

Fig. 2. Overview of teaching LLM to be aware of performance for GEMM.

design the performance model to cover the modeling from microarchitecture to system level to make LLM fully aware of performance. At the microarchitecture level, we establish the relationship between throughput and design parameters (e.g., bit width, pipeline stages) by analyzing array critical paths and incorporating DC synthesis data. At the system level, inspired by the roofline model [6], we use operation intensity (OI) and throughput as metrics. The relationship between OI and parameters like matrix dimensions, GEMM array size, and pipeline stages is shown in Fig. 2. The LLMPM is integrated into the framework via in-context learning with prompts and an internal code interpreter.

III. EXPERIMENTAL RESULTS

In this section, we demonstrate the effectiveness of our LLM4GV framework for automatically generating the optimal Verilog code for GEMM. We use GPT-4 as the default baseline and evaluate the correctness and performance of the Verilog code generated by the frameworks with Synopsys Design Compiler and AMD (formerly Xilinx) Vivado.

To validate the effectiveness of our multi-agent framework with the proposed template and critical code prompt (TCCP) in code generation correctness for variants of GEMM, we design an ablation and comparative experiment with the following baseline and variants: GPT-4 (baseline), GPT-4 + function description prompt (FP), and GPT-4 + FP + TCCP. We also add the open-source LLMs-based frameworks CodeV-7B [4] and VeriGen-16B [3] as references for comparison. Table. I summarizes the quantitative evaluation of both syntax and functionality correctness of different LLMs and different methods. We show the number of syntax-correct and functionality-correct codes for each of the five codes, respectively.

To validate the effectiveness of the proposed LLMPM, we design an comparative experiment with the following baselines: GPT-4 + performance-aware prompt (PAP), GPT-4 + performance model (PM) in [5], Manual design, GPT-4 + LLMPM (ours). As shown in Fig. 3(a), we utilize the three commonly used neural networks as baselines to evaluate the performance of the GEMM code. To validate that our method can achieve the best performance under different hardware resource constraints and system requirements, we visualize the trade-off between DSPs and latency based on VGG-16 with an input resolution of $224 \times 224 \times 3$, as shown in the Fig. 3(b). The results show that the code generated by our proposed framework demonstrates significant advantages in both correctness and performance.

IV. CONCLUSION

Recent works on LLMs for Verilog code generation show great potential but suffer from a lack of automation, correct-

TABLE I
THE SYNTAX AND FUNCTIONALITY CORRECTNESS VERIFICATION FOR DIFFERENT LLMs AND DIFFERENT METHODS

GEMM module	VeriGen-16B [3]		CodeV-7B [4]		GPT-4 (baseline)		GPT4+FP		GPT-4+FP +TCCP (ours)	
	synt. ^{*1}	func. ^{*2}	synt.	func.	synt.	func.	synt.	func.	synt.	func.
PE Array-1	5	0	5	0	4	0	3	1	5	5
PE Array-2	5	0	5	0	2	0	3	0	5	5
PE Array-3	0	0	4	0	4	0	4	2	5	5
DataReOrganize-1	5	0	2	0	5	0	2	0	5	4
DataReOrganize-2	5	0	5	0	5	0	2	0	4	4
DataReOrganize-3	5	0	5	0	5	2	3	0	4	3
FIFO	0	0	1	0	3	0	1	0	5	5
Requantization	0	0	5	0	3	0	5	1	5	5
success rate	62.5%	0%	80%	0%	77.5%	5%	57.5%	10%	95%	90%

*1: The number of codes successfully compiled by DC in five trials.

*2: The number of codes that passed functional testing in five trials.

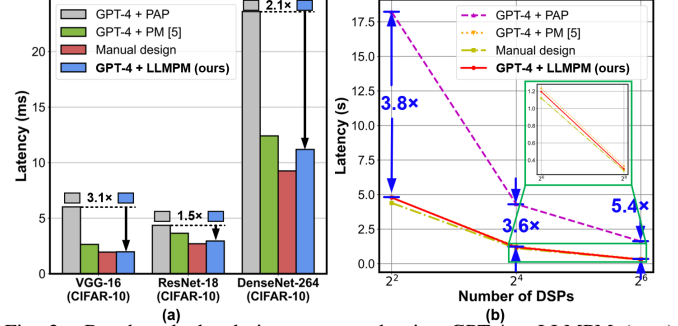


Fig. 3. Benchmark the designs generated using GPT-4 + LLMPM (ours) parameters against those generated using GPT-4 + PAP, GPT-4 + PM [5], and manually design. (a) For different models. (b) Under certain resource constraints.

ness, flexibility, and awareness of performance when supporting GEMM modules. To address these problems, we propose a flexible performance-aware framework called LLM4GV, which utilizes in-context learning based on GPT-4 to generate the optimized Verilog code for the GEMM module with a large design space according to system requirements. In future work, we plan to develop a well-annotated, high-quality Verilog dataset of the GEMM module with this framework and use it to fine-tune an open-source LLM.

V. ACKNOWLEDGEMENT

This work was supported in part by the National Key R&D Program of China under Grant 2022YFB4400600 and the National Natural Science Foundation of China under Grant 62404256.

REFERENCES

- W. Li, A. Hu, N. Xu and G. He, "Quantization and Hardware Architecture Co-Design for Matrix-Vector Multiplications of Large Language Models," in IEEE Transactions on Circuits and Systems I: Regular Papers, vol. 71, no. 6, pp. 2858-2871, June 2024.
- Fu, Yonggan, et al. "Gpt4aigchip: Towards next-generation ai accelerator design automation via large language models." 2023 IEEE/ACM International Conference on Computer Aided Design (ICCAD). IEEE, 2023.
- Thakur, Shailja, et al. "Verigen: A large language model for verilog code generation." ACM Transactions on Design Automation of Electronic Systems 29.3 (2024): 1-31.
- Zhao, Yang, et al. "CodeV: Empowering LLMs for Verilog Generation through Multi-Level Summarization." arXiv preprint arXiv:2407.10424 (2024).
- C. Zhang, et al. "Optimizing FPGA-based accelerator design for deep convolutional neural networks," 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA), Monterey, CA, USA, 2015, pp. 161-170.
- S. Williams, A. Waterman, and D. Patterson. Roofline: An insightful visual performance model for multicore architectures. Commun. ACM, 52(4):65-76, Apr. 2009.