# Collaborative Inference and Learning between Edge SLMs and Cloud LLMs: A Survey of Algorithms, Execution, and Open Challenges

SENYAO LI, School of Computer Science and Technology, Huazhong University of Science and Technology, China

HAOZHAO WANG, School of Computer Science and Technology, Huazhong University of Science and Technology, China

WENCHAO XU, Division of Integrative Systems and Design, Hong Kong University of Science and Technology, China

RUI ZHANG, School of Computer Science and Technology, Huazhong University of Science and Technology, China

SONG GUO, Department of Computer Science and Engineering, Hong Kong University of Science and Technology, China

JINGLING YUAN, Hubei Key Laboratory of Transportation Internet of Things, Wuhan University of Technology, China

XIAN ZHONG, Hubei Key Laboratory of Transportation Internet of Things, Wuhan University of Technology, China

TIANWEI ZHANG, College of Computing and Data Science, Nanyang Technological University, Singapore

RUIXUAN LI, School of Computer Science and Technology, Huazhong University of Science and Technology, China

As large language models (LLMs) evolve, deploying them solely in the cloud or compressing them for edge devices has become inadequate due to concerns about latency, privacy, cost, and personalization. This survey explores a collaborative paradigm in which cloud-based LLMs and edge-deployed small language models (SLMs) cooperate across both inference and training. We present a unified taxonomy of edge-cloud collaboration strategies. For inference, we categorize approaches into task assignment, task division, and mixture-based collaboration at both task and token granularity, encompassing adaptive scheduling, resource-aware offloading, speculative decoding, and modular routing. For training, we review distributed adaptation techniques, including parameter alignment, pruning, bidirectional distillation, and small-model-guided optimization. We further summarize datasets, benchmarks, and deployment cases, and highlight privacy-preserving methods and vertical applications. This survey provides the first systematic foundation for LLM-SLM collaboration, bridging system and algorithm co-design to enable efficient, scalable, and trustworthy edge-cloud intelligence.

CCS Concepts: • **Computing methodologies** → **Distributed artificial intelligence**; *Neural networks*; • **Computer systems organization** → Embedded and cyber-physical systems; Cloud computing.

Authors' Contact Information: Senyao Li, School of Computer Science and Technology, Huazhong University of Science and Technology, Wuhan, China; Haozhao Wang, School of Computer Science and Technology, Huazhong University of Science and Technology, Wuhan, China; Wenchao Xu, Division of Integrative Systems and Design, Hong Kong University of Science and Technology, Hong Kong, China; Rui Zhang, School of Computer Science and Technology, Huazhong University of Science and Technology, Wuhan, China; Song Guo, Department of Computer Science and Engineering, Hong Kong University of Science and Technology, Hong Kong, China; Jingling Yuan, yjl@whut.edu.cn, Hubei Key Laboratory of Transportation Internet of Things, Wuhan University of Technology, Wuhan, China; Xian Zhong, zhongx@whut.edu.cn, Hubei Key Laboratory of Transportation Internet of Things, Wuhan University of Technology, Wuhan, China; Tianwei Zhang, College of Computing and Data Science, Nanyang Technological University, Singapore; Ruixuan Li, rxli@hust.edu.cn, School of Computer Science and Technology, Huazhong University of Science and Technology, Wuhan, China.

## 1 Introduction

Large language models (LLMs) [139–142] have demonstrated remarkable proficiency across a broad spectrum of natural language processing tasks. However, their substantial computational and memory demands make on-device deployment on resource-constrained edge devices prohibitive [24, 25, 187]. With the rapid advancement of big data, cloud computing, and edge computing, edge-cloud collaborative intelligence has emerged as a promising paradigm for unlocking data value and delivering ubiquitous AI [35, 55, 163, 255]. This paradigm combines the cloud's computational power and generalization ability with the edge's responsiveness and adaptability [5, 176].

Rather than merely compressing large models for edge deployment, we focus on the systematic co-optimization of architecture, algorithms, execution, and privacy across cloud-based LLMs and edge-deployed small language models (SLMs) [232, 233, 235, 236]. By leveraging feature sharing, task partitioning, and knowledge transfer, LLMs and SLMs collaborate to deliver efficient and reliable intelligent services in heterogeneous environments [224–226]. Existing cloud-centric methods exploit centralized LLMs for strong generalization but incur high latency, privacy risks, and poor adaptation to localized contexts [57, 197, 247]. In contrast, edge-centric approaches offer fast response and personalization but are limited by on-device compute and generality [25, 118, 237–239]. Prior surveys have examined cloud-centric, edge-centric, and device-to-device cooperative paradigms, yet each faces key limitations in real-world deployment [162]. The LLM-SLM collaborative paradigm unifies these strengths by enabling hierarchical cooperation between large cloud models and small on-device models [6, 14].

Unlike federated learning [47, 279] or model compression [20], this approach treats LLMs and SLMs as distinct but interactive agents [72, 261], enabling more dynamic and modular collaboration across the edge-cloud continuum. However, realizing this vision poses several challenges across both inference and training. From the inference perspective, architectural heterogeneity between LLMs and SLMs complicates unified scheduling and deployment; varying task granularities, resource budgets, and latency constraints make coordinated model execution difficult [228, 230]; and network instability limits high-frequency interactions between models [21], such as model swapping [22] or cross-attention fusion, which are essential for maintaining semantic consistency and responsiveness. On the training side, heterogeneity in data distributions [23], task formulations, and model architectures hinders effective knowledge transfer, particularly for distillation and adaptation methods that assume structural or objective alignment. Moreover, the prevalence of non-IID edge data [159, 279] and the need for personalized model behavior exacerbate this issue, as local fine-tuning may lead to overfitting, while centralized updates risk diluting critical edge-specific or causally relevant patterns [11, 71]. These intertwined factors collectively underscore the need for more principled designs of collaborative learning and inference frameworks between heterogeneous models.

These limitations motivate a dedicated LLM-SLM collaboration framework that is structurally aware, distributionally robust, and capable of delivering efficient, reliable, and generalizable intelligence under dynamic, resource-constrained conditions.

In 2021, Alibaba DAMO Academy and Zhejiang University conducted the first large-scale survey on cloud-edge collaboration [262], identifying three synergy paradigms. In 2022, Alibaba's Top Ten Technology Trends highlighted "co-evolution between cloud-based large models and edge-side small models", emphasizing intelligent, privacy-aware end-device services [242]. Qualcomm echoed this in 2023 with Hybrid AI white papers [243, 245], advocating intelligent workload partitioning between LLMs and SLMs to address latency, efficiency, privacy, and personalization. These developments reflect a shift from centralized LLM pipelines to hybrid edge-cloud architectures, supported by emerging industrial-scale frameworks. Walle [264] offers an end-to-end system for development, deployment, and runtime, enabling cloud-edge collaboration across 300+ tasks with over 10 billion daily invocations, powering personalized recommendations, multi-modal understanding, and real-time 3D rendering via a lightweight edge inference engine. Luoxi [265] adopts a "slow-fast" learning strategy, where cloud LLMs generate latent representations to assist fast, personalized inference on edge devices with real-time
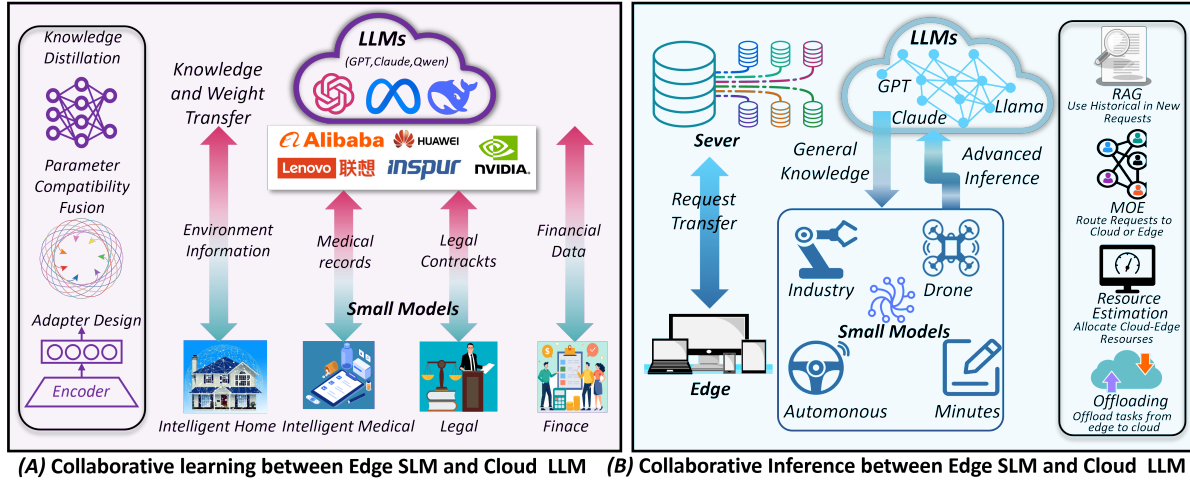
Fig. 1. Overview of cloud-edge collaboration workflows: (a) training and (b) inference.

feedback. InfiGUIAgent [266] demonstrates hierarchical edge reasoning through two-stage fine-tuning for GUI interaction in multi-modal scenarios. Real-world applications further underscore this paradigm's necessity: in vision and graphics (urban video analytics [184], autonomous driving [55, 56], XR), in language tasks (virtual assistants [240, 241], input methods, dialogue agents), and in personalized recommendation [30, 63, 221]. Edge intelligence is also expanding into multi-modal and sensor-based domains such as medical diagnosis [225], smart homes [204], and industrial monitoring, where energy efficiency and responsiveness are critical. Together, these academic and industrial efforts highlight both the technical feasibility and practical urgency of cloud-edge LLM-SLM collaboration, motivating a comprehensive survey of existing paradigms, system designs, and open challenges.

## 1.1 Related Surveys and Their Scope

Xu *et al.* [25] provide the first systematic overview of the full AIGC service lifecycle, from data collection and training to inference, and propose a collaborative cloud-edge-end infrastructure for mobile networks. Qu *et al.* [4] review LLM deployment at the mobile network edge, emphasizing resource-efficient techniques, architectural design, and edge LLM caching. In contrast, Xu *et al.* [246] prioritize system-level optimization, such as model compactification and token pruning, to accelerate inference and enhance deployment efficiency. For task offloading and resource allocation, Wang *et al.* [250] model task migration in mobile-cloud offloading, while [249] apply reinforcement learning to scheduling in CETCN scenarios. Xu *et al.* [248] survey decentralized continual learning on distributed devices, categorizing three algorithmic strategies for mitigating catastrophic forgetting and distributional shift. Building on this, Niu *et al.* [163] comprehensively review collaborative learning mechanisms between device-side small models and cloud-based large models from system, algorithmic, and application perspectives. Under constrained communication, Zhang *et al.* [24] propose air-ground collaborative inference strategies, and Lin *et al.* [187] introduce a multi-agent LLM framework for natural language tasks in 6G networks. Chen *et al.* [162] systematically summarize LLM-SLM collaboration paradigms, pipelining, routing, distillation, and fusion, as foundational strategies for efficient, adaptive intelligent systems. Existing surveys tend to focus narrowly on system deployment or isolated techniques (*e.g.*, distillation, pipelining) and lack a unified algorithmic

**Collaborative Inference and Learning between Edge and Cloud Large Language Models**

**Related Foundations (§1)**
- Related Surveys (§1.1)
- Structure Overview (§1.2)
- Related Foundations (§ 1.3)

**Collaborative Inference (§2)**
- Task Assignment(§2.1)
  - Resource, Uncertainty-Aware (§2.1.1)
  - Modular Collaboration Architectures (§2.1.2)
    - MoE-based approaches
    - Agent-based methods
- Task Division(§2.2)
  - Routing and Forwarding (§2.2.1)
    - Reward and Cost Bandit
    - Dynamic and Semantic
  - Computation Offloading (§2.2.2)
    - Structural Model Partitioning
    - Runtime-Adaptive Scheduling
  - Early-Exit (§2.2.3)
  - Communication Optimization (§2.2.4)
- Mixture: Task-Level(§2.3)
  - Task-Level Scheduling (§2.3.1)
  - Historical Utilization(§2.3.2)
  - Retrieval-Augmented Generation (§2.3.3)
    - Practical Applications
  - Mixture: Token-Level SD (§2.4)
    - Edge Draft, Cloud Validation (§2.4.1)
    - Self-speculative decoding (§2.4.2)
    - Semantic, Skeleton Completion (§2.4.3)
    - Token-Level Tree Validation (§2.4.4)

**Collaborative learning (§3)**
- Pruning, Quantization (§3.1)
- Distillation, Low-Rank (§3.2)
  - Task, Domain-Adaptive Distillation
  - Architectural Decoupling, Feedback-Guided Distillation
- Parameter compatibility, Model Convergence (§3.3)
- Adapter Modular (§3.4)
- Bidirectional Collaborative Learning (§3.5)
  - Dynamic Knowledge Transfer (§3.5.1)
  - Using small model to guide (§3.5.2)

**Benchmarks, Datasets (§4)**
- Datasets (§4.1)
- Benchmark (§4.2)

**Open Challenges (§5)**
- Privacy Preserving, Secure (§5.1)
- Application in vertical field (§5.2)

**Conclusion (§6)**
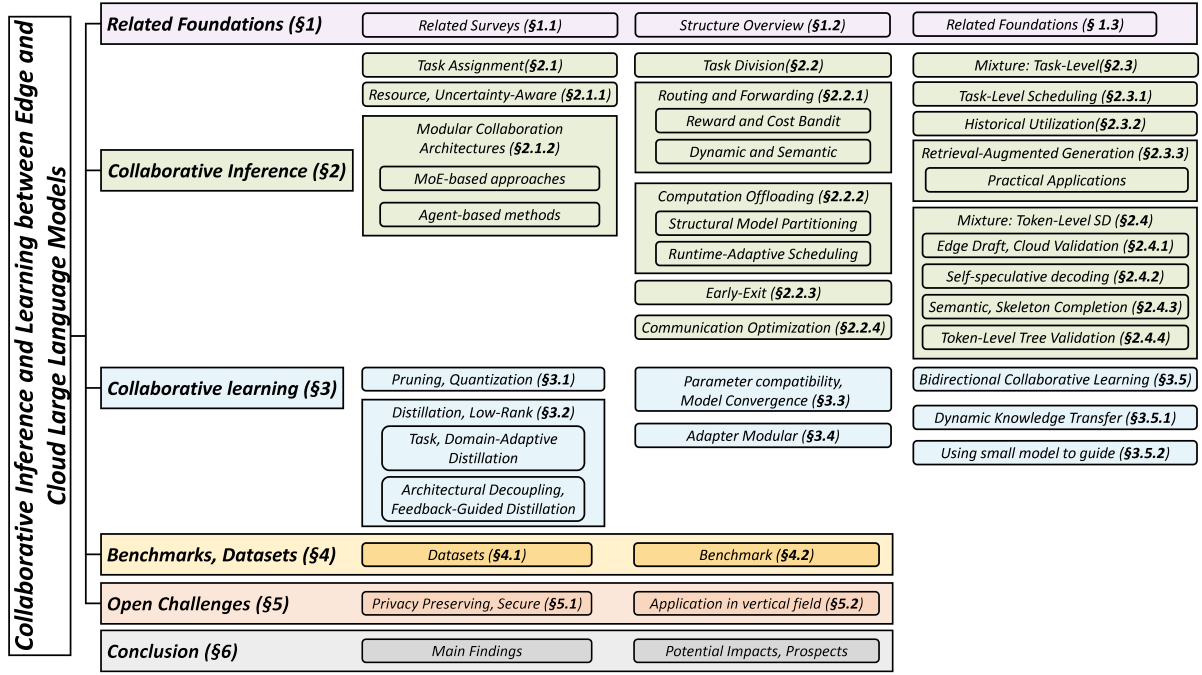- Main Findings
- Potential Impacts, Prospects

Fig. 2. Overview of collaborative paradigms between edge-cloud large and SLMs, structured along two axes: inference cooperation and training collaboration.

perspective on collaborative large-small model training and inference. Crucially, none abstract collaboration into a structured design space or offer a taxonomy capturing functional roles and interaction boundaries between heterogeneous models in edge-cloud settings.

**In contrast, this survey presents the first comprehensive review that jointly considers inference-time collaboration and training-time coordination between LLMs and SLMs.** As shown in Fig. 1 and Table 1, we propose a unified taxonomy of collaboration paradigms, encompassing task assignment, task division, and mixture (with task-level and token-level granularities), and analyze their alignment with algorithmic principles and system constraints. On the training side, we summarize collaborative adaptation methods, including bidirectional distillation, quantization, pruning, and low-rank approximation, which facilitate efficient SLM deployment without sacrificing performance. This survey places particular emphasis on the interaction between LLMs and SLMs across edge and cloud environments. By jointly analyzing collaboration across both inference and training phases, our review bridges algorithmic design and deployment needs, offering methodological insights and practical implications for future edge-cloud LLM systems. Moreover, we ground our discussion in recent literature published since 2023, ensuring both coverage and relevance. By jointly addressing both inference and training dimensions, our survey fills a key methodological gap and establishes a generalizable framework for designing large-small model cooperation in heterogeneous, resource-constrained environments.

## 1.2 Research Questions and Survey Structure

This survey systematically reviews the development, challenges, and future directions of edge-cloud collaborative inference and training with large and SLMs. We organize the discussion around the following core research questions (RQs):

Table 1. Comparison of this survey with related reviews

| Reference | Scope | Focus Area | Distinguishing Features |
|---|---|---|---|
| Xu *et al.* [25] | AIGC service lifecycle | Mobile network infrastructure design | Proposes end-cloud-edge architecture; omits model-level inference collaboration. |
| Xu *et al.* [246] | LLM system deployment | Inference optimization (compression, pruning) | Emphasizes hardware/system acceleration; ignores LLM-SLM interaction logic. |
| Wang *et al.* [249, 250] | Mobile-cloud task offloading | Scheduling and resource allocation | Models offloading strategies via RL; does not address collaborative inference. |
| Xu *et al.* [248] | Decentralized continual learning | Distributed algorithmic adaptation | Surveys lifelong learning; lacks explicit LLM-SLM collaboration structure. |
| Niu *et al.* [163] | LLM-SLM collaborative training | System, algorithm, application layers | Reviews collaborative training; does not cover inference across modalities/tasks. |
| Zhang *et al.* [24, 187] | LLMs in constrained networks (6G) | Multi-agent and aerial-ground collaboration | Focuses on special network conditions; proposes communication-specific inference strategies. |
| Chen *et al.* [162] | LLM-SLM collaboration mechanisms | Pipelining, routing, distillation, fusion | Summarizes interaction methods; lacks a unified algorithmic design space for inference. |
| This Survey | Heterogeneous LLM-SLM collaboration | Inference and training algorithms | First to cover end-cloud collaborative inference (task assignment/division/mixture) and training (parameter fusion, knowledge transfer). |

- **RQ1** *What is the edge-cloud collaborative inference paradigm with large and small models? What are its fundamental concepts and system architecture?* We focus on the basic paradigm, collaboration roles, and system-level design.
- **RQ2** *What are the major paradigms and collaboration patterns in edge-cloud inference?* We cover task scheduling, mixture cooperation at task and token levels, and speculative decoding. Examines how collaboration reduces latency, adapts to uncertainty, and improves efficiency.
- **RQ3** *Why and how should we study collaborative training between large and small models in edge-cloud environments?* We explore the necessity of training collaboration under heterogeneity, summarizing paradigms such as distillation, modular tuning, and parameter alignment for cross-device adaptation.
- **RQ4** *Why is it important to review benchmarks, privacy-preserving methods, and vertical applications in the context of edge-cloud collaboration?* We address fair evaluation standards, privacy challenges in training and inference, and practical requirements across real-world domains.

