models. These systems increasingly integrate expert routing with paradigms like retrieval, instruction tuning, and agent-based control, indicating that MoE now plays a central role in large-scale AI development. The emergence of models such as MoE-LLaVA, MM1 [43], and Omni-SMoLA [70] further reflects a shift toward multimodal and grounded reasoning architectures.

Open-source initiatives have accelerated MoE adoption by releasing advanced models such as Jamba, Qwen1.5-MoE, and Mistral-8x22B. Ongoing improvements to routing and load balancing in Qwen3-MoE, Claude 3.5 Sonnet, and Llama 4 show that innovation now centers on efficiency, controllability, and multi-task generalization rather than scale alone.

Diversity of designs across modalities and tasks. The MoE paradigm has evolved into a diverse architectural space, as shown in Table I. Its variants differ not only in scale and number of experts but also in routing strategies and domain-specific design choices. Early implementations such as Switch Transformer and GLaM [16, 13] focused on language modeling, where token-level gating activates only a subset of experts per token. This form of sparse activation enabled scaling to hundreds of billions of parameters with minimal compute overhead, while preserving sequence-level modeling fidelity.

In contrast, translation-oriented models like GShard and

TABLE I

TAXONOMY OF REPRESENTATIVE MoE ARCHITECTURES ACROSS SIX APPLICATION DOMAINS, HIGHLIGHTING EXPERT COUNTS, ROUTING STRATEGIES, DEVELOPMENT PERIODS (2017–2024), AND KEY INNOVATIONS OR USE CASES.

Category	Model (#Experts)	Routing	Year	Key Innovation / Use Case
Language LLM	Switch Transformer (64) [16] / GLaM (64) [13]	Token-choice	2021–22	Trillion-parameter LLMs with only 1/64 active parameters per token.
Translation	GShard MoE (128) [33] / DeepSpeed-MoE (256) [57]	Token-choice	2020–21	Auto sharding + pipeline parallelism for cross-lingual MT at TB scale.
Multimodal	Omni-SMoLA (16) [71] / T-REX2 (32) [28]	Cross-modal gate	2023–24	Low-rank experts plus dual vision-text prompts for open-set detection.
Computer Vision	MoCaE-DET (8) [52] / Deep-MoE (32) [14]	Attention gate	2017–23	Calibration-aware fusion; boosts COCO AP by 2.5 vs. single detector.
Param-Efficient	LoRA-MoE (4) [80] / Nexus (8) [20]	Frozen router	2023–24	<1 % parameter update (PEFT); “upcycles” dense checkpoints to adaptive MoE.
Hierarchical	H-MoE (32) [48] / MixER (10) [51]	2-level / Top-1	2024	Coarse-to-fine quadratic gating; K-means routing for dynamical systems.

DeepSpeed-MoE [33, 57] emphasize system-level throughput. These models combine MoE with pipeline parallelism and expert sharding to process multilingual data at terabyte scale. Multimodal variants including Omni-SMoLA [71] and T-REX2 [28] route using both visual and textual cues, enabling open-set recognition and grounded captioning where context spans modalities.

Other architectures diverge along structural dimensions. Hierarchical models such as H-MoE [48] and MixER [51] implement multi-stage routing, often using clustering or coarse-to-fine attention to improve interpretability and modularity. Parameter-efficient frameworks like LoRA-MoE [80] and Nexus [20] focus on minimizing update costs during finetuning, by freezing the router and only adapting a few expert weights. These designs are especially useful in settings where retraining large models is computationally prohibitive.

Routing complexity and the stability efficiency tradeoff.

The routing algorithm lies at the heart of every MoE system. Early designs employed simple top- k selection using softmax over expert logits [61], but recent variants increasingly explore entropy-based [85], load-balanced [16], or differentiable attention-based routing mechanisms [34, 85, 59]. The goal is twofold: ensure expert diversity while minimizing redundant computation.

Yet, these methods often suffer from instability or underutilization. For example, overly confident gates can collapse to a few dominant experts, whereas uniform allocation reduces specialization. Techniques like auxiliary load balancing losses or stochastic gating [16] have been proposed, but they introduce additional hyperparameters and training complexity. Designing robust, general-purpose routing remains one of the most technically challenging aspects of MoE development.

Deployment constraints. Despite their computational advantages, MoE models face non-trivial barriers to deployment. Sparse expert activation introduces irregular memory access patterns and frequent cross-device communication, resulting in elevated inference latency and hardware underutilization. Moreover, the stochastic nature of routing leads to unstable batching, fragmented workloads, and poor reproducibility. These issues are particularly problematic in low-latency or memory-constrained environments.

To mitigate these issues, newer designs adopt more deployment-aware strategies. LoRA-MoE [80] and Nexus [20], for example, freeze routing weights or apply low-rank adapters to stabilize expert usage during inference. These methods

reduce routing variance, simplify caching, and enable efficient fine-tuning. This marks a shift from MoE systems optimized purely for scale toward architectures designed with practical usability and system compatibility in mind.

Scope of this work. This paper presents a structured synthesis of recent advances in MoE architectures. It begins by analyzing core components such as expert modules, gating mechanisms, and load balancing strategies. The discussion then extends to domain-specific adaptations in natural language processing, computer vision, and multimodal learning. Key technical challenges are examined, including routing instability, expert underutilization, and scalability limits. Finally, the survey highlights emerging directions such as sparse mixture fusion, expert replay, and lifelong modularity. By consolidating architectural innovations and comparing their design principles, this work clarifies the current landscape and outlines future paths for developing scalable and efficient MoE-based AI systems.

II. CORE ARCHITECTURES AND ROUTING MECHANISMS

This section examines the foundational and advanced architectural components of MoE models, with a particular focus on gating mechanisms, expert routing strategies, and architectural variations that promote specialization and efficiency. To provide a structured overview of these developments, **Figure 2** presents a comprehensive taxonomy of MoE systems. The taxonomy covers seven dimensions, mapping techniques, implementations, and representative systems across language, multimodal, routing, and deployment.

Core concepts and mathematical principles. A MoE layer routes each input x to a sparse subset of k expert networks among a larger pool of N experts. Instead of processing the input through all experts, the model computes the output as a weighted combination of only the selected ones:

$$y = \sum_{i=1}^N g_i(x) E_i(x), \quad (1)$$

where $E_i(x)$ denotes the output of expert i , and $g_i(x)$ is the gating function that assigns nonzero weights to at most $k \ll N$ experts. The gating function is typically implemented via Noisy Top- k routing [61], which introduces Gaussian noise into the expert scores before selecting the top k :

$$H(x)_i = (x \cdot W_g)_i + \mathcal{N}(0, \sigma^2) \quad (2)$$

This stochasticity prevents early expert collapse and encourages exploration, while top- k sparsity ensures linear compute

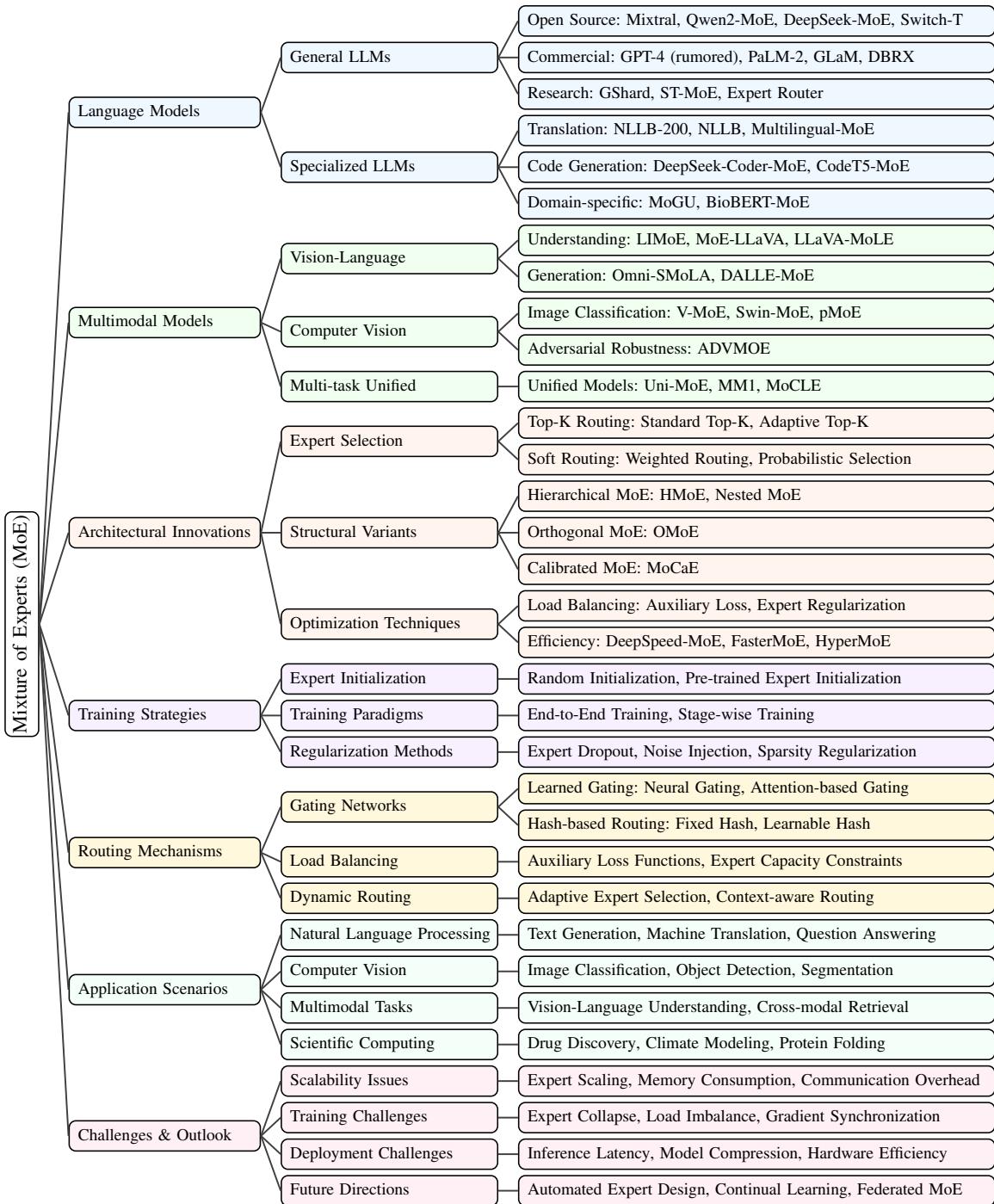


Fig. 2. A comprehensive taxonomy of Mixture of Experts (MoE) models, organizing methodologies into seven key categories: language models, multimodal models, architectural innovations, training strategies, routing mechanisms, application scenarios, and challenges. Each category encompasses specific techniques, implementations, and representative models.

cost with respect to the number of selected experts, not total expert pool size. The resulting sparse activations can be efficiently parallelized across devices, enabling scalable training and inference.

