# Leveraging LLMs to Automate Energy-Aware Refactoring of Parallel Scientific Codes

**Matthew T. Dearing** [1]  **Yiheng Tao** [1]  **Xingfu Wu** [2]  **Zhiling Lan** [1]  **Valerie Taylor** [2]

## Abstract

While large language models (LLMs) are increasingly used for generating parallel scientific codes, most efforts emphasize functional correctness, often overlooking performance, especially energy efficiency. We propose LASSI-EE, an automated LLM-based refactoring framework that generates energy-efficient parallel codes through a multi-stage, iterative approach integrating runtime power profiling, energy-aware prompting, self-correcting feedback loops, and an LLM-as-a-Judge agent for automated screening of code solutions. We introduce energy-reduction@$k$, a novel metric that quantifies expected energy reduction when generating $k$ code candidates and selecting the most energy-efficient, enabling systematic evaluation of multi-attempt generation strategies. Evaluating 20 HeCBench applications and two miniApps on NVIDIA A100 and AMD MI100 GPUs, a single run ($k = 1$) with LASSI-EE delivers refactored parallel codes with an average 29% expected energy reduction at an 81% pass rate, representing a 2.8× improvement over vanilla LLM prompting. Multiple runs ($k = 3$) achieve an average 48% expected energy reduction at a 97% pass rate. These results are consistent across devices, demonstrating LASSI-EE's effectiveness across diverse hardware architectures.

## 1 Introduction

High-quality code generation is important to advance scientific discovery. Applications developed for high-performance computing (HPC) drive discovery across scientific domains, but their parallel implementations are complex, performance-sensitive, and increasingly constrained by energy considerations. Current exascale systems require between 24 and 39 MW of power (Strohmaier et al., 2025), while hyperscale technology companies are seeking facilities supporting gigawatts of power, introducing challenges in expanding existing transmission networks. As energy demands strain both computing infrastructure and the environment, pursuing energy efficiency opportunities across the HPC stack—from hardware to the software delivering scientific solutions—has become essential.

Large language models (LLMs) are now widely used for code generation, with recent attention on parallel scientific codes. Existing work demonstrates that LLMs can generate functionally correct parallel scientific codes and translate code between parallel programming frameworks such as CUDA and OpenMP (Dearing et al., 2024) and others (Nichols et al., 2024; TehraniJamsaz et al., 2024;

Ranasinghe et al., 2025; Zhu et al., 2024). However, most LLM-based code generation efforts prioritize functional correctness over performance and energy efficiency. While generating correct code is fundamental, very little work addresses performant parallel code generation, which involves both execution time and power consumption. This study addresses this gap by developing an effective technique for energy-efficient code generation that considers the target execution platform.

In this paper, we propose LASSI-EE (LLM-based Automated Self-correcting pipeline for generating parallel ScIentific codes for Energy Efficiency), an automated LLM-based refactoring framework designed to generate energy-efficient parallel codes tailored to target execution platforms. LASSI-EE achieves improved energy efficiency through a multi-stage, self-correcting, iterative approach that considers the computing platform used for execution. The framework integrates real-time power profiling, energy-aware context generation, self-correction feedback loops for autonomous error resolution, and an LLM-as-a-Judge agent that provides automated screening of functional equivalence between the source and refactored codes.

Generating energy-efficient parallel codes at scale could supply valuable training data for LLMs, but verifying correctness and efficiency remains difficult because traditional methods require costly tests, annotations, and expertise for complex HPC codes. While the LLM-as-a-Judge agent is not a substitute for formal verification, it enables practi-

[1]Department of Computer Science, College of Engineering, University of Illinois Chicago, Chicago, Illinois USA [2]Argonne National Laboratory, Lemont, Illinois, USA. Correspondence to: Matthew T. Dearing <mdear2@uic.edu>.

cal, automated filtering at scale, offering a critical preliminary step toward robust verification frameworks for LLM-generated scientific codes.

LASSI-EE builds upon the original LASSI framework (Dearing et al., 2024), which demonstrated that LLMs can translate code between CUDA and OpenMP by augmenting an LLM with domain context, self-prompting, and iterative self-correction feedback. While the original LASSI focused on cross-language translation and relied on heuristic human review for validation, LASSI-EE extends this foundation with novel components for energy optimization and automated evaluation: (1) real-time power profiling integration, (2) energy-aware context generation and refactoring plans with iterative prompting strategies tailored to energy optimization, and (3) an LLM-as-a-Judge agent for automated screening.

