

# Towards Edge General Intelligence: Knowledge Distillation for Mobile Agentic AI

Yuxuan Wu, Linghan Ma, Ruichen Zhang, Yinqiu Liu, Dusit Niyato, *Fellow, IEEE*,  
 Shunpu Tang, Zehui Xiong, *Senior Member, IEEE*, Zhu Han, *Fellow, IEEE*, Zhaozhi Yang,  
 Kaibin Huang, *Fellow, IEEE*, Zhaoyang Zhang, *Senior Member, IEEE*, Kai-Kit Wong, *Fellow, IEEE*

**Abstract**—Edge General Intelligence (EGI) represents a paradigm shift in mobile edge computing, where intelligent agents operate autonomously in dynamic, resource-constrained environments. However, the deployment of advanced agentic AI models on mobile and edge devices faces significant challenges due to limited computation, energy, and storage resources. To address these constraints, this survey investigates the integration of Knowledge Distillation (KD) into EGI, positioning KD as a key enabler for efficient, communication-aware, and scalable intelligence at the wireless edge. In particular, we emphasize KD techniques specifically designed for wireless communication and mobile networking, such as channel-aware self-distillation, cross-model Channel State Information (CSI) feedback distillation, and robust modulation/classification distillation. Furthermore, we review novel architectures natively suited for KD and edge deployment, such as Mamba, RWKV (Receptance, Weight, Key, Value) and Cross-Architecture distillation, which enhance generalization capabilities. Subsequently, we examine diverse applications in which KD-driven architectures enable EGI across vision, speech, and multimodal tasks. Finally, we highlight the key challenges and future directions for KD in EGI. This survey aims to provide a comprehensive reference for researchers exploring KD-driven frameworks for mobile agentic AI in the era of EGI.

**Index Terms**—Edge General Intelligence (EGI), mobile agentic ai, Knowledge Distillation (KD), Wireless Edge Intelligence

## I. INTRODUCTION

### A. Background

The Mobile Edge Computing market is undergoing substantial expansion, with projections indicating its value will surge from approximately USD 1.65 billion in 2024 to over USD

Y. Wu, L. Ma, S. Tang, Z. Yang and Z. Zhang are with the College of Information Science and Electronic Engineering, Zhejiang University, Hangzhou, Zhejiang 310027, China(e-mail: {3230100234, 3230100339, tangshunpu, yang\_zhaohui, zhzy }@zju.edu.cn).

R. Zhang, Y. Liu, and D. Niyato are with the College of Computing and Data Science, Nanyang Technological University, Singapore (e-mail: ruichen.zhang@ntu.edu.sg, yinqiu001@e.ntu.edu.sg, dniyato@ntu.edu.sg).

Z. Xiong is with the School of Electronics, Electrical Engineering and Computer Science, Queen's University Belfast, Belfast, BT7 1NN, U.K. (e-mail: z.xiong@qub.ac.uk).

Z. Han is with the University of Houston, Houston, TX 77004, USA, and also with the Department of Computer Science and Engineering, Kyung Hee University, Seoul 446701, South Korea (e-mail: hanzhu22@gmail.com).

K. Huang is with the Department of Electrical and Electronic Engineering, The University of Hong Kong, Hong Kong SAR, China (e-mail: huangkb@eee.hku.hk).

K. K. Wong are with the Department of Electronic and Electrical Engineering, University College London, WC1E 7JE, London, United Kingdom (e-mail: kai-kit.wong@ucl.ac.uk). K. K. Wong is also affiliated with Yonsei Frontier Laboratory, Yonsei University, Seoul, 03722, Republic of Korea.

Y. Wu and L. Ma contributed equally to the work.

13.5 billion by 2032<sup>1</sup>. This growth is propelled by the increasing demand for low-latency computing and enhanced quality of experience (QoE) for a rapidly growing number of mobile and IoT devices. This evolving landscape is fostering a significant technological transformation known as “agentification,” where edge devices are being endowed with autonomous capabilities through the integration of Large language models (LLMs) and other advanced AI modules. This shift transforms passive edge nodes into proactive mobile agentic AI, systems capable of perceiving their environment, reasoning, and executing complex, multi-step tasks without direct human intervention.

The benefits of this agent-centric paradigm are manifold. mobile agentic AI [1] enables a higher degree of automation and personalization, as agents can learn user preferences and adapt to dynamic situations in real-time, directly on the device. This on-device processing significantly reduces latency and enhances data privacy by minimizing reliance on centralized cloud servers. Furthermore, these agents can proactively anticipate issues, optimize workflows, and coordinate across different systems, leading to greater operational efficiency, faster issue resolution, and improved service agility. By embedding sophisticated cognitive abilities at the network’s edge, mobile agentic AI paves the way for more intelligent, responsive, and secure applications, marking a critical step toward the realization of EGI, defined as the ability of edge devices to perform general-purpose reasoning and problem solving comparable to cloud AI under strict resource and latency constraints.[2].

The realization of EGI depends on the deployment of mobile agentic AI. LLMs provide the cognitive engine for such agents, demonstrating remarkable capabilities in planning, reasoning, and tool use. However, the immense computational, memory, and energy requirements of these LLMs are fundamentally incompatible with the resource-constrained nature of mobile and edge devices.[3] This “deployment chasm” is the primary obstacle to achieving the EGI vision.

To bridge this gap, KD has emerged as a key approach for compressing large models. KD refers to the process of training a smaller “student” model to mimic the behavior of a larger, more capable “teacher” model. This allows the student to preserve the teacher’s advanced capabilities in a compact form suitable for deployment on resource-constrained edge hardware [4].

<sup>1</sup><https://www.maximizemarketresearch.com/market-report/global-mobile-edge-computing-market/6951/>

TABLE I  
SUMMARY OF RELATED SURVEYS

scope	Ref.	Kernel	overview	EGI	KD	Agentic AI
EGI	[5]	Edge Intelligence	A comprehensive survey with a 4 pillar framework about edge intelligence	✓	✗	✗
	[2]	EGI via LLMs	A survey categorizing LLM-empowered EGI into centralized, hybrid, and decentralized systems and reviewing their implementations	✓	✗	✗
	[6]	Foundation model towards EGI	A survey introducing foundation models as the key toward General Edge Intelligence, outlining research directions to tackle challenges	✓	✗	✗
	[7]	EGI with World Models	A survey introducing how world models can empower agentic AI with proactive planning and reasoning capabilities at the edge	✓	✗	✓
Agentic AI	[8]	Agentic AI	A survey reviewing the transformative role of agentic AI in organizations, highlighting its core attributes and the strategic shift	✗	✗	✓
	[1]	Mobile Agentic AI	A survey about multimodal mobile agents, categorizing and comparing their deployment effectiveness on mobile devices	✗	✗	✓
Other techniques	[9]	Knowledge distillation	A foundational survey on knowledge distillation, detailing its core components	✗	✓	✗
Ours		EGI with KD for mobile agentic AI	A comprehensive survey introducing KD as the key for mobile agentic AI towards EGI	✓	✓	✓

### B. Comparisons with Related Surveys and Contributions

The framework of KD for advancing agentic AI and enabling its deployment in EGI has demonstrated substantial potential, thereby reducing computational overhead, and enhancing adaptability in dynamic wireless environments. This paper aims to provide a comprehensive survey on the fundamentals of KD, along with its applications in agentic AI, and to highlight the enabling technologies that open new avenues for EGI deployment. Table I presents an in-depth comparison of related surveys, emphasizing KD, agentic AI, and their implementation in EGI. These surveys have primarily focused on summarizing the evolution and optimization of KD. For example, Gou *et al.* [9] offered a foundational survey on KD, detailing its core components such as knowledge types, training schemes, and teacher-student architectures.

Additionally, several surveys have investigated the application of agentic AI. Hosseini *et al.* [1] review the transformative role of agentic AI in organizations, highlighting its core attributes and the strategic shift from human-assisted “Copilot” to “autonomous Autopilot” models. Moreover, Wu *et al.* [8] survey the landscape of multimodal mobile agents, categorizing them into prompt-based and training-based methods and comparing their deployment effectiveness on mobile devices.

Furthermore, several studies have demonstrated the growing significance of agentic ai in shaping the evolution of EGI. Xu *et al.* [5] provided a comprehensive survey on edge intelligence. As the development of powerful LLMs, the exploration of EGI have surged in recent times. Chen *et al.* [2] provide a foundational overview by categorizing LLM-empowered EGI into centralized, hybrid, and decentralized systems and reviewing their implementations. He *et al.* [6] proposed that the integration of foundation models is the key to evolving toward EGI, outlining research directions to tackle challenges. Focusing on the cognitive core of such systems, Zhao *et al.* [7] offer a comprehensive analysis of how world models can empower agentic AI with proactive planning and reasoning capabilities at the edge.

However, existing studies lack sufficient investigation into

systematic KD methodologies tailored to the constraints of EGI, as well as comprehensive analysis of resource limitations and model adaptation challenges faced by agentic AI in edge environments. This paper addresses these gaps by providing an in-depth examination of how advanced model adaptation techniques, such as wireless distillation [10]–[13], and some architectures designed for the edge [14]–[17]. Furthermore, we demonstrate how KD can be effectively adapted to diverse EGI deployment scenarios, such as autonomous vehicles [18], [19], unmanned aerial vehicle (UAV) [20], [21], robotics [22], [23], and other Internet of Things [24]–[26]. The main contributions are summarized as follows.

- We provide a dedicated review of KD techniques in wireless communication scenarios, highlighting how KD enhances channel estimation, feedback compression, and resource-efficient model deployment at the wireless edge.
- We systematically review existing KD techniques and their advantages, highlighting their potential integration with novel architectures. Afterward, we comprehensively discuss the benefits and opportunities of applying KD alongside these architectures in EGI.
- We analyze the limitations present in existing models. Based on these shortcomings, we introduce architectures beyond Transformer and investigate tuning techniques combined with knowledge distillation to create conditions suitable for deployment on edge devices.
- We further analyze the challenges currently faced by EGI from both technical and ethical dimensions, proposing potential future development trends and solutions.

### C. The Structure of Paper

The structure of this survey is outlined in Fig. 1 Section II introduces the fundamental concepts of EGI, mobile agentic AI, and KD, and further elucidates the interrelationships among these three technologies. Subsequently, Section III systematically elaborates on novel architectures and wireless distillation that integrate KD into EGI, and presents the current advancements achieved by existing models. Section IV delineate the role

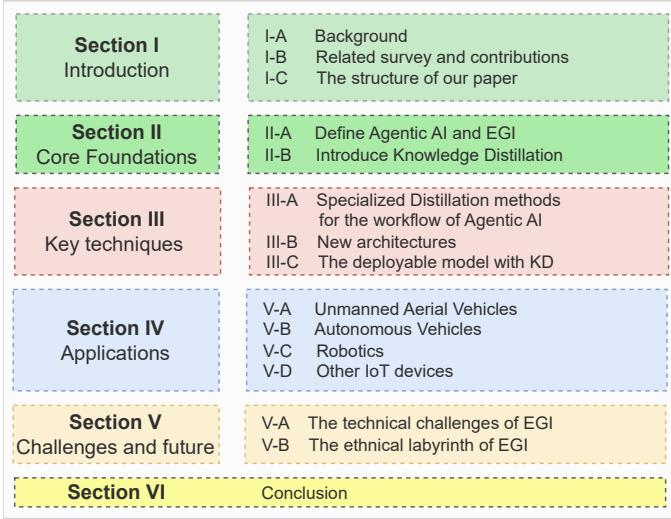


Fig. 1. The structure of this survey

of KD in EGI and its application within specialized domains. Finally, Section V elucidates the lessons learned from the review and challenges faced in this domain from both technical and macro-level perspectives, while outlining potential future research directions.

## II. CORE METHODOLOGICAL FOUNDATIONS

Agentic AI provides an autonomous perception–planning–action–memory loop, while EGI extends this paradigm to the network edge by enabling generalized cognitive capabilities under strict resource constraints. Together, they outline a unified vision for intelligent, adaptive, and autonomous edge systems.

### A. The Rise of Agentic AI and the Vision of Edge General Intelligence

1) *Mobile Agentic AI*: Agentic AI refers to autonomous systems that perceive, reason, plan, and act to achieve goals with minimal human supervision [8]. These systems operate in a continuous cycle, often termed the agentic loop, which is comprised of four core components: Perception, Planning, Action, and Memory [8].

- **Perception:** The perception module integrates multimodal information to form a coherent understanding of the environment. As depicted in Fig. 2, this process begins with data acquisition from physical sensory streams or digital channels. Subsequently, these raw data are processed to extract key features and identify relevant entities and their attributes.

- **Planning:** The planning module devises a sequence of actions to fulfill a high-level goal, frequently leveraging LLMs as a cognitive core. As Fig. 2 shows, agentic planning shifts from static, predefined action sequences to dynamic plans that are continuously updated based on environmental feedback, ensuring robust adaptability.

- **Action:** The agent executes the chosen action, interacting with and affecting its environment. Actions can range from mimicking human interactions with graphical user interfaces (GUIs), such as clicking or typing [27] in Fig. 2, to leveraging API calls for deeper system integration, such as modifying device settings or automating app navigation. In multi-agent systems, communication itself becomes a pivotal action, enabling agents to coordinate and delegate tasks [27].

- **Memory:** The memory module allows agents to retain information over time, providing essential context for coherent interactions and facilitating continual learning. As shown in Fig. 2 This is typically achieved in two tiers: Short-term memory preserves the context of the current session, often managed within the LLM’s context window. Long-term memory stores knowledge across sessions, such as learned facts and past experiences, commonly implemented using external vector databases. Agents can query this knowledge base via similarity search, a core mechanism of Retrieval-Augmented Generation (RAG).

Moreover, Coordination and adaptation are critical capabilities that enable Agentic AI to manage complex, long-horizon tasks in dynamic environments [28], [29]. Coordination is achieved through dynamic multi-agent collaboration, where autonomous agents interact to reach shared objectives [29]. Complementing this, adaptation is the capacity of an agent to learn from its environment and modify its behavior in real-time to pursue high-level goals [28]. This ability is fundamentally underpinned by integrated memory systems and learning mechanisms such as reinforcement learning [30] and lifelong learning [31]. The synergy between robust coordination protocols and memory-driven adaptation distinguishes modern Agentic AI, paving the way for more resilient and autonomous systems.

When deployed in wireless and edge computing environments, each component of the agentic loop must operate under constraints of limited computation, storage, and communication bandwidth. This challenge has catalyzed a shift in communication protocol design, moving from a focus on raw throughput to “semantic efficiency” and prioritizing the transmission of concise, high-value information over verbose raw data [32]. Emerging frameworks for 6G envision “content-aware” networks that can intelligently prioritize semantically critical data, blurring the lines between the application and network layers and transforming the network from a passive conduit into an active participant in information exchange [32].

2) *Edge General Intelligence*: EGI [2] represents a transformative paradigm designed to endow edge devices with general-purpose cognitive capabilities, enabling them to perceive, reason, and act autonomously in dynamic environments. The ultimate vision of EGI is to achieve artificial general intelligence (AGI) at the network edge, where systems can attain high-level cognitive abilities such as comprehension, reasoning, planning, and learning from experience at or even beyond human-level proficiency.

