

# New Solutions on LLM Acceleration, Optimization, and Application

Yingbing Huang<sup>1</sup>, Lily Jiaxin Wan<sup>1</sup>, Hanchen Ye<sup>1</sup>, Manvi Jha<sup>1</sup>, Jinghua Wang<sup>1</sup>, Yuhong Li<sup>1</sup>,  
Xiaofan Zhang<sup>2</sup>, Deming Chen<sup>1</sup>

{yh21, wan25, hanchen8, manvij2, jinghua3, leeyh, dchen}@illinois.edu, xiaofanz@google.com

<sup>1</sup>University of Illinois Urbana-Champaign, <sup>2</sup>Google

## Abstract

Large Language Models (LLMs) have become extremely potent instruments with exceptional capacities for comprehending and producing human-like text in a wide range of applications. However, the increasing size and complexity of LLMs present significant challenges in both training and deployment, leading to substantial computational and storage costs as well as heightened energy consumption. In this paper, we provide a review of recent advancements and research directions aimed at addressing these challenges and enhancing the efficiency of LLM-based systems. We begin by discussing algorithm-level acceleration techniques focused on optimizing LLM inference speed and resource utilization. We also explore LLM-hardware co-design strategies with a vision to improve system efficiency by tailoring hardware architectures to LLM requirements. Further, we delve into LLM-to-accelerator compilation approaches, which involve customizing hardware accelerators for efficient LLM deployment. Finally, as a case study to leverage LLMs for assisting circuit design, we examine LLM-aided design methodologies for an important task: High-Level Synthesis (HLS) functional verification, by creating a new dataset that contains a large number of buggy and bug-free codes, which can be essential for training LLMs to specialize on HLS verification and debugging. For each aspect mentioned above, we begin with a detailed background study, followed by the presentation of several novel solutions proposed to overcome specific challenges. We then outline future research directions to drive further advancements. Through these efforts, we aim to pave the way for more efficient and scalable deployment of LLMs across a diverse range of applications.

## 1 Introduction

In recent years, Large Language Models (LLMs) have emerged as powerful tools across various domains, revolutionizing natural language processing, information retrieval, LLM-aided design, and others. The ability of LLMs to understand, generate, and manipulate human language has propelled them to the forefront of research and applications in various industries. These models, trained on vast amounts of

text data, demonstrate unparalleled proficiency in tasks such as text generation, translation, summarization, sentiment analysis, and more [47][36]. Additionally, there are ongoing efforts to train LLMs with groundbreaking multimodal capabilities, encompassing both visual and speech understanding [27, 28, 55, 56]. They have been successfully applied in various applications, including virtual assistants, content generation, question-answering systems, and recommendation systems. The versatility and effectiveness of LLMs have made them indispensable tools in various industries, driving innovation and accelerating progress in artificial intelligence.

In domains such as LLM-aided design, large language models have been utilized for a variety of tasks, including high-level synthesis, hardware description generation, and functional verification, significantly streamlining the design process and reducing time-to-market for hardware designs [38]. For instance, ChipNeMo [32] enhances LLaMA2 with domain-specific optimizations for more efficient hardware design. AutoChip [48] focuses on automating HDL generation using feedback from LLMs, thereby improving the iterative design process. ChatEDA [53] leverages LLMs to create an autonomous agent for EDA, while VeriGen [49] specializes in generating Verilog code. Additionally, DIVAS [41] provides an end-to-end framework for SoC security analysis, and another approach utilizes LLMs to fix hardware security bugs [1]. Moreover, large language models are also being utilized for automated code generation for information technology tasks in YAML, further showcasing their versatility and efficiency [42].

However, the widespread adoption of these models has been hindered by their demanding computational requirements, which often result in slow response times and high costs for hardware and energy. Addressing these challenges is crucial to fully harnessing the potential of LLMs and unlocking their benefits in real-world applications. To address these challenges, this research paper explores a comprehensive approach to optimize LLMs at algorithm, hardware, compiler, and design-automation levels, aiming to enhance their efficiency and performance across diverse applications.

Previous works explore various algorithmic strategies aimed at decreasing the inference latency of LLMs. We begin by examining various methods for optimizing parameter utilization in large language models. These methods include techniques such as early exiting and layer skipping

All the authors contributed equally.

The presented work is an expanded and more comprehensive study based on our invited DAC'24 paper with the same title and co-authors.

[13], which help reduce computational overhead, as well as contextual sparsity, which dynamically prunes irrelevant parameters during inference [33]. Additionally, previous works explore a Mixture of Experts (MoE) approach, which distribute computation across multiple sub-models to enhance efficiency and scalability [12]. We then delve into optimization techniques for Key-Value (KV) cache, which is crucial for complex tasks like chain-of-thought reasoning and information retrieval. Additionally, we discuss advancements in parallel decoding [29], including techniques for aligning small and large model predictions, multi-token predictions, and parallel processing capabilities. Building upon this background study, we propose two novel approaches: a parallel decoding framework called Medusa [5], which employs multiple decoding heads coupled with an optimized tree-based decoding strategy, and SnapKV, a method for effectively reducing KV cache size [31]. Our experimental results demonstrate significant speedups in inference time without compromising generation quality, along with improved memory efficiency. Finally, we outline future directions for tailored algorithmic optimization, advancements in KV compression, and tackling the computational load from speculative decoding, aiming to boost LLM efficiency and effectiveness in practical applications.

LLM-hardware co-design aims to customize hardware architectures to meet the demands of LLMs while providing insights to optimize LLM architectures [15]. Previously, we proposed an LLM-hardware co-design framework called AutoDistill [64], which integrates model compression, transformer architecture exploration, and multi-objective optimization to produce student models with lower inference latency and smaller sizes while maintaining high accuracy. Moreover, a pruning-aware quantization strategy that combines two effective LLM compression methods, pruning, and quantization, to optimize LLM architectures for hardware efficiency has been proposed [51]. Furthermore, we explore the potential of reconfigurable and heterogeneous hardware for LLMs, aiming to dynamically adjust hardware architectures to accommodate the latest LLM advancements and model compression methods, thereby enhancing both model quality and hardware efficiency.

The demand for efficient hardware accelerators for deep neural networks has led to a new direction of using High-Level Synthesis (HLS) frameworks [7] [9] to quickly translate model architectures into hardware implementations. However, exploring the vast design space effectively to achieve optimal solutions remains a significant challenge. We summarize two novel compilation frameworks published previously by us: ScaleHLS [21, 57] and HIDA [60]. ScaleHLS leverages the MLIR infrastructure [26] for scalable High-Level Synthesis, optimizing hardware designs with a Design Space Exploration engine by performing multi-level transformations. As far as we know, ScaleHLS was the first flow that could take a PyTorch model and transform it into synthesizable C

code that can then be translated into RTL code for hardware implementation. HIDA, built on top of the ScaleHLS framework, automates the conversion of algorithmic hardware descriptions into efficient dataflow architectures, also directly generating HLS accelerators from PyTorch models. Looking forward, we discuss future directions, including spatial architecture exploration, runtime reconfiguration for scalability, and heterogeneous computing solutions, to further enhance the efficiency and scalability of hardware accelerators for LLMs. Through these advancements, we aim to address the computational and memory challenges associated with LLM acceleration, ultimately improving the performance and energy efficiency of hardware implementations.

