

DSL-MONKEYS: SELF-GENERATED IN-CONTEXT EXAMPLES FOR LOW-RESOURCE GPU DSL KERNELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Domain-specific languages (DSLs) for GPU kernels such as TileLang and ThunderKittens remain challenging for LLMs due to scarce training data and limited documentation. Standard test-time techniques (repeated sampling, iterative refinement, evolutionary search) treat each kernel independently, missing a key insight: successful implementations encode reusable DSL usage patterns directly applicable to new tasks. We introduce DSL-Monkeys, a batched test-time approach that considers an entire collection of kernels jointly rather than individually. DSL-Monkeys maintains a growing archive of successfully-implemented kernels and retrieves them as in-context demonstrations for remaining tasks, resembling a form of test-time reinforcement learning—improving on future tasks by curating a retrieval database rather than updating weights. Starting from just 3 seed examples, DSL-Monkeys achieves 71% and 18% correctness on 200 KernelBench Level 1-2 tasks for TileLang and ThunderKittens respectively, compared to baseline pass@1 rates of < 20% and < 1%. Secondly, with curriculum decomposition, DSL-Monkeys solves complex linear attention variants that approach or match human-expert-written Triton performance. Thirdly, we explore scaling DSL-Monkeys to 10K PyTorch programs from GitHub, producing thousands of DSL implementations that address DSL data scarcity. We show that fine-tuning improves Qwen’s TileLang implementation ability via both distillation and self-generated data, a potential path towards improving LLMs’ code generation ability with emerging DSLs.

1 INTRODUCTION

Many software domains and applications rely on specialized libraries and domain-specific languages (DSLs) that are poorly represented in LLM training corpora (Joel et al., 2025; Khandpur et al., 2024; Liu et al., 2023; Cassano et al., 2022). In these “low-resource” settings, LLMs lack knowledge about target DSLs’ syntax and common programming idioms, leading to poor code generation performance.

DSLs for authoring high-performance GPU kernels for ML such as CuTe (NVIDIA, 2026), TileLang (Wang et al., 2025), ThunderKittens (Spector et al., 2024b), and Triton-TLX (Experimental, 2025) exemplify this challenge. They offer the promise of productivity and performance, but are only sparsely documented (<1 million public tokens) and introduce syntax and programming models that are rare in training corpora. For these reasons, although there is significant interest in harnessing the benefits of these frameworks in production settings (FlashAttention-4 and DeepSeek’s Sparse Attention (Dao-AILab, 2024a; DeepSeek-AI, 2025) are successful recent examples), LLMs struggle to automatically generate correct kernels using these frameworks. For pass@1, frontier models correctly implement less than < 20% of KernelBench’s Level 1-2 (Ouyang et al., 2025a) kernels, which test basic kernel generation and optimization from a PyTorch reference, when targeting TileLang, and less than 1% for ThunderKittens.

Since scarce DSL training data limits improvement via fine-tuning directly, scaling test-time compute is the natural approach given the presence of verifier. For instance, standard approaches include repeated sampling Brown et al. (2024), iterative refinement based on feedback from compiler or profiler feedback, or evolutionary search (Sharma, 2025a), which boosts correct implementation rate on KernelBench’s Level 1-2 to 60% (TileLang) and 17% (ThunderKittens). We argue that these prior test-time approaches are limited because they force the LLM to focus on implementing a *single kernel* at a time. However, deep learning kernels share related operators and implementation patterns:

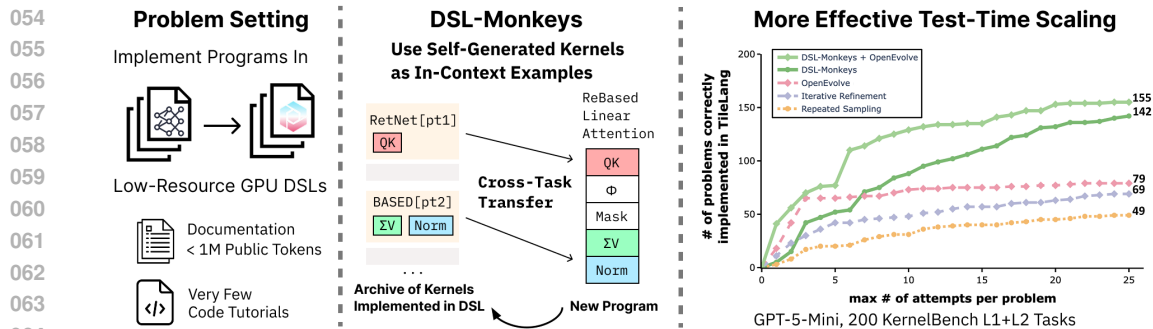


Figure 1: DSL-Monkeys implements kernels in low-resource DSLs via cross-task transfer.

tiling strategies, memory layouts, and reductions. Therefore, when an LLM successfully implements a kernel in the target DSL, that success is likely to be a helpful example for future implementations. Given this observation, we ask: *when implementing kernels in a new DSL, can successful solutions from related kernels be leveraged at test-time to solve the remaining ones?*

Rather than developing advanced techniques for implementing a single kernel (and repeating this procedure individually for all kernels in a collection), our solution, **DSL-Monkeys**, adopts a simple approach that considers the entire collection as a whole: it repeatedly iterates over the collection, attempting to reimplement a batch of kernels in the target DSL. Successfully implemented kernels are logged in an archive. In future iterations, when attempting previously failed kernels, the system retrieves solutions for similar kernels from the archive and provides them as in-context examples. This formulation resembles label propagation: given unlabeled data (kernels to port) and initial labels (a small seed set from DSL documentation), the system propagates knowledge to similar kernels via in-context demonstrations, using a verifier (compilation + tests) for validation. We call this phenomenon *cross-task transfer*: using solved examples to enable solutions for new, related tasks.

Overall, we demonstrate that our test-time compute approach is a simple method that effectively implements ML kernels in low-resource DSLs. Specifically, we show:

1. **Self-Generated In-Context Examples Help Write Kernels in New DSLs.** Starting from just 3 seed examples, DSL-Monkeys using GPT-5-mini correctly implements 71% of the 200 KernelBench Level 1–2 tasks in TileLang and 18% in ThunderKittens, outperforming test-time methods that attempt each kernel independently (Figure 1). Combined with OpenEvolve (Sharma, 2025a), it achieves 77.5% in TileLang.
2. **Curriculum Decomposition Enables Solving Complex Kernels.** For challenging kernels like linear attention variants, we propose automated decomposition into simpler constituent parts. This allows DSL-Monkeys to bootstrap solutions for novel attention variants where other test-time methods fail, achieving parity with expert-written Triton implementations.
3. **Preliminary Findings from Scaling Up Data Generation for Rare DSLs.** We scale DSL-Monkeys to 10K PyTorch programs from GitHub, successfully implementing 4.7K in TileLang and 1.1K in ThunderKittens. Fine-tuning Qwen-30B-A3B on this dataset improves KernelBench L1+2 accuracy from 17% to 61% when used with DSL-Monkeys, pointing toward a data flywheel for emerging DSLs (Appendix B).

2 PROBLEM STATEMENT

Implementing Programs in Novel DSLs. We consider implementing a collection of PyTorch programs in a target DSL, a common challenge when new DSLs emerge. Each DSL has only a few seed examples and a verifier (compiler + correctness tests). Unlike standard benchmarks that treat each program independently, our goal is to maximize successful implementations across programs.

The Low-Resource DSL Challenge. Emerging GPU DSLs exemplify this challenge. Table 1 shows frameworks like TileLang and ThunderKittens have fewer than 1M publicly available tokens. Figure 2 demonstrates the impact: frontier models achieve below 20% pass@1 on KernelBench Level 1-2 for TileLang and below 2% for ThunderKittens.

Cross-Task Transfer Opportunity. Standard test-time strategies (Table 2) apply methods independently to each kernel. While execution feedback helps (OpenEvolve: 39.5% vs repeated sampling: 24.5% with GPT-5-Mini), these approaches discard successful solutions after verification rather than

DSL	Language	Prog. Model	# Tokens
TileLang	Python	Block	733K
ThunderKittens	C++	Warp	395K
CuTe-DSL	Python	Warp	237K
Triton TLX	Python	Warp	182K
CUDA	C++	Thread	361M
Triton	Python	Block	3M

Table 1: **Estimates of Public Tokens for Emerging GPU DSLs.** We characterize GPU DSLs (Appendix D) by syntactic language, conceptual programming level, and number of high-quality tokens. Details in Appendix E.

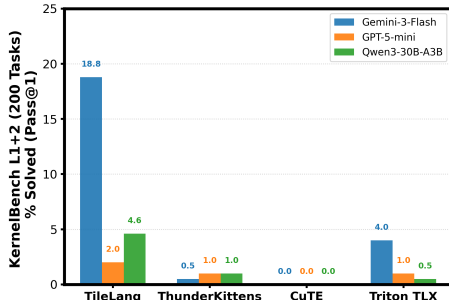


Figure 2: **Models struggle to generate kernels in emerging GPU DSLs.** We measure pass@1 (N=5) on KernelBench Level 1+2 tasks across DSLs. Even frontier models achieve <20% on TileLang and <2% on ThunderKittens.

leveraging them for subsequent tasks. We hypothesize that successful generations contain reusable knowledge applicable to unsolved tasks. To test this, we analyze the 200 KernelBench tasks where repeated sampling (N=25) with Gemini-3-Flash solves 90 kernels. For each of the 110 failures, we check whether required DSL operators and patterns appear in solved kernels. **70% of failed kernels have all required components present in other solved kernels** (91% average coverage). While 2% require novel patterns, the knowledge needed to solve most failures already exists within the batch—it is simply inaccessible when treating tasks independently. This motivates maintaining an archive of successful solutions and retrieving relevant examples for unsolved tasks.

3 DSL-MONKEYS: SELF-GENERATED EXAMPLES FOR DSL KERNELS

Standard test-time approaches treat each problem in isolation, failing to leverage knowledge discovered in related tasks. We introduce DSL-Monkeys, which maintains a persistent archive of verified solutions and retrieves relevant examples for new problems, enabling cross-task knowledge transfer.

DSL-Monkeys Algorithm. DSL-Monkeys follows an iterative bootstrapping process (Algorithm 1 in App. A). We initialize an archive with 3 seed examples from DSL documentation and maintain a set of unsolved problems. At each iteration, we process unsolved problems in mini-batches. For each problem, we retrieve 3 examples from the archive, alternating between cosine similarity and random selection to balance exploitation and exploration (App. K.1). These serve as in-context demonstrations for kernel generation. Generated kernels are validated; passing kernels are added to the archive and marked as solved. Critically, newly solved kernels immediately become retrieval candidates for subsequent mini-batches, enabling intra-iteration bootstrapping.

DSL-Monkeys Effectively Leverages Test-Time Compute. We evaluate DSL-Monkeys against all baselines under identical conditions: 3 seed examples, 25 iterations, temperature 1.0. Table 2 shows results on TileLang and ThunderKittens. Cross-task retrieval substantially outperforms repeated sampling: on TileLang, out of 200 KernelBench tasks, DSL-Monkeys correctly implements 140 tasks with Gemini-3-Flash versus 90 for repeated sampling, and 142 with GPT-5-Mini versus 49. On ThunderKittens, DSL-Monkeys achieves 32 versus 7 for repeated sampling. Cross-task retrieval also outperforms methods using execution feedback on individual problems—on TileLang with Gemini-3-Flash, DSL-Monkeys (140) surpasses Iterative Refinement (102) and OpenEvolve (120), suggesting that retrieving relevant examples from related tasks provides better guidance than compiler errors from the current task alone. Finally, the 155 tasks solved by combining DSL-Monkeys and OpenEvolve (I) show that cross-task retrieval and same-task execution feedback are complementary.

How Retrieval Enables Bootstrapping. We analyze how cross-task retrieval enables knowledge transfer by examining code borrowing patterns and novel example impact. Using LLM-based analysis (Gemini-3-Flash, App. N.1), we identify conceptual code patterns borrowed from examples. On first attempts, successful kernels borrow 60% of code from examples vs. 51% for failures. By iteration 3, successes maintain 60% borrowing while failures show 49%. The fixed-seed baseline shows 49% for both, suggesting when examples lack relevant patterns, borrowing cannot increase even when helpful. Success rate increases from 2.0% when no retrieved examples change to 20.8% when all three are replaced, validating that failures persist due to missing relevant examples (App. M). Remaining failures are mostly synthetic compositions unrepresentative of real-world kernels (App. O).

162
163
164
165
166
167
168
169
170

Test-Time Method	Seeds	TileLang		ThunderKittens	
		gemini-3-flash	gpt-5-mini	gemini-3-flash	gpt-5-mini
Repeated Sampling	3	90	49	7	27
Iterative Refinement	3	102	69	17	25
OpenEvolve	3	120	79	36	34
DSL-Monkeys	3	140	142	32	36
DSL-Monkeys	8	140	152	41	35
DSL-Monkeys+OpenEvolve	3	155	155	35	39

Table 2: **DSL-Monkeys improves DSL kernel synthesis. Results on implementing 200 KernelBench L1-2 tasks correctly (25 attempts).** DSL-Monkeys outperforms baselines (142/200 with GPT-5-Mini, 3 seeds) and combines well with OpenEvolve (155/200 on TileLang).

174
175
176
177
178
179
180
181
182
183

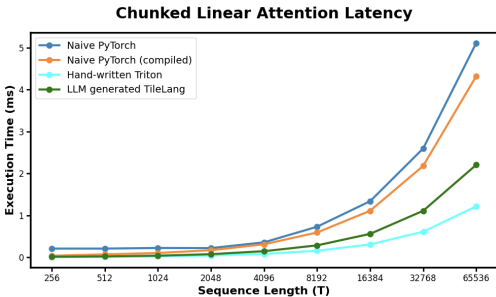


Figure 3: **DSL-Monkeys successful implements Linear Attention Variants that surpass torch.compile.** Chunkwise Linear Attention (Yang & Zhang, 2024a) in TileLang approaches competitive performance with hand-written Triton.

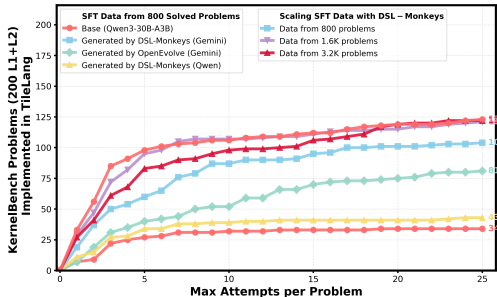


Figure 4: **DSL-Monkeys-generated data from KernelBook teaches Qwen to leverage cross-task transfer.** We show that fine-tuning Qwen3-30B-A3B improves its downstream performance on KernelBench, via both distillation and self-generated data (Appendix B).

184
185
186
187
188
189
190
191
192

4 CASE STUDY: SYNTHESIZING LINEAR ATTENTION KERNELS

Having demonstrated DSL-Monkeys’s effectiveness on KernelBench, we evaluate its transfer to real-world kernels. We focus on 12 linear attention variants (Katharopoulos et al., 2020) from the Flash Linear Attention repository (Yang & Zhang, 2024a), including Based (Arora et al., 2025), Retention (Sun et al., 2023), and Chunked formulations—critical primitives in foundation models like Kimi Linear (Team et al., 2025) and Qwen-next (Qwen Team). These kernels are substantially more complex than KernelBench, requiring intricate recurrent formulations, stateful computations, and expert-level fusion patterns. TileLang has shown promise in this setting (DeepSeek-AI, 2025).