A. Foundational MoE Architectures

Sparse activation via gating networks. Each MoE layer comprises N experts $\{E_1, \dots, E_N\}$ and a gating function $G(x)$ that selects a sparse subset of $k \ll N$ experts for each input. As illustrated in **Figure 3**, these experts are typically independent feedforward modules. The gating function computes a relevance score $H(x)_i$ for each expert, perturbed by Gaussian noise scaled by a soft activation term:

$$H(x)_i = (x \cdot W_g)_i + \mathcal{N}(0, 1) \cdot \text{Softplus}((x \cdot W_n)_i), \quad (3)$$

$$G(x) = \text{Softmax}(\text{TopK}(H(x), k)), \quad (4)$$

where TopK masks all but the k largest scores with $-\infty$ before applying softmax. This Noisy Top- k routing mechanism [61] promotes exploration, mitigates early expert collapse, and ensures sparse yet balanced expert activation.

By evaluating only a small subset of experts per token, MoE models achieve substantial computational savings and parallel scalability. Large-scale implementations such as Switch Transformer [16], GShard [33], and DeepSpeed-MoE [57] demonstrate that this design enables efficient training of models with hundreds of billions of parameters.

Load balancing objectives. A key challenge in MoE training is expert collapse—where only a small subset of experts receive the majority of inputs, leaving others underutilized. To mitigate this, a load-balancing objective is added to the training loss:

$$L_{\text{balance}} = \alpha \sum_{i=1}^N f_i \cdot P_i, \quad (5)$$

where f_i denotes the fraction of tokens assigned to expert i , and P_i is its average gate probability [16]. The product $f_i \cdot P_i$ penalizes deviations between expected and actual expert usage, encouraging uniform load distribution.

This auxiliary loss promotes better resource utilization and prevents training instability due to expert starvation. However, it introduces a tension with the primary routing objective: enforcing balance may degrade expert specialization or routing accuracy. As a result, selecting the appropriate α and balancing regularization strength remains an open trade-off in large-scale MoE design.

From prototypes to production. Transitioning MoE models from conceptual sparsity mechanisms to production-grade systems required addressing significant engineering constraints—particularly in communication overhead, memory fragmentation, and inference latency. Early models like GShard [33] demonstrated the feasibility of expert parallelism via auto-sharded tensor computation, but required substantial infrastructure support.

Recent advances focus on improving throughput and stability during both training and deployment. For instance, Mixtral [27] employs static top-2 routing and fused attention layers

to minimize communication overhead. DBRX [21] introduces fused MoE kernels and low-overhead memory prefetching, while Qwen2 [63] and DeepSeek-v3 [36] integrate quantized MoE layers and expert dropout to reduce inference cost without degrading accuracy.

These systems increasingly adopt fixed-capacity experts and static load balancing schemes, making them compatible with modern serving frameworks. This convergence of sparsity, hardware alignment, and inference regularization marks MoE’s shift from research prototype to scalable, production-ready backbone.

Theoretical capacity and scaling behavior. MoE models benefit from a composite hypothesis space defined by a mixture of subnetworks. Let \mathcal{H}_i denote the hypothesis class of expert E_i , and \mathcal{G} the gating space. Then the overall capacity satisfies:

$$\mathcal{H}_{\text{MoE}} = \bigcup_{g \in \mathcal{G}} \left\{ \sum_{i=1}^N g_i(x) E_i(x) \mid E_i \in \mathcal{H}_i \right\} \quad (6)$$

where \mathcal{H}_i is the hypothesis class of expert E_i , and \mathcal{G} denotes the space of all gating functions.

While MoEs offer high capacity through modular specialization, their lack of shared inductive bias can hinder generalization when routing is unstable. Training dynamics are sensitive to expert overlap and gate smoothness, and unlike Bayesian ensembles, MoEs delegate via hard routing rather than aggregate over uncertainty—potentially increasing variance. Nonetheless, large-scale evaluations show that MoEs achieve comparable performance to dense models on language tasks while activating $\sim 10\times$ fewer parameters per token [19], demonstrating an effective trade-off between compute efficiency and representational specialization.

B. Advanced Architectural Variants

Orthogonal training. To reduce redundancy and encourage specialization among experts, Orthogonal MoE (OMoE) introduces a regularization term that enforces pairwise weight orthogonality:

$$L_{\text{orth}} = \sum_{i \neq j} \langle W_i, W_j \rangle^2, \quad (7)$$

penalizing alignment between expert parameters and promoting functional diversity [17].

Alternatively, mutual distillation encourages knowledge sharing by minimizing pairwise KL divergence across expert outputs:

$$L_{\text{distill}} = \sum_{i=1}^N \sum_{j \neq i} \text{KL}(E_i(x) \| E_j(x)), \quad (8)$$

which improves robustness but risks homogenization if not properly balanced [75].

Parameter efficient tuning. MoE architectures can be tuned with high efficiency by restricting updates to lightweight expert-specific components. This allows models to maintain the benefits of sparse computation while drastically reducing

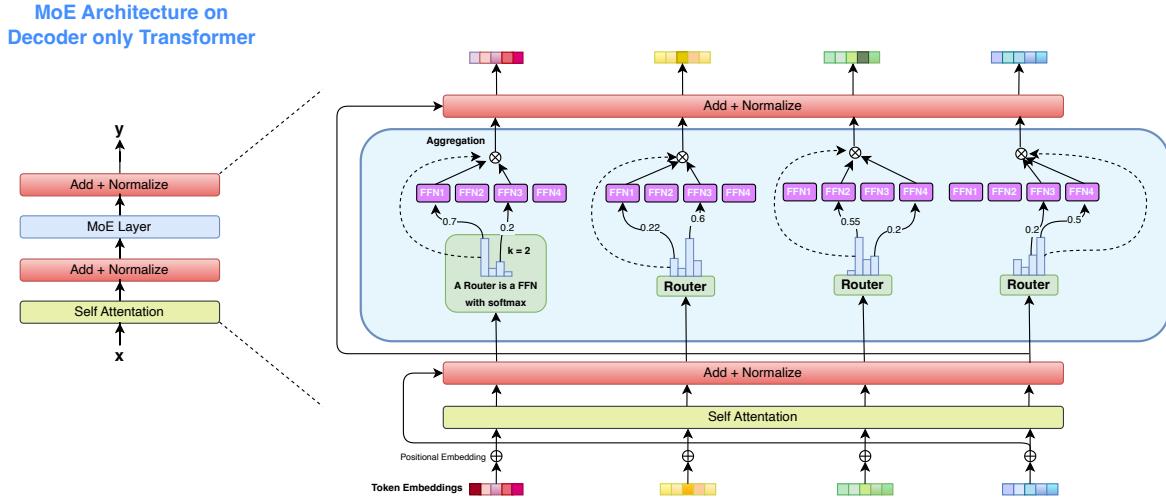


Fig. 3. A brief illustration of sparsely gated Mixture of Experts (MoE) architecture on decoder only transformer. In this figure, the top- k routing mechanism is configured with $k=2$, meaning the gating function selects the two highest-scoring FFN experts for each token based on the router’s softmax probabilities. The selected experts are evaluated in parallel, with their outputs aggregated using weighted summation.

the number of trainable parameters. Formally, the parameter update ratio ρ is minimized:

$$\rho = \frac{\|\Delta\theta_{\text{train}}\|_0}{\|\theta_{\text{full}}\|_0} \ll 1, \quad (9)$$

where $\Delta\theta_{\text{train}}$ denotes updated parameters and θ_{full} is the total parameter set. Methods such as MoCE-IR [80] and Adamix [68] implement this by freezing shared layers and updating only expert heads or low-rank adapters, often achieving $< 1\%$ parameter updates with minimal performance drop.

These strategies align with the goals of Low-Rank Adaptation (LoRA) and adapter-based fine-tuning, treating MoE as a modular scaffold that enables expressive adaptation under strict compute budgets.

Hierarchical and multi-head extensions. To further enhance specialization and routing flexibility, hierarchical MoEs (HMoEs) introduce a two-stage gating process [48]. A coarse gate $G^{(1)}(x)$ selects a super-expert group, while a secondary gate $G^{(2)}(x)$ routes the input within that group:

$$y = \sum_{i \in \mathcal{G}^{(1)}(x)} G_i^{(2)}(x) E_i(x), \quad (10)$$

where $\mathcal{G}^{(1)}(x)$ denotes the selected coarse cluster. This structure enables specialization at multiple abstraction levels and supports large expert pools without routing overhead explosion.

Multi-head MoEs [72], in contrast, assign separate expert subsets to different input dimensions or tasks, allowing input features to be processed in parallel across heads. This setup has shown benefits in vision and speech domains where spatial or temporal decompositions are naturally aligned with the modular structure of MoE.

Heterogeneous and adaptive experts. Traditional MoEs assume homogeneous experts with uniform capacity and archi-

tecture, but recent advances explore heterogeneity to improve routing flexibility and efficiency [65]. In this setup, each expert E_i is characterized by its own computational profile ϕ_i (e.g., depth, width, modality), allowing the gating function to assign complex inputs to high-capacity experts and simpler ones to lightweight modules:

$$g(x) = \arg \max_i [S_i(x) - \lambda \cdot \text{Cost}(\phi_i)], \quad (11)$$

where $S_i(x)$ is the score assigned to expert i and $\text{Cost}(\phi_i)$ reflects its computational expense. This cost-aware routing improves both specialization and hardware efficiency.

Moreover, knowledge integration from unselected experts has emerged as a mechanism for expanding capacity without increasing runtime cost. HyperMoE [83], for instance, aggregates intermediate signals from unselected experts to refine final predictions:

$$y = \sum_{i \in \mathcal{A}(x)} g_i(x) E_i(x) + \gamma \sum_{j \notin \mathcal{A}(x)} h_j(x), \quad (12)$$

where $\mathcal{A}(x)$ is the active expert set, and $h_j(x)$ encodes side information from inactive experts. This promotes cross-task generalization and has shown particular promise in multitask or low-resource regimes.

C. Routing Strategies and Specialization Patterns

Token choice vs. expert choice. Sparse routing in MoE layers can follow either *Token Choice* or *Expert Choice*, as shown in **Figure 4**. In Token Choice, each token x_t is routed to its top- k experts:

$$y_t = \sum_{i \in \mathcal{A}_t} g_i(x_t) E_i(x_t), \quad \mathcal{A}_t = \text{TopK}_{i=1}^N [g_i(x_t)]. \quad (13)$$

For example, in the left panel, token "\We" selects experts 1 and 3, while "\love" selects experts 2 and 4, each independently routed based on its own gate scores.