Unlike traditional edge intelligence, which primarily deploys static, task-specific models for individual tasks, EGI emphasizes versatility, adaptability, and autonomous cognitive reasoning [5]. To achieve this, EGI leverages foundation models, enabling

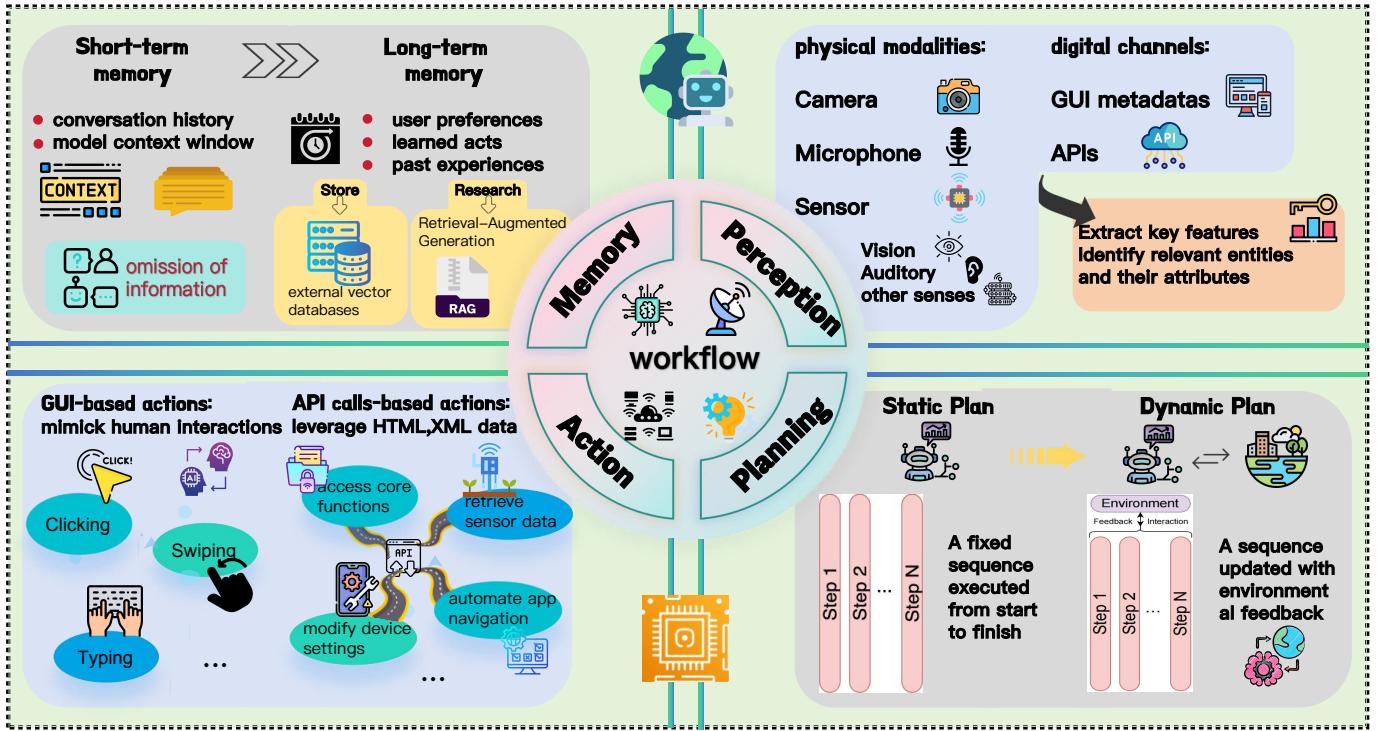


Fig. 2. An overview of the workflow of Agentic AI. *Perception* gathers and interprets multi-model information. *Planning* devises a sequence of actions to achieve a high-level goal. *Action* interacts with and affects its environment. *Memory* enables an agent to retain information over time.

devices to perform multiple diverse tasks without frequent retraining and to dynamically adapt to varying contexts and environments in real-time.

This intelligence is fundamentally rooted in the knowledge, a concept that encompasses not just factual information but also complex reasoning patterns, contextual understanding, and decision-making abilities learned from vast datasets, embedded within AI models. EGI architectures exist on a spectrum defined by how this knowledge is distributed:

- In a centralized model, a single, powerful, cloud-hosted LLM possesses all the comprehensive knowledge. It handles all complex reasoning and planning, while edge devices, possessing minimal local knowledge, serve as simple data collectors and command executors.
- Conversely, a fully decentralized approach aims to imbue each edge device with its own substantial knowledge base by equipping it with a capable medium language model or small language model (SLM), enabling on-device inference and peer-to-peer collaboration.[33]
- A hybrid architecture balances these extremes, using local SLMs with specialized or essential knowledge for routine and latency-sensitive operations, while offloading tasks requiring broader or more profound knowledge to the cloud LLM.

The potential applications for EGI are vast and transformative, promising to redefine human-machine interaction across numerous sectors, including autonomous vehicles, industrial automation, smart cities, and personalized healthcare.

## B. Knowledge Distillation: Core Techniques and Paradigms

KD provides an effective framework for transferring knowledge from large models to compact ones, enabling efficient deployment on resource-constrained devices. Modern KD methods based techniques—capture increasingly richer forms of teacher knowledge, forming the foundation of high-performance edge and agentic AI systems.

**1) Basic Concepts and Working Principles:** Hinton *et al.* [4] firstly proposed the teacher-student architecture, constituting the foundational paradigm of KD. This architecture typically integrates two components: (i) a high-capacity teacher model, frequently implemented as a state-of-the-art neural network or model ensemble, that generates rich knowledge representations and (ii) a compact student model optimized for efficient deployment. Through a specialized distillation process, the student model learns to replicate the teacher's functional mapping by approximating its knowledge outputs.

The core of this framework is the knowledge transfer mechanism [4]. The teacher model provides not just a single correct answer but a full probability distribution across all possible classes. These nuanced probabilities, often called “soft targets” or “dark knowledge,” [34] encode rich information about how the teacher generalizes and perceives similarities between classes, providing far richer supervision than the “hard,” one-hot labels used in standard training. To produce these effective soft targets, a temperature scaling parameter  $T$  [35] is applied to the teacher’s output layer, so that the probability

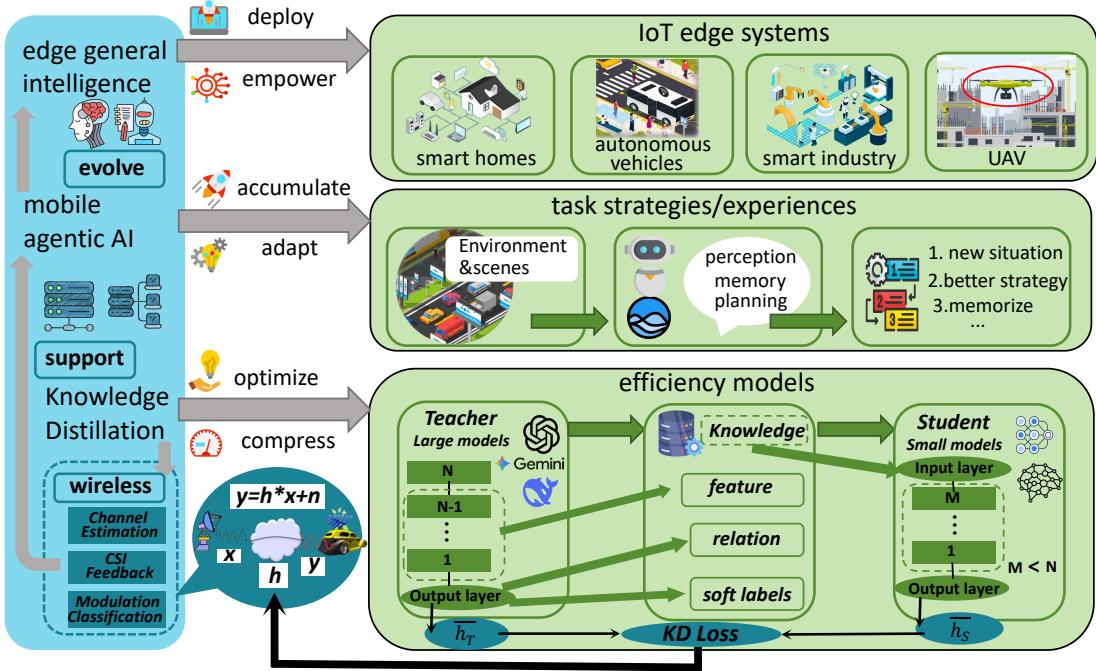


Fig. 3. Knowledge Distillation compresses large models into lightweight ones for deployment, enabling Mobile Agentic AI to accumulate experiences and adapt strategies. These capabilities collectively foster Edge General Intelligence, which empowers IoT edge systems such as UAVs, autonomous vehicles, and smart healthcare.

assigned to class  $i$  is given by

$$q_i = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}, \quad (1)$$

where  $q_i$  denotes the probability assigned to class  $i$ ,  $z_i$  represents its corresponding logit output, and  $T$  represents the temperature.

The overall training objective of the student model consists of two components. (i) loss function with soft targets: The student model employs the same high temperature  $T$  as the teacher in its softmax layer to compute the cross-entropy loss between its output distribution and the soft targets provided by the teacher. (ii) loss function with hard targets: In parallel, the student is also trained to predict the ground-truth labels from the original training data, i.e., the hard targets. For this component, the softmax temperature is set to 1 to reflect standard classification behavior. The combined loss can be written as

$$\mathcal{L} = \alpha \cdot \mathcal{L}_{\text{soft}} + (1 - \alpha) \cdot \mathcal{L}_{\text{hard}}, \quad (2)$$

where

$$\mathcal{L}_{\text{soft}} = L \left( \sigma \left( \frac{\mathbf{z}^S}{T} \right), \sigma \left( \frac{\mathbf{z}^T}{T} \right) \right), \quad (3)$$

and

$$\mathcal{L}_{\text{hard}} = L \left( \sigma(\mathbf{z}^S), \mathbf{y} \right). \quad (4)$$

Here,  $\mathbf{z}^S$  and  $\mathbf{z}^T$  denote the logits of the student and teacher models, respectively, and  $\mathbf{y}$  is the one-hot encoded ground-truth label. Let  $\mathbf{p}^S = \sigma(\mathbf{z}^S/T)$  and  $\mathbf{p}^T = \sigma(\mathbf{z}^T/T)$  represent the softened probability distributions, where  $\sigma(\cdot)$  denotes the Softmax function and  $T$  is the temperature parameter.

The distillation objective  $L$  is defined as the Kullback–Leibler (KL) divergence, which measures the discrepancy between the

teacher's and student's distributions:

$$L_{\text{KL}}(\mathbf{p}^T \| \mathbf{p}^S) = \sum_{i=1}^C p_i^T \log \left( \frac{p_i^T}{p_i^S} \right). \quad (5)$$

Alternatively, when adopting the Mean Squared Error (MSE), the loss is formulated as:

$$L_{\text{MSE}}(\mathbf{p}^T, \mathbf{p}^S) = \frac{1}{C} \|\mathbf{p}^T - \mathbf{p}^S\|_2^2. \quad (6)$$

Finally,  $\alpha \in [0, 1]$  is a hyperparameter that balances the weight between the ground-truth cross-entropy loss and the distillation loss term  $L$ .

2) *A taxonomy of modern KD techniques:* KD methods based on knowledge sources can be broadly classified into three types: response-based distillation, feature-based distillation and relation-based distillation.

a) **Response-Based Distillation:** Response-based KD transfers knowledge by using the teacher's final outputs to supervise the student model. The earliest teacher–student formulation of KD was inherently response-based [4].

Building on the KD framework, several variants have been developed to address its different limitations, forming a progressive line of research. Together, these works illustrate a spectrum of strategies from bridging model capacity gaps, to replacing or removing teachers, to designing more robust and adaptive loss functions. Mirzadeh *et al.* [38] extended the teacher–student paradigm with an assistant model to mitigate the capacity gap between a large teacher and a small student. Moreover, their framework outperformed other strong baselines—including FitNets, Attention Transfer (AT) and Mutual Learning. This yielded more effective transfer under constrained resources, though at the cost of extra training stages. Building on the

TABLE II  
A COMPARATIVE ANALYSIS OF THREE KNOWLEDGE DISTILLATION PARADIGMS

Category	Representative Method	ref	Distillation Target	Layer Level	Pros	Cons
Response-based	teacher-student	[1]	softened logits (class probability)	output layer	<ul style="list-style-type: none"> <li>• Conceptually simple</li> <li>• Easy to implement</li> </ul>	<ul style="list-style-type: none"> <li>• Restricted to output-layer information</li> <li>• Neglect intermediate feature representations</li> <li>• Limited capacity to transfer structural or semantic knowledge</li> </ul>
Feature-based	FitNets	[36]	intermediate feature maps	hidden layer	<ul style="list-style-type: none"> <li>• Leverages rich intermediate representations</li> <li>• Enables the transfer of structural and semantic information</li> </ul>	<ul style="list-style-type: none"> <li>• Challenging alignment between heterogeneous architectures</li> <li>• High computational and memory costs</li> </ul>
Relation-based	Relational Knowledge Distillation(RKD)	[37]	pairwise/triplet relational distance	across multiple layers	<ul style="list-style-type: none"> <li>• Captures higher-order dependencies among samples</li> <li>• Enhances generalization</li> </ul>	<ul style="list-style-type: none"> <li>• Computationally intensive and complex</li> <li>• Sensitive to data distribution or noisy relations</li> </ul>

idea of removing external teachers, Pham *et al.* [39] proposed self-distillation, where a model distills knowledge from its own intermediate predictions. This design is attractive for continual on-device learning, though its effectiveness hinges on prediction quality. Empirical results show that self-distillation consistently improves over teacher baselines, with student models achieving up to 1–2.5 percentage points higher accuracy on CIFAR-10/100 under comparable settings. Taken together, these works trace a clear trajectory, from bridging teacher–student gaps, to eliminating the need for teachers altogether, to strengthening robustness under noisy conditions. The collective insight suggests future research will likely integrate these dimensions, seeking hybrid methods that are resource-efficient, adaptive, and robust, while addressing scalability and parameter sensitivity. When viewed through the lens of agentic AI and mobile deployment, these directions highlight KD’s potential at the edge.

b) **Feature-Based Distillation:** Feature-based knowledge distillation leverages intermediate feature representations of the teacher network to guide student training [40]. The basic paradigm typically consists of two stages [36], [40]–[42]:

The first stage is Intermediate Feature Alignment. In this stage, a teacher layer is selected as the hint layer, and a corresponding student layer as the guided layer. A regressor is introduced to project the features from the student’s guided layer into the representation space of the teacher’s hint layer. The loss for this stage is defined as

$$L_{HT} = L(u_h(x) - r(v_g(x))), \quad (7)$$

where  $L_{HT}$  is the loss function,  $L$  is a custom metric,  $u_h(x)$  is the output of the teacher’s hint layer,  $v_g(x)$  is the output of the student’s guided layer, and  $r$  is the regressor function.

The second stage is KD Phase. In this stage, the aligned student is further trained with ground-truth labels and the teacher’s soft labels, while parameters learned in the first stage remain fixed.

Several extensions have further enriched the KD paradigm by deepening teacher–student interactions. Gou et al. [43] proposed Reciprocal Teacher–Student Learning, where knowledge transfer is bidirectional. The student not only learns from the teacher but also provides feedback for adaptive teaching. This design increases adaptability but requires careful coordination of the mutual updates. Chen *et al.* [44] advanced supervision quality through Knowledge Review, which leverages shallower or multiple teacher layers to supply more diverse guidance signals, thereby alleviating the risk of overfitting to a sin-

gle teacher representation. Complementing these supervision strategies, Guo *et al.* [45] developed attention-based distillation, encouraging the student to mimic both feature activations and teacher attention maps, thus capturing structural dependencies beyond raw outputs. The progression of these works highlights a growing trend of moving beyond output-level mimicry toward exploiting deeper representational structures for more effective knowledge transfer.

c) **Relation-Based Distillation:** Relation-based distillation transfers pairwise similarity information that captures the structural relationships among samples [46]. Its core design relies on a relation potential function, typically defined using distance-based potentials or angle-based measures [47], [48].

Several variants further refine the KD paradigm by focusing on structural information and feature relationships. Tian *et al.* [49] introduced a contrastive learning framework, where student and teacher embeddings are aligned via positive and negative sample pairs. This shifts distillation toward modeling structural relationships rather than solely mimicking outputs. Building on the idea of preserving structure, Liu *et al.* [50] developed Inter-Channel Correlation Distillation, which enforces channel-wise correlation consistency. This balances feature diversity and stability, but may introduce computational overhead in high-dimensional channels. Extending structural preservation to the relational domain, Xin *et al.* [51] proposed Neighborhood Relation KD (NRKD), which distills local neighborhood relationships instead of relying on global similarity, thereby improving robustness to noisy or imbalanced data. This progression reflects a broader shift in KD research—from simple output matching toward more nuanced modeling of representation structures at multiple levels of granularity.

### C. Knowledge Distillation as the Bridge to Mobile Agentic AI and Edge General Intelligence

Agentic AI is central to EGI, enabling edge devices to act as proactive, goal-driven systems rather than passive executors. However, agentic capabilities are inherently procedural (“how to think”), not merely static knowledge (“what to say”), creating a mismatch with conventional model compression methods that focus on reproducing output probabilities. Such techniques imitate final predictions but fail to capture the underlying reasoning process required for agentic behavior.

KD emerges as a key technology for reconciling the computational demands of large-scale AI with the limited resources of edge devices. As illustrated in Fig. 3, it acts as a conceptual

bridge enabling both mobile agentic AI and EGI by making advanced intelligence feasible on edge platforms. Through distilling massive, multi-billion-parameter “teacher” models into compact “student” models, KD substantially reduces model size, memory usage, and computation cost. This compression also improves operational efficiency: lightweight models deliver lower inference latency, essential for real-time EGI tasks such as autonomous navigation and industrial robotics, and significantly reduce energy consumption, extending the battery life of edge devices. Moreover, on-device inference enhances privacy and reliability by keeping sensitive data local and reducing dependence on unstable network connections.