193
194
195
196
197
198
199

Standard approaches fail: Repeated sampling solves 1/12 tasks; OpenEvolve solves 0/12. Even using the 140 TileLang solutions from Section 3 as a starting archive, DSL-Monkeys solves only 2/12; these kernels are out-of-distribution from KernelBench’s standalone primitives. While 85% of attempts compile, they fail due to errors in recurrent state propagation and complex fusion logic.

200
201
202
203
204
205
206
207
208

Curriculum learning via task decomposition. We automatically decompose each kernel into 3–9 interpretable subproblems (Appendix P), such as ‘compute recurrent state update,’ or ‘fuse causal masking with projection,’ to bridge the distribution gap. DSL-Monkeys solves each subproblem using the evolving archive as stepping stones, enabling solving 4/12 variants (Based, Retention, Chunked, ReBased). These implementations exceed 200 lines and coordinate 7+ optimization techniques.

209
210
211
212
213
214
215

Generated kernels achieve expert-level performance. The solved kernels demonstrate sophisticated optimizations including tensor core utilization and pipelining. Figure 3 shows our Chunked Linear Attention achieves 2.1× speedup over Naive PyTorch, outperforms compiled PyTorch, and approaches hand-tuned Triton. DSL-Monkeys outperforms human code by 15% on some formulations (Figure 6), discovering optimizations beyond the training corpus. While direct transfer is limited for OOD problems, DSL-Monkeys with targeted curriculum construction scales to production kernels that challenge expert programmers (Appendix P).

REFERENCES

- 216
217
218 Yaroslav Aksenov, Nikita Balagansky, Sofia Maria Lo Cicero Vaina, Boris Shaposhnikov, Alexey
219 Gorbatoovski, and Daniil Gavrilov. Linear transformers with learnable kernel functions are better
220 in-context models, 2024. URL <https://arxiv.org/abs/2402.10644>.
- 221 Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algo-
222 rithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*,
223 2022.
- 224 Jason Ansel, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznesensky, Bin
225 Bao, Peter Bell, David Berard, Evgeni Burovski, Geeta Chauhan, Anjali Chourdia, Will Constable,
226 Alban Desmaison, Zachary DeVito, Elias Ellison, Will Feng, Jiong Gong, Michael Gschwind, Brian
227 Hirsh, Sherlock Huang, Kshiteej Kalambarkar, Laurent Kirsch, Michael Lazos, Mario Lezcano,
228 Yanbo Liang, Jason Liang, Yinghai Lu, C. K. Luk, Bert Maher, Yunjie Pan, Christian Puhersch,
229 Matthias Reso, Mark Saroufim, Marcos Yukio Siraichi, Helen Suk, Shunting Zhang, Michael Suo,
230 Phil Tillet, Xu Zhao, Eikan Wang, Keren Zhou, Richard Zou, Xiaodong Wang, Ajit Mathews,
231 William Wen, Gregory Chanan, Peng Wu, and Soumith Chintala. Pytorch 2: Faster machine
232 learning through dynamic python bytecode transformation and graph compilation. In *Proceedings*
233 *of the 29th ACM International Conference on Architectural Support for Programming Languages*
234 *and Operating Systems, Volume 2*, ASPLOS '24, pp. 929–947, New York, NY, USA, 2024.
235 Association for Computing Machinery. ISBN 9798400703850. doi: 10.1145/3620665.3640366.
236 URL <https://doi.org/10.1145/3620665.3640366>.
- 237 Simran Arora, Sabri Eyuboglu, Michael Zhang, Aman Timalsina, Silas Alberti, Dylan Zinsley,
238 James Zou, Atri Rudra, and Christopher Ré. Simple linear attention language models balance the
239 recall-throughput tradeoff, 2025. URL <https://arxiv.org/abs/2402.18668>.
- 240 Ian Barber. Cute-dsl. <https://ianbarber.blog/2025/07/04/cute-dsl/>, 2025. URL
241 <https://ianbarber.blog/2025/07/04/cute-dsl/>. Ian’s Blog, accessed January
242 2026.
- 243 Carlo Baronio, Pietro Marsella, Ben Pan, Simon Guo, and Silas Alberti. Kevin: Multi-turn rl for
244 generating cuda kernels, 2025. URL <https://arxiv.org/abs/2507.11948>.
- 245 Dan Biderman, Jacob Portes, Jose Javier Gonzalez Ortiz, Mansheej Paul, Philip Greengard, Connor
246 Jennings, Daniel King, Sam Havens, Vitaliy Chiley, Jonathan Frankle, Cody Blakeney, and John P.
247 Cunningham. Lora learns less and forgets less, 2024. URL <https://arxiv.org/abs/2405.09673>.
- 248
249
- 250 Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and
251 Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling,
252 2024. URL <https://arxiv.org/abs/2407.21787>.
- 253 Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinck-
254 ney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, Arjun Guha, Michael
255 Greenberg, and Abhinav Jangda. Multipl-e: A scalable and extensible approach to benchmarking
256 neural code generation, 2022. URL <https://arxiv.org/abs/2208.08227>.
- 257 Audrey Cheng, Shu Liu, Melissa Pan, Zhifei Li, Bowen Wang, Alex Krentsel, Tian Xia, Mert Cemri,
258 Jongseok Park, Shuo Yang, Jeff Chen, Lakshya Agrawal, Aditya Desai, Jiarong Xing, Koushik
259 Sen, Matei Zaharia, and Ion Stoica. Barbarians at the gate: How ai is upending systems research,
260 2025. URL <https://arxiv.org/abs/2510.06189>.
- 261
262 Christopher Choy. Cutedsl basics. <https://chrischoy.org/posts/cutedsl-basics/>,
263 2025. URL <https://chrischoy.org/posts/cutedsl-basics/>. Online technical
264 note, accessed January 2026.
- 265 Colfax International Research. A user’s guide to flexattention in flashat-
266 tention & cutedsl. [https://research.colfax-intl.com/](https://research.colfax-intl.com/a-users-guide-to-flexattention-in-flash-attention-cute-dsl/)
267 [a-users-guide-to-flexattention-in-flash-attention-cute-dsl/](https://research.colfax-intl.com/a-users-guide-to-flexattention-in-flash-attention-cute-dsl/),
268 2025. URL [https://research.colfax-intl.com/](https://research.colfax-intl.com/a-users-guide-to-flexattention-in-flash-attention-cute-dsl/)
269 [a-users-guide-to-flexattention-in-flash-attention-cute-dsl/](https://research.colfax-intl.com/a-users-guide-to-flexattention-in-flash-attention-cute-dsl/).
Colfax International Research blog, accessed January 2026.

- 270 Dao-AILab. Flash-attention. <https://github.com/Dao-AILab/flash-attention>,
271 2024a. URL <https://github.com/Dao-AILab/flash-attention>.
272
- 273 Dao-AILab. Quack. <https://github.com/Dao-AILab/quack>, 2024b. URL <https://github.com/Dao-AILab/quack>. GitHub repository.
274
- 275 DeepSeek-AI. Deepseek-v3.2-exp: Boosting long-context efficiency with deepseek sparse attention,
276 2025.
277
- 278 Ryan Ehrlich, Bradley Brown, Jordan Juravsky, Ronald Clark, Christopher Ré, and Azalia Mirhoseini.
279 Codemonkeys: Scaling test-time compute for software engineering, 2025. URL <https://arxiv.org/abs/2501.14723>.
280
- 281 Facebook Experimental. Triton tlx: Tensor language extensions, 09 2025. URL <https://github.com/facebookexperimental/triton/tree/tlxa>.
282
283
- 284 Zacharias V. Fisches, Sahan Paliskara, Simon Guo, Alex Zhang, Joe Spisak, Chris Cummins, Hugh
285 Leather, Gabriel Synnaeve, Joe Isaacson, Aram Markosyan, and Mark Saroufim. KernelLLM:
286 Making kernel development more accessible, 6 2025. URL <https://huggingface.co/facebook/KernelLLM>.
287
- 288 Etash Guha, Ryan Marten, Sedrick Keh, Negin Raoof, Georgios Smyrnis, Hritik Bansal, Marianna
289 Nezhurina, Jean Mercat, Trung Vu, Zayne Sprague, Ashima Suvarna, Benjamin Feuer, Liangyu
290 Chen, Zaid Khan, Eric Frankel, Sachin Grover, Caroline Choi, Niklas Muennighoff, Shiye Su,
291 Wanjia Zhao, John Yang, Shreyas Pimpalgaonkar, Kartik Sharma, Charlie Cheng-Jie Ji, Yichuan
292 Deng, Sarah Pratt, Vivek Ramanujan, Jon Saad-Falcon, Jeffrey Li, Achal Dave, Alon Albalak,
293 Kushal Arora, Blake Wulfe, Chinmay Hegde, Greg Durrett, Sewoong Oh, Mohit Bansal, Saadia
294 Gabriel, Aditya Grover, Kai-Wei Chang, Vaishaal Shankar, Aaron Gokaslan, Mike A. Merrill,
295 Tatsunori Hashimoto, Yejin Choi, Jenia Jitsev, Reinhard Heckel, Maheswaran Sathiamoorthy,
296 Alexandros G. Dimakis, and Ludwig Schmidt. Openthoughts: Data recipes for reasoning models,
297 2025. URL <https://arxiv.org/abs/2506.04178>.
- 298 HazyResearch. Thunderkittens. <https://github.com/HazyResearch/ThunderKittens>,
299 2024. URL <https://github.com/HazyResearch/ThunderKittens>. GitHub repository.
300
301
- 302 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
303 and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL <https://arxiv.org/abs/2106.09685>.
304
- 305 Shengran Hu, Cong Lu, and Jeff Clune. Automated design of agentic systems, 2025. URL <https://arxiv.org/abs/2408.08435>.
306
307
- 308 Sathvik Joel, Jie JW Wu, and Fatemeh H. Fard. A survey on llm-based code generation for low-
309 resource and domain-specific programming languages, 2025. URL <https://arxiv.org/abs/2410.03981>.
310
- 311 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are
312 rns: Fast autoregressive transformers with linear attention, 2020. URL <https://arxiv.org/abs/2006.16236>.
313
314
- 315 Kabir Khandpur, Kilian Lieret, Carlos E. Jimenez, Ofir Press, and John Yang. Swe-bench multilin-
316 gual, 2024. URL [https://huggingface.co/datasets/SWE-bench/SWE-bench_](https://huggingface.co/datasets/SWE-bench/SWE-bench_Multilingual)
317 [Multilingual](https://huggingface.co/datasets/SWE-bench/SWE-bench_Multilingual).
- 318 kiui.moe. Cute — cuda dsl documentation. <https://note.kiui.moe/cuda/cute/#print>,
319 2025. URL <https://note.kiui.moe/cuda/cute/#print>. Online technical note, ac-
320 cessed January 2026.
321
- 322 Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Carlos Muñoz Ferrandis,
323 Yacine Jernite, Margaret Mitchell, Sean Hughes, Thomas Wolf, Dzmitry Bahdanau, Leandro von
Werra, and Harm de Vries. The stack: 3 tb of permissively licensed source code. *Preprint*, 2022.

- 324 Itay Levy, Ben Bogin, and Jonathan Berant. Diverse demonstrations improve in-context compositional
325 generalization. *arXiv preprint arXiv:2212.06800*, 2022.
326
- 327 Jianling Li, Shangzhan Li, Zhenye Gao, Qi Shi, Yuxuan Li, Zefan Wang, Jiacheng Huang, Haojie
328 Wang, Jianrong Wang, Xu Han, Zhiyuan Liu, and Maosong Sun. Tritonbench: Benchmarking
329 large language model capabilities for generating triton operators, 2025. URL <https://arxiv.org/abs/2502.14752>.
330
- 331 Gang Liao, Hongsen Qin, Ying Wang, Alicia Golden, Michael Kuchnik, Yavuz Yetim, Jia Jiunn Ang,
332 Chunli Fu, Yihan He, Samuel Hsia, Zewei Jiang, Dianshi Li, Uladzimir Pashkevich, Varna Puvvada,
333 Feng Shi, Matt Steiner, Ruichao Xiao, Nathan Yan, Xiayu Yu, Zhou Fang, Abdul Zainul-Abedin,
334 Ketan Singh, Hongtao Yu, Wenyuan Chi, Barney Huang, Sean Zhang, Noah Weller, Zach Marine,
335 Wyatt Cook, Carole-Jean Wu, and Gaoxiang Liu. Kernelevolve: Scaling agentic kernel coding
336 for heterogeneous ai accelerators at meta, 2025. URL <https://arxiv.org/abs/2512.23236>.
337
- 338
- 339 Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. What
340 makes good in-context examples for gpt-3? *arXiv preprint arXiv:2101.06804*, 2021.
341
- 342 Mingjie Liu, Nathaniel Pinckney, Brucek Khailany, and Haoxing Ren. Verilogeval: Evaluating large
343 language models for verilog code generation, 2023. URL <https://arxiv.org/abs/2309.07544>.
344
- 345 Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane
346 Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov,
347 Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul,
348 Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii,
349 Nii Osaе Osaе Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan
350 Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muhtasham Oblokulov,
351 Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri
352 Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, Julian McAuley, Han Hu, Torsten
353 Scholak, Sebastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostofa
354 Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes,
355 Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder 2 and the stack v2:
356 The next generation, 2024.
- 357 Thinking Machines. Announcing tinker: An llm-fine-tuning ecosystem. [https://](https://thinkingmachines.ai/blog/announcing-tinker/)
358 thinkingmachines.ai/blog/announcing-tinker/, 2025. Accessed: 2026-01-22.
359
- 360 Jean-Baptiste Mouret and Jeff Clune. Illuminating search spaces by mapping elites, 2015. URL
361 <https://arxiv.org/abs/1504.04909>.
- 362 Alexander Novikov, Ngàn Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wag-
363 ner, Sergey Shirobokov, Borislav Kozlovskii, Francisco J. R. Ruiz, Abbas Mehrabian, M. Pawan
364 Kumar, Abigail See, Swarat Chaudhuri, George Holland, Alex Davies, Sebastian Nowozin, Push-
365 meet Kohli, and Matej Balog. Alphaevolve: A coding agent for scientific and algorithmic discovery,
366 2025. URL <https://arxiv.org/abs/2506.13131>.
367
- 368 NVIDIA. Cutlass python dsl documentation. [https://docs.nvidia.com/cutlass/](https://docs.nvidia.com/cutlass/latest/media/docs/pythonDSL/cute_dsl.html)
369 [latest/media/docs/pythonDSL/cute_dsl.html](https://docs.nvidia.com/cutlass/latest/media/docs/pythonDSL/cute_dsl.html), 2025. URL https://docs.nvidia.com/cutlass/latest/media/docs/pythonDSL/cute_dsl.html.
370 NVIDIA official documentation.
371
- 372 NVIDIA. Cutlass: Cuda templates and python dsls for high-performance linear algebra, 2026.
373 URL <https://github.com/NVIDIA/cutlass>. includes CuTe DSL — a Python domain-
374 specific language for GPU code generation.
375
- 376 Anne Ouyang, Simon Guo, Simran Arora, Alex L. Zhang, William Hu, Christopher Ré, and Azalia
377 Mirhoseini. Kernelbench: Can llms write efficient gpu kernels?, 2025a. URL <https://arxiv.org/abs/2502.10517>.

