

Resource-efficient Algorithms and Systems of Foundation Models: A Survey

MENGWEI XU*, Beijing University of Posts and Telecommunications, China

DONGQI CAI*, Beijing University of Posts and Telecommunications, China

WANGSONG YIN*, Peking University, China

SHANGGUANG WANG, Beijing University of Posts and Telecommunications, China

XIN JIN, Peking University, China

XUANZHE LIU, Peking University, China

Large foundation models, including large language models, vision transformers, diffusion, and LLM-based multimodal models, are revolutionizing the entire machine learning lifecycle, from training to deployment. However, the substantial advancements in versatility and performance these models offer come at a significant cost in terms of hardware resources. To support the growth of these large models in a scalable and environmentally sustainable way, there has been a considerable focus on developing resource-efficient strategies. This survey delves into the critical importance of such research, examining both algorithmic and systemic aspects. It offers a comprehensive analysis and valuable insights gleaned from existing literature, encompassing a broad array of topics from cutting-edge model architectures and training/serving algorithms to practical system designs and implementations. The goal of this survey is to provide an overarching understanding of how current approaches are tackling the resource challenges posed by large foundation models and to potentially inspire future breakthroughs in this field.

CCS Concepts: • **Computing methodologies** → **Natural language processing**; **Computer vision**.

Additional Key Words and Phrases: Resource-efficiency, foundation models, algorithm and system optimization

ACM Reference Format:

Mengwei Xu*, Dongqi Cai*, Wangsong Yin*, Shangguang Wang, Xin Jin, and Xuanzhe Liu. 2024. Resource-efficient Algorithms and Systems of Foundation Models: A Survey. 1, 1 (November 2024), 35 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

In the rapidly evolving field of artificial intelligence (AI), a paradigm shift is underway. We are witnessing the transition from specialized, fragmented deep learning models to versatile, one-size-fits-all foundation models (FMs). These advanced AI systems are capable of operating in an open-world context, interacting with open vocabularies and image pixels for unseen AI tasks, i.e., zero-shot abilities. They are exemplified by (1) Large Language Models (LLMs) such as GPTs [26] that can ingest almost every NLP task in the form as a prompt; (2) Vision Transformers Models (ViTs) such as Masked Autoencoder [96] that can handle various downstream vision tasks; (3) Latent Diffusion Models (LDMs) such as Stable Diffusion [220] that generate high-quality images with arbitrary text-based prompts; (4) Multimodal

Mengwei Xu, Dongqi Cai, and Wangsong Yin contributed equally to this survey.

Authors' addresses: Mengwei Xu*, mwx@bupt.edu.cn, Beijing University of Posts and Telecommunications, China; Dongqi Cai*, Beijing University of Posts and Telecommunications, China; Wangsong Yin*, Peking University, China; Shangguang Wang, Beijing University of Posts and Telecommunications, China; Xin Jin, Peking University, China; Xuanzhe Liu, Peking University, China.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

models such as CLIP [214] and ImageBind [79] that map different modal data into the same latent space and are widely used as backbone for cross-modality tasks like image retrieval/search and visual-question answering. Such flexibility and generality mark a significant departure from the earlier era of AI, setting a new standard for how AI interfaces with the world.

The success of these FMs is deeply rooted in their scalability: unlike their predecessors, these models' accuracy and generalization ability can continuously expand with more data or parameters, without altering the underlying algorithms and architectures. An impressive evidence is the scaling law [119]: it describes how the performance of transformer-based models can predictably improve with more model size and data volume; until today, the scaling law stands still. This scalability is not just a matter of model size; it extends to their ability to tackle increasingly complex tasks, making them a cornerstone in the journey towards artificial general intelligence (AGI).

However, the scalability comes at a cost of huge resource demand. Foundation models, by their very nature, are resource-hungry for training and deployment. These resources encompass not only the computing processors like GPUs and TPUs, but also the memory, energy, and network bandwidth. For example, the pre-training of LLaMa-2-70B takes $1.7 \times$ millions of GPU hours and consumes 2.5×10^{12} Joules of energy. The estimated total emissions were 291 tons of CO₂ equivalent. Beyond training, the data processing, experimentation, and inference stages consume comparable or even more electricity according to Meta AI [267]. A recent analysis [50] reveals that, to satisfy the continuation of the current trends in AI capacity and adoption, NVIDIA needs to ship 1.5 million AI server units per year by 2027. These servers, running at full capacity, would consume at least 85.4 terawatt-hours of electricity annually – more than what many countries like New Zealand and Austria use in a whole year. Since FMs proceed growth in size and complexity, their resource requirements escalate, posing a significant challenge in their development and deployment.

The huge resource footprint of large FM also hinders its democratization. Till the end of 2023, there are only a few major players capable of training and deploying the state-of-the-art FMs, who thereby have powerful control over the public and can potentially manipulate them in a way they prefer. The models are served on clouds instead of devices as many lightweight DNNs do [277, 303]; it makes data privacy preservation almost impossible. Though recently, smartphone vendors have been boasting about running large FMs locally and some pioneering engines are developed for on-device LLMs [6, 7, 78], the models demonstrated are limited to relatively small scale (e.g., <10B) and have not yet seen real-world deployment.

Thereby, a significant amount of research has been dedicated to enhance the efficiency of these FMs. These efforts span a wide range of approaches, from optimizing algorithms to system-level innovations, focusing on reducing the resource footprint of these models without compromising their performance. This survey aims to delve into these research efforts, exploring the diverse strategies employed to make FMs more resource-efficient. We will examine advancements in algorithmic efficiency, system optimizations, and the development of novel architectures that are less resource-intensive. The survey also spans from clouds to edge and devices, where the large FMs gain dramatic attentions as well. Through this exploration, we aim to provide a comprehensive understanding of the current state and future directions of resource-efficient algorithms and systems in the realm of FMs.

Scope and rationales. (i) We survey only algorithm and system innovations; we exclude a huge body of work at hardware design; (ii) The definition of resource in this survey is limited to mainly physical ones, including computing, memory, storage, bandwidth, etc; we exclude training data (labels) and privacy that can also be regarded as resources; (iii) We mainly survey papers published on top-tier CS conferences, i.e., those included in CSRankings. We also manually pick related and potentially high-impact papers from arXiv. (iv) We mainly survey papers published after the year of 2020, since the innovation of AI is going fast with old knowledge and methods being overturned frequently.

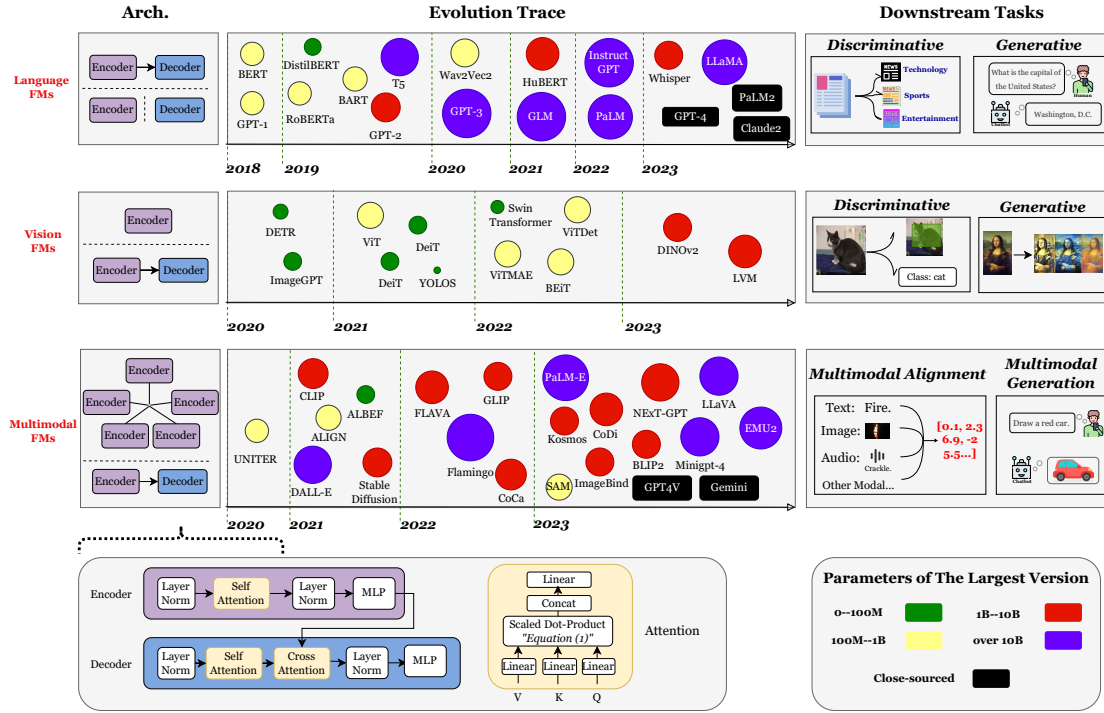


Fig. 1. The evolutionary trace of foundation models.

Organization. Section 2 overviews the classical foundation models and their runtime cost. Section 3 investigates the architectural innovations that revise or replace the existing foundation model architectures. Section 4 and Section 5 examine the algorithm-level and system-level literature towards more resource-efficient foundation models. Section 6 concludes the survey and presents potential future directions.

Comparison to relevant surveys. Concurrent to this work, there are a few (not yet peer-reviewed) surveys about efficient large language models, spanning from compression [318], algorithms [56], system-algorithm [184, 248], and hardware [125]. As comparison, this work is the first comprehensive survey towards resource-efficient foundation models, including not only large language models, but also multimodal ones that are equally important such as diffusion and ViT models. An extended version of this survey is available at [279].

2 FOUNDATION MODEL OVERVIEW

Fig. 1 illustrates the evolutionary trace of popular foundation models (FMs) up to Jan. 2024. In general, there are three types of FMs: language-based, vision-based, and multimodal FMs.

Language FMs typically employ attention-based transformer architecture [246]. The process initiates by converting input words into high-dimensional vectors through an embedding layer. During processing, attention mechanisms assign varying weights to different segments of these input vectors. Following attention, layer normalization is applied to the output, ensuring stabilization and standardization of the activations. Subsequently, each position-wise vector undergoes transformation through a feedforward network, introducing non-linearity and enabling the model to capture complex data patterns. Through multiple layers that incorporate these components, the Transformer learns hierarchical

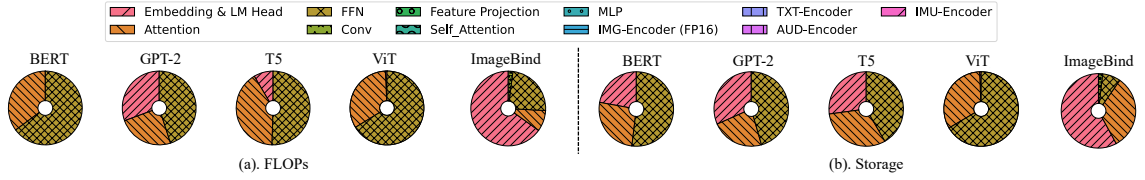


Fig. 2. Empirical computation and storage comparison across different FMs.

representations of the input data. In the final stage, the output from the last Transformer layer is directed into a linear layer, culminating in the final prediction.

Vision FMs In this paper, we use the term “Vision FMs” to refer to the foundation models that only involve the pure vision modality in their main pipeline. Vision FMs (e.g., SAM, seggpt [127, 256]) typically employ Vision Transformer (ViT) architecture [58], a transformer-based visual information processing block. As such, efficient vision FMs (as well as those multimodal ones that rely on ViT) often benefit from the efficient ViT designs. Given an input image, ViT firstly splits image into fixed-size patches (i.e., tokens) by a convolutional embedding layer. For instance, a standard size RGB image input (i.e., $3 \times 224 \times 224$) will be splitted to 14×14 patches with 16×16 pixels. This embedding overhead is almost negligible compared to the following compute-intensive transformer encoder (e.g., less than 5%). Besides, an extra learnable classification token ([CLS]) is added to the token sequence in order to perform classification. After that, positional embeddings are added into each token, and tokens are fed to a standard Transformer encoder. Depending on the specific downstream tasks, the hidden states generated by the Transformer encoder are finally fed into different heads, such as classification, detection, etc.

Multimodal FMs are used in two specific goals: encoding input data in different modalities into the same latent space; or generating output data in different modalities. The two lines of research have convergence, e.g., multimodal-to-multimodal (or even any-to-any) generation. To ingest and align multimodal input data, existing model architectures like CLIP [214] typically consist of multiple transformer encoders, with each modality having its own set of transformer encoders. Notably, these encoders are generally trained from scratch, utilizing paired data with the aligned modalities and current modality. To generate multimodal data, FMs can either (i) reuse the LLM to generate text; (ii) or diffusion models [220] to generate high-quality image pixels. The diffusion module primarily consists of two components: an image encoder/decoder and a denoising network. There are also variants of diffusion model that replace the convolution with Transformer, e.g., DiTs [206], as well as foundation models [293] that involve richer modalities like IMU or audio. Yet such modalities are mainly embedded with only a dedicated embedding layer and reuses the same transformer architecture. Thereby, we do not discuss these models in isolation.

Applications of FMs: In real-world applications, language foundation models like GPT-4 [198] have transformed tasks such as content generation [103], code assistance [144], and natural language understanding [317] across multiple industries. These advancements enable chatbots and personal agents to better understand user queries and provide more meaningful responses. In the case of vision FMs, models such as SAM [127] are widely applied in medical imaging, allowing healthcare professionals to accurately segment and analyze images with minimal manual intervention, significantly improving diagnostic accuracy [11, 180]. Multimodal FMs, including CLIP [214] and Stable Diffusion [220], are transforming the creative industries by enabling artists to generate artwork from simple text prompts [31], thereby expanding creative possibilities while reducing manual effort.

Cost Analysis of transformer Since most FMs are based on transformer architecture, we briefly analyze the resource cost of it. The attention mechanism in large FMs faces significant computational bottlenecks primarily due to its quadratic complexity. This complexity stems from calculating attention scores for every pair of positions within the

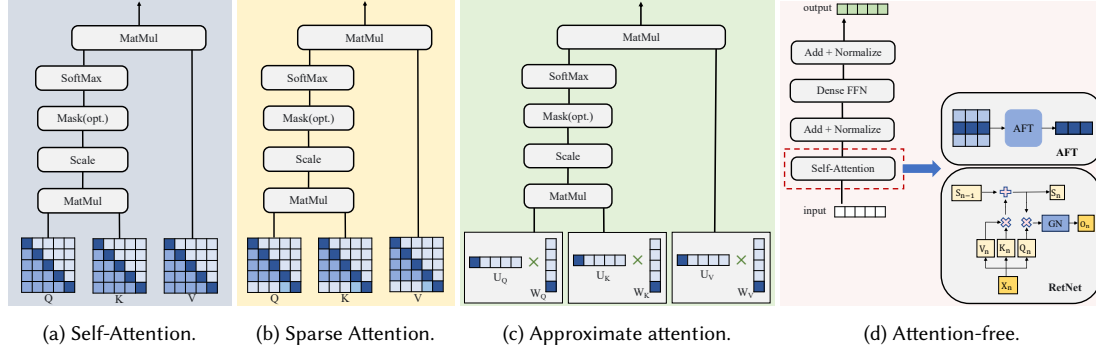


Fig. 3. Illustrations of efficient attentions architectures.

input sequence, posing challenges in managing long sequences and impacting both training and inference efficiency. Additionally, beyond the attention mechanism, the computation complexity of the FFN scales linearly with input length but quadratically with the model's dimension. An increase in the length of the input sequence causes a substantial rise in computational demand, attributable to the quadratic nature of the attention mechanism. In quantitative terms, the computation complexity of attention is $O(T^2D)$, while that of FFN is $O(TD^2)$, where T represents the sequence length and D the hidden state dimension of the model [161]. The decoder's attention mechanism, similar to that in the encoder, also experiences quadratic scaling with token length. This aspect becomes particularly significant in autoregressive decoding tasks, where each token's generation depends on the preceding ones, intensifying computational requirements. The implementation of a KV cache in the decoder can substantially mitigate computational costs by reusing key and value vectors across various positions [134].

We empirically analyze the resource costs of different FMs by comparing their demands in terms of FLOPs and storage, as shown in Fig. 2. For language models such as BERT, GPT-2, and T5, the embedding layer and LM head contribute significantly to storage. However, these components require minimal computational FLOPs. The FFN layer is the most computationally intensive component. Similar trends are observed in vision and speech models, such as Wav2Vec2 and ViTs, where convolution is not dominant. Instead, MLP and self-attention layers consume the most resources. In multimodal models like ImageBind, the IMG-Encoder is the most resource-demanding, while other encoders require significantly fewer resources.

3 RESOURCE-EFFICIENT ARCHITECTURES

3.1 Efficient Attention

As summarized in Fig. 3, numerous efforts have been invested to mitigate the huge resource cost of attention-based transformer architecture. The time and space complexity comparison is shown in Table 1.

Model	Time	Space	Model	Time	Space
Transformer [246]	$O(T^2d)$	$O(T^2 + Td)$	AFT [299]	$O(T^2d)$	$O(Td)$
Reformer [128]	$O(T \log Td)$	$O(T \log T + Td)$	Hyena [210]	$O(T \log Td)$	$O(Td)$
SSM [84]	$O(T \log Td)$	$O(Td)$	Linear Transformers [120]	$O(Td^2)$	$O(Td + d^2)$
RetNet [231]	$O(Td)$	$O(Td)$	RWKV [207]	$O(Td)$	$O(d)$

Table 1. The time and space complexity comparison, where T represents sequence length, and d represents hidden dimension.

3.1.1 Sparse Attention. Motivated by graph sparsification, sparse attention aims to build a sparse attention matrix. This approach aims to retain the empirical advantages of a fully quadratic self-attention scheme while employing a reduced

number of inner products. For instance, Longformer [68], ETC [161], and BIGBIRD [295] decompose conventional attention into local windowed attention and task-specific global attention, effectively reducing self-attention complexity to linear. HEPOS [104] introduces head-wise positional strides, allowing each attention head to concentrate on a specific subset of the input sequence. MATE [60] transforms attention into a multi-view format, efficiently addressing either rows or columns in a table. TDANet [145] emulates the human brain’s top-down attention mechanism to selectively focus on the most relevant information, thereby enhancing speech separation efficiency.