As shown in Fig. 2, we structure the survey as follows:

- **Section 1.3** introduces the fundamentals of edge-cloud collaboration with large and SLMs, covering system architectures, edge device constraints, and trends in large-model development (addresses RQ1).
- **Section 2** examines task assignment and inference strategies, with emphasis on dynamic scheduling and draft-refine verification frameworks, illustrating how edge and cloud models execute inference collaboratively and efficiently (addresses RQ2).
- **Section 3** surveys cloud-edge training architectures and deployment strategies, including quantization, pruning, parameter-efficient fine-tuning, and bidirectional knowledge transfer, to support efficient SLM deployment without sacrificing performance (addresses RQ3).
- **Section 4 - Section 5** summarize benchmarks, evaluation metrics, and domain-specific deployments to facilitate model comparison and practical implementation, and highlight open challenges such as generalizability, privacy, sustainability, and multi-agent collaboration (addresses RQ4).
- **Section 6** concludes with key insights and future research directions.

## 1.3 Related Concepts and Foundations

*1.3.1 Definition and Characteristics of LLMs.* LLMs typically refer to transformer-based architectures with billions to trillions of parameters. These models are pretrained on massive-scale corpora using strategies such as autoregressive learning (*e.g.*, GPT [141, 240]) or masked language modeling (*e.g.*, Qwen [235], DeepSeek [289]). Architecturally, LLMs stack tens to hundreds of transformer blocks, each comprising multi-head attention, feed-forward layers, residual connections, and normalization. As scale grows, LLMs exhibit emergent abilities, such as symbolic reasoning, instruction following, and multimodal understanding, that do not scale linearly with parameter count. In edge-cloud collaboration, LLMs are typically cloud-hosted due to their intensive computational and memory needs. While offering strong generalization and zero-/few-shot transfer across tasks, they also introduce significant communication overhead, as each inference often requires multiple cloud interactions, increasing latency and bandwidth load. Additionally, their training and fine-tuning involve high I/O costs and poor device adaptability, limiting their use in real-time or privacy-sensitive edge scenarios.

*1.3.2 Definition and Characteristics of SLMs.* SLMs, typically comprising millions to a few hundred million parameters, are built for resource-constrained environments while retaining core language modeling abilities. They simplify architecture via shallower transformer stacks, smaller hidden sizes, and fewer attention heads, greatly reducing memory and power use, enabling efficient deployment on mobile devices, microcontrollers, and edge nodes. SLMs emphasize responsiveness and local availability over broad reasoning. Though they lack LLMs' emergent abilities, they perform well in context-specific tasks like smart input, offline assistants, and embedded systems. Through knowledge distillation, LoRA fine-tuning, and instruction tuning, SLMs can approximate LLM outputs in narrow domains while maintaining low latency and data locality. In collaborative setups, they act as first-response agents or draft generators, offloading complex reasoning to cloud LLMs only when needed, thus optimizing communication cost and model efficiency.

*1.3.3 Complementarity of hardware architectures and application scenarios.* The foundation for cloud-edge LLM-SLM collaboration stems from hardware-driven capability partitioning: Cloud servers leverage multi-GPU clusters (*e.g.*, NVIDIA H100) with high-bandwidth interconnects (NVLink/RoCE) to run trillion-parameter models (*e.g.*, LLaMA-3 405B) for complex reasoning and global data processing [193], while resource-constrained edge devices, spanning smartphones (6-16GB RAM), XR headsets, and embedded systems (*e.g.*, Jetson Orin), deploy billion-parameter SLMs (*e.g.*, TinyBERT [157]) for latency-critical tasks (<10ms response) and local data optimization [192]. This architectural complementarity enables hierarchical task allocation: SLMs execute real-time perception (keyword detection, speech-to-text), privacy-sensitive preprocessing (medical data anonymization), and intent filtering (90% accuracy in customer service), offloading complex generation/reasoning to LLMs.

Table 2. Comparison of major paradigms in cloud-edge collaborative inference

| Category | Definition & Characteristics | Advantages | Limitations |
|---|---|---|---|
| **Task assignment § 2.1** | Routes requests to SLM or LLM based on confidence, resource, or task type. | Simple and fast; minimal communication overhead. | Hard switching may misroute uncertain inputs; lacks joint reasoning. |
| **Task division § 2.2** | Splits tasks into modules (*e.g.*, early exit, routing) processed by SLM and LLM. | Supports modular, adaptive inference with fine-grained control. | Requires explicit segmentation; incurs coordination overhead. |
| **Mixture: Task-level § 2.3** | Combines assignment and division across stages or roles via communication. | Flexible decomposition; leverages complementary strengths. | Dependent on partitioning quality and coordination efficacy. |
| **Mixture: Token-level § 2.4** | Collaborates at token generation (*e.g.*, speculative decoding). | Low-latency with accurate output; efficient cloud fallback. | Sensitive to draft quality; complex fusion and validation. |

Technical synergies like knowledge distillation (TinyBERT achieving 96% of BERT's performance [157]) and parameter-efficient tuning (LoRA boosting medical QA accuracy by 20%) further bridge capability gaps, with federated learning ensuring privacy when SLMs transmit anonymized features to cloud LLMs for multimodal fusion, collectively enabling 40% efficiency gains in applications like in-vehicle voice systems (200ms latency) and healthcare diagnostics.

## 2 Overview of Collaborative Inference

Collaborative inference between edge-deployed SLMs and cloud-based LLMs seeks to optimize latency and service efficiency under bandwidth and responsiveness constraints [288]. As shown in Table 2, recent work has explored three complementary collaboration paradigms. It is important to note that the classification presented here focuses specifically on inference-time collaboration between heterogeneous models. Although some mechanisms, such as mixture-of-experts (MoE), are applicable to both inference and training phases, this survey confines its scope to inference-oriented scenarios. That is, we summarize and compare collaborative strategies that are either designed for or primarily evaluated under inference workloads. We categorize and review existing works from the following perspectives:

- **Task assignment**, which routes queries to the most suitable model (*e.g.*, via agent-based or MoE strategies);
- **Task division**, which decomposes execution across models or system components (*e.g.*, offloading, request routing, or early-exit mechanisms);
- **Mixture strategies**, which integrate assignment and division at task-level or token-level granularity using techniques such as speculative decoding and multi-stage validation.

## 2.1 Task Assignment: Dynamic Model Selection for Cost-Quality Trade-offs

Task assignment determines whether a request is processed entirely by an edge-side SLM or a cloud-based LLM, in contrast to collaborative decoding, which involves inter-model interaction. To minimize latency and energy consumption under quality-of-service (QoS) constraints, prior methods employ lightweight scorers, calibrated reward estimators, or bandit-based controllers to select the optimal execution path before generation. Some approaches further incorporate user preference vectors or learned model-capability representations for
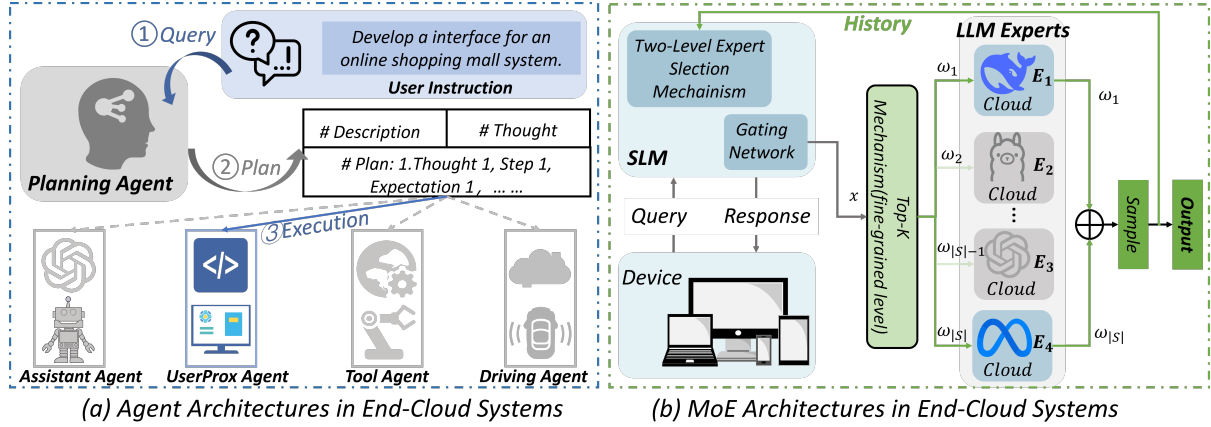
Fig. 3. Comparison of task assignment in edge-cloud systems: (a) Agents decompose user instructions to multi-stage plans executed by specialized agents; (b) MoE-based selection and fusion of top-$k$ cloud LLM experts via fine-grained gating.

adaptive, runtime routing. FS-GEN [84] introduces *collaboration frequency* to quantify the trade-off between large and small models in generation tasks: by aligning input spaces and defining a unified output-fusion objective, FS-GEN measures how often the cloud LLM intervenes in frequency-dominated sequences. Broadly, task assignment strategies fall into three architectural paradigms, resource- and uncertainty-aware assignment, agent-based scheduling and mixture-of-experts (MoE) frameworks (see Fig. 3), each reflecting different coordination granularities and model-selection philosophies.

### 2.1.1 Resource- and Uncertainty-Aware Task Assignment.
To leverage heterogeneous models while respecting edge-cloud capacity constraints (see Table 3), recent work explores diverse task allocation mechanisms that handle runtime uncertainty, system heterogeneity, and model capacity: FS-GEN [34] adopts a dual-system architecture where fast but uncertain SLMs serve as *System 1*, and reliable LLMs act as *System 2*, invoked when ambiguity arises. Fang *et al.* [10] propose a feedback-driven framework that adjusts routing decisions dynamically based on runtime uncertainty, model confidence, and system state, avoiding static thresholds. EdgeLLM [66] uses a value-density-first algorithm to rank tasks by cost-effectiveness, supporting preemption and batched execution via adaptive thresholds. Yang *et al.* [15] introduce fine-grained module-level partitioning, offloading heavy components (*e.g.*, attention, linear layers) to the cloud while retaining lightweight operations at the edge. U-VPA [202] employs uncertainty-guided sampling within a teacher-student framework, selectively uploading domain-informative samples to enhance generalization under non-stationary distributions. A hybrid cost function [123] identifies branch points for edge-cloud transitions, with fault-tolerant re-execution of failed branches locally. Ye *et al.* [204] dynamically select intermediate INT8 representations for edge deployment via automatic tuning. KDSL [85] leverages LLMs to generate logical rules via UCT-guided search. These rules are verified and deployed on the edge, enabling symbolic reasoning to refine cloud predictions.

### 2.1.2 Modular Collaboration via MoE and Agent Architectures.

*MoE-based approaches.* Mixture-of-Experts (MoE) architectures enable scalable, adaptive inference in edge-cloud settings. A lightweight on-device SLM routes queries via a gating or routing network to activate a sparse set of cloud LLM experts (see Fig. 3), whose responses are fused for soft assignment and diversity. MoE$^2$ [18] and EdgeMoE [50] guide expert routing using statistical priors or learned policies to reduce memory and I/O overhead. LiteMoE [251] further reduces cost by identifying and merging critical experts without retraining.

Table 3. Resource- and uncertainty-aware task assignment strategies

| Reference | Key Idea | Advantage | Limitation |
|---|---|---|---|
| PerLLM [115] | Personalized scheduling with constraint satisfaction | Jointly optimizes deployment and resource usage | Complex constraint modeling |
| EdgeLLM [66] | Value-density-first scheduling with preemption and batching | Cost-effective, real-time task prioritization | Preemption can destabilize under dynamic loads |
| FS-GEN [34] | Dual-system inference with uncertainty-triggered LLM fallback | Latency-efficient by conditional large-model invocation | Accuracy depends on uncertainty-estimation precision |
| Yang et al. [15] | Operator-level scheduling based on compute intensity | Fine-grained control of module placement | Requires precise operator profiling |
| Stammler et al. [123] | Hybrid cost-aware branching with fault-tolerant rollback | Reduces re-execution cost during cloud failures | Overhead in maintaining state consistency |
| U-VPA [202] | Uncertainty-guided sampling in a teacher-student setup | Bandwidth-efficient selection of informative samples | Needs robust domain-shift detection |
| Li et al. [125, 204] | INT8 model partitioning with automatic tuning | Efficient, low-latency dynamic deployment | Precision loss under aggressive quantization |
| KDSL [85] | LLM-based rule generation with feedback verification | Closed-loop correction with minimal cloud calls | Complexity in evaluating rule quality |

Extending selection to structuring, DoT [49] uses a task decomposition module and dependency-graph scheduler to partition complex queries into subtasks, enabling adaptive granularity and critical-path prioritization. Tian *et al.* [229] propose role-aware MoE routing for multi-turn dialogue, assigning segments to specialized experts and fusing outputs for contextually rich generation. CoEL [272] adds resource-adaptive coordination across edge devices, supporting elastic deployment and efficient scaling under varying budgets.

*Agent-based methods.* These methods use modular execution pipelines in which a planning agent interprets instructions, formulates structured plans, and delegates subtasks to specialized agents [283] (*e.g.*, assistants, tools, drivers) for context-aware, goal-driven coordination [275]. ARAG [274] integrates four agents, user understanding, natural language inference (NLI), context summarization, and item ranking, into a RAG pipeline. AgentVerse [95] deploys an edge-based coordinator to assemble cloud expert teams by task intent, supporting horizontal aggregation (*e.g.*, voting) and vertical communication for complex reasoning. Salve *et al.* [198] propose a specialized multi-agent framework with a central module orchestrating context-aware query generation and multi-source retrieval. In physically grounded scenarios, WebAgent [230] fuses LLMs with an embodied VirtualHome agent for goal-directed exploration and physical reasoning. ChatEval [228] introduces a debate-style framework where LLM-based judges collaboratively assess generation quality. EcoAgent [254] implements a closed-loop system with a cloud planner and two edge agents, one for execution and one for result verification, that triggers replanning upon failure. Finally, the MADRL framework [257] enables centralized cloud training and decentralized edge execution, reducing training overhead and inference latency in multi-agent reinforcement learning [255].

Table 4. Comparison of collaborative inference paradigms: routing, offloading, early exit, and communication optimization

| Paradigm | Core Innovation | Typical Scenario | Limitations |
|---|---|---|---|
| **Routing and Forwarding § 2.2.1** | Combines confidence, semantic planning, and contextual bandits for query routing. | Budget- or latency-constrained query routing. | Depends on accurate confidence estimates; fallback may add overhead or misroute. |
| **Computation Offloading § 2.2.2** | Supports dynamic scheduling and token-/layer-wise partitioning of models. | Latency-sensitive tasks under fluctuating resource availability. | Performance degrades in highly dynamic or non-stationary environments. |
| **Early Exit § 2.2.3** | Enables token- or layer-level early termination via dropout and pipeline reuse. | Low-latency generation or streaming applications. | Accuracy can drop at exit points; modeling token-dependent dynamics is challenging. |
| **Communication Optimization § 2.2.4** | Uses entropy compression, bandit routing, and adaptive switching to reduce transfer. | Resource-constrained settings with frequent edge-cloud interactions. | Balancing compression and semantic fidelity remains difficult; sensitive to distribution shifts. |

## 2.2 Task Division: Cooperative Subtask Decomposition Between Large and Small Models

Task division breaks hierarchical or modular tasks into semantic or functional parts, enabling LLMs and SLMs to jointly execute complementary subtasks. This cooperation improves parallelism and responsiveness under edge constraints. Unlike the categorization [287], we identify three fundamentally distinct paradigms of task division in edge-cloud collaborative inference, routing, offloading, and early exit, as illustrated in Fig. 4. Each paradigm offers unique advantages and is suited to different application scenarios, as summarized in Table 4.

*2.2.1 Routing and Forwarding Mechanisms.* Routing-based methods dynamically select models at inference time to balance latency, cost, and accuracy:

*Trust- and Semantic-Aware Routing.* These approaches use runtime confidence estimation and semantic planning for adaptive model selection. FrugalGPT [167] reduces API costs via prompt compression and cascaded routing. Tabi [168] employs calibrated confidence scores to choose between lightweight models and powerful LLMs. Dekoninck *et al.* [268] integrate pre- and post-inference quality estimation under a unified framework, improving routing accuracy and fallback success. Kag *et al.* [53, 205] introduce backup-block preloading for robust failover across heterogeneous devices. RouteLLM [106] organizes models hierarchically and applies dynamic programming to minimize tiered performance disparity.

*Reward- and Cost-Aware Bandit Routing.* These methods use online learning to optimize routing based on reward signals and cost-performance trade-offs. HybridLLM [97] defines a quality gap and uses a lightweight encoder to predict routing probabilities. ZOOTER [166] leverages query-level rewards for expert selection. RouterDC [94] uses dual contrastive loss to align queries with high-performing models. LLM Bandit [269] introduces "model identity vectors" for preference-conditioned routing, while MixLLM [169] and CITER [183] employ contextual bandits to estimate cost and quality. RouteT2I [165] selects between edge and cloud models for text-to-image tasks using multi-dimensional quality metrics.
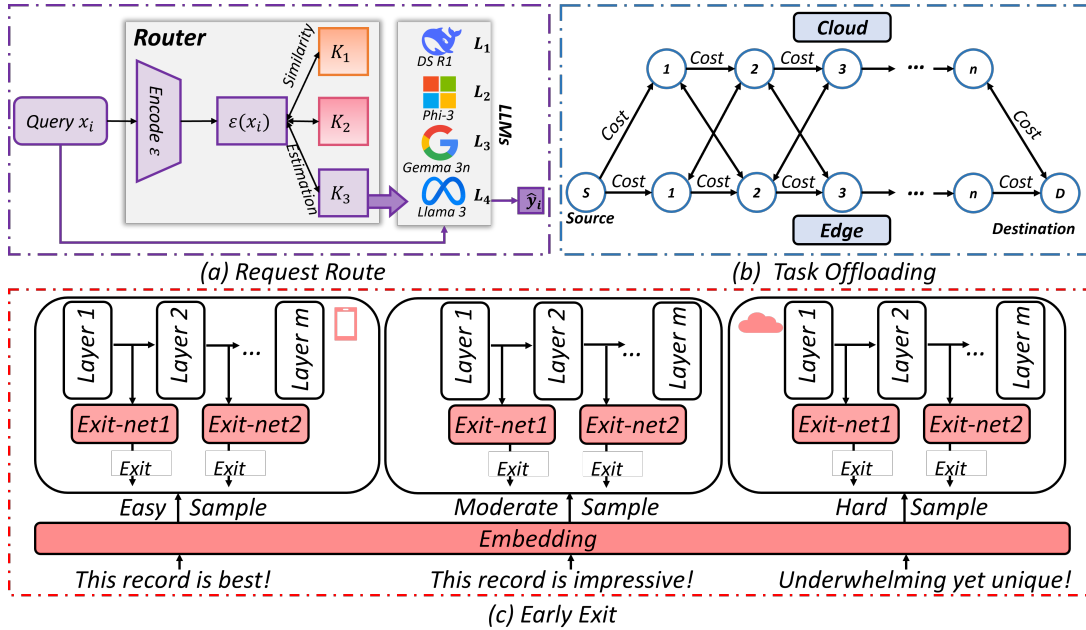
Fig. 4. Illustration of three execution-phase task-division strategies: (a) request routing, dispatching queries to optimal LLM experts based on semantics and cost; (b) task offloading, splitting computation graphs across edge and cloud for cost-aware load balancing; and (c) early exit, multi-exit networks that terminate early for simple inputs to reduce latency and energy.