- 378 Anne Ouyang, Azalia Mirhoseini, and Percy Liang. Surprisingly fast ai-generated kernels we didn't
379 mean to publish (yet). <https://crfm.stanford.edu/2025/05/28/fast-kernels.html>, May 2025b.
- 380
381
- 382 Sahar Paliskara and Mark Saroufim. Kernelbook, 5 2025. URL <https://huggingface.co/datasets/GPUMODE/KernelBook>.
- 383
- 384 Chengwei Qin, Aston Zhang, Chen Chen, Anirudh Dagar, and Wenming Ye. In-context learning with
385 iterative demonstration selection. *arXiv preprint arXiv:2310.09881*, 2023.
- 386
- 387 Qwen Team. Qwen3-coder-next technical report. Technical report. URL https://github.com/QwenLM/Qwen3-Coder/blob/main/qwen3_coder_next_tech_report.pdf. Accessed: 2026-02-03.
- 388
389
- 390 Vishnu Sarukkai, Zhiqiang Xie, and Kayvon Fatahalian. Self-generated in-context examples improve
391 LLM agents for sequential decision-making tasks. In *The Thirty-ninth Annual Conference on*
392 *Neural Information Processing Systems*, 2025. URL <https://openreview.net/forum?id=WdL3058gde>.
- 393
- 394 John Schulman and Thinking Machines Lab. Lora without regret. *Thinking Machines Lab: Connectionism*, 2025. doi: 10.64434/tml.20250929. <https://thinkingmachines.ai/blog/lora/>.
- 395
396
- 397 Asankhaya Sharma. Openevolve: an open-source evolutionary coding agent, 2025a. URL <https://github.com/algorithmicsuperintelligence/openevolve>.
- 398
399
- 400 Asankhaya Sharma. Automated discovery of high-performance gpu kernels with openevolve. <https://huggingface.co/blog/codelion/openevolve-gpu-kernel-discovery>,
401 June 27 2025b. Hugging Face blog post.
- 402
- 403 Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and
404 Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023. URL
405 <https://arxiv.org/abs/2303.11366>.
- 406
- 407 Benjamin Spector, Aaryan Singhal, Simran Arora, and Christopher Ré. Gpus go brrr. <https://hazyresearch.stanford.edu/blog/2024-05-12-tk>, 2024a. URL <https://hazyresearch.stanford.edu/blog/2024-05-12-tk>. Hazy Research blog post,
408 May 12, 2024; introduces the ThunderKittens embedded DSL.
- 409
410
- 411 Benjamin F. Spector, Simran Arora, Aaryan Singhal, Daniel Y. Fu, and Christopher Ré. Thunderkittens: Simple, fast, and adorable ai kernels, 2024b. URL <https://arxiv.org/abs/2410.20399>.
- 412
413
- 414 Stuart Sul. 1.5x faster moe training with custom mxfp8 kernels. <https://cursor.com/blog/kernels>, August 2025. URL <https://cursor.com/blog/kernels>. Achieving a 3.5x
415 MoE layer speedup with a complete rebuild for Blackwell GPUs.
- 416
417
- 418 Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and
419 Furu Wei. Retentive network: A successor to transformer for large language models, 2023. URL
420 <https://arxiv.org/abs/2307.08621>.
- 421
- 422 Garrett Tanzer, Mirac Suzgun, Eline Visser, Dan Jurafsky, and Luke Melas-Kyriazi. A benchmark
423 for learning to translate a new language from one grammar book, 2024. URL <https://arxiv.org/abs/2309.16575>.
- 424
- 425 Gemini Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context,
426 2024. URL <https://arxiv.org/abs/2403.05530>.
- 427
- 428 Kimi Team, Yu Zhang, Zongyu Lin, Xingcheng Yao, Jiayi Hu, Fanqing Meng, Chengyin Liu, Xin
429 Men, Songlin Yang, Zhiyuan Li, Wentao Li, Enzhe Lu, Weizhou Liu, Yanru Chen, Weixin Xu,
430 Longhui Yu, Yejie Wang, Yu Fan, Longguang Zhong, Enming Yuan, Dehao Zhang, Yizhi Zhang,
431 T. Y. Liu, Haiming Wang, Shengjun Fang, Weiran He, Shaowei Liu, Yiwei Li, Jianlin Su, Jiezhong Qiu, Bo Pang, Junjie Yan, Zhejun Jiang, Weixiao Huang, Bohong Yin, Jiacheng You, Chu Wei, Zhengtao Wang, Chao Hong, Yutian Chen, Guanduo Chen, Yucheng Wang, Huabin Zheng, Feng

- 432 Wang, Yibo Liu, Mengnan Dong, Zheng Zhang, Siyuan Pan, Wenhao Wu, Yuhao Wu, Longyu
433 Guan, Jiawen Tao, Guohong Fu, Xinran Xu, Yuzhi Wang, Guokun Lai, Yuxin Wu, Xinyu Zhou,
434 Zhilin Yang, and Yulun Du. Kimi linear: An expressive, efficient attention architecture, 2025.
435 URL <https://arxiv.org/abs/2510.26692>.
- 436
437 Qwen Team. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- 438 TileLang Developers. Tilelang. <https://tilelang.com/>, 2025. URL [https://](https://tilelang.com/)
439 tilelang.com/. Official project website.
- 440
441 Philippe Tillet, H. T. Kung, and David Cox. Triton: an intermediate language and compiler for
442 tiled neural network computations. In *Proceedings of the 3rd ACM SIGPLAN International*
443 *Workshop on Machine Learning and Programming Languages*, MAPL 2019, pp. 10–19, New
444 York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450367196. doi:
445 10.1145/3315508.3329973. URL <https://doi.org/10.1145/3315508.3329973>.
- 446 Lukas Veitner. An applied introduction to cutedsl. [https://veitner.bearblog.](https://veitner.bearblog.dev/an-applied-introduction-to-cutedsl/)
447 [dev/an-applied-introduction-to-cutedsl/](https://veitner.bearblog.dev/an-applied-introduction-to-cutedsl/), 2024. URL [https://veitner.](https://veitner.bearblog.dev/an-applied-introduction-to-cutedsl/)
448 [bearblog.dev/an-applied-introduction-to-cutedsl/](https://veitner.bearblog.dev/an-applied-introduction-to-cutedsl/). Bear Blog, accessed
449 January 2026.
- 450 Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev,
451 Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In
452 *International Conference on Machine Learning*, pp. 35151–35174. PMLR, 2023.
- 453
454 Lei Wang, Yu Cheng, Yining Shi, Zhengju Tang, Zhiwen Mo, Wenhao Xie, Lingxiao Ma, Yuqing
455 Xia, Jilong Xue, Fan Yang, and Zhi Yang. Tilelang: A composable tiled programming model for ai
456 systems, 2025. URL <https://arxiv.org/abs/2504.17577>.
- 457 Jiin Woo, Shaowei Zhu, Allen Nie, Zhen Jia, Yida Wang, and Youngsuk Park. Tritonrl: Training
458 llms to think and code triton without cheating, 2025. URL [https://arxiv.org/abs/2510.](https://arxiv.org/abs/2510.17891)
459 [17891](https://arxiv.org/abs/2510.17891).
- 460
461 John Yang, Kilian Lieret, Carlos E. Jimenez, Alexander Wettig, Kabir Khandpur, Yanzhe Zhang,
462 Binyuan Hui, Ofir Press, Ludwig Schmidt, and Diyi Yang. Swe-smith: Scaling data for software
463 engineering agents, 2025. URL <https://arxiv.org/abs/2504.21798>.
- 464 Songlin Yang and Yu Zhang. Fla: A triton-based library for hardware-efficient implementations
465 of linear attention mechanism, January 2024a. URL [https://github.com/fla-org/](https://github.com/fla-org/flash-linear-attention)
466 [flash-linear-attention](https://github.com/fla-org/flash-linear-attention).
- 467
468 Songlin Yang and Yu Zhang. Fla: A triton-based library for hardware-efficient implementations
469 of linear attention mechanism, January 2024b. URL [https://github.com/fla-org/](https://github.com/fla-org/flash-linear-attention)
470 [flash-linear-attention](https://github.com/fla-org/flash-linear-attention).
- 471 Mert Yuksekgonul, Daniel Kocaja, Xinhao Li, Federico Bianchi, Jed McCaleb, Xiaolong Wang, Jan
472 Kautz, Yejin Choi, James Zou, Carlos Guestrin, and Yu Sun. Learning to discover at test time.
473 *arXiv preprint arXiv:2601.16175*, 2026. URL <https://arxiv.org/abs/2601.16175>.
- 474 Alex L. Zhang, Tim Kraska, and Omar Khattab. Recursive language models, 2025a. URL [https:](https://arxiv.org/abs/2512.24601)
475 [//arxiv.org/abs/2512.24601](https://arxiv.org/abs/2512.24601).
- 476
477 Genghan Zhang, Weixin Liang, Olivia Hsu, and Kunle Olukotun. Adaptive self-improvement llm
478 agentic system for ml library development. *arXiv preprint arXiv:2502.02534*, 2025b.
- 479
480 Jenny Zhang, Shengran Hu, Cong Lu, Robert Lange, and Jeff Clune. Darwin godel machine: Open-
481 ended evolution of self-improving agents, 2025c. URL [https://arxiv.org/abs/2505.](https://arxiv.org/abs/2505.22954)
482 [22954](https://arxiv.org/abs/2505.22954).
- 483
484
485

A ALGORITHM

Here we show the algorithm for DSL-Monkeys.

Algorithm 1 DSL-Monkeys: In-Context Bootstrapping for Kernel Generation

```

1:  $\mathcal{A} \leftarrow \text{SEEDEXAMPLES}()$ ,  $\mathcal{U} \leftarrow \text{ALLPROBLEMS}()$  ▷ Init archive, unsolved
2: for iteration  $t = 1$  to  $T_{\max}$  do
3:   for each mini-batch  $B \subseteq \mathcal{U}$ , problem  $p \in B$  do
4:     if  $t \bmod 2 = 1$  then
5:        $E \leftarrow \text{RETRIEVESIMILAR}(\mathcal{A}, p, k = 3)$  ▷ Cosine sim. retrieval on odd iterations
6:     else
7:        $E \leftarrow \text{RETRIEVERANDOM}(\mathcal{A}, p, k = 3)$  ▷ Random retrieval on even iterations
8:      $kernel \leftarrow \text{GENERATEKERNEL}(\text{LLM}, p, E)$ 
9:     if  $\text{TESTKERNEL}(kernel)$  then
10:       $\mathcal{A} \leftarrow \mathcal{A} \cup \{(p, kernel)\}$  ▷ Check correctness, add to archive
11:       $\mathcal{U} \leftarrow \mathcal{U} \setminus \{p\}$  ▷ Mark solved
12: return  $\mathcal{A}$ 

```

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

B SCALING DATA UNLOCKS CROSS-TASK LEARNING

Generating DSL Programs at Scale. Previous sections demonstrated the effectiveness of DSL-Monkeys in implementing collections of programs in novel DSLs—200 KernelBench tasks (Section 3) and 10 linear attention variants (Section 4), showing the value of cross-task transfer paired with curriculum. Can we scale to more tasks, capturing greater variety, and generate corresponding DSL code? We turn to the massively available PyTorch programs on the internet, such as KernelBook (Paliskara & Saroufim, 2025), a collection of PyTorch programs from The Stack (Kocetkov et al., 2022) that are executable through `torch.compile`. We run DSL-Monkeys on this collection of 10,082 valid KernelBook programs (after decontamination against KernelBench) with up to 5 attempts per task. For TileLang, we solve 4.7K programs with Gemini-3-Flash and 1.1K with Qwen3-30B-A3B. For the more challenging ThunderKittens, we solve 500 with Gemini-3-Flash.

Furthermore, a byproduct of this process is a growing **dataset** of high-quality, annotated code samples for these low-resource DSLs: each successfully implemented program provides a trajectory containing (PyTorch Specification, Implementation Plan, verified DSL kernel). As emerging DSLs each have $< 1M$ public tokens (Section 2), these trajectories represent a significant addition to publicly available code data, totaling 14M tokens of verified code paired with reasoning traces.

Enable Weaker Model to Leverage Cross-Task Learning. Leveraging cross-task transfer is highly effective in this low-resource regime, as demonstrated with frontier models in Section 3. However, weaker models struggle to leverage this: Qwen3-30B-A3B solves only 34 KernelBench Level 1+2 problems and plateaus quickly, even after 25 test-time iterations with DSL-Monkeys (Figure 4). How could we enable weaker models to leverage cross-task transfer and unlock test-time compute effectively?

The data generated at scale from KernelBook can bridge this gap. Fine-tuning Qwen (see Appendix L) on expert (Gemini) trajectories from just 800 solved problems improves test-time performance from 34 \rightarrow 104. Fine-tuning on the same number of OpenEvolve-generated solutions also helps (34 \rightarrow 81), but DSL-Monkey trajectories are more effective, as they explicitly teach model to take advantage of cross-task transfer beyond improving its ability in implementing in DSL. Further data scaling reaches 120+, approaching the teacher model’s 140. Training Qwen on its own DSL-Monkeys-generated trajectories (1.6K samples) also yields gains (34 \rightarrow 43). While smaller, this shows models can bootstrap DSL capability without a stronger teacher, pointing toward a data flywheel as open-source models continue to improve.

C RELATED WORK AND LIMITATIONS

C.1 AI FOR GPU KERNEL GENERATION

GPU kernel optimization is central to deep learning performance, yet it is challenging for experts to write, motivating recent interest in AI-assisted kernel generation Ouyang et al. (2025a); Li et al. (2025). Prior approaches have shown promise on more established framework such as CUDA and Triton, with evolutionary search Novikov et al. (2025); Liao et al. (2025), reinforcement learning Baronio et al. (2025); Woo et al. (2025), and test-time training Yuksekogonul et al. (2026). In particular, they require strong enough prior in the target language and face significant challenges with data scarcity in newer DSLs as shown in our evaluation in 3.

As GPU kernel data is rare, synthetic data generation has been explored, notably through leveraging compiler systems as Paliskara & Saroufim (2025); Fisches et al. (2025) generate (PyTorch, Triton) pairs via Torch Inductor Ansel et al. (2024); however such transpilers only exist for limited amount of DSLs. Our work differs by focusing on DSL kernel generation from minimal examples and leveraging the model’s inherent in-context learning capability, making it applicable to more DSLs.

C.2 BOOTSTRAPPING CAPABILITY IN DATA-SCARCE DOMAINS

Our approach addresses capability bootstrapping in low-resource settings through two mechanisms: in-context learning and synthetic data generation leveraging verifiers.