3.1.2 Approximate Attention. Approximate attention mainly includes low-rank approximations of the self-attention matrix and innovative reformulations of the self-attention. Linformer [252] effectively decomposes the attention matrix into a low-rank matrix. It involves projecting the length dimensions of keys and values into a lower-dimensional space, resulting in a significant reduction in memory complexity. Reformer [128] utilizes locality-sensitive hashing to replace the conventional dot-product attention. Katharopoulos et al. [120] introduced a kernel-based alternative to self-attention, leveraging the associative property of matrix multiplication for computing self-attention weights. Polysketchformer [118] employs polynomial functions and sketching techniques to approximate softmax attention outputs. Mega [179], featuring a single-head gated attention mechanism, incorporates exponential moving average. Deformable Attention [270] proposes a data-aware, deformable attention mechanism, contributing to improved performance within the ViT architecture. CrossViT [37] introduces linear cross-attention, empowering the ViT architecture to efficiently handle variably-sized input tokens while mitigating computational costs.

3.1.3 Attention-Free Approaches. Despite the dominance of attention-based transformer architectures in large FMs, several works have put forth innovative architectures that hold the potential to replace the traditional transformer model. For instance, Hyena [210] introduces an architecture that interleaves implicitly parametrized long convolutions with data-controlled gating. This design provides a subquadratic alternative to attention in large-scale language models, thereby enhancing efficiency in processing long sequences. Another notable trend is the substitution of the attention mechanism with state space models (SSMs), as explored in [46, 84, 199]. Mamba [83] seamlessly integrates selective SSMs into a streamlined neural network architecture, eliminating attention and MLP blocks. This model achieves a notable $5\times$ speed increase over traditional transformers and exhibits linear scaling with sequence length. Recurrent-Style Transformers [27, 28] adopts an recurrent neural network-based architecture, replacing attention with an RNN to achieve linear complexity. RWKV [207] combines the efficient parallelizable training of Transformers with the effective inference capabilities of RNNs. RetNet [231] introduces an architecture that replaces multi-head attention with a multi-scale retention mechanism. During training, RetNet demonstrates a 25-50% memory saving and a $7\times$ acceleration compared to the standard Transformer.

3.2 Dynamic Neural Network

3.2.1 Mixture of Experts. Mixture-of-Experts (MoE), as illustrated in Fig. 4(b), represents an efficient and sparse approach for training and deploying large FMs with extensive parameter sets. This model utilizes routed sparse parameters during inference. Switch Transformer [65] introduces a switch routing algorithm, leading to models with improved efficiency and reduced computational and communication costs. Switch Transformer demonstrates the scalability and effectiveness of MoE framework by managing up to one trillion parameters, with as many as 2,048 experts. GLaM [59], a family of decoder-only language models, leverages a sparsely activated MoE design. V-MoE [219] presents a sparse adaptation of the ViT, scaling to 15 billion parameters, and achieves performance matching dense models while requiring less training time. LIMoE [190] represents the first multimodal model to incorporate sparse MoE,

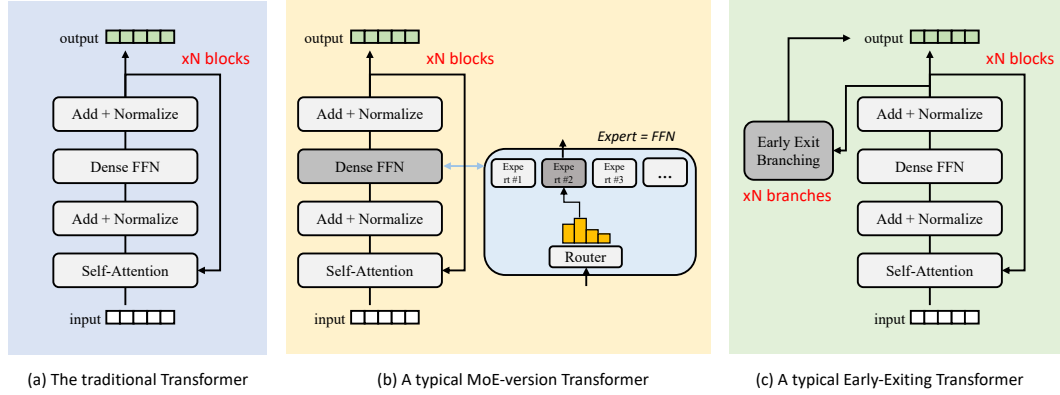


Fig. 4. Traditional and typical dynamic transformers.

significantly outperforming CLIP in various tasks. Mistral AI introduces Mistral¹, an MoE model comprising 8 experts, each with 7 billion parameters. This model outperforms the performance of LLaMA2-70B model [240]. MoEfication [306] converts a model into its MoE variant with equivalent parameters. Sparse upcycling [129] initializes sparsely activated MoE from dense checkpoints, reducing about 50% of the original dense pretraining costs. FFF [24] divides the feed-forward layer into separate leaves instead of copying the entire feed-forward layer as an expert, being up to 220× faster than the original feed-forward layer with about 5% accuracy loss. §5.1 will detail systematic optimizations applied to MoE models.

3.2.2 Early Exiting. As illustrated in Fig. 4(c), early exiting optimization is a strategy that allows a model to terminate its computational process prematurely when it attains high confidence in the prediction or encounters resource constraints. He et al. [95] investigates modifications to the standard transformer block, aiming for simpler yet efficient architectures without sacrificing performance. M4 [293] introduces a multi-path task execution framework, enabling elastic fine-tuning and execution of foundational model blocks for different training and inference tasks. FREE [21] proposes a shallow-deep module that synchronizes the decoding of the current token with previously processed early-exit tokens. SkipDecode [51] is designed for batch inferencing and KV caching, overcoming previous limitations by establishing a unique exit point for each token in a batch at every sequence position. PABEE [316] enhances the efficiency of pre-trained language models by integrating internal classifiers at each layer. The inference process halts when predictions stabilize for a set number of steps, facilitating quicker predictions with reduced layer usage. DeeBERT [273] augments BERT's inference efficiency by incorporating early exit points. DeeBERT allows instances to terminate at intermediate layers based on confidence levels, effectively reducing computational demands and accelerating inference. Bakhtiarnia et al. [23] proposes 7 distinct architectural designs for early-exit branches suitable for dynamic inference in ViTs backbones. LGViT [276] presents an early-exiting framework tailored for general ViTs, featuring diverse exiting heads, such as local perception and global aggregation heads, to balance efficiency and accuracy. This approach achieves competitive performance with an approximate 1.8× speedup.

3.3 Diffusion-specific Optimization

Generating images through diffusion models typically involves iterative process with numerous denoising steps. Recent research has focused on accelerating the denoising process and reducing the resource requirements during image

¹<https://mistral.ai/>

generation, which fall into three main categories: (1) efficient sampling, (2) diffusion in latent space, and (3) diffusion architecture variants.

3.3.1 Efficient Sampling. To enhance the denoising process of diffusion model while maintaining or improving sample quality, many efforts have been made to improve the sampling process. These works emphasize resource and time efficiency in their architectures. Nichol et al. [194] made strides in enhancing the traditional DDPM by focusing on resource efficiency. Their improved model not only competes in log-likelihoods but also enhances sample quality. This efficiency is achieved by learning the variances of the reverse diffusion process and employing a hybrid training objective. This methodology requires fewer forward passes, and shows improved scalability in terms of model capacity and computational power. DDIM [227] represents a significant improvement in time efficiency for diffusion models. By introducing a non-Markovian, deterministic approach to sampling, DDIM accelerates the generation process, allowing for faster sampling without compromising sample quality. PNDM [166] enhances the efficiency of DDPM in generating high-quality samples. The approach treats the diffusion process as solving differential equations on manifolds, greatly accelerating the inference process. DPM-Solver [176] utilizes a high-order solver that exploits the semi-linear structure of diffusion ODEs, facilitating fast and high-quality sample generation. Remarkably, DPM-Solver achieves this with as few as 10-20 denoising steps, highlighting the latency efficiency in sample generation.

3.3.2 Diffusion in Latent Space. In traditional diffusion models, operations are usually performed within the pixel space of images. However, this approach proves to be inefficient for high-resolution images because of the considerable computational demands and significant memory requirements. In response to these challenges, researchers proposed a shift towards conducting diffusion processes in latent space through VAEs. This paradigm results in substantial memory-efficient advancements, allowing for the generation of high-resolution images with reduced computational resources. LDM [220], also known as stable diffusion, serves as a notable example of memory-efficient image generation. By performing diffusion processes within a latent space derived from pixel data through a VAE, LDM effectively tackles scalability issues present in earlier diffusion models. LD-ZNet [209] leverages the memory-efficient properties of LDM for image segmentation tasks. This approach capitalizes on the deep semantic understanding inherent in LDM’s internal features, providing a nuanced bridge between real and AI-generated imagery. SALAD [132] introduces a memory-efficient methodology for 3D shape generation and manipulation with a cascaded diffusion model.

3.3.3 Diffusion Architecture Variants. Another method for enhancing diffusion models involves the adoption of more efficient model architectures. This strategy focuses on refining the structural framework of diffusion models to optimize their performance. SnapFusion [150] introduces an optimized text-to-image diffusion model for mobile devices, featuring a resource-efficient network architecture. This model overcomes the computational and latency limitations of existing models through a redesigned network architecture and improved step distillation. It generates high-quality 512×512 images in under 2 seconds with fewer denoising steps. ScaleCrafter [98] addresses the generation of ultra-high-resolution images using pre-trained diffusion models with an innovative and resource-efficient network design. ScaleCrafter incorporates techniques like “re-dilation”, “dispersed convolution”, and “noise-damped classifier-free guidance” to dynamically adjust convolutional perception fields during inference. ERNIE-ViLG [67] introduces a novel text-to-image diffusion model that integrates fine-grained textual and visual knowledge into a highly efficient network architecture. With a mixture-of-denoising-experts mechanism and scaling up to 24B parameters, ERNIE-ViLG outperforms the existing models on MS-COCO with a remarkable zero-shot FID-30k score of 6.75. Mobile diffusion [312] conducts a

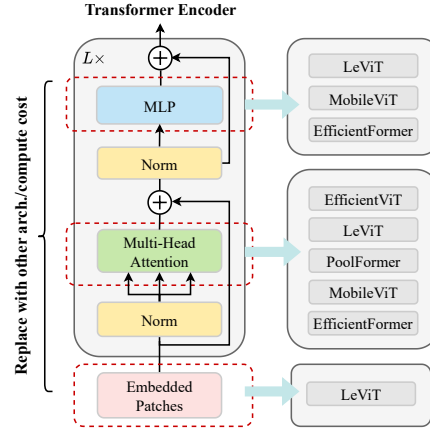


Fig. 5. A summary of resource-efficient ViT variants.

comprehensive examination of model architecture design to minimize model size and FLOPs. Besides, the authors also optimize the sampling steps, making one-step sampling compatible to downstream applications.

3.4 ViT-specific Optimizations

As a Transformer variant, ViT benefits from general optimizations aforementioned; yet, there exists also ViT-specific architecture optimizations as summarized in Fig. 5. LeViT [82] is a hybrid neural network designed for efficient image classification. Its main backbone features a pyramid architecture, progressively reducing the dimensionality of features while concurrently increasing the number of attention heads. MobileViT [181] adheres to the idea of utilizing CNNs to construct a more lightweight transformer architecture. Through the design of a convolution-like MobileViT block, the model achieves a lightweight and low-latency implementation, specifically tailored for practical hardware platforms. EfficientFormer [155] designs a lightweight CNN-Transformer hybrid architecture, achieving more efficient on-device inference. EfficientViT [29] introduces a linear attention mechanism to alleviate the computational cost linked with the high overhead of softmax in non-linear attention. In the domain of super-resolution, EfficientViT achieves a speedup of up to $6.4 \times$ compared to Restormer [296]. FastViT [244] introduces a token mixing operator that uses structural reparameterization to lower the memory access cost by removing the skip-connections in the network. Efficientvit [169] identifies that the speed of existing transformer models is commonly bounded by memory inefficient operations, especially the tensor reshaping and element-wise functions in MHSA. In response, the authors reduce the MHSA by a sandwiched structure. LightViT [105] presents several learning-based optimizations of pure convolution-free ViT architecture. EdgeViT [201] enables attention based vision models to compete with the best light-weight CNNs in the trade off between accuracy and on-device efficiency.

4 RESOURCE-EFFICIENT ALGORITHMS

This section focuses on resource-efficient large FMs techniques at the algorithm level. Compared to traditional DNNs, large FMs exhibit new characteristics such as its huge parameter set and autoregressive inference. This disparity has led to the emergence of numerous resource-efficient algorithms, which are categorized based on the lifecycle of FMs: pre-training, fine-tuning, serving algorithms, and model compression.

4.1 Pre-training Algorithms

Pre-training for large FMs relies on a substantial amount of computation resources. For instance, GPT-3-175B consumes 3.14×10^{23} FLOPs and LLaMA-70B takes 1.7×10^6 GPU hours. Consequently, optimizing the utilization of computational resources is crucial for the efficient pre-training of FMs. Resource-efficient algorithms can be categorized into training data deduction, neural architecture search, progressive learning, and mixed precision training.

4.1.1 Training Data Quality Control. A portion of work focus on controlling the quality of training data. DataComp [75] proposes a novel paradigm of locking the model/hyperparameters and refining the pretraining data. DFN [64] uses a proxy network as a modeling of the pretraining dataset. It recognizes that a better performance of the proxy network does not necessarily translate to the higher performance of the to-be-trained network. DataCompDR [245] of MobileCLIP leverages knowledge transfer from an image captioning model and an ensemble of strong CLIP encoders to improve the accuracy of efficient models.

4.1.2 Training Data Reduction. Pre-training for large FMs needs a dataset at the trillion-scale, exemplified by 0.3 trillion tokens for GPT-3-175B [26] and 2 trillion tokens for LLaMa-2-70B [240]. More data indicates more resource expenditure. Thereby, prior literature resort to reduce vast training data through two aspects: deduplicate text datasets and image patch removal.

Deduplicating text datasets [139] shows training data has redundancy caused by near-duplicate examples and long repetitive substrings. The reduction of repetitions can lead to fewer training steps without compromising performance.

Image patch removal is achieved by either reducing the number of patch inputs to the model or reorganizing image tokens based on modified model architectures. For instance, TRIPS [114] employs a patch selection layer to reduce image patches. This layer computes attentive image tokens through text guidance, resulting in a 40% reduction in computation resources, compared to previous pre-training vision-language models. Masked autoencoders (MAE) [96] masks image patches in pre-training phrase, but the large masking ratio brings significant computation resource wastage. MixMAE [164] introduces a method for mixing multiple images at the patch level, thereby avoiding the need for introducing “[MASK]” symbols. COPA [115] introduces an auxiliary pre-training task called patch-text alignment. This patch-level alignment strategy aims to decrease redundancy in image patches. PatchDropout [171] introduces the concept of patch dropout to enhance both computation and memory efficiency. This method involves the random sampling of a subset of original image patches to effectively shorten the length of token sequences.

4.1.3 Progressive Learning. Progressive learning is a training strategy that begins by training a small model and then gradually increases the model size, throughout the training process. This approach optimizes computational resource usage by reusing the computations from the previous stage. Inspired by the insight that knowledge can be shared across models of different depths, stackingBERT [80] introduces a progressive stacking algorithm. This algorithm cost-effectively trains a large model with no performance degradation by sequentially stacking attention layers from smaller models. CompoundGrow [85] identifies the similarity between progressive training algorithms and NAS. Staged training [223] adopts a strategy where a small model is pre-trained initially, and subsequently, the depth and width of the model are increased, continuing the training process. Knowledge inheritance [213] suggests employing existing pre-trained language models as teacher models to provide guidance during the training of larger models. The supplementary auxiliary supervision offered by the teacher model can effectively enhance the training speed of the larger model. The progressive training algorithm in AutoProg [143] is for the vision Transformer. AutoProg automatically adjusts the growth schedule to achieve lossless performance and make training resource consumption

minimal. LiGO [251] introduces small model parameters to initialize the large model through a trainable parameter linear map. LiGO achieves this by factorizing the growing transformation into a composition of linear operators at width and depth dimensions.

4.1.4 Mixed Precision Training. Mixed precision training often utilizes half-precision floating-point data representation instead of single precision. This approach significantly reduces memory requirements, approximately halving the storage space needed for weights, activations, and gradients. Mesa [203] proposes the combination of activation compressed training [33] with mixed precision training to further reduce the memory used by activations. The method quantifies activation based on the distribution of multi-head self-attention layers to minimize the approximation error. GACT [170] introduces a dynamically adjusted compression ratio based on the importance of each gradient.

4.2 Finetuning Algorithms

Efficient fine-tuning algorithms are designed to reduce the workload to adapt a pre-trained FM to downstream tasks. As summarized in Fig. 6, these techniques can be categorized into three groups: additive tuning, selective tuning, and re-parameter tuning.

4.2.1 Additive Tuning. Large FMs can achieve high performance with low costs by incorporating additional parameters and fine-tuning them for new tasks. In particular, this additive tuning process in Large FMs can be categorized into three main classes: adapter tuning, prompt tuning, and prefix tuning.

Adapter tuning aims to reduce training costs by introducing adapter modules to specific layers (or all layers) of a pre-trained large FMs. During tuning, the backbone of the pre-trained model remains frozen, and Adapter modules are utilized to acquire task-specific knowledge. Some works [62, 202, 236] focus on designing adapters for multi-task or multi-modal extensions. ADA [62] and MetaTroll [236] concentrate on incrementally extending pre-trained Transformers' capabilities across multiple tasks. This approach helps alleviate catastrophic forgetting during learning while simultaneously reducing computational expenses. ST-Adapter [202] introduces built-in spatiotemporal reasoning abilities, allowing pre-trained models to significantly reduce the number of parameters that need to be updated in cross-modal tasks. HiWi [158] improves inference speed by applying adapters to pre-trained parameters rather than hidden representations. AdaMix [257] designs a combined mechanism that merges the weights of different adapters into a single adapter at each Transformer layer. This innovation significantly reduces the additional storage cost introduced by multiple adapters. MEFT [159] designs a method for inserting adapters into LLM by modifying LLM to its reversible variant, reducing activation memory and thus improving the memory efficiency of fine-tuning. Residual Adapters [238] utilizes personalized residual adapters to address the issue of performance degradation in automatic speech recognition caused by non-standard speech. AutoProg [143] achieves lossless acceleration by automatically increasing the training overload on-the-fly. Such a procedure is done by progressively growth of subnets.