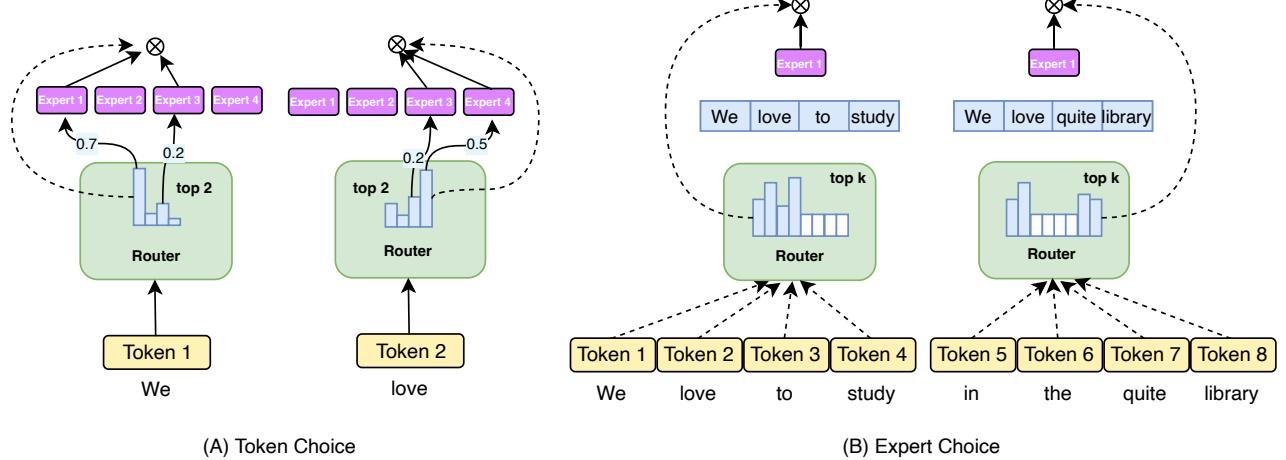


Fig. 4. **Comparison of routing strategies in Token Choice and Expert Choice Mixture-of-Experts architectures.** (A) Token Choice routing: Each token is processed by selecting the most suitable experts based on computed affinity scores, with tokens "We" being routed to Expert 1 and Expert 3, while "Like" being routed to and Expert 3 and Expert4, respectively with their corresponding probability weights. (B) Expert Choice routing: Experts maintain fixed computational budgets and select their preferred tokens from the input sequence, where Expert 1 processes tokens ["We", "Love", "To", "Study"] and Expert 2 handles ["We", "Love", "Quite", "Library"], enabling balanced workload distribution across experts while allowing tokens to be processed by multiple experts when beneficial.

In contrast, Expert Choice lets each expert E_i selects a subset of tokens \mathcal{T}_i to process under a fixed budget:

$$y_t = \sum_{i:x_t \in \mathcal{T}_i} \tilde{g}_i(x_t) E_i(x_t), \quad \mathcal{T}_i = \text{TopB}_{t=1}^T [s_i(x_t)]. \quad (14)$$

In the right panel, Expert 1 selects tokens "\We", "\love", "\to", and "\study", while Expert 2 handles "\quiet" and "\library", enabling token grouping and controlled load balancing.

These paradigms differ fundamentally in control flow: Token Choice assigns experts per-token, while Expert Choice assigns tokens per-expert. The latter improves expert utilization and coherence, particularly in vision or structured input tasks [38].

Learned routers vs. Fixed routers. Although trainable gating networks are widely used, recent work [15] shows that randomly initialized fixed routers can yield comparable or even superior performance. This challenges the assumption that learned gating $g_\theta(x)$ necessarily improves specialization. In fixed routing, the gate $g(x)$ becomes a deterministic mapping or a sampled sparse mask:

$$g_i(x) = \begin{cases} 1, & i \in \mathcal{A}_{\text{fixed}}(x) \\ 0, & \text{otherwise} \end{cases} \quad (15)$$

where $\mathcal{A}_{\text{fixed}}(x)$ is a static assignment determined at initialization. This approach avoids routing instability and eliminates variance from gradient updates, especially in early training.

Empirical studies further show that routing granularity influence the emergence of expert functions. Sequence-level gating (i.e., one routing decision per input) typically groups experts by topic or discourse structure, whereas token-level gating yields more fine-grained specialization that frequently

aligns with syntactic categories such as nouns or verbs [15].

Emergent linguistic structure in expert assignments. Probing experiments [2] reveal that MoE layers implicitly cluster inputs by part-of-speech and morphological role, even without explicit supervision. Let $S_i = \{x_t \mid g_i(x_t) > 0\}$ be the token set for expert i ; statistical analysis shows that S_i has high mutual information with specific POS categories:

$$I(\text{POS}; \text{Expert}_i) = H(\text{POS}) - H(\text{POS} \mid S_i). \quad (16)$$

This implies that specialization arises from the interplay between architecture and training dynamics, not solely from design choices. In practice, this interpretability enhances modular debugging, domain adaptation, and controllable generation.

Adaptive expert selection. In adaptive MoE architectures, the gating function $g(x)$ is augmented with input-dependent capacity control, enabling dynamic scaling of expert participation. For input x , the number of activated experts $k(x)$ is computed as:

$$k(x) = \min(K_{\max}, \lfloor \tau \cdot \|x\| \rfloor), \quad (17)$$

where τ is a learnable or fixed scaling coefficient and $\|x\|$ denotes input complexity (e.g., norm, entropy, or proxy task difficulty) [89].

This enables the router to allocate deeper or wider sub-networks for semantically rich or ambiguous inputs, while routes simpler tokens to lightweight modules. In multi-modal settings, adaptive expert counts have shown measurable gains in sample efficiency and robustness under domain shift, particularly when expert allocation aligns with modality boundaries (e.g., image tokens vs. text tokens). Unlike static MoEs, these models respond to context shifts by reshaping expert usage dynamically, without retraining.

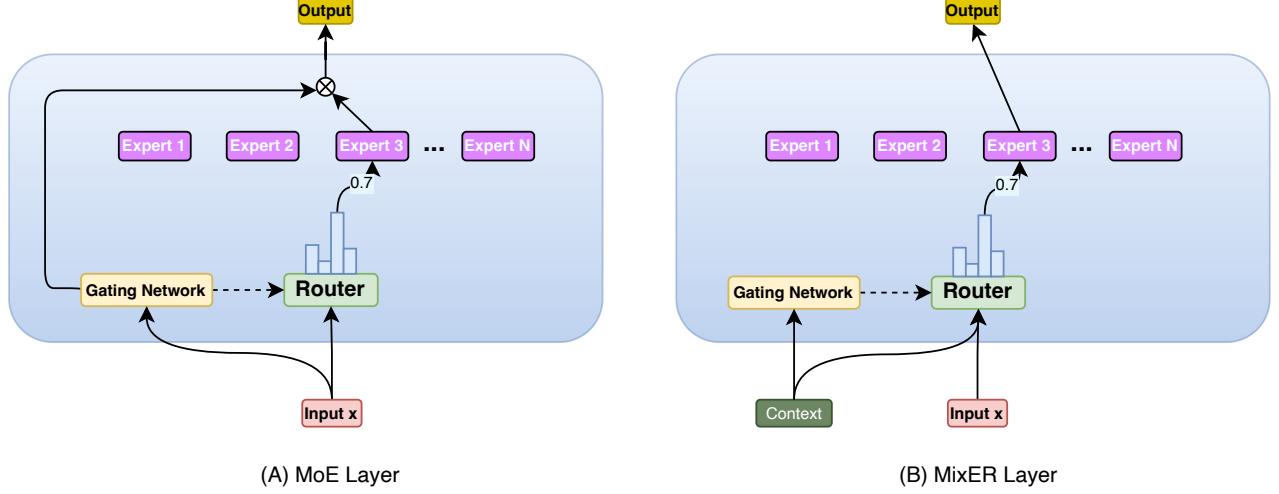


Fig. 5. Architectural comparison between standard MoE and MixER layers. (A) Standard MoE architecture where input x flows through the gating network to generate routing decisions, directing computation to selected experts via the router mechanism. (B) Enhanced MixER layer design that incorporates an additional context vector ξ alongside input x for routing decisions. The gating network leverages both inputs to compute expert selection probabilities, while the MixER approach eliminates the traditional softmax-weighted output combination used in conventional MoE implementations

In continual learning, dynamic routing helps preserve old knowledge while incorporating new, mitigating catastrophic forgetting [42]. It also improves robustness to distribution shifts, adapting expert usage to evolving input domains [5].

III. META-LEARNING AND KNOWLEDGE TRANSFER IN MOE

While the architectural foundations we've explored provide the building blocks for effective MoE systems, and the advanced variants of MoE model expertizes in dealing with a variety of critical challenges, a key question remains: how can these models learn to learn and transfer knowledge across different domains? This section examines meta-learning's development with MoE capabilities, which enables rapid adaptation and efficient knowledge transfer.

A. Meta-Learning Framework Design

Meta-learning enhances MoE systems by enabling rapid generalization across diverse tasks without retraining from scratch [6, 40]. Instead of learning task-specific routing independently, meta-MoE architectures optimize a routing policy θ across a distribution of tasks \mathcal{T} , enabling fast adaptation to unseen tasks with limited support data:

$$\theta_{\mathcal{T}_{\text{new}}} = \theta - \eta \nabla_{\theta} \mathcal{L}_{\text{support}}(\theta), \quad (18)$$

where $\mathcal{L}_{\text{support}}$ denotes the support loss induced by the sparse expert outputs.

Hierarchical meta-learning with MixER. To handle nested dynamical systems, the MixER model (Mixture of Expert Reconstructors) introduced in [51] augments classical MoE layers by feeding an additional context vector ξ alongside the input x into the router. As shown in **Figure 5**, MixER

bypasses softmax gating and instead performs discrete expert selection via a K-means-inspired objective:

$$z(x, \xi) = \arg \min_j \|f_{\theta}(x, \xi) - \mu_j\|^2, \quad (19)$$

where f_{θ} maps (x, ξ) to a latent space, and μ_j denotes the prototype of expert j .

This configuration enables top-1 routing with interpretable cluster assignments and avoids the overhead of differentiable soft selection. It demonstrates strong performance on sparse reconstruction tasks, such as parametric ODE systems. However, when contextual hierarchies are weak or missing, specialization degrades due to overlapping token-to-expert assignments.

Meta-distillation for domain adaptation. To address domain shift, the Meta-DMoE framework proposed in [84] formulates test-time adaptation as a meta-distillation problem. A set of domain-specific experts $\{E_i\}$ are pre-trained on disjoint source domains $\{\mathcal{D}_i\}$, and their predictions are aggregated through a transformer-based aggregator \mathcal{A} to supervise a lightweight student model S :

$$\mathcal{L}_{\text{meta}} = \text{KL}(S(x) \parallel \mathcal{A}(E_1(x), \dots, E_N(x))). \quad (20)$$

The aggregator learns to combine the expert outputs based on inter-domain dependencies, while meta-optimization ensures that the knowledge is transferable to unseen target domains. This approach improves generalization in scenarios where domain labels are unavailable or where training a unified model on heterogeneous domains is suboptimal.