In-context learning from self-generated examples. Prior work has established that the choice of in-context examples significantly impacts model performance Liu et al. (2021); Levy et al. (2022); Qin et al. (2023), with theoretical characterizations of in-context learning Akyürek et al. (2022); Von Oswald et al. (2023). This is particularly valuable in low-resource domains such as rare languages Tanzer et al. (2024); Team (2024). Recent work has explored using self-generated examples to build growing libraries: Sarukkai et al. (2025) show self-generated examples improve LLM agent performance, and Zhang et al. (2025b) leverage in-context learning from an expanding library of functions for ML library development. We apply these principles to test-time kernel generation, where each attempt at any problem can produce verified solutions that serve as examples for future attempts at different problems, enabling cross-task knowledge transfer in data-scarce DSL settings.

Synthetic data leveraging verifiers Solutions generated through test-time problem-solving can also serve as training data for domain specialization. Recent work has shown synthetic data improves model capabilities in various domains, including repository-level bug fixing Yang et al. (2025), and self-generated examples have been used for fine-tuning in agent tasks Sarukkai et al. (2025). In Section B, We demonstrate this pathway is particularly effective for emerging DSLs where human-written training data is scarce, as DSL-Monkeys enables generation via in-context learning and execution-based verifier as way to ensure data quality.

C.3 LIMITATIONS

We note here that our system prioritizes simplicity over specialized method engineering. For instance, we intentionally omit building in profiler feedback to iterative refinement, which could improve kernel quality but adds more complexity. Also, while our recipe (self-generated examples + verification) should generalize to other DSLs, it requires a reliable execution-based verifier and a task distribution where solved examples serve as useful demonstrations. Finally, our approach requires API credit compute (approx. \$40 USD per run with Gemini-3-Flash / \$100 with GPT-5-Mini with the settings shown in Figure 1) from iterative generation and evaluation, which we consider justified for GPU kernels where small gains save substantial developer time and deployment costs.

D GPU DSL FUNDAMENTALS

Modern GPUs have a hierarchical programming model. On NVIDIA GPUs, a *thread* is the fundamental unit of execution, 32 threads form a *warp*, and up to 32 warps (1024 threads) form a *block*. Blocks execute on Streaming Multiprocessors (SMs) with dedicated computational and memory resources. Recent DSLs for GPU kernel development differ in which level of this hierarchy they expose, their host language (Python or C++), and the amount of publicly available code.

We focus on recently-introduced DSLs with minimal public code representation. TileLang Wang et al. (2025) is a Python-based DSL that operates at higher abstraction and relies on compiler optimizations, enabling DeepSeek’s Sparse Attention DeepSeek-AI (2025). ThunderKittens Spector et al. (2024b) is a C++ embedded DSL with opinionated primitives operating over 16×16 tiles, enabling practitioners to author high-performance kernels Sul (2025). Both were released within the last year. Table 1 characterizes these DSLs by their syntax, programming model, and volume of publicly available documentation and code samples. Each provides minimal documentation with 3-8 code examples, which we use as seed examples for our experiments.

E COUNTING TOKENS

We determine the number of high-quality, publicly-available tokens for each DSL by running the Qwen (Team, 2025) tokenizer through various sources of data:

- **TileLang.** We scrape the TileLang documentation (TileLang Developers, 2025) and paper (Wang et al., 2025).
- **ThunderKittens.** We scrape the TK Github repository (HazyResearch, 2024), paper (Spector et al., 2024b), blog (Spector et al., 2024a), and a publicly available onboarding document linked from the ThunkerKittens repo.
- **CuTe DSL.** We aggregate pages from Nvidia’s CuTe DSL documentation (NVIDIA, 2025), Github repositories with CuTe DSL kernels (Dao-AILab, 2024b), and various online blog posts (kiui.moe, 2025; Barber, 2025; Veitner, 2024; Colfax International Research, 2025; Choy, 2025).
- **Triton TLX.** We scrape the Triton TLX Github repository (Experimental, 2025).
- **Triton.** We count the number of tokens from GPUMODE’s Triton KernelBook (Paliskara & Saroufim, 2025).
- **CUDA.** We analyze the *the-stack-v2-dedup* dataset on HuggingFace (Lozhkov et al., 2024) and count the total number of CUDA bytes. Based on previous byte to token ratios, we divide the number of bytes by 4 to arrive at our token count.

F EVALUATION SETUP

F.1 KERNELBENCH TASK FORMAT

We evaluate on KernelBench tasks Ouyang et al. (2025a) from the most recent v0.2 version. KernelBench is a standardized benchmark where each task consists of: (1) a PyTorch reference implementation that defines the operator’s functionality, (2) input-generation routines that specify tensor dimensions and batch sizes, and (3) correctness and performance testing infrastructure. The benchmark is organized into four difficulty levels, from basic single-kernel operators (Level 1) to complex fused patterns (Level 2) and full model architectures (Levels 3-4). To ensure reliable benchmarking, the v0.2 version scales tensor dimensions and batch sizes for Level 1–2 tasks so each test runs in roughly 1–15ms on an H100, avoiding cases where kernel launch overhead dominates the timing.

F.2 EVALUATION SETTING

We run evaluation on NVIDIA H100 GPUs on Modal. The evaluation container uses CUDA 12.8.0 and PyTorch 2.9.0. For each task, we compile and execute the generated kernel and require agreement with the PyTorch reference across 5 correctness trials, using the benchmark’s dtype-specific tolerances. We measure runtime using Triton’s `do_bench` Tillet et al. (2019) with a warmup budget of 25 ms and a repetition budget of 100 ms.

F.3 EVALUATION RESULTS

Table 3: Model Performance Across Different Tasks (N=5 for both effort levels). Note: Only genuine successes are counted; reward hacking attempts are excluded.

Task	Model	Level 1			Level 2		
		Acc.	Cov.	Pass@1	Acc.	Cov.	Pass@1
TileLang	Gemini-3-flash-preview	38%	58%	35.6%	2%	7%	2%
	GPT-5-mini	6%	18%	4%	0.4%	2%	0%
	Qwen3-30B-A3B	9%	13%	9.2%	0%	0%	0%
ThunkerKittens	Gemini-3-flash-preview	0.8%	3%	1%	0%	0%	0%
	GPT-5-mini	1.4%	5%	2%	0%	0%	0%
	Qwen3-30B-A3B	0.6%	2%	2%	0%	0%	0%
CuTe-DSL	Gemini-3-flash-preview	0.4%	2%	0%	0%	0%	0%
	GPT-5-mini	0.2%	1%	0%	0%	0%	0%
	Qwen3-30B-A3B	0%	0%	0%	0%	0%	0%
TLX	Gemini-3-flash-preview	6%	16%	8%	0.2%	1%	0%
	GPT-5-mini	1.4%	7%	2%	0%	0%	0%
	Qwen3-30B-A3B	0.6%	3%	1%	0%	0%	0%

810 G CONTEXT CONSTRUCTION

811
812 Models without DSL-specific information perform poorly, so we provide strong baseline context for
813 all methods. For each DSL, we provide: (1) a DSL-specific guide synthesized from documentation
814 using recursive language models Zhang et al. (2025a), and (2) 3 representative in-context examples
815 (see Table 4) drawn from each DSL’s documentation. These 3 examples serve as our seed set.
816 Following prior work Zhang et al. (2025c); Hu et al. (2025); Ouyang et al. (2025b), we use a 2-
817 stage generation setup where the model first generates an implementation plan, then produces code
818 conditioned on the plan.

819 G.1 DSL LANGUAGE GUIDE

820
821 Because frontier models are unlikely to have been trained on newly released DSLs, we supply them
822 with the relevant language knowledge at inference time. To compress the large volume of scraped
823 documentation into a single prompt-friendly reference, we synthesize domain-specific language
824 guides using recursive language models (RLMs) to handle the long context Zhang et al. (2025a).
825 Specifically, we run an RLM instance of GPT-5.1 with max recursive depth 7. The prompt to the
826 RLM is shown below.

827 DSL Language Guide RLM Prompt

830 You are given extensive documentation, papers, blogposts, and code
831 for a GPU **<DSL>**. Your task is to synthesize this material into a
832 structured guide that teaches a performance engineer how to *think*
833 in this DSL; not just use it, but learn its abstraction so they can
834 translate any PyTorch operation into an efficient kernel plan, which
835 targets the DSL features.

836 Do not use bullet points; explain in clear natural language. The
837 goal is to convey the DSL’s philosophy as accurately as possible, not
838 necessarily its syntax.

839 Start the prompt with: "You are an expert GPU kernel engineer
840 specializing in translating PyTorch operations into **<DSL>**
841 implementation plans. Your task is to analyze PyTorch code and
842 produce detailed, structured natural language plans that describe how
843 to implement the operations efficiently in **<DSL>**. You do not write
844 **<DSL>** code. You produce implementation plans that a **<DSL>** programmer
845 could follow to produce the kernel that leverages the **<DSL>** feature.
846 Your plans must demonstrate a deep understanding of **<DSL>**
847 abstractions, philosophy, and the hardware realities of modern GPUs."

848 The rest of your prompt should address the following:

849 ### 1. What is **<DSL>** and what problem does it solve?

850 What gap does it fill between PyTorch and raw CUDA? What is the main
851 abstraction that the **<DSL>** provides?

852 ### 2. What is the fundamental unit of thought? What does it mean
853 to write **<DSL>**?

854 Every DSL has a primary abstraction that defines how programmers
855 should decompose problems. Identify it precisely.

856 ### 3. What does the compiler handle vs. what must the programmer
857 specify?

858 State and explain what the **<DSL>** handles automatically, exposes for
859 optional control, or requires explicit specification. The goal
860 should be that a programmer should know what to specify and what to
861 leave alone.

862 ### 4. How does **<DSL>** model the memory hierarchy?

863 GPU DSLs typically address global -> shared -> register data movement.
Describe: - What primitives or types represent each memory tier

```

864
865 (actually use the <DSL> documentation to find the names of the
866 primitives) - The pattern for staging data through the hierarchy -
867 How reuse is expressed - Any DSL-specific features (swizzling, TMA
868 integration, typed tiles, etc.)
869
870 ### 5. What are the core compute primitives?
871
872 List the key operations the DSL provides and what hardware they map
873 to: - Matrix multiply / MMA primitives (and tensor core
874 requirements) - Copy/load primitives (and async capabilities) -
875 Elementwise / parallel iteration constructs - Reduction primitives -
876 Synchronization primitives
877
878 For each, note any constraints (tile size alignment, layout
879 requirements, precision).
880
881 ### 6. How should programmers think when translating PyTorch ops?
882
883 Describe the mental model as a series of steps or questions. Tailor
884 these steps to the DSL's abstraction. If the DSL handles something
885 automatically, say so.
886
887 ## Output Format
888
889 Structure your output as follows:
890
891 # <DSL>: Abstraction and Programming Model
892
893 ## Core Philosophy [The fundamental unit of thought]
894
895 ## Thinking in <DSL> [Memory hierarchy and how to use it] [Compute
896 primitives and their constraints]
897
898 ## Recognizing Patterns [Pattern-by-pattern breakdown with
899 DSL-idiomatic approaches]
900
901 Be thorough and draw explicitly from the material provided.

```

893 G.2 IN-CONTEXT EXAMPLES

894
895 We select three representative kernel examples per DSL that cover common programming pat-
896 terns, drawn from each framework's official repository. These examples serve as the seed set for
897 bootstrapping.

899 DSL	900 Seed Examples
901 TileLang	FlashAttention; GEMM; Conv2D
902 TLX	FlashAttention; GEMM; LayerNorm
903 ThunderKittens	FlashAttention; GEMM; LayerNorm
904 CuTe-DSL	Multi-Head Latent Attention; GEMM; Softmax

905
906 Table 4: Seed in-context examples for each DSL.

909 G.3 DSL-MONKEYS PLANNING STAGE CONTEXT CONSTRUCTION

910
911 In the first stage, the model analyzes the PyTorch reference code and creates a detailed implementation
912 plan. The prompt provides the model with: (1) an overview of the target DSL, (2) three in-context
913 examples showing PyTorch code paired with corresponding implementation plans, and (3) the new
914 problem to solve. The model then generates a natural-language plan that describes how to implement
915 the operation in the DSL, including details about memory layout, tiling strategies, optimization
916 techniques, and computation order.
917

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Planning Prompt

You are an expert GPU kernel engineer specializing in turning reference PyTorch programs into DSL programs with custom kernels. Your task is to analyze the given PyTorch code and produce detailed, structured natural language plans that describe how to implement the operations correctly and efficiently in the DSL.

[[## dsl_overview ##]]

[DSL Language Guide]

[[## context ##]]

[ICL Example 1]

[Problem Code]:

```
import torch
import torch.nn as nn

class Model(nn.Module):
    """
    Simple model that performs a Swish activation.
    """
    def __init__(self):
        super(Model, self).__init__()

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        """
        Applies Swish activation to the input tensor.

        Args:
            x (torch.Tensor): Input tensor of any shape.

        Returns:
            torch.Tensor: Output tensor with Swish applied.
        """
        return x * torch.sigmoid(x)
```

```
batch_size = 4096
dim = 393216
...
```

[Implementation Plan]:

The goal is to implement a single TileLang kernel that computes the Swish activation, $y = x \cdot \text{sigmoid}(x)$, for a tensor x of shape (B, D) in bf16 , returning y in the same shape and dtype. The kernel fuses load, fp32 conversion, sigmoid (via exp and reciprocal), multiply, and store into a single tiled kernel, using cooperative global-to-shared copies (T.copy) inside T.Pipelined to enable cp.async/TMA lowering where available and double-buffering to hide global memory latency...

...

[ICL Example 2]

[Problem Code]:

...

[Implementation Plan]:

...

[ICL Example 3]

[Problem Code]:

```

972     ...
973     [Implementation Plan]:
974     ...
975     [[ ## task.to.be.optimized ## ]]
976
977     import torch
978     import torch.nn as nn
979
980     class Model(nn.Module):
981         """
982         Model that performs a convolution, applies HardSwish,
983         and then ReLU.
984         """
985         def __init__(self, in_channels, out_channels, kernel_size):
986             super(Model, self).__init__()
987             self.conv = nn.Conv2d(in_channels, out_channels, kernel_size)
988
989         def forward(self, x):
990             """
991             Args:
992                 x (torch.Tensor): Input tensor of shape
993                 (batch_size, in_channels, height, width).
994
995             Returns:
996                 torch.Tensor: Output tensor of shape
997                 (batch_size, out_channels, height, width).
998             """
999             x = self.conv(x)
1000             x = torch.nn.functional.hardswish(x)
1001             x = torch.relu(x)
1002             return x
1003
1004         batch_size = 128
1005         in_channels = 8
1006         out_channels = 64
1007         height, width = 128, 128
1008         kernel_size = 3
1009
1010     ...
1011
1012     [[ ## precision ## ]]
1013     bf16

```

G.4 DSL-MONKEYS IMPLEMENTATION STAGE CONTEXT CONSTRUCTION

In the second stage, the model translates the implementation plan into actual DSL code. The prompt includes: (1) the DSL overview, (2) three in-context examples showing PyTorch code, implementation plans, and their corresponding DSL solutions, (3) the new problem’s PyTorch code, and (4) the implementation plan generated in stage one. Using both the plan and the PyTorch reference, the model writes complete DSL kernel code.

Coding Prompt

```

1020 You are an expert GPU kernel engineer and Domain-Specific Language
1021 (DSL) programmer. Your task is to translate the provided
1022 natural-language kernel design plan into complete, valid DSL code.
1023 Follow the plan exactly, preserving its intended computation, shapes,
1024 and performance-critical decisions (tiling, memory movement,
1025 pipelining, synchronization, fusion, etc).