2.2.2 *Computation Offloading and System-Level Optimization.* Offloading divides inference tasks across the device and cloud based on runtime conditions: rather than statically assigning tasks from the outset (as in routing), offloading makes dynamic, stage-wise decisions during execution, evaluating current workload against system capacity to maximize performance and utilize resources cost-effectively. We categorize offloading into two complementary types:

*Structural Model Partitioning.* This strategy decomposes inference by structurally partitioning models across device and edge/cloud at selected layers or token stages according to computational cost and depth [292]. ADAS [55] leverages internet of things (IoT) networks and proposes an enhanced task offloading algorithm based on DDPG, where a diffusion model is employed to generate noise and determine the optimal execution location for each task (*i.e.*, cloud, edge, or local). Liu *et al.* [51] assign lower layers to local devices and offload higher reasoning to edge servers. CE-CoLLM [13] uses token-level confidence to decide which tokens to process locally. Li *et al.* [125] automatically select INT8-quantized intermediate layers as partition points.

*Runtime-Adaptive Scheduling.* This group of methods dynamically schedules offloading decisions by adapting to real-time metrics, such as latency, confidence, and resource availability, and distributes tasks across local devices, control units, and auxiliary vehicles. Hao *et al.* [12, 124] use device metrics and confidence scores for fine-grained control. He *et al.* [52] introduce a reward-free policy based on latent states. Enhanced Hybrid Inference [65] incorporates user-aware utility models under bandwidth constraints. AVA [132] combines federated and multi-agent reinforcement learning for multi-tier distribution. While these approaches excel at adaptive, condition-aware scheduling, they often struggle to generalize in highly dynamic or heterogeneous environments.

*2.2.3 Early Exit.* Early-exit is a hybrid inference paradigm in which execution can terminate at intermediate layers based on system load or model confidence, reducing computation through dynamic exit decisions. LITE [259] introduces a confidence-guided early-exit strategy that halts inference without sacrificing generation quality. LayersKip [130] applies progressively increasing layer dropout and an early-exit loss across Transformer layers to support reliable early termination. EE-LLM [261], designed for 3D-parallelized LLMs, integrates token-level exits via key-value recomputation and pipeline parallelism [252], making it compatible with both large-scale training and inference. EESD [62] uses the initial $N$ layers of a base model and appends a single-layer Transformer to efficiently generate high-quality draft tokens. FREE [258] introduces temporally coherent parallelism by synchronizing shallow processing units with precomputed residuals, accelerating token generation with minimal depth. Collectively, these methods reduce latency and cost but still struggle to maintain consistent accuracy across exit points and coordinate exits under token-dependent dynamics.

Efficient communication is critical for collaborative inference in edge-cloud systems, enabling fine-grained reasoning through joint task division and communication-aware scheduling. Hu *et al.* [267] propose a hybrid architecture that dynamically delegates inference between lightweight edge models and powerful cloud LLMs, achieving 91% diagnostic accuracy with a 28.6% reduction in energy consumption. LLMCascades [111] uses a multi-stage framework in which edge predictions are either accepted or escalated to cloud models via voting and verification, enhancing trustworthiness with minimal redundancy. EdgeShard [117] reduces unnecessary cloud queries by forwarding only inference-critical token features, while entropy-based compression [17] further minimizes transmission cost without accuracy loss. For throughput optimization, PipeEdge [116] partitions LLMs across devices using pipeline parallelism and automated scheduling, reducing latency while preserving model capacity. PerLLM [115] frames coordination as a constrained multi-armed bandit problem, employing an upper confidence bound algorithm to adaptively select the most cost-effective execution path. Blending [271] enables turn-wise model switching in multi-turn interactions, leveraging model heterogeneity to improve response quality while maintaining cost efficiency. Despite these advances, minimizing overhead in dynamic environments and preserving semantic coherence across partitioned reasoning remain open challenges.

*2.2.4 Communication Optimization.*

## 2.3 Mixture: Task-Level Orchestration and Delegation in LLM-SLM Collaborative Inference

The *mixture paradigm* hybrids task assignment and execution division, allowing SLMs and LLMs to cooperatively handle a request through staged responsibility sharing. As shown in Fig. 5, we distinguish two granularities: *1) task-level mixture*, *2) token-level mixture*, where generation is shared step-by-step across models [129, 196, 203].

*2.3.1 Task-Level Decomposition, Scheduling, and Orchestration.* Task-level mixture uses semantic cues to identify subtask boundaries and assign them to edge or cloud models based on computational requirements, yielding interpretable, modular workflows. For example, MinionS [138] pipelines edge-side decomposition with cloud aggregation, while HybridSD [7] offloads high-level structural reasoning to the cloud and refines perceptual details on-device. IntellectReq [75] generates abstract intents in the cloud for lightweight edge execution. BAIM [64] integrates multiple edge-model outputs via a gating network for multi-modal collaboration, and OT-GAH [118] uses online tree pruning and assignment heuristics to maximize batch throughput. To streamline chain-of-thought reasoning, HAWKEYE [284] applies length-penalized reinforcement learning in the LLM and offloads detailed expansions to the SLM.

*2.3.2 Historical-Enhancement Collaborative Inference.* Users' queries to cloud APIs often mirror those sent to local models, motivating methods that leverage historical user-cloud interactions to improve on-device inference. Existing approaches include retrieval-augmented generation, cache-based scheduling, and collaborative learning. For instance, SlimPLM [174] uses heuristic confidence scores to trigger multi-stage retrieval only when the local

*(a) Task-Level Collaboration of LLM and SLM*      *(b) Token-Level Collaboration of LLM and SLM*
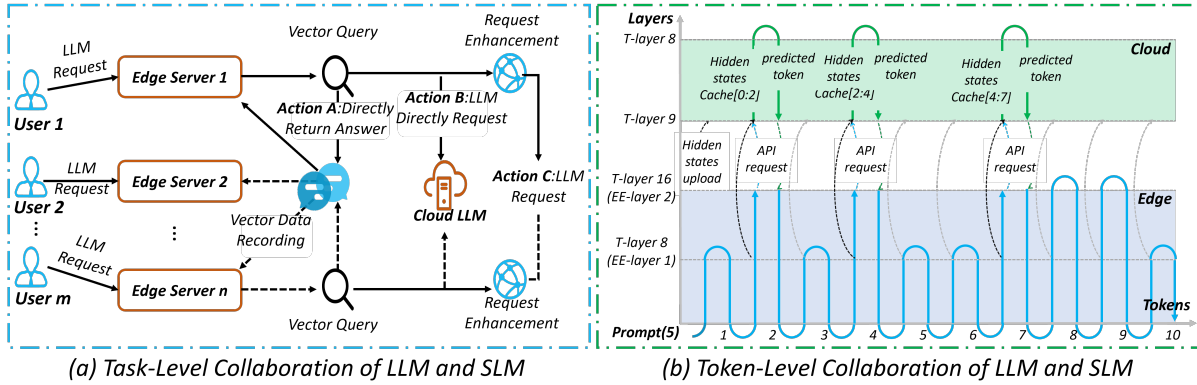
Fig. 5. Mixture strategies for edge-cloud model collaboration: (a) task-level, edge servers decide to respond locally, query the cloud, or forward to LLMs based on context and history; (b) token-level, edge SLMs generate easy tokens locally and offload harder ones to cloud LLMs via hidden-state sharing for fine-grained inference.

response is unreliable. VELO [45] caches prior request embeddings on edge servers and schedules execution paths based on vector similarity. Ding *et al.* [44] store interaction histories for nearest-neighbor retrieval, using subset selection to limit storage overhead. Hybrid-RACA [114] combines a cloud retriever with a lightweight edge predictor by transmitting compressed memory units for enhanced local inference without continuous cloud access. Xu *et al.* [83, 185, 186] introduce an iterative Direct Preference Optimization (DPO) mechanism, where a large model continually guides updates to a small on-device model.

*2.3.3 Retrieval-Augmented Generation with Self-Assessment and Feedback.* As shown in Fig. 6, these methods leverage external knowledge and reflective reasoning to enhance generation in complex tasks. Qin *et al.* [188, 282] analyze trade-offs among fine-tuning, RAG, data scale, model size, and task difficulty, demonstrating that RAG can outperform fine-tuning under resource constraints and that compressed LLMs benefit more from limited personalized data. Self-Knowledge Retrieval [170] prioritizes internal knowledge, invoking external retrieval only when needed via prompt learning or confidence heuristics. Self-RAG [171] introduces reflection tokens, special control signals that let the model explicitly assess retrieved passages. CRAG [172] employs a lightweight assessor to evaluate retrieval quality and trigger secondary retrieval or web expansion as necessary. Jiang *et al.* [200] extend introspective signals by allowing the model to either incorporate or recalibrate retrieved content. RA-ISF [173] decomposes failed queries into subquestions via a three-stage framework, self-assessment, retrieval, and query decomposition, and integrates subanswers into the final response. SlimRAG [273] builds a lightweight entity-to-fragment index on the cloud for salient entity detection, reducing graph complexity and communication. SpeculativeRAG [211] clusters documents for parallel draft generation.

In practice, Glocker [80] applies timestamp-enhanced RAG in autonomous home object management, enabling robots to track past actions while executing high-level instructions. Zhu *et al.* [190] introduce sparse context selection, parallelizing document encoding and decoding only the most relevant cached content via control tokens. SparseRAG [231] integrates document evaluation and response generation to minimize context loading, accelerating inference and improving output quality for both short- and long-form tasks.

## 2.4 Mixture: Token-Level Speculation and Verification in Edge-Cloud LLM-SLM Collaboration

Most LLMs generate tokens autoregressively, producing each token sequentially for the next prediction. While this yields high-quality output, it incurs substantial latency, which limits real-time deployment on resource-constrained
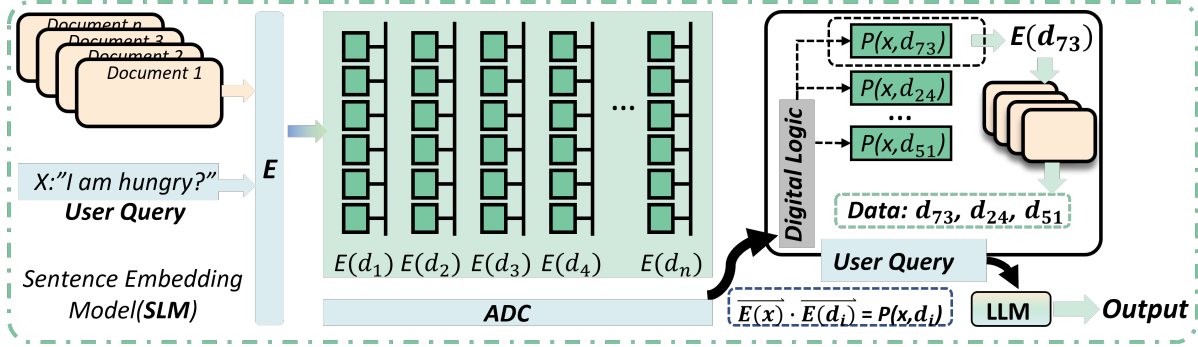
Fig. 6. Hybrid SLM-LLM collaboration in RAG: the SLM encodes queries and retrieves documents via approximate nearest-neighbor search, and the LLM generates responses using the query and retrieved context.

devices. To address this, recent work [244] proposes a "lightweight drafting + precise verification" paradigm: the edge-deployed SLM rapidly generates high-confidence drafts, and the cloud-based LLM verifies semantic consistency, corrects uncertain tokens, and enriches content with additional detail [276]. As shown in Fig. 7, we survey speculative decoding architectures and techniques in edge-cloud collaboration, covering draft-refine frameworks, draft-completion strategies, skeleton-completion mechanisms, and token-level verification [14].

*2.4.1 Edge Draft and Cloud Validation.* This line of work divides into two categories: *1)* vanilla speculative decoding frameworks with algorithmic enhancements and *2)* system-level designs that improve parallelism and reduce latency under diverse deployment constraints.

*Vanilla Speculative Decoding and Algorithmic Enhancements.* Early works such as [32, 40, 208] introduced the basic speculative decoding paradigm. A fallback mechanism [207] lets the lightweight model defer uncertain tokens to the larger model, while a rollback strategy corrects inaccuracies. Building on this, RSD [37] and SpecExec [42] incorporate controlled bias to prioritize high-reward outputs. AutoMix [92] integrates symbolic reasoning with reinforcement learning to mitigate hallucinations [213] and reasoning errors, using kernel density estimation to close the feedback loop between on-device verification and cloud routing. Fu *et al.* [215] improve efficiency by alternating the draft and target models as proposer and verifier via a learned verification distribution.

*Parallel and Low-Latency Speculative Inference Frameworks.* Interactive applications (*e.g.*, question answering, voice assistants, real-time translation) demand low Time-To-First-Token (TTFT) and Token-By-Token Time (TBT). DiSCo [41] estimates the acceptance probability of each speculative token using the drafting model's logits. To avoid mutual stalling, where the draft model waits for validation, SpecDec [100, 109] performs parallel token verification with soft discrimination. PEARL [191] uses adaptive draft lengths and early verification to accelerate decoding, and SEED [210] manages drafts with a round-robin FCFS queue.

*Multi-Layer and Lightweight Verification Pipelines.* To reduce cloud load, [135, 195, 219, 223] propose multi-layer pipelines: a lightweight verifier filters token blocks before final cloud validation [217]. Falcon [194] boosts intra-block token dependency via coupled sequential scanning, improving speculative accuracy while maintaining a compact two-layer transformer draft model suited for lightweight, multi-layer verification. Spector [128] adds an extra draft stage using a simpler model to pre-filter low-confidence hypotheses. [291] uses layer parallelism via early exits for mid-verification token drafting. Clients pre-draft tokens using idle time before final validation, boosting edge-cloud parallelism. This achieves a 21% speedup over cloud decoding on the Unitree Go 2 robot.
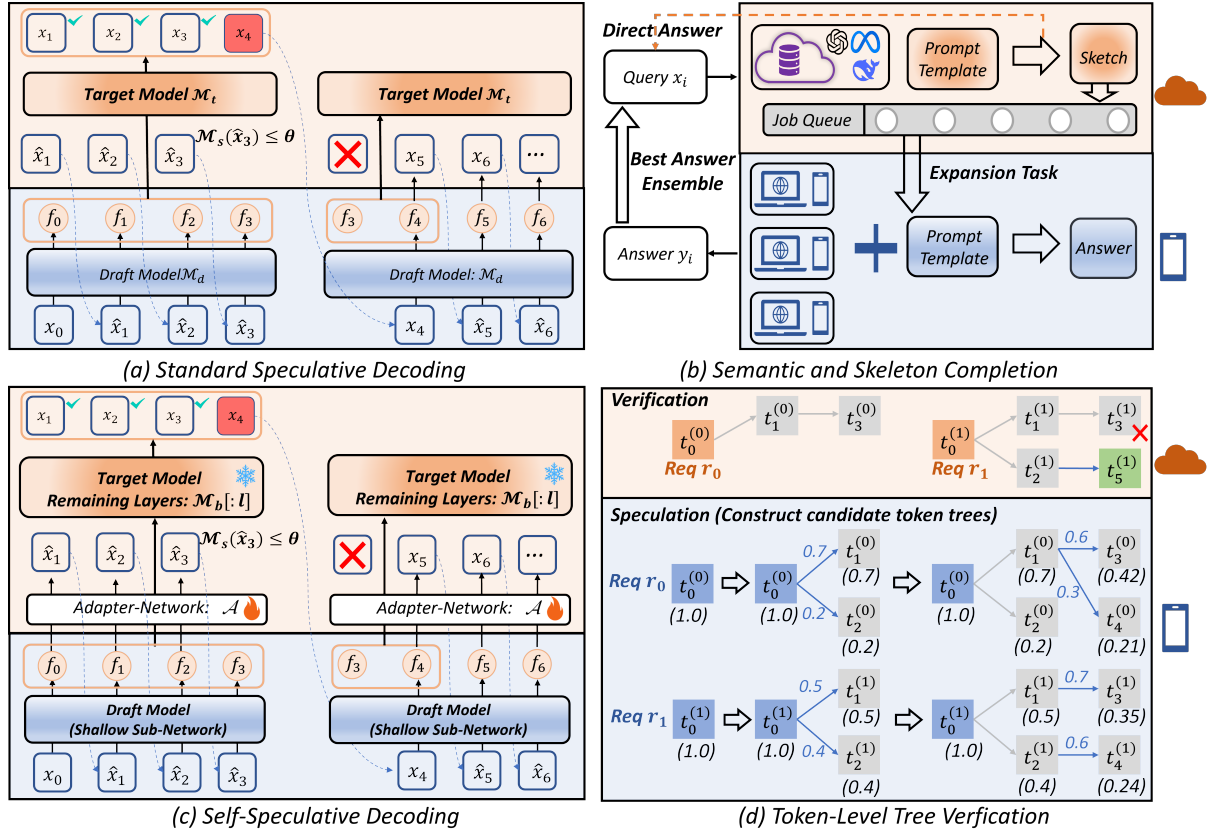
Fig. 7. (a) Vanilla: the cloud verifies edge-generated drafts and rejects or regenerates invalid tokens; (b) Semantic Skeleton: the cloud generates high-level semantic skeletons via prompts, which the edge completes with identical or refined prompts; (c) Self-Speculative: the model is split into shallow and deep subnetworks, using adapters for lightweight on-device verification; (d) Token Tree: token trees are expanded probabilistically, and the cloud verifies optimal branches per request.

*Draft Quality Enhancement and Hardware-Aware Scheduling.* Fuzzy Speculative Decoding (FSD) [27] relaxes strict match constraints, allowing minor mismatches to boost speed [221]. Zhao *et al.* [126, 136] propose phrase-level reuse, caching validated phrases for rapid generation. DistillSpec [137] distills the target model into the draft model offline using self-sampled data. DuoDecoding [28] assigns drafting to the CPU and validation to the GPU, dynamically adjusting draft lengths based on system load [218]. SPIN [290] uses a learning-based SSM selection without request priors and accelerates execution via GPU pipelining of speculation and verification. BanditSpec [214] formulated adaptive hyperparameter selection problem when generating text as a multi arm Bandit problem and designed two bandit based hyperparameter selection algorithms.

*2.4.2 Self-Speculative Decoding without Auxiliary Draft Models.* Some studies [61, 216] note that adding a separate draft model increases communication, storage, and training costs [253]. Kangaroo [38] reuses the target LLM's shallow subnetwork and LM head for self-drafting, significantly reducing overhead. SWIFT [61] exploits context-aware layer skipping for dynamic computation, and ASD [216] extends speculative sampling to diffusion models

Table 5. Comparison of cloud-to-edge and edge-to-cloud skeleton collaboration paradigms

| Paradigm | Ref. | Key Ideas | Advantages | Limitations |
|---|---|---|---|---|
| Cloud-to-Edge | PICE [43] | Progressive inference: cloud LLM drafts, edge SLM refines | Reduces latency and supports incremental, real-time reasoning | Requires fine-grained synchronization between LLM and SLM |
| | CoGenesis [33] | Sketch planning by LLM with local completion; distribution-only collaboration via logits | Enables privacy-preserving inference and adapts to bandwidth constraints | Coarse plans may limit expressiveness; logit protocols are model-specific |
| | NEST [201] | Cloud performs full retrieval, edge conducts token-level neighbor search for caching | Reduces corpus storage and improves retrieval latency | Requires careful tuning of retrieval-token matching; introduces multi-stage complexity |
| Edge-to-Cloud | SlimPLM [174] | Triggers LLM only when edge model uncertainty is high | Minimizes cloud calls and provides fast edge responses | Depends on accurate uncertainty detection; may miss subtle errors |
| | Hao et al. [14] | Token correction: LLM refines output by editing a few edge tokens | Efficient corrections with minimal cloud usage | Requires reliable error localization; not suited for structural rewrites |
| | Probe Sampling [209] | Measures draft-target similarity to decide cloud fallback | Lowers cloud overhead while maintaining output fidelity | Relies on precise similarity estimation; risk of under-correction |

by leveraging DDPM exchangeability for parallel inference. These methods maintain speed gains without extra model training.