In the context of Agentic AI, it is crucial to move beyond the traditional view of KD as mere compression and redefine it as a framework for capability transfer. The ultimate objective is not simply to reduce model size but to comprehensively transfer the complex cognitive skills essential for agentic behavior [1]. Key challenges arise in this process, including retaining generalization under resource-constrained mobile environments[52], preserving robustness against noisy or dynamic data, and reducing sensitivity to hyperparameters. Emerging solutions involve leveraging KD in wireless scenarios[11] to enable efficient over-the-air learning, designing lightweight and adaptive frameworks tailored for edge deployment, and developing tuning techniques that align well with KD. These lessons highlight KD’s evolving role from compression to capability transfer, paving the way toward resource-efficient and resilient agentic intelligence on mobile platforms.[53]

### III. DISTILLING AGENTIC CAPABILITIES FOR MOBILE DEPLOYMENT

KD establishes a layered relationship with EGI. At the capability level [54], KD equips EGI with fundamental capabilities, such as perception, planning, action and memory, that are required for robust operation and completion. At the methodological level [55], a range of frameworks and techniques aligned with KD enhance computational efficiency and reduce resource consumption [56], while simultaneously fostering generalizable representations that better align with diverse downstream tasks in EGI. At the toolkit level, KD itself offers reliable compact models that directly support the deployment of EGI. [57] This section proceeds to elaborate on how KD can be leveraged for EGI deployment across the three levels.

#### A. Specialized Distillation Frameworks for Capabilities of Agentic AI

To address these limitations, recent research has focused on developing specialized distillation frameworks tailored for agentic systems. A key principle unifying these advanced methods is a shift from imitating raw outputs toward transferring more abstract, structured, and reusable knowledge components.

1) *Capability in “Perception”*: KD is crucial for establishing the perception capabilities of Agentic AI. An agent’s awareness of the wireless environment depends on precise CSI estimation, accurate signal classification and so on. To overcome the

difficulty of large model for direct deployment, KD enables efficient, accurate, and robust environmental sensing on resource-constrained edge devices.[68], [69]

- **Channel Estimation:** Knowledge-driven deep learning models, such as the network developed by Yang *et al.* [70] that integrates the knowledge of prior communication systems into its architecture, are increasingly being used for complex channel estimation. To address the practical deployment challenges, KD has become a key enabling technique. Catak *et al.* [58] employed a defensive distillation framework to improve a model’s resilience to malicious signal perturbations, smoothing the decision surface and drastically reduces the success rate of attacks. This method is outstanding with its security enhancement by lowering the Attack Success Ratio (ASR) from approximately 0.9 down to 0.06. Kong *et al.* [59] used KD to solve a fundamental training problem where a non differentiable binarization layer for channel quantization provides inaccurate gradients and hinders learning. An auxiliary “teacher” network was introduced to bypass this layer and provide “lossless gradients”, guiding the main “student” network to a better solution. This technique gained a novel solution to a core optimization challenge, achieving a sum rate of approximately 9.0 at 30 dB Signal-to-Noise Ratio (SNR), compared to 8.2 without KD.
- **CSI Feedback:** To deploy complex Channel State Information (CSI) feedback models on resource-limited devices, Tang *et al.* [60] used KD to create lightweight models without difficult manual redesign. The method improved the student’s performance by up to 7.6 dB in Normalized Mean-Squared Error (NMSE) while reducing the encoder’s computational complexity to 25.50%-43.46% of the teacher’s. However, a limitation is that the student’s performance cannot surpass the teacher’s. Cui *et al.* [11] advanced this by proposing an “Encoder KD” framework to reduce the high training overhead of full autoencoder distillation, with the lightweight student encoder trained via KD before being combined with the powerful teacher decoder. This reduced total training time by 66.8% and further improved feedback performance, with NMSE dropping from -9.54 dB to -9.75 dB. Beyond efficiency, Gong *et al.* [61] used KD to enhance robustness against imperfect CSI in MIMO semantic communication to solve a problem where the ambiguous relationship between estimated and true CSI hinders model convergence. A teacher with perfect CSI guided a student trained on imperfect CSI, helping it learn a resilient representation. This KD stage alone contributed to an overall improvement of 0.40–0.54 dB Peak Signal-to-Noise Ratio (PSNR) over benchmarks including DeepJSCC [71], DeepJSCC-V [72], DeepJSCC-MIMO [73], and SwinJSCC [74], which are pre-trained with accurate CSI and then finetuned with estimated channel matrix.
- **Modulation Classification:** To deploy complex Automatic Modulation Classification (AMC) models on edge devices, Yang *et al.* [62] utilized KD to improve the classification accuracy of the computationally simple student model

TABLE III  
KD FOR WIRELESS APPLICATIONS

Application Task	Role of KD	Improvements	Implications for EGI
Channel Estimation	Teacher networks provide “lossless gradients” and Defensive distillation enhances robustness	<ul style="list-style-type: none"> <li>• Higher estimation accuracy</li> <li>• Improved security (lower ASR from 0.9 to 0.06 [58])</li> <li>• Overcomes non-differentiable training bottlenecks (A sum rate of 9.0 at 30 dB SNR, compared to 8.2 without KD [59])</li> </ul>	Provides reliable and robust channel modeling foundation
CSI Feedback	To deploy complex CSI feedback models on resource-limited devices	<b>Improved model's performance</b> <ul style="list-style-type: none"> <li>• Up to 7.6 dB in NMSE [60]</li> <li>• NMSE dropping from -9.54 dB to -9.75 dB [11]</li> <li>• improvement of 0.40–0.54 dB PSNR [61])</li> </ul> <b>Reduced training time and computational complexity</b> <ul style="list-style-type: none"> <li>• Reduced computational complexity to 25.50%–43.46% [60]</li> <li>• Reduced total training time by 66.8% [11]</li> </ul>	Enables efficient inference at edge nodes
Modulation classification	Compresses complex detection networks while retaining generalization	<b>Higher accuracy</b> <ul style="list-style-type: none"> <li>• Boosted the student's maximum accuracy by 2.4% [62]</li> <li>• 77.86% accuracy under PGD attack at 30 dB SNR [63]</li> </ul>	Enables efficient inference at edge nodes
Beam prediction	Transfers knowledge from multi-sensor model to lightweight and enhances beam prediction model's robustness and reliability	<ul style="list-style-type: none"> <li>• Achieved a lower Mean Squared Error under attack [64]</li> <li>• Achieves 86.33% of the teacher's accuracy with only 12.4% of the parameters [65]</li> </ul>	Enables the deployment of a secure and robust beam prediction model on base stations
Resource Allocation	Transfers prior knowledge via teacher models to reduce search space	<ul style="list-style-type: none"> <li>• Reduced Energy Consumption by an average of 94% compared to standard Federated Learning [66]</li> <li>• Lower Latency and Cost [67]</li> <li>• Improved Accuracy with the algorithm converging faster to accuracy level of approximately 0.80 [67]</li> </ul>	Facilitates efficient resource coordination for EGI

without increasing its complexity. The method boosted the student's maximum accuracy by 2.4% while the student's floating point operations remained less than half of the teacher's. Further enhancing reliability, Xu *et al.* [63] proposed an adversarial robust distillation strategy to transfer the resilience of a large, adversarially-trained teacher to a compact student. This was necessary because lightweight models are more vulnerable to attacks. Their proposed KD method enabled the student model to achieve 77.86% accuracy under a Projected Gradient Descent (PGD) attack at 30 dB SNR, which was 4.6% higher than AT and about 30% higher than other robust distillation techniques.

This line of work is central to the Perception stage of an agentic AI workflow. Through KD, a compact student model can be deployed on edge devices to enable real-time and reliable perception of the wireless environment. Accurate channel-state understanding forms the foundation for the agent's higher-level cognitive functions, supporting autonomous and secure operation in dynamic settings.

2) *Capability in “Planning”*: By using knowledge distillation, the complex and computationally intensive planning processes of a large teacher model, which encompass procedural reasoning and high-level strategies, can be compressed into a compact and efficient student model suitable for real-world deployment.

Kuzlu *et al.* [64] applied a defensive distillation framework to transfer a teacher model's robustness to a compact student, improving the beam prediction model's resilience against adversarial attacks. Park *et al.* [65], [75] introduced cross-modal relational KD to transfer knowledge from a complex multimodal teacher with the use of LiDAR, radar, GPS, and RGB camera to a lightweight radar-only student, enabling resource-efficient beam prediction without relying on expensive sensor data. Moreover, Yang *et al.* [76] developed a knowledge-driven approach by unrolling the domain knowledge of the iterative

algorithm into a Graph Neural Network (GNN), creating a scalable and interpretable model for efficient resource allocation. Furthermore, Gong *et al.* [77] employed distillation methods to transfer a policy learned by a computationally expensive deep reinforcement learning (DRL) agent into a gradient boosting decision tree model, which achieves faster and more cost-effective resource optimization.

A robust and efficient planning faculty acts as the agent's cognitive engine, translating its goals and perceived state into a structured and executable course of action. This is the foundational layer that enables the agent to operate with foresight and intentionality, facilitating sophisticated, goal-directed behavior in dynamic and complex environments.

a) **Policy Distillation:** Policy distillation operates on the core principle of behavioral mimicry. The “knowledge” transferred is the teacher's decision-making policy, typically represented as a probability distribution over all possible actions for a given state [78].

To dynamically manage resources in network slicing, Li *et al.* [79] used DRL to create policies that adapt to volatile user demand without a traffic model. This approach could learn effective policies for complex environments, but is limited by the high training cost. To efficiently implement these policies in distributed settings, policy distillation has emerged as a critical technique to optimize decision-making models in wireless networks, enabling efficient and collaborative resource management. Zhou *et al.* [80] developed a federated bidirectional distillation mechanism to reduce the high communication overhead inherent in Federated Reinforcement Learning. By distilling knowledge to and from a central server, the method cut the required communication rounds by nearly 50% in non-IID scenarios. Moving beyond simple compression, Mensah *et al.* [81] introduced a federated mutual policy distillation scheme for collaborative resource trading among heterogeneous agents. This was necessary to overcome learning instability and non-

convergence caused by conflicting objectives and non-IID data in multi-agent systems. By allowing agents to learn from each other’s policies, the method improved the total system utility by 12.5% compared to baseline DRL techniques.

Despite its effectiveness, policy distillation exhibits inherent limitations. As a quintessential “closed-box” method, it transfers final decisions while the underlying reasoning and causal logic are lost [82]. Furthermore, if the environment changes, a mismatch between the distilled policy and the value function can lead to incorrect judgments, training divergence, and catastrophic forgetting [83].

**b) Interpretable Strategy Teaching:** The central idea of Interpretable Strategy Teaching is to extract and transfer the knowledge embedded in a teacher agent in the form of explicit, interpretable, natural language strategies [84].

To enhance the interpretability of AI in wireless communications, Zaidi *et al.* [85] pioneered an intrinsic approach by developing a domain knowledge-guided neural network whose architecture directly reflects the physics of radio propagation, making its decisions inherently transparent and robust. In contrast, Chen *et al.* [86] adopted a post-hoc strategy for complex DRL agents, proposing a framework that analyzes an agent’s behavior to infer comprehensible decision heuristics, thereby building trust and providing actionable insights for engineers.

This paradigm offers several advantages. Natural language strategies improve interpretability and enable expert auditing [87]. It is also compatible with closed-box models, requiring only API access rather than internal parameters [87]. By separating a general-purpose model from an editable strategy library, it yields a modular hybrid intelligence framework that supports lifelong learning and enhances the safety and trustworthiness of agentic AI systems.

**c) Chain-of-Thought (CoT) Distillation:** CoT prompting [88] represents a pivotal technique for generating the reasoning component of an agent’s trajectory. CoT guides a large language model to emulate a human cognitive process by generating intermediate reasoning steps prior to reaching a final answer [89]. The technique augments few-shot exemplars in a prompt with explicit reasoning chains, formatting examples as (input, CoT, output) rather than simple (input, output) pairs. Wang *et al.* [90] proposed a framework where a CoT enabled LLM translates high-level user intents into executable wireless network control policies. They used CoT to overcome the key limitations of standard LLMs in wireless applications, such as their poor multi-step reasoning, lack of interpretability, and tendency to generate incorrect or “hallucinated” outputs. This resulted in significant, measurable gains, including a 27.2% increase in network sum rate in a Unmanned Aerial Vehicle (UAV) case study compared to a non-CoT baseline.

CoT-based KD is a methodology designed to transfer the complex, multi-step reasoning capabilities of large “teacher” language models to smaller, more efficient “student” models, enabling the model to decompose complex problems into manageable sub-problems for sequential resolution [91].

Based on the analysis above, the following strategic recommendations can be provided for different application scenarios in Fig. 4.

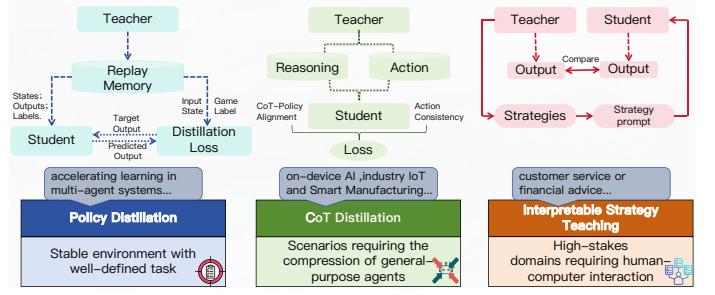


Fig. 4. Strategic Recommendations of Policy Distillation, Interpretable Strategic Teaching and CoT-Based KD. Policy Distillation for stable environments with well-defined tasks [78]; Interpretable Strategy Teaching for high-stakes domains requiring human-computer interaction [92]; CoT-Based KD for scenarios requiring the compression of general-purpose agents [91]

**3) Capability in “Action” :** In the context of EGI, action represents the system’s capacity to translate perception and reasoning into concrete operations within its environment. It encompasses both digital and physical interactions, ranging from executing database queries or invoking local tools to coordinating with external devices or robotic actuators. KD transcends its origin as a mere model compression technique, evolving into a fundamental mechanism that endows EGI with a suite of critical, action-oriented capabilities. Specifically, KD enhances *performance-optimized decision making* [93] by guiding lightweight models toward robust and adaptive policies—for instance, enabling mobile agents to perform low-latency beam selection in wireless communications. It enriches *action reliability under data scarcity* [64] [11] through pseudo-labeling and synthetic data generation, which can stabilize handover decisions when limited user trajectory data are available. It enables *cross-domain and multi-task generalization*, ensuring that actions remain transferable and coordinated across heterogeneous settings such as UAV-assisted networks [94] and RIS-enabled environments [95], [96]. Finally, it facilitates *privacy-preserving, efficient coordination* among distributed agents, exemplified by federated KD frameworks for collaborative spectrum sensing or interference management. Together, these dimensions illustrate how KD serves as a unifying toolkit to strengthen the effectiveness, generalizability, and trustworthiness of actions in EGI.

**a) Performance-Optimized Decision Making:** Traditionally, decision-making at the edge is constrained by limited compute and stringent latency requirements, often forcing agents to trade off accuracy for efficiency. Contrary to this limitation, recent developments in KD show that lightweight students can even outperform their teachers, as demonstrated in Born-Again Networks [97] and self-distillation frameworks [39]. Mechanistically, KD regularizes the optimization trajectory toward flatter minima [98], [99], which enhances robustness under dynamic and uncertain environments. Moreover, extensions such as bias-corrected distillation [100], [101] mitigate flawed teacher guidance, enabling edge agents not only to replicate but also to refine decision policies. For EGI, this transforms action from resource-limited execution into adaptive and resilient control in real-time contexts.

**b) Action Enrichment through Data Augmentation:**

EGI agents frequently work under data-scarce or noisy conditions, where direct supervision or high-quality data is difficult to obtain. KD addresses this limitation by leveraging soft targets, pseudo-labeling, and synthetic data generation to enrich training signals, thereby enabling action policies to remain reliable even in uncertain scenarios. In data-scarce settings, Li *et al.* [102] and Wu *et al.* [103] demonstrated that KD serves as a powerful amplifier of data value, primarily through pseudo-labeling. The teacher's soft outputs provide richer, more regularized signals than one-hot labels, making KD particularly effective in semi-supervised and unsupervised learning [104]. Beyond generic vision or language tasks, such mechanisms are also promising in wireless applications. For instance, pseudo-labeling can support automatic modulation classification under limited labeled samples, enhance channel state information feedback with synthetic training data, or improve spectrum sensing reliability in low-SNR environments.[105] Mechanistically, techniques such as confidence thresholding and explicit denoising frameworks [102], [106] enhance robustness against noisy supervision. Beyond pseudo-labeling, generative models such as GANs, a class of models in which a generator and a discriminator are trained in opposition to produce realistic synthetic samples, have been employed to synthesize new data, which are subsequently labeled by the teacher to train the student, aligning with the paradigm of Data-Centric AI that emphasizes improving data quality over altering model architectures [107]. For EGI, these advances translate into action policies that can maintain confidence and adaptability under severe data scarcity, strengthening agents' ability to act effectively in unpredictable edge environments.