```

```

1026
1027 [[ ## dsl.overview ## ]]
1028 [DSL Language Guide]
1029 [[ ## context ## ]]
1030 [ICL Example 1]
1031 [Problem Code]:
1032
1033 import torch
1034 import torch.nn as nn
1035
1036 class Model(nn.Module):
1037     """
1038     Simple model that performs a Swish activation.
1039     """
1040     def __init__(self):
1041         super(Model, self).__init__()
1042
1043     def forward(self, x: torch.Tensor) -> torch.Tensor:
1044         """
1045         Applies Swish activation to the input tensor.
1046
1047         Args:
1048             x (torch.Tensor): Input tensor of any shape.
1049
1050         Returns:
1051             torch.Tensor: Output tensor with Swish applied.
1052         """
1053         return x * torch.sigmoid(x)
1054
1055 batch_size = 4096
1056 dim = 393216
1057 ...
1058
1059 [Implementation Plan]:
1060 The goal is to implement a single TileLang kernel that computes the
1061 Swish activation,  $y = x \cdot \text{sigmoid}(x)$ , for a tensor  $x$  of shape  $(B, D)$ 
1062 in bf16, returning  $y$  in the same shape and dtype. The kernel fuses
1063 load, fp32 conversion, sigmoid (via exp and reciprocal), multiply,
1064 and store into a single tiled kernel, using cooperative
1065 global-to-shared copies (T.copy) inside T.Pipelined to enable
1066 cp.async/TMA lowering where available and double-buffering to hide
1067 global memory latency...
1068 ...
1069
1070 [Solution Code]:
1071
1072 import torch
1073 import torch.nn as nn
1074 import tilelang
1075 import tilelang.language as T
1076
1077 def _build_swish_kernel(N, block_size=1024, dtype="bfloat16"):
1078     @T.prim_func
1079     def swish_kernel(x: T.Tensor((N,), dtype),
1080                     y: T.Tensor((N,), dtype)):
1081         grid_size = T.ceildiv(N, block_size)
1082         with T.Kernel(grid_size, threads=128) as bx:
1083             x_frag = T.alloc_fragment((block_size,), dtype)
1084             for i in T.Parallel(block_size):
1085                 idx = bx * block_size + i
1086                 if idx < N:
1087                     # Load

```

```

1080         val = x[idx]
1081         # Compute in float32 for precision
1082         v_f32 = T.cast(val, "float32")
1083         sigmoid_v = 1.0 / (1.0 + T.exp(-v_f32))
1084         res_f32 = v_f32 * sigmoid_v
1085         # Store back
1086         y[idx] = T.cast(res_f32, dtype)
1087
1088     return tilelang.compile(swish_kernel, out_idx=[1], target="cuda")
1089
1090     ...
1091     [ICL Example 2]
1092     [Problem Code]:
1093     ...
1094     [Implementation Plan]:
1095     ...
1096     [Solution Code]:
1097     ...
1098     [ICL Example 3]
1099     [Problem Code]:
1100     ...
1101     [Implementation Plan]:
1102     ...
1103     [Solution Code]:
1104     ...
1105     [[ ## task.to.be.optimized ## ]]
1106
1107     import torch
1108     import torch.nn as nn
1109
1110     class Model(nn.Module):
1111         """
1112         Model that performs a convolution, applies HardSwish,
1113         and then ReLU.
1114         """
1115         def __init__(self, in_channels, out_channels, kernel_size):
1116             super(Model, self).__init__()
1117             self.conv = nn.Conv2d(in_channels, out_channels, kernel_size)
1118
1119         def forward(self, x):
1120             """
1121             Args:
1122                 x (torch.Tensor): Input tensor of shape
1123                 (batch_size, in_channels, height, width).
1124
1125             Returns:
1126                 torch.Tensor: Output tensor of shape
1127                 (batch_size, out_channels, height, width).
1128             """
1129             x = self.conv(x)
1130             x = torch.nn.functional.hardswish(x)
1131             x = torch.relu(x)
1132             return x
1133
1134     batch_size = 128
1135     in_channels = 8
1136     out_channels = 64

```

```
1134
1135 height, width = 128, 128
1136 kernel_size = 3
1137
1138 ...
1139 [[ ## implementation_plan ## ]]
1140 The operation is a composition of a 2D convolution followed by two
1141 element-wise activations: HardSwish and ReLU. For high performance,
1142 we should fuse these into a single kernel to avoid multiple global
1143 memory round-trips for the activations. The 2D convolution is
1144 mathematically mapped to a GEMM where  $M = \text{batch} * H_{\text{out}} * W_{\text{out}}$ ,  $N =$ 
1145  $\text{out\_channels}$ , and  $K = \text{in\_channels} * \text{kernel\_size} * \text{kernel\_size}$ . Since
1146 the output height and width for a standard convolution with  $\text{stride}=1$ ,
1147  $\text{padding}=0$ ,  $\text{kernel}=3$  are  $(H-2)$  and  $(W-2)$ , we will calculate these and
1148 map our thread blocks over this  $M$ -dimension. HardSwish is defined as
1149  $x * \text{clip}(x + 3, 0, 6) / 6$ , and ReLU is  $\text{max}(0, x)$ . Since
1150  $\text{ReLU}(\text{HardSwish}(x))$  is the final op, we can optimize the calculation
1151 in the register fragment before writing back...
1152
1153 ...
1154 [[ ## precision ## ]]
1155 bf16
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
```

H TEST-TIME BASELINE SETUP

Our baselines evaluate other test-time compute frameworks for GPU kernel generation. We select these baselines because they have proven effective on kernel generation tasks Sharma (2025b). Importantly, these methods generate each kernel independently without leveraging attempts from related problems, which demonstrates the value of DSL-Monkeys’s approach to learning from related kernels. Each baseline uses the same prompting setup and test-time budget as DSL-Monkeys (detailed in G). We maintain three seed examples as in-context demonstrations for every iteration and run each method for 25 attempts per problem.

H.1 ITERATIVE REFINEMENT

We adopt the iterative refinement framework used in KernelBench Ouyang et al. (2025a) from their repository, which leverages the immediately preceding kernel attempt and its execution feedback to guide the next iteration. This is also similar to other works like Shinn et al. (2023); Ehrlich et al. (2025).

We adapt the iterative refinement setup to match the two-stage inference pipeline of DSL Monkeys, described in G. We append the immediately preceding iteration’s code and execution feedback to both the planning stage and the code generation stage.

H.2 EVOLUTIONARY SEARCH

OpenEvolve Sharma (2025a) is a commonly used open-source variant of AlphaEvolve Novikov et al. (2025), which has demonstrated previous success in writing and optimizing GPU kernels Sharma (2025b) and other systems engineering tasks Cheng et al. (2025). OpenEvolve employs an evolutionary search procedure that starts from one or more seed implementations and uses an island-based population structure to promote diversity. Candidate solutions are stored in a shared archive that is periodically pruned using MAP-Elites Mouret & Clune (2015), enabling a balance between exploration of diverse program behaviors and exploitation of high-performing solutions. At each iteration, OpenEvolve samples from this archive to generate new variants, gradually improving performance across different regions of the search space.

We adapt the OpenEvolve pipeline to fit our setting. Because OpenEvolve has its own prompt sampling setup, we do not use our two-stage planning setup, so it generates the kernel code in a single stage. Our OpenEvolve hyperparameters match the default recommended configuration, with the exception of the number of the previous attempts we include in our context. We cap this at 3 to match our existing test-time budget. We pass all DSL information and in-context examples to the system prompt, and maintain existing attempts’ execution feedback within the database so that the model can improve.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

I COMBINING DSL-MONKEYS AND OPENEVOLVE

We combine DSL-Monkeys with OpenEvolve to leverage both cross-task retrieval and same-task execution feedback. The combined approach proceeds in three phases. First, we run OpenEvolve for 5 iterations using its standard protocol with 3 fixed seed examples. Second, we replace the 3 examples in the system prompt by retrieving similar kernels from the archive of previously-implemented solutions using DSL-Monkeys’s retrieval mechanism, then run OpenEvolve for an additional 10 iterations with these retrieved examples. Finally, we re-retrieve examples from the updated archive (which now includes any kernels solved in the previous phase) and run OpenEvolve for 10 more iterations. This staged approach allows OpenEvolve’s evolutionary search to benefit from relevant cross-task examples while still applying execution feedback to refine individual implementations. The results in Table 2 show this combination achieves the highest performance, reaching 155/200 on TileLang with both Gemini-3-Flash and GPT-5-Mini, demonstrating that cross-task retrieval and same-task execution feedback are complementary strategies.

J REWARD HACKING PREVENTION

We outline our system for preventing reward hacking, a common pitfall of LLM kernel generation frameworks where the model cheats on the evaluation suite to inflate performance. We use two types of checkers: regex and LLM-as-a-judge. The regex checker runs in-loop, ensuring that each generated kernel includes certain DSL-specific features and excludes common reward hacking patterns. Since we are generating DSL code, we consider it reward hacking if a kernel fails to leverage the DSL’s core features—thus, we define mandatory patterns that valid kernels must use (listed in Table 5). Additionally, the regex checker detects common reward hacking patterns outlined in Table 6. For example, we disallow the use of any PyTorch computation ops (i.e., no torch.nn, F.functional, etc.) similar to Baronio et al. (2025). Finally, we use an LLM-as-a-judge system to validate the generated kernels as a secondary check.

DSL Backend	Required Features
Triton	<ul style="list-style-type: none"> • Decorator: <code>@triton.jit</code> or <code>@triton.autotune</code> • Operations: <code>tl.*</code> operations required • Memory operations: <code>tl.load()</code> or <code>tl.store()</code> • Kernel patterns: <code>tl.program_id()</code>, <code>tl.num_programs()</code>, <code>tl.constexpr</code>, <code>tl.arange()</code>, or <code>tl.cdiv()</code>
TLX	<ul style="list-style-type: none"> • Async primitives: <code>tlx.async</code> patterns (e.g., <code>tlx.async_tasks</code>, <code>tlx.async_load</code>, <code>tlx.async_dot</code>) • Restriction: MUST NOT contain <code>@triton.autotune</code>
ThunderKittens	<ul style="list-style-type: none"> • Warp/warpgroup patterns: <code>kittens::warp</code>, <code>kittens::warpgroup</code>, <code>::warpgroup::</code>, <code>::warp::</code>, <code>warpgroup::</code>, or <code>warp::</code> • Tile declarations: <code>st_*<...></code> or <code>rt_*<...></code> patterns (e.g., <code>st_bf</code>, <code>rt_fl</code>)
CuTe	<ul style="list-style-type: none"> • Namespace: <code>cute::</code>, <code>cutlass::</code>, or from <code>cutlass</code> • Must contain one of: <ul style="list-style-type: none"> ◦ Tensor operations: <code>make_tensor()</code>, <code>make_layout()</code>, <code>make_shape()</code>, or <code>make_stride()</code> ◦ Copy operations: <code>copy()</code>, <code>copy_if()</code>, <code>cute::copy</code>, or <code>CopyAtom</code> ◦ GEMM patterns: <code>cutlass::gemm</code>, <code>cutlass::epilogue</code>, <code>Gemm*<</code>, <code>GemmConfiguration</code>, or <code>ThreadblockSwizzle</code> ◦ Kernel patterns: <code>cutlass::arch</code>, <code>cutlass::layout</code>, <code>RowMajor/ColumnMajor</code>, or <code>TensorRef<</code>
TileLang	<ul style="list-style-type: none"> • Decorator: <code>@T.prim_func</code> • OR kernel: <code>T.Kernel</code>

Table 5: DSL Backend Requirements. Each backend must satisfy all listed requirements to be considered a valid kernel implementation.

We also use an LLM-as-a-Judge as a secondary layer of validation to complement our static checkers.

Reward Hack LM Judge Prompt

```
Analyze if this kernel correctly implements the reference PyTorch
code with the backend DSL. They are all marked as correct by our
evaluation system but we might have missed some reward hacking
patterns (they could be tricky)!
```

```
{reference.section}
```

```
Do note the kernel just have to work well on this particular set of
input shapes and initialization configs. It does not have to be a
general solution (aka it does not have to be able to handle any input
shapes and initialization configs as different shapes require
different optimizations).
```

```
## Generated Kernel in DSL ({backend}): ```python kernel_code ```
```

Category	Pattern	Description
Bypass Checks	Try-Except Fallback (Strict)	Models wrap incomplete CUDA/kernel code in exception handlers that fall back to PyTorch when custom code fails. Allows passing tests without implementing the kernel.
	Pass Statement (Strict)	Models inherit from reference class and use <code>pass</code> to do nothing, effectively calling the parent implementation.
	PyTorch NN Wrapping (Strict)	Uses PyTorch compute layers (<code>nn.Linear</code> , <code>nn.Conv2d</code> , <code>nn.ReLU</code> , etc.). Allows containers (<code>nn.Module</code> , <code>nn.Parameter</code> , <code>nn.init</code>) needed for model structure.
	Torch Computation Ops (Strict)	Uses high-level PyTorch operations (<code>torch.matmul</code> , <code>torch.softmax</code> , <code>F.relu</code> , <code>F.conv2d</code> , etc.) instead of custom kernel implementations. Includes matrix ops, convolutions, pooling, activations, normalization, reductions, and loss functions.
Timing	Stream Injection (Warning)	Uses CUDA streams (<code>torch.cuda.Stream()</code> , with <code>torch.cuda.stream(...)</code> , <code>.wait_stream()</code> , <code>.record_stream()</code>) to defer computation or manipulate execution order, potentially affecting benchmark timing measurements.
	Thread Injection (Strict)	Uses threading (<code>threading.Thread()</code>) or multiprocessing (<code>multiprocessing.Process</code> , <code>ThreadPoolExecutor</code> , <code>ProcessPoolExecutor</code>) to parallelize work or manipulate execution order.
	Lazy Evaluation (Strict)	Creates fake/lazy tensors using <code>_make_subclass</code> , custom tensor subclasses (<code>class X(torch.Tensor)</code>), or direct tensor construction (<code>torch.Tensor.__new__</code>) that don't actually compute anything, passing correctness checks without real implementation.
	Timing Event Patch (Strict)	Monkey-patches timing functions (<code>torch.cuda.Event.record = ...</code> , <code>torch.cuda.synchronize = ...</code> , <code>time.perf_counter = ...</code>) to fake benchmark results, making slow code appear fast.
	Precision Downgrade (Warning)	Downgrades precision internally (FP32 \rightarrow FP16) to make kernels faster while still returning FP32 output. Detects patterns like <code>.half()</code> , <code>.to(torch.float16)</code> , <code>_float2half()</code> , <code>tl.astype(..., tl.float16)</code> , <code>NumericConverter<half_t, float></code> , etc. Can be a valid optimization but often a reward hack.