There has recently been a growing interest in leveraging LLMs to enhance Electronic Design Automation (EDA) processes, offering significant potential for improving design productivity, code generation, and verification [53]. Existing research in this domain encompasses various applications, including assistant chatbots for design workflow enhancement, Verilog and script generation, and Hardware Description Language (HDL) verification and analysis. Despite these advancements, several challenges persist, notably the scarcity of high-quality, domain-specific datasets and the need for more specialized LLMs tailored to grasp the intricacies of electronic design languages and processes. As a case study to leverage LLMs for assisting circuit design, we focus on an important task: High-Level Synthesis (HLS) functional verification. We pursue this task through the construction of the Chrysalis dataset, an extensive collection of HLS designs embedded with diverse sets of realistic bugs, and the development of an HLS-specific debugging assistant. A debugging assistant can be built by training an LLM fine-tuned on the Chrysalis dataset which aims to significantly expedite the verification process, enhance productivity, and reduce time-to-market for hardware designs. Additionally, we outline future research directions, including LLM-aided formal verification and the integration of LLMs into multi-agent systems for hardware/software design automation, offering a transformative approach to streamlining the design, verification, and debugging processes in EDA.

In the rest of the paper, Section 2 delves into algorithm-level acceleration for LLMs, while Section 3 provides an overview of hardware co-design tailored for LLMs. Section 4 focuses on the compiler for mapping LLMs to accelerators, and Section 5 explores LLM aided designs. Finally, Section 6 presents the conclusion of the study.

## 2 LLM Algorithm-Level Acceleration

### 2.1 Background Study

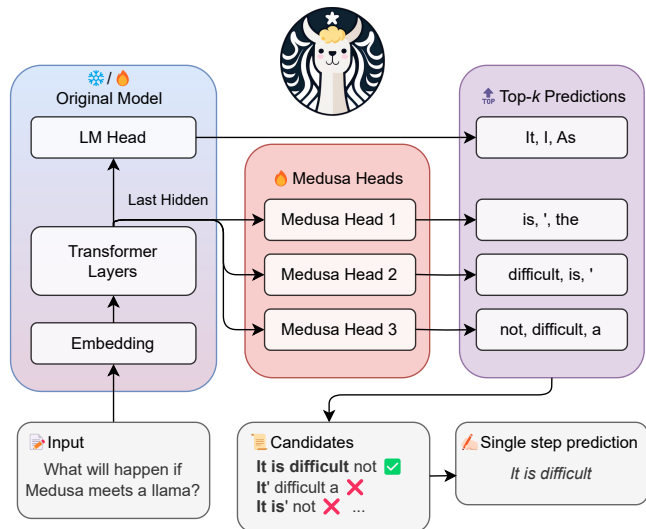
LLMs excel in tasks such as coding copilot, question answering, and summarization. However, their autoregressive nature, where each token depends on all previous ones, makes

decoding memory-intensive and hard to parallelize. This results in significant latency due to the massive size of LLM parameters, impacting applications requiring fast responses like chatbots. Addressing the challenge of reducing inference latency in LLMs is becoming increasingly critical. This section primarily explores previous methods aimed at decreasing the inference latency of LLMs from an algorithmic standpoint, which could facilitate straightforward implementation and integration across various platforms.

**2.1.1 Efficient Parameter Utilization.** The early study [45] shows only a necessary subset of parameters is used per input token to reduce language model inference latency while preserving accuracy. The concepts of early exiting and layer skipping [13, 44] in decoder architectures, allow for an efficient generation of sequences by potentially exiting the decoding process early based on certain criteria, thus saving on computational resources while maintaining output quality. On another perspective, contextual sparsity, as investigated by Liu et al. [33], leverages the insight that a significant portion of model weights can be dynamically omitted without affecting performance, capitalizing on the variability of importance across different weights for different inputs. Lastly, the Mixture of Experts (MoE) [12, 19, 68] approach decouples model size from computational demands, enabling significant scaling of model capacity with minimal impact on inference efficiency, offering a pathway to enhancing model performance without proportional increases in computational burden.

**2.1.2 KV Cache Optimization.** KV Cache in LLMs stores previously computed attention values and keys. This caching proves particularly effective for complex tasks like chain-of-thought [52] reasoning or information retrieval [30]. However, it introduces overheads like setup time, extra memory for cache storage, and the complexity of managing cache validity when the sequence length or batch size increases. Several strategies have been developed to enhance KV Cache efficiency. One key approach is through advanced KV Cache management techniques. For instance, the VLLM [25] introduces PagedAttention which stores keys and values in segmented memory blocks, allowing for more efficient retrieval during attention calculations. Additionally, solutions like Hydragen [22] employ a Shared-prefix KV Cache strategy, greatly improving cache reuse rates by leveraging common sequences. Another significant advancement is the use of KV Cache compression [65], which implements eviction policies to selectively retain tokens in the cache, guided by a scoring function based on cumulative attention.

**2.1.3 Parallel Decoding.** Parallel decoding presents a unique approach by executing multiple decoding steps simultaneously to reduce the overall steps needed, diverging from traditional methods. It typically involves a smaller "draft" model predicting several upcoming words, which are then collectively assessed by the main LLM. This technique



**Figure 1.** The proposed parallel decoding framework Medusa. During inference, each head generates multiple top predictions for its designated position. These predictions are assembled into candidates processed in parallel using a *tree-based attention* mechanism. Then the framework verifies the candidates and accepts a continuation [4].

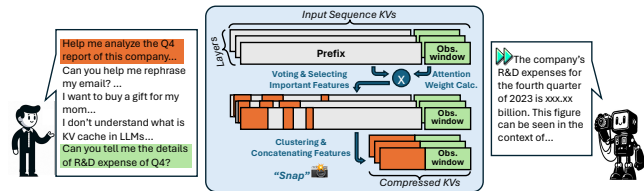
is specifically adapted for LLM efficiency. Recent advancements in parallel decoding for LLMs include techniques by Leviathan et al. [29] and Chen et al. [6], which introduce a re-sampling strategy to align small and large model predictions with the LLM’s distribution, ensuring output consistency. Stern et al. [46] explore multi-token predictions from a single forward pass using a linear projection layer and a tree-based decoding structure to improve decoded sequence acceptance. Additionally, Santilli et al. [43] and Fu et al. [14] adapt Jacobi and Gaussian-Seidel algorithms [40] for parallelizing decoding, incorporating n-gram reuse and attention masks to enhance LLM efficiency.

## 2.2 Proposed Works

LLM inference is predominantly memory-bandwidth-bound with the main latency bottleneck stemming from accelerators’ memory bandwidth rather than arithmetic computations. This bottleneck is inherent to the sequential nature of auto-regressive decoding, where each forward pass requires transferring the complete model parameters from High-Bandwidth Memory (HBM) to the accelerator’s cache. This process, which generates only a single token, underutilizes the arithmetic computation potential of modern accelerators, leading to inefficiency. In our proposed work [4], named as Medusa, we revisit the concept of parallel decoding with a new perspective, noting that current research primarily aims to boost generation speed through draft models. Yet, obtaining an appropriate draft model either from scratch or

from distillation is non-trivial. Also, hosting dual-sized models on a server presents challenges, and it’s even harder to integrate the draft model into a distributed system. To tackle this, we present a novel approach (shown in Fig. 1) using multiple decoding heads as the adapter for prediction, coupled with an optimized tree-based decoding strategy, enhancing the efficiency of the method. Our proposed technique does not need a separate draft model and allows for seamless integration into existing LLM systems. Our experiments demonstrate that limited-resource fine-tuning can achieve over 2.2× speedup without compromising generation quality, while full fine-tuning further improves the speedup to 2.3-2.8× on models with various sizes. Furthermore, parallel decoding improves resource utilization due to increased matrix operations for multi-token validation per step. By employing an optimized tree-based attention mechanism, we strive to minimize the overhead introduced by parallel decoding. Our focus on optimization of both fine-tuning and inference with the decoding adapter in the context of speculative decoding presents a novel direction for enhancing LLM performance.

