Gen-AI in a Bottle Experiments with LLMs to generate HPC kernels Upasana Sridhar, Elliott Binder, Tze Meng Low

Carnegie Mellon University



Writing Libraries is a Recurring Task

- Library implementations must be rewritten for different:
 - Datatypes float, double, int8, etc
 - ISAs x86 AVX, AVX2, ARM NEON, etc
 - Memory hierarchies: cach size, associativity, levels



Writing Libraries is a Recurring Task

- Library implementations must be rewritten for different:
 - Datatype precisions float, double, int8, etc
 - ISAs x86 AVX, AVX2, ARM NEON, etc
 - Memory hierarchies: cach size, associativity, levels
- Repeated for each API in the library





Per µArch

How can we make library-writing more productive?

© 2025 Upasana Sridhar



Specialization can be Simplified

- Portable library frameworks such as BLIS^[1] reduce specialization effort
- Construct an abstraction that is shared across:
 - API calls,
 - µArchs, and
 - datatypes



[1] Field G. Van Zee, Tyler M. Smith, Bryan Marker, Tze Meng Low, Robert A. Van De Geijn, Francisco D. Igual, Mikhail Smelyanskiy, Xianyi Zhang, Michael Kistler, Vernon Austel, John A. Gunnels, and Lee Killough. 2016. The BLIS Framework: Experiments in Portability. ACM Trans. Math. Softw. 42, 2, Article 12 (June 2016), 19 pages. https://doi.org/10.1145/2755561



Specialization can be Simplified

- Portable library frameworks such as BLIS reduce specialization effort
- Construct an abstraction that is shared across:
 - API calls,
 - µArchs, and
 - datatypes



ISA specialization is still necessary



Can we use LLMs to specialize libraries?



© 2025 Upasana Sridhar

Using LLMs for HPC Code-Generation

1. LLMs for Code-Generation: Tips and Tricks

2. Designing High-Performance GEMM Microkernels

3. Encoding Microkernel Design Principles with Prompt Engineering



Microkernels present a solvable problem

- Restricted scope has benefits:
 - Fixed problem size
 - Microkernel code does not have to handle many edge cases
 - Commonly used example
- Complexity still scales with datatype precision and ISA
 - LLM needs to retrieve the correct "form" for specialized instruction



Goal: Generate specialized microkernels using LLMs



Performance Expertise is Latent

Prompt:

Write a high-performance C-function to implement double-precision matrix multiplication with A: 6x256, B:256x8 and C: 6x8 Function signature: `dgemm(double const * A, const double * B, double * C)`

Does not automatically use any ISA-specific code



Providing ISA is helpful but insufficient

Prompt:

Write a high-performance C-function to implement double-precision matrix multiplication with A: 6x256, B:256x8 and C: 6x8 Function signature: `dgemm(double const * A, const double * B, double * C)`

Context Architecture: AVX2 vector instructions Data type: double void dgemm(double const * A, const double * B, double * C) { int i, j, k; m256d a0, b0, c0, ..., a6, b6, c6; m256d t0, ..., t14, t15; // Loop over the rows of C for (i = 0; i < 6; i += 8) { for (j = 0; j < 8; j += 16) {</pre> // Initialize the accumulators a0 = mm256 setzero pd();b0 = mm256 loadu pd(B + j);c0 = mm256 setzero pd();



Providing ISA is helpful but insufficient

Prompt:

Write a high-performance C-function to implement double-precision matrix multiplication with A: 6x256, B:256x8 and C: 6x8 Function signature: `dgemm(double const * A, const double * B, double * C)`

Context

Architecture: AVX2 vector instructions Data type: double void dgemm(double const * A, const double * B, double * C) { int i, j, k; m256d a0, b0, c0, ..., a6, b6, c6; m256d t0, ..., t14, t15; // Loop over the rows of C for (i = 0; i < 6; i += 8) { for (j = 0; j < 8; j += 16) { // Initialize the accumulators a0 = mm256 setzero pd();Using specialized b0 = mm256 loadu pd(B + j);datatypes c0 = mm256 setzero pd();



Providing ISA is helpful but insufficient

Prompt:

Write a high-performance C-function to implement double-precision matrix multiplication with A: 6x256, B:256x8 and C: 6x8 Function signature: `dgemm(double const * A, const double * B, double * C)`

Context Architecture: AVX2 vector instructions

Data type: double



Model needs more guidance for how to write code



A Primer on Priming LLMs

- HPC-code is a tiny fraction of training data
 - Hard to retrieve without context

Language Representation on GitHub



Electrical & Computer ENGINEERING

Scaffolded Learning

- Prompt can be used to "prime" LLM to use the right subset of its knowledge
- Context could contain
 - text-based information
 - "Worked examples"



What should be encoded in the microkernel prompt?