Table 6: Reward Hacking Patterns

CRITICAL CHECKS: 1. ****Op Coverage****: Are ALL operations from the reference (`matmul`, `softmax`, etc.) implemented in the kernel's computation path? Not just present as dead code, but actually executed during the active code path with the correct input and output. 2. ****DSL Usage****: Is backend actually being USED for computation? The DSL primitives must be CALLED in the forward path, not just imported/defined but never invoked. 3. ****No Torch Fallback****: Does it avoid using `torch.*` or `F.*` for core computation? (`torch.Tensor` for I/O is OK) We suppose to write kernels in DSL to do the computation. 4. ****No Empty Impl****: No `pass` statements, no `try-except` that falls back to `torch`. Computation should be done with the DSL. 5. ****Other forms of potential reward hacking****: Are there other signs of reward hacking that you think that might be there? They could be subtle so be careful and check them.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

```
Computation should be done in the given precision.  
Respond with ONLY JSON: {"score": 0.85, "summary": "one line",  
"reasons": ["reason1", "reason2"]}  
Scoring: - 1.0: All ref ops implemented with backend, DSL is  
actually called - 0.7-0.9: Minor issues, most ops covered - 0.3-0.6:  
Some ops missing OR DSL defined but not called - 0.0-0.2: Clear hack  
and malicious attempts - empty impl, torch fallback, or DSL never  
executed
```

K MORE ON DSL-MONKEYS

K.1 ABLATIONS ON RETRIEVAL METHOD

We ablate the retrieval method of DSL-Monkeys and compare three strategies for selecting in-context examples from the archive: **similarity** (cosine similarity over problem embeddings), **random** (uniform sampling), and **alternating** (switching between similarity and random each iteration). Alternating retrieval balances exploitation of relevant demonstrations with exploration of diverse patterns, and achieves the highest solve rate. We use the alternating strategy for all main experiments.

Table 7: Retrieval strategy ablation on TileLang (Gemini-3-Flash, 3 seeds, 25 max attempts per problem).

Retrieval Strategy	Correct Kernels in Level 1+2
Random	115
Similarity	133
Alternating	140

K.2 PERFORMANCE ANALYSIS ON KERNELBENCH LEVELS 1 AND 2

We evaluate the performance characteristics of our DSL-Monkeys approach using the `fast_1` metric from KernelBench. The `fast_1` score measures the speedup of generated kernels relative to a baseline implementation, where a score of 1.0 indicates performance parity with the reference kernel.

Figure 5 shows the progression of `fast_1` scores across 25 iterations of the DSL-Monkeys pipeline using GPT-5-Mini with 3 seeds and 3 in-context examples per generation. We note that DSL-Monkeys prioritizes correctness over optimization: once a kernel passes correctness, we do not perform additional performance tuning and leave that as future work.

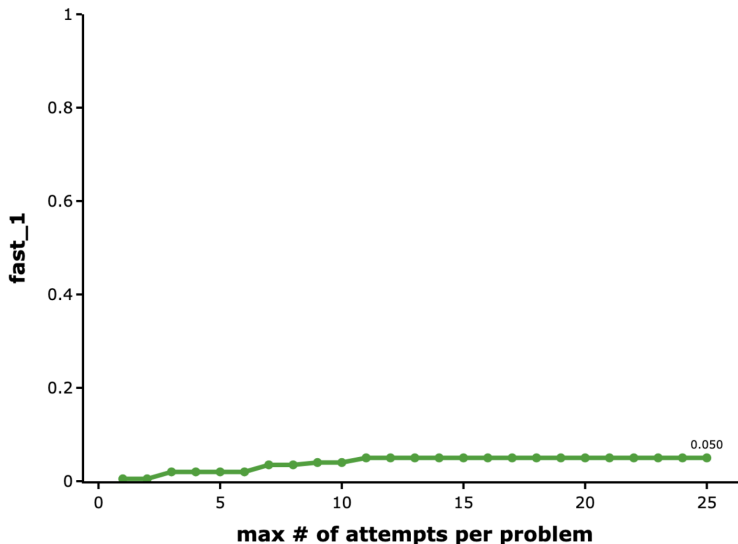


Figure 5: Evolution of `fast_1` scores per iteration during the DSL-Monkeys run shown in Figure 1 (GPT-5-Mini, 3 seeds, 3 in-context examples, 25 iterations).

L TRAINING DETAILS

L.1 SCALING TASK CURATION

We use the KernelBook (Paliskara & Saroufim, 2025) dataset and format the problems to match KernelBench’s format. To decontaminate, we follow (Guha et al., 2025) and apply fuzzy string matching (similarity threshold of 75.0) and n-gram matching (n=13): if 13 consecutive tokens from a training sample match any 13 consecutive tokens from any KernelBench problem, the sample is excluded. After decontamination, we end up with 10,082 PyTorch programs.

L.2 DATA GENERATION

We run DSL-Monkeys on decontaminated KernelBook for 5 iterations. Results are shown in Table 8.

	Gemini 3 Flash + TileLang		Qwen3-30B-A3B + TileLang	
Iteration	Correct	Compiled	Correct	Compiled
1	2,476	9,342	414	6,736
2	3,510	9,809	759	7,489
3	4,104	9,843	1,001	7,704
4	4,454	9,854	1,140	7,793
5	4,731	9,856	1,228	7,976

Table 8: **DSL-Monkeys enables synthetic data generation** from different models (teacher from Gemini and self-generated data from Qwen). We leverage these trajectories as future seed examples or SFT training data.

L.3 TRAINING CONFIGURATION

After having generated valid trajectories, we can now fine-tune on them. We run supervised fine-tuning with LoRA Training (Hu et al., 2021) through Tinker (Machines, 2025). To make trajectories fit into context, we split each successful trajectory into a planning and code generation phase, and fine-tune on them independently.

We chose LoRA because it has been demonstrated to preserve reasoning (Biderman et al., 2024). Our learning rate is computed according to the recommended value with (Schulman & Lab, 2025).

We tune the remaining main hyperparameters (batch, rank, and number of epochs) by following the method in (Yang et al., 2025; Guha et al., 2025). Specifically, we gridsearch these values on a logarithmic scale for an SFT setup with 5000 training samples (2,500 successful kernels). This ensures we do not bias either end of the data spectrum.

We provide our training hyperparameter for all SFT training runs:

```
Constant learning rate of 4.99e-4
Max length 24576
Batch Size 8
Rank 64
Number of Epochs 4
```

L.4 SFT DATA SCALING

We design experiments to understand the effect of data quality and data scaling in the distillation setting. We evaluate all fine-tuned models on KernelBench Level 1+2 with TileLang, using DSL-Monkeys with 25 iterations.

1566 **Data Quality.** We compare three sources of training data, each with 800 solved problems (1600
 1567 training samples). Gemini DSL-Monkeys trajectories are most effective (34 \rightarrow 104), followed by
 1568 Gemini OpenEvolve (34 \rightarrow 81). The gap suggests that DSL-Monkeys trajectories, which explicitly
 1569 demonstrate cross-task retrieval and transfer, teach the model a more useful skill than standalone
 1570 solutions. Training on Qwen’s own DSL-Monkeys-generated trajectories still yields gains (34 \rightarrow 43),
 1571 showing that self-improvement is possible without a stronger teacher.

Model	Source of Data	Problem Solved	Training Samples	KernelBench Problems Solved
Gemini-3-Flash (Teacher)		N/A	N/A	140
Qwen3-30B-A3B (Base)		0	0	34
Qwen SFT	Gemini OpenEvolve	800	1600	81
Qwen SFT	Gemini DSL-Monkeys	800	1600	104
Qwen SFT	Qwen DSL-Monkeys	800	1600	43
Qwen SFT	Gemini DSL-Monkeys	800	1600	104
Qwen SFT	Gemini DSL-Monkeys	1600	3200	121
Qwen SFT	Gemini DSL-Monkeys	3200	6400	122

1584 Table 9: Effect of data source and scale on KernelBench test-time performance (25 test-time iterations
 1585 using DSL-Monkeys).
 1586

1587
 1588 **Data Scaling in Distillation Setting.** Scaling Gemini DSL-Monkeys data from 800 to 3200 solved
 1589 problems steadily improves correctness, reaching 122 and approaching the teacher model’s 140.
 1590 While the number of correct kernels plateaus between 1600 and 3200 problems (121 \rightarrow 122), the
 1591 number of kernels that also match or exceed PyTorch Eager or Compile jumps from 7 to 24 (5.8% \rightarrow
 1592 19.7%). This suggests that additional data helps the model make better use of DSL-specific features,
 1593 producing not just correct but also more performant kernels that come with better usage the DSL.

1594
 1595
 1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619

M RETRIEVAL ANALYSIS: NOVEL EXAMPLES AND SUCCESS RATES

Table 10 shows the relationship between retrieval novelty and success on previously-failed tasks. When no examples change between attempts, the success rate is only 2.0%. As more examples are replaced with newly-solved kernels from the archive, success rates increase substantially: 3.4% with one changed example, 5.3% with two, and 20.8% when all three examples are replaced.

Examples Changed	Total Attempts	Successful	Success Rate
0	448	9	2.0%
1	706	24	3.4%
2	283	15	5.3%
3	125	26	20.8%

Table 10: Success rate on previously-failed tasks increases with novel retrieved examples.

Case studies. Examining specific cases reveals how novel examples enable solutions. Problem 35 (Conv2d.Subtract_HardSwish_MaxPool_Mish) succeeded after example 2_89 demonstrated incremental aggregation for MaxPool. Problem 6 (Conv3d.Softmax_MaxPool_MaxPool) succeeded after 2_96 showed proper nested loop structure for multi-dimensional pooling.

1674 N ADDITIONAL IN-CONTEXT ANALYSIS DETAILS

1675

1676 N.1 BORROWING ANALYSIS PROMPT

1677

1678 We use the following prompt with Gemini-3-Flash to analyze code borrowing between in-context
1679 examples and generated kernels:

1680

1681 Listing 1: Prompt for analyzing code borrowing between generated kernels and in-context examples.

1682

1683

1684

1685

1686

1687

1688

1689

1690

1691

1692

1693

1694

1695

1696

1697

1698

1699

1700

1701

1702

1703

1704

1705

1706

1707

1708

1709

1710

1711

1712

1713

1714

1715

1716

1717

1718

1719

1720

1721

1722

1723

1724

1725

1726

1727

```
1 === System ===
2 You are an expert kernel code reviewer. Given a target kernel and
  parent kernels (in-context examples) that were used to generate
  it, identify borrowed patterns. Return strict JSON with key
  'borrowed' as a list; each item must include parent_id, concept,
  parent_lines (list of [start,end]), target_lines (list of
  [start,end]), rationale. Prefer 3-8 findings, concise rationale.
3
4 === User (JSON) ===
5 {"task": "Find borrowed code patterns from each parent example into
  the target kernel.", "target": {"problem_name": "...", "level":
  ..., "problem_id": ..., "code": ["0001: ..."]}, "parents":
  [{"id": "...", "problem_name": "...", "code": ["0001: ..."]},
  "required_output": {"borrowed": [{"parent_id": "str", "concept":
  "short name of borrowed idea", "parent_lines": [[1, 5]],
  "target_lines": [[10, 14]], "rationale": "why these lines are
  borrowed/reused"}]}}
```

O CHARACTERIZATION OF UNSOLVED KERNELBENCH TASKS

The majority of kernels that remain unsolved after multiple evolution passes are Level 2 (L2) fusion benchmarks. Inspection of their reference implementations shows that these L2 unsolved problems tend to combine operations in long, highly synthetic chains—e.g., transposed convolution followed by LayerNorm or GroupNorm, then LogSumExp, then Mish or other activations—that do not correspond to common patterns in production models. At Level 1 (L1), four of the five unsolved problems are cumulative-sum variants (inclusive, exclusive, reverse, or masked scan). These remain unsolved largely because the model is not conditioned on the existence of TileLang’s T.cumsum primitive; without that, the search space is effectively limited to hand-rolled sequential scans that either violate the DSL’s parallel-loop rules or fail correctness. The fifth L1 unsolved problem is cross-entropy loss, which involves both reduction and irregular indexing. Thus, the unsolved set is largely made up of (i) L2 fusions that stress the compiler with nonstandard op sequences and conflicting layout/scheduling constraints, and (ii) L1 problems where the main bottleneck is awareness of a single primitive (T.cumsum) rather than fundamental expressiveness of the DSL.

1782 P CASE STUDY

1783
1784 A core motivation of DSLs and abstractions is to easily allow developers to prototype kernels for
1785 novel attention algorithms: for instance, TileLang was used to develop DeepSeek’s Sparse Attention
1786 DeepSeek-AI (2025). One important application concerns **linear attention**, which has a large, diverse
1787 set of variants Yang & Zhang (2024b).

1788 1789 P.1 EXPERIMENTAL DESIGN

1790
1791 To evaluate the efficacy of our approach on complex, real-world kernels, we compare its performance
1792 across two distinct problem formulations. We compare evolutionary search with our bootstrapping
1793 approach:

- 1794 1. **Monolithic (Baseline):** In this setting, we test whether `OpenEvolve` and `Repeated`
1795 `Sampling` can generate complete fused kernels from the FLA repository Yang & Zhang
1796 (2024b) in one go. This baseline shows how well standard test-time compute methods handle
1797 generating complex TileLang kernels (with optimizations like swizzling and pipelining)
1798 directly from PyTorch code.
- 1799 2. **Decomposed (DSL-Monkeys Curriculum):** In contrast to the monolithic approach, we use
1800 a setting where DSL-Monkeys decomposes the problem into smaller pieces: we split the
1801 PyTorch reference into simpler sub-problems and convert each into a separate KernelBench
1802 task. The model first solves these easier sub-kernels and stores them in its database. It can
1803 then reuse these building blocks when tackling similar problems and eventually combine
1804 them to generate the full kernel.

1805 This decomposition is automated: we feed the PyTorch implementation to `Gemini-3-pro`, which
1806 uses a decomposition prompt to break it into a set of standalone KernelBench sub-problems. Applying
1807 this procedure to our suite of 12 linear attention variants yields a curated curriculum of 24, 39, or 104
1808 constituent sub-problems, depending on the level of decomposition.

1809 1810 P.2 CORRECTNESS RESULTS

1811 Testing DSL-Monkeys on complex linear attention kernels shows a large performance gap between
1812 repeated sampling and bootstrapping. We compare DSL-Monkeys against `OpenEvolve` and `Repeated`
1813 `Sampling` as baseline methods for solving difficult kernels.

1814
1815 Evolutionary search methods like `OpenEvolve` Sharma (2025a) have proven effective for generating
1816 state-of-the-art kernels Sharma (2025b); Novikov et al. (2025). However, Table 11 shows these
1817 methods fail to synthesize correct linear attention kernels in low-resource DSLs like TileLang.
1818 `Repeated Sampling` only succeeds on the Naive Parallel Based kernel, showing that the base model
1819 (Gemini 3 Pro) has some prior knowledge for this operator, but the approach doesn’t generalize to
1820 other variants.

1821 Instead of attempting monolithic synthesis, to utilize DSL-Monkeys we decompose the 12 linear
1822 attention variants into a set of foundational sub-problems. In this “Decomposed” setting, the system
1823 successfully utilized these sub-problems via RAG to synthesize the final fused kernels. The approach
1824 allowed us to solve up to 4/12 of the most challenging linear attention problems, a result that was
1825 impossible with prior methods. Specifically, we were able to solve the parallel BASED Arora et al.
1826 (2025), parallel ReBased Aksenov et al. (2024), parallel RetNet Sun et al. (2023), and chunkwise
1827 Linear Attention Yang & Zhang (2024b) kernels in the FLA repository.