*2.4.3 Semantic and Skeleton Completion.* The draft-refine paradigm splits generation into two stages: edge-side drafting and cloud-side refinement (see Table 5).

*Cloud-to-Edge Skeletons.* This class of methods adopts a top-down strategy, where the cloud LLM first produces a high-level semantic "skeleton", and the edge-deployed SLM complements or adapts the content based on local context or constraints. PICE [43] generates a task sketch for concurrent edge refinement, reducing latency. CoGenesis [33] offers *1)* sketch-based planning with local completion and *2)* logit-based privacy-preserving inference. NEST [201] uses hybrid retrieval in the cloud and token-level neighbor search on the edge via lightweight caches.

*Edge-to-Cloud Drafting.* This paradigm starts with fast, low-cost draft generation at the edge, followed by refinement from a more capable cloud model. SlimPLM [174] triggers LLM engagement based on local knowledge sufficiency. Hao *et al.* [14] perform token-level corrections, selectively replacing ambiguous tokens. Probe Sampling [209] filters prompts based on draft-target similarity, minimizing cloud usage. While collaboration granularity varies, from document-chunk retrieval to token-level correction, these methods exemplify two complementary coordination paradigms: cloud-to-edge sketching for semantic decomposition and edge-to-cloud refinement for lightweight, responsive drafting.

Table 6. Comparison of collaborative training paradigms in edge-cloud LLM-SLM systems

| Category | Definition & Characteristics | Advantages | Limitations |
|---|---|---|---|
| Distillation-Based Collaboration | Uses logits, features, or hidden states for LLM-to-SLM supervision with task- and domain-adaptive schemes | Flexible; enables effective model compression and personalization | Sensitive to domain shift and teacher-student mismatch; data quality-dependent; unidirectional by default |
| Multi-SLM Parameter Fusion | Integrates heterogeneous SLMs via stitching, subnet assembly and capacity-guided transfer | Enhances edge personalization; enables cross-domain parameter reuse; supports one-pass deployment | Structural mismatch sensitivity; requires functional/meta-guidance; faces stability challenges |
| Adapter-Based Modular Training | Inserts LoRA-like adapters for parameter-efficient fine-tuning and federated updates | Scalable across tasks and devices; reduces communication costs; preserves base model integrity | Depends on shared adapter architectures; limited in cross-modal or non-aligned scenarios |
| SLM-Driven LLM Supervision | Employs SLM CoT to refine LLM outputs in noisy/low-resource contexts | Improves interpretability; enhances reliability on simple or repetitive tasks | Constrained by SLM capacity; risk of reinforcing local biases |
| Cloud-Guided Capability Injection | LLMs generate task-specific modules or behavioral updates for SLMs via meta-learning or distillation | Enables few-shot adaptation; supports evolving edge demands; modular and task-aware | Requires costly cloud computation; may induce concept drift; on-device validation is challenging |

*2.4.4 Token Tree Verification.* In speculative decoding, optimal verification timing is crucial: early verification may waste computation, while delayed checks risk irreversible errors (see Fig. 7). To mitigate this, the token tree structure allows each node to branch into multiple candidate paths, enabling broader output exploration and reducing waste from single-path errors. Along this line, LLMCad [26] and Traversal Verification [133] introduce a non-autoregressive token-tree verifier capable of simultaneously validating and correcting all branches within a single iteration. Building on this foundation, AdaServe [39] constructs a candidate token tree for each request and dynamically selects the optimal token branch to maximize system throughput under predefined service level objectives (SLOs) [29, 131]. OPT-Tree [60] addresses the inefficiency of fixed-structure token trees by proposing a dynamic tree construction method. It greedily builds an expectation-optimal candidate tree at each decoding step and prunes branches via probabilistic modeling, thus maximizing valid tokens under a node budget. Traversal Verification [133] further enhances tree utilization by introducing bottom-up, sequence-level verification. This method evaluates complete token paths rather than individual tokens, reducing the risk of prematurely discarding useful subsequences. Sequoia [127] builds on this by formulating token tree construction as a dynamic programming problem. It avoids draft model resampling and integrates hardware-aware optimizations to adaptively select tree depth and size for target platforms.

## 3 Overview of Collaborative learning Architectures

Collaborative training facilitates ongoing knowledge exchange between edge SLMs and cloud LLMs under heterogeneous conditions. Unlike latency-focused inference cooperation, training stresses structural compatibility, non-IID adaptation, and privacy-preserving transfer. Traditional distillation and transfer learning assume aligned
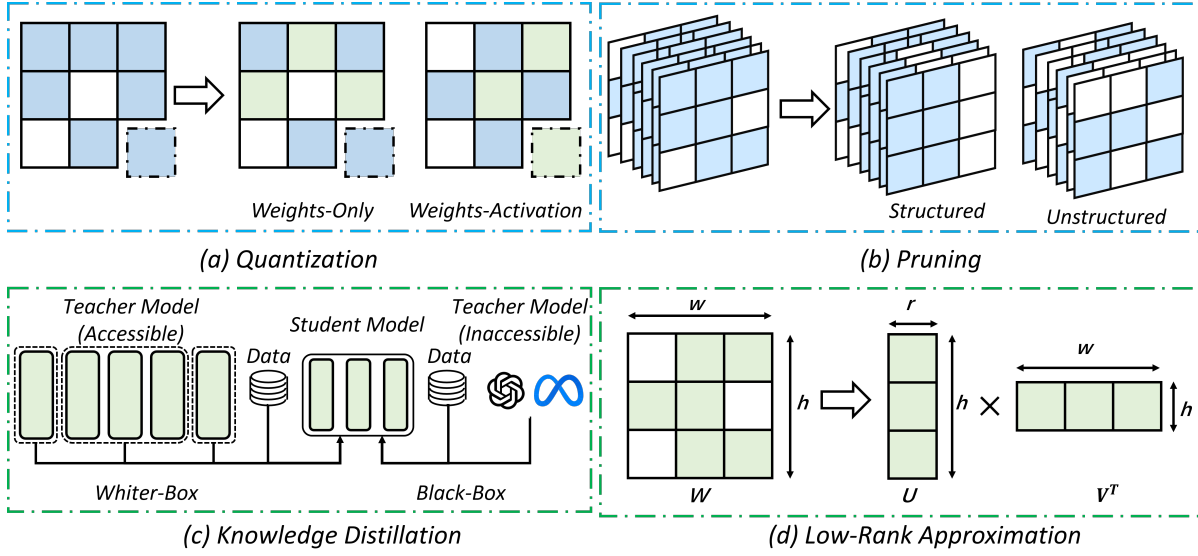
Fig. 8. Model compression techniques for efficient cloud-edge collaborative training: (a) quantization of weights and activations; (b) pruning of parameters (structured or unstructured sparsity); (c) knowledge distillation for teacher-student transfer; and (d) low-rank approximation of weight matrices.

architectures or tasks, limiting generality. In practice, structural gaps and data heterogeneity cause biases and spurious correlations that harm aggregation and generalization. Liu *et al.* [90] formalize large-small model collaboration under privacy, exposure, and cost constraints, decomposing it into three stages: *1)* LLM-to-SLM transfer, *2)* SLM-to-LLM feedback, and *3)* closed-loop adaptation. To unify this landscape, we propose a taxonomy of collaborative training across five paradigms: distillation-based collaboration, adapter-based modular training, bidirectional learning, SLM-driven supervision, and cloud-guided capability injection. As shown in Table 6, our classification highlights each paradigm's functional role and progressive interdependence, from unidirectional transfer to modular tuning, iterative feedback, and reversed supervision. Grounded in knowledge flow, coupling strength, structural assumptions, and deployment constraints, this taxonomy provides a principled basis for selecting suitable strategies in edge-cloud contexts and inspires more robust, generalizable training designs.

## 3.1 Pruning and Quantization Strategies for Efficient Device-Cloud Collaboration

In device-cloud collaborative training frameworks, pruning and quantization are key techniques for constructing lightweight models that reconcile edge devices' limited resources with cloud models' accuracy demands. Li *et al.* [120] introduce a sparsity-aware channel pruning method that evaluates feature distribution deviations to remove globally unimportant channels; its soft-mask mechanism allows selective reactivation. EfficientLLM [121] applies progressive structural pruning during pretraining to eliminate redundant substructures. Complementing pruning, Jiang *et al.* [122] propose a split-transformer design in which a lightweight edge encoder produces quantized intermediate embeddings for cloud decoding. MergeNet [67] addresses structural incompatibility via low-rank decomposition of LLM and SLM parameters, followed by attention-based fusion, enabling expressive knowledge transfer without altering the edge architecture. As illustrated in Fig. 8, these methods co-optimize pruning and quantization across training, inference, and scheduling dimensions. For task-granularity mixtures in

edge-cloud inference, the OT-GAH algorithm [118] leverages online tree pruning and generalized assignment heuristics to achieve near-optimal throughput under diverse latency and accuracy constraints.

## 3.2 Distillation and Low-Rank Approximation

Edge-cloud distillation techniques fall into three broad categories: *1)* task- and domain-adaptive distillation, *2)* architectural decoupling with proxy learning, and *3)* bidirectional, feedback-driven distillation.

*Task- and Domain-Adaptive Distillation.* ATKD redefines distillation along task- and diversity-oriented axes [112], introducing an uncertainty coefficient to quantify token-level learning difficulty and revealing that high LLM output certainty suppresses diversity. Methods such as SLMREC [91] and GKT [16] align intermediate hidden states between lightweight on-device SLMs and cloud LLMs for latent-space transfer. Domain-robust techniques [76, 81] construct a pseudo-sample space, using latent mapping, masking, and Mixup, to align student representations without shared data. DDK [101] further improves cross-domain robustness via domain-guided sampling and factor-smoothing mechanism to facilitate efficient cloud-to-edge distillation across diverse domains.

*Architectural Decoupling and Feedback-Guided Distillation.* To decouple deployment between edge and cloud, DC-CCL [9] vertically partitions a base model into cloud and edge components, training a lightweight proxy via distillation to mimic cloud behavior with minimal communication. Starodubcev *et al.* [108] adaptively invoke teacher corrections only when student outputs fail quality checks. Collin *et al.* [93] integrate weak supervision and high-capacity teachers into a closed-loop framework, while Co-Supervised Learning [102] dynamically allocates teachers and filters noisy signals via posterior consistency. Inverted supervision methods like SALT [107] and SKD [212] have edge SLMs teach LLMs early in training, gradually reversing roles to focus LLMs on complex samples. Multi-modal, feedback-enhanced pipelines, CD-CCA [48] and LLM-QAT [103], upload utility-selected samples for cloud optimization and return quantized updates to the edge, enabling efficient, compressible, and adaptive model enhancement. Collectively, these methods demonstrate a progression from static partitioning to dynamic, bidirectional, and multi-modal knowledge transfer in edge-cloud collaborative training.

*Leveraging Low-Rank Approximation for Edge-side SLMs Deployment.* . Beyond knowledge distillation, low-rank approximation has emerged as a lightweight and effective approach for optimizing LLM deployment in resource-constrained edge environments. QLLMS [294] addresses edge performance unpredictability by reconstructing the AQS matrix from partial samples via a low-rank attribute-driven recovery method. DP-LoRA [293] reduces transmission overhead in distributed training through low-rank weight updates. Moreover, [295] enhances stability-aware fine-tuning using fractional programming and an iterative-level penalty (IRP) to jointly optimize resource allocation and user–edge association.

## 3.3 Parameter Compatibility and Model Convergence

Traditional unified model-transfer methods struggle with parameter mismatches and personalized edge requirements. To address fragmented domain knowledge in collaborative learning, recent studies propose efficient parameter-coordination mechanisms: Graft [277] introduces a compatibility-aware stitching strategy, using local functional attribution and global information-theoretic signals to integrate only compatible parameters. CKI [68] evaluates a source model's information capacity and fuses compatible parameters into the target via a two-stage transfer and stitching process. Forward-OFA [69] uses real-time edge demands to assemble sub-networks through behavior-to-structure mapping, resolving gradient conflicts and allowing one-pass adaptation without back-propagation. DIET [70] maintains a unified backbone across devices, while the cloud generates personalized "diets" (subnetworks) based on each device's usage history. FedMKT [46] explores two-way knowledge exchange between cloud LLMs and client SLMs, improving adaptation under non-IID data. ModelGPT [74] leverages LLMs to automatically generate customized small models from user-provided samples or task descriptions.

*(a) Dynamic injection of knowledge capability*    *(b) Bidirectional Collaborative Learning*
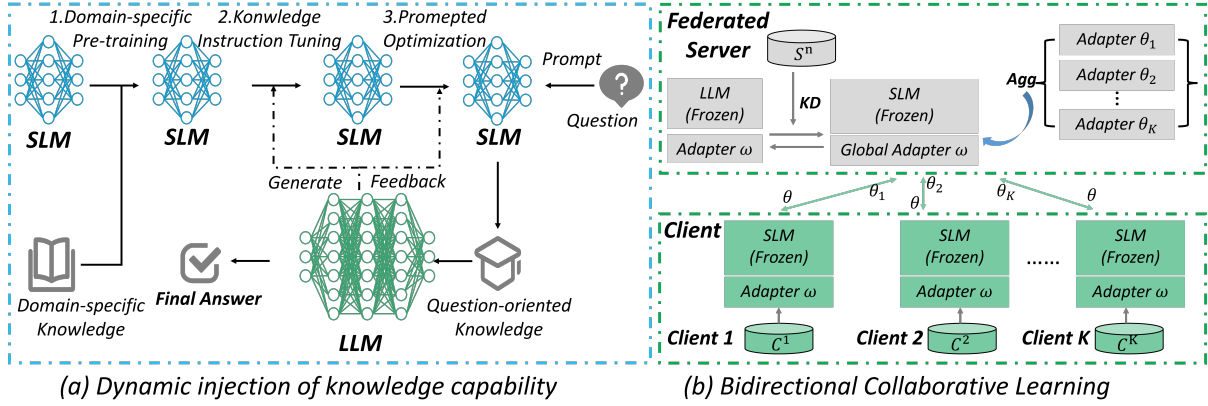
Fig. 9. Collaborative training with bidirectional knowledge transfer: (a) the SLM acquires task knowledge through domain-specific pretraining, instruction tuning, and prompt optimization; (b) frozen client SLMs augmented with local adapters train collaboratively under centralized coordination.

## 3.4 Adapter-Based Modular Designs

In collaborative edge-cloud systems, integrating knowledge and weight parameters between SLMs and LLMs remains a core challenge. To address this, adapter methods (*e.g.*, LoRA) enable scalable, parameter-efficient tuning across heterogeneous edge-cloud systems: PEFT [79] uses lightweight LoRA modules as communication bridges: clients update only adapter weights during federated training, while the server performs bidirectional knowledge extraction via supervised fine-tuning and KL-divergence regularization. Lu *et al.* [280] insert bottleneck adapters between Transformer layers to encode domain knowledge into pre-trained models. PLURALISM [88] integrates community LLMs into a foundation model via LoRA, enabling edge SLMs to fine-tune locally and the cloud LLM to manage integration. HETLoRA [96] combines high- and low-rank LoRA modules; clients rank modules by capability, and the server applies rank-aware pruning and sparsity-weighted aggregation. CDC-MMPG [36] introduces a Fast Domain Adapter that uses historical multi-modal data to train a global model and generates personalized parameters for real-time device inputs. Collectively, these approaches mark a shift toward flexible, modular, and data-aware tuning mechanisms that seamlessly integrate knowledge across end-cloud collaborative frameworks.

## 3.5 Bidirectional Collaborative Learning with Heterogeneous Models

Most distribution-shift methods rely on on-device personalization via online training or gradient updates, which incur high latency and overhead. To address this, recent approaches favors collaborative optimization across heterogeneous edge-cloud models, emphasizing iterative, bidirectional knowledge transfer over static, one-way distillation (see Fig. 9): DUET [104] separates static and dynamic layers, using a Universal Meta-Network and hierarchical HyperNetworks to generate adaptive parameters for personalized edge inference. Hu *et al.* [270] develop a knowledge-transfer prompting technique that lets different LLMs share and integrate complementary knowledge effectively. CROSSLM [87] establishes a pipeline where edge models train locally, and the cloud generates and filters pseudo-data using edge feedback, enabling mutual enhancement. SLM [89] applies inverse reinforcement learning on the edge to produce structured, task-aligned samples, while cloud LLMs perform parallel-decoding distillation to align global and local objectives.

*3.5.1 Dynamic Knowledge Transfer from Cloud LLMs to On-Device SLMs.* This research line treats small models as pluggable knowledge modules that supplement or enhance LLM performance during inference. Knowledge Card [176] uses a pretrained SLM to generate "knowledge cards", a parameterized knowledge base dynamically queried during inference. BLADE [181] employs a plugin-style architecture where domain-specific SLMs capture expert knowledge and fuse it with general-purpose LLMs via Bayesian optimization to improve reasoning quality. Progressive Distillation [99] targets unlabeled data by having cloud LLMs generate task labels and rationales to train lightweight edge models, jointly optimizing label and rationale losses for improved accuracy and interpretability. The HEF framework [180] integrates a small-scale Empathy Model to detect emotions and triggers, guiding LLM generation for more nuanced dialogue. FedCFA [73] addresses non-IID generalization in federated settings by constructing a global feature dataset and applying causality-aware counterfactual replacement, aligning SLMs with cloud semantics. AcKnowledge [281] leverages user-question-driven meta-learning to acquire external knowledge and continually refine SLMs with user feedback when distortions are detected.

*3.5.2 Adaptive Capability Injection into Small Models via Cloud Guidance.* This category reverses the typical distillation flow, using SLMs as "teachers" or behavioral priors to guide LLM training and fine-tuning. ECLM [256] decomposes cloud models into composable modules, assembling task-specific submodels for diverse devices and periodically integrating new edge knowledge back into the cloud. Mitchell *et al.* [105] encode fine-tuning-induced behavioral changes as a compact "logit delta", which a lightweight SLM computes and transfers to an LLM to emulate these adaptations without full retraining. Tang *et al.* [179, 206] present an entity-relation extraction system where a small PLM captures head-class knowledge and generates intermediate predictions serving as chain-of-thought guidance for the LLM. Purify-LLM [177] filters noisy data for LLMs using a trusted SLM and the CP-$\delta$ algorithm, integrating only aligned knowledge. These methods highlight SLMs' role in constraining and guiding LLM learning, creating a positive feedback loop that optimizes model behavior along explicit paths.

## 4 Benchmarks, Datasets, and Evaluation Protocols

Evaluating collaborative LLM-SLM systems is more complex than centralized models due to a lack of standardized benchmarks for edge-cloud collaboration. Traditional NLP benchmarks assume IID data and centralized execution, unlike real deployments where SLMs manage personalized, non-IID workloads. Thus, new benchmarks must enable user- or device-level partitioning to reflect deployment heterogeneity.