B. Knowledge Transfer Mechanisms

Sparse-to-Dense knowledge integration. Transferring knowledge from sparse MoE models to dense architectures

presents significant challenges, prompting researchers to propose several effective solutions. Inspired by human pedagogical frameworks, a multi-teacher distillation strategy has been proposed, wherein multiple expert models collaboratively supervise a single student model. This setup enables the student to integrate diverse knowledge sources, improving generalization across tasks [78]. Their framework tackles some real problems with sparse MoE models, which tend to overfit, are tricky, and often don't play nice with existing hardware. The core idea is straightforward: Instead of trying to compress everything into one dense model directly, you let multiple experts teach a single student model. This knowledge integration approach sidesteps many of the headaches that come with sparse MoE deployment.

The proposed framework encompasses knowledge gathering and knowledge distillation phases. Four distinct knowledge gathering methods are investigated: summation, averaging, Top-K Knowledge Gathering (Top-KG), and Singular Value Decomposition Knowledge Gathering (SVD-KG). The dense student model (OneS) preserves 61.7% of MoE benefits on ImageNet, achieving 78.4% top-1 accuracy with only 15M parameters. On natural language processing datasets, OneS obtains 88.2% of MoE benefits while outperforming baselines by 51.7%. The approach achieves 3.7x inference speedup compared to MoE counterparts due to reduced computation and hardware-friendly architecture.

Mutual distillation among experts. A known limitation of MoE architectures is their narrow learning scope: individual experts often struggle to generalize due to restricted exposure to diverse training samples. This challenge has been addressed by incorporating mutual distillation mechanisms. The MoDE (Mixture-of-Distilled-Expert) framework [75] introduces moderate mutual distillation among experts, facilitating knowledge sharing and enhancing task awareness.

The MoDE framework addresses the fundamental issue where gate routing mechanisms restrict experts to limited sample exposure, thereby constraining generalization ability improvements. Through moderate mutual distillation, each expert acquires features learned by other experts, gaining more accurate perceptions of their originally allocated sub-tasks. Extensive experiments across tabular, natural language processing, and computer vision datasets demonstrate MoDE's effectiveness, universality, and robustness. The approach is validated through innovative "expert probing" studies. These validations demonstrate that moderate knowledge distillation improves individual expert performance, thus enhanced the overall MoE performance.

C. System Platform Support

The deployment of meta-learning-based MoE systems demands not only model-level optimization but also end-to-end systems engineering support. AwesomeMeta+ [67] addresses this need by introducing a standardized prototyping platform that encapsulates core meta-learning components as reusable and configurable modules. This modular design abstracts recurrent patterns such as task conditioning, gradient aggrega-

tion, and adaptation loops into composable units, simplifying integration into broader training and inference pipelines.

A key contribution of AwesomeMeta+ lies in its ability to translate theoretical meta-learning constructs into deployable MoE workflows. Traditional meta-learning implementations tend to be highly task-specific, limiting their extensibility and reproducibility. AwesomeMeta+ mitigates this through a layered architecture comprising: (i) a declarative model interface that maps task descriptors to expert selectors, (ii) a scheduler that optimizes expert instantiation under resource constraints, and (iii) an evaluation monitor that tracks transferability metrics such as few-shot accuracy or expert stability across tasks.

To assess platform effectiveness, both automated benchmarking and user studies were conducted. Feedback from over 50 researchers indicated that the framework improves understanding of meta-adaptation logic and accelerates system assembly. Moreover, measured performance across meta-dataset benchmarks confirms that the platform incurs negligible overhead while enabling consistent deployment of otherwise fragmented designs. By bridging theoretical innovation and practical engineering, AwesomeMeta+ demonstrates that platform-level standardization is essential for scalable, adaptable, and maintainable MoE-based meta-learning systems.

IV. MIXTURE OF EXPERTS APPLICATIONS AND DOMAIN-SPECIFIC MODELS

Having reviewed the theoretical framework and learning mechanisms of MoE, we now examine the applications of it in real world. In this section, we will examine the adoption of MoE across recommendation systems, search, computer vision, NLP, healthcare, etc. These models are transforming diverse fields by taking advantage of their specialized expertise to solve real-world challenges.

A. Recommendation Systems and Search

In large-scale recommendation and search applications, MoE architectures are increasingly employed to handle the inherent complexity of multi-domain, multi-task personalization [82, 22]. Conventional models often struggle to balance domain-specific signal extraction and shared knowledge transfer, particularly under varying user contexts and item scenarios.

To address these limitations, M3oE [82] proposes a modular MoE framework that jointly models domain- and task-level heterogeneity. It deploys three parallel expert modules to capture shared user preferences, domain-specific behavior, and task-specific patterns respectively. These modules are fused through a hierarchical gating design that allows adaptive aggregation of relevant signals across heterogeneous learning objectives. By integrating AutoML techniques for structure search, M3oE dynamically adjusts its expert composition, improving alignment with evolving personalization targets and ensuring performance scalability under real-world constraints.

As illustrated in **Figure 6**, AESM2 begins with a shared embedding lookup over input features, scenario indicators,

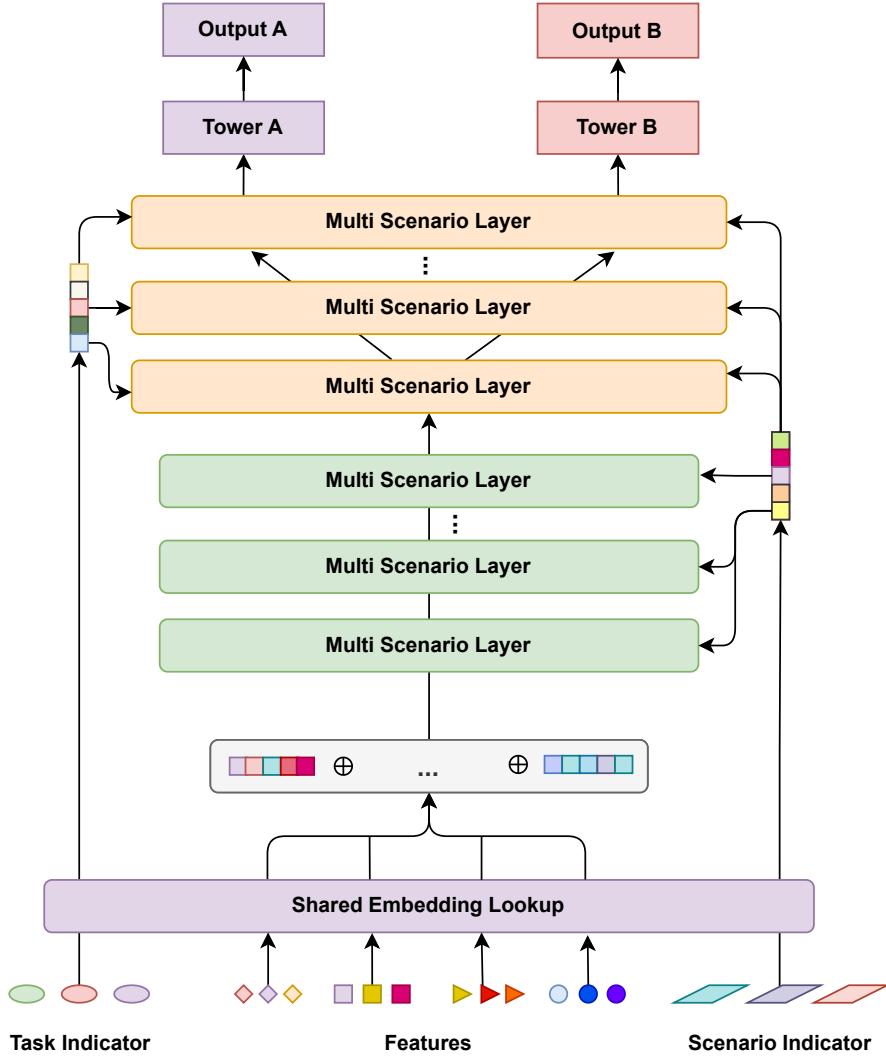


Fig. 6. **The Architecture of AESM² for Multi-Task Learning.** The framework is composed of major modules including Shared Embedding Layer used to cast the raw categorical and numerical features into continuous embeddings, and Multi-Scenario Layer used for expert selection, as well as Multi-Task Layer for multi task learning.

and task indicators. The encoded representations are processed through stacked multi-scenario layers, where each layer refines contextual understanding across traffic segments or environments. On top of these layers, scenario-informed representations are passed into multi-task-specific towers to generate differentiated outputs. The architecture employs a hierarchical routing mechanism that automatically selects experts at both the scenario and task levels. It does so without manual intervention, allowing the model to allocate computation to scenario-shared or scenario-specialized experts as needed. This structure enables explicit control over knowledge transfer across scenarios and tasks while maintaining model flexibility. Evaluations show that AESM² consistently improves retrieval quality and training stability under dynamic traffic distributions, outperforming static multi-gate MoE variants that lack automatic structure adaptation.

B. Multimodal and Multitask Learning

MoE architectures have proven effective in addressing the inherent complexity of multimodal and multitask learning, where diverse input modalities and objective functions must be handled simultaneously. Recent advancements in models such as MoVA [88], DeepSeek-VL2 [73], Omni-SMoLA [71], T-REX2 [28], MoME [62], and MoTE [66] collectively reflect a growing consensus that fine-tuning alone is insufficient for achieving robust generalization across heterogeneous tasks and data types.

Omni-SMoLA addresses a core limitation in large multimodal models: task interference from shared representations. By integrating low-rank expert modules specialized for different modalities and tasks, it enables modular specialization while preserving generalist capabilities. Empirical studies show that compared to standard LMM fine-tuning, SMoLA's sparse expert routing leads to more stable conver-

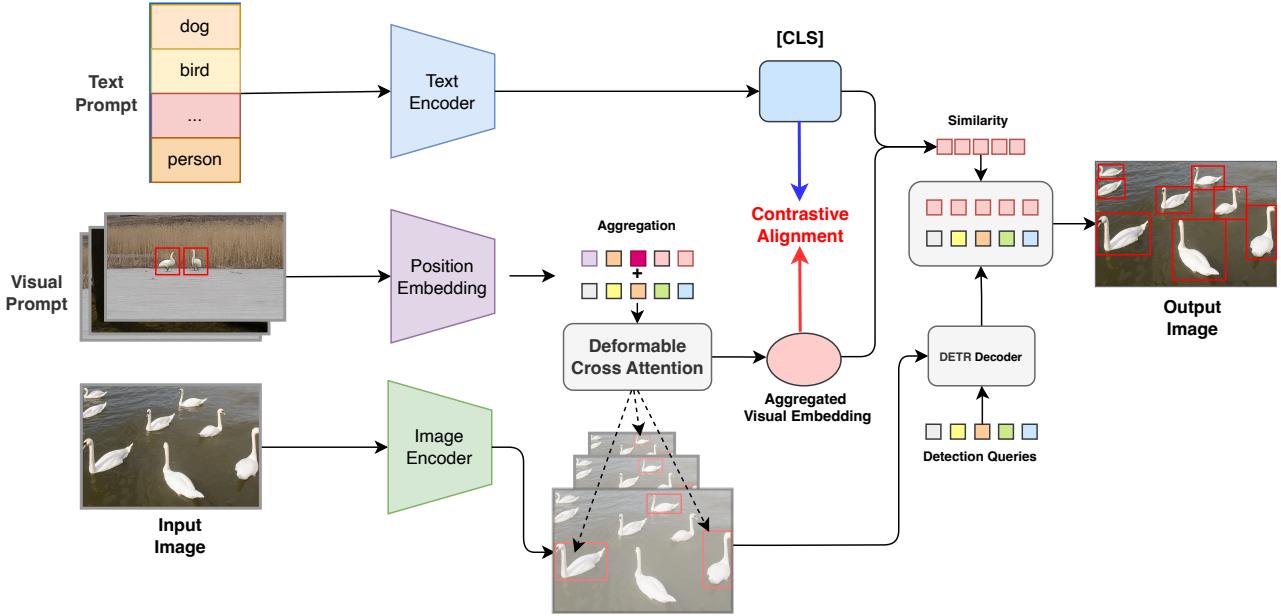


Fig. 7. Architectural overview of the T-Rex2 framework. The model adopts DETR-based design methodology for end-to-end object detection [7]. Visual and textual prompts are processed via deformable cross-attention [87] mechanisms and CLIP [56] text encoding respectively, with multimodal alignment achieved through contrastive learning strategies.

gence and better performance across diverse tasks, suggesting that architectural modularity is key to scalable multimodal generalization [71].