Prompt tuning involves designing a task-specific prompt for each task, with the aim of replacing the traditional fine-tuning of pre-trained large FMs parameters. By tuning the input prompts instead, this method significantly reduces the resources and time required for the fine-tuning. Some works [18, 140, 241] focus on improving the efficient scalability of prompts in multi-task settings. For example, PromptTuning [140], ATTEMPT [18], and BioInstruct [241] investigate how the utilization of mixed soft prompts can efficiently transfer knowledge across different tasks. These approaches help mitigate parameter update costs by reusing the frozen pre-trained large model. Furthermore, some works [38, 285] focus on minimizing prompt fine-tuning costs for specific tasks. For instance, DualPL [285] designs two prompts and separately captures the relevant knowledge of both tasks. This approach addresses the high cost associated with

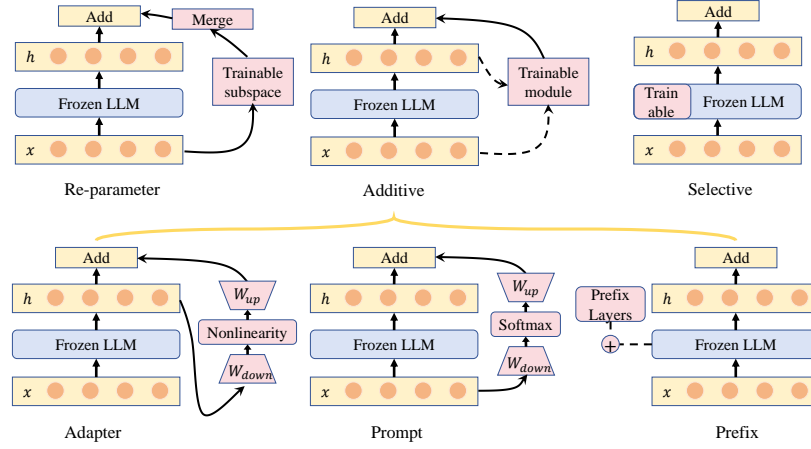


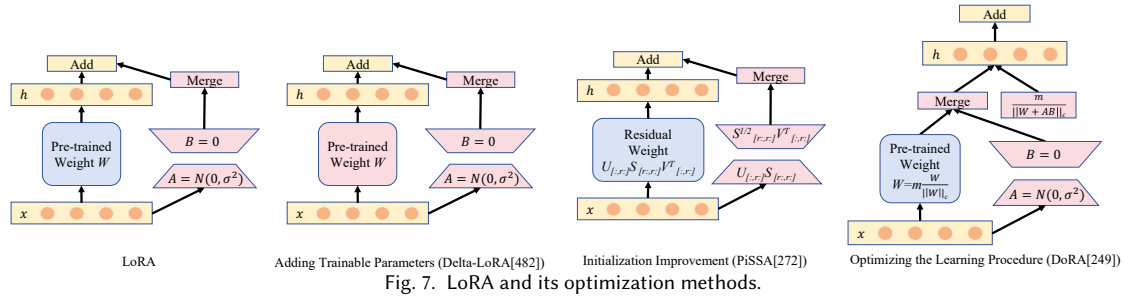
Fig. 6. A summary of various fine-tuning algorithms.

collecting state labels for slots and values in dialogue state tracking systems. In machine reading comprehension tasks, MPrompt [38] introduces task-specific multi-level prompt tuning to enhance the understanding of input semantics at different granularities while reducing the number of parameter updates.

Prefix tuning introduces a trainable, task-specific prefix part to each layer of large FMs. This technique aims to reduce the tuning cost by limiting the updates to the parameters in this prefix. Some works [162, 191, 247, 254, 309] focus on enhancing the performance of prefix tuning in specific domains. For example, UAPT [254] and Prefix-diffusion [162] address the issue of limited diversity in generating captions for images. These approaches extract image features from large FMs and design prefixes to enhance performance while reducing additional overhead. DOP [309] and DAPA [191] concentrate on domain-generalization problems in abstract summarization. These approaches design prefixes for each source domain to improve the model’s generalization capabilities. PIP [247] focuses on syntactic control in paraphrase generation and reduces training costs by designing parsing-indicating prefixes.

4.2.2 Selective Tuning. Selective tuning aims to maintain high performance on new tasks with low training costs by freezing the majority of parameters in large FMs and selectively updating only a small portion of the parameters. Some works focus on optimizing the performance of selective tuning. For example, SAM [74] explores how the choice of tunable parameters affects tuning. By proposing a second-order approximation method, it tunes fewer parameters to achieve better model performance. SmartFRZ [148] focuses on improving the efficiency of layer freezing by introducing an adaptive layer freezing technique based on different network structures. This innovation enhances system accuracy and training speed. FiSH-DiP [48] explores the effectiveness of tuning with limited data by introducing a sample-aware dynamic sparse tuning strategy. This approach selectively tunes partial parameters using sample feedback to enhance the model’s generalization in resource-constrained situations. Token Mixing [172] and VL-PET [102] enhance fine-tuning efficiency of visual-language tasks by adjusting and selecting a subset of trainable parameters.

4.2.3 Re-parameter Tuning. Re-parameter tuning adapts large FMs by targeting a significantly smaller subspace than the original, expansive training space. This approach involves fine-tuning low-rank matrix parameters, a technique that effectively reduces the overall training cost. The majority of existing research centers on reparameterization tuning through the implementation of the low-rank adapter design. For example, EfficientDM [97], QLoRA [53], PEQA [123], QALoRA [280] and LoftQ [153] incorporate quantization techniques, building upon the foundation of



LoRA. GLoRA [34] enhances LoRA’s generality, improving model transferability, few-shot capabilities, and domain generalization. PELA [89] derives inspiration from LoRA and devises a low-rank approximation compression method. LongLoRA [40] extends the capabilities of LoRA by incorporating context expansion through shift short attention. For ViT’s linear layers, LBP-WHT [284] diminishes the computational costs of matrix multiplication by employing low-rank backward propagation based on the Walsh-Hadamard transform. Additionally, DSEE [39] investigates the application of sparse-aware low-rank updates on pre-trained model weights. Dynamic-Pooling [193] mechanisms are designed to predict inference boundaries through autoregressive prediction.

LoRA, as the most popular parameter-efficient fine-tuning method, still exhibits performance gaps when compared to full fine-tuning. To address this, various methods have been developed to enhance LoRA’s performance, as shown in Fig. 7. Delta-LoRA [320] aims to bridge the performance gap by updating the pre-trained weights through the product of low-rank matrices A and B, thus adding trainable parameters without incurring additional memory overhead. On the other hand, PiSSA [182] identifies an issue where LoRA initializes low-rank matrices with Gaussian random values and zeros, resulting in very small initial gradient values and slow convergence. Lastly, DoRA [168] and LoRA+ [94] focus on enhancing the learning process itself to further improve efficiency and effectiveness. DoRA decomposes the pre-trained weights into their magnitude and directional components, and fine-tunes the directional matrix. LoRA+ sets the unbalanced learning rate for different blocks, accelerating convergence and improving fine-tuning performance.

4.3 Inference Algorithms

4.3.1 Opportunistic Decoding. Autoregressive mechanism significantly hinders the inference efficiency of large FMs. To address this, various approaches aim to replace autoregressive decoding with more efficient non-autoregressive techniques. Speculative decoding has been widely acknowledged as an effective method to accelerate autoregressive decoding. It involves generating sequences autoregressively with a cost-efficient small model, followed by parallel token verification using a larger model. Yaniv et al. [141] report a 2–3× improvement in performance using speculative decoding on the T5X model, while a concurrent study [36] demonstrates similar speedups on a 70B Chinchilla model. SpecTr [232] further enhances speculative decoding by increasing the number of candidate tokens and improving the draft selection process, resulting in a 2.13× improvement in wall clock speed and an additional 1.37× speedup on standard benchmarks. ProphetNet [281] introduces a sequence modeling architecture that predicts future tokens, partially reducing the reliance on autoregression. In the draft stage, Draft & Verify [301] skips certain intermediate layers, achieving a 1.73× speedup when tested on Llama-2. Medusa [30] offers another non-autoregressive decoding architecture that requires no auxiliary model, predicting multiple tokens by pre-training heads for different time steps and verifying them concurrently. Look-ahead decoding [72] accelerates inference in large FMs without relying on a draft model or data store, reducing decoding steps in proportion to $\log(\text{FLOPs})$. Additionally, speculative decoding is

the foundation for various inference systems, such as SpecInfer [185], which uses multiple draft models in the cloud, and LLMcad [274], deployed at the edge.

4.3.2 Input Filtering and Compression. This method includes directly filtering raw data, i.e., Prompt filtering, or filtering hidden activations of FMs, i.e., Token pruning.

Prompt compression. Computations can be effectively reduced by compressing the prompt to the model. LLMIn-gua [116] introduces a prompt compression approach from a coarse-to-fine perspective. Jiang et al. [264] investigate the feasibility, applicability, and potential of compressing natural language for large FMs while preserving semantics. EntropyRank [242] presents an unsupervised approach for extracting keywords and keyphrases from textual data. This method leverages a pre-trained language large FM and incorporates Shannon’s information maximization. LLMZip [243] employs LLaMA-7B for compressing natural language. Experimental results demonstrate that LLMZip outperforms cutting-edge text compression methods, including BSC, ZPAQ, and paq8h. AutoCompressors [42] utilizes large FMs to compress natural language into compact summary vectors. These vectors can then serve as soft prompts for large FMs usage. ICAE [77] utilizes the capabilities of large FMs to condense an extensive context into concise memory slots. These memory slots are directly adaptable by the large FMs for diverse purposes. Nugget 2D [212] introduces a prompt compression method specifically designed to handle long contexts. CoT-Max [106] is a context pruner, aiming to enhance the Chain-of-Thought ability of large FMs.

Token Pruning. Research has also explored the pruning of input sequences for transformers, often involving the incremental removal of less important tokens during inference. PoWER-BERT [81] proposes the direct learning of token pruning configurations. Length-Adaptive Transformer [122] extends this idea by introducing LengthDrop, a technique that entails training the model with various token pruning configurations, followed by an evolutionary search. TR-BERT [288] formulates token pruning as a multi-step token selection problem and addresses it through reinforcement learning. DynamicViT [216] hierarchically prunes redundant tokens based on their importance scores. AdaViT [183] and A-ViT [291] employ adaptive token reduction mechanisms and select different tokens for different images. AdaViT dynamically determines the usage of patches, self-attention heads, and transformer blocks based on the input. A-ViT discards tokens in vision transformers during inference, adapting the token retention based on the complexity of the input images. SPViT [131] devises an adaptive instance-wise token selector and introduces a soft pruning technique. PuMer [32] combines similar textual and visual tokens during inference for large-scale vision language models.

4.3.3 Key-Value Cache. Optimizing memory for the KV cache is a crucial aspect of the autoregressive decoder-based model inference process.

Memory efficient sparse attention. An alternative approach involves leveraging sparse attention. However, it’s noteworthy that most sparse attention designs, which primarily target the reduction of computational complexity [25, 295], do not necessarily lead to a reduction in KV cache memory consumption. This is because achieving a reduced memory footprint for the KV cache necessitates a more stringent sparsity pattern. Specifically, tokens that are sparsified should not be dynamically accessed in subsequent steps. To address this, H2O [307] introduces a KV cache eviction strategy designed for optimal memory efficiency. This strategy employs attention scores to identify and select the least important KV cache tokens in the current state for eviction. When compared to robust baselines, H2O demonstrates the capability to reduce latency by up to 1.9× and increase throughput by 29×. Dynamic Context Pruning [17] learns a memory-efficient KV cache eviction strategy during the pre-training phase. This approach has demonstrated the ability to achieve up to a 2× increase in inference throughput and even greater memory savings. Scissorhands [173] utilizes

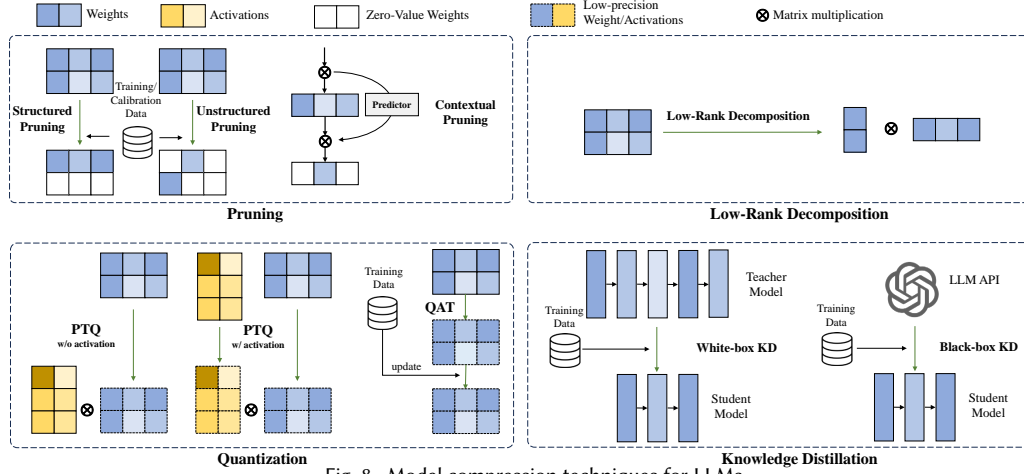


Fig. 8. Model compression techniques for LLMs.

Method	Categories	Unique challenge
Pruning	Structured Pruning [178, 269, 302], Unstructured Pruning [70, 226, 230], Contextual Pruning [175, 229]	Massive re-pretraining, Unique transformer structures
Knowledge distillation	White-box KD [41, 188], Black-box KD [86, 237]	Massive re-pretraining
Quantization	Quantization-Aware Training [174, 234], Post-Training Quantization [52, 69]	Quantization outliers Per-tensor quantization
Low-Rank Decomposition	[121, 154, 278]	/

Table 2. Model compression methods and their unique challenges.

an innovative compact KV cache and results in a notable reduction in KV cache inference memory usage, achieving up to a 5× reduction while maintaining model quality. By employing a landmark token to demarcate a token block, Landmark Attention [189] optimizes KV cache storage. This approach enables the storage of most KV cache in a slower but larger capacity memory, resulting in reduced memory requirements without compromising performance.

4.3.4 Long Context. To effectively process long sequences, transformers need to adapt their positional encoding to enhance their capability to capture long-range information. Due to the quadratic computational cost associated with attention mechanisms, various resource-efficient optimizations have been proposed to handle long inputs. LM-Infinite [91] introduces a Λ -shaped attention mechanism to handle long contexts efficiently. Characterized by computational efficiency with $O(n)$ time and space complexity, LM-Infinite consistently demonstrates fluency and quality in text generation for sequences as long as 128k tokens on ArXiv and OpenWebText2 datasets. StreamingLLM [272] facilitates large FMs trained with a finite-length attention window to generalize to infinite stream decoding without the need for any fine-tuning. PCW [217] segments a long context into chunks or “windows”, constrains the attention mechanism to operate solely within each window, and reuses positional embeddings across the windows. LongNet [55] introduces dilated attention, expanding the attentive field exponentially as the distance increases. This innovation allows LongNet to scale Transformers efficiently, enabling them to handle sequences of up to 1B tokens. SLED [110], short for SLiding-Encoder and Decoder, repurposes and capitalizes on well-validated short-text pre-trained language models. Despite competing effectively with specialized models that are up to 50× larger, SLED does not require a dedicated and expensive pretraining step.

4.4 Model Compression

As summarized in Fig. 8, model compression refers to a set of techniques aimed at reducing the model size without significant performance degradation, categorized into: pruning, knowledge distillation, quantization, and low-rank decomposition. While compression has been extensively studied in pre-LLM era [92, 93], compressing FMs faces unique challenges such as weight outliers and extensive training efforts, as discussed in Table 2.

4.4.1 Pruning. Pruning technique removes redundant or non-essential connections, neurons, or layers from a neural network. The primary objective is to reduce the model size, subsequently decreasing computational and storage costs, while maintaining model accuracy. Structured pruning and unstructured pruning target weight reduction without modifying sparsity during inference. In contrast, Contextual Pruning dynamically selects activated neurons or layers during inference based on the sparsity of the model.

Structured Pruning compresses large foundational models by eliminating entire structural components, such as groups of consecutive parameters or hierarchical structures. Examples of these structural components include channels or blocks of the model’s weights. It is often combined with fine-tuning to mitigate accuracy loss. LLM-Pruner [178] is a task-agnostic structured pruning algorithm that utilizes a small amount of data to assess the importance of coupled structure weights. The method selectively removes non-essential model structures based on gradient information. LLM-Pruner incorporates LoRA to recover the model’s accuracy after pruning. LoRAPrune [302] is another structured pruning approach based on LoRA, leveraging LoRA’s weights and gradients for importance estimation. This method iteratively eliminates excess channels and attention heads, achieving superior results compared to LLM-Pruner. Lagunas et al. [135] improved structured pruning techniques by incorporating blocks of variable sizes. This integration is applied within the movement pruning framework during fine-tuning, resulting in the removal of entire model components, such as attention heads. It achieves a 2.4× speedup and is 74% smaller compared to the original BERT.

Structured pruning is also employed in the training of large foundational models as well. Sheared LLaMA [269] adopts an end-to-end approach to remove channels, encompassing layers, attention heads, intermediate layers, and hidden layers. Sheared LLaMA demonstrates the capability to prune the LLaMA2-7B model down to 1.3B parameters. AdaPrune[108] accelerates neural network training using transposable masks, resulting in a 2× speed-up in matrix multiplications during both inference and training. GUM [222] considers neuron specificity and introduces pruning through network component-based global mobility and local uniqueness scores. This approach aims to simultaneously maximize sensitivity and uniqueness, effectively reducing redundant parameters in large FMs weights. PLATON [304] tackles the uncertainty in importance scores during model pruning by employing the upper confidence bound of importance estimation. This approach ensures stability in training and leads to improved generalization.

Unstructured Pruning does not consider the inherent structure of the model. Typically, it removes neurons with weights below a threshold, thereby compressing the model. When deploying unstructured pruning, specialized techniques are required to implement model storage compression. SparseGPT [70] treats the pruning framework as a generalized sparse regression problem and employs an approximate sparse regression solver, achieving 60% unstructured pruning on large GPT models like 175B. Wanda [230] leverages the observation of emergent large-magnitude features in large FMs. Wanda introduces sparsity by pruning weights with the smallest magnitudes multiplied by corresponding input activations, on a per-output basis. UPop [226] serves as a universal vision-language Transformer compression framework, which incorporates unifiedly multimodal subnets and progressively searching/retraining. SIGE [226] is proposed to convert computation reduction into latency reduction on standard hardware, achieving notable accelerations for models like DDPM, Stable Diffusion, and GauGAN with minimal edits.