## 4.1 Datasets for Federated Edge-Cloud Collaboration

Realistic, heterogeneous, and reproducible benchmarks are crucial for evaluating collaborative training and inference in edge-cloud LLM-SLM systems (see Table 7). The LEAF framework [145] provides modular tools and real-world, user-partitioned datasets to assess performance under statistical heterogeneity, resource constraints, and personalization. For large-scale visual tasks, Zhu *et al.* [146] introduce iNaturalist-User-120k (120K images, 1.2K species, 9.3K users) and Landmarks-User-160k (164K images, 2K categories, 1.3K uploaders) with non-IID partitions. In dialogue modeling, PersonalDialog [147] offers 56 M utterances from 8.5 M speakers annotated with demographic attributes, while PERSONA-CHAT [148] provides persona-based consistency benchmarks. LiveChat [149] supplies 1.33 M multi-party Chinese dialogues with 351 roles, enabling response generation and recipient identification under temporal dynamics. FedScale [150] covers 18 real-world tasks across modalities and scales from hundreds to millions of clients, with its FAR platform supporting asynchronous training, latency simulation, and metrics for accuracy, cost, and privacy. FedNLP [151] benchmarks federated NLP tasks (*e.g.*, 20Newsgroup, MRQA) using Dirichlet splits to simulate non-IID distributions, facilitating reproducible evaluation and real-world optimization of collaborative LLM-SLM systems.

Table 7. Representative datasets for edge-cloud collaboration

| Dataset | Task | Volume & Classes | Partitioning | Metrics & Features |
|---------|------|------------------|--------------|--------------------|
| LEAF [145] | Multi-task federated learning | Sent140: 660k tweets; FEMNIST: 805k images | User-level splits reflecting real-world heterogeneity | Supports personalization and transfer learning |
| iNaturalist-User-120k [146] | Image classification | 120,300 images from 9,275 users; 1,203 species | Uploader-ID splits (user-level) | Large label space with strong user-level statistical skew |
| Landmarks-User-160k [146] | Landmark recognition | 164,172 images from 1,262 photographers; 2,028 landmarks | Region-aware S2-cell splits via GPS metadata | Captures spatial distribution heterogeneity |
| PersonalDialog [147] | Personalized dialogue | 56.3 M utterances from 8.5 M speakers | User-ID splits with demographic annotations | Supports sociolinguistic and attribute personalization |
| PERSONA-CHAT [148] | Persona-grounded chit-chat | 164,356 utterances across 10,981 dialogues | Crowdworker pairs on 1,155 personas | Next-utterance prediction and personal consistency |
| LiveChat [149] | Multi-party live dialogue | 1.33 M utterances from 351 streamers | ASR-aligned utterances with recipient matching | Single-turn dialogue with personas; addressee recognition |
| FedScale [150] | Multi-domain federated learning | 20 datasets (CV, NLP, audio) with millions of clients | Non-IID splits by user ID or trace logs | Benchmarks statistical; System heterogeneous, dynamics |
| FedNLP [151] | Federated NLP (TC, NER, QA, Seq2Seq) | 20News: 11.3k, 20 classes; OntoNotes: 50k, 37 tags | Dirichlet non-IID splits and cross-silo partitions | Metrics: accuracy, F1, ROUGE; federated evaluation |

## 4.2 Benchmarks for Dual-Model Edge-Cloud Learning

Edge-cloud collaborative learning demands dual evaluation: global generalization by the cloud LLM and local adaptation by edge SLMs, under non-IID data across devices. Recent benchmarks create non-IID test environments via user-level splits. LEAF extends CIFAR-10 and Shakespeare into federated versions with device-ID partitions [145]. Persona-Chat [148] and LiveChat [149] exploit dialogue roles for personalized generation. FedMulti-modal [152] integrates ten datasets across eight modalities, simulating missing modalities and label noise. FederatedScope-GNN [160] offers heterogeneous graph benchmarks (*e.g.*, FedDBLP) with node-, edge-, and graph-level splits. pFL-Bench [153] unifies partition strategies across 10+ datasets and introduces device heterogeneity and sparsity configurations for realism.

Hierarchical evaluation metrics are vital: image classification employs weighted global accuracy (by device data volume), local SLM accuracy, and distillation loss to assess knowledge transfer [158]. For NLG, global BLEU/ROUGE scores complement user-level perplexity. In recommendations, global AUC weighted by user data and local CTR correction reflect personalization [154], with OCPC evaluating end-to-end traffic-allocation efficiency [155].

Open-source platforms support reproducible edge-cloud experiments: FedML [159] offers on-device, single-machine, and distributed modes with modular communication-training decoupling; Flower [161] scales to millions of virtual clients across languages and backends; Google's TensorFlow Federated [158] runs on hundreds of millions of devices. For LLM-specific operations, LLMOps frameworks [156, 164] introduce tailored monitoring and safety protocols. SpecBench [244] evaluates speculative decoding across latency, accuracy, and compute cost. MessageRewriteVal [227] standardizes mobile text rewriting with high-quality, human-annotated scenarios to assess a model's ability to rewrite messages based on natural language instructions. Together, these benchmarks and platforms establish a foundation for fair, rigorous evaluation of collaborative LLM-SLM systems under realistic edge-cloud conditions.

## 5 Open Challenges

### 5.1 Privacy-Preserving and Secure Collaboration Mechanisms

Traditional optimizations for inference efficiency often overlook privacy and security risks in collaborative training and updating [113]. Achieving efficient edge-cloud cooperation while preserving data locality has thus become a critical research direction, with secure multi-party computation encrypting gradient updates to prevent sensitive information leakage during aggregation.

*Training-level protections.* Liu *et al.* [57] integrate edge-side prompt adaptation with cloud-based optimization via lightweight adapters, enabling privacy-preserving collaboration by keeping sensitive data local while leveraging cloud resources for efficiency. DCPR [63] adopts a hierarchical diffusion paradigm for personalized recommendation, where general patterns are learned in the cloud, regional preferences are adapted at the edge, and fine-grained personalization occurs on-device. Luo *et al.* [278] utilize federated learning to train models across decentralized clients, ensuring privacy by transmitting only encrypted updates. Building upon these advances, future efforts will increasingly center on uncertainty-guided supervision for selective cloud involvement, fair aggregation under skewed and non-*i.i.d.* client distributions, and generalizable knowledge transfer across heterogeneous tasks, modalities, and model scales. Furthermore, counterfactual representation learning is expected to play a pivotal role in debiasing cloud-side updates before aggregation, enhancing both utility and fairness in collaborative training.

*Inference-level protections.* RemoteRAG [247] introduces a semantic differential-privacy metric in embedding space and a dynamic secure retrieval protocol that only returns document indices when similarity exceeds a threshold. SuperICL [58] uses plugin-style small models to inject local knowledge via prompt augmentation, guiding the LLM's reasoning without exposing raw data. Pan *et al.* [35] propose a hybrid framework where the cloud builds proxy models with adapters and compression, and the edge applies symbolic masking to sensitive content, generating substitute data via cGANs to balance utility and privacy. POST [222] accelerates privacy-preserving speculative decoding by offloading draft generation to a public GPT model and optimizing cryptographic operations. Chen *et al.* [197] orchestrate a hybrid Kubernetes-based system that retains sensitive data and vector stores (*e.g.*, Chroma) in private clouds while elastically offloading LLM inference to public clouds under strict confidentiality constraints. These efforts provide initial approaches to privacy-preserving collaboration but reveal open challenges. Future work should scale secure, efficient speculative decoding across diverse model hierarchies, balance draft quality and verification cost in multi-model setups, ensure robustness amid noisy feedback and unstable connections, and enable emergent capabilities through self-reflective SLMs within larger collaborative systems.

Table 8. Industrial cloud-edge LLM-SLM collaboration frameworks and application scenarios

| Category | Domain | Collaboration Strategy | Advantages & Limitations |
|---|---|---|---|
| Industrial Edge-Cloud Frameworks | Walle [264] | End-to-end deployment pipeline handling development, runtime, and scaling. | Manages 300+ tasks and 10B+ daily calls; constrained by tight infrastructure coupling. |
| | Luoxi [265] | Slow-fast learning: cloud LLM generates aux. reps; edge SLM fast inference with feedback. | Enhances personalization via feedback loop; requires robust bidirectional communication. |
| | InfiGUIAgent [266] | Two-stage hierarchical reasoning with on-device fine-tuning for multi-modal GUI interaction. | Enables on-device GUI reasoning; may face edge memory and compute limitations. |
| Vertical Applications | Autonomous driving [55, 56] | Edge SLM handles perception; cloud LLM performs high-level planning and reasoning. | Balances fast perception with complex strategy; requires efficient cloud triggering. |
| | Livestream product recognition [184] | Edge extracts keyframe features; cloud LLM conducts multi-modal classification. | Reduces bandwidth via selective uploads; risks missing critical content in keyframes. |
| | Cultural heritage restoration [77] | Edge proposes fragment matches; cloud LLM ranks and aligns results. | Supports time-sensitive, high-accuracy restoration; depends on fragment proposal quality. |
| | Product modeling[78] | Edge injects geometric priors; cloud completes texture synthesis. | Improves partial-view generation; limited by cross-device alignment noise. |
| | Virtual assistants [240, 241] | Edge SLM manages routine interactions; cloud LLM handles complex or fallback queries. | Ensures responsiveness; cloud fallback may disrupt conversational continuity. |
| | Personalized recommendation [30, 63, 221] | Cloud LLM generates candidate items; edge re-ranks using real-time user context. | Adapts swiftly to user intent; candidate diversity hinges on LLM quality. |
| | Mobile task automation [224] | Cloud LLM parses app semantics; edge executes tasks via dynamic analysis. | Enables cross-app automation without manual setup; depends on accurate semantic parsing. |
| | Healthcare [225] | Biomedical LLM fine-tuned from general LLM for on-premise or hybrid deployment. | Offers domain-specific open-source model; needs careful tuning for sensitive settings. |
| | Web interaction [230] | Cloud LLM decomposes tasks and generates executable scripts. | Automates multi-step web tasks; relies on precise decomposition. |

## 5.2 Application of Edge-Cloud Large-Small Models in Vertical Domains

Edge-cloud LLM-SLM collaborations demonstrate strong potential across diverse application domains (see Table 8). In autonomous driving, ADAS [55] uses a lightweight multi-modal model (CogVLM2) at the edge for latency-critical tasks (object detection, obstacle avoidance) while a cloud LLM (ChatGPT-4o) handles high-level route planning. EC-Drive [56] extends this with an event-driven design: the edge LLM manages routine driving, triggering the cloud LLM upon distribution shifts (*e.g.*, novel obstacles) for commonsense reasoning and cross-modal understanding.

Live streaming challenges real-time product recognition with low latency, high concurrency, and costly multi-modal inference. Recent frameworks tackle this by using uni-modal edge models to selectively upload keyframes for cloud multi-modal processing [184]. For example, LLMCO4MR [77] treats manuscript restoration as a combinatorial task, using on-device neural solvers for fragment selection and cloud LLMs for confidence scoring and alignment. BoFiCap [286] decouples image captioning into boundary detection on-device and non-autoregressive filling in the cloud. VITA-1.5 [285] employs a three-stage training pipeline (vision-language pretraining, audio alignment, speech decoding) to avoid modality conflicts, while MPOD123 [78] injects geometric priors at the edge and offloads complex visual generation to the cloud.

Recommendation systems such as cenarios [263], LSC4Rec [30] use cloud LLMs to generate diverse candidate lists, with lightweight edge ranking models adapting in real time to user preferences. AutoDroid [224] automates mobile workflows by combining LLM commonsense reasoning with dynamic analysis of app interfaces. In healthcare, BioMistral [225] specializes a general LLM via continued pretraining on biomedical texts to support domain-specific tasks. For web-based interactions, WebAgent [230] addresses open-domain variability and the lack of HTML inductive biases by decomposing user instructions into structured subtasks, summarizing lengthy documents, and generating executable Python scripts to perform complex operations. Together, these vertical applications illustrate the transformative impact of hierarchical edge-cloud LLM-SLM architectures across real-world scenarios.

## 6 Conclusion

This survey offers the first comprehensive review of edge-cloud collaboration between large and small language models (SLMs), addressing both inference-time cooperation and training-time coordination. We propose a unified taxonomy of inference paradigms, task assignment, task division, and hybrid approaches, further refined into task- and token-level granularities. On the training side, we identify key enablers such as bidirectional knowledge transfer, quantization, pruning, and low-rank adaptation, which support efficient deployment without compromising performance. By grounding our analysis in recent advances, we bridge algorithmic techniques with system-level requirements, providing methodological insights and practical guidance for building scalable, low-latency, resource-aware LLM-SLM systems. This unified perspective establishes a solid foundation for future research and deployment in heterogeneous environments.

**Future Prospects**. In edge-cloud collaborative inference, many systems rely on uncertainty estimation to support dynamic decisions such as early exit, fallback, or model switching, typically by comparing uncertainty between edge-side SLMs and cloud LLMs. However, mainstream methods based on sampling consistency or token-level probabilities face inherent limitations. Normalized softmax probabilities obscure the raw evidential strength encoded in logits, often failing to distinguish whether a model is confident, uncertain, or unfamiliar with the input. This misalignment is particularly problematic in open-ended language generation, where multiple valid continuations may exist. We advocate a future shift toward evidence-based uncertainty estimation, which leverages unnormalized logits to preserve the model's accumulated training experience. By decomposing uncertainty into epistemic and aleatoric components, *e.g.*, via Dirichlet-based formulations, edge models can more accurately characterize their own reliability. Such evidence-aware strategies offer a finer-grained basis for trust calibration,

response escalation, and offloading in collaborative LLM-SLM deployments, enabling more robust and semantically aligned inference.

Based on evidence-driven uncertainty estimation, future collaborative frameworks will shift from fixed rules to adaptive, intelligent policies. Reinforcement and meta-learning can optimize collaboration using task needs, resource limits, and real-time feedback. Cooperation will deepen from simple task splitting to joint optimization and shared representations, enabling finer coordination. In complex fields like multimodal and embodied intelligence, lightweight edge models will handle real-time sensing while cloud LLMs perform high-level reasoning, balancing latency and decision complexity. Challenges remain in ensuring robust collaboration, reducing communication costs, and enabling personalization. Future systems must balance performance, flexibility, and privacy for truly effective edge-cloud LLM-SLM collaboration.

## References

[1] Yu Qiu, Junbin Liang, Victor CM Leung, et al. Online security-aware and reliability-guaranteed ai service chains provisioning in edge intelligence cloud. *IEEE Trans. Mob. Comput.*, 2023.

[2] Othmane Friha, Mohamed Amine Ferrag, Burak Kantarci, et al. Llm-based edge intelligence: A comprehensive survey on architectures, applications, security and trustworthiness. *IEEE Open J. Commun. Soc.*, 2024.

[3] Yiming Wang, Yu Lin, Xiaodong Zeng, et al. Privatelora for efficient privacy preserving llm. *arXiv preprint arXiv:2311.14030*, 2023.

[4] Guanqiao Qu, Qiyuan Chen, Wei Wei, et al. Mobile edge intelligence for large language models: A contemporary survey. *IEEE Commun. Surv. Tutor.*, 2025.

[5] Mulei Ma, Chenyu Gong, Liekang Zeng, et al. Multi-tier multi-node scheduling of LLM for collaborative AI computing. In *Proc. IEEE INFOCOM*, pages 1–10, 2025.

[6] Zhongkai Yu, Shengwen Liang, Tianyun Ma, et al. Cambricon-llm: A chiplet-based hybrid architecture for on-device inference of 70b LLM. In *Proc. 57th IEEE/ACM Int. Symp. Microarchitecture, MICRO*, pages 1474–1488, 2024.

[7] Chenqian Yan, Songwei Liu, Hongjian Liu, et al. Hybrid sd: Edge-cloud collaborative inference for stable diffusion models. *arXiv preprint arXiv:2408.06646*, 2024.

[8] Xuchao Zhang, Menglin Xia, Camille Couturier, et al. Hybrid retrieval-augmented generation for real-time composition assistance. *arXiv preprint arXiv:2308.04215*, 2023.

[9] Yucheng Ding, Chaoyue Niu, Fan Wu, et al. Dc-ccl: Device-cloud collaborative controlled learning for large vision models. *arXiv preprint arXiv:2303.10361*, 2023.

[10] Jingcheng Fang, Ying He, F Richard Yu, et al. Large language models (llms) inference offloading and resource allocation in cloud-edge networks: An active inference approach. In *Proc. 98th IEEE Veh. Technol. Conf., VTC2023-Fall*, pages 1–5. IEEE, 2023.

[11] Jihwan Bang, Juntae Lee, Kyuhong Shim, et al. Crayon: Customized on-device llm via instant adapter blending and edge-server hybrid inference. In *Proc. 62nd Annu. Meet. Assoc. Comput. Linguist. (ACL), Vol. 1, ACL*, pages 3720–3731, 2024.

[12] Ruslan Svirschevski, Avner May, Zhuoming Chen, et al. Specexec: Massively parallel speculative decoding for interactive llm inference on consumer devices. *arXiv preprint arXiv:2406.02532*, 2024.

[13] Hongpeng Jin and Yanzhao Wu. Ce-collm: Efficient and adaptive large language models through cloud-edge collaboration. *arXiv preprint arXiv:2411.02829*, 2024.

[14] Zixu Hao, Huiqiang Jiang, Shiqi Jiang, et al. Hybrid slm and llm for edge-cloud collaborative inference. In *Proc. Workshop Edge Mobile Found. Models, EdgeFM*, pages 36–41, 2024.

[15] Fan Yang, Zehao Wang, Haoyu Zhang, et al. Efficient deployment of large language model across cloud-device systems. In *IEEE Int. Syst.-on-Chip Conf., SOCC*, pages 1–6, 2024.

[16] Yao Yao, Zuchao Li, and Hai Zhao. Gkt: A novel guidance-based knowledge transfer framework for efficient cloud-edge collaboration llm deployment. *arXiv preprint arXiv:2405.19635*, 2024.

[17] Yifei Shen, Jiawei Shao, Xinjie Zhang, et al. Large language models empowered autonomous edge ai for connected intelligence. *IEEE Commun. Mag.*, 62(10):140–146, 2024.

[18] Lyudong Jin, Yanning Zhang, Yanhan Li, et al. Moe$^2$: Optimizing collaborative inference for edge large language models. *arXiv preprint arXiv:2501.09410*, 2025.

[19] Wei Chen, Zhiyuan Li, Zhen Guo, et al. Octo-planner: On-device language model for planner-action agents. *arXiv preprint arXiv:2406.18082*, 2024.

[20] Xunyu Zhu, Jian Li, Yong Liu, et al. A survey on model compression for large language models. *Trans. Assoc. Comput. Linguistics*, 12:1556–1577, 2024.

[21] Jindong Li, Tenglong Li, Guobin Shen, et al. Pushing up to the limit of memory bandwidth and capacity utilization for efficient llm decoding on embedded fpga. In *Des. Autom. Test Eur. Conf., DATE*, pages 1–7. IEEE, 2025.

[22] Fucheng Jia, Zewen Wu, Shiqi Jiang, et al. Scaling up on-device llms via active-weight swapping between dram and flash. *arXiv preprint arXiv:2504.08378*, 2025.

[23] Yiheng Liu, Hao He, Tianle Han, et al. Understanding llms: A comprehensive overview from training to inference. *Neurocomputing*, 620:129190, 2025.

[24] Shuhang Zhang, Qingyu Liu, Ke Chen, et al. Large models for aerial edges: An edge-cloud model evolution and communication paradigm. *IEEE J. Sel. Areas Commun.*, 43(1):21–35, 2025.

[25] Minrui Xu, Hongyang Du, Dusit Niyato, et al. Unleashing the power of edge-cloud generative ai in mobile networks: A survey of aigc services. *IEEE Commun. Surv. Tutor.*, 26(2):1127–1170, 2024.

[26] Daliang Xu, Wangsong Yin, Xin Jin, et al. Llmcad: Fast and scalable on-device large language model inference. *arXiv preprint arXiv:2309.04255*, 2023.

[27] Maximilian Holsman, Yukun Huang, and Bhuwan Dhingra. Fuzzy speculative decoding for a tunable accuracy-runtime tradeoff. *arXiv preprint arXiv:2502.20704*, 2025.