Electrical & Compute

Designing GEMM Microkernels

<pre>#define FLOAT_CONV_TILE_C(step, a, b, W_ob, C_ob)</pre>	λ.
<pre>float const * a_ptr = a;</pre>	۸
a_reg = _mm256_broadcast_ss(a_ptr + 0 * step);	λ.
b0 = _mm256_load_ps(b);	λ.
<pre>b1 = _mm256_load_ps(b + FLOAT_SIMD);</pre>	λ.
c12 = _mm256_broadcast_ss(a_ptr + 1 * step);	λ.
c0 = _mm256_fmadd_ps(a_reg, b0, c0);	λ.
c1 = _mm256_fmadd_ps(a_reg, b1, c1);	Λ.
a_reg = _mm256_broadcast_ss(a_ptr + 2 * step);	λ.
a_ptr += 3*step;	N
c2 = _mm256_fmadd_ps(c12, b0, c2);	N
c3 = _mm256_fmadd_ps(c12, b1, c3);	λ.
c12 = _mm256_broadcast_ss(a_ptr + 0 * step);	λ.
	N
c4 = _mm256_fmadd_ps(a_reg, b0, c4);	۸
c5 = _mm256_fmadd_ps(a_reg, b1, c5);	λ.
a_reg = _mm256_broadcast_ss(a_ptr + 1 * step);	١.

x86 AVX2 float Kernel

a_0 = vld1q_dup_f32(a + 0 * step + 0 * FLOAT_SIMD);\ b_0 = vld1q_f32(bb + 0*C_ob + (0 * 4 + 0)*FLOAT_SIMD);\ asm volatile ("fmla %0.4s, %1.4s, %2.s[0]" : "+w"(c 0 0) : "w"(b $b_1 = vld1q_{f32}(bb + 0*C_{ob} + (0 * 4 + 1)*FLOAT_{SIMD});$ ___asm__ volatile ("fmla %0.4s, %1.4s, %2.s[0]" : "+w"(c_0_1) : "w"(b_ $b_2 = vld1q_{f32}(bb + 0*C_{ob} + (0 * 4 + 2)*FLOAT_{SIMD});$ ___asm__ volatile ("fmla %0.4s, %1.4s, %2.s[0]" : "+w"(c_0_2) : "w"(b_2) b 3 = vld1q f32(bb + 0*C ob + (0 * 4 + 3)*FLOAT SIMD);___asm__ volatile ("fmla %0.4s, %1.4s, %2.s[0]" : "+w"(c_0_3) : "w"(b_ a_1 = vld1q_dup_f32(a + 1 * step + 0 * FLOAT_SIMD);\ asm volatile ("fmla %0.4s, %1.4s, %2.s[0]" : "+w"(c 1 0) : "w"(b __asm__ volatile ("fmla %0.4s, %1.4s, %2.s[0]" : "+w"(c_1_1) : "w"(b_ __asm__ volatile ("fmla %0.4s, %1.4s, %2.s[0]" : "+w"(c_1_2) : "w"(b_ a_2 = vld1q_dup_f32(a + 2 * step + 0 * FLOAT_SIMD);\ asm volatile ("fmla %0.4s, %1.4s, %2.s[0]" : "+w"(c 2 0) : "w"(b _asm__ volatile ("fmla %0.4s, %1.4s, %2.s[0]" : "+w"(c_2_1) : "w"(b_