Contextual Pruning selects the sparse state of each layer, making it hardware-optimization friendly. DeJa Vu [175] dynamically predicts the sparsity of the next layer using the activations of the previous layer. It determines which neurons of MLP blocks and the heads of attention blocks need to be retained. To mitigate the overhead of this predictor, DeJa Vu asynchronously predicts the next layer. PowerInfer [229] utilizes the sparsity of activation to dynamically predict the hot-activated neurons of the next layer and computes them on the GPU, while other cold-activated neurons are computed on the CPU. In comparison to llama.cpp [78], PowerInfer achieves up to 11× acceleration, enabling the 40B model to output ten tokens per second on a personal computer.

4.4.2 Knowledge Distillation. Knowledge Distillation (KD) transfers knowledge from a complex, heavy model (i.e., teacher model) to a simpler corresponding model (i.e., student model) for model compression. In general, there are two ways to apply KD to large FMs based on whether the internal structure of the teacher model is considered: white-box knowledge distillation and black-box knowledge distillation.

Black-box Knowledge Distillation. Assuming that the internal structure of the teacher’s large base model is not visible, this approach fine-tunes the student model using prompt-response pairs generated by large FMs’ API. The goal is to imbue the student model with the capabilities of the teacher model. For large FMs, the insights gained due to the increased parameter count contribute to strong generalization abilities. Therefore, techniques such as In-Context Learning (ICL) [57] and Chain-of-Thought (CoT) [259] can be utilized to enable the student model to thoroughly learn the capabilities of the large FMs. ICL distillation transfers few-shot learning and language model capabilities from the teacher model to the student model by integrating in-context learning objectives with traditional language modeling objectives. In Meta-ICL [188] and Metal-ICL [41], language models undergo meta-training on diverse tasks using in-context learning objectives. This process enables them to fine-tune for unseen tasks through in-context learning. Multitask-ICL [107] introduces the concept of in-context learning distillation, fine-tuning models with ICL objectives and examples from target tasks. CoT introduces intermediate reasoning steps in prompts, guiding language models to solve complex reasoning tasks step by step. Fu et al. [73] enhance the mathematical reasoning capabilities of smaller models by instructing them through CoT distilled from LLM teachers. Distilling Step-by-Step [101] extracts rationales from large FMs using CoT in a multi-task framework, providing additional guidance for training smaller models in a multi-task environment. Fine-tune-CoT [99] uses zero-shot CoT prompting techniques, employing random sampling to generate multiple reasoning solutions from large FMs to guide the training of student models.

White-box Knowledge Distillation. In contrast to black-box knowledge distillation, white-box knowledge distillation not only has access to the output results of the teacher model but also to its structure and intermediate results. Therefore, white-box knowledge distillation can better leverage the structure of the teacher model, enabling smaller student models to replicate and learn the capabilities of larger teacher models.

Timiryasov et al. [237] train an ensemble consisting of a GPT-2 and small LLaMA models on the developmentally plausible BabyLM dataset. Subsequently, they distilled it into a small LLaMA model with 58 million parameters, surpassing in performance both of its teachers as well as a similar model trained without distillation. MiniLLM [86] distills smaller language models from generative larger language models. This approach replaces the forward Kullback-Leibler Divergence (KLD) objective in the standard KD approaches with reverse KLD, which is more suitable for KD on generative language models, to prevent the student model from overestimating the low-probability regions of the teacher distribution. Instead of solely relying on a fixed set of output sequences, GKD [12] trains the student model using self-generated output sequences. TED [157] employs task-aware filters to align the hidden representations of the

student and the teacher at each layer. These filters are designed to select task-relevant knowledge from the hidden representations.

4.4.3 Quantization. Quantization is a well-established model compression method to mitigate the storage and computational demands. Compared to traditional DNNs, LLMs exhibit a higher frequency of activation outliers, which are crucial for maintaining model accuracy. Standard quantization often removes these outliers, leading to a significant performance drop.

Quantization-aware training (QAT) involves training a quantized model in such a way that it adapts its parameters to the lower precision introduced by quantization. The primary objective of this process is to mitigate the accuracy loss that occurs as a result of quantization. LLM-QAT tackles the issue of obtaining training data for LLMs by leveraging pre-trained models to generate samples through data-free distillation. Concurrently, it quantizes weights, activations, and KV cache, thereby improving training throughput. QuantGPT [234] achieves this by incorporating contrastive distillation from a full-precision teacher model and distilling logit information to a quantized student model during autoregressive pretraining. BitNet [249] pioneers QAT for 1-bit language models, training the language model with 1-bit weights and activations. Due to the substantial parameter count in large models often reaching tens or hundreds of billions, the training cost of QAT remains considerable. On the one hand, QAT for large FMs is often combined with knowledge distillation to reduce the training cost, as seen in approaches such as LLM-QAT and QuantGPT. On the other hand, quantization is frequently employed in the fine-tuning process of large models, such as in PEQA [268] and QLoRA [53].

Post-training quantization (PTQ) converts a trained full-precision model to a low-precision model without retraining. The advantage of PTQ lies in compressing models without altering the model structure or necessitating retraining, thereby reducing the storage and computational costs of models. Due to its low deployment cost, PTQ is also the most easily deployable and widely applicable technique in model compression. However, unlike QAT and distillation, PTQ lacks the feedback loop for adjusting precision through training. Research related to PTQ often focuses on efficiently preserving relevant information in weights/activations while compressing models. PTQ can be categorized into two groups: weight-only quantization and weight-activation co-quantization.

(1) *Weight-only quantization.* It only quantizes the model weights. There are two primary methods for mitigating quantization errors in the weight quantization of large FMs.

The first category involves identifying outliers and important weights in weights that significantly contribute to accuracy and treating these outliers specially. For instance, SpQR [54] identifies outlier weights and maintains them with high precision, while quantizing the rest of the weights. LLM.int8() [52] employs vectorized quantization and mixed-precision decomposition to handle outlier values for efficient inference. LLM.int8() utilizes 8-bit quantization for matrix multiplication, effectively reducing GPU memory usage during inference. AWQ [160] reduces quantization error by protecting the top 1% important weights in the model, utilizing per-channel scaling to determine the optimal scaling factor. OWQ [137] analysis suggests that abnormal activations amplify quantization errors, and it employs a mixed-precision scheme, applying higher precision quantization to weights with a significant impact from activated outlier values. SqueezeLLM [124] observes that sensitive weights determine the final model’s quantization performance and proposes a non-uniform quantization approach to minimize quantization errors in these sensitive weights.

The second category of quantization reduction methods is based on the second-order information updated weights. GPTQ [71] employs layer-wise quantization with OBQ [69], utilizing inverse Hessian information to update weights. GPTQ reduces the bit-width of each weight to 3 or 4 bits, allowing quantization of GPT models with 175 billion

parameters with minimal accuracy loss. QuIP [35] uses an adaptive rounding process, minimizing a second-order proxy objective for quantization.

(2) *Weights-activation co-quantization.* Quantizing both weights and activation facilitates deployment on hardware accelerators. SmoothQuant [271] takes advantage of the similarity in the channel-wise activations of different tokens and performs quantization on both weight and activation using per-channel scaling transforms. RPTQ [294] recognizes the substantial range differences across different channels, reordering the channels for quantization and integrating them into Layer Normalization and linear layer weights. OliVe [87] adopts outlier-victim pair quantization and locally processes outliers. Outlier Suppression+ [260] builds upon Outlier Suppression [261], discovering that harmful outliers exhibit an asymmetric distribution mainly concentrated in specific channels. Considering the asymmetry of outliers and quantization errors from the weights of the next layer, this approach performs channel-level translation and scaling operations. QLLM [163] addresses the issue of activation outliers through an adaptive channel reassembly method and mitigates the information loss caused by quantization using calibration data. LLM-FP4 [286] quantizes weights into 4-bit float points, proposes per-channel activation quantization, and reparameters additional scaling factors as exponential biases of weights. ZeroQuant [287] combines layer-wise knowledge distillation and optimized quantization support to achieve 8-bit quantization. FlexRound [138] updates the quantization scale of weights and activations by minimizing the error between the quantized values and the full-precision values. ATOM [311] significantly boosts serving throughput by using low-bit operators and considerably reduces memory consumption via low-bit quantization.

There is also extensive quantization research for backbone networks in FMs like ViT and BERT. For instance, BinaryBERT [197] and I-BERT [22] have achieved higher accuracy for BERT under low-precision quantization. Wang et al. [255] exploit the operator fusion [196], PTQ techniques, and structured pruning [135] to reduce the memory cost. They also reduce the number of computation operations of DeiT-Tiny [239]. Q-ViT [152], I-ViT [156], and OFQ [167] also achieve high accuracy for ViT under low-precision quantization. Q-Diffusion [149] compresses the noise estimation network to expedite the generation process of diffusion models.

4.4.4 Low-Rank Decomposition. Low-rank decomposition (LoRD) approximates weight matrix in large FMs by decomposing a given weight matrix into two or more smaller matrices. As mentioned in §4.2.3, low-rank decomposition has been widely applied in large FMs finetuning methods like LoRA. LoRD has also shown substantial compression capabilities with minimal impact on performance, highlighting its potential for large FMs compression [121]. To reduce the dimensionality of high-dimensional token embeddings underpinning large FMs, TensorGPT [278] proposes an approach based on the tensor-train decomposition, where each token embedding is treated as a matrix product state that can be efficiently computed in a distributed manner. Through TensorGPT, the embedding layer can be compressed by a factor of up to 38.40×. LoSparse [154] employs low-rank approximation to compress the coherent and expressive elements. The method uses iterative training to assess the significance scores of column neurons for the pruning process, showcasing superior performance compared to traditional iterative pruning techniques. Saha et al. [221] compress matrices through randomized low-rank and low-precision factorization, achieving compression ratios as aggressive as one bit per matrix coordinate while surpassing or maintaining the performance of traditional compression techniques. ViTALiTy [49] is an algorithm-hardware codesigned framework to enhance the inference efficiency of ViTs. It achieves approximation of the dot-product softmax operation with first-order Taylor attention, utilizing row-mean centering as the low-rank component to linearize the cost of attention blocks.

Name	Descriptions	Tags	Name	Descriptions	Tags
DeepSpeed [1]	An open-sourced Python library proposed by Microsoft. Supports MoE, long-sequence training, RLHF, ZeRO optimizations, and model compression.	C/T/I	Megatron [192]	The first cloud training system that introduces tensor parallelism to distributed training models like GPT, Bert, and T5. It is proposed by NVIDIA.	C/T
Alpa [313]	An automatic FM parallelization engine from UCB.	C/T/I	FairScale [63]	A new scaling library from Meta.	C/T/I
Colossal AI [146]	From HPC-AI Tech. Supports common parallelism strategies and heterogeneous memory management.	C/T/I	FlexFlow [185]	A cloud FM training and serving compiler from CMU and Stanford University. Automatic parallelization.	C/T/I
PyTorch FSDP [310]	A cloud large-scale training system atop PyTorch. It shards parameters, optimizer states, and gradients.	C/T/I	HF PEFT [2]	An efficient fine-tuning system from HuggingFace. It supports a set of PEFT methods like LoRA, p-tuning.	C/T
MII [1]	A library from DeepSpeed. Supports FastGen.	C/I	vLLM [134]	A serving engine from UC Berkeley. PagedAttention.	C/I
LightLLM [4]	A framework for token-wise's KV cache management.	C/I	Ray LLM [9]	A multiple LLMs serving solution from Anyscale.	C/I
TGI [3]	A high-performance serving engines from HuggingFace. It supports tensor parallelism, quantization with bitsandbytes and GPT-Q, and PagedAttention.	C/I	TRT-LLM [8]	A TensorRT toolbox for optimized LLM inference. It supports AWQ, GPTQ, SmoothQuant, speculative decoding, pipeline/tensor parallelism, PagedAttention.	C/I
llama.cpp [78]	A popular on-device LLM serving engine supporting mixed F16 / F32 precision and 2/3/4/5/6/8-bits int quantization. Mainly for LLaMA-based LLMs.	C/E/I	MNN-LLM [7]	An edge LLMs serving engine proposed by Alibaba and inherited from mnn. It optimizes the inference procedure separately in the prefill/decoding phase.	E/I
mlml [6]	A versatile and efficient on-device multimodal engine.	E/I	MLC-LLM [5]	Natively deploy LLMs with compiler accelerated APIs.	C/E/I

Table 3. Popular open-source tools for training and deploying large FMs. “C”: Cloud; “E”: Edge; “T”: Training; “I”: Inference.

5 RESOURCE-EFFICIENT SYSTEMS

Training and serving systems are key to practical large FMs. This section investigates the system research to enable resource-efficient large FMs, notable at four aspects: (1) distributed training; (2) hardware-aware optimizations; (3) serving in cloud, and (4) serving in edge. Table 3 summarizes widely-used open-source frameworks in this domain.

5.1 Distributed Training

Distributed training systems serve as the foundation for training large FMs, encompassing pretraining and fine-tuning phases. Pretraining, involving intensive computation and communication, demands substantial resources compared to other large FMs processes. Fine-tuning is widely used to transform a general-purpose model into a specialized model for particular use cases. Considering the large scale and new execution pattern of large FMs, designing resource-efficient systems for FMs has drawn great attention from the community. We categorize techniques for optimizing distributed training systems, covering aspects such as resilience, parallelism, communication, storage, and heterogeneous GPUs. Additionally, MoE has emerged as a trend in training extremely large models, for which several approaches are tailored. These specialized methods are detailed at the end of this subsection.

Resilience. The increasing size and duration of training for large FMs have led to a rise in failures, emphasizing the importance of resilient training [263]. Fault tolerance approaches for large FMs primarily manifest in four forms. First, Varuna and Gemini [19, 258] facilitate resilient training by implementing checkpoints to restart training. Varuna [19] is designed for training in commodity clusters with low-bandwidth networks, frequent pre-emptions, and user-friendly features. On the other hand, Gemini [258] expedites failure recovery through in-memory checkpoints. Second, Bamboo [235] utilizes redundant computations where one node performs computations for both itself and its neighbors. Bamboo avoids the overhead of recovering but introduces the overhead during training. Third, activation checkpointing [133, 308], which avoids storing the activation and recomputes it when needed, falls between the checkpointing and redundant computation approaches. The fourth approach involves recovering partial layers, as demonstrated by Oobleck [111]. In the event of a failure, the affected pipeline can be restored using partial layers from other replicas, incurring less overhead than employing the entire checkpoint.

Parallelism. Parallelism plays a crucial role in distributed training, especially for large FMs. Three types of parallelism are commonly employed for training large FMs. *Data parallelism* (DP) involves distributing the data across workers to scale up distributed training. DeepSpeed ZeRO [215] optimizes memory usage by splitting the model states. *Model parallelism* (MP) partitions the model in intra-layer paradigm (Tensor parallelism [192]) or inter-layer paradigm (Pipeline

parallelism [136, 200]). Tensor parallelism (TP) improves the training speed while leading to more communication. Pipeline parallelism (PP) improves GPU utilization by filling the bubbles. Breadth-first pipeline parallelism [136] designs a looping placement and breadth-first schedule to achieve both high GPU utilization and low cost. PipeFisher [200] assigns extra work to the bubbles for further benefits. Mobius [66] is designed for fine-tuning with a novel PP scheme and heterogeneous memory. FTPipe [61] partitions the model into finer-grained blocks rather than layers for flexible execution and low resource demand. *Sequence parallelism* (SP) [133, 147] is designed for the trend of long sequence training where training one sentence exceeds the memory capacity of one worker. SP divides the long sequence into multiple chunks and puts them on different workers. In practice, these parallelisms are usually used in a hybrid way. Galvaton [187] can automatically determine the most efficient hybrid parallelism strategy.

Communication. The large scale and complex parallelism lead to significant communication overhead. We summarize the optimization of communication into two categories: reducing the communication time directly and hiding the communication. Some work explores parallelism-aware communication compression [228] and heterogeneity-aware traffic reduction [308]. Existing work usually overlaps the communication with computation, by unifying the abstraction of computation and communication [112], decomposing the original communication collective [253], or designing a novel pipelining schedule [319].

Storage. Large FMs require a significant amount of storage resources, e.g., GPU memory for model states, host memory for model analysis, and disk for dataset and checkpoint. Various approaches have been proposed to alleviate the storage constraints for efficiency. Offloading is a common way to reduce the stress of GPU memory. ZeRO-Offload [218] offloads data and computations to CPU to train large models on a single GPU. FlashNeuron [20], on the other hand, offloads selective data to the SSD for higher throughput. Additionally, Behemoth [126] replaces low-capacity, high-performance HBM with high-capacity, low-performance NAND flash to enable data-parallel training for large FMs.

Heterogeneous GPUs. Training on specialized high-performance GPU clusters is impossible for most people or enterprises. Moreover, heterogeneous GPUs commonly exist even in specialized GPU clusters. Therefore, some efforts try to train large FMs on heterogeneous GPUs. Hetpipe [204] accelerates training with low-performance GPUs and Wave Synchronous Parallel to synchronize parameters among heterogeneous GPUs. Whale [113] introduces a hardware-aware load-balancing algorithm to speed up training.

MoE. MoE is an efficient approach to scaling up DNN models. The goals of optimizing MoE training systems are mainly efficiency and scalability. Existing work mainly optimizes the dynamism-related mechanisms, parallelism, and communication in MoE training. MegaBlocks [76] leverages sparse primitives to handle dynamic routing and load-imbalanced computation. Brainstorm [43] is a framework for dynamic DNNs by abstracting the dynamism and profile-based optimization. FlexMoE [195] focuses on the dynamic expert management and device placement problem. Additionally, Tutel [109] designs dynamic adaptive parallelism and pipelining strategies. SmartMoE [298] optimizes the parallelism strategy for efficient MoE training with a combination of offline and online mechanisms. Janus [165] changes communication from an expert-centric paradigm to a data-centric paradigm for faster communication in MoE training. MoE-Mamba [208] integrates MoE with Mamba [83] to enable selective state space models, reaching the same performance as Mamba in 2.35× fewer training steps.