[28] Kai Lv, Honglin Guo, Qipeng Guo, and Xipeng Qiu. Duodecoding: Hardware-aware heterogeneous speculative decoding with dynamic multi-sequence drafting. *arXiv preprint arXiv:2503.00784*, 2025.

[29] Penghui Yang, Cunxiao Du, Fengzhuo Zhang, et al. Longspec: Long-context speculative decoding with efficient drafting and verification. *arXiv preprint arXiv:2502.17421*, 2025.

[30] Zheqi Lv, Tianyu Zhan, Wenjie Wang, et al. Collaboration of large language models and small recommendation models for device-cloud recommendation. In *Proc. ACM SIGKDD Conf. Knowl. Discov. Data Min., KDD*, pages 962–973, 2025.

[31] Jihwan Bang, Juntae Lee, Kyuhong Shim, et al. Crayon: Customized on-device LLM via instant adapter blending and edge-server hybrid inference. In *Proc. 62nd Annu. Meet. Assoc. Comput. Linguist. (ACL), Vol. 1, ACL*, pages 3720–3731, 2024.

[32] Sehoon Kim, Karttikeya Mangalam, Suhong Moon, et al. Speculative decoding with big little decoder. In *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 2023.

[33] Kaiyan Zhang, Jianyu Wang, Remo Hua, et al. Cogenesis: A framework collaborating large and small language models for secure context-aware instruction following. In *Proc. 62nd Annu. Meet. Assoc. Comput. Linguist. (ACL), Vol. 1, ACL*, pages 4295–4312, 2024.

[34] Kaiyan Zhang, Jianyu Wang, Ning Ding, et al. Fast and slow generating: An empirical study on large and small language models collaborative decoding. *arXiv preprint arXiv:2406.12295*, 2024.

[35] Yanghe Pan, Zhou Su, Yuntao Wang, et al. Cloud-edge collaborative large model services: Challenges and solutions. *IEEE Netw.*, 2024.

[36] Wei Ji, Li Li, Zheqi Lv, et al. Backpropagation-free multi-modal on-device model adaptation via cloud-device collaboration. *ACM Trans. Multim. Comput. Commun. Appl.*, 21(2):69:1–69:17, 2025.

[37] Baohao Liao, Yuhui Xu, Hanze Dong, et al. Reward-guided speculative decoding for efficient llm reasoning. *arXiv preprint arXiv:2501.19324*, 2025.

[38] Fangcheng Liu, Yehui Tang, Zhenhua Liu, et al. Kangaroo: Lossless self-speculative decoding for accelerating llms via double early exiting. In *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 2024.

[39] Zikun Li, Zhuofu Chen, Remi Delacourt, et al. Adaserve: Slo-customized llm serving with fine-grained speculative decoding. *arXiv preprint arXiv:2501.12162*, 2025.

[40] Yujin Kim, Euiin Yi, Minu Kim, et al. Guiding reasoning in small language models with llm assistance. *arXiv preprint arXiv:2504.09923*, 2025.

[41] Ting Sun, Penghan Wang, and Fan Lai. Disco: Device-server collaborative llm-based text streaming services. *arXiv preprint arXiv:2502.11417*, 2025.

[42] Ruslan Svirschevski, Avner May, Zhuoming Chen, et al. Specexec: Massively parallel speculative decoding for interactive llm inference on consumer devices. In *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 2024.

[43] Huiyou Zhan, Xuan Zhang, Haisheng Tan, et al. Pice: A semantic-driven progressive inference system for llm serving in cloud-edge networks. *arXiv preprint arXiv:2501.09367*, 2025.

[44] Yucheng Ding, Chaoyue Niu, Fan Wu, et al. Enhancing on-device llm inference with historical cloud-based llm interactions. In *Proc. ACM SIGKDD Conf. Knowl. Discov. Data Min., KDD*, pages 597–608, 2024.

[45] Zhi Yao, Zhiqing Tang, Jiong Lou, et al. Velo: A vector database-assisted cloud-edge collaborative llm qos optimization framework. In *IEEE Int. Conf. Web Serv., ICWS*, pages 865–876, 2024.

[46] Tao Fan, Guoqiang Ma, Yan Kang, et al. Fedmkt: Federated mutual knowledge transfer for large and small language models. In *Proc. Int. Conf. Comput. Linguist., COLING*, pages 243–255.

[47] Chuantao Li, Bruce Gu, Zhigang Zhao, et al. Federated transfer learning for on-device llms efficient fine tuning optimization. *Big Data Min. Anal.*, 8(2):430–446, 2025.

[48] Guanqun Wang, Jiaming Liu, Chenxuan Li, et al. Cloud-device collaborative learning for multimodal large language models. In *IEEE/CVF Conf. Comput. Vis. Pattern Recognit., CVPR*, pages 12646–12655, 2024.

[49] Chenyang Shao, Xinyuan Hu, Yutang Lin, et al. Division-of-thoughts: Harnessing hybrid language model synergy for efficient on-device agents. In *Proc. ACM Web Conf., WWW*, pages 1822–1833, 2025.

[50] Rongjie Yi, Liwei Guo, Shiyun Wei, et al. Edgemoe: Fast on-device inference of moe-based large language models. *arXiv preprint arXiv:2308.14352*, 2023.

[51] Chang Liu and Jun Zhao. Resource allocation for stable llm training in mobile edge computing. In *Proc. ACM Int. Symp. Mobile Ad Hoc Netw. Comput., MOBIHOC*, pages 81–90, 2024.

[52] Ying He, Jingcheng Fang, F. Richard Yu, et al. Large language models (llms) inference offloading and resource allocation in cloud-edge computing: An active inference approach. *IEEE Trans. Mob. Comput.*, 23(12):11253–11264, 2024.

[53] Xinyi Hu, Zihan Chen, Kun Guo, et al. Adaptlink: A heterogeneity-aware adaptive framework for distributed mllm inference. In *Proc. AAAI 2025 Workshop Artif. Intell. Wireless Commun. Netw., AI4WCN*, 2025.

[54] Bhaskarjit Sarmah, Dhagash Mehta, Benika Hall, et al. Hybridrag: Integrating knowledge graphs and vector retrieval augmented generation for efficient information extraction. In *Proc. 5th ACM Int. Conf. AI Finance, ICAIF*, pages 608–616, 2024.

[55] Yaqi Hu, Dongdong Ye, Jiawen Kang, et al. A cloud-edge collaborative architecture for multimodal llm-based advanced driver assistance systems in iot networks. *IEEE Internet Things J.*, 12(10):13208–13221, 2025.

[56] Jiao Chen, Suyan Dai, Fangfang Chen, et al. Edge-cloud collaborative motion planning for autonomous driving with large language models. In *Proc. IEEE ICCT, ICCT*, pages 185–190, 2024.

[57] Ya Liu, Kai Yang, Yu Zhu, et al. Grey-box prompt optimization and fine-tuning for cloud-edge LLM agents, 2024.

[58] Canwen Xu, Yichong Xu, Shuohang Wang, et al. Small models are valuable plug-ins for large language models. *arXiv preprint arXiv:2305.08848*, 2023.

[59] Chi-Heng Lin, Shikhar Tuli, James Seale Smith, et al. Slim: Speculative decoding with hypothesis reduction. In *Proc. 2024 Conf. North Am. Chap. Assoc. Comput. Linguist.: Hum. Lang. Technol. (Vol. 1: Long Papers), NAACL*, pages 1005–1017, 2024.

[60] Jikai Wang, Yi Su, Juntao Li, et al. Opt-tree: Speculative decoding with adaptive draft tree structure. *Trans. Assoc. Comput. Linguistics*, 13:188–199, 2025.

[61] Heming Xia, Yongqi Li, Jun Zhang, et al. Swift: On-the-fly self-speculative decoding for llm inference acceleration. In *Proc. 13th Int. Conf. Learn. Represent., ICLR*, 2025.

[62] Jiahao Liu, Qifan Wang, Jingang Wang, et al. Speculative decoding via early-exiting for faster llm inference with thompson sampling control mechanism. In *Findings Assoc. Comput. Linguist., ACL*, pages 3027–3043, 2024.

[63] Jing Long, Guanhua Ye, Tong Chen, et al. Diffusion-based cloud-edge-device collaborative learning for next poi recommendations. In *Proc. ACM SIGKDD Conf. Knowl. Discov. Data Min., KDD*, pages 2026–2036, 2024.

[64] Yuqing Tian, Zhaoyang Zhang, Yuzhi Yang, et al. An edge-cloud collaboration framework for generative AI service provision with synergetic big cloud model and small edge models. *IEEE Netw.*, 38(5):37–46, 2024.

[65] Teresa Peng, Liam Liu, Maya Gupta, et al. Enhanced hybrid inference techniques for scalable on-device llm personalization and cloud integration. *contexts*, 8(17):18–25, 2024.

[66] Fenglong Cai, Dong Yuan, Zhe Yang, et al. Edge-llm: A collaborative framework for large language model serving in edge computing. In *IEEE Int. Conf. Web Serv., ICWS*, pages 799–809, 2024.

[67] Kunxi Li, Tianyu Zhan, Kairui Fu, et al. Mergenet: Knowledge migration across heterogeneous models, tasks, and modalities. In *Proc. AAAI Conf. Artif. Intell., AAAI*, volume 39, pages 4824–4832, 2025.

[68] Zheqi Lv, Keming Ye, Zishu Wei, et al. Optimize incompatible parameters through compatibility-aware knowledge integration. In *Proc. AAAI Conf. Artif. Intell., AAAI*, pages 19233–19241, 2025.

[69] Kairui Fu, Zheqi Lv, Shengyu Zhang, et al. Forward once for all: Structural parameterized adaptation for efficient cloud-coordinated on-device recommendation. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD*, pages 318–329, 2025.

[70] Kairui Fu, Shengyu Zhang, Zheqi Lv, et al. Diet: Customized slimming for incompatible networks in sequential recommendation. In *Proc. ACM SIGKDD Conf. Knowl. Discov. Data Min., KDD*, pages 816–826, 2024.

[71] Minsu Kim, Pinyarash Pinyoanuntapong, Bong-Ho Kim, et al. Edge vs cloud: How do we balance cost, latency, and quality for large language models over 5g networks? In *Proc. IEEE Wireless Commun. Netw. Conf., WCNC*, pages 1–6, 2025.

[72] Jing Liu, Yao Du, Kun Yang, et al. Edge-cloud collaborative computing on distributed intelligence and model optimization: A survey. *arXiv preprint arXiv:2505.01821*, 2025.

[73] Zhonghua Jiang, Jimin Xu, Shengyu Zhang, et al. Fedcfa: Alleviating simpson's paradox in model aggregation with counterfactual federated learning. In *Proc. AAAI Conf. Artif. Intell., AAAI*, pages 17662–17670, 2025.

[74] Zihao Tang, Zheqi Lv, Shengyu Zhang, et al. Modelgpt: Unleashing llm's capabilities for tailored model generation. *arXiv preprint arXiv:2402.12408*, 2024.

[75] Zheqi Lv, Wenqiao Zhang, Zhengyu Chen, et al. Intelligent model update strategy for sequential recommendation. In *Proc. ACM Web Conf., WWW*, pages 3117–3128, 2024.

[76] Zihao Tang, Zheqi Lv, Shengyu Zhang, et al. Aug-kd: Anchor-based mixup generation for out-of-domain knowledge distillation. In *Proc. 12th Int. Conf. Learn. Represent., ICLR*, 2024.

[77] Yuqing Zhang, Hangqi Li, Shengyu Zhang, et al. Llmco4mr: Llms-aided neural combinatorial optimization for ancient manuscript restoration from fragments with case studies on dunhuang. In *Proc. Eur. Conf. Comput. Vis., ECCV*, volume 15133, pages 253–269, 2024.

[78] Jimin Xu, Tianbao Wang, Tao Jin, et al. Mpod123: One image to 3d content generation using mask-enhanced progressive outline-to-detail optimization. In *IEEE/CVF Conf. Comput. Vis. Pattern Recognit., CVPR*, pages 10682–10692, 2024.

[79] Tao Fan, Yan Kang, Guoqiang Ma, et al. Fedcollm: A parameter-efficient federated co-tuning framework for large and small language models. *arXiv preprint arXiv:2411.11707*, 2024.

[80] Marc Glocker, Peter Hönig, Matthias Hirschmanner, et al. Llm-empowered embodied agent for memory-augmented task planning in household robotics. *arXiv preprint arXiv:2504.21716*, 2025.

[81] Yuxian Gu, Li Dong, Furu Wei, et al. Minillm: Knowledge distillation of large language models. In *Proc. 12th Int. Conf. Learn. Represent., ICLR*, 2024.

[82] M. Mehdi Mojarradi, Lingyi Yang, Robert McCraith, et al. Improving in-context learning with small language model ensembles. *arXiv preprint arXiv:2410.21868*, 2024.

[83] Ran Xu, Wenqi Shi, Yuchen Zhuang, et al. Collab-rag: Boosting retrieval-augmented generation for complex question answering via white-box and black-box llm collaboration. *arXiv preprint arXiv:2504.04915*, 2025.

[84] Kaiyan Zhang, Jianyu Wang, Ning Ding, et al. Fast and slow generating: An empirical study on large and small language models collaborative decoding. *arXiv preprint arXiv:2406.12295*, 2024.

[85] Bin Chen, Yu Zhang, Hongfei Ye, et al. Knowledge-decoupled synergetic learning: An MLLM based collaborative approach to few-shot multimodal dialogue intention recognition. In *Proc. ACM Web Conf., WWW*, pages 3044–3048, 2025.

[86] Yu-Neng Chuang, Leisheng Yu, Guanchu Wang, et al. Confident or seek stronger: Exploring uncertainty-based on-device llm routing from benchmarking to generalization. *arXiv preprint arXiv:2502.04428*, 2025.

[87] Yongheng Deng, Ziqing Qiao, Ju Ren, et al. Mutual enhancement of large and small language models with cross-silo knowledge transfer. *arXiv preprint arXiv:2312.05842*, 2023.

[88] Shangbin Feng, Taylor Sorensen, Yuhan Liu, et al. Modular pluralism: Pluralistic alignment via multi-llm collaboration. In *Proc. Conf. Empir. Methods Nat. Lang. Process., EMNLP*, pages 4151–4171, 2024.

[89] Yu-Chen Lin, Sanat Sharma, Hari Manikandan, et al. Efficient multitask learning in small language models through upside-down reinforcement learning. *arXiv preprint arXiv:2502.09854*, 2025.

[90] Yang Liu, Bingjie Yan, Tianyuan Zou, et al. Towards harnessing the collaborative power of large and small models for domain tasks. *arXiv preprint arXiv:2504.17421*, 2025.

[91] Wujiang Xu, Qitian Wu, Zujie Liang, et al. Slmrec: Distilling large language models into small for sequential recommendation. In *Proc. 13th Int. Conf. Learn. Represent., ICLR*, 2025.

[92] Pranjal Aggarwal, Aman Madaan, Ankit Anand, et al. Automix: Automatically mixing language models. In *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 2024.

[93] Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, et al. Weak-to-strong generalization: Eliciting strong capabilities with weak supervision. In *Proc. Int. Conf. Mach. Learn., ICML*, 2024.

[94] Shuhao Chen, Weisen Jiang, Baijiong Lin, et al. Routerdc: Query-based router by dual contrastive learning for assembling large language models. In *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 2024.

[95] Weize Chen, Yusheng Su, Jingwei Zuo, et al. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors. In *Proc. 12th Int. Conf. Learn. Represent., ICLR*, 2024.

[96] Yae Jee Cho, Luyang Liu, Zheng Xu, et al. Heterogeneous lora for federated fine-tuning of on-device foundation models. In *Proc. Conf. Empir. Methods Nat. Lang. Process., EMNLP*, pages 12903–12913, 2024.

[97] Dujian Ding, Ankur Mallick, Chi Wang, et al. Hybrid llm: Cost-efficient and quality-aware query routing. In *Proc. 12th Int. Conf. Learn. Represent., ICLR*, 2024.

[98] Dong Chen, Yueting Zhuang, Shuo Zhang, et al. Data shunt: Collaboration of small and large models for lower costs and better performance. In *Proc. AAAI Conf. Artif. Intell., AAAI*, volume 38, pages 11249–11257, 2024.

[99] Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, et al. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. In *Findings Assoc. Comput. Linguist., ACL*, pages 8003–8017, 2023.

[100] Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast inference from transformers via speculative decoding. In *Proc. Int. Conf. Mach. Learn., ICML*, 2023.

[101] Jiaheng Liu, Chenchen Zhang, Jinyang Guo, et al. Ddk: Distilling domain knowledge for efficient large language models. In *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 2024.

[102] Yuejiang Liu and Alexandre Alahi. Co-supervised learning: Improving weak-to-strong generalization with hierarchical mixture of experts. *arXiv preprint arXiv:2402.15505*, 2024.

[103] Zechun Liu, Barlas Oguz, Changsheng Zhao, et al. Llm-qat: Data-free quantization aware training for large language models. In *Findings Assoc. Comput. Linguist., ACL*, pages 467–484, 2024.

[104] Zheqi Lv, Wenqiao Zhang, Shengyu Zhang, et al. Duet: A tuning-free device-cloud collaborative parameters generation framework for efficient device model generalization. In *Proc. ACM Web Conf., WWW*, pages 3077–3085, 2023.

[105] Eric Mitchell, Rafael Rafailov, Archit Sharma, et al. An emulator for fine-tuning large language models using small language models. In *Proc. 12th Int. Conf. Learn. Represent., ICLR*, page 16, 2024.

[106] Isaac Ong, Amjad Almahairi, Vincent Wu, et al. Routellm: Learning to route llms from preference data. In *Proc. 13th Int. Conf. Learn. Represent., ICLR*, page 16, 2025.

[107] Ankit Singh Rawat, Veeranjaneyulu Sadhanala, Afshin Rostamizadeh, et al. A little help goes a long way: Efficient llm training by leveraging small lms. *arXiv preprint arXiv:2410.18779*, 2024.

[108] Nikita Starodubcev, Dmitry Baranchuk, Artem Fedorov, et al. Your student is better than expected: Adaptive teacher-student collaboration for text-conditional diffusion models. In *IEEE/CVF Conf. Comput. Vis. Pattern Recognit., CVPR*, pages 9275–9285, 2024.

[109] Heming Xia, Tao Ge, Peiyi Wang, et al. Speculative decoding: Exploiting speculative execution for accelerating seq2seq generation. In *Proc. Conf. Empir. Methods Nat. Lang. Process., EMNLP*, pages 3909–3925, 2023.

[110] Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, et al. Sheared llama: Accelerating language model pre-training via structured pruning. In *Proc. 12th Int. Conf. Learn. Represent., ICLR*, 2024.

[111] Murong Yue, Jie Zhao, Min Zhang, et al. Large language model cascades with mixture of thought representations for cost-efficient reasoning. In *Proc. 12th Int. Conf. Learn. Represent., ICLR*, 2024.

[112] Qihuang Zhong, Liang Ding, Li Shen, et al. Revisiting knowledge distillation for autoregressive language models. In *Proc. 62nd Annu. Meet. Assoc. Comput. Linguist. (ACL), Vol. 1, ACL*, pages 10900–10913, 2024.

[113] Daihang Chen, Yonghui Liu, Mingyi Zhou, et al. Llm for mobile: An initial roadmap. *ACM Trans. Softw. Eng. Methodol.*, 34(5), 2025.

[114] Menglin Xia, Xuchao Zhang, Camille Couturier, et al. Hybrid-raca: Hybrid retrieval-augmented composition assistance for real-time text prediction. In *Proc. Conf. Empir. Methods Nat. Lang. Process., EMNLP*, pages 120–131, 2024.