To quantify the energy optimization capabilities of LLM-based code generation, we introduce energy-reduction@$k$, a new metric inspired by the established pass@$k$ metric (Chen et al., 2021), to evaluate beyond code correctness with a focus on energy efficiency. By providing quantitative guidance on the relationship between generation effort ($k$) and expected energy savings, this metric enables developers and researchers to make informed decisions about balancing computational costs of code generation against potential energy benefits.

We evaluate LASSI-EE using 22 representative applications from science and engineering domains on both NVIDIA A100 and AMD MI100 GPUs. Our key findings show that, across all tested applications and target systems, LASSI-EE achieves an average energy reduction of 35% in its successful code generations. LASSI-EE produces codes with average pass@1 and pass@3 rates of 81.4% and 97.3%, respectively. Of practical importance, the expected energy reductions achieved by LASSI-EE are, on average, 2.8× greater than those of a vanilla LLM when generating parallel codes in a single attempt ($k = 1$). For $k = 3$ repeated attempts, LASSI-EE delivers up to 48% lower energy consumption on average for both A100 and MI100 GPUs.

Overall, the key contributions of this paper include:

- We present **LASSI-EE**, an LLM-based, automated, self-correcting framework for generating energy-efficient parallel scientific codes. LASSI-EE integrates runtime power profiling, energy-aware prompting, and an automated code screening mechanism using an LLM-as-a-Judge agent (§2).

- We introduce **energy-reduction@$k$**, a new metric that quantifies the expected energy reduction when generating $k$ code candidates, enabling systematic and fair comparison across energy optimization methods (§3).

- We analyze LASSI-EE's effectiveness through **cross-platform evaluation**, using 22 applications on both NVIDIA A100 and AMD MI100 GPUs. LASSI-EE achieves nearly 50% average energy reduction when generating $k = 3$ code candidates (§ 4-5).

## 2 LASSI-EE FRAMEWORK

Figure 1 illustrates the overall pipeline, highlighting the iterative stages of the LASSI-EE framework. Given *a parallel code* executed on a measurable accelerator device (e.g., a GPU like NVIDIA A100 or AMD MI100), LASSI-EE conducts code performance and power profiling on *the target system* (Stage 0), baseline code refactoring using an LLM (Stage 1), context-driven and self-prompted refactoring (Stages 2 and 3), and outputs the final energy-efficient code (Stage 4). The framework is designed to be LLM-agnostic and extensible to multiple architectures and programming models.

### 2.1 Stage 0: Code Profiling

The initial stage captures baseline performance metrics of parallel source code on the target system. This code is compiled, executed, and profiled for power consumption, with metrics stored in memory for comparison with refactored versions.

LASSI-EE collects real-time power measurements using platform-specific monitoring libraries: NVIDIA's Management Library (NVML) on NVIDIA GPUs and ROCm System Management Interface (ROCm-SMI) on AMD GPUs. Power samples are captured at 10-millisecond intervals during execution, then aggregated to calculate average power and total energy consumed. Unlike static analysis approaches that predict performance, LASSI-EE takes direct, runtime measurements to provide accurate feedback for iterative refinement. An example power profile for the `jacobi` application is in Appendix B Figure 5.

### 2.2 Stage 1: Baseline Energy-Aware Refactoring

Stage 1 establishes another baseline by evaluating an out-of-the-box LLM's inherent capability to generate energy-efficient code, serving as a comparison point for the value added by the context-building, iterative refinement, and self-correction in later stages. We refer to this enhanced approach as the *vanilla LLM*.