5.2 Hardware-aware Optimizations

Some hardware-aware methods are also proposed to optimize FM. For instance, edgebert [233] proposes an in-depth algorithm-hardware co-design for latency-aware energy optimization for multi-task NLP. Its core is an entropy-based early exit prediction for dynamic DVFS at a sentence granularity. FlightLLM [297] is an end-to-end LLM inference

mapping flow on FPGAs. Its core is the computation and memory overhead of LLMs can be solved by utilizing FPGA-specific resources (e.g., DSP48 and heterogeneous memory hierarchy). SpAtten [250] proposes a sparse attention mechanism with cascade token and head pruning. It designs a novel top-k engine to rank token and head importance scores with high throughput, and along with other careful optimizations like progressive quantization. A3 [90] makes a key insight that the attention mechanism is semantically a content-based search where a large portion of computations ends up not being used. Recognizing that, it proposes an architecture with algorithmic approximation and hardware specialization.

5.3 Serving on Cloud

FM serving has two main phases: the prefill phase and the decoding phase. The prefill phase often processes a long sequence of input tokens in parallel, which is compute-intensive and can lead to potential bottlenecks if resources are not carefully allocated. In contrast, the decoding phase generates one token at a time, making it more bandwidth-bound [315]. Therefore, a series of optimizations for FM serving systems have been introduced to accelerate this process.

Inference Accelerating To accelerate the computation in a single accelerator, kernel optimization is a common approach. FlashAttention [45] and FlashAttention-2 [44] design for FM training can be simply used to accelerate the prefill phase. However, due to the unique characteristics of the decoding phase, Flash-Decoding [47] proposes a specific NVIDIA CUDA kernel to accelerate the decoding phase. FlashDecoding++ [100] further improves the performance of Flash-Decoding by optimizing the softmax operation and flat GEMM operation in the decoding phase and provides additional AMD GPU support. DeepSpeed-Inference [16], ByteTransformer [300], and Google’s PaLM serving system [211] also optimize GPU/TPU optimizations for small batch size scenarios, which is common in FM serving but rare in FM training. When scaling FM inference to numerous GPUs at a large scale, many works [16, 211] exploit combinations of various parallelism strategies, such as data parallelism, pipeline parallelism, tensor parallelism, and expert parallelism. These works efficiently serve FM inference on multiple modern accelerators, such as GPUs/TPUs.

Given the auto-regressive nature of FMs, various requests may feature distinct lengths of input tokens and output tokens. To address this issue, request batching and scheduling constitute another set of methods to enhance the computational efficiency of request processing. Orca [292] proposes selective batching and iteration-level scheduling to batch requests of different lengths at the granularity of iterations to increase the maximum batch size. FlexGen [225] proposes a request scheduling algorithm to mitigate the impact of offloading on the performance of latency-insensitive FM serving in a single GPU. FastServe [265] proposes an iteration-level preemptive scheduling and proactive KV cache swapping to mitigate the impact of head-of-line blocking on the performance of distributed FM serving. SARATHI [14] and DeepSpeed-FastGen [1] split the computation of the prefill phase into small chunks and schedule these chunks with the decoding phase to mitigate the impact of the prefill phase on the performance of large FMs serving. Splitwise [205] splits the prefill phase and the decoding phase onto different machines according to their different computation and memory requirements. Sarathi-Serve [13] introduces chunked-prefills scheduler which splits a prefill request into near equal sized chunks and creates stall-free schedules that adds new requests in a batch without pausing ongoing decodes. dLoRA [266] dynamically merges/unmerges adapters with the base model and migrating requests/adapters between worker replicas, significantly improving the serving throughput.

Memory Saving An FM consumes a large amount of memory during the serving process. To reduce the memory consumption of FM serving, many works propose various memory management techniques. As for FMs’ parameters and activations, DeepSpeed-Inference [16] and FlexGen [225] offload activations or model parameters to the DRAM or NVMe memories when the GPU memory is insufficient.

KV cache is another important memory component in FM serving. To reduce the memory consumption of KV cache, vLLM [134] adopts a block-level on-demand memory allocation mechanism, which only allocates memory to intermediate states when needed. vLLM also proposes a new operator, Paged Attention, to support attention operation when using this memory allocation mechanism. S-LoRA [224] extends this idea to Unified Paging to manage multiple LoRA adapters at the same time. SGLang [314] further exposes prompt programming primitives to users to enable more complex KV cache management among all requests with the help of RadixAttention.

Emerging Platforms Typical FM serving systems are usually deployed on data centers equipped with plenty of homogeneous high-performance servers. Due to the scarcity and cost of these high-performance servers, there are also some FM serving systems specifically designed for other deployment platforms. SpotServe [186] tries to serve FMs on spot instances, which are low-cost but unreliable cloud instances. SpotServe dynamically adjusts its parallelism strategy to accommodate the impact of spot instance preemption. As for FM serving on heterogeneous GPUs, HexGen [117] uses an evolutionary algorithm to search for high-performance FM placement on heterogeneous GPUs.

5.4 Serving on Edge

Large FMs have been widely adopted in many real-world mobile applications, such as search engines [10], chatbots [283], and intelligent agents [151]. With ever-increasing data privacy concerns and the stringent response latency requirement, running large FM on mobile devices locally (i.e., on-device inference) has recently attracted attention from both academia and industry. While small language models [177, 290] have been developed for on-device deployment, the runtime efficiency (decoding speed, memory footprint, energy consumption, etc) still remains a key challenge. Thereby, many on-device inference optimization techniques have been introduced.

Edge-cloud collaboration. A common strategy to tackle the scarce resources on mobile devices is to speed up the intensive inference with a powerful edge/cloud server collaboration. For instance, EdgeFM [282] queries and adapts the large FMs to the specific edge models with customized knowledge and architectures so that the dynamic edge model can ensure both low latency and close accuracy to the original large FMs.

On-device MoE models are proposed to only execute in routed sparse parameters during inference, which can decrease computation (detailed in §3.2). EdgeMoe [289] identifies the problem that experts have to be dynamically loaded into memory during inference. To tackle this issue, this approach proposes expert-wise bit-width adaptation to reduce the size of expert parameters with acceptable accuracy loss, saving parameters loading time. PC-MoE [130] is based on a crucial observation that expert activations are subject to temporal locality. Based on this observation, PC-MoE proposes Parameter Committee, which intelligently maintains a subset of crucial experts in use to reduce resource consumption.

Memory optimization. Since large FMs often rely on large parameter sizes and on-device memory resources are scarce (e.g., 8GB), inferring large FMs on devices faces the challenge of “memory wall”. To tackle this issue, LLMcad [274] utilizes speculative decoding [141] which can offload most workloads to a smaller memory-resident draft model. PowerInfer [229] relies on large FMs runtime sparsity, i.e., only hot neurons are consistently activated across inputs. To that end, PowerInfer preloads hot-activated neurons onto the GPU for fast access, while cold-activated neurons are computed on the CPU, thus significantly reducing GPU memory demands and CPU-GPU data transfers.

I/O optimization. As parameter size increasing speed is larger than edge devices’ memory increasing speed, dynamically loading parameters from disks to memory is avoidable. STI [88] identifies that loading parameters time is highly longer than computation time. To address this problem, STI proposes dynamically adapting weights bit-width during the loading procedure according to parameters importance, minimizing loading overhead under maximum

inference accuracy. LLM in a flash [15] solves this problem by fine-grained management of flash storage to reduce the volume of data transferred from flash to memory as well as reading data in larger, more contiguous chunks.

Kernel optimization. Computing resources are also crucial while limited resources on the devices. Prior study [305] implements the first 32-bit integer-based edge kernel for vision transformers with post-training integer-only quantization to speedup inference process. This method also introduces a range-constrained quantization technique for activation and normalization operators in transformers to trade-off data range and inference accuracy. llm.npu [275] offloads most of the LLM inference computation to hardware accelerator (NPU) to significantly improve the runtime efficiency.

6 CONCLUSIONS AND FUTURE DIRECTIONS

This survey provides a holistic, systematic overview of recent literature towards resource-efficient large FMs. We first present the preliminary background and cost analysis of the popular FMs, including language, vision, and multimodal. We then dive into the model architecture, algorithm, and system designs to enable more resource-efficient large FM lifecycle. In the future, the research of this domain will continue to be (or even more) crucial since the scaling law guarantees a promising future of more powerful AI with larger and larger models. Such research is also highly interdisciplinary, involving various CS communities such as machine learning, NLP/CV/Speech, networking, cloud computing, edge computing, etc.

The research opportunity of resource-efficient large FM is extremely large, notably:

(1) **Cloud-edge hybrid deployment.** To enable ubiquitous, privacy-preserving, and highly available general intelligence, many FMs will ultimately sink to near-user devices. Preliminary efforts have been already conducted to bring LLaMA-7B to smartphones and PCs. The killer applications include personal assistants/agents [151, 262], multimodal information retrieval [142], etc. In the future, at what size and speed the FMs can run on devices will become a key competitive force in the business model of hardware vendors.

(2) **Exploiting the model sparsity.** With model being larger, the activated ratio of model will go smaller for a given task. Recent literature [175] finds that even a densely trained non-MoE model exhibits runtime activation sparsity, which can be exploited to reduce inference time and memory footprint. We believe that exploiting the model and activation sparsity will be a promising direction towards sustainable model size scaling. More efficient sparse architectures other than MoE could emerge.

(3) **Large FM as a service.** On both clouds and devices, large FMs are unifying the DNN ecosystem [293]. Ultimately, it becomes a universal service to be invoked just as today's Web and Database. On the one hand, it opens the opportunity for highly hardware-algorithm co-design and optimizations; meanwhile, it poses new challenges in system and infrastructure design for scheduling, load balancing, and security&isolation.

(4) **Agent as a holistic system to optimize.** In the future, FMs especially LLMs will be used as a key building block for establishing agents [151, 262]. Its efficiency shall not be considered as in a standalone LLM service; instead, the algorithm and system designs need to cater to the specific agent workflow. For example, an agent system might require multiple FMs to cooperate, where there exists inherent logic dependency. In this process, the design space of selecting the proper FMs for each task and scheduling them on a given set of hardware resources to maximize the agent performance is huge.

(5) **Practical privacy-preserving FM.** As the volume of user data uploaded to the cloud for FM processing continues to increase, the severity of privacy concerns correspondingly escalates. Existing methods include federated learning²,

²You can find a brief literature survey of resource-efficient federated learning in the Appendix.

homomorphic encryption, and disentanglement learning. While being theoretically sound, those methods still confront significant performance challenges, hindering their large-scale in-the-wild deployment. A promising direction involves the development of innovative privacy-preserving techniques specifically designed for large FMs, or the refinement of existing methods, to effectively balance privacy with performance.

ACKNOWLEDGEMENT

Mengwei Xu was supported by NSFC 62102045.

REFERENCES

- [1] 2023. DeepSpeed-FastGen: High-throughput Text Generation for LLMs via MII and DeepSpeed-Inference. <https://github.com/microsoft/DeepSpeed/tree/master/blogs/deepspeed-fastgen>.
- [2] 2023. Huggingface PEFT. <https://github.com/huggingface/peft>.
- [3] 2023. HuggingFace Text Generation Inference. <https://github.com/huggingface/text-generation-inference>.
- [4] 2023. LightLLM. <https://github.com/ModelTC/lightllm>.
- [5] 2023. mlc-llm. <https://github.com/mlc-ai/mlc-llm>.
- [6] 2023. mllm. <https://github.com/UbiquitousLearning/mllm>.
- [7] 2023. mnn-llm. <https://github.com/wangzhaode/mnn-llm>.
- [8] 2023. NVIDIA TensorRT-LLM. <https://github.com/NVIDIA/TensorRT-LLM>.
- [9] 2023. Ray LLM. <https://github.com/ray-project/ray-llm>.
- [10] 2024. Microsoft Recall. <https://support.microsoft.com/en-us/windows/retrace-your-steps-with-recall-aa03f8a0-a78b-4b3e-b0a1-2eb8ac48701c>.
- [11] Fernando Acosta and others. 2023. Medical diagnostics with Segment Anything Model (SAM): A case study in tumor segmentation. *arXiv preprint arXiv:2307.01234* (2023).
- [12] Rishabh Agarwal, Nino Vieillard, Piotr Stanczyk, and others. 2023. GKD: Generalized Knowledge Distillation for Auto-regressive Sequence Models. *arXiv preprint arXiv:2306.13649* (2023).
- [13] Amey Agrawal, Nitin Kedia, Ashish Panwar, and others. 2024. Taming {Throughput-Latency} Tradeoff in {LLM} Inference with {Sarathi-Serve}. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*. 117–134.
- [14] Amey Agrawal, Ashish Panwar, Jayashree Mohan, and others. 2023. SARATHI: Efficient LLM Inference by Piggybacking Decodes with Chunked Prefills. *arXiv preprint arXiv:2308.16369* (2023).
- [15] Keivan Alizadeh, Iman Mirzadeh, Dmitry Belenko, and others. 2023. LLM in a flash: Efficient Large Language Model Inference with Limited Memory. *arXiv:2312.11514* [cs.CL]
- [16] Reza Yazdani Aminabadi, Samyam Rajbhandari, Ammar Ahmad Awan, and others. 2022. DeepSpeed-inference: enabling efficient inference of transformer models at unprecedented scale. In *SC22: International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 1–15.
- [17] Sotiris Anagnostidis, Dario Pavlo, Luca Biggio, and others. 2023. Dynamic Context Pruning for Efficient and Interpretable Autoregressive Transformers. *arXiv:2305.15805* [cs.CL]
- [18] Akari Asai, Mohammadreza Salehi, Matthew E Peters, and Hannaneh Hajishirzi. 2022. Attempt: Parameter-efficient multi-task tuning via attentional mixtures of soft prompts. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. 6655–6672.
- [19] Sanjith Athlur, Nitika Saran, Muthian Sivathanu, and others. 2022. Varuna: scalable, low-cost training of massive deep learning models. In *Proceedings of the Seventeenth European Conference on Computer Systems*. 472–487.
- [20] Jonghyun Bae, Jongsung Lee, Yunho Jin, and others. 2021. FlashNeuron:SSD-Enabled Large-Batch Training of Very Deep Neural Networks. In *19th USENIX Conference on File and Storage Technologies*. 387–401.
- [21] Sangmin Bae, Jongwoo Ko, Hwanjun Song, and Se-Young Yun. 2023. Fast and Robust Early-Exiting Framework for Autoregressive Language Models with Synchronized Parallel Decoding. *arXiv preprint arXiv:2310.05424* (2023).
- [22] Haoli Bai, Wei Zhang, Lu Hou, and others. 2020. Binarybert: Pushing the limit of bert quantization. *arXiv preprint arXiv:2012.15701* (2020).
- [23] Arian Bakhtiarnia, Qi Zhang, and Alexandros Iosifidis. 2021. Multi-Exit Vision Transformer for Dynamic Inference. (2021).
- [24] Peter Belcak and Roger Wattenhofer. 2023. Fast Feedforward Networks. *arXiv:2308.14711* [cs.LG]
- [25] Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150* (2020).
- [26] Tom Brown, Benjamin Mann, Nick Ryder, and others. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [27] Aydar Bulatov, Yuri Kuratov, and Mikhail Burtsev. 2022. Recurrent memory transformer. *Advances in Neural Information Processing Systems* 35 (2022), 11079–11091.
- [28] Aydar Bulatov, Yuri Kuratov, and Mikhail S Burtsev. 2023. Scaling Transformer to 1M tokens and beyond with RMT. *arXiv preprint arXiv:2304.11062* (2023).

- [29] Han Cai, Junyan Li, Muyan Hu, and others. 2024. EfficientViT: Multi-Scale Linear Attention for High-Resolution Dense Prediction. *arXiv:2205.14756* [cs.CV] <https://arxiv.org/abs/2205.14756>
- [30] Tianle Cai, Yuhong Li, Zhengyang Geng, and others. 2023. Medusa: Simple Framework for Accelerating LLM Generation with Multiple Decoding Heads. <https://github.com/FasterDecoding/Medusa>.
- [31] Liang Cao and others. 2023. Text-to-image diffusion models for art generation: Capabilities and challenges. *arXiv preprint arXiv:2306.12567* (2023).
- [32] Qingqing Cao, Bhargavi Paranjape, and Hannaneh Hajishirzi. 2023. PuMer: Pruning and Merging Tokens for Efficient Vision Language Models. *arXiv:2305.17530* [cs.CV]
- [33] Ayan Chakrabarti and Benjamin Moseley. 2019. Backprop with approximate activations for memory-efficient network training. *Advances in Neural Information Processing Systems* 32 (2019).
- [34] Arnav Chavan, Zhuang Liu, Deepak Gupta, and others. 2023. One-for-All: Generalized LoRA for Parameter-Efficient Fine-tuning. *arXiv preprint arXiv:2306.07967* (2023).
- [35] Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher De Sa. 2023. QuIP: 2-Bit Quantization of Large Language Models With Guarantees. *arXiv preprint arXiv:2307.13304* (2023).
- [36] Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, and others. 2023. Accelerating Large Language Model Decoding with Speculative Sampling. *ArXiv abs/2302.01318* (2023). <https://api.semanticscholar.org/CorpusID:256503945>
- [37] Chun-Fu Chen, Quanfu Fan, and Rameswar Panda. 2021. CrossViT: Cross-Attention Multi-Scale Vision Transformer for Image Classification. *arXiv:2103.14899* [cs.CV]
- [38] Guoxin Chen, Yiming Qian, Bowen Wang, and Liangzhi Li. 2023. MPrompt: Exploring Multi-level Prompt Tuning for Machine Reading Comprehension. *arXiv preprint arXiv:2310.18167* (2023).
- [39] Xuxi Chen, Tianlong Chen, Yu Cheng, and others. 2021. Dsee: Dually sparsity-embedded efficient tuning of pre-trained language models. *arXiv preprint arXiv:2111.00160* (2021).
- [40] Yukang Chen, Shengju Qian, Haotian Tang, and others. 2023. Longlora: Efficient fine-tuning of long-context large language models. *arXiv preprint arXiv:2309.12307* (2023).
- [41] Yanda Chen, Ruiqi Zhong, Sheng Zha, and others. 2021. Meta-learning via language model in-context tuning. *arXiv preprint arXiv:2110.07814* (2021).
- [42] Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and Danqi Chen. 2023. Adapting Language Models to Compress Contexts. *arXiv:2305.14788* [cs.CL]
- [43] Weihao Cui, Zhenhua Han, Lingji Ouyang, and others. 2023. Optimizing Dynamic Neural Networks with Brainstorm. In *17th USENIX Symposium on Operating Systems Design and Implementation*. 797–815.
- [44] Tri Dao. 2023. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv preprint arXiv:2307.08691* (2023).
- [45] Tri Dao, Dan Fu, Stefano Ermon, and others. 2022. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems* 35 (2022), 16344–16359.
- [46] Tri Dao, Daniel Y Fu, Khaled K Saab, and others. 2022. Hungry hungry hippos: Towards language modeling with state space models. *arXiv preprint arXiv:2212.14052* (2022).
- [47] Tri Dao, Daniel Haziza, Francisco Massa, and Grigory Sizov. 2023. Flash-decoding for long-context inference.
- [48] Sarkar Snigdha Sarathi Das, Ranran Haoran Zhang, Peng Shi, and others. 2023. Unified Low-Resource Sequence Labeling by Sample-Aware Dynamic Sparse Finetuning. *arXiv preprint arXiv:2311.03748* (2023).
- [49] Jyotikrishna Dass, Shang Wu, Huihong Shi, and others. 2023. Vitality: Unifying low-rank and sparse approximation for vision transformer acceleration with a linear taylor attention. In *2023 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 415–428.
- [50] Alex de Vries. 2023. The growing energy footprint of artificial intelligence. *Joule* 7, 10 (2023), 2191–2194.
- [51] Luciano Del Corro, Allie Del Giorno, Sahaj Agarwal, and others. 2023. Skipdecode: Autoregressive skip decoding with batching and caching for efficient llm inference. *arXiv preprint arXiv:2307.02628* (2023).
- [52] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. 2022. LLM.int8(): 8-bit Matrix Multiplication for Transformers at Scale. *arXiv:2208.07339* [cs.LG]
- [53] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314* (2023).
- [54] Tim Dettmers, Ruslan Svirschevski, Vage Egiazarian, and others. 2023. SpQR: A Sparse-Quantized Representation for Near-Lossless LLM Weight Compression. *arXiv preprint arXiv:2306.03078* (2023).
- [55] Jiayu Ding, Shuming Ma, Li Dong, and others. 2023. LongNet: Scaling Transformers to 1,000,000,000 Tokens. *arXiv:2307.02486* [cs.CL]
- [56] Tianyu Ding, Tianyi Chen, Haidong Zhu, and others. 2023. The efficiency spectrum of large language models: An algorithmic survey. *arXiv preprint arXiv:2312.00678* (2023).
- [57] Qingxiu Dong, Lei Li, Damai Dai, and others. 2022. A survey for in-context learning. *arXiv preprint arXiv:2301.00234* (2022).
- [58] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, and others. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).