Multimodal specialization via hybrid prompts. T-REX2 further advances open-set object detection by exploiting the complementary strengths of textual and visual prompts. As depicted in **Figure 7**, the model incorporates dual encoders, with one responsible for processing abstract text categories (e.g., “dog”, “bird”), and the other extracting instance-level features from visual exemplars. These streams are fused through a deformable cross-attention module that integrates position-aware visual embeddings with context from prompts. A contrastive alignment loss links the aggregated visual embedding with the [CLS] representation from the text encoder, enforcing semantic consistency across modalities.

This hybrid prompting design enables robust generalization to unseen classes, where text offers conceptual grounding and visual prompts provide fine-grained instance cues. Unlike single-modal methods, T-REX2 flexibly accommodates different prompt combinations and achieves high zero-shot detection accuracy across diverse domains [28].

Expert reuse via hypernet-based modulation. HyperMoE addresses a central limitation in sparse expert models: the underutilization of unselected experts during forward inference. In standard MoE architectures, only the top- k experts are activated per input, leaving the rest unused. This underutilization of capacity can impair generalization, especially in multitask or low-resource scenarios.

To overcome this, HyperMoE introduces a hypernetwork that leverages the hidden states of inactive experts to generate

lightweight modulation signals. These signals are injected into the output paths of the active experts, enabling implicit expert collaboration without evaluating the full expert set. This mechanism retains routing sparsity while enriching the active computation with global knowledge from the broader model.

Empirical results show that such modulation enhances downstream performance, particularly on tasks with limited data, and improves expert diversity without increasing computational cost. The approach exemplifies how structural innovations can balance expressiveness and efficiency in scalable multitask systems.

C. Healthcare and Life Sciences

MoE architectures are increasingly applied in healthcare, addressing key challenges in patient care, clinical decision-making, and system efficiency. Representative models include Med-MoE [54], BiMediX [54], and LoRA-based medical MoEs [8]. Given the safety-critical nature of medical applications, these systems emphasize accuracy, modularity, and interpretability, aiming to support diagnostic reasoning while remaining compatible with clinical constraints.

A significant frontier is Embodied Intelligence in healthcare, where robotic systems assist with elderly care, rehabilitation, and clinical procedures. As shown in **Figure 8**, embodied agents leverage perception, actuation, planning, and memory across tasks ranging from bedside assistance to surgical support [39]. However, barriers such as limited integration into existing workflows, simulation-to-reality gaps, and the absence of standardized evaluation benchmarks persist, limiting widespread deployment.

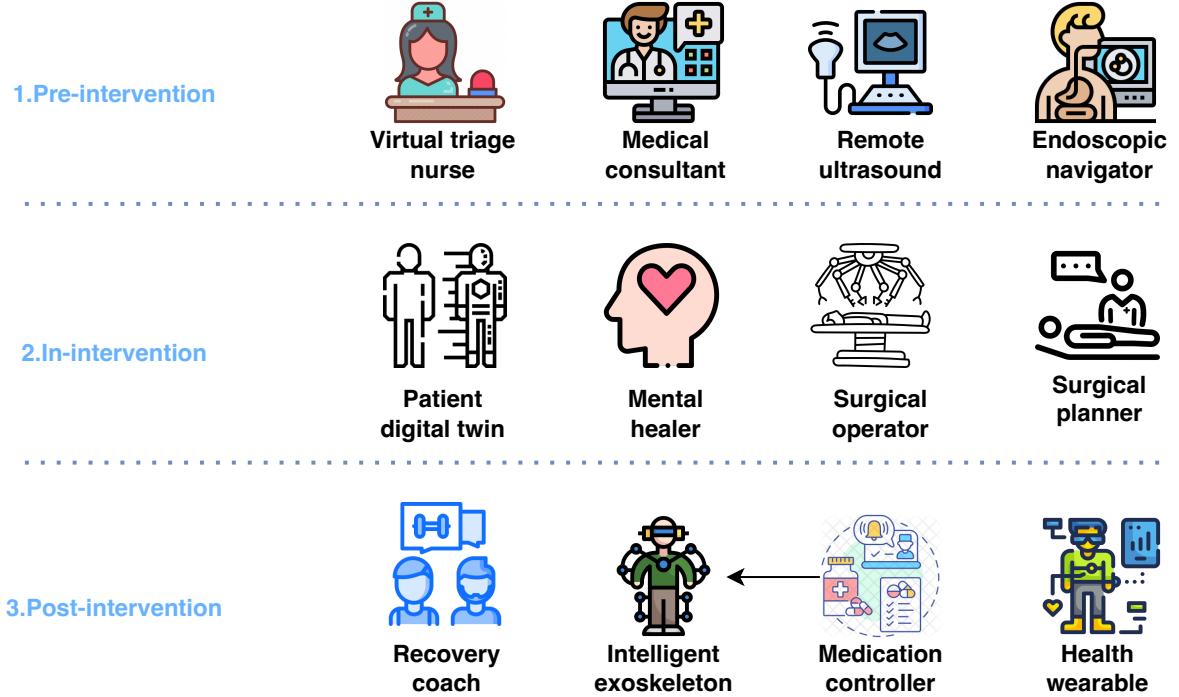


Fig. 8. **Applications of embedded AI in healthcare.** As the figure shows, embodied AI has been used in health care scenarios including but not limited to: pre-intervention (virtual triage nurse, medical consultant, remote ultrasound, endoscopic navigator), in-intervention (patient digital twin, mental healer, surgical operator, surgical planner), and post-intervention (recovery coach, intelligent exoskeleton, medication controller, health wearable).

To address the scarcity of medical data, Syn-Mediverse [45] introduces a large-scale synthetic dataset comprising over 48,000 hyper-realistic images and 1.5 million annotations across five vision tasks, enabling robust visual perception in complex healthcare environments. Building on this, AT-MoE [35] enhances interpretability and specialization through LoRA-tuned expert layers and grouped adaptive routing, which dynamically fuses task-relevant modules to support controllable, transparent decision-making. These advancements highlight the need to balance technological innovation with safety, interpretability, and ethical considerations in life-critical medical applications.

D. Computer Vision and Image Processing

Modern computer vision has progressed from CNN-based pipelines to Transformer and diffusion-based architectures. MoE designs are increasingly integrated into this evolution to address the growing complexity of tasks such as object detection, image classification, and scene understanding. Representative applications include AdaMV-MoE [9], GNT-MOVE [11], and recent efforts in image classification using expert-based decomposition [64].

In object detection, the Mixture of Calibrated Experts (MoCaE) [52] framework introduces a principled approach to aggregating predictions from multiple detectors. Traditional ensemble-based detectors suffer from miscalibrated confidence outputs, where dominant experts may overwhelm the consensus even in uncertain regions. MoCaE addresses this by calibrating each expert's output based on empirical performance, leading to more reliable prediction fusion. Evaluated on COCO and related benchmarks, MoCaE achieves up to +2.5 AP improvement, establishing new performance baselines.

The issue of effective expert specialization has also been tackled through refined gating architectures and regularization [29]. Early MoE models often failed to disentangle task-relevant features, resulting in homogeneous expert behavior. To address this, attention-like gates combined with entropy-minimizing regularizers enable low-overlap, semantically aligned expert selection. Experiments across classification datasets (e.g., MNIST [31], CIFAR [30], FashionMNIST [74]) demonstrate that this gating scheme improves both accuracy and interpretability of expert routing.

Deeper variants such as the Deep Mixture of Experts model [14] propose hierarchical expert compositions with stacked routing layers. This architecture separates “where” experts in early spatial layers from “what” experts in deeper semantic stages, allowing task-conditional specialization without inflating parameter count. These designs reflect a broader trend in visual MoEs: balancing model sparsity with fine-grained expressiveness.

By enabling modular learning over hierarchical visual signals, MoE frameworks have demonstrated efficacy in managing multi-scale feature hierarchies, adaptive capacity allocation, and task-specific specialization within scalable visual

pipelines.

E. Natural Language Processing and Large Language Models

Among all domains, NLP and LLMs have witnessed the most impactful and widespread adoption of MoE architectures. The fundamental motivation is clear: enabling scalable capacity without proportional increases in inference or training cost [3]. This has led to numerous architectural advances that have redefined efficiency paradigms in the tech industry and open-source communities alike.

A key limitation of early MoE systems is their incompatibility with parameter-efficient fine-tuning (PEFT), due to the need to store and update a full set of experts. To overcome this, a recent study proposes an extremely parameter-efficient MoE framework [79], which replaces dense expert networks with lightweight modules, achieving performance comparable to full fine-tuning while updating less than 1% of parameters in an 11B-scale model. This makes fine-tuning feasible even in constrained environments, without sacrificing task-specific adaptability.

Complementing this, flexible composition frameworks for domain adaptation have emerged. A toolkit for building low-cost Mixture-of-Domain-Experts (MoDE) [32] enables combining trained adapters or full models into expert pools tailored for specific domains. This method supports modular domain composition without retraining from scratch, offering practical deployment value. The framework includes guidance for optimal configuration and has demonstrated efficacy in multi-domain scenarios under constrained compute.

On the theoretical side, recent work on MoE hypothesis construction [60] deepens our understanding of their representational properties. Unlike Bayesian ensembles, which aggregate across uncertainty distributions, MoEs select hypotheses via discrete routing. This mechanism effectively enables abductive reasoning over the hypothesis space. Under mild assumptions, it is shown that MoE models exhibit higher functional capacity and can outperform Bayesian alternatives in certain regimes, even when relying on weaker inductive priors.