Furthermore, our method, SnapKV [31], effectively reduces KV cache size, addressing the computational and memory bottlenecks in scenarios involving long sequence inputs. Our findings demonstrate the consistent attention allocation patterns of important features in prompts used throughout the generation process, independent of prompt formats. This observation highlights the potential on KV cache compression for long sequence input, which could reduce the computational and memory overhead in attention calculation during generation steps. Leveraging this insight, our innovative approach intelligently identifies these attention allocation patterns by using the window of features at the end of long sequence input, as shown in Fig. 2, and compresses the KV cache accordingly. This proposed work achieves consistent decoding speeds, significantly enhancing generation speed by 3.6x and improving memory efficiency by 8.2x compared to the baseline, when processing inputs of 16k tokens.



**Figure 2.** The graph shows the simplified workflow of SnapKV, where the orange area represents the group of positions per head clustered and selected by SnapKV.

## 2.3 Future Directions

**2.3.1 Enhanced Versatility in Parallel Decoding.** With the growth in the size of LLMs and their deployment across

both cloud and edge devices, accurately predicting the performance of parallel decoding models has become increasingly complex. A universal solution for LLM performance optimization remains elusive. Without accurate prediction before training and deployment, there is a risk of wasting computational resources and failing to meet performance targets. We focus on modeling and predicting various scenarios, utilizing an analytic model to achieve precise performance estimations. Our objective is to create predictive frameworks that aid in selecting optimal parallel decoding algorithms and their hyperparameters tailored to model size, task complexity, and performance goals. This aims to enhance efficiency, adaptability, and overall model performance.

**2.3.2 Combining KV Compression and Parallel Decoding.** Leveraging KV compression, we see opportunities for notable improvements in tasks with large input prompts, like summarization and multi-round chats, where precise prompt compression will be crucial for maintaining retrieval accuracy and understanding. In long-context scenarios, directly processing the entire prompt and performing inference with parallel decoding introduces significant inference overhead due to the increased computational complexity and memory requirements. To address this, we explore effective attention mechanisms such as Group Query Attention [2] and techniques like quantizing the KV cache to reduce computational load. These refinements are intended to boost LLM efficiency and effectiveness in practical uses while maintaining the generation quality.

## 3 LLM-Hardware co-design

### 3.1 Background Study

The exceptional capabilities of LLMs are countered by their significant memory and computational overheads. Addressing these, LLM-hardware co-design, inspired by the DNN-accelerator co-design methodology [15] [16], customizes hardware to meet these demands and provides insights to optimize LLM architectures. Specialized accelerators, like GPUs, TPUs, and FPGAs, enhance parallel processing and memory capacity and provide efficient LLM execution. At the same time, software strategies, such as model distillation, pruning, and quantization, can effectively reduce LLM size and complexity, making them adaptable to hardware constraints.

We have seen customized hardware accelerators are built for LLM workloads. For example, Tensor Processing Units (TPUs) [20] are designed to efficiently handle matrix operations which are fundamental for LLM attention and linear layers. These accelerator designs enhance LLM support by integrating High Bandwidth Memory (HBM), reconfigurable high-speed interconnects, and multi-type parallel computation support, offering cost-effective LLM training and serving solutions. Beyond ASICs, FPGA-based accelerators are being actively investigated for their potential to provide more

flexible and faster turnaround solutions. For example, DFX [17] utilizes model parallelism and enables rapid concurrent execution of transformer-based workloads, while FlightLLM [62] introduces a configurable compute unit and a LLM mapping flow to support LLM inference.

For LLM designs, researchers have investigated hardware-aware model compression technologies to optimize LLM architecture. FlashAttention [10] reduces the number of High Bandwidth Memory (HBM) accesses by using tiling techniques in attention computations and extends it to block-sparse attention. PagedAttention [25] divides the KV cache into blocks and manages blocks as pages in OS’s virtual memory, reducing the internal and external fragmentation and thus increasing the efficiency within a single request. In addition, model distillation, pruning, and quantization have proven to improve hardware efficiency for LLM deployment. MLFS [24] freezes a base model and stores many small low-rank adapter matrices, which maintains high quality encoder models on edge applications and reduces training time. LLM.int8 [11] develops a two-part vector-wise quantization procedure and a new mixed-precision decomposition scheme, enabling models like OPT-175B on a single server with consumer GPUs. SmoothQuant [54] uses a per-channel smoothing factor to handle outliers in activations and achieves up to 1.56x speedup and 2x memory reduction. ViTCoD [61] prunes attention maps to either dense or sparse patterns and designs an accelerator that coordinates between these two workloads to boost hardware utilization while integrating on-chip encoder and decoder engines.

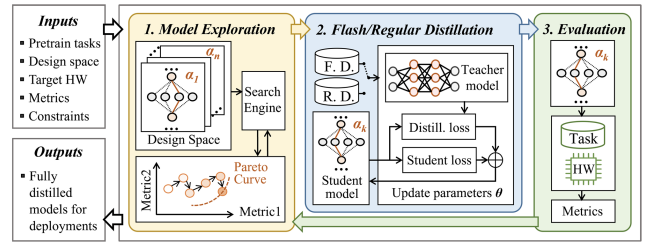
### 3.2 Proposed Works

Following the LLM-hardware co-design method, we propose AutoDistill [64], an end-to-end model distillation framework to deliver hardware-efficient models. As shown in Fig. 3, AutoDistill introduces a three-stage solution, which includes model architecture exploration, model compression, and model evaluation to deliver efficient models given the target hardware and hardware-software tradeoff requirements. To facilitate the hardware/software co-design process, these stages are tightly connected and continuously iterated in a quality-performance space. During the evaluation stage, model quality and its hardware performance results are passed back to the model exploration to guide the search engine for finding a better model architecture that could fulfill both software and hardware requirements. Results show that AutoDistill can efficiently produce student models with lower inference latency and smaller sizes, while maintaining high accuracy on multiple downstream tasks, such as SQuAD for question answering and reading comprehension in Table 1.