- [59] Nan Du, Yanping Huang, Andrew M Dai, and others. 2022. Glam: Efficient scaling of language models with mixture-of-experts. In *International Conference on Machine Learning*. PMLR, 5547–5569.
- [60] Julian Eisenschlos, Maharshi Gor, Thomas Mueller, and William Cohen. 2021. MATE: Multi-view Attention for Table Transformer Efficiency. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. 7606–7619.
- [61] Saar Eliad, Ido Hakimi, Alon De Jagger, and others. 2021. Fine-tuning giant neural networks on commodity hardware with automatic pipeline model parallelism. In *USENIX Annual Technical Conference*. 381–396.
- [62] Beyza Ermiş, Giovanni Zappella, Martin Wistuba, and others. 2022. Memory efficient continual learning with transformers. *Advances in Neural Information Processing Systems* 35 (2022), 10629–10642.
- [63] FairScale authors. 2021. FairScale: A general purpose modular PyTorch library for high performance and large scale training. <https://github.com/facebookresearch/fairscale>.
- [64] Alex Fang, Albin Madappally Jose, Amit Jain, and others. 2023. Data Filtering Networks. arXiv:2309.17425 [cs.AI] <https://arxiv.org/abs/2309.17425>
- [65] William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *The Journal of Machine Learning Research* 23, 1 (2022), 5232–5270.
- [66] Yangyang Feng, Minhui Xie, Zijie Tian, and others. 2023. Mobius: Fine tuning large-scale models on commodity gpu servers. In *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*. 489–501.
- [67] Zhida Feng, Zhenyu Zhang, Xintong Yu, and others. 2023. ERNIE-ViLG 2.0: Improving text-to-image diffusion model with knowledge-enhanced mixture-of-denoising-experts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10135–10145.
- [68] Quentin Fournier, Gaétan Marceau Caron, and Daniel Aloise. 2023. A practical survey on faster and lighter transformers. *Comput. Surveys* 55, 14s (2023), 1–40.
- [69] Elias Frantar and Dan Alistarh. 2022. Optimal brain compression: A framework for accurate post-training quantization and pruning. *Advances in Neural Information Processing Systems* 35 (2022), 4475–4488.
- [70] Elias Frantar and Dan Alistarh. 2023. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning*. PMLR, 10323–10337.
- [71] Elias Frantar, Saleh Ashkboos, Torsten Hoefer, and Dan Alistarh. 2022. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323* (2022).
- [72] Yichao Fu, Peter Bailis, Ion Stoica, and Hao Zhang. 2023. Breaking the Sequential Dependency of LLM Inference Using Lookahead Decoding. <https://lmsys.org/blog/2023-11-21-lookahead-decoding/>
- [73] Yao Fu, Hao Peng, Litu Ou, and others. 2023. Specializing Smaller Language Models towards Multi-Step Reasoning. *arXiv preprint arXiv:2301.12726* (2023).
- [74] Zihao Fu, Haoran Yang, Anthony Man-Cho So, and others. 2023. On the effectiveness of parameter-efficient fine-tuning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 12799–12807.
- [75] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, and others. 2023. DataComp: In search of the next generation of multimodal datasets. arXiv:2304.14108 [cs.CV] <https://arxiv.org/abs/2304.14108>
- [76] Trevor Gale, Deepak Narayanan, Cliff Young, and Matei Zaharia. 2023. MegaBlocks: Efficient Sparse Training with Mixture-of-Experts. *Proceedings of Machine Learning and Systems* (2023).
- [77] Tao Ge, Jing Hu, Lei Wang, and others. 2023. In-context Autoencoder for Context Compression in a Large Language Model. arXiv:2307.06945 [cs.CL]
- [78] Georgi Gerganov. 2023. llama.cpp: A C++ implementation of the Large Language Models. <https://github.com/ggerganov/llama.cpp>.
- [79] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, and others. 2023. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15180–15190.
- [80] Linyuan Gong, Di He, Zhuohan Li, and others. 2019. Efficient training of BERT by progressively stacking. In *International conference on machine learning*. PMLR, 2337–2346.
- [81] Saurabh Goyal, Anamitra R. Choudhury, Saurabh M. Raje, and others. 2020. PoWER-BERT: Accelerating BERT Inference via Progressive Word-vector Elimination. arXiv:2001.08950 [cs.LG]
- [82] Ben Graham, Alaaeldin El-Nouby, Hugo Touvron, and others. 2021. LeViT: a Vision Transformer in ConvNet’s Clothing for Faster Inference. arXiv:2104.01136 [cs.CV]
- [83] Albert Gu and Tri Dao. 2023. Mamba: Linear-Time Sequence Modeling with Selective State Spaces. *arXiv preprint arXiv:2312.00752* (2023).
- [84] Albert Gu, Karan Goel, and Christopher Ré. 2022. Efficiently modeling long sequences with structured state spaces. (2022).
- [85] Xiaotao Gu, Liyuan Liu, Hongkun Yu, and others. 2021. On the Transformer Growth for Progressive BERT Training. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 5174–5180.
- [86] Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. 2023. Knowledge Distillation of Large Language Models. *arXiv preprint arXiv:2306.08543* (2023).
- [87] Cong Guo, Jiaming Tang, Weiming Hu, and others. 2023. Olive: Accelerating Large Language Models via Hardware-friendly Outlier-Victim Pair Quantization. In *Proceedings of the 50th Annual International Symposium on Computer Architecture (ISCA ’23)*. ACM. <https://doi.org/10.1145/3579371.3589038>
- [88] Liwei Guo, Wonkyo Choe, and Felix Xiaozhu Lin. 2023. STI: Turbocharge NLP Inference at the Edge via Elastic Pipelining. In *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*. 791–803.

- [89] Yangyang Guo, Guangzhi Wang, and Mohan Kankanhalli. 2023. PELA: Learning Parameter-Efficient Models with Low-Rank Approximation. *arXiv preprint arXiv:2310.10700* (2023).
- [90] Tae Jun Ham, Sung Jun Jung, Seonghak Kim, and others. 2020. A³: Accelerating Attention Mechanisms in Neural Networks with Approximation. *arXiv:2002.10941* [cs.DC] <https://arxiv.org/abs/2002.10941>
- [91] Chi Han, Qifan Wang, Wenhan Xiong, and others. 2023. LM-Infinite: Simple On-the-Fly Length Generalization for Large Language Models. *arXiv:2308.16137* [cs.CL]
- [92] Song Han, Huizi Mao, and William J. Dally. 2016. Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding. *arXiv:1510.00149* [cs.CV] <https://arxiv.org/abs/1510.00149>
- [93] Song Han, Jeff Pool, John Tran, and William J. Dally. 2015. Learning both Weights and Connections for Efficient Neural Networks. *arXiv:1506.02626* [cs.NE] <https://arxiv.org/abs/1506.02626>
- [94] Soufiane Hayou, Nikhil Ghosh, and Bin Yu. [n. d.]. LoRA+: Efficient Low Rank Adaptation of Large Models. ([n. d.]).
- [95] Bobby He and Thomas Hofmann. 2023. Simplifying Transformer Blocks. *arXiv preprint arXiv:2311.01906* (2023).
- [96] Kaiming He, Xinlei Chen, Saining Xie, and others. 2022. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 16000–16009.
- [97] Yefei He, Jing Liu, Weijia Wu, and others. 2023. EfficientDM: Efficient Quantization-Aware Fine-Tuning of Low-Bit Diffusion Models. *arXiv preprint arXiv:2310.03270* (2023).
- [98] Yingqing He, Shaoshu Yang, Haoxin Chen, and others. 2023. ScaleCrafter: Tuning-free Higher-Resolution Visual Generation with Diffusion Models. *arXiv preprint arXiv:2310.07702* (2023).
- [99] Namgyu Ho, Laura Schmid, and Se-Young Yun. 2022. Large language models are reasoning teachers. *arXiv preprint arXiv:2212.10071* (2022).
- [100] Ke Hong, Guohao Dai, Jiaming Xu, and others. 2023. FlashDecoding++: Faster Large Language Model Inference on GPUs. *arXiv preprint arXiv:2311.01282* (2023).
- [101] Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, and others. 2023. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. *arXiv preprint arXiv:2305.02301* (2023).
- [102] Zi-Yuan Hu, Yanyang Li, Michael R Lyu, and Liwei Wang. 2023. VL-PET: Vision-and-Language Parameter-Efficient Tuning via Granularity Control. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3010–3020.
- [103] Jane Huang and Kevin Williams. 2023. GPT-4 for creative writing: A case study of content generation in digital media. *arXiv preprint arXiv:2305.11234* (2023).
- [104] Luyang Huang, Shuyang Cao, Nikolaus Parulian, and others. 2021. Efficient attentions for long document summarization. *arXiv preprint arXiv:2104.02112* (2021).
- [105] Tao Huang, Lang Huang, Shan You, and others. 2022. LightViT: Towards Light-Weight Convolution-Free Vision Transformers. *arXiv:2207.05557* [cs.CV] <https://arxiv.org/abs/2207.05557>
- [106] Xijie Huang, Li Lyna Zhang, Kwang-Ting Cheng, and Mao Yang. 2023. Boosting LLM Reasoning: Push the Limits of Few-shot Learning with Reinforced In-Context Pruning. *arXiv:2312.08901* [cs.CL]
- [107] Yukun Huang, Yanda Chen, Zhou Yu, and Kathleen McKeown. 2022. In-context Learning Distillation: Transferring Few-shot Learning Ability of Pre-trained Language Models. *arXiv preprint arXiv:2212.10670* (2022).
- [108] Itay Hubara, Brian Chmiel, Moshe Island, and others. 2021. Accelerated Sparse Neural Training: A Provable and Efficient Method to Find N:M Transposable Masks. *arXiv:2102.08124* [cs.AI]
- [109] Changho Hwang, Wei Cui, Yifan Xiong, and others. 2023. Tutel: Adaptive mixture-of-experts at scale. *Proceedings of Machine Learning and Systems* (2023).
- [110] Maor Ivgi, Uri Shaham, and Jonathan Berant. 2023. Efficient long-text understanding with short-text models. *Transactions of the Association for Computational Linguistics* 11 (2023), 284–299.
- [111] Insu Jang, Zhenning Yang, Zhen Zhang, and others. 2023. Ooblock: Resilient Distributed Training of Large Models Using Pipeline Templates. In *Proceedings of the 29th Symposium on Operating Systems Principles*. 382–395.
- [112] Abhinav Jangda, Jun Huang, Guodong Liu, and others. 2022. Breaking the computation and communication abstraction barrier in distributed machine learning workloads. In *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*. 402–416.
- [113] Xianyan Jia, Le Jiang, Ang Wang, and others. 2022. Whale: Efficient giant model training over heterogeneous GPUs. In *USENIX Annual Technical Conference*. 673–688.
- [114] Chaoya Jiang, Haiyang Xu, Chenliang Li, and others. 2022. TRIPS: Efficient Vision-and-Language Pre-training with Text-Relevant Image Patch Selection. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. 4084–4096.
- [115] Chaoya Jiang, Haiyang Xu, Wei Ye, and others. 2023. COPA: Efficient Vision-Language Pre-training through Collaborative Object-and Patch-Text Alignment. In *Proceedings of the 31st ACM International Conference on Multimedia*. 4480–4491.
- [116] Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, and others. 2023. LLMingua: Compressing Prompts for Accelerated Inference of Large Language Models. *arXiv:2310.05736* [cs.CL]
- [117] Youhe Jiang, Ran Yan, Xiaozhe Yao, and others. 2023. HexGen: Generative Inference of Foundation Model over Heterogeneous Decentralized Environment. *arXiv preprint arXiv:2311.11514* (2023).

- [118] Praneeth Kacham, Vahab Mirrokni, and Peilin Zhong. 2023. Polysketchformer: Fast transformers via sketches for polynomial kernels. *arXiv preprint arXiv:2310.01655* (2023).
- [119] Jared Kaplan, Sam McCandlish, Tom Henighan, and others. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361* (2020).
- [120] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. 2020. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International conference on machine learning*. PMLR, 5156–5165.
- [121] Ayush Kaushal, Tejas Vaidhya, and Irina Rish. 2023. LORD: Low Rank Decomposition Of Monolingual Code LLMs For One-Shot Compression. *arXiv preprint arXiv:2309.14021* (2023).
- [122] Gyuwan Kim and Kyunghyun Cho. 2021. Length-Adaptive Transformer: Train Once with Length Drop, Use Anytime with Search. *arXiv:2010.07003* [cs.CL]
- [123] Jeonghoon Kim, Jung Hyun Lee, Sungdong Kim, and others. 2023. Memory-Efficient Fine-Tuning of Compressed Large Language Models via sub-4-bit Integer Quantization. *arXiv preprint arXiv:2305.14152* (2023).
- [124] Sehoon Kim, Coleman Hooper, Amir Gholami, and others. 2023. SqueezeLLM: Dense-and-Sparse Quantization. *arXiv preprint arXiv:2306.07629* (2023).
- [125] Sehoon Kim, Coleman Hooper, Thanakul Wattanawong, and others. 2023. Full stack optimization of transformer inference: a survey. *arXiv preprint arXiv:2302.14017* (2023).
- [126] Shine Kim, Yunho Jin, Gina Sohn, and others. 2021. Behemoth: a flash-centric training accelerator for extreme-scale DNNs. In *19th USENIX Conference on File and Storage Technologies*. 371–385.
- [127] Alexander Kirillov, Eric Mintun, Nikhila Ravi, and others. 2023. Segment Anything.
- [128] Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. 2020. Reformer: The Efficient Transformer. In *International Conference on Learning Representations*.
- [129] Aran Komatsuzaki, Joan Puigcerver, James Lee-Thorp, and others. 2023. Sparse Upcycling: Training Mixture-of-Experts from Dense Checkpoints. *arXiv:2212.05055* [cs.LG]
- [130] Rui Kong, Yuanchun Li, Qingtian Feng, and others. 2023. Serving MoE Models on Resource-constrained Edge Devices via Dynamic Expert Swapping. *arXiv preprint arXiv:2308.15030* (2023).
- [131] Zhenglun Kong, Peiyan Dong, Xiaolong Ma, and others. 2022. SPViT: Enabling Faster Vision Transformers via Soft Token Pruning. *arXiv:2112.13890* [cs.CV]
- [132] Juil Koo, Seungwoo Yoo, Minh Hieu Nguyen, and Minhyuk Sung. 2023. Salad: Part-level latent diffusion for 3d shape generation and manipulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 14441–14451.
- [133] Vijay Anand Korthikanti, Jared Casper, Sangkug Lym, and others. 2023. Reducing activation recomputation in large transformer models. *Proceedings of Machine Learning and Systems* (2023).
- [134] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, and others. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*. 611–626.
- [135] François Lagunas, Ella Charlaix, Victor Sanh, and Alexander M Rush. 2021. Block pruning for faster transformers. *arXiv preprint arXiv:2109.04838* (2021).
- [136] Joel Lamy-Poirier. 2023. Breadth-First Pipeline Parallelism. *Proceedings of Machine Learning and Systems* (2023).
- [137] Changhun Lee, Jungyu Jin, Taesu Kim, and others. 2023. OWQ: Lessons learned from activation outliers for weight quantization in large language models. *arXiv preprint arXiv:2306.02272* (2023).
- [138] Jung Hyun Lee, Jeonghoon Kim, Se Jung Kwon, and Dongsoo Lee. 2023. FlexRound: Learnable Rounding based on Element-wise Division for Post-Training Quantization. *arXiv preprint arXiv:2306.00317* (2023).
- [139] Katherine Lee, Daphne Ippolito, Andrew Nystrom, and others. 2022. Deduplicating Training Data Makes Language Models Better. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 8424–8445.
- [140] Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691* (2021).
- [141] Yaniv Leviathan, Matan Kalman, and Y. Matias. 2022. Fast Inference from Transformers via Speculative Decoding. In *International Conference on Machine Learning*. <https://api.semanticscholar.org/CorpusID:254096365>
- [142] Chunyuan Li, Zhe Gan, Zhengyuan Yang, and others. 2023. Multimodal Foundation Models: From Specialists to General-Purpose Assistants. *arXiv preprint arXiv:2309.10020* (2023).
- [143] Changlin Li, Bohan Zhuang, Guangrun Wang, and others. 2022. Automated Progressive Learning for Efficient Training of Vision Transformers. *arXiv:2203.14509* [cs.CV] <https://arxiv.org/abs/2203.14509>
- [144] David Li, Erik Nijkamp, and et al. 2022. Competition-level code generation with AlphaCode. *Science* 378, 6624 (2022), 1107–1114.
- [145] Kai Li, Runxuan Yang, and Xiaolin Hu. 2022. An efficient encoder-decoder architecture with top-down attention for speech separation. *arXiv preprint arXiv:2209.15200* (2022).
- [146] Shenggui Li, Hongxin Liu, Zhengda Bian, and others. 2023. Colossal-AI: A Unified Deep Learning System For Large-Scale Parallel Training. In *Proceedings of the 52nd International Conference on Parallel Processing*. Association for Computing Machinery, New York, NY, USA, 766–775.
- [147] Shenggui Li, Fuzhao Xue, Chaitanya Baranwal, and others. 2023. Sequence parallelism: Long sequence training from system perspective. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*. 2391–2404.