**c) Cross-Domain and Multi-Task Actions:**

In real-world EGI systems, actions often span heterogeneous domains such as communication, sensing, and control, requiring policies that generalize across tasks and modalities. KD has emerged as a powerful enabler of such cross-domain and multi-task action integration. A central challenge lies in the heterogeneity of feature spaces, which KD mitigates through feature space alignment techniques [108], [109], often enhanced by contrastive learning strategies [110]. In multi-task learning (MTL), KD addresses negative transfer and asynchronous convergence by facilitating knowledge sharing across tasks [111]. For example, Li *et al.* [112] proposed leveraging multiple task-specific expert teachers to guide a unified multi-task student, thereby balancing specialization and generalization. More broadly, distilled knowledge can function as a shared intermediate representation, allowing diverse models to be modularly composed into larger AI systems. In wireless networks, these advances enable models to transfer knowledge across heterogeneous tasks, such as modulation classification [113], channel estimation [59], and interference management, thereby enhancing adaptability and scalability in dynamic communication environments.

**d) Privacy-Preserving and Efficient Action Coordination:**

EGI relies on distributed agents collaborating across heterogeneous devices and networks, where privacy protection and communication efficiency are critical. KD provides an effective solution by enabling the exchange of distilled outputs (e.g., logits) instead of raw parameters or sensitive data, thereby

reducing both privacy risks and communication overhead. Li *et al.* [114] demonstrated that output-sharing on a public proxy dataset enhances privacy [115] while improving communication efficiency [116], a significant improvement over standard federated averaging (FedAvg) [117]. However, dependence on public datasets introduces new challenges such as dataset selection and potential leakage. To address this, Lu *et al.* [118] and Zhu *et al.* [119] employed generative models to synthesize proxy data, avoiding reliance on external datasets. Further advances include decentralized peer-to-peer (P2P) and co-distillation frameworks [120]–[122], where clients directly exchange knowledge with neighbors [123], often using local data as a temporary proxy for peer evaluation and weighted updates [124]. For EGI, these approaches enable privacy-preserving and resource-efficient action coordination, allowing edge agents to collaborate securely and effectively in distributed environments.

**4) Capability in “Memory” :** In the context of EGI, memory can be understood across three complementary dimensions. Model memory concerns the ability to store compressed or distilled models on resource-constrained edge devices, reflecting how efficiently knowledge is represented and retained for long-term use. Operational memory refers to the dynamic allocation of memory during inference and training, which directly impacts the agent's capacity to execute real-time tasks under limited hardware resources. Finally, knowledge memory captures the system's capability to preserve and recall past knowledge, for example retaining information about channel quality variations [125], user mobility patterns [126], or interference conditions [127] in wireless networks, which can be leveraged in continual learning scenarios [128] or through the sharing of distilled global representations in federated settings.

Recent works have explored how KD can effectively enhance memory efficiency at the edge through storage reduction, runtime optimization, and knowledge retention. Chen *et al.* [129] introduced a computation offloading framework that integrates Deep Imitation Learning with KD to reduce task latency in heterogeneous edge–cloud settings. By distilling offline-learned optimal strategies into lightweight student models, their method enables real-time decision-making on mobile devices and offers a principled alternative to heuristic offloading approaches. Li *et al.*[130] tackled catastrophic forgetting which appears because Deep neural networks tend to forget the learned knowledge of previous tasks when sequentially trained on a stream of tasks. Unlike rehearsal-based methods that demand high storage and risk privacy leakage, their data-free distillation embeds past knowledge into gradient-based regularization, allowing models to recall prior tasks without retaining raw data. This shift highlights KD's role in strengthening functional memory for lifelong edge learning. Extending from knowledge preservation to distributed personalization, Pan *et al.*[131] introduced FedCache 2.0, where dataset distillation replaces parameter exchange. Clients contribute compact synthetic datasets to a central cache, reducing communication and storage burdens while enabling adaptive personalization, though sensitivity to client heterogeneity remains a challenge.

Together, these works illustrate a layered trajectory: from optimizing runtime memory, to safeguarding knowledge memory, to redefining storage and communication memory, and finally to

enhancing functional memory and robustness. KD thus emerges as a versatile paradigm for strengthening memory in diverse edge intelligence scenarios, forming a natural bridge toward the broader vision of Edge General Intelligence.

### B. Beyond Transformers: Architectures Natively Designed for the Edge

The inherent quadratic complexity of the Transformer architecture remains a fundamental bottleneck for edge deployment. This has spurred the development of novel architectures engineered for computational efficiency. However, these emerging models, such as Mamba[141] and RWKV [16], often lack the extensive pre-training on vast datasets that underpins the powerful capabilities of their Transformer-based predecessors. KD, particularly in the form of cross-architecture distillation, emerges as a critical and capital-efficient strategy to bridge this gap. By transferring the rich, generalized knowledge from a large, pre-trained Transformer "teacher" to a more efficient "student" model like Mamba or RWKV, it is possible to imbue these lightweight architectures with advanced capabilities without incurring the prohibitive costs of training them from scratch. This symbiotic relationship between efficient novel architectures and sophisticated distillation techniques is paving the way for the next generation of deployable, high-performance models for Edge General Intelligence.

*1) Mamba: Architecture with Linear-Time Sequence:* The Mamba architecture can be regarded as a significant evolution of the Structured State-Space Model (S4) [141]–[143], inheriting its mathematical formulation while introducing novel mechanisms to enhance flexibility and efficiency. Mamba also incorporates an efficient algorithm with parallel scan and kernel fusion, streamlining its architecture by unifying the separate attention and MLP layers of a Transformer into a homogeneous Mamba block [141]. Mamba maintains the latent state formulation, where the hidden state evolves over time based on the input and parameterized dynamics.

Recent research has established the Mamba architecture as a highly efficient and versatile tool for wireless systems. For network optimization, *Mehrabian and Wong* [133] proposed the A-Gamba model for cellular traffic prediction, achieving a  $7\times$  faster training time and a  $681\times$  reduction in multiply-accumulate operations (MACs) over strong baselines . In core communication functions, *Zhang et al.* [132] demonstrated that Mamba can enhance semantic decoding, while also enabling scalable resource allocation such as energy-efficient beamforming with near-constant inference time as user numbers grow. Moreover, in wireless human sensing, *Liu et al.* [144] developed TF-Mamba that reduces MACs by up to 98.9% while achieving high accuracy. Similarly, *Huang et al.* [134] created SenseMamba, which achieves state-of-the-art accuracy with only 0.021M parameters. Extending this to multimodal perception, *Yang et al.* [145] showed that a Mamba-based framework could improve performance when sensor data is missing by fine-tuning less than 0.3% of its parameters. Collectively, these attributes underscore Mamba's potential as a key enabling technology for decentralized, edge-native intelligence in next-generation wireless networks and even the possible way towards mobile agentic AI.

*2) RWKV: RNN-Transformer Hybrid Architectures:* The core purpose of RWKV [16] is to integrate the high performance of Transformers with the computational efficiency of recurrent neural networks (RNNs). Its core innovation lies in a linear-complexity attention mechanism: the hidden state  $h(t)$  depends only on  $h(t-1)$  and the current input  $x(t)$ , enabling step-by-step computation with constant memory and linear time complexity.

RWKV is composed of four fundamental components: (i) R(Receptance): a gating mechanism regulating how much past information influences the current state. (ii) W(Weight):a channel-specific, learnable time decay vector. (iii) K(Key): analogous to attention keys, computed via linear interpolation of current and past tokens. (iv) V(Value): analogous to attention values, also derived through linear interpolation of current and past tokens. Together, these define the WKV operator, which unifies transformer-style parallel training with RNN-style efficient inference.

Recent studies have leveraged the RWKV architecture to create highly efficient solutions. For instance, *Fu et al.* [136] developed an RWKV-based channel estimator for Vehicle-to-Everything communications that reduces computational complexity by 24.9% compared to LSTM-based methods. In the domain of network security, *Xie et al.* [135] proposed an RWKV-based scheme to detect selective forwarding attacks in wireless sensor networks, achieving an average false detection rate of just 0.09% in harsh mobile environments. Critically, the architecture's suitability for the wireless edge has been bolstered by the work of *Choe et al.* [146], who introduced a suite of compression techniques that reduce the RWKV model's memory footprint by  $3.4\times$  to  $5\times$  with only negligible accuracy degradation, making it feasible for resource-constrained devices.

#### 3) Cross-Architecture Distillation:

New architectures such as Mamba and RWKV promise superior efficiency but lack the trillions of tokens of pre-training data and massive computational investment that have made Transformers so capable [147]. Cross-architecture distillation breaks this cycle by providing a capital-efficient pathway [148].

Unlike same-architecture distillation, heterogeneous architectures exhibit significant feature divergence [149].This mismatch is particularly acute when distilling non-causal models like Transformers into causal models like Mamba. Consequently, naive feature-based distillation methods often fail, necessitating more sophisticated alignment techniques [150]. To overcome these challenges, *Yao et al.* [151] constructed a hybrid model composed of 86% of Mamba blocks and 14% of self-attention blocks, leveraging the computational efficiency of Mamba blocks while strategically incorporating self-attention blocks to handle global interactions. Another direction in cross-architecture distillation involves the use of specially designed "projectors" to align feature spaces. *Liu et al.* [152] and *Hao et al.* [150] introduced methods that utilize these projectors to map the student's features into the teacher's feature or attention space.

As discussed previously, the high efficiency and less model parameters achieved through cross-architecture distillation make it a highly advantageous technique for edge deployment. To enable the deployment of an accurate retinal disease classifier on an edge device, *Yilmaz and Aiyengar* [153] used a

TABLE IV  
A COMPREHENSIVE COMPARATIVE ANALYSIS OF SEQUENCE MODELING ARCHITECTURES:

Feature	Mamba	RWKV	Transformer	RNN (LSTM/GRU)
Core Mechanism	Selective State Space Model	Linear Attention with Time Decay	Self-Attention	Gated Recurrence
Training Time	$O(B \cdot L \cdot D \cdot N)$	$O(L \cdot d^2)$	$O(L^2 \cdot d)$	$O(L \cdot d^2)$
Training Space	$O(B \cdot L \cdot D \cdot N)$	$O(L \cdot d^2)$	$O(L^2 + L \cdot d)$	$O(d^2)$
Inference (per token) Time	$O(1)$ : Constant time	$O(1)$ : Constant time	$O(L)$ : Linear in sequence length	$O(1)$ : Constant time
Inference (state) Space	$O(1)$ : Constant space.	$O(1)$ : Constant space.	$O(L)$ : Linear in sequence length	$O(1)$ : Constant space.
Key Advantages	<ul style="list-style-type: none"> <li>State-of-art performance</li> <li>Linear time complexity</li> <li>Extremely fast inference</li> <li>Excellent long-context scaling</li> </ul>	<ul style="list-style-type: none"> <li>Excellent inference efficiency</li> <li>Parallelizable training</li> <li>No quadratic bottleneck</li> <li>Simple recurrent formulation</li> </ul>	<ul style="list-style-type: none"> <li>Best-in-class performance at moderate scales</li> <li>Parallelizable training</li> <li>well-established ecosystem</li> <li>No information bottleneck</li> </ul>	<ul style="list-style-type: none"> <li>Extremely low inference compute cost</li> <li>Conceptually simple</li> <li>Memory footprint during inference</li> </ul>
Key Disadvantages	<ul style="list-style-type: none"> <li>Less explored architecture</li> <li>May underperform LTI SSMs on continuous data</li> <li>Relies on custom CUDA kernels for efficiency</li> </ul>	<ul style="list-style-type: none"> <li>Information bottleneck limit performance on tasks requiring high-fidelity recall</li> <li>Highly sensitive to prompt structure/information order</li> </ul>	<ul style="list-style-type: none"> <li>Quadratic scaling issues in time and memory</li> <li>Slow and memory-intensive inference due to KV cache</li> <li>Prohibitively expensive for very long sequences</li> </ul>	<ul style="list-style-type: none"> <li>Really difficult to parallelize training</li> <li>Vanishing gradient problem</li> <li>Poor at capturing long-range dependencies</li> </ul>
Applications in Wireless area	<ul style="list-style-type: none"> <li>General-purpose applications in wireless communications and networking [132], [133]</li> <li>Lightweight, efficient real-time human activity and gesture recognition with wireless signals [134]</li> </ul>	<ul style="list-style-type: none"> <li>Detecting malicious node attacks in wireless sensor networks under harsh environmental conditions [135]</li> <li>Channel estimator for Vehicle-to-Everything communications [136]</li> </ul>	<ul style="list-style-type: none"> <li>Vision Transformer-based semantic communication for efficient and robust image transmission [137]</li> <li>Deep learning-based hybrid beamforming in millimeter-wave MIMO-OFDM systems [138]</li> </ul>	<ul style="list-style-type: none"> <li>Proactive resource optimization in optical backbone networks through deep learning-based traffic prediction [139]</li> <li>Low-complexity CSI prediction in massive MIMO (mMIMO) systems [140]</li> </ul>

L: sequence length    D(d): model dimension    B: batch size    N: state dimension

cross-architecture distillation framework to transfer diagnostic knowledge from a large Vision Transformer teacher model to a lightweight CNN student model. Their framework, featuring specialized projectors, dramatically improved the student's performance with 97.4% reduction in the model parameters (from 85.8M to 2.2M), which resulted in a compact 8.79 MB model that retained 93% of the teacher's diagnostic accuracy on the edge device. The primary strength of this approach is the successful creation of an efficient, high-accuracy model for resource-constrained settings. However, its limitations include being trained on a relatively small dataset and for a limited number of epochs due to computational constraints [153].

### C. Utilize KD to obtain compact and deployable edge models

1) *Language Models as a Foundation: From BERT to TinyBERT*: Language models (LMs) are central to the core functions of modern intelligent agents, operationalizing intelligence as the capability to interact with environments, utilize tools, and achieve goals. However, the prevailing architecture, which relies on Large Language Model (LLM) API endpoints hosted on centralized cloud infrastructure, is fundamentally incompatible with the resource-constrained nature of the edge due to computational cost, memory footprint, and energy consumption.

Consequently, a new vision is emerging where general intelligence is realized not as a singular cloud entity but as the aggregate capability of numerous specialized SLMs-powered agents operating autonomously in local contexts. The true potential of SLMs for edge intelligence is unlocked by knowledge distillation. In this section, we use BERT as a case study to illustrate this point.

The deployment challenges posed by large-scale models like Bidirectional Encoder Representations from Transformers (BERT), whose exceptional performance is linked to its immense size, are systematically addressed by knowledge distillation.

Sanh *et al.* [154] create a miniaturized general-purpose version of BERT, which serves as a seminal example of successfully applying knowledge distillation. Its core contribution lies in demonstrating that distillation can be performed during the computationally intensive pre-training phase. DistilBERT's compression effects are remarkably significant, successfully striking an excellent balance between model efficiency and performance. Compared to BERT-Base, DistilBERT reduces the parameter count by 40% and achieves a 60% increase in inference speed. Moreover, DistilBERT retains 97% of BERT-Base's language understanding capabilities. On the GLUE benchmark, which includes 9 different tasks, its average score is very close to that of BERT-Base.

If DistilBERT represents the pioneering application of KD to BERT, then Jiao *et al.* [155] embodies a more profound and refined distillation strategy. TinyBERT pioneers a "Transformer distillation" method that delves deep into the internal mechanisms of the Transformer, aiming for a fine-grained imitation of the teacher model's intermediate representations at every layer.