A context-free zero-shot prompt instructs the vanilla LLM to refactor the source code for energy efficiency without context, examples, or self-correction assistance Table 5 in Appendix A. The generated code is compiled and executed once. If either fails, then the vanilla LLM attempt is marked as unsuccessful. If successful, then a power profile is captured and the code is evaluated by the LLM-as-a-Judge agent
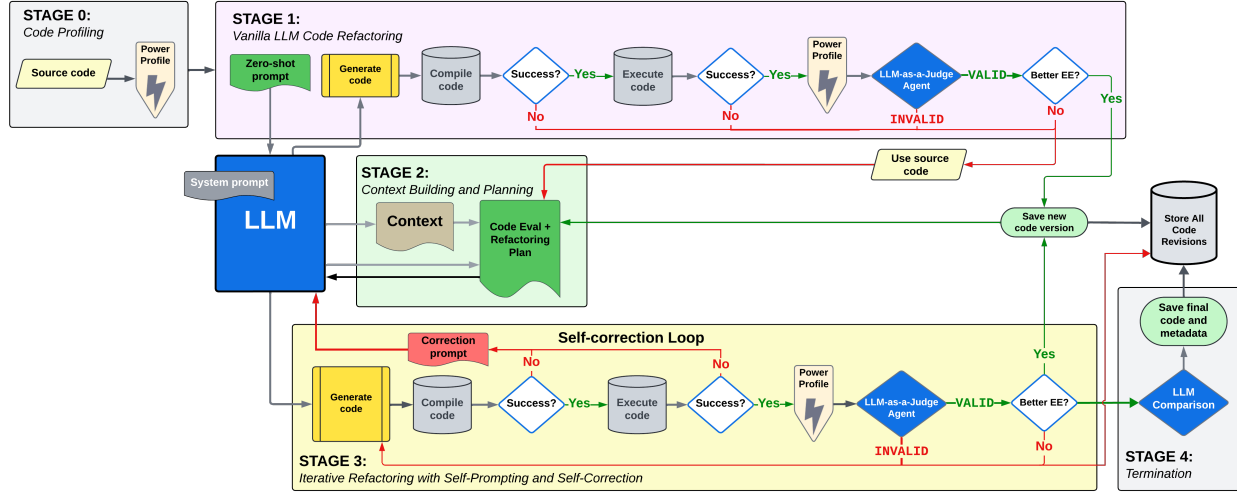
*Figure 1.* Diagram of the **LASSI-EE pipeline**.

(§ 2.3). Code that passes screening and demonstrates energy savings becomes a candidate for the final solution and seeds the subsequent context-driven refinements in Stages 2 and 3. Otherwise, the source code passes to Stage 2.

The same system prompt is used for all inference steps to ensure differences in results stem from LASSI-EE mechanisms rather than inconsistent prompting. The system prompt instructs the LLM to act as an expert C++/CUDA (or HIP) engineer refactoring code to improve energy efficiency (see Appendix A Table 6).

## 2.3 LLM-as-a-Judge: Candidate Code Screening

A fundamental challenge with automating code generation is verifying functional equivalence of generated and source code. The original LASSI framework for code translation (Dearing et al., 2024) relied on manual human review, limiting scalability. While formal verification methods exist to rigorously prove functional equivalence (Abadi et al., 2019; Siegel et al., 2015; 2006; Blom & Huisman, 2014), they are computationally expensive, require expert configuration, and may not readily apply to all code transformations.

LASSI-EE takes an incremental step toward automated validation by incorporating an LLM-as-a-Judge agent that provides automated screening of likely functionally equivalent candidates. This agent replaces manual human review with a practical filtering mechanism that operates at scale and enables the metrics in § 3.1. The Judge agent receives the source code, refactored code, standard outputs from both, and energy metrics, then returns a binary verdict: VALID (passes screening, suitable for energy comparison) or INVALID (failed generation task, filtered from energy metrics).

The Judge prompt (Appendix A Table 8) instructs the LLM to accept minor output formatting differences (e.g., whitespace, layout, and numerical precision within tolerance) while flagging significant algorithmic changes, missing functionality, or substantial output discrepancies as INVALID. This automated screening filters problematic refactorings at scale before they contribute to energy metric calculations. Future work will explore formal verification integration for stronger correctness guarantees.

## 2.4 Stage 2: Context Building and Planning

LASSI-EE's core prompting strategy uses an enhanced context-building phase to prepare for refactoring. After Stage 1, the pipeline curates LLM-authored summaries of background knowledge and the incoming code (source or Stage 1 output), with the LLM also identifying inefficiencies in the code and proposing specific optimizations.