We also propose a pruning-aware quantization strategy by combining two of the most effective LLM compression

**Table 1.** The results on SQuADv1.1. Ours-1 and Ours-2 denote two models designed with Autodistill [64].

	# Param	Latency	F1	EM
BERT <sub>BASE</sub>	109 M	-	88.5	80.8
DistilBERT	67 M	-	85.8	77.1
DistilBERT*	67 M	-	86.9	79.1
TinyBERT <sub>6</sub> *	67 M	-	87.5	79.7
NAS-BERT*	60 M	-	88.0	80.5
NAS-BERT*†	60 M	-	88.4	81.2
MobileBERT	25.3 M	0.65 ms	<b>90.0</b>	<b>82.9</b>
MobileBERT‡	25.3 M	0.65 ms	87.7	80.0
Ours-1	22.8 M	0.59 ms	88.4	80.8
Ours-2	<b>20.6 M</b>	<b>0.49 ms</b>	88.1	80.5



**Figure 3.** The proposed AutoDistill framework [64].

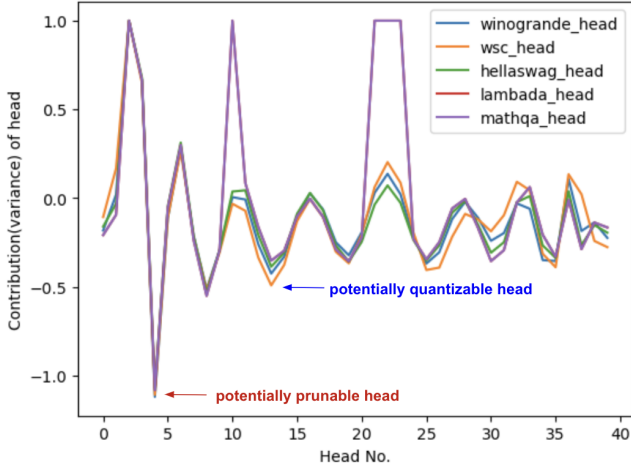
methods, pruning and quantization, for LLM-hardware co-design [51]. We observe a similar sparsity distribution pattern of attention heads across various datasets as shown in Fig. 4, which could potentially be used to smartly choose either completely pruning (0 bit) or different quantization precision (4, 8, 16 bits) for attention heads on individual layers without additional overhead, based on the hardware level objective. Moreover, pruning-aware quantization method could also be combined with other state-of-the-art hardware-aware LLM acceleration frameworks, such as flash attention in Fig. 5, and give higher throughput.

### 3.3 Future Directions

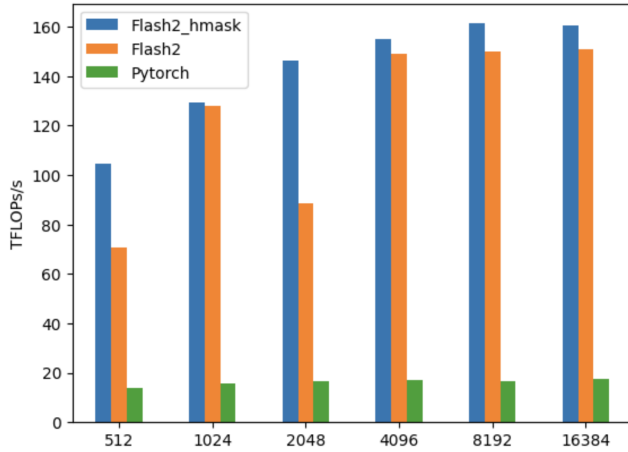
**3.3.1 System-aware algorithmic optimization.** As LLMs continue to grow in size and complexity, it is becoming increasingly critical to design them with hardware system in mind. It means the model developers should consider the hardware system configurations, including accelerator architecture, compute power, memory capacity, system topology, network bandwidth, model parallelism strategies, in addition to the LLM design parameters.

This future direction will create multiple optimization opportunities to explore the combined design space consisting of hardware and software configurations. Model quality will not be the only optimization objective while hardware performance and efficiency metrics, such as Queries per second (QPS) and latency, will also be considered and explored as parts of the LLM designs. With such an enhanced design





**Figure 4.** The profiling results on the activity of heads across different datasets by measuring each head’s contribution based on its variance over the input sequence. Heads that show low variance are considered inactive, leading to contextual sparsity.



**Figure 5.** The preliminary result from forward throughput improvement. Flash2\_hmask is the result from the combination of FlashAttention2 [10] and our pruning-aware quantization approach [51].

space, LLMs can be specifically tailored to the underlying system and hardware architecture and we anticipate further innovations in model size reduction, efficient sharding strategies, optimized data layouts, and other techniques to fully utilize the full potential of target systems.

**3.3.2 Reconfigurable and heterogeneous hardware for LLMs.** Reconfigurable hardware, such as FPGAs, is a promising solution to address the continuously evolving LLM designs. It offers the ability to adapt to the specific computational patterns of different LLM workloads, which allows a

fast development of hardware accelerators to efficiently handle key LLM operations, including matrix multiplication and attention mechanisms. Additionally, it can be combined with heterogeneous hardware to explore new compute paradigms, such as adopting in-memory computing, to address memory-bound operations. This direction enables trade-offs between model quality and hardware efficiency.

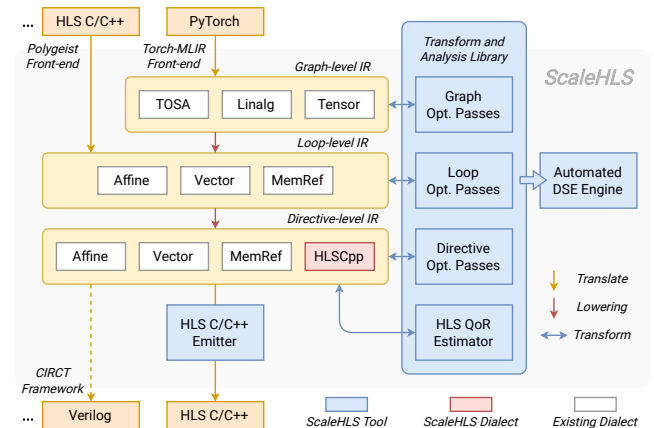
**3.3.3 Co-design for edge LLM applications.** The co-design for edge LLM applications is crucial, given the intricate challenges posed by edge computing’s energy and resource limitations. LLM-hardware co-design emerges as a promising solution to these challenges, aiming to harmonize software and hardware to optimize LLM performance on edge devices. Future research will focus on creating tailored architectures and algorithms that efficiently manage computational resources, ensuring that the quality of LLM services remains high. This could involve exploring adaptive power management techniques, optimizing memory usage, and enhancing processing speeds without sacrificing the accuracy or responsiveness of LLM applications.

## 4 LLM-to-Accelerator Compiler

### 4.1 Background Study

High-Level Synthesis (HLS) [7][9] is vital for rapidly developing efficient, high-density hardware accelerators, enabling quick evaluation of different algorithmic choices. The challenge of enabling scalable compiler from LLM models to HLS-based accelerators lies in effectively exploring the vast design space, which can lead to sub-optimal solutions if not done well, undermining the productivity benefits of HLS. To tackle this challenge, in this section, we will introduce two compilation frameworks, ScaleHLS and HIDA, which can generate HLS accelerators directly from PyTorch models.

### 4.2 Proposed Works



**Figure 6.** ScaleHLS framework architecture [57].

**4.2.1 ScaleHLS.** ScaleHLS [21, 57–59] is a scalable HLS framework based on Multi-Level Intermediate Representation (MLIR) [26]. Fig. 6 shows the overall architecture. ScaleHLS supports C/C++ and PyTorch as design entries. Once the inputs are parsed, ScaleHLS supports three levels of IR, *graph*, *loop*, and *directive*, to apply the HLS-oriented optimizations progressively. At the graph and loop level, graph optimizations (e.g., node fusion and coarse-grained pipelining) and loop optimizations can be performed efficiently. At the lowest directive level, HLS-specific optimizations are applied to fine-tune the hardware micro-architecture.