These findings collectively underscore the strength of MoE models in language representation tasks, combining modularity, adaptability, and efficiency. As a result, MoE architectures continue to play a central role in building scalable LLM systems that are both performant and cost-effective.

F. Methodological Innovation and Theoretical Foundations

The success of MoE architectures is grounded in a growing body of theoretical insights and methodological advancements that address scalability, convergence properties, expert diversity, and model integration.

Scaling behavior under memory constraints has been systematically explored through unified scaling laws for dense and sparse models [41]. By incorporating active parameter counts, dataset sizes, and expert configurations, the study demonstrates that MoE models can surpass dense models in memory efficiency, contrary to earlier assumptions. Empirical

evaluations over parameter regimes up to 5B validate this advantage, offering guidance for efficient training in constrained environments.

The convergence behavior of MoE parameter estimation has been formalized in recent work on Gaussian-gated models [49]. By analyzing maximum likelihood estimation (MLE) under gating covariates, the study introduces Voronoi-based loss functions to characterize non-uniform convergence rates. The analysis reveals how different configurations of location parameters lead to distinct solution spaces governed by polynomial systems, providing a deeper understanding of optimization dynamics in sparse settings.

Expert diversity and inter-expert knowledge transfer are addressed through several architectural strategies. Orthogonal MoE (OMoE) [17] enforces expert specialization via Gram-Schmidt orthogonalization of expert weights, while mutual distillation in MoDE and related frameworks [75, 23] encourages knowledge sharing through pairwise KL divergence. These mechanisms enhance feature coverage and mitigate redundancy. Nexus [20] introduces adaptive routing and parameter reuse, enables the conversion of dense models into sparse expert systems without retraining from scratch.

To improve generalization and model merging, hypernetwork-based strategies have emerged. HMoE [55] leverages hypernetworks to generate expert parameters dynamically, supporting low-dimensional alignment across domains. Concurrently, parameter merging frameworks [86] address conflicts in heterogeneous expert integration, introducing alignment and reparameterization strategies that preserve performance while minimizing interference.

Collectively, these developments reinforce the theoretical robustness of MoE models and expand their applicability to complex training regimes, where capacity, adaptability, and diversity must be simultaneously managed.

V. EVALUATIONS, CHALLENGES AND FUTURE DIRECTIONS

As MoE models gain wider adoption in real-world applications, critical questions about their evaluation and limitations are becoming increasingly important. Although traditional benchmarks used to evaluate LLMs, such as LLM-Perf Leaderboard [25], Mlperf inference benchmark [58] and Mmbench [76], are very popular and powerful tools, they are not well-suited in evaluating MoE models. A methodology to evaluate MoE models and provide guidance in the design choice, is one of the urgent needs. This section examines the development of methods for evaluating MoE performance, the remaining challenges, and the future directions for research and development.

A. Evaluation Framework and Methodology

Theoretical foundation and evaluation principles. The evaluation of MoE architectures presents unique challenges that stem from their fundamental design principles. Unlike traditional dense models, MoE systems are conditional computation, which only activates specific expert subsets for each

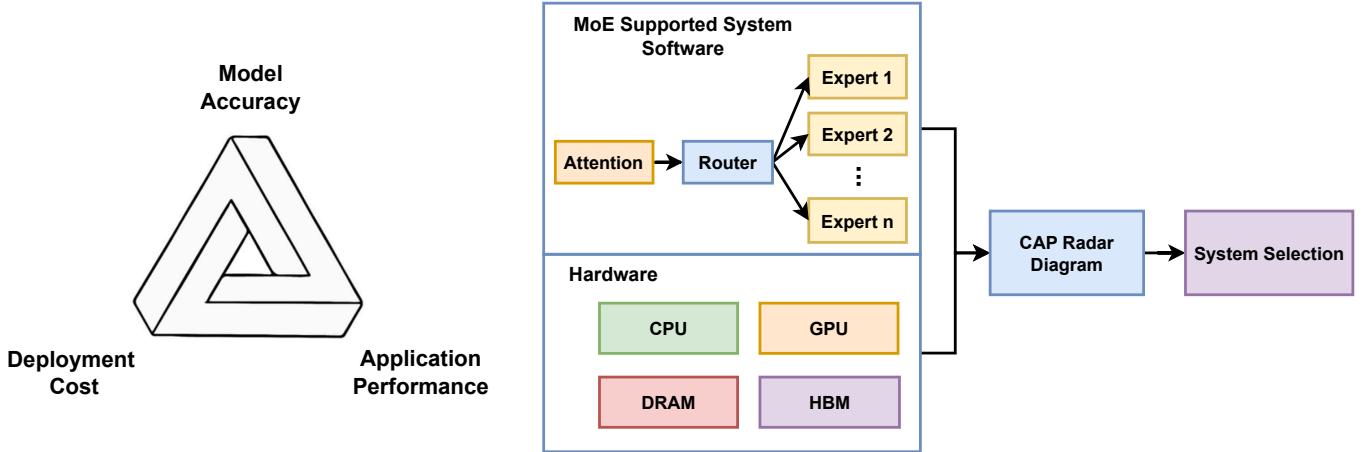


Fig. 9. **Framework illustration of MoE-CAP methodology.** Left: The triangular relationship demonstrates the balance among deployment cost, model accuracy, and application performance in MoE system design. Right: The MoE-CAP framework employs novel sparsity-aware evaluation metrics and CAP radar visualization to provide comprehensive assessment capabilities for MoE architectures, facilitating informed decisions regarding system architecture and hardware configuration selection.

input, and require special attention. In these architectures, there are many factors that affect the system's performance, including the interplay between expert specialization, routing. This evaluation frameworks need to be aware of those factors, and accurately capture their effects.

To be more specific, MoE's divide-and-conquer principle is more complex than dense models, as the problem space partition's effectiveness directly impacts model performance. Apart from evaluation of the final output, evaluation frameworks also needs to assess the intermediate processes of expert assignment, load balancing, and knowledge distribution across the expert ensemble, which is a end to end process.

Standardized evaluation framework. There is a critical need for comprehensive benchmarking platforms, since advancing MoE research has encountered a blocker: the lack of standard measure, and guidance. Recently researchers have been proposing studies on new benchmarks, some notable ones are Mixtral 8x7B [44] and LibMoE [50]. LibMoE is especially noteworthy since it presents a modular framework, covering streamlined research, training, and evaluation steps of MoE algorithms, basically the full lifecycle.

To validate the framework's benchmarking capabilities, researchers executed experiments and were astonished by the findings: by executing systematic evaluation over the five state-of-the-art MoE algorithms across three different large language models and eleven datasets under zero-shot settings [50], LibMoE have yielded critical insights. One notable finding is: although the evaluated MoE algorithms were developed on different task purposes, all of them achieved similar average performance across a wide range of tasks. This finding suggests that the choice of the MoE algorithm may be less critical than previously assumed.

System-level multi-dimensional evaluation. The deployment of MoE architectures in real-world systems requires evaluation beyond model accuracy. As shown in Figure 9, the MoE-CAP framework [18] introduces a triadic evaluation

paradigm encompassing *model accuracy*, *application performance*, and *deployment cost*. This CAP triangle highlights the inherent trade-offs among the three, where optimizing two often comes at the expense of the third.

To operationalize this evaluation, MoE-CAP integrates software- and hardware-level profiling. The software stack includes attention mechanisms, routers, and expert networks, while the hardware layer incorporates compute and memory components (e.g., CPU, GPU, DRAM, HBM). By analyzing the interaction between routing sparsity and hardware utilization, the framework generates CAP radar plots to support system comparison under constraints such as latency or budget. This enables structured, deployment-aware architecture selection aligned with practical feasibility rather than accuracy alone.

Specific evaluation method examples. The Mixture of Calibrated Experts (MoCaE) framework [52] offers a practical enhancement to MoE evaluation by addressing a persistent issue: *prediction miscalibration*. In conventional MoE models, expert outputs are often fused based on confidence scores that do not reliably reflect true accuracy, resulting in suboptimal predictions dominated by overconfident experts.

MoCaE mitigates this by introducing calibration procedures prior to expert output aggregation. Instead of directly averaging raw predictions, each expert's output is first adjusted to better reflect its empirical reliability. This calibrated fusion improves ensemble robustness and reduces overfitting to dominant experts. Empirical results on the COCO benchmark demonstrate performance gains of up to 2.5 AP, establishing MoCaE as a new state-of-the-art approach across multiple object detection tasks.

Expert diversity and representation learning challenges. A fundamental limitation in MoE architectures lies in the lack of expert specialization. Empirical studies have shown that experts often converge to nearly identical representations, with similarity scores exceeding 99% across diverse inputs [37].

Notably, this phenomenon is not confined to underperforming or poorly regularized models, but also emerges in high-performing configurations, suggesting it is a systemic issue.

Such representational homogeneity directly undermines the divide-and-conquer principle that motivates the MoE paradigm. Without sufficient diversity, experts fail to develop complementary capabilities, leading to inefficient parameter utilization and degraded task generalization. This challenge highlights the need for architectural and training mechanisms that explicitly promote expert differentiation.

Architectural design and integration challenges. The integration of shared layers within MoE architectures has been observed to degrade performance in certain configurations [53]. One plausible explanation is that experts learn redundant or conflicting representations when exposed to identical shared features, thereby reducing specialization and increasing interference. This suggests that naïve parameter sharing may hinder the decomposition of tasks and limit the expressiveness of individual experts.

In addition, dynamic expert expansion in incremental learning settings poses nontrivial challenges. When new experts are added post hoc [1], the system must resolve conflicts arising from inconsistent outputs among parallel experts. These inconsistencies can destabilize training or lead to suboptimal predictions. Addressing this requires the development of conflict-aware routing or mediation strategies that ensure coherence across dynamically evolving expert pools.

Routing mechanisms and specialization challenges. The necessity and efficacy of learned routing mechanisms in MoE systems remain an open and debated question. Empirical studies have shown that frozen, randomly initialized routers can perform on par with learned routing strategies across several benchmarks [15]. These findings challenge the commonly held assumption that adaptive routing is critical for MoE performance and suggest that increased routing complexity may not always yield proportional benefits. This raises important questions regarding the trade-off between routing expressiveness and architectural simplicity, particularly in low-resource or latency-constrained environments.

Challenges in Theoretical Grounding of MoE Architectures. Despite the empirical success of MoE architectures, particularly within natural language processing tasks, their theoretical foundations remain underdeveloped [10, 46]. Most existing designs rely on experimental heuristics rather than principled models, leaving unresolved key questions about the relationship between expert diversity, specialization dynamics, and system-level generalization.

A quantitative framework linking expert diversity to generalization and modular efficiency would guide principled MoE design, setting expert count, garsity and gating for any task or data. Formal analysis of routing and selection, grounded in information or learning theory, can replace costly empirical tuning and shrink the design space.