Context-building involves multiple self-prompting steps. First, the LLM summarizes relevant programming language documentation (e.g., CUDA (NVIDIA Corporation, 2024) or HIP (Advanced Micro Devices, Inc., 2024) programming guide excerpts). Second, the LLM reviews the current code version and generates a description of its functionality and potential energy inefficiencies. These summaries become enriched context for subsequent refactoring stages (prompts in Appendix A Table 9).

Following context building, LASSI-EE prompts the LLM to develop a structured refactoring plan (Appendix A Table 10). identifying specific optimization strategies, such as reducing memory transfers, optimizing kernel launch configurations, eliminating redundant computations, improving memory access patterns, or leveraging hardware-specific features. The plan provides a roadmap for each Stage 3 refactoring

attempt.

## 2.5 Stage 3: Iterative Self-Prompted Refactoring

This stage represents the core iterative refinement loop of LASSI-EE. Using the context and refactoring plan from Stage 2, Stage 3 repeatedly generates, tests, and evaluates energy-optimized code versions. The LLM responds with a description of changes and a code block for evaluation (Appendix A Table 11).

*Self-Correcting Feedback Loop.* A critical innovation in LASSI-EE is autonomous error correction without human intervention. When compilation or runtime errors occur, an iterative correction process continues until all errors are eliminated (see Self-correction Loop in Figure 1). LASSI-EE captures compiler errors and re-prompts the LLM to generate revised code. Following successful compilation, any runtime errors from code execution are provided to the LLM with the buggy code to generate corrections. This cycles until the generated code runs error-free, with prompts ensuring the LLM maintains function signatures, kernel logic, and original behavior while fixing errors (Appendix A Table 7). A configurable maximum threshold of self-correction attempts prevents runaway iterations. None of the trials reported in § 4 terminated prematurely due to unrepairable code errors. This self-correction mechanism is not applied to the vanilla LLM generation in Stage 1, allowing for fair comparison between an unassisted LLM and the LASSI-EE framework.

*Energy-Driven Iteration and Screening.* Once generated code passes self-correction, power profiles are collected during execution, energy measures calculated, and the code evaluated by the LLM-as-a-Judge (§ 2.3). If screening fails, this version is archived but filtered from energy comparisons, and Stage 3 restarts with the same refactoring plan. If screening passes, energy metrics are compared to current "best" code (source, vanilla LLM version, or any Stage 3 version), defined as code that compiles, executes, passes LLM-as-a-Judge screening, and consumes the least energy. If energy savings improves, this revision becomes the new "best," prompting the LLM to write a *new* refactoring plan based on the latest optimizations, initiating another Stage 3 iteration.

*Adaptive Temperature Tuning.* During Stage 3, adaptive adjustment is applied to the LLM `temperature` parameter when refactorings do not yield additional energy savings. With default `temperature` set to 0.2, code generations tend to be more deterministic. Increasing this value allows broader token selection, emulating more "creativity" in LLM responses. Following iterations without energy savings, the pipeline increments `temperature` by 0.2 before another code generation with the *same* refactoring plan. A stopping condition triggers when `temperature` exceeds

0.8, suggesting exhausted optimization capabilities. Some reasoning LLMs, (e.g., o3, o4-mini, and GPT-5) do not support the sampling parameter of `temperature`. In these cases, the same stopping mechanism still triggers after a fixed number of unsuccessful iterations. Additionally, a configurable maximum iteration count for Stage 3 prevents runaway computation.

## 2.6 Stage 4: Termination and Output

Upon termination (triggered when `temperature` exceeds 0.8 or maximum Stage 3 iterations are reached), the framework identifies the "best" code version from all candidates. If no generated code passes screening and reduces energy, the source code is retained as the "best" solution and LASSI-EE reports no energy improvement.

When the "best" code achieves energy reduction (such as the `jacobi` example in Appendix B Figure 6), LASSI-EE performs a final LLM-based comparison analysis identifying code transformations, energy-relevant optimizations, and potential reasons for the measured energy reduction (Appendix A Table 12). This analysis generates a natural language summary documenting changes. While no automated decisions are based on this comparison, it enables human validation, supports post-hoc analysis, and documents which code refactoring tactics yield energy savings on the target platform.

The final "best" code is saved along with comprehensive metadata: all intermediate code versions, energy profiles, Judge agent verdicts, refactoring plans, iteration counts, and (when applicable) the LLM-generated comparison analysis. This metadata archive supports future LLM fine-tuning for energy-aware code generation and retrospective analysis of the LASSI-EE optimization process.