On top of each level of IR, ScaleHLS provides a set of transform passes to optimize HLS designs. By performing each transform pass at the "correct" level of abstraction, ScaleHLS is able to leverage the intrinsic hierarchy of HLS designs and reduce the algorithmic complexity of transforms. Meanwhile, we propose a Design Space Exploration (DSE) engine to automatically optimize the configurable design parameters and search for the Pareto-dominating design points in the latency-resource utilization space. Finally, the optimized IR is emitted as synthesizable HLS C/C++ code. Experimental results show that, comparing to the baseline designs without manual directives insertion and code-rewriting, that are only synthesized by Vitis HLS, ScaleHLS improves the performances with up to 3825.0× on representative neural network models.

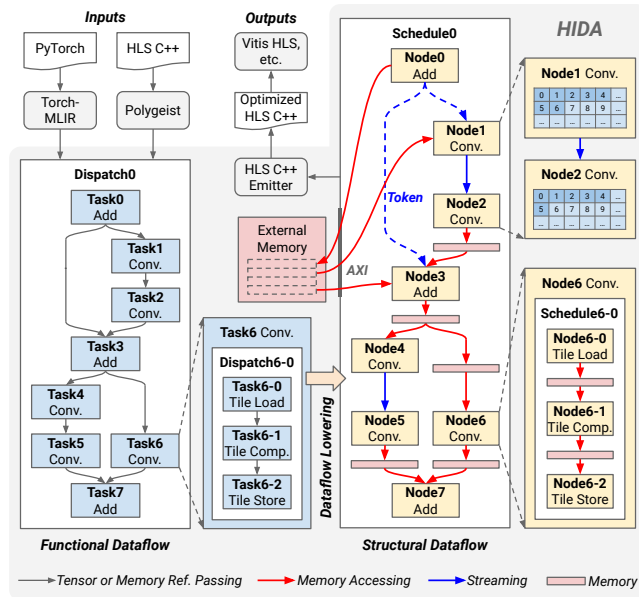


Figure 7. HIDA framework architecture [60].

**4.2.2 HIDA.** HIDA [60] is an HLS framework built upon ScaleHLS with hierarchical dataflow intermediate representations (IR) and optimizations, enabling the automated transformation of algorithmic hardware descriptions to efficient dataflow architectures. Fig. 7 shows the HIDA’s overall architecture. The core of HIDA is its novel dataflow IR, named

HIDA-IR, which is designed for modeling dataflow at two distinct abstraction levels: *Functional* and *Structural*. This dual-level approach is critical for capturing the dataflow’s characteristics and its multi-level hierarchy, thereby facilitating effective optimizations.

Another important aspect of HIDA is the introduction of HIDA-OPT, a new dataflow optimizer. This optimizer utilizes a pattern-driven task fusion algorithm coupled with an intensity- and connection-aware dataflow parallelization algorithm, which can capture the computation complexity and interconnection topology of the dataflow nodes during the parallelization process. Furthermore, HIDA is designed to be end-to-end and extensible, supporting both PyTorch and C++ inputs. This flexibility empowers users to rapidly experiment with various design parameters and prototype new dataflow architectures, broadening the framework’s applicability and ease of use. Despite being fully automated and able to handle various applications, HIDA achieved throughputs that were 8.54× and 1.29× higher than those of ScaleHLS and RTL-based neural network accelerators [63], respectively.

### 4.3 Future Directions

**4.3.1 Spatial Architecture.** Due to the substantial volume of parameters and intermediate computations involved in LLMs, the bottleneck in hardware acceleration frequently resides in the external memory bandwidth. Contrary to the Von Neumann architecture, which consistently battles the "memory wall," spatial architectures can leverage on-chip communication among tasks to minimize frequent external memory accesses. By overlapping the execution of distinct LLM layers and buffering only a subset of intermediate results on-chip, spatial architectures can markedly decrease on-chip memory requirements and overall latency. This approach presents a compelling solution for LLM inference. Nonetheless, the automatic generation of spatial architectures remains challenging, opening vast avenues for innovation in compilation and architecture design.

**4.3.2 Runtime Reconfiguration.** To achieve spatial parallelization, tasks must be instantiated simultaneously on-chip. However, due to constrained computational and on-chip memory resources, it is infeasible to simultaneously map all layers of emerging LLMs on-chip, which significantly limits the scalability of spatial architectures. Consequently, runtime reconfiguration emerges as a crucial strategy for enabling scalable spatial solutions. The main challenge lies in automating the balance between spatial and sequential execution – that is, addressing the scheduling problem – to optimize the performance-energy trade-off.

**4.3.3 Heterogeneous Computation.** Accelerating LLMs presents a unique challenge due to their dual nature, being both computation-bound and memory-bound. The prefill phase of LLMs is dominated by General Matrix Multiply

(GEMM) operators, making it computation-intensive. In contrast, the generation phase is dominated by General Matrix-Vector (GEMV) operations, demanding substantial memory bandwidth to keep the computation units engaged (refer to Section 2.2). This dual nature of LLMs unveils significant opportunities for heterogeneous computing solutions, where compilers assume an important role in the code generation for heterogeneous platforms and facilitating efficient communication between them.

**4.3.4 Advanced HIDA.** Although the HIDA framework can conduct effective dataflow-aware design space exploration, optimizing streaming buffers in LLM accelerators remains a formidable challenge due to the self-attention mechanism and complex inter-layer connections in LLMs. Enhancements in the HIDA framework could address more complicated stream optimizations to reduce on-chip memory consumption. Additionally, recent works [8, 67] have demonstrated the ability to generate efficient kernel designs through customized scheduling. We propose to integrate these highly-optimized kernels into the HIDA explorer to further improve the efficiency of LLM accelerators. We also propose to enhance the code generation of HIDA to support more hardware platforms with dataflow architecture, such as AMD Versal ACAP [3].

## 5 LLM-aided design

### 5.1 Background Study

The existing related work regarding leveraging LLMs in the field of EDA could be divided into three categories [66]: (1) Assistant Chatbot for Enhanced Design Workflow: Chip-Nemo [32] leverages domain adaptation techniques such as custom tokenizers, domain-adaptive continued pretraining, and Supervised Fine-Tuning (SFT) atop the foundation model LLaMA2 [37]. This integration facilitates instant access to knowledge, streamlines the query-response process, and diminishes reliance on conventional search methods and associated delays. (2) HDL and Script Generation: LLMs, such as those in AutoChip [48] and VeriGen [49], have shown their effectiveness in generating Verilog codes and EDA tool scripts from natural language instructions; (3) HDL Verification and Analysis: RTLFixer [50] exemplifies this by introducing a framework aimed at correcting Verilog code errors utilizing tools like OpenAI GPT-3.5/GPT-4, supplemented by Retrieval Augmented Generation (RAG) and ReAct prompting techniques. Additional efforts in this area focus on generating SystemVerilog Assertions (SVA) for security purposes [23][35][41], illustrating the wide-ranging potential of LLMs in bolstering HDL verification and analysis processes. CHI-RAAG [34] is proposed to generate SVA assertions from natural language specification based on GPT-4. For those assertions with syntax error or simulation error, LLMs could receive the automatic feedback of log file and then regenerate the SVA for retest. Orenes-Vera [39] proposed an iterative

methodology where LLMs, particularly GPT-4, are prompted with refined rules and RTL code to generate SVA, which is then evaluated for correctness and completeness through a series of testbench simulations and revisions.