ARM NEON Float Kernel

There is a shared underlying structure

© 2025 Upasana Sridhar



Both are a Vectorized Outer Product



We need encode this pattern into the prompt

Steps to write a Microkernel

Transformation Class	Transformation Required				
Indexing	column-major A				
Loop Order	kij				
Vectorization	Vectorize j-loop use bcast, vload and fma				
Loop Unrolling	Unroll i and j loops				
Loop Invariant Analysis	Hoist B-loads				
Redundant Store Removal	Delay C stores				
Register Allocation	Reuse registers for A broadcast				

```
// naive loops
for (i = 0; i < 6; i++) {
  for (j = 0; j < 8; j++) {
    C[i * 8 + j] = 0.0;
    for (k = 0; k < 256; k++) {
        C[i * 8 + j] += A[i * 256 + k] * B[k * 8 + j];
    }
}</pre>
```



Create prompts with these steps in the context

CodeLLaMA as a Kernel Generator

- Small model, runs locally on a mini-pc
- Trained specifically to generate code
- Weights are available to refine/finetune

Name	codellama/CodeLlama-7b-Instruct-hf				
architecture	LlamaForCausalLM				
Embedding Size	4096				
Max Seq. Length	16384				
model_type	llama				
Heads	32				
Layers	32				
DataType	bfloat16				
Vocabulary Size	32016				



The Prompt Types

• We tried 5 different types



Least Prescriptive, Most Generalizable Most Prescriptive, Least Generalizable

19

Electrical & Computer

Focus on generating entire microkernel



Least Prescriptive, Most Generalizable Most Prescriptive, Least Generalizable

20

Electrical & Computer



Context-Encoding: Entire Microkernel

- 1. Explicitly encode kernel pattern as transformations from naive loops.
- 2. Implicitly encode the pattern with a related example





Generating the entire microkernel is hard

Prompt Strategy	Indexing	Loop	Vectorize	Unroll	Loop	Redunda	nt Register	Errors	%Peak
		Order			Invari- ant	Store	Alloc.		
Simple Prompt (Scalar)	X	X (ikj)	Х	X	Х	Х	n/a	None	7%
Simple Prompt (Vector)	X	X (ijk)	\checkmark	Х	\checkmark	Х	\checkmark	No reduction within vector register	-
Step-by-step Transformations	~	X (ijki)	~	X (not j)	\checkmark	\checkmark	\checkmark	Undeclared variables. Extra i-loop. Wrong store indexing	-
Style Transfer	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	X (20)	Undeclared variables	99%

Create prompts with these steps in the context

What Worked

- CodeLLaMA can reproduce patterns very well
 - Patterns shown in the context are applied even if they are not often seen in training data
- Breaking down each step in a complex reasoning process improved performance
 - When asked to emit a precise piece of code, Code-LLaMA was the most successful
- Style-transfer works well, but is surface level
 - General structure is copied which is the intent in this case. But requires at least one "golden kernel" to be used as a seed



What Didn't

- When an arithmetic calculation was left implicit, CodeLLaMA often guessed incorrectly
 - e.g. the stride of a vectorized loop is reproduced exactly from the example
- Inter-connected constraints leads to confused code-generation
 - When a loop is involved in unrolling as well as vectorization, the generated code muddles the transformations



Towards Fully Automatic Kernel Generation

- Preliminary success with using off-the-shelf LLMs for writing specialized GEMM code.
 - Good reasons to do a hybrid approach with traditional code-gen and LLMs
- A good prompt shows success with a small model
 - could be run natively on the target device for the library
- Systematic steps to write microkernels is impactful
 - Serves as a metric to evaluate the performance of other models
 - Can be used in finetuning
 - Integrated into guardrails for provably correct code-generation



Thank you!



https://linktr.ee/upasanasridhar



🔿 Meta



© 2025 Upasana Sridhar

Prompt Strategies: Only Vectorization

Single Step:

Objective Vectorize the following code snippet

Vectorization Principles

Vectorization Examples Original code Vectorized code Analysis (text)

Prompt Vectorize the following code snippet

Table-filling

You are a vectorization expert...

Example

Objective Identify the vector instructions to fill out the table below. Use the AVX2 vector ISA with float scalars

```
"Vector":
   {"load": "_mm256_load_ps"
        <a list of vector instructions>
```

Prompt

Identify the vector instructions to fill out the table below. Use the AVX2 vector ISA with <u>double</u> scalars