## 3 METRICS AND EXPERIMENTAL SETUP

### 3.1 Evaluation Metrics

***Pass@k.*** To measure the likelihood that LASSI-EE generates code that passes automated screening by the LLM-as-a-Judge, we adopt the pass@$k$ metric (Chen et al., 2021). This metric quantifies the probability that at least one code solution passes screening when generating $k$ independent refactorings from a given source code. Formally, for a task with $n$ total generation attempts ($n \geq k$), let $c$ denote the number of codes that pass candidate screening:

$$\text{pass@}k = \mathbb{E}\left[1 - \frac{\binom{n-c}{k}}{\binom{n}{k}}\right], \quad (1)$$

where the expectation is taken over all evaluated tasks and we focus on $k \in \{1, 3, 5\}$ for practical evaluation.

***Energy-reduction@k.*** While pass@$k$ addresses the prob-

ability of generating passable code, we introduce energy-reduction@$k$, a novel metric for energy-aware code generation that quantifies expected optimization quality: "If I generate $k$ code candidates and select the most energy-efficient one that passes screening, what energy reduction can I expect?" Formally, let $\bar{E}_{\text{source}}$ denote the average energy consumption of the source code across all trial executions, and let $E_i$ represent the energy consumption of candidate $i$ from the subset of screened generated codes:

**energy-reduction@$k$** :=
$$\text{pass@}k \times \text{energy-reduction@}k^{\text{valid}} \quad (2)$$

energy-reduction@$k^{valid}$ :=
$$\mathbb{E}\left[\max\left\{\left(\frac{\bar{E}_{\text{source}} - E_i}{\bar{E}_{\text{source}}}\right)^+ \mid E_i \in c_k\right\}\right] \quad (3)$$

where $c_k$ represents $k$ candidates sampled from generated codes that passed screening among $n$ independent runs; $E_i$ are energy measurements from candidates that received `VALID` screening verdicts; $\mathbb{E}[\cdot]$ denotes expectation (estimated by Monte Carlo); the ratio represents the relative energy reduction achieved by candidate $i$ compared to the source code. The algorithm for calculating energy-reduction@$k$ is listed in Appendix C.

The energy-reduction@$k$ metric represents the *expected energy reduction per generation attempt*, accounting for the probability that generation may fail to produce a candidate code. This provides realisitic assessment of practical utility: if a system has low pass@$k$ but high energy-reduction@$k^{valid}$, then the actual expected benefit will be diminished by the high failure rate.

In our experiments, we set $n = 30$ independent generation attempts per application per device. Each run executes both a vanilla LLM generation (without LASSI-EE enhancements) and the complete LASSI-EE pipeline, producing paired results. For LASSI-EE, a code sample passes screening (contributing to the count $c$ in Equation 1) if it: (i) compiles, (ii) executes without errors, and (iii) receives a `VALID` verdict from the LLM-as-a-Judge.

### 3.2 Applications

Table 1 lists the 22 applications tested in this study. The HeCBench suite (Jin & Vetter, 2023) provides a robust collection of heterogeneous computing implementations of scientific kernels and applications. We selected 20 source codes from HeCBench spanning multiple categories and sizes to capture a broad range of problem types. We also introduce two miniApps to evaluate LASSI-EE with longer, more complex codes, miniMDock (Thavappiragasam et al., 2020) and XSBench (Tramm et al., 2014). Each application contains multiple core parallel kernels referenced

from the main program script. We merged these separate files and modified the Makefile to build the app using this merged script, enabling extended codes with 1,472 lines from XSBench and 2,394 for miniMDock, considerably longer than the maximum HeCBench option of 1,079 lines (`lid-driven-cavity`). We verified the applications executed as expected from this merged source code.