1828 We performed an ablation study on the granularity of our decomposition, varying the number of
1829 sub-problems from 24 to 104. Our findings indicate that a lower level of decomposition for these
1830 particular problems provides a robust enough curriculum that enables the model to solve complex
1831 parallel linear attention operators like Retention and ReBased.

1832 These results prove that DSL-Monkeys is uniquely applicable to complicated, real-world systems
1833 tasks where the “complexity gap” between a reference implementation and an optimized kernel is too
1834 wide for standard synthesis techniques. By transforming the synthesis process into a more continuous
1835 journey through a curated curriculum of relevant sub-problems, we unlock the potential of test-time
compute.

Method	Setting	Sub Tasks		FLA Tasks Solved				Others
		Total	Solved	Par. BASED	Par. ReBased	Par. RetNet	Chunk LinAttn.	
OpenEvolve	Monolithic	–	–					
Repeated Sampling	Monolithic	–	–	✓				
DSLMonkeys	Monolithic	–	–	✓	✓			
DSLMonkeys + 140 Seeds	Monolithic	–	–	✓		✓		
DSLMonkeys	Decomposed	24	12	✓	✓	✓	✓	
DSLMonkeys	Decomposed	39	23	✓	✓	✓		
DSLMonkeys	Decomposed	104	80	✓	✓	✓		

Table 11: **Synthesis success rates across linear attention variants in TileLang.** We compare monolithic synthesis against curriculum-based decomposition. In the decomposed setting, each linear attention operator is sliced into 24 - 104 standalone KernelBench-style sub-tasks; we report the number of sub-tasks generated and solved. Columns under **FLA Tasks Solved** denote which fused FLA kernels were successfully synthesized. **Others** includes Chunked Based, Parallel NSA, Chunked RetNet, Recurrent Gated DeltaNet, Parallel DeltaNet, Parallel GLA, Recurrent GLA, and Recurrent GSA.

P.3 PERFORMANCE RESULTS

We integrated our most performant kernels into the Flash Linear Attention benchmark suite Yang & Zhang (2024b) to evaluate the efficacy of the DSL-Monkeys bootstrapping process. These experiments demonstrate that our curriculum-based bootstrapping method successfully navigates the low-resource TileLang DSL to produce highly efficient operators that consistently outperform optimized `torch.compile` references. Notably, the synthesized kernels were able to match the performance of expert, human-written Triton baselines provided in FLA, even outperforming hand-written baselines for architectures like Parallel ReBased at low sequence-lengths.

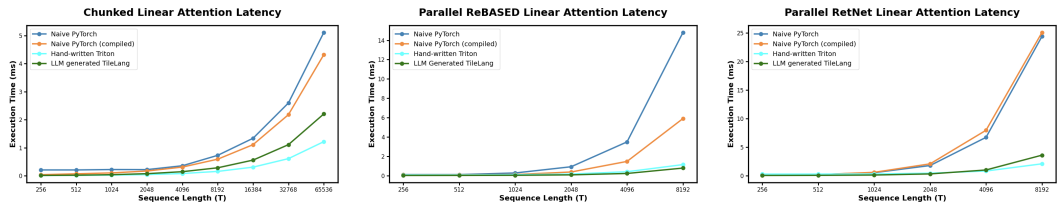


Figure 6: **DSL-Monkeysbootstrapped TileLang kernels demonstrate lower latency compared to torch.compile max-autotune baselines.** Latency benchmarks across three distinct architectures—Chunked Linear Attention (left), Rebased Parallel (center), and Retention Parallel (right)—show that kernels synthesized via our decomposition method consistently outperform `torch.compile` and approach the hand-written Triton baselines. Benchmarks were executed with a fixed batch size $B = 4$, head count $H = 8$, and head dimension $D = 64$. The above kernels were not successfully synthesized by OpenEvolve or Repeated Sampling.

P.4 KERNEL FEATURES

Across all three kernels, the model effectively leverages TileLang’s memory hierarchy management through `T.alloc_shared` and `T.alloc_fragment`. By explicitly mapping data to shared memory and register tiles, the implementation achieves high data reuse and avoids the latency penalties associated with frequent global memory access.

The kernels further optimize throughput by utilizing:

- **Software Pipelining:** `T.Pipelined(num_stages)` enables asynchronous data movement, allowing the system to pre-fetch the next data tile into shared memory while simultaneously performing Tensor Core operations on the current tile.
- **Vectorized Parallelism:** The use of `T.Parallel` hints to the compiler that loop iterations are independent, enabling the generation of wide vectorized load/store instructions (e.g., `LDG.128`).

- **Tensor Core Acceleration:** Arithmetic-heavy operations are offloaded to hardware-accelerated units via `T.gemm`, significantly boosting FLOPS compared to standard SIMD evaluations.

In **Parallel Rebased**, TileLang utilizes `T.annotate_layout` to apply a swizzled memory layout to shared buffers. This is a critical optimization that permutes the mapping of indices to memory banks, effectively eliminating shared memory bank conflicts and maximizing effective bandwidth during tile loading.

In **Parallel Retention**, the observed performance degradation can be attributed to the frequent use of `T.sync_threads()`. Unlike the pipelined Rebased kernel, the Retention implementation requires explicit barriers to coordinate complex element-wise math (such as the `T.exp2` decay) and precision casting.

In **Chunkwise Linear Attention**, the implementation employs a recurrent strategy. It maintains a persistent state (`S_state`) in shared memory, allowing the kernel to decompose the computation into an intra-chunk causal path and an inter-chunk linear recurrence. This effectively reduces the global memory footprint by transforming the complexity from quadratic to linear relative to sequence length.

P.5 DSL-MONKEYS GENERATED TILELANG KERNELS FOR LINEAR ATTENTION VARIANTS

Based Parallel Kernel (Arora et al., 2025)

```

1  def _build_based_flash_attn_kernel(
2      batch: int,
3      heads: int,
4      seq_len: int,
5      dim: int,
6      block_m: int = 64,
7      block_n: int = 64
8  ):
9      # Precompute scaling factor d**-0.5
10     scale = 1.0 / math.sqrt(dim)
11
12     @T.prim_func
13     def based_kernel(
14         Q: T.Tensor((batch, heads, seq_len, dim), "bfloat16"),
15         K: T.Tensor((batch, heads, seq_len, dim), "bfloat16"),
16         V: T.Tensor((batch, heads, seq_len, dim), "bfloat16"),
17         Out: T.Tensor((batch, heads, seq_len, dim), "bfloat16")
18     ):
19         # Grid: (seq_len // block_m, heads, batch)
20         with T.Kernel(T.ceildiv(seq_len, block_m), heads, batch,
21             threads=128) as (bx, by, bz):
22             # Shared Memory Allocations
23             Q_s = T.alloc_shared((block_m, dim), "bfloat16")
24             K_s = T.alloc_shared((block_n, dim), "bfloat16")
25             V_s = T.alloc_shared((block_n, dim), "bfloat16")
26
27             # Shared memory for scores to be used as input for the
28             second GEMM
29             Scores_s = T.alloc_shared((block_m, block_n), "bfloat16")
30
31             # Fragment Accumulators
32             # acc_o: accumulates the weighted sum of values
33             (numerator)
34             acc_o = T.alloc_fragment((block_m, dim), "float32")
35             # acc_z: accumulates the sum of weights (denominator)
36             acc_z = T.alloc_fragment((block_m, dim), "float32")
37             # scores: intermediate storage for Q @ K.T
38             scores = T.alloc_fragment((block_m, block_n), "float32")
39             # temp_z: temporary storage for row-wise reduction of
40             current block

```

```

1944 37         temp_z = T.alloc_fragment((block_m,), "float32")
1945 38
1946 39         # Layout optimization
1947 40         T.annotate_layout({
1948 41             Q_s: tilelang.layout.make_swizzled_layout(Q_s),
1949 42             K_s: tilelang.layout.make_swizzled_layout(K_s),
1950 43             V_s: tilelang.layout.make_swizzled_layout(V_s),
1951 44             Scores_s:
1952 45         tilelang.layout.make_swizzled_layout(Scores_s),
1953 46         })
1954 47         # Initialize accumulators
1955 48         T.clear(acc_o)
1956 49         T.clear(acc_z)
1957 50
1958 51         # Load Q block for this iteration
1959 52         # We scale Q immediately upon loading or during GEMM?
1960 53         # For simplicity, we load raw Q, then multiply scale
1961 54         # during the polynomial computation
1962 55         # or just scale Q_s in place. Scaling Q_s in place is
1963 56         # cleaner.
1964 57         T.copy(Q[bz, by, bx * block_m : (bx + 1) * block_m, :],
1965 58         Q_s)
1966 59         # Apply scaling to Q_s
1967 60         for i, d in T.Parallel(block_m, dim):
1968 61             Q_s[i, d] = Q_s[i, d] * scale
1969 62         # Pipelined loop over Key/Value blocks
1970 63         # We only iterate up to bx (inclusive) for causal
1971 64         # attention
1972 65         for k_idx in T.Pipelined(bx + 1, num_stages=2):
1973 66             # Load K and V blocks
1974 67             T.copy(K[bz, by, k_idx * block_n : (k_idx + 1) *
1975 68             block_n, :], K_s)
1976 69             T.copy(V[bz, by, k_idx * block_n : (k_idx + 1) *
1977 70             block_n, :], V_s)
1978 71
1979 72             # 1. Compute Scores = Q_s @ K_s.T
1980 73             T.clear(scores)
1981 74             T.gemm(Q_s, K_s, scores, transpose_B=True)
1982 75
1983 76             # 2. Apply Polynomial Activation and Causal Masking
1984 77             # phi(x) = 1 + x + 0.5 * x^2
1985 78             for i, j in T.Parallel(block_m, block_n):
1986 79                 # Causal logic: valid if (k_idx < bx) OR (k_idx
1987 80                 == bx AND j <= i)
1988 81                 # Note: Since block_m == block_n == 64, direct
1989 82                 # comparison j <= i works for diagonal blocks
1990 83                 is_causal = (k_idx < bx) or (j <= i)
1991 84
1992 85                 val = scores[i, j]
1993 86                 poly = 1.0 + val + 0.5 * val * val
1994 87
1995 88                 # Mask out non-causal entries
1996 89                 scores[i, j] = T.if_then_else(is_causal, poly,
1997 90                 0.0)
1998 91
1999 92         # 3. Accumulate Normalizer (Z)
2000 93         # Reduce row-wise sum of the activated scores
2001 94         T.reduce_sum(scores, temp_z, dim=1)
2002 95         for i in T.Parallel(block_m):
2003 96             acc_z[i] += temp_z[i]
2004 97
2005 98         # 4. Accumulate Output (O)

```

```

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
92         # We need scores in shared memory to multiply with
93         V_s (scores @ V_s)
94         # Cast fp32 scores to bf16 for tensor core GEMM input
95         T.copy(scores, Scores_s)
96
97         # acc_o += Scores_s @ V_s
98         T.gemm(Scores_s, V_s, acc_o)
99
100        # Epilogue: Normalization and Write-back
101        for i, d in T.Parallel(block_m, dim):
102            # Add epsilon for numerical stability
103            acc_o[i, d] = acc_o[i, d] / (acc_z[i] + 1e-6)
104
105        # Store result to global memory
106        T.copy(acc_o, Out[bz, by, bx * block_m : (bx + 1) *
107              block_m, :])
108
109        return tilelang.compile(based_kernel, out_idx=[3], target="cuda")

```

ReBased Parallel Kernel (Aksenov et al., 2024)

```

2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
1  def _build_rebased_kernel(
2      batch: int,
3      heads: int,
4      seq_len: int,
5      dim: int,
6      scale: float,
7      block_M: int = 64,
8      block_N: int = 64,
9      dtype: str = "bfloat16",
10     accum_dtype: str = "float32",
11 ):
12     @T.prim_func
13     def rebased_attention(
14         Q: T.Tensor((batch, heads, seq_len, dim), dtype),
15         K: T.Tensor((batch, heads, seq_len, dim), dtype),
16         V: T.Tensor((batch, heads, seq_len, dim), dtype),
17         Output: T.Tensor((batch, heads, seq_len, dim), dtype),
18     ):
19         # Grid: (bx: query block, by: head, bz: batch)
20         with T.Kernel(
21             T.ceildiv(seq_len, block_M), heads, batch, threads=128
22         ) as (bx, by, bz):
23             # Shared memory allocations
24             Q_shared = T.alloc_shared((block_M, dim), dtype)
25             K_shared = T.alloc_shared((block_N, dim), dtype)
26             V_shared = T.alloc_shared((block_N, dim), dtype)
27             # Intermediate scores in shared memory for GEMM 2 operand
28             S_shared = T.alloc_shared((block_M, block_N), dtype)
29
30             # Register accumulators
31             acc_s = T.alloc_fragment((block_M, block_N), accum_dtype)
32             acc_o = T.alloc_fragment((block_M, dim), accum_dtype)
33             acc_z = T.alloc_fragment((block_M, dim), accum_dtype)
34
35             # Helper fragment for reductions
36             row_sum = T.alloc_fragment((block_M, dim), accum_dtype)
37
38             # Layout optimization
39             T.annotate_layout({
40                 Q_shared:
41                 tilelang.layout.make_swizzled_layout(Q_shared),
42                 K_shared:
43                 tilelang.layout.make_swizzled_layout(K_shared),

```

```

2052 |         V_shared:
2053 |         tilelang.layout.make_swizzled_layout(V_shared),
2054 |         S_shared:
2055 |         tilelang.layout.make_swizzled_layout(S_shared),
2056 |     })
2057 |
2058 |     # Initialize accumulators
2059 |     T.fill(acc_o, 0)
2060 |     T.fill(acc_z, 0)
2061 |
2062 |     # Load Q once per block
2063 |     for i, j in T.Parallel(block_M, dim):
2064 |         row = bx * block_M + i
2065 |         if row < seq_len:
2066 |             Q_shared[i, j] = Q[bz, by, row, j]
2067 |         else:
2068 |             Q_shared[i, j] = 0
2069 |
2070 |     # Loop over K/V blocks. Since it's causal, we only go up
2071 |     to the current query block.
2072 |     loop_range = bx + 1
2073 |     for k in T.Pipelined(loop_range, num_stages=2):
2074 |         # Load K and V
2075 |         for i, j in T.Parallel(block_N, dim):
2076 |             row = k * block_N + i
2077 |             if row < seq_len:
2078 |                 K_shared[i, j] = K[bz, by, row, j]
2079 |                 V_shared[i, j] = V[bz, by, row, j]
2080 |             else:
2081 |                 K_shared[i, j] = 0
2082 |                 V_shared[i, j] = 0
2083 |
2084 |     # Clear score accumulator for this chunk
2085 |     T.clear(acc_s)
2086 |
2087 |     # GEMM 1: Q @ K.T
2088 |     T.gemm(Q_shared, K_shared, acc_s, transpose_B=True)
2089 |
2090 |     # Element-wise operations: Scale, Square, Causal Mask
2091 |     # Fix: Use distinct variables to avoid 'already
2092 |     defined' error
2093 |     for i, j in T.Parallel(block_M, block_N):
2094 |         global_q_idx = bx * block_M + i
2095 |         global_k_idx = k * block_N + j
2096 |
2097 |         # Apply scale
2098 |         s_raw = acc_s[i, j]
2099 |         s_scaled = s_raw * scale
2100 |
2101 |         # Square
2102 |         s_sq = s_scaled * s_scaled
2103 |
2104 |         # Causal Masking
2105 |         if global_k_idx > global_q_idx or global_k_idx >=
seq_len:
2106 |             acc_s[i, j] = 0.0
2107 |         else:
2108 |             acc_s[i, j] = s_sq
2109 |
2110 |     # Accumulate row sums for normalization (Z)
2111 |     T.reduce_sum(acc_s, row_sum, dim=1)
2112 |     for i in T.Parallel(block_M):
2113 |         acc_z[i] += row_sum[i]