TinyBERT has achieved astonishing results in model compression with its advanced two-stage, multi-level distillation framework. Taking the 4-layer TinyBERT4 as an example, its parameter count is merely 14.5 million. Compared to BERT-Base, its model size is reduced by 7.5 times, while its inference speed is boosted by 9.4 times. While achieving such extreme

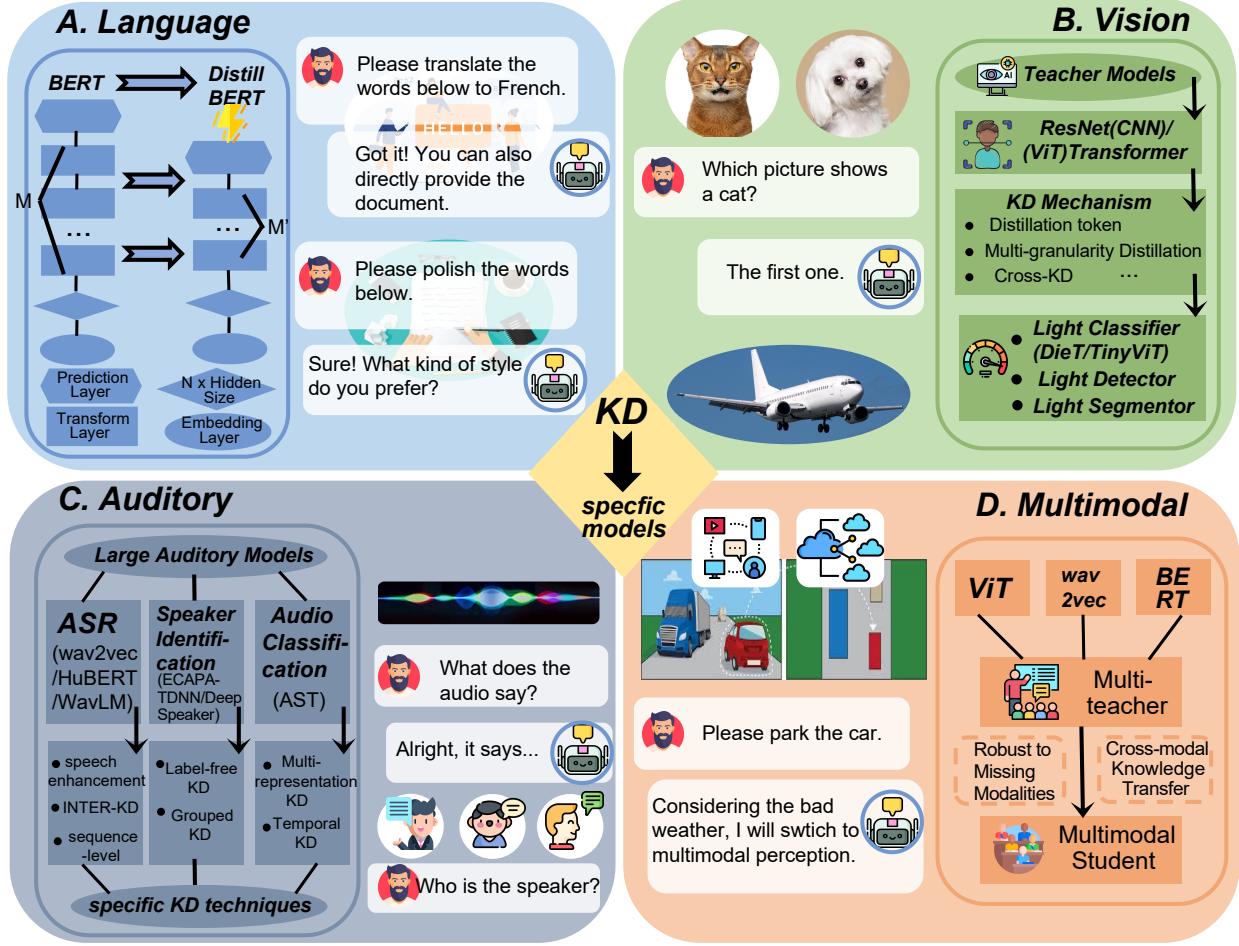


Fig. 5. An overview of KD techniques across different modalities. (A) Language: distillation from large language models such as BERT into compact models (e.g., DistilBERT). (B) Vision: KD mechanisms including distillation tokens, multi-granularity distillation, and cross-KD applied to light classifiers, detectors, and segmentors. (C) Auditory: specific KD strategies such as sequence-level KD, label-free KD, grouped KD, multi-representation KD, and temporal KD for tasks like ASR, speaker identification, and audio classification. (D) Multimodal: integration of KD across modalities with multi-teacher settings, robust handling of missing modalities, and cross-modal knowledge transfer for multimodal student models.

compression, TinyBERT4 retains 96.8% of BERT-Base’s performance on the GLUE benchmark, with minimal performance loss. Even more impressively, a 6-layer TinyBERT6 performs nearly on par with the full BERT-Base teacher model on GLUE. The emergence of TinyBERT provides a powerful and reliable option for deploying high-performance NLP models designed to process and understand human language on edge devices with extremely limited resources.

2) *Distilled Models for Wireless Tasks:* In wireless communication scenarios, the development of compact distilled models has not yet reached the same level of standardization as in NLP, where models such as DistilBERT and TinyBERT serve as well-established benchmarks. Nevertheless, in several wireless-specific tasks—such as CSI feedback and reconstruction, there exist a variety of efficient and task-tailored knowledge distillation frameworks specifically designed to address the unique characteristics of wireless systems.

#### • CRNet-SE [11]

In recent years, various KD-based models have been proposed to enhance CSI feedback and edge inference efficiency in wireless communication systems. Among

them, CRNet-SE focuses on the CSI feedback problem in large-scale MIMO systems. Its main objective is to compress high-dimensional CSI on user equipment (UE) with limited computational resources. The teacher model, CRNet, employs a multi-resolution encoder to extract rich channel features, while the student model simplifies the encoder into only one convolutional and one fully connected layer for lightweight deployment. Despite this reduction, the KD-trained student achieves a significantly lower NMSE than a non-distilled counterpart, illustrating the efficacy of KD in preserving reconstruction accuracy under severe model compression.

#### • CSI-ALM-Light [12]

While CRNet-SE emphasizes compression for feedback efficiency, CSI-ALM-Light extends the KD paradigm toward temporal prediction in high-mobility communication scenarios. Instead of reconstructing current CSI, CSI-ALM-Light predicts future downlink CSI using historical uplink information. Its teacher, CSI-ALM, builds on a pre-trained GPT-2 backbone and integrates a Modality Alignment Mechanism to align numerical CSI features with

language embeddings. The distilled student replaces these complex components with a lightweight Transformer and learnable soft prompts, while employing a self-attention relation-based distillation strategy that aligns multi-level attention relations (Query-Key, Query-Query, Key-Key, and Value-Value). Remarkably, CSI-ALM-Light achieves near-teacher performance using only 10% of the data and less than 1/200 of the parameters, highlighting KD's potential to bridge performance gaps across drastically different model scales.

- **KD-Based AIoT Framework [13]**

Complementary to these communication-oriented designs, the KD-Based AIoT Framework applies KD to Wi-Fi CSI-based gesture recognition for intelligent edge devices. Unlike CRNet-SE and CSI-ALM-Light, which focus on CSI reconstruction or prediction, this framework leverages KD to enable efficient human activity recognition (e.g., waving, walking, falling) under strict memory and computation constraints. The student network, containing merely 28K parameters, integrates multiple parallel convolutional branches and a bottleneck fusion layer, and through a joint response-plus-feature distillation scheme, it achieves up to 99.5% accuracy on SignFi datasets—outperforming both response-only and hard-label training.

Overall, these studies collectively demonstrate that KD serves as a versatile mechanism across diverse wireless and AIoT tasks—from channel compression (CRNet-SE), to temporal prediction (CSI-ALM-Light), and edge activity recognition (KD-Based AIoT)—each highlighting a distinct yet complementary facet of efficient model design under hardware constraints.

#### D. Lesson learned

A lesson from this analysis is that KD is not a uniform methodology but a collection of adaptive strategies whose effectiveness varies by the target agentic capability. In NLP, KD primarily preserves architectural fidelity and semantic representation, while in wireless communication tasks, domain constraints such as limited data and noisy, resource-restricted edge environments shift the focus toward task-oriented adaptation, as demonstrated by models and the KD-Based AIoT Framework [13]. However, existing wireless KD approaches remain mostly task-specific prototypes and lack generalized distillation protocols comparable to those in NLP [93]. Moreover, KD efficiency depends heavily on the teacher's domain adaptability and the quality of intermediate representations [156], [157], which are harder to define for non-textual modalities. Future directions include establishing standardized KD benchmarks for edge wireless intelligence, developing modality-agnostic distillation mechanisms transferable across heterogeneous inputs [158], and enabling collaborative KD across multiple edge nodes to approximate cloud-level intelligence [159].

## IV. APPLICATION

This section explores the specialized architectures and technologies for specific domains such as Unmanned Aerial Vehicles (UAVs), Autonomous Vehicles(AV), Robotics and Other IoT Devices.

### A. UAV

KD is a critical technique for deploying complex models on resource-constrained UAVs, with researchers developing specialized frameworks to enhance on-device efficiency and performance across different operational domains. These approaches can be broadly categorized by their primary objective, whether optimizing high-level mission logistics, bolstering on-board cybersecurity, or refining the fundamental distillation process for perception tasks.

To begin with, in the domain of mission optimization, *Sun et al.* [160] created a cooperative co-evolutionary framework to address the challenge of running complex task-planning algorithms on energy-limited UAVs. In their approach, multiple student subnetworks collaboratively learn and exchange information to compress a large Deep Neural Network (DNN) to just 1.3% of its original parameters. The use of KD here drastically reduces task completion delay by up to 95% in large-scale scenarios. This work illustrates that efficiency is only one side of the equation, security remains another vital concern for autonomous UAVs operating in dynamic and adversarial environments.

The efficiency gains from KD are foundational, not merely incremental, for the evolution of UAVs. By enabling complex cognitive capabilities such as on-device reasoning and planning, KD provides the crucial bridge to transform UAVs from automated tools into autonomous, goal-driven Agentic UAVs [20]. This latest evolution, the Agentic UAV, combines cognitive AI models, multimodal sensing, and edge computing to support real-time perception and reasoning. To address cloud latency and limited onboard resources, Agentic UAVs employ compact edge AI platforms for tasks such as semantic segmentation and path reconfiguration [20]. The capabilities of these agents are further extended by Vision-Language Models (VLMs), which enable the execution of high-level user commands through semantic grounding and zero-shot generalization. Thus, KD serves as the enabling foundation upon which higher-order autonomy and cognitive generalization are being built in modern UAV systems.

However, the story of UAV intelligence does not end at the single-agent level. Given the persistent limitations on battery and processing power[161], UAVs often offload compute-intensive tasks, such as 3D mapping and large-scale inference, to edge servers [21]. This shift from individual to distributed intelligence marks the transition toward swarm-based and EGI-powered UAV ecosystems, where multiple agents collaboratively sense, reason, and act. In summary, while edge computing techniques like KD greatly enhance UAV autonomy by overcoming inherent physical limitations, the evolution toward distributed cognitive systems like EGI-powered swarms introduces new challenges. The opacity of deep learning creates an “explainability-accountability crisis,” complicating error analysis, while physical safety remains a critical, underexplored area for ensuring public acceptance and safe deployment [21].

### B. Autonomous Vehicles

Achieving Level-5 autonomous driving, where a vehicle can operate in all environments without any human intervention, requires extremely robust, low-latency, and high-reliability AI capabilities. Such systems must continuously process large

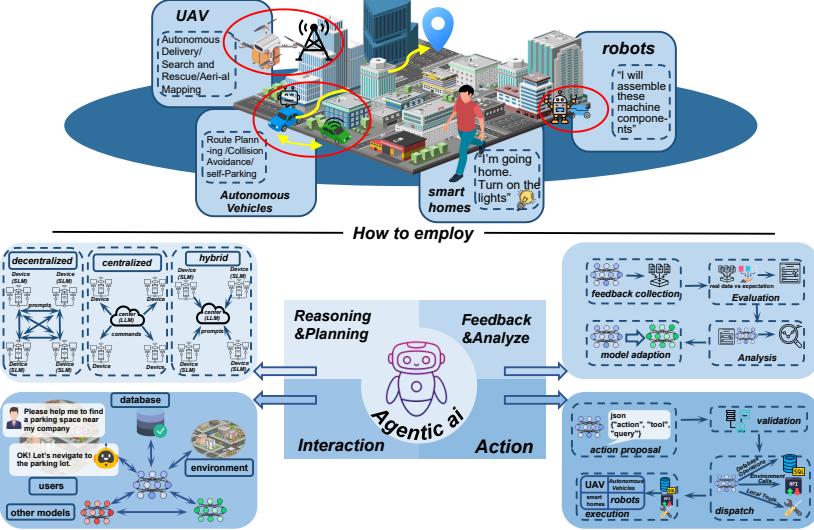


Fig. 6. Illustration of KD in enabling EGI within mobile agentic AI ecosystems. Use cases span UAVs (e.g., autonomous delivery, search and rescue, aerial mapping), autonomous vehicles (e.g., route planning, collision avoidance, self-parking), robots (e.g., task assembly, household assistance), and smart homes. The figure highlights how agentic AI leverages both LLMs in the cloud and smaller student models (SLMs) at the edge to achieve reasoning, planning, feedback, and interaction. KD facilitates the transfer of memory and learned knowledge, enabling mobile agentic AI to perform efficient, real-time, and context-aware decision-making at the edge.

volumes of high-dimensional, multi-modal sensor data (e.g., camera, LiDAR, radar, GPS, IMU) and make real-time decisions for perception, planning, and control [18].

Current approaches to supporting autonomous vehicles can be broadly categorized into five types: vehicle-based, vehicle-to-vehicle collaborative computing, cloud-assisted, edge-assisted, and cloud-edge-assisted approaches [19]. To realize EGI for autonomous vehicles (AVs), general methods like KD are complemented by specialized techniques such as Bird's-Eye-View (BEV) perception and Collaborative Edge Intelligence (CEI).

BEV perception transforms and fuses data from multiple cameras, LiDAR, or millimeter-wave radar sensors to generate a unified, top-down spatial representation. This approach mitigates challenges like perspective distortion and occlusion inherent in single sensors, providing downstream path planning systems with more stable and information-rich inputs.

CEI is a distributed computing paradigm that integrates geo-distributed edge resources into a federated pool. It enables both vertical collaboration (cloud-edge-end) and horizontal collaboration (edge-to-edge), creating a decentralized intelligent network that can dynamically orchestrate resources to meet the stringent demands of autonomous driving [19].

### C. Robotics

The strict real-time and security requirements of robotic Operational Technologies (OT) are often poorly met by centralized cloud computing, making Edge Computing a critical complementary approach that provides computation closer to the system [22]. Several frameworks facilitate this, including FogROS2 [23] for task offloading, the EI-for-AMR SDK [162] for containerized deployment and hardware-accelerated inference, and the Edge-Driving Robotics Platform (EDRP) [163] for 5G-enabled collaborative mapping.

General-purpose robotics often involves perceiving and manipulating complex objects. Several studies have sought to address such challenges. For example, Chen *et al.* [54] adopt an innovative teacher-student learning paradigm, in which a 'privileged agent' equipped with complete information about the state of the cloth serves as the teacher, guiding the robot through complex cloth manipulation tasks using vision-based reinforcement learning. Similarly, for high-precision industrial tasks, Zhao *et al.* [164] proposes a novel semi-supervised knowledge distillation framework aimed at enhancing the robustness and accuracy of vision-based guidance in intelligent manufacturing settings.

Policies optimized in simulation often fail when deployed on physical robots due to discrepancies between simulated and real-world dynamics. A common solution involves training a teacher policy in simulation with access to privileged information (e.g., precise object velocities, friction coefficients), while the student policy learns to replicate its behavior using only the sensor data available to a physical robot (e.g., visual input, joint encoder readings) [165].

To manage diverse skills or collaborative multi-robot scenarios, specialized expert teacher policies can be trained for individual tasks. KD is then used to consolidate their capabilities into a unified student network. This process distills the collective behaviors of all teachers, allowing the student to learn cross-task commonalities and form a compact, generalizable decision-making model [166], [167].

### D. Other IoT Devices

The Internet of Things(IoT) features a massive, growing network of interconnected devices that generate immense volumes of heterogeneous, real-time data [24].

The Internet of Agents (IoA) provides the architectural blueprint to realize EGI, offering a scalable, agent-centric infras-

ture. Its hierarchical, layered design is engineered to manage the complexity of a global agent network [25]. The foundational infrastructure layer supplies essential resources, including AI models, diverse computing environments, multimodal data, and high-reliability communication networks.

Effective multi-agent communication in IoT environments requires balancing efficiency, reliability, and semantic richness. While lightweight protocols like MQTT [26] and CoAP are efficient for simple sensor data, they lack the semantic depth for complex agent interactions [168]. To address this, emerging agent-centric standards like Google’s Agent-to-Agent (A2A) protocol and Anthropic’s Model Context Protocol (MCP) enable goal-oriented collaboration. A2A uses “Agent Cards” for dynamic discovery and task coordination, while MCP unifies access to external tools and real-time data to enhance contextual awareness [169].

Agentic AI and the IoA are transitioning from theory to practice, with applications emerging across several domains:

- **Smart Homes:** Systems like Sasha [170] and SAGE [171] demonstrate advanced goal reasoning and personalization, while decentralized device-as-agent architectures enable proactive, privacy-preserving collaboration at the edge [25].
- **Urban and Industrial IoT:** Frameworks such as CityGPT [172] and CASIT [173] showcase large-scale multi-agent coordination. CityGPT augments human decision-making with spatiotemporal data analysis, and CASIT autonomously manages remote environmental monitoring by using hierarchical agent structures to optimize bandwidth.[174]
- **Smart Transportation:** Connected and Autonomous Vehicles (CAVs) [175] operate as cooperative agents via V2V and V2I communication [176] for applications like cooperative merging and platooning. Internally, multi-agent hierarchies manage planning, tactical decisions, and energy optimization.

#### E. Lesson learned

Across domains such as UAVs [20], autonomous vehicles [18], robotics [162], and IoT devices [24], KD consistently enables EGI by compressing large teacher models into lightweight, deployable student models for resource-constrained environments [9]. These applications highlight KD’s role in bridging high-capacity models and edge agents, supporting cross-modal knowledge transfer and real-time, context-aware decision-making [177]. However, the opacity of distilled models raises explainability and accountability concerns [178], and KD alone does not address physical safety in critical systems [9]. Promising directions include safety-aware distillation and extending KD from policy transfer to semantic understanding [179], enabling more transparent, robust, and efficient multi-agent coordination.