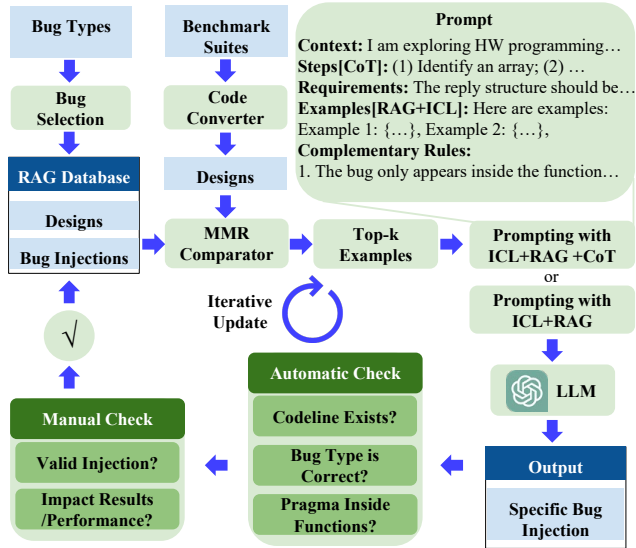
Despite the promising developments in LLM-aided design within EDA, several challenges remain: (1) Data Quality and Availability: The efficacy of LLMs in EDA critically hinges on the availability of high-quality, domain-specific datasets for their training and refinement. Unfortunately, the proprietary nature of many electronic designs and the tools used for EDA significantly limit access to such datasets. The bulk of detailed hardware design codes, primarily developed within corporate settings, are not made public. This restriction leads to a scarcity of accessible, high-grade datasets, thus hindering the development and optimization of LLMs specifically engineered for EDA applications; (2) Development of Specialized LLMs: There is a critical need for the development of LLMs that are specifically tailored to grasp the complexities of electronic design languages and processes. The generic models, while useful, often lack the nuanced understanding required to effectively generate, verify, and analyze hardware code and to interact with EDA tools at a level that matches human experts. This necessitates a concerted effort to create more specialized models that can comprehend and manipulate the intricate details of electronic designs with a high degree of accuracy and efficiency.

### 5.2 Proposed Works

One use case of LLM-aided design is to harness LLMs for enhancing the verification and debugging processes for HLS code development. HLS, with its higher level of abstraction, can significantly improve design productivity as explained in Section 4.1.

**5.2.1 Chrysalis Dataset.** The cornerstone of our proposed work is the Chrysalis dataset, an extensive collection of HLS designs embedded with a diverse set of realistic bugs [51]. This dataset is meticulously curated from a wide range of open-source HLS benchmark suites, featuring over 1,500 designs, with both the version embedded with bugs and the corresponding bug-free version. Fig. 8 outlines our methodology for constructing the Chrysalis dataset. We begin by gathering open-source benchmark suites and difficult bug types (bugs also include non-ideal HLS Pragma insertions), which compilers often struggle to identify. These suites are then converted to the function-level designs. Using the Maximal Marginal Relevance (MMR) algorithm, we select the top-k similar designs from the RAG database for bug injection prompts. The prompt generation chooses one strategy based on bug type statistics: one combining In-Context Learning (ICL), Retrieval-augmented Generation (RAG), and Chain-of-Thoughts (CoT); the other using just RAG and ICL. After integration, the prompts are processed by an LLM (GPT-4 in our case) to generate bug (or non-ideal Pragma) injection solutions, which are then validated through both automatic



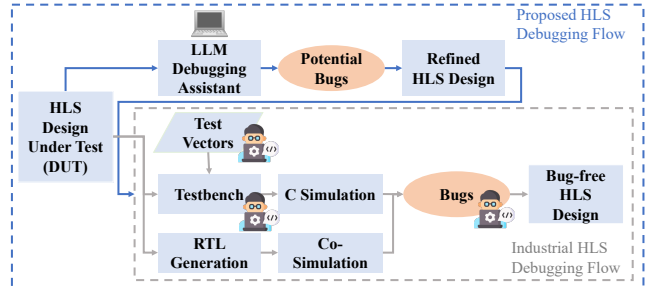


**Figure 8.** Overview of the *Chrysalis* Dataset Construction and Iterative Upgrade Process: For each check iteration, it involves evaluating the dataset’s validity and expanding the RAG dataset accordingly. Through these iterations, the quality of the dataset progressively improves.

procedures as well as manual checks by hardware engineers. Successful solutions are added to the RAG database to enhance its diversity and volume, improving future solutions. Our evaluations show 88% validity for all the bugs. This dataset serves not only as a tool for training our domain-specific LLM, but also as a benchmark for evaluating the model’s proficiency in identifying and suggesting fixes for common and complex HLS bugs.

**5.2.2 HLS-specific Debugging Assistant.** Building upon the *Chrysalis* dataset, our next step involves the creation of an HLS-specific debugging assistant, as Fig. 9 shows. Engineers typically design test vectors and create test benches manually, then perform C simulations and co-simulations to analyze and identify potential bugs, which is time-consuming. To improve the efficiency of the debugging process, we proposed a novel flow leveraging the capabilities of LLMs on top of the traditional HLS debugging flow. This LLM will be fine-tuned to understand the intricacies of HLS code, enabling it to identify errors, suggest corrections, and provide guidance directly within the developers’ workflow. The assistant aims to integrate seamlessly with popular development environments, offering real-time support to engineers as they navigate the complexities of HLS designs. By providing context-aware suggestions and corrections, the debugging assistant will significantly expedite the verification process, enhancing productivity and reducing the time-to-market for hardware designs.

The entire methodology could be adapted to RTL debugging as well, starting from the bug injection stage, using



**Figure 9.** LLM-based HLS Debugging Flows Working Together with Traditional Flows: By incorporating our LLM debugging assistant, the number of bugs requiring verification by test cases can be significantly reduced.

open-source LLMs and developing a domain-specific RTL debugger through fine-tuning. To effectively transition to this new application, we must tackle diverse bug types specific to RTL, such as those related to timing constraints, race conditions, and synthesis-simulation mismatches. Particularly, the inherent timing characteristics of RTL designs can lead to more complex bugs, often manifesting as issues in timing analysis that are not present in higher-level abstractions. Given the complexity of RTL code, one ambitious goal is to reduce manual debugging and verification effort by building an advanced and automated RTL verification and debugging flow. This would involve enriching our dataset to include RTL designs (can be generated through HLS working with our *Chrysalis* dataset) together with test benches and test vectors. A flow similar to Fig. 8 can be developed to assess and improve the validity of such bug injections. Furthermore, seamless integration with EDA tools is crucial to enable real-time analysis and correction within the existing design frameworks.