[115] Zheming Yang, Yuanhao Yang, Chang Zhao, et al. Perllm: Personalized inference scheduling with edge-cloud collaboration for diverse llm services. *arXiv preprint arXiv:2405.14636*, 2024.

[116] Yang Hu, Connor Imes, Xuanang Zhao, et al. Pipeedge: Pipeline parallelism for large-scale model inference on heterogeneous edge devices. In *25th Euromicro Conf. Digit. Syst. Des., DSD*, pages 298–307, 2022.

[117] Mingjin Zhang, Xiaoming Shen, Jiannong Cao, et al. Edgeshard: Efficient llm inference via collaborative edge computing. *IEEE Internet Things J.*, 12(10):13119–13131, 2025.

[118] Xinyuan Zhang, Jiangtian Nie, Yudong Huang, et al. Beyond the cloud: Edge inference for generative large language models in wireless networks. *IEEE Trans. Wireless Commun.*, 24(1):643–658, 2025.

[119] Yaqi Hu, Dongdong Ye, Jiawen Kang, et al. A cloud-edge collaborative architecture for multimodal llm-based advanced driver assistance systems in iot networks. *IEEE Internet Things J.*, 12(10):13208–13221, 2025.

[120] Mingran Li, Xuejun Zhang, Jiasheng Guo, et al. Cloud-edge collaborative inference with network pruning. *Electronics*, 12(17):3598, 2023.

[121] Xingrun Xing, Zheng Liu, Shitao Xiao, et al. Efficientllm: Scalable pruning-aware pretraining for architecture-agnostic edge language models. *arXiv preprint arXiv:2502.06663*, 2025.

[122] Penghao Jiang, Ke Xin, Chunxi Li, et al. High-efficiency device-cloud collaborative transformer model. In *IEEE/CVF Conf. Comput. Vis. Pattern Recognit., CVPR*, pages 2204–2210, 2023.

[123] Matthias Stammler, Vladimir Sidorenko, Fabian Kreß, et al. Context-aware layer scheduling for seamless neural network inference in cloud-edge systems. In *IEEE 16th Int. Symp. Embedded Multicore/Many-core Syst.-on-Chip,MCSoC*, pages 97–104, 2023.

[124] Hao Luo, Hui Tian, Peng Zhang, et al. Cloud-edge collaborative intelligent inference based on distributed neural networks in power distribution networks. In *Int. Conf. Space-Air-Ground Comput., SAGC*, pages 129–136, 2021.

[125] Guangli Li, Lei Liu, Xueying Wang, et al. Auto-tuning neural network quantization framework for collaborative inference between the cloud and edge. In *Int. Conf. Artif. Neural Netw. Mach. Learn., ICANN*, volume 11139, pages 402–411, 2018.

[126] Weilin Zhao, Yuxiang Huang, Xu Han, et al. Ouroboros: Generating longer drafts phrase by phrase for faster speculative decoding. In *Proc. Conf. Empir. Methods Nat. Lang. Process., EMNLP*, pages 13378–13393, 2024.

[127] Zhuoming Chen, Avner May, Ruslan Svirschevski, et al. Sequoia: Scalable and robust speculative decoding. In *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 2024.

[128] Benjamin Spector and Chris Re. Accelerating llm inference with staged speculative decoding. *arXiv preprint arXiv:2308.04623*, 2023.

[129] Ning Chen, Zhipeng Cheng, Xuwei Fan, et al. Towards integrated fine-tuning and inference when generative ai meets edge intelligence. *arXiv preprint arXiv:2401.02668*, 2024.

[130] Mostafa Elhoushi, Akshat Shrivastava, Diana Liskovich, et al. Layerskip: Enabling early exit inference and self-speculative decoding. In *Proc. 62nd Annu. Meet. Assoc. Comput. Linguist. (ACL), Vol. 1, ACL*, pages 12622–12642, 2024.

[131] Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, et al. Specinfer: Accelerating large language model serving with tree-based speculative inference and verification. In *Proc. 29th ACM Int. Conf. Archit. Support Program. Lang. Oper. Syst., vol. 3, ASPLOS*, pages 932–949, 2024.

[132] Huaming Wu, Anqi Gu, and Yonghui Liang. Federated reinforcement learning-empowered task offloading for large models in vehicular edge computing. *IEEE Trans. Veh. Technol.*, 74(2):1979–1991, 2025.

[133] Yepeng Weng, Qiao Hu, Xujie Chen, et al. Traversal verification for speculative tree decoding. *arXiv preprint arXiv:2505.12398*, 2025.

[134] Zhuoming Chen, Avner May, Ruslan Svirschevski, et al. Sequoia: Scalable and robust speculative decoding. In *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 2024.

[135] Bradley McDanel, Sai Qian Zhang, Yunhai Hu, et al. Pipespec: Breaking stage dependencies in hierarchical llm decoding. *arXiv preprint arXiv:2505.01572*, 2025.

[136] Jingyu Liu and Ce Zhang. Hamburger: Accelerating llm inference via token smashing. *arXiv preprint arXiv:2505.20438*, 2025.

[137] Yongchao Zhou, Kaifeng Lyu, Ankit Singh Rawat, et al. Distillspec: Improving speculative decoding via knowledge distillation. In *Proc. 12th Int. Conf. Learn. Represent., ICLR*, 2024.

[138] Avanika Narayan, Dan Biderman, Sabri Eyuboglu, et al. Minions: Cost-efficient collaboration between on-device and cloud language models. *arXiv preprint arXiv:2502.15964*, 2025.

[139] Gemini Team, Rohan Anil, Sebastian Borgeaud, et al. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2025.

[140] Marah Abdin, Jyoti Aneja, Hany Awadalla, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.

[141] OpenAI, Aaron Hurst, Adam Lerer, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.

[142] S Anthropic. Model card addendum: Claude 3.5 haiku and upgraded claude 3.5 sonnet. *URL https://api. semanticscholar. org/CorpusID*, 273639283.

[143] Protection Regulation. General data protection regulation. *Intouch*, 25:1–5, 2018.

[144] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, et al. Palm: Scaling language modeling with pathways. *J. Mach. Learn. Res.*, 24:240:1–240:113, 2023.

[145] Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, et al. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097*, 2019.

[146] Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Federated visual classification with real-world data distribution. In *Proc. Eur. Conf. Comput. Vis., ECCV*, volume 12355, pages 76–92, 2020.

[147] Yinhe Zheng, Guanyi Chen, Minlie Huang, et al. Personalized dialogue generation with diversified traits. *arXiv preprint arXiv:1901.09672*, 2019.

[148] Saizheng Zhang, Emily Dinan, Jack Urbanek, et al. Personalizing dialogue agents: I have a dog, do you have pets too? In *Proc. 56th Annu. Meet. Assoc. Comput. Linguist., ACL*, pages 2204–2213, 2018.

[149] Jingsheng Gao, Yixin Lian, Ziyi Zhou, et al. Livechat: A large-scale personalized dialogue dataset automatically constructed from live streaming. In *Proc. 61st Annu. Meet. Assoc. Comput. Linguist. (ACL), Vol. 1, ACL*, pages 15387–15405, 2023.

[150] Fan Lai, Yinwei Dai, Sanjay Sri Vallabh Singapuram, et al. Fedscale: Benchmarking model and system performance of federated learning at scale. In *Proc. Int. Conf. Mach. Learn., ICML*, volume 162, pages 11814–11827, 2022.

[151] Bill Yuchen Lin, Chaoyang He, Zihang Ze, et al. FedNLP: Benchmarking federated learning methods for natural language processing tasks. In *Findings Assoc. Comput. Linguist., NAACL*, pages 157–175, 2022.

[152] Tiantian Feng, Digbalay Bose, Tuo Zhang, et al. Fedmultimodal: A benchmark for multimodal federated learning. In *Proc. ACM SIGKDD Conf. Knowl. Discov. Data Min., KDD*, page 4035–4045, 2023.

[153] Daoyuan Chen, Dawei Gao, Weirui Kuang, et al. pfl-bench: A comprehensive benchmark for personalized federated learning. In *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 2022.

[154] Yucheng Ding, Chaoyue Niu, Fan Wu, et al. On-device model fine-tuning with label correction in recommender systems. *arXiv preprint arXiv:2211.01163*, 2022.

[155] Han Zhu, Junqi Jin, Chang Tan, et al. Optimized cost per click in taobao display advertising. In *Proc. ACM SIGKDD Conf. Knowl. Discov. Data Min., KDD*, page 2191–2200, 2017.

[156] Eduardo Zimelewicz, Marcos Kalinowski, Daniel Méndez, et al. Ml-enabled systems model deployment and monitoring: Status quo and problems. In *16th Int. Conf. Softw. Qual., SWQD*, volume 505, pages 112–131, 2024.

[157] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, et al. Tinybert: Distilling BERT for natural language understanding. In *Proc. Conf. Empir. Methods Nat. Lang. Process., EMNLP*, pages 4163–4174, 2020.

[158] Kallista A. Bonawitz, Hubert Eichner, Wolfgang Grieskamp, et al. Towards federated learning at scale: System design. In *Proc. 2nd Conf. Mach. Learn. Syst., SysML*, 2019.

[159] Chaoyang He, Songze Li, Jinhyun So, et al. Fedml: A research library and benchmark for federated machine learning. 2020.

[160] Zhen Wang, Weirui Kuang, Yuexiang Xie, et al. Federatedscope-gnn: Towards a unified, comprehensive and efficient package for federated graph learning. In *Proc. ACM SIGKDD Conf. Knowl. Discov. Data Min., KDD*, pages 4110–4120, 2022.

[161] Daniel J. Beutel, Taner Topal, Akhil Mathur, et al. Flower: A friendly federated learning research framework. *arXiv preprint arXiv:2007.14390*, 2022.

[162] Yi Chen, JiaHao Zhao, and HaoHao Han. A survey on collaborative mechanisms between large and small language models. *arXiv preprint arXiv:2505.07460*, 2025.

[163] Chaoyue Niu, Yucheng Ding, Junhui Lu, et al. Collaborative learning of on-device small model and cloud-based large model: Advances and future directions. *arXiv preprint arXiv:2504.15300*, 2025.

[164] Saurabh Pahune and Zahid Akhtar. Transitioning from mlops to llmops: Navigating the unique challenges of large language models. *Information*, 16(2), 2025.

[165] Zewei Xin, Qinya Li, Chaoyue Niu, et al. Edge-cloud routing for text-to-image model with token-level multi-metric prediction. *arXiv preprint arXiv:2411.13787*, 2024.

[166] Keming Lu, Hongyi Yuan, Runji Lin, et al. Routing to the expert: Efficient reward-guided ensemble of large language models. In *Proc. 2024 Conf. North Am. Chap. Assoc. Comput. Linguist.: Hum. Lang. Technol. (Vol. 1: Long Papers), NAACL*, pages 1964–1974, 2024.

[167] Lingjiao Chen, Matei Zaharia, and James Zou. Frugalgpt: How to use large language models while reducing cost and improving performance. *Trans. Mach. Learn. Res.*, 2024.

[168] Yiding Wang, Kai Chen, Haisheng Tan, et al. Tabi: An efficient multi-level inference system for large language models. In *Proc. 18th Eur. Conf. Comput. Syst., EuroSys*, pages 233–248, 2023.

[169] Xinyuan Wang, Yanchi Liu, Wei Cheng, et al. MixLLM: Dynamic routing in mixed large language models. In *Proc. 2025 Conf. North Am. Chap. Assoc. Comput. Linguist.: Hum. Lang. Technol. (Vol. 1: Long Papers), NAACL*, pages 10912–10922, 2025.

[170] Yile Wang, Peng Li, Maosong Sun, et al. Self-knowledge guided retrieval augmentation for large language models. In *Proc. Conf. Empir. Methods Nat. Lang. Process., EMNLP*, pages 10303–10315, 2023.

[171] Akari Asai, Zeqiu Wu, Yizhong Wang, et al. Self-RAG: Learning to retrieve, generate, and critique through self-reflection. In *Proc. 12th Int. Conf. Learn. Represent., ICLR*, 2024.

[172] Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, et al. Corrective retrieval augmented generation. *arXiv preprint arXiv:2401.15884*, 2024.

[173] Yanming Liu, Xinyue Peng, Xuhong Zhang, et al. Ra-isf: Learning to answer and understand from retrieval augmentation via iterative self-feedback. In *Findings Assoc. Comput. Linguist., ACL*, pages 4730–4749, 2024.

[174] Jiejun Tan, Zhicheng Dou, Yutao Zhu, et al. Small models, big insights: Leveraging slim proxy models to decide when and what to retrieve for llms. In *Proc. 62nd Annu. Meet. Assoc. Comput. Linguist. (ACL), Vol. 1, ACL*, pages 4420–4436, 2024.

[175] Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, et al. Toolformer: language models can teach themselves to use tools. In *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 2023.

[176] Shangbin Feng, Weijia Shi, Yuyang Bai, et al. Knowledge card: Filling llms' knowledge gaps with plug-in specialized language models. In *Proc. 11th Int. Conf. Learn. Represent., ICLR*, 2023.

[177] Tianlin Li, Qian Liu, Tianyu Pang, et al. Purifying large language models by ensembling a small language model. *arXiv preprint arXiv:2402.14845*, 2024.

[178] Dheeraj Mekala, Alex Nguyen, and Jingbo Shang. Smaller language models are capable of selecting instruction-tuning training data for larger language models. In *Findings Assoc. Comput. Linguist., ACL*, pages 10456–10470, 2024.

[179] Xuemei Tang and Jun Wang. Small language models as effective guides for large language models in chinese relation extraction. *arXiv preprint arXiv:2402.14373*, 2024.

[180] Zhou Yang, Zhaochun Ren, Wang Yufeng, et al. Enhancing empathetic response generation by augmenting llms with small-scale empathetic models. *arXiv preprint arXiv:2402.11801*, 2024.

[181] Haitao Li, Qingyao Ai, Jia Chen, et al. Blade: Enhancing black-box large language models with small domain-specific models. In *Proc. AAAI Conf. Artif. Intell., AAAI*, volume 39, pages 24422–24430, 2025.

[182] Kun Zhao, Bohao Yang, Chen Tang, et al. SLIDE: A framework integrating small and large language models for open-domain dialogues evaluation. In *Findings Assoc. Comput. Linguist., ACL*, pages 15421–15435, 2024.

[183] Wenhao Zheng, Yixiao Chen, Weitong Zhang, et al. Citer: Collaborative inference for efficient large language model decoding with token-level routing. *arXiv preprint arXiv:2502.01976*, 2025.

[184] Yufei Zhu, Chaoyue Niu, Yikai Yan, et al. Device-unimodal cloud-multimodal collaboration for livestreaming content understanding. In *IEEE Int. Conf. Data Min., ICDM*, pages 1571–1576, 2023.

[185] Jiaxing Li, Chi Xu, Lianchen Jia, et al. Eaco-rag: Towards distributed tiered llm deployment using edge-assisted and collaborative rag with adaptive knowledge update. *arXiv preprint arXiv:2410.20299*, 2025.

[186] Masoomali Fatehkia, Ji Kim Lucas, and Sanjay Chawla. T-rag: Lessons from the llm trenches. *arXiv preprint arXiv:2402.07483*, 2024.

[187] Zheng Lin, Guanqiao Qu, Qiyuan Chen, Xianhao Chen, Zhe Chen, and Kaibin Huang. Pushing large language models to the 6g edge: Vision, challenges, and opportunities, 2023.

[188] Ruiyang Qin, Dancheng Liu, Chenhui Xu, et al. Empirical guidelines for deploying llms onto resource-constrained edge devices. *ACM Trans. Des. Autom. Electron. Syst.*, 2025.

[189] Yun Zhu, Jia-Chen Gu, Caitlin Sikora, et al. Accelerating inference of retrieval-augmented generation via sparse context selection. In *Proc. 13th Int. Conf. Learn. Represent., ICLR*, 2025.

[190] Ruiyang Qin, Zheyu Yan, Dewen Zeng, et al. Robust implementation of retrieval-augmented generation on edge-based computing-in-memory architectures. In *Proc. 43rd IEEE/ACM Int. Conf. Comput.-Aided Des., ICCAD*, 2025.

[191] Tianyu Liu, Yun Li, Qitan Lv, et al. Pearl: Parallel speculative decoding with adaptive draft length. In *Proc. 13th Int. Conf. Learn. Represent., ICLR*, 2025.

[192] Seong Hoon Seo, Junghoon Kim, Donghyun Lee, et al. Facil: Flexible DRAM address mapping for soc-pim cooperative on-device LLM inference. In *IEEE Int. Symp. High Perform. Comput. Archit., HPCA*, pages 1720–1733, 2025.

[193] Weiyi Sun, Mingyu Gao, Zhaoshi Li, et al. Lincoln: Real-time 50~100b LLM inference on consumer devices with lpddr-interfaced, compute-enabled flash memory. In *IEEE Int. Symp. High Perform. Comput. Archit., HPCA*, pages 1734–1750. IEEE, 2025.

[194] Xiangxiang Gao, Weisheng Xie, Yiwei Xiang, et al. Falcon: Faster and parallel inference of large language models through enhanced semi-autoregressive drafting and custom-designed decoding tree. In *Proc. AAAI Conf. Artif. Intell., AAAI*, pages 23933–23941, 2025.

[195] Shensian Syu and Hung-yi Lee. Hierarchical speculative decoding with dynamic window. In *Proc. 2025 Conf. North Am. Chap. Assoc. Comput. Linguist.: Hum. Lang. Technol. (Vol. 1: Long Papers), NAACL*, 2025.

[196] Minsu Kim, Alexander DeRieux, and Walid Saad. A bargaining game for personalized, energy efficient split learning over wireless networks. In *Proc. IEEE Wireless Commun. Netw. Conf., WCNC*, pages 1–6, 2023.

[197] Chia-Chuan Chuang and Kai-Ching Chen. Retrieval augmented generation on hybrid cloud: A new architecture for knowledge base systems. In *16th IIAI Int. Congr. Adv. Appl. Inform.,IIAI-AAI*, pages 68–71, 2024.

[198] Aniruddha Salve, Saba Attar, Mahesh Deshmukh, et al. A collaborative multi-agent approach to retrieval-augmented generation across diverse data. *arXiv preprint arXiv:2412.05838*, 2024.

[199] Pei Chen, Shuai Zhang, and Boran Han. CoMM: Collaborative multi-agent, multi-reasoning-path prompting for complex problem solving. In *Proc. 2024 Conf. North Am. Chap. Assoc. Comput. Linguist.: Hum. Lang. Technol. (Vol. 1: Long Papers), NAACL*, pages 1720–1738, 2024.

[200] Ziyan Jiang, Xueguang Ma, and Wenhu Chen. Longrag: Enhancing retrieval-augmented generation with long-context llms. *arXiv preprint arXiv:2406.15319*, 2024.

[201] Minghan Li, Xilun Chen, Ari Holtzman, et al. Nearest neighbor speculative decoding for llm generation and attribution. *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 37:80987–81015, 2024.

[202] Yulu Gan, Mingjie Pan, Rongyu Zhang, et al. Cloud-device collaborative adaptation to continual changing environments in the real-world. In *IEEE/CVF Conf. Comput. Vis. Pattern Recognit., CVPR*, pages 12157–12166, 2023.

[203] Yuxuan Chen, Rongpeng Li, Xiaoxue Yu, et al. Adaptive layer splitting for wireless llm inference in edge computing: A model-based reinforcement learning approach. *arXiv preprint arXiv:2406.02616*, 2024.

[204] Shengyuan Ye, Jiangsu Du, Liekang Zeng, et al. Galaxy: A resource-efficient collaborative edge ai system for in-situ transformer inference. In *Proc. IEEE INFOCOM*, pages 1001–1010, 2024.