- [148] Sheng Li, Geng Yuan, Yue Dai, and others. 2022. SmartFRZ: An Efficient Training Framework using Attention-Based Layer Freezing. In *The Eleventh International Conference on Learning Representations*.
- [149] Xiuyu Li, Yijiang Liu, Long Lian, and others. 2023. Q-diffusion: Quantizing diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 17535–17545.
- [150] Yanyu Li, Huan Wang, Qing Jin, and others. 2023. SnapFusion: Text-to-Image Diffusion Model on Mobile Devices within Two Seconds. *arXiv preprint arXiv:2306.00980* (2023).
- [151] Yuanchun Li, Hao Wen, Weijun Wang, and others. 2024. Personal LLM Agents: Insights and Survey about the Capability, Efficiency and Security. *arXiv preprint arXiv:2401.05459* (2024).
- [152] Yanjing Li, Sheng Xu, Baochang Zhang, and others. 2022. Q-vit: Accurate and fully quantized low-bit vision transformer. *Advances in Neural Information Processing Systems* 35 (2022), 34451–34463.
- [153] Yixiao Li, Yifan Yu, Chen Liang, and others. 2023. Loftq: Lora-fine-tuning-aware quantization for large language models. *arXiv preprint arXiv:2310.08659* (2023).
- [154] Yixiao Li, Yifan Yu, Qingru Zhang, and others. 2023. LoSparse: Structured Compression of Large Language Models based on Low-Rank and Sparse Approximation. *arXiv preprint arXiv:2306.11222* (2023).
- [155] Yanyu Li, Geng Yuan, Yang Wen, and others. 2022. EfficientFormer: Vision Transformers at MobileNet Speed. *arXiv:2206.01191* [cs.CV]
- [156] Zhikai Li and Qingyi Gu. 2023. I-vit: Integer-only quantization for efficient vision transformer inference. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 17065–17075.
- [157] Chen Liang, Simiao Zuo, Qingru Zhang, and others. 2023. Less is more: Task-aware layer-wise distillation for language model compression. In *International Conference on Machine Learning*. PMLR, 20852–20867.
- [158] Baohao Liao, Yan Meng, and Christof Monz. 2023. Parameter-Efficient Fine-Tuning without Introducing New Latency. *arXiv preprint arXiv:2305.16742* (2023).
- [159] Baohao Liao, Shaomu Tan, and Christof Monz. 2023. Make Your Pre-trained Model Reversible: From Parameter to Memory Efficient Fine-Tuning. *arXiv preprint arXiv:2306.00477* (2023).
- [160] Ji Lin, Jiaming Tang, Haotian Tang, and others. 2023. AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration. *arXiv preprint arXiv:2306.00978* (2023).
- [161] Tianyang Lin, Yuxin Wang, Xiangyang Liu, and Xipeng Qiu. 2022. A survey of transformers. *AI Open* (2022).
- [162] Guisheng Liu, Yi Li, Zhengcong Fei, and others. 2023. Prefix-diffusion: A Lightweight Diffusion Model for Diverse Image Captioning. *arXiv preprint arXiv:2309.04965* (2023).
- [163] Jing Liu, Ruihao Gong, Xiuying Wei, and others. 2023. QLLM: Accurate and Efficient Low-Bitwidth Quantization for Large Language Models. *arXiv:2310.08041* [cs.CL]
- [164] Jihao Liu, Xin Huang, Jinliang Zheng, and others. 2023. MixMAE: Mixed and masked autoencoder for efficient pretraining of hierarchical vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6252–6261.
- [165] Juncai Liu, Jessie Hui Wang, and Yimin Jiang. 2023. Janus: A Unified Distributed Training Framework for Sparse Mixture-of-Experts Models. In *Proceedings of the ACM SIGCOMM 2023 Conference*. 486–498.
- [166] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. 2022. Pseudo numerical methods for diffusion models on manifolds. *arXiv preprint arXiv:2202.09778* (2022).
- [167] Shih-Yang Liu, Zechun Liu, and Kwang-Ting Cheng. 2023. Oscillation-free quantization for low-bit vision transformers. *arXiv preprint arXiv:2302.02210* (2023).
- [168] Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, and others. 2024. DoRA: Weight-Decomposed Low-Rank Adaptation. <http://arxiv.org/abs/2402.09353> *arXiv:2402.09353* [cs].
- [169] Xinyu Liu, Houwen Peng, Ningxin Zheng, and others. 2023. EfficientViT: Memory Efficient Vision Transformer with Cascaded Group Attention. *arXiv:2305.07027* [cs.CV] <https://arxiv.org/abs/2305.07027>
- [170] Xiaoxuan Liu, Lianmin Zheng, Dequan Wang, and others. 2022. GACT: Activation compressed training for generic network architectures. In *International Conference on Machine Learning*. PMLR, 14139–14152.
- [171] Yue Liu, Christos Matsoukas, Fredrik Strand, and others. 2023. Patchdropout: Economizing vision transformers using patch dropout. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 3953–3962.
- [172] Yuqi Liu, Luhui Xu, Pengfei Xiong, and Qin Jin. 2023. Token mixing: parameter-efficient transfer learning from image-language to video-language. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 1781–1789.
- [173] Zichang Liu, Aditya Desai, Fangshuo Liao, and others. 2023. Scissorhands: Exploiting the Persistence of Importance Hypothesis for LLM KV Cache Compression at Test Time. *arXiv:2305.17118* [cs.LG]
- [174] Zechun Liu, Barlas Oguz, Changsheng Zhao, and others. 2023. LLM-QAT: Data-Free Quantization Aware Training for Large Language Models. *arXiv preprint arXiv:2305.17888* (2023).
- [175] Zichang Liu, Jue Wang, Tri Dao, and others. 2023. Deja vu: Contextual sparsity for efficient llms at inference time. In *International Conference on Machine Learning*. PMLR, 22137–22176.
- [176] Cheng Lu, Yuhao Zhou, Fan Bao, and others. 2022. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems* 35 (2022), 5775–5787.

- [177] Zhenyan Lu, Xiang Li, Dongqi Cai, and others. 2024. Small Language Models: Survey, Measurements, and Insights. *arXiv preprint arXiv:2409.15790* (2024).
- [178] Xinyin Ma, Gongfan Fang, and Xinchao Wang. 2023. LLM-Pruner: On the Structural Pruning of Large Language Models. *arXiv preprint arXiv:2305.11627* (2023).
- [179] Xuezhe Ma, Chunting Zhou, Xiang Kong, and others. 2022. Mega: moving average equipped gated attention. *arXiv preprint arXiv:2209.10655* (2022).
- [180] Maciej A Mazurowski and others. 2023. SAM (Segment Anything Model) for Medical Image Segmentation: Applications, Opportunities and Limitations. *arXiv preprint arXiv:2305.02652* (2023).
- [181] Sachin Mehta and Mohammad Rastegari. 2022. MobileViT: Light-weight, General-purpose, and Mobile-friendly Vision Transformer. *arXiv:2110.02178* [cs.CV]
- [182] Fanxu Meng, Zhaohui Wang, and Muhan Zhang. 2024. PiSSA: Principal Singular Values and Singular Vectors Adaptation of Large Language Models. <http://arxiv.org/abs/2404.02948> *arXiv:2404.02948* [cs].
- [183] Lingchen Meng, Hengduo Li, Bor-Chun Chen, and others. 2021. AdaViT: Adaptive Vision Transformers for Efficient Image Recognition. *arXiv:2111.15668* [cs.CV]
- [184] Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, and others. 2023. Towards efficient generative large language model serving: A survey from algorithms to systems. *arXiv preprint arXiv:2312.15234* (2023).
- [185] Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, and others. 2023. SpecInfer: Accelerating Generative Large Language Model Serving with Speculative Inference and Token Tree Verification. *arXiv:2305.09781* [cs.CL]
- [186] Xupeng Miao, Chunan Shi, Jiangfei Duan, and others. 2023. SpotServe: Serving Generative Large Language Models on Preemptible Instances. *arXiv preprint arXiv:2311.15566* (2023).
- [187] Xupeng Miao, Yujie Wang, Youhe Jiang, and others. 2023. Galvatron: Efficient transformer training over multiple gpus using automatic parallelism. *Proceedings of the VLDB Endowment* (2023).
- [188] Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2021. Metaicl: Learning to learn in context. *arXiv preprint arXiv:2110.15943* (2021).
- [189] Amirkeivan Mohtashami and Martin Jaggi. 2023. Landmark Attention: Random-Access Infinite Context Length for Transformers. *arXiv preprint arXiv:2305.16300* (2023).
- [190] Basil Mustafa, Carlos Riquelme, Joan Puigcerver, and others. 2022. Multimodal contrastive learning with limoe: the language-image mixture of experts. *Advances in Neural Information Processing Systems* 35 (2022), 9564–9576.
- [191] Pranav Ajit Nair, Sukomal Pal, and Pradeepika Verm. 2023. Domain Aligned Prefix Averaging for Domain Generalization in Abstractive Summarization. *arXiv preprint arXiv:2305.16820* (2023).
- [192] Deepak Narayanan, Mohammad Shoeybi, Jared Casper, and others. 2021. Efficient large-scale language model training on gpu clusters using megatron-lm. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. 1–15.
- [193] Piotr Nawrot, Jan Chorowski, Adrian Łańcucki, and Edoardo M Ponti. 2022. Efficient Transformers with Dynamic Token Pooling. *arXiv preprint arXiv:2211.09761* (2022).
- [194] Alexander Quinn Nichol and Prafulla Dhariwal. 2021. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*. PMLR, 8162–8171.
- [195] Xiaonan Nie, Xupeng Miao, Zilong Wang, and others. 2023. FlexMoE: Scaling Large-scale Sparse Pre-trained Model Training via Dynamic Device Placement. *Proceedings of the ACM on Management of Data* (2023), 1–19.
- [196] Wei Niu, Jiexiong Guan, Yanzhi Wang, and others. 2021. DNNFusion: accelerating deep neural networks execution with advanced operator fusion. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*. 883–898.
- [197] Georgii Sergeevich Novikov, Daniel Bershtsky, Julia Gusk, and others. 2023. Few-bit backward: Quantized gradients of activation functions for memory footprint reduction. In *International Conference on Machine Learning*. PMLR, 26363–26381.
- [198] OpenAI. 2023. GPT-4 Technical Report. *arXiv:2303.08774* [cs.CL]
- [199] Antonio Orvieto, Samuel L Smith, Albert Gu, and others. 2023. Resurrecting recurrent neural networks for long sequences. *arXiv preprint arXiv:2303.06349* (2023).
- [200] Kazuki Osawa, Shigang Li, and Torsten Hoefler. 2023. PipeFisher: Efficient Training of Large Language Models Using Pipelining and Fisher Information Matrices. *Proceedings of Machine Learning and Systems* (2023).
- [201] Junting Pan, Adrian Bulat, Fuwen Tan, and others. 2022. EdgeViTs: Competing Light-weight CNNs on Mobile Devices with Vision Transformers. *arXiv:2205.03436* [cs.CV] <https://arxiv.org/abs/2205.03436>
- [202] Junting Pan, Ziyi Lin, Xiatian Zhu, and others. 2022. St-adapter: Parameter-efficient image-to-video transfer learning. *Advances in Neural Information Processing Systems* 35 (2022), 26462–26477.
- [203] Zizheng Pan, Peng Chen, Haoyu He, and others. 2021. Mesa: A memory-saving training framework for transformers. *arXiv preprint arXiv:2111.11124* (2021).
- [204] Jay H Park, Gyeongchan Yun, M Yi Chang, and others. 2020. HetPipe: Enabling large DNN training on (whimpy) heterogeneous GPU clusters through integration of pipelined model parallelism and data parallelism. In *USENIX Annual Technical Conference*. 307–321.
- [205] Pratyush Patel, Esha Choukse, Chaojie Zhang, and others. 2023. Splitwise: Efficient generative LLM inference using phase splitting. *arXiv preprint arXiv:2311.18677* (2023).

- [206] William Peebles and Saining Xie. 2023. Scalable Diffusion Models with Transformers. *arXiv:2212.09748* [cs.CV] <https://arxiv.org/abs/2212.09748>
- [207] Bo Peng, Eric Alcaide, Quentin Anthony, and others. 2023. RWKV: Reinventing RNNs for the Transformer Era. *arXiv preprint arXiv:2305.13048* (2023).
- [208] Maciej Pióro, Kamil Ciebiera, Krystian Król, and others. 2024. Moe-mamba: Efficient selective state space models with mixture of experts. *arXiv preprint arXiv:2401.04081* (2024).
- [209] Koutilya Pnvr, Bharat Singh, Pallabi Ghosh, and others. 2023. LD-ZNet: A latent diffusion approach for text-based image segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4157–4168.
- [210] Michael Poli, Stefano Massaroli, Eric Nguyen, and others. 2023. Hyena hierarchy: Towards larger convolutional language models. *arXiv preprint arXiv:2302.10866* (2023).
- [211] Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, and others. 2023. Efficiently scaling transformer inference. *Proceedings of Machine Learning and Systems* 5 (2023).
- [212] Guanghui Qin, Corby Rosset, Ethan C. Chau, and others. 2023. Nugget 2D: Dynamic Contextual Compression for Scaling Decoder-only Language Models. *arXiv:2310.02409* [cs.CL]
- [213] Yujia Qin, Yankai Lin, Jing Yi, and others. 2022. Knowledge Inheritance for Pre-trained Language Models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 3921–3937.
- [214] Alec Radford, Jong Wook Kim, Chris Hallacy, and others. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- [215] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. 2020. Zero: Memory optimizations toward training trillion parameter models. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 1–16.
- [216] Yongming Rao, Wenliang Zhao, Benlin Liu, and others. 2021. DynamicViT: Efficient Vision Transformers with Dynamic Token Sparsification. *arXiv:2106.02034* [cs.CV]
- [217] Nir Ratner, Yoav Levine, Yonatan Belinkov, and others. 2023. Parallel context windows for large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 6383–6402.
- [218] Jie Ren, Samyam Rajbhandari, Reza Yazdani Aminabadi, and others. 2021. ZeRO-Offload: Democratizing Billion-Scale model training. In *USENIX Annual Technical Conference*. 551–564.
- [219] Carlos Riquelme, Joan Puigcerver, Basil Mustafa, and others. 2021. Scaling vision with sparse mixture of experts. *Advances in Neural Information Processing Systems* 34 (2021), 8583–8595.
- [220] Robin Rombach, Andreas Blattmann, Dominik Lorenz, and others. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10684–10695.
- [221] Rajarshi Saha, Varun Srivastava, and Mert Pilanci. 2023. Matrix Compression via Randomized Low Rank and Low Precision Factorization. *arXiv preprint arXiv:2310.11028* (2023).
- [222] Michael Santacrose, Zixin Wen, Yelong Shen, and Yuanzhi Li. 2023. What Matters In The Structured Pruning of Generative Language Models? *arXiv preprint arXiv:2302.03773* (2023).
- [223] Sheng Shen, Pete Walsh, Kurt Keutzer, and others. 2022. Staged training for transformer language models. In *International Conference on Machine Learning*. PMLR, 19893–19908.
- [224] Ying Sheng, Shiyi Cao, Dacheng Li, and others. 2023. S-lora: Serving thousands of concurrent lora adapters. *arXiv preprint arXiv:2311.03285* (2023).
- [225] Ying Sheng, Lianmin Zheng, Binhang Yuan, and others. 2023. Flexgen: High-throughput generative inference of large language models with a single gpu. In *International Conference on Machine Learning*. PMLR, 31094–31116.
- [226] Dachuan Shi, Chaofan Tao, Ying Jin, and others. 2023. Upop: Unified and progressive pruning for compressing vision-language transformers. *arXiv preprint arXiv:2301.13741* (2023).
- [227] Jiaming Song, Chenlin Meng, and Stefano Ermon. 2020. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502* (2020).
- [228] Jaeyong Song, Jinkyu Yim, Jaewon Jung, and others. 2023. Optimus-CC: Efficient Large NLP Model Training with 3D Parallelism Aware Communication Compression. In *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*. 560–573.
- [229] Yixin Song, Zeyu Mi, Haotong Xie, and Haibo Chen. 2023. PowerInfer: Fast Large Language Model Serving with a Consumer-grade GPU. *arXiv:2312.12456* [cs.LG]
- [230] Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. 2023. A Simple and Effective Pruning Approach for Large Language Models. *arXiv preprint arXiv:2306.11695* (2023).
- [231] Yutao Sun, Li Dong, Shaohan Huang, and others. 2023. Retentive network: A successor to transformer for large language models. *arXiv preprint arXiv:2307.08621* (2023).
- [232] Ziteng Sun, Ananda Theertha Suresh, Jae Hun Ro, and others. 2023. SpecTr: Fast Speculative Decoding via Optimal Transport. In *Workshop on Efficient Systems for Foundation Models @ ICML2023*. <https://openreview.net/forum?id=d0mGsaheuT>
- [233] Thierry Tambe, Coleman Hooper, Lillian Pentecost, and others. 2021. EdgeBERT: Sentence-Level Energy Optimizations for Latency-Aware Multi-Task NLP Inference. *arXiv:2011.14203* [cs.AR] <https://arxiv.org/abs/2011.14203>
- [234] Chaofan Tao, Lu Hou, Wei Zhang, and others. 2022. Compression of generative pre-trained language models via quantization. *arXiv preprint arXiv:2203.10705* (2022).