Technical method innovation. Efforts to improve expert diversity and routing precision in MoE architectures have led to a range of architectural innovations. Methods such as

DeepSeekMoE [12], TableMoE [81], and Pre-gated MoE [24] introduce input-aware expert allocation and gating preconditioning to enhance functional specialization. Orthogonal Mixture-of-Experts (OMoE) [17], in particular, applies orthogonality constraints on expert weights to reduce representational redundancy, thereby encouraging more disentangled expert behavior.

While such structural designs improve modularity, their potential can be further enhanced by incorporating feedback-based optimization. For instance, reinforcement learning techniques, including Reinforcement Learning from Human Feedback (RLHF), can guide expert selection and adjust routing policies using reward signals aligned with human preferences. These observations indicate that hybrid approaches combining architectural regularization with adaptive learning strategies hold significant promise for building more robust and generalizable MoE systems.

VI. CONCLUSION

This survey provides a comprehensive overview of recent advances in Mixture of Experts architectures. We first traced the evolution of MoE from its theoretical origins to large-scale implementations, followed by an in-depth analysis of core architectural components and design principles. We then examined advanced variants, meta-learning frameworks, and knowledge transfer mechanisms, as well as domain-specific applications across industry and academia. Furthermore, we discussed current evaluation methodologies, highlighted major challenges such as routing stability and expert specialization, and identified promising directions for future research. We hope this work serves as a valuable resource for researchers and practitioners, and contributes to the continued development of scalable and efficient MoE-based systems.

REFERENCES

- [1] Sebastian Agethen and Winston H. Hsu. Mediated experts for deep convolutional networks. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2016, Shanghai, China, March 20-25, 2016*, pages 2687–2691. IEEE, 2016.
- [2] Elie Antoine, Frédéric Béchet, and Philippe Langlais. Part-of-speech sensitivity of routers in mixture of experts models. In Owen Rambow, Leo Wanner, Marianna Apidianaki, Hend Al-Khalifa, Barbara Di Eugenio, and Steven Schockaert, editors, *Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025, Abu Dhabi, UAE, January 19-24, 2025*, pages 6467–6474. Association for Computational Linguistics, 2025.
- [3] Mikel Artetxe, Shruti Bhosale, Naman Goyal, Todor Miaylov, Myle Ott, Sam Shleifer, Xi Victoria Lin, Jingfei Du, Srinivasan Iyer, Ramakanth Pasunuru, et al. Efficient large scale language modeling with mixtures of experts. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11699–11732, 2022.

- [4] Emmanuel Bengio, Pierre-Luc Bacon, Joelle Pineau, and Doina Precup. Conditional computation in neural networks for faster models. *arXiv preprint arXiv:1511.06297*, 2015.
- [5] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48, 2009.
- [6] Rafael Bischof and Michael A Kraus. Mixture-of-experts-ensemble meta-learning for physics-informed neural networks. In *Proceedings of 33. forum bauinformatik*, 2022.
- [7] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020.
- [8] Qian Chen, Lei Zhu, Hangzhou He, Xinliang Zhang, Shuang Zeng, Qiushi Ren, and Yanye Lu. Low-rank mixture-of-experts for continual medical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 382–392. Springer, 2024.
- [9] Tianlong Chen, Xuxi Chen, Xianzhi Du, Abdullah Rashwan, Fan Yang, Huizhong Chen, Zhangyang Wang, and Yeqing Li. Adamv-moe: Adaptive multi-task vision mixture-of-experts. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17346–17357, 2023.
- [10] Zixiang Chen, Yihe Deng, Yue Wu, Quanquan Gu, and Yuanzhi Li. Towards understanding mixture of experts in deep learning. *arXiv preprint arXiv:2208.02813*, 2022.
- [11] Wenyan Cong, Hanxue Liang, Peihao Wang, Zhiwen Fan, Tianlong Chen, Mukund Varma, Yi Wang, and Zhangyang Wang. Enhancing nerf akin to enhancing llms: Generalizable nerf transformer with mixture-of-view-experts. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3193–3204, October 2023.
- [12] Damai Dai, Chengqi Deng, Chenggang Zhao, Rx Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai Yu, Y Wu, et al. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1280–1297, 2024.
- [13] Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. Glam: Efficient scaling of language models with mixture-of-experts. In *International conference on machine learning*, pages 5547–5569. PMLR, 2022.
- [14] David Eigen, Marc’Aurelio Ranzato, and Ilya Sutskever. Learning factored representations in a deep mixture of experts. In Yoshua Bengio and Yann LeCun, editors, *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Workshop Track Proceedings*, 2014.
- [15] Dongyang Fan, Bettina Messmer, and Martin Jaggi. Towards an empirical understanding of moe design choices. *arXiv preprint arXiv:2402.13089*, 2024.
- [16] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022.
- [17] Jinyuan Feng, Zhiqiang Pu, Tianyi Hu, Dongmin Li, Xiaolin Ai, and Huimu Wang. Omoe: Diversifying mixture of low-rank adaptation by orthogonal finetuning. *arXiv preprint arXiv:2501.10062*, 2025.
- [18] Yao Fu, Yinsicheng Jiang, Yeqi Huang, Ping Nie, Zhan Lu, Leyang Xue, Congjie He, Man-Kit Sit, Jilong Xue, Li Dong, et al. Moe-cap: Cost-accuracy-performance benchmarking for mixture-of-experts systems. *arXiv preprint arXiv:2412.07067*, 2024.
- [19] Wensheng Gan, Zhenyao Ning, Zhenlian Qi, and Philip S. Yu. Mixture of experts (moe): A big data perspective. *arXiv preprint arXiv:2501.16352*, 2025.
- [20] Nikolas Gritsch, Qizhen Zhang, Acyr Locatelli, Sara Hooker, and Ahmet Üstün. Nexus: Specialization meets adaptability for efficiently training mixture of experts. In *Adaptive Foundation Models: Evolving AI for Personalized and Efficient Learning*, 2024.
- [21] Nikhil Gupta and Jason Yip. Dbrx: Creating an llm from scratch using databricks. In *Databricks Data Intelligence Platform: Unlocking the GenAI Revolution*, pages 311–330. Springer, 2024.
- [22] Yupeng Hou, Shanlei Mu, Wayne Xin Zhao, Yaliang Li, Bolin Ding, and Ji-Rong Wen. Towards universal sequence representation learning for recommender systems. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 585–593, 2022.
- [23] Xinyi Huang, Yingyi Wu, Danyang Zhang, Jiacheng Hu, and Yujian Long. Improving academic skills assessment with nlp and ensemble learning. In *2024 IEEE 7th International Conference on Information Systems and Computer Aided Education (ICISCAE)*, pages 37–41. IEEE, 2024.
- [24] Ranggi Hwang, Jianyu Wei, Shijie Cao, Changho Hwang, Xiaohu Tang, Ting Cao, and Mao Yang. Pre-gated moe: An algorithm-system co-design for fast and scalable mixture-of-expert inference. In *2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA)*, pages 1018–1031. IEEE, 2024.
- [25] Régis Pierrard Ilyas Moutawakil. Llm-perf leaderboard. <https://huggingface.co/spaces/optimum/llm-perf-leaderboard>, 2023.
- [26] Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991.
- [27] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, De-

- vendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- [28] Qing Jiang, Feng Li, Zhaoyang Zeng, Tianhe Ren, Shilong Liu, and Lei Zhang. T-rex2: Towards generic object detection via text-visual prompt synergy. In *European Conference on Computer Vision*, pages 38–57. Springer, 2024.
- [29] Yamuna Krishnamurthy, Chris Watkins, and Thomas Gaertner. Improving expert specialization in mixture of experts. *arXiv preprint arXiv:2302.14703*, 2023.
- [30] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [31] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 2002.
- [32] Rhui Dih Lee, Laura Wynter, and Raghu Kiran Ganti. Flexible and effective mixing of large language models into a mixture of domain experts. *arXiv preprint arXiv:2408.17280*, 2024.
- [33] Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, De-hao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional computation and automatic sharding. In *International Conference on Learning Representations*, 2021.
- [34] Mike Lewis, Shruti Bhosale, Tim Dettmers, Naman Goyal, and Luke Zettlemoyer. Base layers: Simplifying training of large, sparse models. In *International Conference on Machine Learning*, pages 6265–6274. PMLR, 2021.
- [35] Xurui Li and Juanjuan Yao. At-moe: Adaptive task-planning mixture of experts via lora approach. *arXiv preprint arXiv:2410.10896*, 2024.
- [36] 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.
- [37] Boan Liu, Liang Ding, Li Shen, Keqin Peng, Yu Cao, Dazhao Cheng, and Dacheng Tao. Diversifying the mixture-of-experts representation for language models with orthogonal optimizer. In Ulle Endriss, Francisco S. Melo, Kerstin Bach, Alberto José Bugarín Díz, Jose María Alonso-Moral, Senén Barro, and Fredrik Heintz, editors, *ECAI 2024 - 27th European Conference on Artificial Intelligence, 19-24 October 2024, Santiago de Compostela, Spain - Including 13th Conference on Prestigious Applications of Intelligent Systems (PAIS 2024)*, volume 392 of *Frontiers in Artificial Intelligence and Applications*, pages 2966–2973. IOS Press, 2024.
- [38] Tianlin Liu, Mathieu Blondel, Carlos Riquelme Ruiz, and Joan Puigcerver. Routers in vision mixture of experts: An empirical study. *Trans. Mach. Learn. Res.*, 2024, 2024.
- [39] Yihao Liu, Xu Cao, Wei-Ting Chen, Yankai Jiang, Junjie You, Minghua Wu, Xiaosong Wang, Mengling Feng, Yaochu Jin, and Jintai Chen. From screens to scenes: A survey of embodied AI in healthcare. *Inf. Fusion*, 119:103033, 2025.
- [40] Zhaoyang Liu, Xijun Wang, Chenyuan Feng, Xinghua Sun, Wen Zhan, and Xiang Chen. Meta-reinforcement learning with mixture of experts for generalizable multi access in heterogeneous wireless networks. *arXiv preprint arXiv:2412.03850*, 2024.
- [41] Jan Ludziejewski, Maciej Pióro, Jakub Krajewski, Maciej Stefaniak, Michał Krutul, Jan Małański, Marek Cygan, Piotr Sankowski, Kamil Adamczewski, Piotr Miłoś, and Sebastian Jaszczerz. Joint moe scaling laws: Mixture of experts can be memory efficient. *arXiv preprint arXiv:2502.05172*, 2025.
- [42] James L McClelland, Felix Hill, Maja Rudolph, Jason Baldridge, and Hinrich Schütze. Placing language in an integrated understanding system: Next steps toward human-level performance in neural language models. *Proceedings of the National Academy of Sciences*, 117(42):25966–25974, 2020.
- [43] Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruti Shah, Xianzhi Du, Futang Peng, Anton Belyi, et al. Mml: methods, analysis and insights from multimodal llm pre-training. In *European Conference on Computer Vision*, pages 304–323. Springer, 2024.
- [44] MLCommons. Mixtral 8x7b: a new mlperf inference benchmark for mixture of experts. Technical report, MLCommons, 2024.
- [45] Rohit Mohan, José Arce, Sassan Mokhtar, Daniele Cataneo, and Abhinav Valada. Syn-mediverse: A multimodal synthetic dataset for intelligent scene understanding of healthcare facilities. *IEEE Robotics Autom. Lett.*, 9(8):7094–7101, 2024.
- [46] Siyuan Mu and Sen Lin. A comprehensive survey of mixture-of-experts: Algorithms, theory, and applications. *arXiv preprint arXiv:2503.07137*, 2025.
- [47] Basil Mustafa, Carlos Riquelme, Joan Puigcerver, Rodolphe Jenatton, and Neil Houlsby. Multimodal contrastive learning with limoe: the language-image mixture of experts. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, pages 9564–9576, 2022.
- [48] Huy Nguyen, Xing Han, Carl Harris, Suchi Saria, and Nhat Ho. On expert estimation in hierarchical mixture of experts: Beyond softmax gating functions. *arXiv preprint arXiv:2410.02935*, 2024.
- [49] Huy Nguyen, TrungTin Nguyen, Khai Nguyen, and Nhat Ho. Towards convergence rates for parameter estimation in gaussian-gated mixture of experts. In Sanjoy Dasgupta, Stephan Mandt, and Yingzhen Li, editors, *International Conference on Artificial Intelligence and Statistics, 2-4 May 2024, Palau de Congressos, Valencia, Spain*, volume 238 of *Proceedings of Machine Learning Research*, pages 2683–2691. PMLR, 2024.
- [50] Nam V Nguyen, Thong T Doan, Luong Tran, Van