[205] Anil Kag, Igor Fedorov, Aditya Gangrade, et al. Efficient edge inference by selective query. In *Proc. 11th Int. Conf. Learn. Represent., ICLR*, 2023.

[206] Zhanhui Zhou, Zhixuan Liu, Jie Liu, et al. Weak-to-strong search: align large language models via searching over small language models. In *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 2025.

[207] Zihua Wang, Ruibo Li, Haozhe Du, et al. Flash: Latent-aware semi-autoregressive speculative decoding for multimodal tasks. *arXiv preprint arXiv:2505.12728*, 2025.

[208] Mukul Gagrani, Raghavv Goel, Wonseok Jeon, et al. On speculative decoding for multimodal large language models. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops, CVPRW*, pages 8285–8289, 2024.

[209] Yiran Zhao, Wenyue Zheng, Tianle Cai, et al. Accelerating greedy coordinate gradient and general prompt optimization via probe sampling. In *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 2025.

[210] Zhenglin Wang, Jialong Wu, Yilong Lai, et al. Seed: Accelerating reasoning tree construction via scheduled speculative decoding. In *Proc. Int. Conf. Comput. Linguist., COLING*, pages 4920–4937, 2025.

[211] Zilong Wang, Zifeng Wang, Long Le, et al. Speculative rag: Enhancing retrieval augmented generation through drafting. In *Proc. 13th Int. Conf. Learn. Represent., ICLR*, 2025.

[212] Wenda Xu, Rujun Han, Zifeng Wang, et al. Speculative knowledge distillation: Bridging the teacher-student gap through interleaved sampling. In *Proc. 13th Int. Conf. Learn. Represent., ICLR*, 2025.

[213] Valentin De Bortoli, Alexandre Galashov, Arthur Gretton, et al. Accelerated diffusion models via speculative sampling. In *Proc. Int. Conf. Mach. Learn., ICML*, 2025.

[214] Cunxiao Du Yunlong Hou, Fengzhuo Zhang et al. Banditspec: Adaptive speculative decoding via bandit algorithms. In *Proc. Int. Conf. Mach. Learn., ICML*, 2025.

[215] Jiale Fu, Yuchu Jiang, Junkai Chen, et al. Fast large language model collaborative decoding via speculation. In *Proc. Int. Conf. Mach. Learn., ICML*, 2025.

[216] Hengyuan Hu, Aniket Das, Dorsa Sadigh, et al. Diffusion models are secretly exchangeable: Parallelizing ddpms via autospeculation. In *Proc. Int. Conf. Mach. Learn., ICML*, 2025.

[217] Ziteng Sun, Ananda Theertha Suresh, Jae Hun Ro, et al. Spectr: fast speculative decoding via optimal transport. In *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 2023.

[218] Cunxiao Du, Jing Jiang, Xu Yuanchen, et al. Glide with a cape: a low-hassle method to accelerate speculative decoding. In *Proc. Int. Conf. Mach. Learn., ICML*, 2024.

[219] Haifeng Qian, Sujan Kumar Gonugondla, Sungsoo Ha, et al. Bass: Batched attention-optimized speculative sampling. In *Findings Assoc. Comput. Linguist., ACL*, pages 8214–8224, 2024.

[220] Heming Xia, Yongqi Li, Jun Zhang, et al. Swift: On-the-fly self-speculative decoding for llm inference acceleration. In *Proc. 13th Int. Conf. Learn. Represent., ICLR*, 2025.

[221] Xinyu Lin, Chaoqun Yang, Wenjie Wang, et al. Efficient inference for large language model-based generative recommendation. In *Proc. 13th Int. Conf. Learn. Represent., ICLR*, 2025.

[222] Zhengyi Li, Yue Guan, Kang Yang, et al. An efficient private gpt never autoregressively decodes. In *Proc. Int. Conf. Mach. Learn., ICML*, 2025.

[223] Ziyao Wang, Muneeza Azmat, Ang Li, et al. Speculate, then collaborate: Fusing knowledge of language models during decoding. In *Proc. Int. Conf. Mach. Learn., ICML*, 2025.

[224] Hao Wen, Yuanchun Li, Guohong Liu, et al. Autodroid: Llm-powered task automation in android. In *Proc. 30th Annu. Int. Conf. Mobile Comput. Netw., MobiCom*, page 543–557, 2024.

[225] Yanis Labrak, Adrien Bazoge, Emmanuel Morin, et al. Biomistral: A collection of open-source pretrained large language models for medical domains. In *Findings Assoc. Comput. Linguist., ACL*, pages 5848–5864, 2024.

[226] Jiannan Xiang, Tianhua Tao, Yi Gu, et al. Language models meet world models: Embodied experiences enhance language models. *Proc. Int. Conf. Neural Inf. Process. Syst., NeurIPS*, 36:75392–75412, 2023.

[227] Yun Zhu, Yinxiao Liu, Felix Stahlberg, et al. Towards an on-device agent for text rewriting. In *Proc. 2024 Conf. North Am. Chap. Assoc. Comput. Linguist.: Hum. Lang. Technol. (Vol. 1: Long Papers), NAACL*, pages 2535–2552, 2024.

[228] Chi-Min Chan, Weize Chen, Yusheng Su, et al. Chateval: Towards better LLM-based evaluators through multi-agent debate. In *Proc. 12th Int. Conf. Learn. Represent., ICLR*, 2024.

[229] Yuanhe Tian, Fei Xia, and Yan Song. Dialogue summarization with mixture of experts based on large language models. In *Proc. 62nd Annu. Meet. Assoc. Comput. Linguist. (ACL), Vol. 1, ACL*, pages 7143–7155, 2024.

[230] Izzeddin Gur, Hiroki Furuta, Austin V Huang, et al. A real-world webagent with planning, long context understanding, and program synthesis. In *Proc. 12th Int. Conf. Learn. Represent., ICLR*, 2024.

[231] Yun Zhu, Jia-Chen Gu, Caitlin Sikora, et al. Accelerating inference of retrieval-augmented generation via sparse context selection. In *Proc. 13th Int. Conf. Learn. Represent., ICLR*, 2025.

[232] Volodymyr Kuleshov. Minillm: Large language models on consumer gpus. https://github.com/kuleshov/minillm, 2023.

[233] Yuxian Gu, Li Dong, Furu Wei, et al. Minillm: Knowledge distillation of large language models. In *Proc. 12th Int. Conf. Learn. Represent., ICLR*, 2024.

[234] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, et al. TinyBERT: Distilling BERT for natural language understanding. In *Proc. Conf. Empir. Methods Nat. Lang. Process., EMNLP*, pages 4163–4174, 2020.

[235] Qwen, An Yang, Baosong Yang, et al. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2025.

[236] Sachin Mehta, Mohammad Hossein Sekhavat, Qingqing Cao, et al. Openelm: An efficient language model family with open training and inference framework. *arXiv preprint arXiv:2404.14619*, 2024.

[237] Zhengxiao Du, Yujie Qian, Xiao Liu, et al. GLM: General language model pretraining with autoregressive blank infilling. In *Proc. 60th Annu. Meet. Assoc. Comput. Linguist. (ACL), Vol. 1, ACL*, pages 320–335, 2022.

[238] Ross Taylor, Marcin Kardas, Guillem Cucurull, et al. Galactica: A large language model for science. *arXiv preprint arXiv:2211.09085*, 2022.

[239] Gemma Team, Morgane Riviere, Shreya Pathak, et al. Gemma 2: Improving open language models at a practical size, 2024.

[240] OpenAI, Josh Achiam, Steven Adler, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2024.

[241] Team GLM. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *arXiv preprint arXiv:2406.12793*, 2024.

[242] Alibaba Damo Academy. Top ten technology trends of damo academy, 2022.

[243] Qualcomm. The future of ai is hybrid; part i: Unlocking the generative ai future with on-device and hybrid ai. https://www.qualcomm.com/content/dam/qcomm-martech/dm-assets/documents/Whitepaper-The-future-of-AI-is-hybrid-Part-1-Unlocking-the-generative-AI-future-with-on-device-and-hybrid-AI.pdf, 2024.

[244] Heming Xia, Zhe Yang, Qingxiu Dong, et al. Unlocking efficiency in large language model inference: A comprehensive survey of speculative decoding. In *Findings Assoc. Comput. Linguist., ACL*, pages 7655–7671, 2024.

[245] Qualcomm. The future of ai is hybrid; part ii: Qualcomm is uniquely positioned to scale hybrid ai. https://www.qualcomm.com/content/dam/qcomm-martech/dm-assets/documents/Whitepaper-The-future-of-AI-is-hybrid-Part-2-Qualcomm-is-uniquely-positioned-to-

scale-hybrid-AI.pdf, 2024.

[246] Wenchao Xu, Jinyu Chen, Peirong Zheng, et al. Deploying foundation model powered agent services: A survey. *IEEE Commun. Surv. Tutor.*, pages 1–1, 2025.

[247] Yihang Cheng, Lan Zhang, Junyang Wang, et al. Remoterag: A privacy-preserving llm cloud rag service. *arXiv preprint arXiv:2412.12775*, 2024.

[248] Minrui Xu, Hongyang Du, Dusit Niyato, et al. Unleashing the power of edge-cloud generative ai in mobile networks: A survey of aigc services. *IEEE Commun. Surv. Tutor.*, 26(2):1127–1170, 2024.

[249] Huixian Gu, Liqiang Zhao, Zhu Han, et al. Ai-enhanced cloud-edge-terminal collaborative network: Survey, applications, and future directions. *IEEE Commun. Surv. Tutor.*, 26(2):1322–1385, 2024.

[250] Jianyu Wang, Jianli Pan, Flavio Esposito, et al. Edge cloud offloading algorithms: Issues, methods, and perspectives. *ACM Comput. Surv.*, 52(1):2:1–2:23, 2019.

[251] Yan Zhuang, Zhenzhe Zheng, Fan Wu, et al. Litemoe: Customizing on-device llm serving via proxy submodel tuning. In *Proc. 22nd ACM Conf. Embedded Networked Sensor Syst., SenSys*, page 521–534, 2024.

[252] Haseena Rahmath P, Vishal Srivastava, Kuldeep Chaurasia, et al. Early-exit deep neural network - A comprehensive survey. *ACM Comput. Surv.*, 57(3):75:1–75:37, 2025.

[253] Jie Ou, Yueming Chen, and Prof. Tian. Lossless acceleration of large language model via adaptive n-gram parallel decoding. In *Proc. 2024 Conf. North Am. Chap. Assoc. Comput. Linguist.: Hum. Lang. Technol. (Vol. 1: Long Papers), NAACL*, pages 10–22, 2024.

[254] Biao Yi, Xavier Hu, Yurun Chen, et al. Ecoagent: An efficient edge-cloud collaborative multi-agent framework for mobile automation. *arXiv preprint arXiv:2505.05440*, 2025.

[255] Junhao Feng, Boyang Huang, Xiaodong Zhou, et al. Large-scale multi-agent learning-based cloud–edge collaborative distributed pv data compression and information aggregation for multimodal network in power systems. *IEEE Trans. Consum. Electron.*, 71(1):30–40, 2025.

[256] Yan Zhuang, Zhenzhe Zheng, Yunfeng Shao, et al. Eclm: Efficient edge-cloud collaborative learning with continuous environment adaptation. *arXiv preprint arXiv:2311.11083*, 2023.

[257] Hongjun Gao, Renjun Wang, Shuaijia He, et al. A cloud-edge collaboration solution for distribution network reconfiguration using multi-agent deep reinforcement learning. *IEEE Trans. Power Syst.*, 39(2):3867–3879, 2024.

[258] Sangmin Bae, Jongwoo Ko, Hwanjun Song, et al. Fast and robust early-exiting framework for autoregressive language models with synchronized parallel decoding. In *Proc. Conf. Empir. Methods Nat. Lang. Process., EMNLP*, pages 5910–5924, 2023.

[259] Neeraj Varshney, Agneet Chatterjee, Mihir Parmar, et al. Investigating acceleration of llama inference by enabling intermediate layer decoding via instruction tuning with 'lite'. In *Proc. 2024 Conf. North Am. Chap. Assoc. Comput. Linguist.: Hum. Lang. Technol. (Vol. 1: Long Papers), NAACL*, pages 3656–3677, 2024.

[260] Daniel Rotem, Michael Hassid, Jonathan Mamou, et al. Finding the SWEET spot: Analysis and improvement of adaptive inference in low resource settings. In *Proc. 61st Annu. Meet. Assoc. Comput. Linguist. (ACL), Vol. 1, ACL*, pages 14836–14851, 2023.

[261] Yanxi Chen, Xuchen Pan, Yaliang Li, et al. Ee-llm: Large-scale training and inference of early-exit large language models with 3d parallelism. In *Proc. Int. Conf. Mach. Learn., ICML*, 2024.

[262] Jiangchao Yao, Shengyu Zhang, Yang Yao, et al. Edge-cloud polarization and collaboration: A comprehensive survey for ai. *IEEE Trans. on Knowl. and Data Eng.*, 35(7):6866–6886, 2023.

[263] Xufeng Qian, Yue Xu, Fuyu Lv, et al. Intelligent request strategy design in recommender system. In *Proc. ACM SIGKDD Conf. Knowl. Discov. Data Min., KDD*, page 3772–3782, 2022.

[264] Chengfei Lv, Chaoyue Niu, Renjie Gu, et al. Walle: An end-to-end,general-purpose, and large-scale production system for device-cloud collaborative machine learning. In *16th USENIX Symp. Oper. Syst. Des. Implementation, OSDI*, pages 249–265, 2022.

[265] LuoXi Team. Luoxi models. https://github.com/luoxi-model/luoxi_models, 2024. Accessed: 2025-06-28.

[266] Yuhang Liu, Pengxiang Li, Zishu Wei, et al. Infiguiagent: A multimodal generalist gui agent with native reasoning and reflection. *arXiv preprint arXiv:2501.04575*, 2025.

[267] Liao Hu. Hybrid edge-ai framework for intelligent mobile applications: Leveraging large language models for on-device contextual assistance and code-aware automation. *J. Ind. Eng. Appl. Sci.*, 3(3):10–22, 2025.

[268] Jasper Dekoninck, Maximilian Baader, and Martin Vechev. A unified approach to routing and cascading for llms. *arXiv preprint arXiv:2410.10347*, 2024.

[269] Yang Li. Llm bandit: Cost-efficient llm generation via preference-conditioned dynamic routing. *arXiv preprint arXiv:2502.02743*, 2025.

[270] Jinwu Hu, Yufeng Wang, Shuhai Zhang, et al. Dynamic ensemble reasoning for llm experts. *arXiv preprint arXiv:2412.07448*, 2024.

[271] Xiaoding Lu, Zongyi Liu, Adian Liusie, et al. Blending is all you need: Cheaper, better alternative to trillion-parameters llm. *arXiv preprint arXiv:2401.02994*, 2024.

[272] Ning Li, Song Guo, Tuo Zhang, et al. The moe-empowered edge llms deployment: Architecture, challenges, and opportunities. *arXiv preprint arXiv:2502.08381*, 2025.

[273] Jiale Zhang, Jiaxiang Chen, Zhucong Li, et al. Slimrag: Retrieval without graphs via entity-aware context selection. *arXiv preprint arXiv:2506.17288*, 2025.

[274] Reza Yousefi Maragheh, Pratheek Vadla, Priyank Gupta, et al. Arag: Agentic retrieval augmented generation for personalized recommendation. *arXiv preprint arXiv:2506.21931*, 2025.

[275] Longkun Guo, Jiawei Lin, Xuanming Xu, et al. Algorithmics and complexity of cost-driven task offloading with submodular optimization in edge-cloud environments. *arXiv preprint arXiv:2411.15687*, 2024.

[276] Wei Zhong, Manasa Bharadwaj, Yixiao Wang, et al. Cross-attention speculative decoding. *arXiv preprint arXiv:2505.24544*, 2025.

[277] Yang Dai, Jianxiang An, Tianwei Lin, et al. Graft: Integrating the domain knowledge via efficient parameter synergy for mllms. *arXiv preprint arXiv:2506.23940*, 2025.

[278] Huaiying Luo and Cheng Ji. Federated learning-based data collaboration method for enhancing edge cloud ai system security using large language models. *arXiv preprint arXiv:2506.18087*, 2025.

[279] Ling Li, Lidong Zhu, and Weibang Li. Cloud–edge–end collaborative federated learning: Enhancing model accuracy and privacy in non-iid environments. *Sensors*, 24(24), 2024.

[280] Qiuhao Lu, Dejing Dou, and Thien Huu Nguyen. Parameter-efficient domain knowledge integration from multiple sources for biomedical pre-trained language models. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proc. Conf. Empir. Methods Nat. Lang. Process., EMNLP*, pages 3855–3865, 2021.

[281] Sourav Das, Sanjay Chatterji, and Imon Mukherjee. Acknowledge: Acquired knowledge representation by small language model without pre-training. In *Proc. 1st Workshop Towards Knowl. Lang. Models, KnowLLM*, pages 83–95, 2024.

[282] Chitranshu Harbola and Anupam Purwar. Knowslm: A framework for evaluation of small language models for knowledge augmentation and humanised conversations. *arXiv preprint arXiv:2504.04569*, 2025.

[283] Kun Wang, Guibin Zhang, Zhenhong Zhou, et al. A comprehensive survey in llm(-agent) full stack safety: Data, training and deployment. *arXiv preprint arXiv:2504.15585*, 2025.

[284] Jianshu She, Zhuohao Li, Zhemin Huang, et al. Hawkeye:efficient reasoning with model collaboration. *arXiv preprint arXiv:2504.00424*, 2025.

[285] Chaoyou Fu, Haojia Lin, Xiong Wang, et al. Vita-1.5: Towards gpt-4o level real-time vision and speech interaction. *arXiv preprint arXiv:2501.01957*, 2025.

[286] Zheng Ma, Changxin Wang, Bo Huang, et al. Bounding and filling: A fast and flexible framework for image captioning. In *Nat. Lang. Process. Chin. Comput., NLPCC*, pages 469–481, 2023.

[287] Evan King, Haoxiang Yu, Sahil Vartak, et al. Thoughtful things: Building human-centric smart devices with small language models. *arXiv preprint arXiv:2405.03821*, 2024.

[288] Zhijun Chen, Jingzheng Li, Pengpeng Chen, et al. Harnessing multiple large language models: A survey on llm ensemble. *arXiv preprint arXiv:2502.18036*, 2025.

[289] DeepSeek-AI, Aixin Liu, Bei Feng, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2025.

[290] Fahao Chen, Peng Li, Tom H. Luan, et al. Spin: Accelerating large language model inference with heterogeneous speculative models. In *Proc. IEEE INFOCOM*, pages 1–10, 2025.

[291] Yeshwanth Venkatesha, Souvik Kundu, and Priyadarshini Panda. Fast and cost-effective speculative edge-cloud decoding with early exits, 2025.

[292] Jingling Yuan, Yao Xiang, Yuhui Deng, et al. UPOA: A user preference based latency and energy aware intelligent offloading approach for cloud-edge systems. *IEEE Trans. Cloud Comput.*, 11(2):2188–2203, 2023.

[293] Xiao-Yang Liu, Rongyi Zhu, Daochen Zha, et al. Differentially private low-rank adaptation of large language model using federated learning. *ACM Trans. Manag. Inf. Syst.*, 16(2):1–24, 2025.

[294] Miao Hu, Qi He, and Di Wu. QLLMS: quantization-adaptive LLM scheduling for partially informed edge serving systems. In *Proc. IEEE INFOCOM*, pages 1–10. IEEE, 2025.

[295] Chang Liu and Jun Zhao. Enhancing stability and resource efficiency in llm training for edge-assisted mobile systems. *IEEE Trans. Mob. Comput.*, pages 1–18, 2025.