- [235] John Thorpe, Pengzhan Zhao, Jonathan Eyolfson, and others. 2023. Bamboo: Making Preemptible Instances Resilient for Affordable Training of Large DNNs. In *20th USENIX Symposium on Networked Systems Design and Implementation*. 497–513.
- [236] Lin Tian, Xiuzhen Zhang, and Jey Han Lau. 2023. MetaTroll: Few-shot Detection of State-Sponsored Trolls with Transformer Adapters. In *Proceedings of the ACM Web Conference 2023*. 1743–1753.
- [237] Inar Timiryasov and Jean-Loup Tastet. 2023. Baby Llama: knowledge distillation from an ensemble of teachers trained on a small dataset with no performance penalty. *arXiv preprint arXiv:2308.02019* (2023).
- [238] Katrin Tomanek, Vicky Zayats, Dirk Padfield, and others. 2021. Residual adapters for parameter-efficient ASR adaptation to atypical and accented speech. *arXiv preprint arXiv:2109.06952* (2021).
- [239] Hugo Touvron, M Cord, M Douze, and others. 2012. Training data-efficient image transformers and distillation through attention (2020). doi: 10.48550. *arxiv* (2012).
- [240] Hugo Touvron, Louis Martin, Kevin Stone, and others. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* (2023).
- [241] Hieu Tran, Zhichao Yang, Zonghai Yao, and Hong Yu. 2023. BioInstruct: Instruction Tuning of Large Language Models for Biomedical Natural Language Processing. *arXiv preprint arXiv:2310.19975* (2023).
- [242] Alexander Tsvetkov and Alon Kipnis. 2023. EntropyRank: Unsupervised Keyphrase Extraction via Side-Information Optimization for Language Model-based Text Compression. In *ICML 2023 Workshop Neural Compression: From Information Theory to Applications*.
- [243] Chandra Shekhara Kaushik Valmeekam, Krishna Narayanan, Dileep Kalathil, and others. 2023. LLMZip: Lossless Text Compression using Large Language Models. *arXiv:2306.04050* [cs.IT]
- [244] Pavan Kumar Anasosalu Vasu, James Gabriel, Jeff Zhu, and others. 2023. FastViT: A Fast Hybrid Vision Transformer using Structural Reparameterization. *arXiv:2303.14189* [cs.CV] <https://arxiv.org/abs/2303.14189>
- [245] Pavan Kumar Anasosalu Vasu, Hadi Pouransari, Fartash Faghri, and others. 2024. MobileCLIP: Fast Image-Text Models through Multi-Modal Reinforced Training. *arXiv:2311.17049* [cs.CV] <https://arxiv.org/abs/2311.17049>
- [246] Ashish Vaswani, Noam Shazeer, Niki Parmar, and others. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [247] Yixin Wan, Kuan-Hao Huang, and Kai-Wei Chang. 2023. PIP: Parse-Instructed Prefix for Syntactically Controlled Paraphrase Generation. *arXiv preprint arXiv:2305.16701* (2023).
- [248] Zhongwei Wan, Xin Wang, Che Liu, and others. 2023. Efficient large language models: A survey. *arXiv preprint arXiv:2312.03863* 1 (2023).
- [249] Hongyu Wang, Shuming Ma, Li Dong, and others. 2023. BitNet: Scaling 1-bit Transformers for Large Language Models. *arXiv preprint arXiv:2310.11453* (2023).
- [250] Hanrui Wang, Zhekai Zhang, and Song Han. 2021. SpAtten: Efficient Sparse Attention Architecture with Cascade Token and Head Pruning. In *2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE. <https://doi.org/10.1109/hpca51647.2021.00018>
- [251] Peihao Wang, Rameswar Panda, Lucas Torroba Hennigen, and others. 2022. Learning to Grow Pretrained Models for Efficient Transformer Training. In *The Eleventh International Conference on Learning Representations*.
- [252] Sinong Wang, Belinda Z Li, Madian Khabsa, and others. 2020. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768* (2020).
- [253] Shibo Wang, Jinliang Wei, Amit Sabne, and others. 2023. Overlap communication with dependent computation via decomposition in large deep learning models. In *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1*. 93–106.
- [254] Xuan Wang, Guan hong Wang, Wenhao Chai, and others. 2023. User-Aware Prefix-Tuning is a Good Learner for Personalized Image Captioning. *arXiv preprint arXiv:2312.04793* (2023).
- [255] Xudong Wang, Li Lyna Zhang, Yang Wang, and Mao Yang. 2022. Towards efficient vision transformer inference: A first study of transformers on mobile devices. In *Proceedings of the 23rd Annual International Workshop on Mobile Computing Systems and Applications*. 1–7.
- [256] Xinlong Wang, Xiaosong Zhang, Yue Cao, and others. 2023. SegGPT: Segmenting Everything In Context.
- [257] Yaqing Wang, Sahaj Agarwal, Subhabrata Mukherjee, and others. 2022. AdaMix: Mixture-of-adaptations for parameter-efficient model tuning. *arXiv preprint arXiv:2210.17451* (2022).
- [258] Zhuang Wang, Zhen Jia, Shuai Zheng, and others. 2023. Gemini: Fast failure recovery in distributed training with in-memory checkpoints. In *Proceedings of the 29th Symposium on Operating Systems Principles*. 364–381.
- [259] Jason Wei, Xuezhi Wang, Dale Schuurmans, and others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems* 35 (2022), 24824–24837.
- [260] Xiuying Wei, Yunchen Zhang, Yuhang Li, and others. 2023. Outlier Suppression+: Accurate quantization of large language models by equivalent and optimal shifting and scaling. *arXiv:2304.09145* [cs.CL]
- [261] Xiuying Wei, Yunchen Zhang, Xiangguo Zhang, and others. 2022. Outlier suppression: Pushing the limit of low-bit transformer language models. *Advances in Neural Information Processing Systems* 35 (2022), 17402–17414.
- [262] Hao Wen, Yuanchun Li, Guohong Liu, and others. 2023. Empowering LLM to use Smartphone for Intelligent Task Automation. *arXiv preprint arXiv:2308.15272* (2023).

- [263] Qizhen Weng, Wencong Xiao, Yinghao Yu, and others. 2022. MLaaS in the wild: Workload analysis and scheduling in Large-Scale heterogeneous GPU clusters. In *19th USENIX Symposium on Networked Systems Design and Implementation*. 945–960.
- [264] David Wingate, Mohammad Shoeybi, and Taylor Sorensen. 2022. Prompt Compression and Contrastive Conditioning for Controllability and Toxicity Reduction in Language Models. *arXiv:2210.03162* [cs.CL]
- [265] Bingyang Wu, Yinmin Zhong, Zili Zhang, and others. 2023. Fast Distributed Inference Serving for Large Language Models. *arXiv preprint arXiv:2305.05920* (2023).
- [266] Bingyang Wu, Ruidong Zhu, Zili Zhang, and others. 2024. {dLoRA}: Dynamically Orchestrating Requests and Adapters for {LoRA} {LLM} Serving. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*. 911–927.
- [267] Carole-Jean Wu, Ramya Raghavendra, Udit Gupta, and others. 2022. Sustainable ai: Environmental implications, challenges and opportunities. *Proceedings of Machine Learning and Systems* 4 (2022), 795–813.
- [268] Xiaoxia Wu, Cheng Li, Reza Yazdani Aminabadi, and others. 2023. Understanding INT4 quantization for language models: latency speedup, composability, and failure cases. In *International Conference on Machine Learning*. PMLR, 37524–37539.
- [269] Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. 2023. Sheared llama: Accelerating language model pre-training via structured pruning. *arXiv preprint arXiv:2310.06694* (2023).
- [270] Zhuofan Xia, Xuran Pan, Shiji Song, and others. 2022. Vision Transformer with Deformable Attention. *arXiv:2201.00520* [cs.CV]
- [271] Guangxuan Xiao, Ji Lin, Mickael Seznec, and others. 2023. SmoothQuant: Accurate and Efficient Post-Training Quantization for Large Language Models. *arXiv:2211.10438* [cs.CL]
- [272] Guangxuan Xiao, Yuandong Tian, Beidi Chen, and others. 2023. Efficient Streaming Language Models with Attention Sinks. *arXiv preprint arXiv:2309.17453* (2023).
- [273] Ji Xin, Raphael Tang, Jaejun Lee, and others. 2020. DeeBERT: Dynamic early exiting for accelerating BERT inference. *arXiv preprint arXiv:2004.12993* (2020).
- [274] Daliang Xu, Wangsong Yin, Xin Jin, and others. 2023. LLMcad: Fast and Scalable On-device Large Language Model Inference. *arXiv:2309.04255* [cs.NI]
- [275] Daliang Xu, Hao Zhang, Liming Yang, and others. 2024. Empowering 1000 tokens/second on-device llm prefilling with mllm-npu. *arXiv preprint arXiv:2407.05858* (2024).
- [276] Guanyu Xu, Jiawei Hao, Li Shen, and others. 2023. LGViT: Dynamic Early Exiting for Accelerating Vision Transformer. In *Proceedings of the 31st ACM International Conference on Multimedia*. 9103–9114.
- [277] Mengwei Xu, Jiawei Liu, Yuanqiang Liu, and others. 2019. A first look at deep learning apps on smartphones. In *The World Wide Web Conference*. 2125–2136.
- [278] Mingxue Xu, Yao Lei Xu, and Danilo P Mandic. 2023. Tensorgpt: Efficient compression of the embedding layer in llms based on the tensor-train decomposition. *arXiv preprint arXiv:2307.00526* (2023).
- [279] Mengwei Xu, Wangsong Yin, Dongqi Cai, and others. 2024. A survey of resource-efficient llm and multimodal foundation models. *arXiv preprint arXiv:2401.08092* (2024).
- [280] Yuhui Xu, Lingxi Xie, Xiaotao Gu, and others. 2023. Qa-lora: Quantization-aware low-rank adaptation of large language models. *arXiv preprint arXiv:2309.14717* (2023).
- [281] Yu Yan, Weizhen Qi, Yeyun Gong, and others. 2020. ProphetNet: Predicting Future N-gram for Sequence-to-Sequence Pre-training. *ArXiv abs/2001.04063* (2020). <https://api.semanticscholar.org/CorpusID:210164665>
- [282] Bufang Yang, Lixing He, Neiweng Ling, and others. 2023. EdgeFM: Leveraging Foundation Model for Open-set Learning on the Edge. *arXiv preprint arXiv:2311.10986* (2023).
- [283] Jingfeng Yang, Hongye Jin, Ruixiang Tang, and others. 2023. Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond. *arXiv:2304.13712* [cs.CL]
- [284] Yuedong Yang, Hung-Yueh Chiang, Guihong Li, and others. 2023. Efficient Low-rank Backpropagation for Vision Transformer Adaptation. *arXiv preprint arXiv:2309.15275* (2023).
- [285] Yuting Yang, Wenqiang Lei, Pei Huang, and others. 2023. A Dual Prompt Learning Framework for Few-Shot Dialogue State Tracking. In *Proceedings of the ACM Web Conference 2023*. 1468–1477.
- [286] Shih yang Liu, Zechun Liu, Xijie Huang, and others. 2023. LLM-FP4: 4-Bit Floating-Point Quantized Transformers. *arXiv:2310.16836* [cs.CL]
- [287] Zhewei Yao, Reza Yazdani Aminabadi, Minjia Zhang, and others. 2022. Zeroquant: Efficient and affordable post-training quantization for large-scale transformers. *Advances in Neural Information Processing Systems* 35 (2022), 27168–27183.
- [288] Deming Ye, Yankai Lin, Yufei Huang, and Maosong Sun. 2021. TR-BERT: Dynamic Token Reduction for Accelerating BERT Inference. *arXiv:2105.11618* [cs.CL]
- [289] Rongjie Yi, Liwei Guo, Shiyun Wei, and others. 2023. Edgemoe: Fast on-device inference of moe-based large language models. *arXiv preprint arXiv:2308.14352* (2023).
- [290] Rongjie Yi, Xiang Li, Weikai Xie, and others. 2024. PhoneLM: an Efficient and Capable Small Language Model Family through Principled Pre-training. *arXiv preprint arXiv:2411.05046* (2024).
- [291] Hongxu Yin, Arash Vahdat, Jose Alvarez, and others. 2022. AdaViT: Adaptive Tokens for Efficient Vision Transformer. *arXiv:2112.07658* [cs.CV]
- [292] Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, and others. 2022. Orca: A distributed serving system for {Transformer-Based} generative models. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*. 521–538.

- [293] Jinliang Yuan, Chen Yang, Dongqi Cai, and others. 2024. Mobile Foundation Model as Firmware. In *Proceedings of the 30th Annual International Conference on Mobile Computing and Networking (ACM MobiCom24, Vol. 33)*. ACM. <https://doi.org/10.1145/3636534.3649361>
- [294] Zhihang Yuan, Lin Niu, Jiawei Liu, and others. 2023. RPTQ: Reorder-based Post-training Quantization for Large Language Models. *arXiv:2304.01089* [cs.CL]
- [295] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, and others. 2020. Big bird: Transformers for longer sequences. *Advances in neural information processing systems* 33 (2020), 17283–17297.
- [296] Syed Waqas Zamir, Aditya Arora, Salman Khan, and others. 2022. Restormer: Efficient Transformer for High-Resolution Image Restoration. *arXiv:2111.09881* [cs.CV]
- [297] Shulin Zeng, Jun Liu, Guohao Dai, and others. 2024. FlightLLM: Efficient Large Language Model Inference with a Complete Mapping Flow on FPGAs. *arXiv:2401.03868* [cs.AR] <https://arxiv.org/abs/2401.03868>
- [298] Mingshu Zhai, Jiaao He, Zixuan Ma, and others. 2023. SmartMoE: Efficiently Training Sparsely-Activated Models through Combining Offline and Online Parallelization. In *USENIX Annual Technical Conference*. 961–975.
- [299] Shuangfei Zhai, Walter Talbott, Nitish Srivastava, and others. 2021. An attention free transformer. *arXiv preprint arXiv:2105.14103* (2021).
- [300] Yujia Zhai, Chengquan Jiang, Leyuan Wang, and others. 2023. ByteTransformer: A high-performance transformer boosted for variable-length inputs. In *2023 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*. IEEE, 344–355.
- [301] Jinchao Zhang, Jue Wang, Huan Li, and others. 2023. Draft & Verify: Lossless Large Language Model Acceleration via Self-Speculative Decoding. *ArXiv abs/2309.08168* (2023). <https://api.semanticscholar.org/CorpusID:262013673>
- [302] Mingyang Zhang, Chunhua Shen, Zhen Yang, and others. 2023. Pruning Meets Low-Rank Parameter-Efficient Fine-Tuning. *arXiv preprint arXiv:2305.18403* (2023).
- [303] Qiyang Zhang, Xiangying Che, Yijie Chen, and others. 2023. A comprehensive deep learning library benchmark and optimal library selection. *IEEE Transactions on Mobile Computing* (2023).
- [304] Qingru Zhang, Simiao Zuo, Chen Liang, and others. 2022. Platon: Pruning large transformer models with upper confidence bound of weight importance. In *International Conference on Machine Learning*. PMLR, 26809–26823.
- [305] Zining Zhang, Bingsheng He, and Zhenjie Zhang. 2023. Practical Edge Kernels for Integer-Only Vision Transformers Under Post-training Quantization. *Proceedings of Machine Learning and Systems* 5 (2023).
- [306] Zhengyan Zhang, Yankai Lin, Zhiyuan Liu, and others. 2022. MoEification: Transformer Feed-forward Layers are Mixtures of Experts. In *Findings of the Association for Computational Linguistics: ACL 2022*. 877–890.
- [307] Zhenyu Zhang, Ying Sheng, Tianyi Zhou, and others. 2023. H₂O: Heavy-Hitter Oracle for Efficient Generative Inference of Large Language Models. *arXiv:2306.14048* [cs.LG]
- [308] Zhen Zhang, Shuai Zheng, Yida Wang, and others. 2023. MiCS: Near-linear scaling for training gigantic model on public cloud. *Proceedings of the VLDB Endowment* (2023).
- [309] Lulu Zhao, Fujia Zheng, Weihao Zeng, and others. 2022. Domain-oriented prefix-tuning: Towards efficient and generalizable fine-tuning for zero-shot dialogue summarization. *arXiv preprint arXiv:2204.04362* (2022).
- [310] Yanli Zhao, Andrew Gu, Rohan Varma, and others. 2023. Pytorch FSDP: experiences on scaling fully sharded data parallel. *arXiv preprint arXiv:2304.11277* (2023).
- [311] Yilong Zhao, Chien-Yu Lin, Kan Zhu, and others. 2023. Atom: Low-bit Quantization for Efficient and Accurate LLM Serving. *arXiv:2310.19102* [cs.LG]
- [312] Yang Zhao, Yanwu Xu, Zhisheng Xiao, and others. 2024. MobileDiffusion: Instant Text-to-Image Generation on Mobile Devices. *arXiv:2311.16567* [cs.CV] <https://arxiv.org/abs/2311.16567>
- [313] Lianmin Zheng, Zhuohan Li, Hao Zhang, and others. 2022. Alpa: Automating inter-and {Intra-Operator} parallelism for distributed deep learning. In *16th USENIX Symposium on Operating Systems Design and Implementation*. 559–578.
- [314] Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, and others. 2023. Efficiently Programming Large Language Models using SGLang. *arXiv preprint arXiv:2312.07104* (2023).
- [315] Yinmin Zhong, Shengyu Liu, Junda Chen, and others. 2024. {DistServe}: Disaggregating Prefill and Decoding for Goodput-optimized Large Language Model Serving. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*. 193–210.
- [316] Wangchunshu Zhou, Canwen Xu, Tao Ge, and others. 2020. Bert loses patience: Fast and robust inference with early exit. *Advances in Neural Information Processing Systems* 33 (2020), 18330–18341.
- [317] Lingfeng Zhu and others. 2023. A comprehensive survey of large language models: Application to natural language processing tasks. *arXiv preprint arXiv:2304.12345* (2023).
- [318] Xunyu Zhu, Jian Li, Yong Liu, and others. 2023. A survey on model compression for large language models. *arXiv preprint arXiv:2308.07633* (2023).
- [319] Yonghao Zhuang, Hexu Zhao, Lianmin Zheng, and others. 2023. On optimizing the communication of model parallelism. *Proceedings of Machine Learning and Systems* (2023).
- [320] Bojia Zi, Xianbiao Qi, Lingzhi Wang, and others. 2023. Delta-LoRA: Fine-Tuning High-Rank Parameters with the Delta of Low-Rank Matrices. <http://arxiv.org/abs/2309.02411> *arXiv:2309.02411* [cs].