- Nguyen, and Quang Pham. Libmoe: A library for comprehensive benchmarking mixture of experts in large language models. *arXiv preprint arXiv:2411.00918*, 2024.
- [51] Roussel Desmond Nゾyem, David A. W. Barton, and Tom Deakin. Towards foundational models for dynamical system reconstruction: Hierarchical meta-learning via mixture of experts. *arXiv preprint arXiv:2502.05335*, 2025.
- [52] Kemal Oksuz, Selim Kuzucu, Tom Joy, and Puneet K. Dokania. Mocae: Mixture of calibrated experts significantly improves object detection. *Transactions on Machine Learning Research*, 2024.
- [53] Sejik Park. Learning more generalized experts by merging experts in mixture-of-experts. *arXiv preprint arXiv:2405.11530*, 2024.
- [54] Sara Pieri, Sahal Shaji Mullappilly, Fahad Khan, Rao Anwer, Salman Khan, Timothy Baldwin, and Hisham Cholakkal. Bimedix: Bilingual medical mixture of experts llm. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 16984–17002, 2024.
- [55] Jingang Qu, Thibault Faney, Ze Wang, Patrick Gallinari, Soleiman Yousef, and Jean-Charles de Hemtinne. Hmoe: Hypernetwork-based mixture of experts for domain generalization. *arXiv preprint arXiv:2211.08253*, 2022.
- [56] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [57] Samyam Rajbhandari, Conglong Li, Zhewei Yao, Minjia Zhang, Reza Yazdani Aminabadi, Ammar Ahmad Awan, Jeff Rasley, and Yuxiong He. Deepspeed-moe: Advancing mixture-of-experts inference and training to power next-generation ai scale. In *International conference on machine learning*, pages 18332–18346. PMLR, 2022.
- [58] Vijay Janapa Reddi, Christine Cheng, David Kanter, Peter Mattson, Guenther Schmuelling, Carole-Jean Wu, Brian Anderson, Maximilien Breughe, Mark Charlebois, William Chou, et al. Mlperf inference benchmark. In *2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)*, pages 446–459. IEEE, 2020.
- [59] Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts. *Advances in Neural Information Processing Systems*, 34:8583–8595, 2021.
- [60] Bruce Rushing. Peirce in the machine: How mixture of experts models perform hypothesis construction. *arXiv preprint arXiv:2406.17150*, 2024.
- [61] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *International Conference on Learning Representations*, 2017.
- [62] Leyang Shen, Gongwei Chen, Rui Shao, Weili Guan, and Liqiang Nie. Mome: Mixture of multimodal experts for generalist multimodal large language models. *arXiv preprint arXiv:2407.12709*, 2024.
- [63] Qwen Team. Qwen2 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [64] Mathurin Videau, Alessandro Leite, Marc Schoenauer, and Olivier Teytaud. Mixture of experts in image classification: What’s the sweet spot? *arXiv preprint arXiv:2411.18322*, 2024.
- [65] An Wang, Xingwu Sun, Ruobing Xie, Shuaipeng Li, Jiaqi Zhu, Zhen Yang, Pinxue Zhao, J. N. Han, Zhanhui Kang, Di Wang, Naoaki Okazaki, and Cheng zhong Xu. Hmoe: Heterogeneous mixture of experts for language modeling. *arXiv preprint arXiv:2408.10681*, 2024.
- [66] Hongyu Wang, Jiayu Xu, Ruiping Wang, Yan Feng, Yitao Zhai, Peng Pei, Xunliang Cai, and Xilin Chen. Mote: Mixture of ternary experts for memory-efficient large multimodal models. *arXiv preprint arXiv:2506.14435*, 2025.
- [67] Jingyao Wang, Chuyuan Zhang, Ye Ding, and Yuxuan Yang. Awesome-meta+: Meta-learning research and learning platform. *arXiv preprint arXiv:2304.12921*, 2023.
- [68] Yaqing Wang, Sahaj Agarwal, Subhabrata Mukherjee, Xiaodong Liu, Jing Gao, Ahmed Hassan Awadallah, and Jianfeng Gao. Adamix: Mixture-of-adaptations for parameter-efficient model tuning. *arXiv preprint arXiv:2205.12410*, 2022.
- [69] Steve Waterhouse, David MacKay, and Anthony Robinson. Bayesian methods for mixtures of experts. *Advances in neural information processing systems*, 8, 1995.
- [70] Jialin Wu, Xia Hu, Yaqing Wang, Bo Pang, and Radu Soricut. Omni-smola: Boosting generalist multimodal models with soft mixture of low-rank experts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14205–14215, 2024.
- [71] Jialin Wu, Xia Hu, Yaqing Wang, Bo Pang, and Radu Soricut. Omni-smola: Boosting generalist multimodal models with soft mixture of low-rank experts. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16–22, 2024*, pages 14205–14215. IEEE, 2024.
- [72] Xun Wu, Shaohan Huang, Wenhui Wang, Shuming Ma, Li Dong, and Furu Wei. Multi-head mixture-of-experts. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024.
- [73] Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao

- Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, et al. Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding. *arXiv preprint arXiv:2412.10302*, 2024.
- [74] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashionmnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- [75] Zhitian Xie, Yinger Zhang, Chenyi Zhuang, Qitao Shi, Zhining Liu, Jinjie Gu, and Guannan Zhang. Mode: A mixture-of-experts model with mutual distillation among the experts. In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan, editors, *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2024, February 20-27, 2024, Vancouver, Canada*, pages 16067–16075. AAAI Press, 2024.
- [76] Cheng Xu, Xiaofeng Hou, Jiacheng Liu, Chao Li, Tianhao Huang, Xiaozhi Zhu, Mo Niu, Lingyu Sun, Peng Tang, Tongqiao Xu, et al. Mmbench: Benchmarking end-to-end multi-modal dnns and understanding their hardware-software implications. In *2023 IEEE International Symposium on Workload Characterization (IISWC)*, pages 154–166. IEEE, 2023.
- [77] Lei Xu, Michael Jordan, and Geoffrey E Hinton. An alternative model for mixtures of experts. *Advances in neural information processing systems*, 7, 1994.
- [78] Fuzhao Xue, Xiaoxin He, Xiaozhe Ren, Yuxuan Lou, and Yang You. One student knows all experts know: From sparse to dense. In Krystal Maughan, Rosanne Liu, and Thomas F. Burns, editors, *The First Tiny Papers Track at ICLR 2023, Tiny Papers @ ICLR 2023, Kigali, Rwanda, May 5, 2023*. OpenReview.net, 2023.
- [79] Ted Zadouri, Ahmet Üstün, Arash Ahmadian, Beyza Ermiş, Acyr Locatelli, and Sara Hooker. Pushing mixture of experts to the limit: Extremely parameter efficient moe for instruction tuning. In *The Twelfth International Conference on Learning Representations*, 2024.
- [80] Eduard Zamfir, Zongwei Wu, Nancy Mehta, Yuedong Tan, Danda Pani Paudel, Yulun Zhang, and Radu Timofte. Complexity experts are task-discriminative learners for any image restoration. *arXiv preprint arXiv:2411.18466*, 2024.
- [81] Junwen Zhang, Pu Chen, and Yin Zhang. Tablemoe: Neuro-symbolic routing for structured expert reasoning in multimodal table understanding. *arXiv preprint arXiv:2506.21393*, 2025.
- [82] Zijian Zhang, Shuchang Liu, Jiaao Yu, Qingpeng Cai, Xiangyu Zhao, Chunxu Zhang, Ziru Liu, Qidong Liu, Hongwei Zhao, Lantao Hu, Peng Jiang, and Kun Gai. M³oe: Multi-domain multi-task mixture-of experts recommendation framework. In Grace Hui Yang, Hongning Wang, Sam Han, Claudia Hauff, Guido Zuccon, and Yi Zhang, editors, *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2024, Washington DC, USA, July 14-18, 2024*, pages 893–902. ACM, 2024.
- [83] Hao Zhao, Zihan Qiu, Huijia Wu, Zili Wang, Zhaofeng He, and Jie Fu. Hypermoe: Towards better mixture of experts via transferring among experts. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 10605–10618. Association for Computational Linguistics, 2024.
- [84] Tao Zhong, Zhixiang Chi, Li Gu, Yang Wang, Yuanhao Yu, and Jin Tang. Meta-dmoe: Adapting to domain shift by meta-distillation from mixture-of-experts. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- [85] Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew M Dai, Quoc V Le, James Laudon, et al. Mixture-of-experts with expert choice routing. *Advances in Neural Information Processing Systems*, 35:7103–7114, 2022.
- [86] Yuhang Zhou, Giannis Karamanolakis, Victor Soto, Anna Rumshisky, Mayank Kulkarni, Furong Huang, Wei Ai, and Jianhua Lu. Mergeme: Model merging techniques for homogeneous and heterogeneous moes. *arXiv preprint arXiv:2502.00997*, 2025.
- [87] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable {detr}: Deformable transformers for end-to-end object detection. In *International Conference on Learning Representations*, 2021.
- [88] Zhuofan Zong, Bingqi Ma, Dazhong Shen, Guanglu Song, Hao Shao, Dongzhi Jiang, Hongsheng Li, and Yu Liu. Mova: Adapting mixture of vision experts to multimodal context. *CoRR*, 2024.
- [89] Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and William Fedus. St-moe: Designing stable and transferable sparse expert models. *arXiv preprint arXiv:2202.08906*, 2022.