*Table 1.* Characteristics of the selected source code applications.

| Category | Application | Lines of Code | Input args |
|---|---|---|---|
| **HeCBench** | | | |
| Bandwidth | randomAccess | 158 | [10] |
| Bioinformatics | all-pairs-distance | 328 | [10000] |
| Computer vision and image processing | colorwheel | 154 | [10000, 8, 1] |
| | marchingCubes* | 574 | [100] |
| Cryptography | chacha20* | 130 | [300000] |
| Data compression and reduction | segment-reduce | 95 | [16384, 100] |
| Data encoding, decoding, or verification | entropy* | 171 | [10000, 1024, 1] |
| | murmurhash3 | 245 | [750000, 500] |
| Graph and Tree | floydwarshall | 295 | [2048, 100, 256] |
| Language and kernel features | layout | 197 | [2500] |
| | threadfence | 142 | [100, 250000000] |
| Machine learning | dense-embedding | 193 | [10000, 8, 1] |
| Math | jacobi† | 235 | None |
| | jaccard | 417 | [1024, 512, 1000] |
| | matrix-rotate | 67 | [30000, 100] |
| Search | bsearch | 279 | [10000, 1] |
| | keogh* | 143 | [256, 22500000, 100] |
| Signal processing | extrema | 349 | [750] |
| Simulation | pathfinder | 286 | [10000, 1000, 1000] |
| | lid-driven-cavity | 1079 | None |
| **miniApps** | | | |
| Particle transport | XSBench | 1472 | [1700000000] |
| Molecular dynamics | miniMDock | 2394 | [7cpa_ligand.pdbqt, -nrun 500] |

\* Code includes a dependency that was not considered during refactoring.
†Because this app does not accept runtime arguments, we modified the thread block size multiplier from 2048 to 16384 in the source code.
‡ App contains parallel and serial components, so the serial portion was removed.

### 3.3 Target Platforms

Our experiments were conducted on Chameleon testbed nodes (Kate Keahey, 2020) at UC and TACC sites, each featuring two AMD EPYC 7763 64-core processors (128 physical cores, 256 threads total). To evaluate LASSI-EE's performance across different GPU architectures, we implemented the pipeline on two node types: **NVIDIA A100** (80 GB GPU memory, NVIDIA driver version 560.35.05, CUDA 12.6) and **AMD MI100** (32 GB GPU memory, ROCm driver version 6.14.0).

Each of the 22 applications was run independently on both GPU platforms, resulting in 44 total experimental trials. For the NVIDIA A100 platform, applications were implemented in CUDA C++, while the AMD MI100 platform

used HIP C++ implementations. The LASSI-EE codebase is developed in Python 3.12.3. Interfacing between LASSI-EE and the system command line was facilitated through the `subprocess` library for calls to gather system state information, compile code, and execute binaries.

*Code Compilation.* Each code was compiled in the same compute environment with consistent compiler flags across all experimental runs so that each generated code benefited from identical baseline compiler optimizations.

For CUDA codes on the NVIDIA A100:

```
nvcc -std=c++14 -Xcompiler -Wall -arch=sm_80 -O3
```

For HIP codes on the AMD MI100:

```
hipcc -std=c++14 -Wall -O3 --amdgpu-target=gfx908
```

*Power Measurement.* Our power measurements focus on GPU power consumption. For the NVIDIA A100, the pipeline utilizes `pyNVML`, the Python bindings to the NVIDIA Management Library, to collect instantaneous power measurements from the GPU. For the AMD MI100, the pipeline uses `rocm-smi` command-line utility with the `--showpower` flag via subprocess calls.

Power usage is sampled every 0.01 seconds on both platforms, beginning with a brief initial period to estimate average idle power. During code execution, power is collected in a separate thread at the same interval. Timestamps are recorded at the launch and completion of code execution to mark the runtime window. Measurements continue for 15 seconds post-execution to capture a stable post-run idle state. Pre- and post-run idle values are averaged to estimate idle system power during execution.

To ensure accurate comparisons, we subtract idle power from all measurements to calculate *net power usage*. This correction accounts for idle power drift, especially during longer sessions. *Net energy* is then estimated by summing $N$ idle-subtracted power samples $P_i$ multiplied by the time interval ($\Delta t = 0.01$ s), approximating the continuous energy integral:

$$E_{net} \approx \Delta t \sum_{i=1}^{N} \max\{P_i, 0\} \qquad (4)$$

Any negative power measurements are clamped to zero to avoid non-physical outcomes due to noise.