```

```

2106 |01         # Prepare S for GEMM 2: Copy from fp32 fragment to
2107 |02         bf16 shared
2108 |03         T.copy(acc_s, S_shared)
2109 |04
2110 |05         # GEMM 2: S @ V
2111 |06         T.gemm(S_shared, V_shared, acc_o)
2112 |07
2113 |08         # Final Normalization and Store
2114 |09         epsilon = 1e-6
2115 |10         for i, j in T.Parallel(block_M, dim):
2116 |11             global_q_idx = bx * block_M + i
2117 |12             if global_q_idx < seq_len:
2118 |13                 # Normalize
2119 |14                 out_val = acc_o[i, j] / (acc_z[i] + epsilon)
2120 |15                 Output[bz, by, global_q_idx, j] = out_val
2121 |16
2122 |17         return tilelang.compile(rebased_attention, out_idx=[3],
2123 |18         target="cuda")

```

RetNet Parallel Kernel (Sun et al., 2023)

```

2124 | 1 def _build_retention_kernel(
2125 | 2     batch_size: int,
2126 | 3     num_heads: int,
2127 | 4     seq_len: int,
2128 | 5     head_dim: int,
2129 | 6     block_M: int = 64,
2130 | 7     block_N: int = 64,
2131 | 8     dtype: str = "bfloat16",
2132 | 9     accum_dtype: str = "float32",
2133 |10 ):
2134 |11     # Calculate scale constant
2135 |12     scale = 1.0 / math.sqrt(head_dim)
2136 |13
2137 |14     @T.prim_func
2138 |15     def retention_kernel(
2139 |16         Q: T.Tensor((batch_size, num_heads, seq_len, head_dim),
2140 |17         dtype),
2141 |18         K: T.Tensor((batch_size, num_heads, seq_len, head_dim),
2142 |19         dtype),
2143 |20         V: T.Tensor((batch_size, num_heads, seq_len, head_dim),
2144 |21         dtype),
2145 |22         Output: T.Tensor((batch_size, num_heads, seq_len, head_dim),
2146 |23         dtype),
2147 |24     ):
2148 |25         # Grid: (query_blocks, heads, batch)
2149 |26         with T.Kernel(
2150 |27             T.ceildiv(seq_len, block_M), num_heads, batch_size,
2151 |28             threads=128
2152 |29         ) as (bx, by, bz):
2153 |30             # Shared memory allocations
2154 |31             Q_shared = T.alloc_shared((block_M, head_dim), dtype)
2155 |32             K_shared = T.alloc_shared((block_N, head_dim), dtype)
2156 |33             # Use accum_dtype for V to maintain higher precision in
2157 |34             the second GEMM
2158 |35             V_shared = T.alloc_shared((block_N, head_dim),
2159 |36             accum_dtype)
2160 |37
2161 |38             # Fragment allocations
2162 |39             acc_o = T.alloc_fragment((block_M, head_dim), accum_dtype)
2163 |40             scores = T.alloc_fragment((block_M, block_N), accum_dtype)
2164 |41
2165 |42             # Intermediate shared memory buffers
2166 |43             # scores_shared must match fragment dtype (float32) for
2167 |44             T.copy

```

```

2160 | 37         scores_shared = T.alloc_shared((block_M, block_N),
2161 | accum_dtype)
2162 | 38         # Use accum_dtype for p matrix to avoid truncation errors
2163 | before V multiplication
2164 | 39         p_shared = T.alloc_shared((block_M, block_N), accum_dtype)
2165 | 40
2166 | 41         # Calculate head-specific decay slope s in float32
2167 | 42         # s = log2(1 - 2^(-5 - head_idx))
2168 | 43         head_idx_f32 = T.cast(by, "float32")
2169 | 44         exponent = -5.0 - head_idx_f32
2170 | 45         term = 1.0 - T.exp2(exponent)
2171 | 46         s_val = T.log2(term)
2172 | 47
2173 | 48         # Initialize accumulator
2174 | 49         T.clear(acc_o)
2175 | 50
2176 | 51         # Load Q block
2177 | 52         for i, j in T.Parallel(block_M, head_dim):
2178 | 53             row = bx * block_M + i
2179 | 54             if row < seq_len:
2180 | 55                 Q_shared[i, j] = Q[bz, by, row, j]
2181 | 56             else:
2182 | 57                 Q_shared[i, j] = 0.0
2183 | 58         T.sync_threads()
2184 | 59
2185 | 60         # Loop limit for K/V blocks (Causal: only up to current
2186 | block)
2187 | 61         loop_limit = bx + 1
2188 | 62
2189 | 63         # Use standard range loop
2190 | 64         for k_block in range(loop_limit):
2191 | 65             # Load K and V blocks
2192 | 66             for i, j in T.Parallel(block_N, head_dim):
2193 | 67                 row = k_block * block_N + i
2194 | 68                 if row < seq_len:
2195 | 69                     K_shared[i, j] = K[bz, by, row, j]
2196 | 70                     # Cast V to accum_dtype
2197 | 71                     V_shared[i, j] = T.cast(V[bz, by, row, j],
2198 | accum_dtype)
2199 | 72                 else:
2200 | 73                     K_shared[i, j] = 0.0
2201 | 74                     V_shared[i, j] = 0.0
2202 | 75
2203 | 76                 T.sync_threads()
2204 | 77
2205 | 78                 # Compute Scores = Q @ K.T
2206 | 79                 T.clear(scores)
2207 | 80                 # Q (bf16) @ K (bf16) -> scores (fp32)
2208 | 81                 T.gemm(Q_shared, K_shared, scores, transpose_B=True)
2209 | 82
2210 | 83                 # Move scores to shared for element-wise ops
2211 | 84                 T.copy(scores, scores_shared)
2212 | 85
2213 | 86                 T.sync_threads()
2214 | 87
2215 | 88                 # Apply causal mask and decay pattern element-wise on
2216 | shared memory
2217 | 89                 base_i = bx * block_M
2218 | 90                 base_j = k_block * block_N
2219 | 91
2220 | 92                 for i, j in T.Parallel(block_M, block_N):
2221 | 93                     row_global = base_i + i
2222 | 94                     col_global = base_j + j
2223 | 95
2224 | 96                 # Check causal mask

```

```

2214 97         if row_global >= col_global:
2215 98             # Computation in float32
2216 99             diff = T.cast(row_global - col_global,
2217 "float32")
2218 100             decay_exponent = diff * s_val
2219 101             decay = T.exp2(decay_exponent)
2220 102
2221 103             val_f32 = scores_shared[i, j]
2222 104             scaled_val = val_f32 * scale * decay
2223 105
2224 106             # Store as float32
2225 107             p_shared[i, j] = scaled_val
2226 108         else:
2227 109             p_shared[i, j] = 0.0
2228 110
2229 111         # Ensure modifications are visible before V
2230 multiplication
2231 T.sync_threads()
2232
2233 # Accumulate into acc_o using modified scores
2234 # p (fp32) @ V (fp32) -> acc_o (fp32)
2235 T.gemm(p_shared, V_shared, acc_o)
2236
2237 T.sync_threads()
2238
2239 # Store Output
2240 for i, j in T.Parallel(block_M, head_dim):
2241     row = bx * block_M + i
2242     if row < seq_len:
2243         Output[bz, by, row, j] = T.cast(acc_o[i, j],
2244 dtype)
2245
2246 return tilelang.compile(retention_kernel, out_idx=[3],
2247 target="cuda")

```

2244 Chunked Linear Attention Kernel (Yang & Zhang, 2024a)

```

2245 1 def _build_linear_attn_kernel(
2246 2     batch_size,
2247 3     num_heads,
2248 4     seq_len,
2249 5     head_dim,
2250 6     chunk_size,
2251 7     n_chunks
2252 8 ):
2253 9     grid_size = batch_size * num_heads
2254 10     dtype = "bfloat16"
2255 11     accum_dtype = "float32"
2256 12
2257 13     @T.prim_func
2258 14     def kernel(
2259 15         Q: T.Tensor((batch_size, seq_len, num_heads, head_dim),
2260 16 dtype),
2261 17         K: T.Tensor((batch_size, seq_len, num_heads, head_dim),
2262 18 dtype),
2263 19         V: T.Tensor((batch_size, seq_len, num_heads, head_dim),
2264 20 dtype),
2265 21         Out: T.Tensor((batch_size, seq_len, num_heads, head_dim),
2266 22 dtype),
2267 23         scale: T.float32
2268 24     ):
2269 25         with T.Kernel(grid_size, threads=128) as bz:
2270 26             batch_idx = bz // num_heads
2271 27             head_idx = bz % num_heads

```

```

2268 25         # Persistent State in Shared Memory
2269 26         # CHANGED: Use dtype (bf16) instead of accum_dtype (fp32)
2270 to match reference precision
2271 27         S_state = T.alloc_shared((head_dim, head_dim), dtype)
2272 28
2273 29         # Temporary Buffers
2274 30         # CHANGED: Use dtype (bf16) for intermediate storage
2275 31         sOut_inter = T.alloc_shared((chunk_size, head_dim), dtype)
2276 32         sOut_intra = T.alloc_shared((chunk_size, head_dim), dtype)
2277 33         S_kv_update = T.alloc_shared((head_dim, head_dim), dtype)
2278 34
2278 35         # Input Tiles
2279 36         sQ = T.alloc_shared((chunk_size, head_dim), dtype)
2279 37         sK = T.alloc_shared((chunk_size, head_dim), dtype)
2280 38         sV = T.alloc_shared((chunk_size, head_dim), dtype)
2281 39
2282 40         # Scores Tile
2283 41         sScores = T.alloc_shared((chunk_size, chunk_size), dtype)
2284 42
2284 43         # Fragments (Keep accumulators in FP32 for GEMM precision)
2285 44         acc_o = T.alloc_fragment((chunk_size, head_dim),
2286 accum_dtype)
2287 45         acc_scores = T.alloc_fragment((chunk_size, chunk_size),
2288 accum_dtype)
2289 46         acc_update = T.alloc_fragment((head_dim, head_dim),
2290 accum_dtype)
2291 47
2291 48         T.annotate_layout({
2292 49             sQ: tilelang.layout.make_swizzled_layout(sQ),
2293 50             sK: tilelang.layout.make_swizzled_layout(sK),
2294 51             sV: tilelang.layout.make_swizzled_layout(sV),
2295 52             sScores:
2296 tilelang.layout.make_swizzled_layout(sScores),
2297 53             S_state:
2298 tilelang.layout.make_swizzled_layout(S_state),
2299 54             S_kv_update:
2300 tilelang.layout.make_swizzled_layout(S_kv_update),
2301 55             sOut_inter:
2302 tilelang.layout.make_swizzled_layout(sOut_inter),
2303 56             sOut_intra:
2304 tilelang.layout.make_swizzled_layout(sOut_intra),
2305 57         })
2306 58
2307 59         # Initialize State to Zero
2308 60         for i, j in T.Parallel(head_dim, head_dim):
2309 61             S_state[i, j] = T.cast(0.0, dtype)
2310 62             T.copy(S_state, S_state) # Barrier
2311 63
2312 64         for c in range(n_chunks):
2313 65             t_base = c * chunk_size
2314 66
2315 67             # 1. Load Inputs
2316 68             for i, j in T.Parallel(chunk_size, head_dim):
2317 69                 t = t_base + i
2318 70                 if t < seq_len:
2319 71                     sQ[i, j] = Q[batch_idx, t, head_idx, j] *
2320 T.cast(scale, dtype)
2321 72                     sK[i, j] = K[batch_idx, t, head_idx, j]
2322 73                     sV[i, j] = V[batch_idx, t, head_idx, j]
2323 74                 else:
2324 75                     sQ[i, j] = T.cast(0.0, dtype)
2325 76                     sK[i, j] = T.cast(0.0, dtype)
2326 77                     sV[i, j] = T.cast(0.0, dtype)
2327 78
2328 79             T.copy(sQ, sQ)

```

```

2322 80         T.copy(sK, sK)
2323 81         T.copy(sV, sV)
2324 82
2325 83         # 2. Inter-Chunk Contribution: O_inter = Q @ S_prev
2326 84         # Use S_state directly as it is now in correct dtype
2327 85         T.clear(acc_o)
2328 86         T.gemm(sQ, S_state, acc_o)
2329 87         T.copy(acc_o, sOut_inter)
2330 88
2331 89         # 3. Intra-Chunk Contribution: O_intra = Mask(Q @
2332 90         K.T) @ V
2333 91         T.clear(acc_scores)
2334 92         T.gemm(sQ, sK, acc_scores, transpose_B=True)
2335 93         T.copy(acc_scores, sScores)
2336 94         T.copy(sScores, sScores)
2337 95
2338 96         # Apply Causal Mask
2339 97         for i, j in T.Parallel(chunk_size, chunk_size):
2340 98             if j > i:
2341 99                 sScores[i, j] = T.cast(0.0, dtype)
2342 100         T.copy(sScores, sScores)
2343 101
2344 102         T.clear(acc_o)
2345 103         T.gemm(sScores, sV, acc_o)
2346 104         # Store to dedicated buffer sOut_intra
2347 105         T.copy(acc_o, sOut_intra)
2348 106         T.copy(sOut_intra, sOut_intra) # Barrier ensures
2349 107         visibility for step 4
2350 108
2351 109         # 4. Final Accumulation and Output Store
2352 110         for i, j in T.Parallel(chunk_size, head_dim):
2353 111             # Safe read: sOut_inter and sOut_intra are
2354 112             distinct and fully written
2355 113             val = sOut_inter[i, j] + sOut_intra[i, j]
2356 114             t = t_base + i
2357 115             if t < seq_len:
2358 116                 Out[batch_idx, t, head_idx, j] = val
2359 117
2360 118         # 5. Update State: S_state += K.T @ V
2361 119         T.clear(acc_update)
2362 120         T.gemm(sK, sV, acc_update, transpose_A=True)
2363 121
2364 122         # Move to shared buffer S_kv_update (casts to bf16)
2365 123         T.copy(acc_update, S_kv_update)
2366 124         T.copy(S_kv_update, S_kv_update) # Barrier
2367 125
2368 126         for i, j in T.Parallel(head_dim, head_dim):
2369 127             S_state[i, j] += S_kv_update[i, j]
2370 128         T.copy(S_state, S_state) # Barrier for next chunk
2371 129
2372 130     iteration
2373 131
2374 132     return tilelang.compile(kernel, out_idx=3, target="cuda")
2375 133

```