### 5.3 Future Directions

**5.3.1 LLM-Aided Debugging.** Our research highlights challenges in using LLMs to inject certain HLS bug types, such as operator precedence errors and incorrect data access patterns for interface pragmas. These difficulties stem from sparse code patterns and the complexities of the existing codebase, necessitating further investigation and refined methodologies for effective bug injection. Additionally, as manual verification of bug injections remains necessary in our current flow, creating an automated flow to estimate performance could speed up the identification and resolution of non-ideal Pragma issues, thus enhancing the quality and quantity of the dataset. Furthermore, for the HLS-specific debugging assistant, we will employ Low-Rank Adaptation (LoRA) [18] for supervised fine-tuning on state-of-the-art open-source LLMs such as LLaMA3 [36], utilizing commercial HLS documentation for design guidelines and rules together with our *Chrysalis* dataset.

**5.3.2 LLM-Aided Formal Verification.** LLMs can enhance the formal verification process in hardware design by generating precise assertions for the proof of correctness. By integrating these assertions into the formal verification workflow, LLMs can substantially increase hardware design productivity. One promising direction is to explore an iterative process: after the initial proof attempts, the theorem prover’s feedback is utilized to refine the LLM’s output. This feedback loop enables the LLM to adjust its generated proofs iteratively until the assertions are fully verifiable. Through this dynamic interaction between LLMs and theorem provers, the generation of program proofs becomes both faster and more achievable. This methodology not only speeds up the verification process but also ensures a higher degree of reliability in hardware design verification.

**5.3.3 LLM-Aided Hardware/Software Design Automation.** In the realm of EDA, employing LLM multi-agent systems promises a transformative approach to streamlining the design, verification, and debugging processes. These sophisticated systems autonomously manage various phases of the workflow, seamlessly transitioning from design to verification and debugging. By deploying multiple specialized LLM agents – each adept in distinct facets of the design process such as code generation, verification, error detection, and performance optimization – a highly efficient pipeline is crafted. This orchestrated integration allows the agents to collaboratively refine and optimize the design iteratively, leveraging real-time feedback and comprehensive verification results. Throughout the process, hardware engineers are only tasked with overseeing the initial specification and periodically reviewing the outputs from the LLMs to ensure that they align with the design intentions and confirm the reliability of the LLMs’ outputs.

## 6 Conclusion

In our study, we focused on optimizing LLMs to reduce inference latency and improve efficiency across various applications. We presented a new method, Medusa, to use multiple decoding heads for prediction, coupled with an optimized tree-based decoding strategy for parallel token processing to speed up the execution of LLMs. We also proposed a novel method, SnapKV, that effectively reduced KV cache size, addressing the computational and memory bottlenecks in scenarios involving long sequence inputs.

We discussed LLM/hardware co-design to integrate both hardware optimization for efficient execution and model architecture exploration for improved system efficiency while maintaining LLM accuracy. HLS frameworks like ScaleHLS and HIDA were explored for accelerating LLMs directly from PyTorch models, envisioning automated generation of spatial architectures and heterogeneous computing solutions.

We also explored the advancements in LLM-aided design for EDA and discussed a novel flow to create the Chrysalis

dataset that can be used to train an LLM-based HLS-specific debugging assistant. A similar strategy can be adopted for building an RTL-specific debugging assistant as well. These methods are promising for streamlining the debugging and verification process of the hardware code development.

For each aspect mentioned above, we also outlined promising future directions for further research and exploration.

## Acknowledgments

This work is supported in part by the IBM-Illinois Discovery Accelerator Institute, AMD Center of Excellence at UIUC, AMD Heterogeneous Adaptive Compute Cluster (HACC) initiative, NSF 2117997 grant through the A3D3 institute, and Semiconductor Research Corporation (SRC) 2023-CT-3175 grant.

## References

- [1] Baleegh Ahmad et al. 2023. Fixing hardware security bugs with large language models. *arXiv:2302.01215*.
- [2] Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. 2023. Gqa: Training generalized multi-query transformer models from multi-head checkpoints. *arXiv preprint arXiv:2305.13245*.
- [3] AMD/Xilinx. [n. d.]. Versal Adaptive Compute Acceleration Platform. <https://www.xilinx.com/products/silicon-devices/acap/versal.html>
- [4] Tianle Cai et al. 2024. Medusa: Simple llm inference acceleration framework with multiple decoding heads. *arXiv:2401.10774*.
- [5] Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, Jason D. Lee, Deming Chen, and Tri Dao. 2024. Medusa: Simple LLM Inference Acceleration Framework with Multiple Decoding Heads. In *ICML*. <https://openreview.net/forum?id=PEpbUobfJv>
- [6] Charlie Chen et al. 2023. Accelerating large language model decoding with speculative sampling. *arXiv:2302.01318*.
- [7] Deming Chen et al. 2005. xPilot: A Platform-Based Behavioral Synthesis System. In *SRC Techcon*.
- [8] Hongzheng Chen et al. 2024. Allo: A Programming Model for Composable Accelerator Design. *arXiv preprint arXiv:2404.04815*.
- [9] Jason Cong, Jason Lau, Gai Liu, Stephen Neuendorffer, Peichen Pan, Kees Vissers, and Zhiru Zhang. 2022. FPGA HLS Today: Successes, Challenges, and Opportunities. *ACM Trans. Reconfigurable Technol. Syst.* 15, 4, Article 51, 42 pages. <https://doi.org/10.1145/3530775>
- [10] Tri Dao et al. 2022. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*.
- [11] Tim Dettmers et al. 2022. Llm.int8(): 8-bit matrix multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339*.
- [12] Nan Du et al. 2022. Glam: Efficient scaling of language models with mixture-of-experts. In *ICML*.
- [13] Maha Elbayad et al. 2019. Depth-adaptive transformer. *arXiv:1910.10073*.
- [14] Yichao Fu et al. 2024. Break the sequential dependency of llm inference using lookahead decoding. *arXiv:2402.02057*.