### 3.4 Other Configurations

LASSI-EE is designed as an LLM-agnostic pipeline. In this study, we employed a multi-model approach to leverage strengths of different LLMs while maintaining consistency across our experimental trials. We used o4-mini for context building and code generation across the 20 HeCBench applications. For the two miniApps, we used GPT-5 as the code generator due to context length constraints in o4-mini, while

retaining o4-mini for context building. We used GPT-4.1 as the LLM-as-a-Judge agent for its larger context window and detailed evaluation capabilities. The o4-mini model does not support temperature tuning, common to OpenAI's o-series reasoning models. While LASSI-EE incorporates temperature adjustment as a stopping condition, this parameter was ignored by o4-mini but was employed for GPT-5 during miniApp generation.

## 4 EXPERIMENTAL RESULTS

We evaluate LASSI-EE on 22 applications, conducting 30 independent trials per application on both NVIDIA A100 and AMD MI100 GPUs. In total, our experiments consist of 1,320 trials (22 × 2 × 30).

### 4.1 Performance Gains

Table 2 presents average measured performance for representative source codes and LASSI-EE-generated codes on both AMD MI100 and NVIDIA A100 GPUs. We list select applications from small to larger lines of code (Table 1) with full data on all 22 applications in Appendix E Table 13. For each application, the consumed energy (J), average instantaneous power (W), and runtime (s) for the source code are averaged over 30 independent trials. The corresponding LASSI-EE measurements are averaged over only those trials that passed the LLM-as-a-Judge screening and consumed less total energy than the source code, with differences shown in parentheses. The 'Best Trial' columns report pipeline metadata for the lowest-energy code generation attempt for the application/GPU pair. 'Self-Corrections' indicates iterations required to fix compilation or execution errors (see the *Self-correction Loop* in Figure 1), and 'Iterations for EE' is the iteration number where the most energy-efficient code was generated. For example, the best `segment-reduce` refactoring on MI100 had no bugs and was generated during the fourth refactoring attempt. Average pipeline runtime per trial is typically around 30 minutes, ranging from 10 minutes to outliers of 80-114 minutes (Appendix E Table 13).

*Cross-Device Consistency.* When LASSI-EE generated energy-efficient code, it achieved average energy savings of 33.9% on NVIDIA A100 and 36.2% on AMD MI100, a difference of only 2.3 percentage points (pp). Average power reductions were 7.9% and 6.5%, respectively, while runtime improvements averaged 26.2% and 29.6%. This consistency demonstrates that LASSI-EE's energy-aware refactoring generalizes effectively across hardware architectures. Overall, average savings are 35.1% in energy, 7.2% in power, and 27.9% in runtime.

*Runtime and Power Trade-offs.* Application-specific energy reductions covered a broad range from 2.5%

*Table 2.* Performance characteristics comparing source code and LASSI-EE-generated code for representative applications. Energy, power, and runtime show LASSI-EE performance with percentage change ($\Delta\%$) in parentheses. Negative values indicate reductions (improvements). These metrics represent conditional averages over trials that compiled, executed successfully, passed LLM-as-a-Judge screening, and consumed less energy than source code. Best-trial statistics report values from the trial with lowest energy per application/GPU. Pipeline time shows average across all successful trials. Complete results for all applications in Appendix E Table 13.