## V. CHALLENGES AND FUTURE DIRECTIONS

### A. Challenges

While these lessons highlight KD’s potential in enabling edge general intelligence, they also reveal a deeper truth: significant

technical and ethical challenges remain before EGI can be deployed safely and reliably at scale<sup>2</sup>.

#### 1) Confronting the Technical Challenges of Edge General Intelligence:

- **Toward New Benchmarks for On-Device Agentic Performance:** Most existing AI benchmarks are cloud-centric and focus only on static task accuracy, making them unsuitable for evaluating mobile agentic AI at the edge. Current edge benchmarks suffer from limited hardware and software coverage, unrealistic testing environments, lack of multi-tenancy support, and weak end-to-end evaluation. Moreover, they typically measure accuracy but ignore latency, energy efficiency, and robustness. Traditional benchmarks like ImageNet [180] and SQuAD [181] rely on pre-collected datasets, while on-device agents must perform within dynamic, interactive feedback loops.
- **Mitigating Hallucinations and Errors in Compressed Agents:** Large models often hallucinate, producing confident but incorrect outputs because training with hard labels enforces deterministic supervision and overconfidence. KD alleviates this by using softened teacher distributions, but still suffers from exposure bias: teacher-generated examples often exhibit high perplexity for the student, making adaptation difficult and reducing the effectiveness of distillation [182].
- **Vulnerabilities in Distillation and Federated Systems:** In federated distillation (FD), the server cannot observe how clients generate logits, making logits-poisoning attacks (LPA) particularly effective. Although FDLA [183] introduces a tailored LPA, its impact across heterogeneous clients and different distillation stages remains insufficiently studied. More sophisticated variants, such as Peak-Controlled FDLA (PCFDLA), manipulate peak logits to evade detection. Beyond these targeted attacks, Byzantine clients may degrade performance through parameter tampering or high-dimensional perturbations. Other attacks, including Local Model Averaging (LMA) and Collusion-based Poisoning (CPA), further obscure adversarial behavior. Privacy is also not guaranteed—logits can leak sensitive user information through membership inference attacks (MIAs) [184], [185].

#### 2) Navigating the Ethical Labyrinth of Edge General Intelligence:

The technical capabilities and limitations of EGI directly give rise to a series of profound ethical challenges. This section will systematically analyze these challenges and reveal their inextricable link to the aforementioned technical difficulties.

- **The Human Element:** Encoding complex and sometimes contradictory human values into EGI systems is an ethical challenge. Simulated empathy without real understanding risks manipulation, while rigid ethical rules fail in situations requiring nuanced judgment (e.g., the trolley problem [186]). The goal is not to hard-code emotions or ethics, but to enable a functional analogue of human wisdom.
- **Algorithmic Integrity and Trust:** Bias, privacy risks, and errors are inherent to data-driven AI systems [2].

<sup>2</sup>The authors would like to thank the Qwen, an advanced large language model, for its valuable academic polish service.

Ambiguity, adversarial inputs, or overfitting may cause harmful mistakes in high-stakes settings. Meanwhile, EGI must balance privacy with the need for authentic data, and biased training sets can lead to unfair or discriminatory decisions, especially in sensitive domains such as law enforcement, defense, and healthcare.

- **Societal and Economic Impacts:** The rise of EGI may trigger large-scale socioeconomic disruption. Its development requires massive investment, while its automation capabilities threaten broad employment sectors. Traditional responses such as retraining may be insufficient—not due to a skills gap alone, but the potential obsolescence of human labor. If EGI surpasses humans in most cognitive tasks, the role of human work in economic production becomes uncertain.

## B. Future Directions

- Safety-aware KD frameworks that explicitly optimize for model robustness and predictable behavior in real-world deployments.
- Modality-agnostic distillation mechanisms capable of transferring both structural and semantic priors across heterogeneous input forms [158].
- Collaborative multi-agent KD, where edge devices distill and exchange knowledge to collectively approximate cloud-level intelligence [159].
- Transitioning from policy transfer to semantic understanding transfer [179], enabling transparent decision-making and reducing the explainability gap.
- Instead of replacing human judgment, embed human oversight into decision-critical loops, ensuring accountability, interpretability, and alignment with societal norms.

## VI. CONCLUSION

The journey towards Edge General Intelligence is not merely a technical problem of model compression but a comprehensive scientific and engineering endeavor. By adopting the principled, synergistic approach detailed in this survey, which combined agent-aware capability transfer, edge-native architectures, and efficient adaptation methods, the research community can successfully bridge the deployment chasm. The challenges that lie ahead are formidable, spanning the technical domains of benchmarking and robustness to the deep socio-ethical questions of transparency, trust, and societal impact. Yet, these challenges define a clear and compelling research agenda for the coming years. This survey has charted a viable path forward, laying the groundwork for a future where intelligent, autonomous agents can operate securely, efficiently, and beneficially at the very edge of our increasingly connected digital world.

## REFERENCES

- [1] S. Hosseini and H. Seilani, “The role of agentic ai in shaping a smart future: A systematic review,” *Array*, p. 100399, May. 2025.
- [2] H. Chen, W. Deng, S. Yang, J. Xu, Z. Jiang, E. C. Ngai, J. Liu, and X. Liu, “Towards edge general intelligence via large language models: Opportunities and challenges,” *IEEE Network*, Feb. 2025.
- [3] Z. Yang, W. Xu, L. Liang, Y. Cui, Z. Qin, and M. Debbah, “On Privacy, Security, and Trustworthiness in Distributed Wireless Large AI Models (WLAM),” *arXiv e-prints*, p. arXiv:2412.02538, 2024.
- [4] G. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” *arXiv preprint arXiv:1503.02531*, 2015.
- [5] D. Xu, T. Li, Y. Li, X. Su, S. Tarkoma, T. Jiang, J. Crowcroft, and P. Hui, “Edge intelligence: Architectures, challenges, and applications,” *arXiv preprint arXiv:2003.12172*, 2020.
- [6] L. He, L. Fan, X. Lei, P. Fan, A. Nallanathan, and G. K. Karagiannidis, “The road toward general edge intelligence: Standing on the shoulders of foundation models,” *IEEE Communications Magazine*, Mar. 2025.
- [7] C. Zhao, G. Liu, R. Zhang, Y. Liu, J. Wang, J. Kang, D. Niyato, Z. Li, Z. Han, S. Sun *et al.*, “Edge general intelligence through world models and agentic ai: Fundamentals, solutions, and challenges,” *arXiv preprint arXiv:2508.09561*, 2025.
- [8] B. Wu, Y. Li, Y. Wei, M. Fang, and L. Chen, “Foundations and recent trends in multimodal mobile agents: A survey,” *arXiv preprint arXiv:2411.02006*, 2024.
- [9] J. Gou, B. Yu, S. J. Maybank, and D. Tao, “Knowledge distillation: A survey,” *International journal of computer vision*, vol. 129, no. 6, pp. 1789–1819, Mar. 2021.
- [10] A. P. Mohamed, A. S. M. M. Jameel, and A. El Gamal, “Knowledge distillation for wireless edge learning,” in *2021 IEEE Statistical Signal Processing Workshop (SSP)*. IEEE, Apr. 2021, pp. 600–604.
- [11] Y. Cui, J. Guo, Z. Cao, H. Tang, C.-K. Wen, S. Jin, X. Wang, and X. Hou, “Lightweight neural network with knowledge distillation for csi feedback,” *IEEE Transactions on Communications*, vol. 72, no. 8, pp. 4917–4929, Aug. 2024.
- [12] Z. Li, Q. Yang, Z. Xiong, Z. Shi, and T. Q. Quek, “Bridging the modality gap: Enhancing channel prediction with semantically aligned llms and knowledge distillation,” *arXiv preprint arXiv:2505.12729*, 2025.
- [13] M. Zeeshan, R. Raj, A. Anand, A. Pandey, M. T. Quasim, J. Torres-Sospedra, S. Kumar, and K. Dev, “Knowledge distillation-based aiot framework for efficient wireless gesture sensing in b5g/6g networks,” *IEEE Network*, pp. 1–1, 2025.
- [14] Q. Anthony, Y. Tokpanov, P. Glorioso, and B. Millidge, “Blackmamba: Mixture of experts for state-space models,” *arXiv preprint arXiv:2402.01771*, 2024.
- [15] J. Cao, Q. Zhang, J. Sun, J. Wang, H. Cheng, Y. Li, J. Ma, K. Wu, Z. Xu, Y. Shao *et al.*, “Mamba policy: Towards efficient 3d diffusion policy with hybrid selective state models,” *arXiv preprint arXiv:2409.07163*, 2024.
- [16] B. Peng, E. Alcaide, Q. Anthony, A. Albalak, S. Arcadinho, S. Biderman, H. Cao, X. Cheng, M. Chung, M. Grella, G. Kranthikiran, X. Du, X. He, H. Hou, P. Kazienko, J. Kocoñ, J. Kong, B. Koptyra, H. Lau, K. S. I. Mantri, F. Mom, A. Saito, X. Tang, B. Wang, J. S. Wind, S. Wozniak, R. Zhang, Z. Zhang, Q. Zhao, P. Zhou, J. Zhu, and R. Zhu, “Rwkv: Reinventing rnns for the transformer era,” in *Conference on Empirical Methods in Natural Language Processing*, May. 2023.
- [17] B. Peng, R. Zhang, D. Goldstein, E. Alcaide, X. Du, H. Hou, J. Lin, J. Liu, J. Lu, W. Merrill *et al.*, “Rwkv-7 ‘goose’ with expressive dynamic state evolution,” *arXiv preprint arXiv:2503.14456*, 2025.
- [18] S. Grigorescu, B. Trasnea, T. Cocias, and G. Macesanu, “A survey of deep learning techniques for autonomous driving,” *Journal of field robotics*, vol. 37, no. 3, pp. 362–386, 2020.
- [19] L. Bai, J. Cao, M. Zhang, and B. Li, “Collaborative edge intelligence for autonomous vehicles: Opportunities and challenges,” *IEEE Network*, vol. 39, no. 2, pp. 52–60, Mar. 2025.
- [20] M. Janßen, T. Pfandzelter, M. Wang, and D. Bermbach, “Supporting uavs with edge computing: A review of opportunities and challenges,” *arXiv preprint arXiv:2310.11957*, 2023.
- [21] R. Sapkota, K. I. Roumeliotis, and M. Karkee, “Uavs meet agentic ai: A multidomain survey of autonomous aerial intelligence and agentic uavs,” *arXiv preprint arXiv:2506.08045*, 2025.
- [22] M. Groshev, G. Baldoni, L. Cominardi, A. de la Oliva, and R. Gazda, “Edge robotics: Are we ready? an experimental evaluation of current vision and future directions,” *Digital Communications and Networks*, vol. 9, no. 1, pp. 166–174, 2023.
- [23] K. E. Chen, Y. Liang, N. Jha, J. Ichnowski, M. Danielczuk, J. Gonzalez, J. Kubiatowicz, and K. Goldberg, “Fogros: An adaptive framework for automating fog robotics deployment,” in *2021 IEEE 17th International Conference on Automation Science and Engineering (CASE)*, Aug. 2021, pp. 2035–2042.
- [24] A. Bhat, A. Mondal, and A. Tripathy, “Llm agents for internet of things (iot) applications,” in *Submitted to CS598 LLM Agent 2025 Workshop*, 2025.
- [25] Y. Wang, S. Guo, Y. Pan, Z. Su, F. Chen, T. H. Luan, P. Li, J. Kang, and D. Niyato, “Internet of agents: Fundamentals, applications, and challenges,” *IEEE Transactions on Cognitive Communications and Networking*, pp. 1–1, 2025.