- [15] Cong Hao et al. 2018. Deep neural network model and FPGA accelerator co-design: Opportunities and challenges. In *IC-SICT*.
- [16] Cong Hao et al. 2019. FPGA/DNN co-design: An efficient design methodology for IoT intelligence on the edge. In *DAC*.
- [17] Seongmin Hong et al. 2022. DFX: A low-latency multi-FPGA appliance for accelerating transformer-based text generation. In *MICRO*.
- [18] Edward J Hu et al. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- [19] Albert Q Jiang et al. 2024. Mixtral of experts. *arXiv:2401.04088*.
- [20] Norm Jouppi et al. 2023. Tpu v4: An optically reconfigurable supercomputer for machine learning with hardware support for embeddings. In *ISCA*.
- [21] Hyegang Jun et al. 2023. AutoScaleDSE: A scalable design space exploration engine for high-level synthesis. In *TRETS*.
- [22] Jordan Juravsky et al. 2024. Hydragen: High-Throughput LLM Inference with Shared Prefixes. *arXiv:2402.05099*.
- [23] Rahul Kande et al. 2023. Llm-assisted generation of hardware assertions. *arXiv:2306.14027*.
- [24] Achintya Kundu et al. 2024. Efficiently Distilling LLMs for Edge Applications. *arXiv preprint arXiv:2404.01353*.
- [25] Woosuk Kwon et al. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In *SOSP*.
- [26] Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, and Oleksandr Zinenko. 2021. MLIR: Scaling compiler infrastructure for domain specific computation. In *2021 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. IEEE, 2–14.
- [27] Zhihong Lei, Ernest Pusateri, Shiyi Han, Leo Liu, Mingbin Xu, Tim Ng, Ruchir Travadi, Youyuan Zhang, Mirko Hannemann, Man-Hung Siu, et al. 2024. Personalization of ctc-based end-to-end speech recognition using pronunciation-driven subword tokenization. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 10096–10100.
- [28] Zhihong Lei, Mingbin Xu, Shiyi Han, Leo Liu, Zhen Huang, Tim Ng, Yuanyuan Zhang, Ernest Pusateri, Mirko Hannemann, Yaqiao Deng, et al. 2023. Acoustic Model Fusion For End-to-End Speech Recognition. In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 1–7.
- [29] Yaniv Leviathan et al. 2023. Fast inference from transformers via speculative decoding. In *ICML*.
- [30] Patrick Lewis et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*.
- [31] Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle Cai, Patrick Lewis, and Deming Chen. 2024. SnapKV: LLM Knows What You are Looking for Before Generation. *arXiv preprint arXiv:2404.14469*.
- [32] Mingjie Liu et al. 2023. Chipnemo: Domain-adapted llms for chip design. *arXiv:2311.00176*.
- [33] Zichang Liu et al. 2023. Deja vu: Contextual sparsity for efficient llms at inference time. In *ICML*.
- [34] Bhabesh Mali et al. 2024. ChIRAAG: ChatGPT Informed Rapid and Automated Assertion Generation. *arXiv preprint arXiv:2402.00093*.
- [35] Xingyu Meng et al. 2023. Unlocking hardware security assurance: The potential of llms. *arXiv:2308.11042*.
- [36] Meta. 2024. Introducing Meta Llama 3: The most capable openly available LLM to date. <https://ai.meta.com/blog/meta-llama-3/>
- [37] Meta. 2024. Llama 2: open source, free for research and commercial use. <https://llama.meta.com/llama2/>
- [38] Md Rakib Hossain Misu et al. 2024. Towards AI-Assisted Synthesis of Verified Dafny Methods. *Proc. ACM Softw. Eng.* 1, FSE. <https://doi.org/10.1145/3643763>
- [39] Marcelo Orenes-Vera et al. 2023. Using llms to facilitate formal verification of rtl. *arXiv e-prints*, arXiv–2309.
- [40] James M Ortega and Werner C Rheinboldt. 2000. Iterative solution of nonlinear equations in several variables. *Classics in Applied Mathematics*.
- [41] Sudipta Paria et al. 2023. Divas: An llm-based end-to-end framework for soc security analysis and policy-based protection. *arXiv:2308.06932*.
- [42] Saurabh Pujar, Luca Buratti, Xiaojie Guo, Nicolas Dupuis, Burn Lewis, Sahil Suneja, Atin Sood, Ganesh Nalawade, Matt Jones, Alessandro Morari, and Ruchir Puri. 2023. Invited: Automated Code generation for Information Technology Tasks in YAML through Large Language Models. In *2023 60th ACM/IEEE Design Automation Conference (DAC)*. 1–4. <https://doi.org/10.1109/DAC56929.2023.10247987>
- [43] Andrea Santilli et al. 2023. Accelerating transformer inference for translation via parallel decoding. *arXiv:2305.10427*.
- [44] Tal Schuster et al. 2022. Confident adaptive language modeling. *Advances in Neural Information Processing Systems*.
- [45] Antoine Simoulin et al. 2021. How many layers and why? An analysis of the model depth in transformers. In *IJCNLP Student Research Workshop*.
- [46] Mitchell Stern et al. 2018. Blockwise parallel decoding for deep autoregressive models. *Advances in Neural Information Processing Systems*.
- [47] Gemini Team et al. 2024. Gemini: A Family of Highly Capable Multimodal Models. *arXiv:2312.11805 [cs.CL]*
- [48] Shailja Thakur et al. 2023. Autochip: Automating hdl generation using llm feedback. *arXiv:2311.04887*.
- [49] Shailja Thakur et al. 2023. Verigen: A large language model for verilog code generation. In *TRETS*.
- [50] YunDa Tsai et al. 2023. Rtlfixer: Automatically fixing rtl syntax errors with large language models. *arXiv:2311.16543*.
- [51] Lily Jiaxin Wan et al. 2024. Software/Hardware Co-design for LLM and Its Application for Design Verification. In *ASP-DAC*.
- [52] Jason Wei et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*.
- [53] Haoyuan Wu et al. 2024. ChatEDA: A Large Language Model Powered Autonomous Agent for EDA. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 1–1. <https://doi.org/10.1109/TCAD.2024.3383347>
- [54] Guangxuan Xiao et al. 2023. Smoothquant: Accurate and efficient post-training quantization for large language models. In *ICML*.
- [55] Mingbin Xu, Alex Jin, Sicheng Wang, Mu Su, Tim Ng, Henry Mason, Shiyi Han, Yaqiao Deng, Zhen Huang, and Mahesh

- Krishnamoorthy. 2023. Conformer-Based Speech Recognition On Extreme Edge-Computing Devices. *arXiv preprint arXiv:2312.10359*.
- [56] Mingbin Xu, Congzheng Song, Ye Tian, Neha Agrawal, Filip Granqvist, Rogier van Dalen, Xiao Zhang, Arturo Argueta, Shiyi Han, Yaqiao Deng, et al. 2023. Training large-vocabulary neural language models by private federated learning for resource-constrained devices. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 1–5.
- [57] Hanchen Ye et al. 2022. ScaleHLS: A new scalable high-level synthesis framework on multi-level intermediate representation. In *HPCA*.
- [58] Hanchen Ye et al. 2022. ScaleHLS: a scalable high-level synthesis framework with multi-level transformations and optimizations. In *DAC*.
- [59] Hanchen Ye et al. 2023. High-level Synthesis for Domain Specific Computing. In *ISPD*.
- [60] Hanchen Ye et al. 2024. HIDA: A Hierarchical Dataflow Compiler for High-Level Synthesis. In *ASPLOS*.
- [61] Haoran You et al. 2023. Vitcod: Vision transformer acceleration via dedicated algorithm and accelerator co-design. In *HPCA*.
- [62] Shulin Zeng et al. 2024. FlightLLM: Efficient Large Language Model Inference with a Complete Mapping Flow on FPGA. *arXiv:2401.03868*.
- [63] Xiaofan Zhang et al. 2018. DNNBuilder: An automated tool for building high-performance DNN hardware accelerators for FPGAs. In *ICCAD*.
- [64] Xiaofan Zhang et al. 2022. AutoDistill: An end-to-end framework to explore and distill hardware-efficient language models. *arXiv:2201.08539*.
- [65] Zhenyu Zhang et al. 2023. H2o: Heavy-hitter oracle for efficient generative inference of large language models. *arXiv:2306.14048*.
- [66] Ruizhe Zhong et al. 2023. LLM4EDA: Emerging Progress in Large Language Models for Electronic Design Automation. *arXiv:2401.12224*.
- [67] Jinming Zhuang, Jason Lau, Hanchen Ye, Zhuoping Yang, Yubo Du, Jack Lo, Kristof Denolf, Stephen Neuendorffer, Alex Jones, Jingtong Hu, et al. 2023. CHARM: Composing Heterogeneous Accelerators for Matrix Multiply on Versal ACAP Architecture. In *Proceedings of the 2023 ACM/SIGDA International Symposium on Field Programmable Gate Arrays*. 153–164.
- [68] Barret Zoph et al. 2022. Designing effective sparse expert models. *arXiv:2202.08906*.