- [26] M. B. Yassein, M. Q. Shatnawi, S. Aljwarneh, and R. Al-Hatmi, "Internet of things: Survey and open issues of mqtt protocol," in *2017 International Conference on Engineering & MIS (ICEMIS)*, May. 2017, pp. 1–6.
- [27] J. Zhang, C. Zhao, Y. Zhao, Z. Yu, M. He, and J. Fan, "Mobileexperts: A dynamic tool-enabled agent team in mobile devices," *arXiv preprint arXiv:2407.03913*, 2024.
- [28] A. K. Shukla, "From ai agents to agentic intelligence: A comparative study of autonomy, adaptation, and ethical design," *Adaptation, and Ethical Design*, May. 2025.
- [29] H. Derouiche, Z. Brahmi, and H. Mazeni, "Agentic ai frameworks: Architectures, protocols, and design challenges," *arXiv preprint arXiv:2508.10146*, 2025.
- [30] L. P. Kaelbling, M. L. Littman, and A. W. Moore, "Reinforcement learning: A survey," *Journal of artificial intelligence research*, vol. 4, pp. 237–285, Apr. 1996.
- [31] O. Poquet and M. De Laat, "Developing capabilities: Lifelong learning in the age of ai," *British Journal of Educational Technology*, vol. 52, no. 4, pp. 1695–1708, May. 2021.
- [32] G. Liu, Y. Liu, R. Zhang, H. Du, D. Niyato, Z. Xiong, S. Sun, and A. Jamalipour, "Wireless agentic ai with retrieval-augmented multimodal semantic perception," *arXiv preprint arXiv:2505.23275*, 2025.
- [33] Y. Tian, Z. Zhang, Y. Yang, Z. Chen, Z. Yang, R. Jin, T. Q. S. Quek, and K.-K. Wong, "An edge-cloud collaboration framework for generative ai service provision with synergetic big cloud model and small edge models," *IEEE Network*, vol. 38, no. 5, pp. 37–46, 2024.
- [34] D. M. Onchis, C. Istin, and I.-V. Samuila, "Optimal knowledge distillation through non-heuristic control of dark knowledge," *Mach. Learn. Knowl. Extr.*, vol. 6, pp. 1921–1935, Aug. 2024.
- [35] Y. Wei and Y. Bai, "Dynamic temperature knowledge distillation," *arXiv preprint arXiv:2404.12711*, 2024.
- [36] J. Li, Z. Guo, H. Li, S. Han, J.-w. Baek, M. Yang, R. Yang, and S. Suh, "Rethinking feature-based knowledge distillation for face recognition," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, Jun. 2023, pp. 20 156–20 165.
- [37] S. Lin, H. Xie, B. Wang, K. Yu, X. Chang, X. Liang, and G. Wang, "Knowledge distillation via the target-aware transformer," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, May. 2022, pp. 10 915–10 924.
- [38] S. I. Mirzadeh, M. Farajtabar, A. Li, N. Levine, A. Matsukawa, and H. Ghasemzadeh, "Improved knowledge distillation via teacher assistant," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 04, New York, USA, Feb. 2020, pp. 5191–5198.
- [39] M. Pham, M. Cho, A. Joshi, and C. Hegde, "Revisiting self-distillation," *arXiv preprint arXiv:2206.08491*, 2022.
- [40] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," 2015. [Online]. Available: <https://arxiv.org/abs/1412.6550>
- [41] Z. Yang, Z. Li, A. Zeng, Z. Li, C. Yuan, and Y. Li, "Vitkd: Practical guidelines for vit feature knowledge distillation," *arXiv preprint arXiv:2209.02432*, 2022.
- [42] B. Heo, J. Kim, S. Yun, H. Park, N. Kwak, and J. Y. Choi, "A comprehensive overhaul of feature distillation," in *Proceedings of the IEEE/CVF international conference on computer vision*, Apr. 2019, pp. 1921–1930.
- [43] J. Gou, Y. Chen, B. Yu, J. Liu, L. Du, S. Wan, and Z. Yi, "Reciprocal teacher-student learning via forward and feedback knowledge distillation," *IEEE transactions on multimedia*, vol. 26, pp. 7901–7916, Mar. 2024.
- [44] P. Chen, S. Liu, H. Zhao, and J. Jia, "Distilling knowledge via knowledge review," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, Apr. 2021, pp. 5008–5017.
- [45] Z. Guo, H. Yan, H. Li, and X. Lin, "Class attention transfer based knowledge distillation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, Apr. 2023, pp. 11 868–11 877.
- [46] F. Tung and G. Mori, "Similarity-preserving knowledge distillation," in *Proceedings of the IEEE/CVF international conference on computer vision*, Jul. 2019, pp. 1365–1374.
- [47] W. Park, D. Kim, Y. Lu, and M. Cho, "Relational knowledge distillation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, Apr. 2019, pp. 3967–3976.
- [48] Y. Liu, J. Cao, B. Li, C. Yuan, W. Hu, Y. Li, and Y. Duan, "Knowledge distillation via instance relationship graph," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2019, pp. 7089–7097.
- [49] Y. Tian, D. Krishnan, and P. Isola, "Contrastive representation distillation," *arXiv preprint arXiv:1910.10699*, 2019.
- [50] L. Liu, Q. Huang, S. Lin, H. Xie, B. Wang, X. Chang, and X. Liang, "Exploring inter-channel correlation for diversity-preserved knowledge distillation," in *Proceedings of the IEEE/CVF international conference on computer vision*, Oct. 2021, pp. 8271–8280.
- [51] X. Xin, H. Song, and J. Gou, "A new similarity-based relational knowledge distillation method," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, Apr. 2024, pp. 3535–3539.
- [52] E. Tanghatari, M. Kamal, A. Afzali-Kusha, and M. Pedram, "Federated learning by employing knowledge distillation on edge devices with limited hardware resources," *Neurocomputing*, vol. 531, pp. 87–99, Feb. 2023.
- [53] X. Pang, J. Hu, P. Sun, J. Ren, and Z. Wang, "When federated learning meets knowledge distillation," *IEEE Wireless Communications*, vol. 31, no. 5, pp. 208–214, 2024.
- [54] W. Chen and N. Rojas, "Trakdis: A transformer-based knowledge distillation approach for visual reinforcement learning with application to cloth manipulation," *IEEE Robotics and Automation Letters*, vol. 9, no. 3, pp. 2455–2462, Mar. 2024.
- [55] R. Zhao, Y. Fan, Y. Li, D. Zhang, F. Gao, Z. Gao, and Z. Yang, "Knowledge distillation-enhanced behavior transformer for decision-making of autonomous driving," *Sensors (Basel, Switzerland)*, vol. 25, no. 1, p. 191, Jan. 2025.
- [56] J. Liu, D. Dong, X. Wang, A. Qin, X. Li, P. Valdoriez, D. Dou, and D. Yu, "Large-scale knowledge distillation with elastic heterogeneous computing resources," *Concurrency and Computation: Practice and Experience*, vol. 35, no. 26, p. e7272, 2023.
- [57] D. Deniz, E. Ros, E. M. Ortigosa, and F. Barranco, "Optimized edge-cloud system for activity monitoring using knowledge distillation," *Electronics*, vol. 13, no. 23, p. 4786, Dec. 2024.
- [58] F. O. Catak, M. Kuzlu, E. Catak, U. Cali, and O. Guler, "Defensive distillation-based adversarial attack mitigation method for channel estimation using deep learning models in next-generation wireless networks," *IEEE Access*, vol. 10, pp. 98 191–98 203, 2022.
- [59] K. Kong, W.-J. Song, and M. Min, "Knowledge distillation-aided end-to-end learning for linear precoding in multiuser mimo downlink systems with finite-rate feedback," *IEEE Transactions on Vehicular Technology*, vol. 70, no. 10, pp. 11 095–11 100, Oct. 2021.
- [60] H. Tang, J. Guo, M. Matthaiou, C.-K. Wen, and S. Jin, "Knowledge-distillation-aided lightweight neural network for massive mimo csi feedback," in *2021 IEEE 94th Vehicular Technology Conference (VTC2021-Fall)*. IEEE, Sep. 2021, pp. 1–5.
- [61] M. Gong, S. Wang, S. Gao, J. Yan, and S. Bi, "Robust mimo semantic communication with imperfect csi via knowledge distillation," *arXiv preprint arXiv:2509.04005*, 2025.
- [62] J. Yang, S. Chang, S. Xu, L. Zhou, S. Huang, and Z. Feng, "Knowledge distillation based lightweight deep neural network for automatic modulation classification," in *2023 9th International Conference on Computer and Communications (ICCC)*. IEEE, Dec. 2023, pp. 1972–1977.
- [63] F. Xu, C. Wang, J. Liang, C. Zuo, K. Yue, and W. Li, "A knowledge distillation strategy for enhancing the adversarial robustness of lightweight automatic modulation classification models," *IET Communications*, vol. 18, no. 14, pp. 827–845, Jun. 2024.
- [64] M. Kuzlu, F. O. Catak, U. Cali, E. Catak, and O. Guler, "Adversarial security mitigations of mmwave beamforming prediction models using defensive distillation and adversarial retraining," *International Journal of Information Security*, vol. 22, no. 2, pp. 319–332, 2023.
- [65] Y. M. Park, S. S. Hassan, W. Saad, and C. S. Hong, "Cross-modal knowledge distillation for efficient radar-only beam prediction in mmwave communications," in *2025 IEEE 26th International Workshop on Signal Processing and Artificial Intelligence for Wireless Communications (SPAWC)*. IEEE, Jul. 2025, pp. 1–5.
- [66] J. Yao, S. Cal, and X. Sun, "Resource allocation for federated knowledge distillation learning in internet of drones," *IEEE Internet of Things Journal*, vol. 12, no. 7, pp. 8064–8074, Apr. 2025.
- [67] C. Li, Y. Zhang, L. Yu, and M. Yang, "Efficient vehicle selection and resource allocation for knowledge distillation-based federated learning in uav-assisted vec," *IEEE Transactions on Intelligent Transportation Systems*, vol. 26, no. 5, pp. 6321–6331, May. 2025.
- [68] Z. Chen, Z. Zhang, and Z. Yang, "Big Alai models for 6G wireless networks: Opportunities, challenges, and research directions," *arXiv preprint arXiv:2308.06250*, 2023.
- [69] W. Xu, Z. Yang, D. W. K. Ng, M. Levorato, Y. C. Eldar, and M. Debbah, "Edge learning for B5G networks with distributed signal processing: Semantic communication, edge computing, and wireless sensing," *IEEE J. Sel. Topics Signal Process.*, vol. 17, no. 1, pp. 9–39, Jan. 2023.

- [70] R. Yang, S. Xu, Z. Zhu, C. Li, Y. Huang, and L. Yang, "Knowledge-driven channel estimation for asymmetrical massive mimo systems," *IEEE Transactions on Vehicular Technology*, vol. 74, no. 1, pp. 911–924, Jan. 2025.
- [71] E. Bourtsoulatze, D. B. Kurka, and D. Gündüz, "Deep joint source-channel coding for wireless image transmission," *IEEE Transactions on Cognitive Communications and Networking*, vol. 5, no. 3, pp. 567–579, 2019.
- [72] W. Zhang, H. Zhang, H. Ma, H. Shao, N. Wang, and V. C. Leung, "Predictive and adaptive deep coding for wireless image transmission in semantic communication," *IEEE Transactions on Wireless Communications*, vol. 22, no. 8, pp. 5486–5501, Aug. 2023.
- [73] H. Wu, Y. Shao, C. Bian, K. Mikolajczyk, and D. Gündüz, "Deep joint source-channel coding for adaptive image transmission over mimo channels," *IEEE Transactions on Wireless Communications*, vol. 23, no. 10, pp. 15 002–15 017, Oct. 2024.
- [74] K. Yang, S. Wang, J. Dai, X. Qin, K. Niu, and P. Zhang, "Swinjsc: Taming swin transformer for deep joint source-channel coding," *IEEE Transactions on Cognitive Communications and Networking*, vol. 11, no. 1, pp. 90–104, Feb. 2025.
- [75] Y. M. Park, Y. K. Tun, W. Saad, and C. S. Hong, "Resource-efficient beam prediction in mmwave communications with multimodal realistic simulation framework," *arXiv preprint arXiv:2504.05187*, 2025.
- [76] H. Yang, N. Cheng, R. Sun, W. Quan, R. Chai, K. Aldubaikh, A. Alqasir, and X. Shen, "Knowledge-driven resource allocation for wireless networks: A wmmse unrolled graph neural network approach," *IEEE Internet of Things Journal*, vol. 11, no. 10, pp. 18 902–18 916, May. 2024.
- [77] Y. Gong, Y. Wei, F. R. Yu, and Z. Han, "Slicing-based resource optimization in multi-access edge network using ensemble learning aided ddpg algorithm," *Journal of Communications and Networks*, vol. 25, no. 1, pp. 1–14, Feb. 2023.
- [78] S. Wadhwania, D.-K. Kim, S. Omidshafiei, and J. P. How, "Policy distillation and value matching in multiagent reinforcement learning," in *2019 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, Mar. 2019, pp. 8193–8200.
- [79] R. Li, Z. Zhao, Q. Sun, C. Yang, X. Chen, M. Zhao, H. Zhang *et al.*, "Deep reinforcement learning for resource management in network slicing," *IEEE Access*, vol. 6, pp. 74 429–74 441, May. 2018.
- [80] X. Zhou, X. Zheng, X. Cui, J. Shi, W. Liang, Z. Yan, L. T. Yang, S. Shimizu, and K. I.-K. Wang, "Digital twin enhanced federated reinforcement learning with lightweight knowledge distillation in mobile networks," *IEEE Journal on Selected Areas in Communications*, vol. 41, no. 10, pp. 3191–3211, Oct. 2023.
- [81] D. Ayepah-Mensah, G. Sun, G. Owusu Boateng, and G. Liu, "Federated policy distillation for digital twin-enabled intelligent resource trading in 5g network slicing," *IEEE Transactions on Network and Service Management*, vol. 22, no. 1, pp. 361–379, Feb. 2025.
- [82] S. Tan, R. Caruana, G. Hooker, and Y. Lou, "Distill-and-compare: Auditing black-box models using transparent model distillation," in *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, 2018, pp. 303–310.
- [83] M. Laskin, L. Wang, J. Oh, E. Parisotto, S. Spencer, R. Steigerwald, D. Strouse, S. S. Hansen, A. Filos, E. A. Brooks, M. Gazeau, H. Sahni, S. Singh, and V. Mnih, "In-context reinforcement learning with algorithm distillation," *ArXiv*, vol. abs/2210.14215, 2022.
- [84] A. Bilal, D. Ebert, and B. Lin, "Llms for explainable ai: A comprehensive survey," *arXiv preprint arXiv:2504.00125*, 2025.
- [85] B. A. Zaidi and S. Ansari, "Designing interpretable ai-driven propagation models for advanced wireless networks," in *2024 IEEE International Symposium on Antennas and Propagation and INC/USNC-USRI Radio Science Meeting (AP-S/INC-USNC-USRI)*. IEEE, Jul. 2024, pp. 1847–1848.
- [86] Q. Chen and B. Heydari, "Dynamic resource allocation in systems-of-systems using a heuristic-based interpretable deep reinforcement learning," *Journal of Mechanical Design*, vol. 144, no. 9, p. 091711, Jul. 2022.
- [87] R. Vinuesa and B. Sirmacek, "Interpretable deep-learning models to help achieve the sustainable development goals," *Nature Machine Intelligence*, vol. 3, no. 11, pp. 926–926, Aug. 2021.
- [88] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," *Advances in neural information processing systems*, vol. 35, pp. 24 824–24 837, Jan. 2022.
- [89] C.-Y. Hsieh, C.-L. Li, C.-K. Yeh, H. Nakhost, Y. Fujii, A. Ratner, R. Krishna, C.-Y. Lee, and T. Pfister, "Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes," *arXiv preprint arXiv:2305.02301*, 2023.
- [90] X. Wang, J. Zhu, R. Zhang, L. Feng, D. Niyato, J. Wang, H. Du, S. Mao, and Z. Han, "Chain-of-thought for large language model-empowered wireless communications," *ArXiv*, vol. abs/2505.22320, 2025.
- [91] H. Chen, S. Wu, X. Quan, R. Wang, M. Yan, and J. Zhang, "Mcc-kd: Multi-cot consistent knowledge distillation," in *Conference on Empirical Methods in Natural Language Processing*, 2023.
- [92] T. Wang, K. Sudhir, and D. Hong, "Using advanced llms to enhance smaller llms: An interpretable knowledge distillation approach," *arXiv preprint arXiv:2408.07238*, 2024.
- [93] C. Liu, Y. Zhou, Y. Chen, and S.-H. Yang, "Knowledge distillation-based semantic communications for multiple users," *IEEE Transactions on Wireless Communications*, vol. 23, no. 7, pp. 7000–7012, Nov. 2023.
- [94] Y. Ding, Z. Yang, Q.-V. Pham, Y. Hu, Z. Zhang, and M. Shikh-Bahaei, "Distributed machine learning for uav swarms: Computing, sensing, and semantics," *IEEE Int. Things J.*, vol. 11, no. 5, pp. 7447–7473, 2024.
- [95] K. Xiong, H. Yu, S. Leng, C. Huang, and C. Yuen, "Multi-hop ris-aided learning model sharing for urban air mobility," *IEEE Transactions on Intelligent Transportation Systems*, vol. 26, no. 3, pp. 3947–3959, Mar. 2025.
- [96] C. Huang, Z. Yang, G. C. Alexandropoulos, K. Xiong, L. Wei, C. Yuen, Z. Zhang, and M. Debbah, "Multi-hop ris-empowered terahertz communications: A drl-based hybrid beamforming design," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 6, pp. 1663–1677, 2021.
- [97] T. Furlanello, Z. Lipton, M. Tschanen, L. Itti, and A. Anandkumar, "Born again neural networks," in *International conference on machine learning*. PMLR, May. 2018, pp. 1607–1616.
- [98] S. Yun, J. Park, K. Lee, and J. Shin, "Regularizing class-wise predictions via self-knowledge distillation," 2020. [Online]. Available: <https://arxiv.org/abs/2003.13964>
- [99] Y. Shen, L. Xu, Y. Yang, Y. Li, and Y. Guo, "Self-distillation from the last mini-batch for consistency regularization," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, Mar. 2022, pp. 11 943–11 952.
- [100] J. Zhang, Y. Gao, R. Liu, X. Cheng, H. Zhang, and S. Chen, "Can students beyond the teacher? distilling knowledge from teacher's bias," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 39, no. 21, Philadelphia, PA, USA, Feb. 2025, pp. 22 434–22 442.
- [101] J. Li, Z. Ji, G. Wang, Q. Wang, and F. Gao, "Learning from students: Online contrastive distillation network for general continual learning," in *IJCAI*, Jul. 2022, pp. 3215–3221.
- [102] J. Li, S. Nag, H. Liu, X. Tang, S. Sarwar, L. Cui, H. Gu, S. Wang, Q. He, and J. Tang, "Learning with less: Knowledge distillation from large language models via unlabeled data," *arXiv preprint arXiv:2411.08028*, 2024.
- [103] Z. Wu, Y. Mo, P. Zhou, S. Yuan, and X. Zhu, "Self-training based few-shot node classification by knowledge distillation," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 38, no. 14, Vancouver, BC, Canada, Feb. 2024, pp. 15 988–15 995.
- [104] S. Kang, D. B. Lee, H. Jang, and S. J. Hwang, "Simple semi-supervised knowledge distillation from vision-language models via dual-head optimization," *arXiv preprint arXiv:2505.07675*, 2025.
- [105] S. Zheng, S. Chen, P. Qi, H. Zhou, and X. Yang, "Spectrum sensing based on deep learning classification for cognitive radios," *China Communications*, vol. 17, no. 2, pp. 138–148, 2020.
- [106] M. Zafar, P. J. Wall, S. Bakkali, and R. Haque, "Confidence-based knowledge distillation to reduce training costs and carbon footprint for low-resource neural machine translation," *Applied Sciences*, vol. 15, no. 14, p. 8091, Jul. 2025.
- [107] C. Buciluă, R. Caruana, and A. Niculescu-Mizil, "Model compression," in *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, Aug. 2006, pp. 535–541.
- [108] P. Sarkar and A. Etemad, "Xkd: Cross-modal knowledge distillation with domain alignment for video representation learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 13, Vancouver, BC, Canada, Feb. 2024, pp. 14 875–14 885.
- [109] Q. Feng, W. Li, T. Lin, and X. Chen, "Align-kd: Distilling cross-modal alignment knowledge for mobile vision-language large model enhancement," in *Proceedings of the Computer Vision and Pattern Recognition Conference*, Jun. 2025, pp. 4178–4188.
- [110] F. Xu, P. Fu, Q. Huang, B. Zou, A. Aw, and M. Wang, "Leveraging contrastive learning and knowledge distillation for incomplete modality rumor detection," in *Findings of the Association for Computational Linguistics: EMNLP 2023*, Dec. 2023, pp. 13 492–13 503.

- [111] W.-H. Li and H. Bilen, "Knowledge distillation for multi-task learning," in *European Conference on Computer Vision*. Springer, Jul. 2020, pp. 163–176.
- [112] G. M. Jacob, V. Agarwal, and B. Stenger, "Online knowledge distillation for multi-task learning," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, Jan. 2023, pp. 2359–2368.
- [113] S. Wang and C. Liu, "Automatic modulation classification with neural networks via knowledge distillation," *Electronics*, vol. 11, no. 19, p. 3018, Sep. 2022.
- [114] Q. Li, B. He, and D. X. Song, "Model-agnostic round-optimal federated learning via knowledge transfer," *ArXiv*, vol. abs/2010.01017, 2020.
- [115] X. Gong, A. Sharma, S. Karamam, Z. Wu, T. Chen, D. Doermann, and A. Innanje, "Preserving privacy in federated learning with ensemble cross-domain knowledge distillation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 11, Jun. 2022, pp. 11 891–11 899.
- [116] Y. Yu, S. Zhu, and J. Hu, "Logit poisoning attack in distillation-based federated learning and its countermeasures," *arXiv preprint arXiv:2401.17746*, 2024.
- [117] A. Mora, I. Tenison, P. Bellavista, and I. Rish, "Knowledge distillation for federated learning: a practical guide," *arXiv preprint arXiv:2211.04742*, 2022.
- [118] Z. Lu, J. Wang, and C. Jiang, "Data-free knowledge filtering and distillation in federated learning," *IEEE Transactions on Big Data*, vol. 11, no. 3, pp. 1128–1143, Jun. 2025.
- [119] Z. Zhu, J. Hong, and J. Zhou, "Data-free knowledge distillation for heterogeneous federated learning," in *International conference on machine learning*. PMLR, May. 2021, pp. 12 878–12 889.
- [120] E. Jeong and M. Kountouris, "Personalized decentralized federated learning with knowledge distillation," in *ICC 2023-IEEE International Conference on Communications*. IEEE, Feb. 2023, pp. 1982–1987.
- [121] S. Mukherjee, A. Simonetto, and H. Jamali-Rad, "Mapl: Model agnostic peer-to-peer learning," *arXiv preprint arXiv:2403.19792*, 2024.
- [122] Q. Jin and H. Ochiai, "Decentralized p2p federated learning on ad-hoc like networks with non-iid dataset," in *2022 18th International Conference on Wireless and Mobile Computing, Networking and Communications (WiMob)*. IEEE, Oct. 2022, pp. 387–393.
- [123] Z. Zhou, F. Sun, X. Chen, D. Zhang, T. Han, and P. Lan, "A decentralized federated learning based on node selection and knowledge distillation," *Mathematics*, vol. 11, no. 14, p. 3162, Jul. 2023.
- [124] H. Hu, A. N. Kothari, and A. Banerjee, "A novel algorithm for personalized federated learning: Knowledge distillation with weighted combination loss," *Algorithms*, vol. 18, no. 5, p. 274, Apr. 2025.
- [125] M. A. Mohsin, M. Umer, A. Bilal, M. I. Qadir, M. A. Jamshed, D. F. Hougen, and J. M. Cioffi, "Channel prediction under network distribution shift using continual learning-based loss regularization," *arXiv preprint arXiv:2509.15192*, 2025.
- [126] M. Yan, S. Li, C. A. Chan, Y. Shen, and Y. Yu, "Mobility prediction using a weighted markov model based on mobile user classification," *Sensors*, vol. 21, no. 5, p. 1740, Mar. 2021.
- [127] Z. Wang, Y. Zhou, Y. Shi, and W. Zhuang, "Interference management for over-the-air federated learning in multi-cell wireless networks," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 8, pp. 2361–2377, Jun. 2022.
- [128] M. Akroot, A. Feriani, F. Bellili, A. Mezghani, and E. Hossain, "Continual learning-based mimo channel estimation: A benchmarking study," in *ICC 2023-IEEE International Conference on Communications*. IEEE, 2023, pp. 2631–2636.
- [129] H. Chen, L. Zeng, S. Yu, and X. Chen, "Knowledge distillation for mobile edge computation offloading," *arXiv preprint arXiv:2004.04366*, 2020.
- [130] X. Li, S. Wang, J. Sun, and Z. Xu, "Memory efficient data-free distillation for continual learning," *Pattern Recognition*, vol. 144, p. 109875, Aug. 2023.
- [131] Q. Pan, S. Sun, Z. Wu, Y. Wang, M. Liu, B. Gao, and J. Wang, "Fecdcache 2.0: Federated edge learning with knowledge caching and dataset distillation," *arXiv preprint arXiv:2405.13378*, 2024.
- [132] R. Zhang, R. Zhang, Y. Lu, W. Chen, B. Ai, and D. Niyato, "Mamba for wireless communications and networking: Principles and opportunities," *arXiv preprint arXiv:2508.00403*, 2025.
- [133] A. Mehrabian and V. W. Wong, "A-gamba: An adaptive graph-mamba model for traffic prediction in wireless cellular networks," *IEEE Wireless Communications Letters*, vol. 14, no. 6, pp. 1801–1805, Jun. 2025.
- [134] Y. Huang, J. Liu, X. Shi, S. Zhao, T. Mi, and R. C. Qiu, "Sensemamba: A general lightweight state-space model for wireless human sensing," *IEEE Sensors Journal*, vol. 25, no. 8, pp. 13 859–13 870, Apr. 2025.
- [135] S. Xie, Y. Li, Y. Ma, and Y. Wu, "Autogmm-rwkv: A detecting scheme based on attention mechanisms against selective forwarding attacks in wireless sensor networks," *IEEE Internet of Things Journal*, vol. 12, no. 4, pp. 4403–4419, Feb. 2025.
- [136] Q. Fu, H. Yuan, Y. Hou, and X. Shen, "Linear attention based channel estimation scheme for v2x communications," in *2024 8th International Conference on Communication and Information Systems (ICCIS)*. IEEE, Oct. 2024, pp. 96–101.
- [137] M. A. Mohsin, M. Jazib, Z. Alam, M. F. Khan, M. Saad, and M. A. Jamshed, "Vision transformer based semantic communications for next generation wireless networks," in *2025 IEEE International Conference on Communications Workshops (ICC Workshops)*, Jun. 2025, pp. 1803–1808.
- [138] F. Jabbarvaziri and L. Lampe, "Attention-based deep learning for hybrid beamforming in ofdm systems with phase noise," *IEEE Transactions on Wireless Communications*, vol. 24, no. 9, pp. 7733–7746, Sep. 2025.
- [139] S. Troia, R. Alvizu, Y. Zhou, G. Maier, and A. Pattavina, "Deep learning-based traffic prediction for network optimization," in *2018 20th International Conference on Transparent Optical Networks (ICTON)*. IEEE, Jul. 2018, pp. 1–4.
- [140] J. P. Lemayian and J. M. Hamamreh, "Recurrent neural network-based channel prediction in mmimo for enhanced performance in future wireless communication," in *2020 International Conference on UK-China Emerging Technologies (UCET)*. IEEE, Aug. 2020, pp. 1–4.
- [141] A. Gu and T. Dao, "Mamba: Linear-time sequence modeling with selective state spaces," *arXiv preprint arXiv:2312.00752*, 2023.
- [142] A. Gu, I. Johnson, K. Goel, K. Saab, T. Dao, A. Rudra, and C. Ré, "Combining recurrent, convolutional, and continuous-time models with linear state space layers," *Advances in neural information processing systems*, vol. 34, pp. 572–585, Oct. 2021.
- [143] A. Gu, K. Goel, and C. Ré, "Efficiently modeling long sequences with structured state spaces," *arXiv preprint arXiv:2111.00396*, 2021.
- [144] J. Liu, Y. Huang, X. Shi, X. Ren, T. Mi, and R. C. Qiu, "Tf-mamba: A lightweight state-space model for wi-fi-based human activity recognition," *IEEE Sensors Journal*, vol. 25, pp. 23 184–23 194, July. 2025.
- [145] W. Yang, X. Zhou, J. Guan, H. Du, and T. Bai, "Gram-mamba: Holistic feature alignment for wireless perception with adaptive low-rank compensation," *arXiv preprint arXiv:2507.13803*, 2025.
- [146] W. Choe, Y. Ji, and F. X. Lin, "Rwkv-lite: Deeply compressed rwkv for resource-constrained devices," *arXiv preprint arXiv:2412.10856*, 2024.
- [147] A. Bick, K. Li, E. Xing, J. Z. Kolter, and A. Gu, "Transformers to ssms: Distilling quadratic knowledge to subquadratic models," *Advances in Neural Information Processing Systems*, vol. 37, pp. 31 788–31 812, Aug. 2024.
- [148] A. Bick, T. Katsch, N. Sohoni, A. Desai, and A. Gu, "Llamba: Scaling distilled recurrent models for efficient language processing," *arXiv preprint arXiv:2502.14458*, 2025.
- [149] L. Galke, I. Cuber, C. Meyer, H. F. Nölscher, A. Sonderecker, and A. Scherp, "General cross-architecture distillation of pretrained language models into matrix embeddings," in *2022 International Joint Conference on Neural Networks (IJCNN)*. IEEE, Jul. 2022, pp. 1–10.
- [150] Z. Hao, J. Guo, K. Han, Y. Tang, H. Hu, Y. Wang, and C. Xu, "One-for-all: Bridge the gap between heterogeneous architectures in knowledge distillation," *Advances in Neural Information Processing Systems*, vol. 36, pp. 79 570–79 582, Oct. 2023.
- [151] Y. Yao, Y. Hong, D. Liu, L. Mai, F. Liu, and J. Luo, "Diffusion transformer-to-mamba distillation for high-resolution image generation," *arXiv preprint arXiv:2506.18999*, 2025.
- [152] Y. Liu, J. Cao, B. Li, W. Hu, J. Ding, and L. Li, "Cross-architecture knowledge distillation," in *Proceedings of the Asian conference on computer vision*, Jul. 2022, pp. 3396–3411.
- [153] B. Yilmaz and A. Aiyengar, "Cross-architecture knowledge distillation (kd) for retinal fundus image anomaly detection on nvidia jetson nano," *arXiv preprint arXiv:2506.18220*, 2025.
- [154] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter," *arXiv preprint arXiv:1910.01108*, 2019.
- [155] X. Jiao, Y. Yin, L. Shang, X. Jiang, X. Chen, L. Li, F. Wang, and Q. Liu, "Tinybert: Distilling bert for natural language understanding," in *Findings*, Sep. 2019.
- [156] J. Tang, R. Shivanna, Z. Zhao, D. Lin, A. Singh, E. H. Chi, and S. Jain, "Understanding and improving knowledge distillation," *arXiv preprint arXiv:2002.03532*, 2020.
- [157] G. Ji and Z. Zhu, "Knowledge distillation in wide neural networks: Risk bound, data efficiency and imperfect teacher," *Advances in Neural Information Processing Systems*, vol. 33, pp. 20 823–20 833, Oct. 2020.

- [158] A. Albaseer and M. Abdallah, "Tailoring semantic communication at network edge: A novel approach using dynamic knowledge distillation," in *ICC 2024-IEEE International Conference on Communications*. IEEE, Jun. 2024, pp. 1455–1460.
- [159] Z. Wu, S. Sun, Y. Wang, M. Liu, X. Jiang, and R. Li, "Survey of knowledge distillation in federated edge learning," *arXiv preprint arXiv:2301.05849*, 2023.
- [160] L. Sun, Z. Liu, L. Wan, Y. Lin, L. Lin, J. Wang, and M. Gen, "Cooperative knowledge-distillation-based tiny dnn for uav-assisted mobile-edge network," *IEEE Internet of Things Journal*, vol. 11, no. 18, pp. 30 204–30 216, Sep. 2024.
- [161] Z. Yang, W. Xu, and M. Shikh-Bahaei, "Energy efficient UAV communication with energy harvesting," *IEEE Trans. Veh. Technol.*, vol. 69, no. 2, pp. 1913–1927, 2019.
- [162] Y. R. Wei, J. Sathiyamoorthi, M. Rahman, and P. Janik, "Mobile robots design for industrial applications in private 5g networks: Essential factors to consider," in *2023 53rd European Microwave Conference (EuMC)*. IEEE, Sep. 2023, pp. 504–507.
- [163] J. Moon, D. Hong, J. Kim, S. Kim, S. Woo, H. Choi, and C. Moon, "Enhancing autonomous driving robot systems with edge computing and ldm platforms," *Electronics*, vol. 13, no. 14, p. 2740, Jul. 2024.
- [164] Z. Zhao, J. Lyu, Y. Chu, K. Liu, D. Cao, C. Wu, L. Qin, and S. Qin, "Toward generalizable robot vision guidance in real-world operational manufacturing factories: A semi-supervised knowledge distillation approach," *Robotics and Computer-Integrated Manufacturing*, vol. 86, p. 102639, Apr. 2024.
- [165] Y. Chebotar, A. Handa, V. Makovychuk, M. Macklin, J. Issac, N. Ratliff, and D. Fox, "Closing the sim-to-real loop: Adapting simulation randomization with real world experience," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 8973–8979.
- [166] E. Parisotto, J. L. Ba, and R. Salakhutdinov, "Actor-mimic: Deep multitask and transfer reinforcement learning," *arXiv preprint arXiv:1511.06342*, 2015.
- [167] R. Wang, M. Lyu, and J. Zhang, "A multi-robot collaborative exploration method based on deep reinforcement learning and knowledge distillation," *Mathematics*, vol. 13, no. 1, p. 173, Jan. 2025.
- [168] S. Al-Sarawi, M. Anbar, K. Alieyan, and M. Alzubaidi, "Internet of things (iot) communication protocols," in *2017 8th International conference on information technology (ICIT)*. IEEE, May. 2017, pp. 685–690.
- [169] Q. Li and Y. Xie, "From glue-code to protocols: A critical analysis of a2a and mcp integration for scalable agent systems," *arXiv preprint arXiv:2505.03864*, 2025.
- [170] E. King, H. Yu, S. Lee, and C. Julien, "Sasha: creative goal-oriented reasoning in smart homes with large language models," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 8, no. 1, pp. 1–38, 2024.
- [171] D. Rivkin, F. Hogan, A. Feriani, A. Konar, A. Sigal, S. Liu, and G. Dudek, "Sage: Smart home agent with grounded execution," *arXiv preprint arXiv:2311.00772*, 2023.
- [172] J. Feng, T. Liu, Y. Du, S. Guo, Y. Lin, and Y. Li, "Citygpt: Empowering urban spatial cognition of large language models," in *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, 2025, pp. 591–602.
- [173] N. Zhong, Y. Wang, R. Xiong, Y. Zheng, Y. Li, M. Ouyang, D. Shen, and X. Zhu, "Casit: Collective intelligent agent system for internet of things," *IEEE Internet of Things Journal*, vol. 11, no. 11, pp. 19 646–19 656, Jun. 2024.
- [174] Z. Zhao, Z. Yang, C. Huang, L. Wei, Q. Yang, C. Zhong, W. Xu, and Z. Zhang, "A joint communication and computation design for distributed risc assisted probabilistic semantic communication in iiot," *IEEE Int. Things J.*, pp. 1–1, 2024.
- [175] J. He, Z. Tang, X. Fu, S. Leng, F. Wu, K. Huang, J. Huang, J. Zhang, Y. Zhang, A. Radford *et al.*, "Cooperative connected autonomous vehicles (cav): Research, applications and challenges," in *2019 IEEE 27th International Conference on Network Protocols (ICNP)*. IEEE, Oct. 2019, pp. 1–6.
- [176] Z. Khan, A. Koubaa, and H. Farman, "Smart route: Internet-of-vehicles (iov)-based congestion detection and avoidance (iov-based cda) using rerouting planning," *Applied Sciences*, vol. 10, no. 13, p. 4541, Jun. 2020.
- [177] S. A. Cañas Ordóñez, J. Samanta, A. L. Suárez-Cetrulo, and R. S. Carbajo, "Intelligent edge computing and machine learning: A survey of optimization and applications," *Future Internet*, vol. 17, no. 9, p. 417, Sep. 2025.
- [178] H. Lee and S. Kim, "Explaining neural networks using attentive knowledge distillation," *Sensors*, vol. 21, no. 4, p. 1280, Feb. 2021.
- [179] A. Parchami-Araghi, M. Böhle, S. Rao, and B. Schiele, "Good teachers explain: Explanation-enhanced knowledge distillation," in *European Conference on Computer Vision*. Springer, Feb. 2024, pp. 293–310.
- [180] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, Jun. 2009, pp. 248–255.
- [181] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, "Squad: 100,000+ questions for machine comprehension of text," in *Conference on Empirical Methods in Natural Language Processing*, Jun. 2016.
- [182] H. Nguyen, Z. He, S. A. Gandre, U. Pasupulety, S. K. Shivakumar, and K. Lerman, "Smoothing out hallucinations: Mitigating llm hallucination with smoothed knowledge distillation," *arXiv preprint arXiv:2502.11306*, 2025.
- [183] Y. Tang, A. Zhang, Z. Wu, B. Gao, T. Wen, Y. Wang, and S. Sun, "Peak-controlled logits poisoning attack in federated distillation," *Discover Computing*, vol. 28, 2024.
- [184] S. Li, Y. Wang, Y. Li, and Y.-a. Tan, "I-leaks: Membership inference attacks with logits," *arXiv preprint arXiv:2205.06469*, 2022.
- [185] C. Xu, S. Liu, Z. Yang, Y. Huang, and K.-K. Wong, "Learning rate optimization for federated learning exploiting over-the-air computation," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 12, pp. 3742–3756, 2021.
- [186] G. Yenduri, R. Murugan, P. Kumar Reddy Maddikunta, S. Bhattacharya, D. Sudheer, and B. Bhushan Savarala, "Artificial general intelligence: Advancements, challenges, and future directions in agi research," *IEEE Access*, vol. 13, pp. 134 325–134 356, 2025.