| Application/ GPU | Energy (J) | | Power (W) | | Runtime (s) | | Best Trial | | Avg Pipeline |
|---|---|---|---|---|---|---|---|---|---|
| | Source | LASSI-EE ($\Delta\%$) | Source | LASSI-EE ($\Delta\%$) | Source | LASSI-EE ($\Delta\%$) | Self-Corrections | Iterations for EE | Time (min) |
| **segment-reduce** | | | | | | | | | |
| MI100 | 710.9 | 468.0 ($-34.2$) | 62.5 | 128.0 ($+104.6$) | 11.5 | 4.1 ($-64.6$) | 0 | 4 | 28.9 |
| A100 | 622.0 | 289.8 ($-53.4$) | 38.8 | 79.3 ($+104.3$) | 16.0 | 3.5 ($-78.1$) | 0 | 5 | 13.4 |
| **entropy** | | | | | | | | | |
| MI100 | 1.6 | 0.3 ($-82.7$) | 0.9 | 0.2 ($-83.3$) | 1.8 | 1.8 ($-0.5$) | 0 | 9 | 19.6 |
| A100 | 1718.2 | 1659.1 ($-3.4$) | 55.8 | 58.1 ($+4.0$) | 31.3 | 31.1 ($-0.6$) | 0 | 5 | 16.3 |
| **jacobi** | | | | | | | | | |
| MI100 | 2200.9 | 1597.2 ($-27.4$) | 125.1 | 156.6 ($+25.2$) | 17.8 | 10.4 ($-41.6$) | 0 | 1 | 82.2 |
| A100 | 1216.3 | 648.1 ($-46.7$) | 150.9 | 131.7 ($-12.8$) | 8.1 | 4.7 ($-42.3$) | 0 | 5 | 22.2 |
| **lid-driven-cavity** | | | | | | | | | |
| MI100 | 2358.1 | 2119.8 ($-10.1$) | 122.1 | 127.8 ($+4.6$) | 19.3 | 16.6 ($-14.0$) | 0 | 9 | 43.9 |
| A100 | 1557.1 | 1192.0 ($-23.5$) | 36.2 | 37.6 ($+3.9$) | 43.5 | 36.7 ($-15.7$) | 0 | 7 | 42.8 |
| **XSBench** | | | | | | | | | |
| MI100 | 5864.5 | 5531.6 ($-5.7$) | 142.9 | 168.6 ($+18.0$) | 41.7 | 33.1 ($-20.7$) | 0 | 3 | 53.8 |
| A100 | 3038.4 | 2753.7 ($-9.4$) | 155.9 | 141.1 ($-9.5$) | 19.5 | 20.7 ($+5.9$) | 0 | 6 | 37.9 |
| **miniMDock** | | | | | | | | | |
| MI100 | 5059.8 | 3625.8 ($-28.3$) | 223.3 | 205.6 ($-7.9$) | 22.7 | 17.5 ($-22.9$) | 0 | 4 | 50.2 |
| A100 | 2627.2 | 1568.3 ($-40.3$) | 133.2 | 163.7 ($+23.0$) | 19.7 | 9.5 ($-51.9$) | 0 | 8 | 47.6 |
| **Cross-app average ($\Delta\%$)** | | | | | | | | | |
| MI100 | | ($-36.2\%$) | | ($-6.5\%$) | | ($-29.6\%$) | | 5.8 | 38.3 |
| A100 | | ($-33.9\%$) | | ($-7.9\%$) | | ($-26.2\%$) | | 5.4 | 22.0 |

(`marchingCubes` on MI100) to 96.6% (`colorwheel` on A100). Eighteen of 22 applications on A100 and 17 of 22 on MI100 showed runtime improvements. Several applications showed increased average power and still achieved energy reduction through faster runtimes. For example, `segment-reduce` had 104.3% average power increase on A100 and 104.6% on MI100, but still achieved 53.4% and 34.2% energy reductions, respectively, through 78.1% and 64.6% speedups. This trade-off suggests LASSI-EE's capability to target total energy consumption optimization strategies based on source application and device.

*Application Code Length.* For the two longer miniAps, `XSBench` (1,472 lines) and `miniMDock` (2,394 lines), LASSI-EE delivered average energy savings of 5.7-9.4% and 28.3-40.3%, respectively, between MI100 and A100 GPUs. **With these codes being approximately 3× longer than most HeCBench codes, LASSI-EE demonstrates scalability beyond typical benchmark kernels toward more complex applications.**

These averages in Table 2 reflect only successful generated codes and do not account for the *probability* of generating energy-reducing code during any given run of LASSI-EE. In the following sections, we report the *expected* improvements from one or more runs using pass@$k$ and energy-reduction@$k$ metrics (§3.1), which quantify LASSI-EE's practical utility for automated energy-efficient parallel code

generation.

### 4.2 pass@$k$ Results

A snapshot of the pass@$k$ metrics for vanilla LLM and LASSI-EE across representative applications on both GPU platforms are shown in Table 3, with full data in Appendix E Table 14. Averaging across both devices and all 22 applications, LASSI-EE achieves pass@1 = 81.4%, pass@3 = 97.3%, and pass@5 = 99.0%, compared to vanilla LLM's 41.0%, 70.9%, and 82.6%, respectively. At $k = 1$, comparative improvements are +40.5 pp on average (+35.8 pp on MI100, +45.2 pp on A100), and at $k = 3$, LASSI-EE delivers +26.4 pp average improvement over vanilla LLM (+24.3 pp on MI100, +28.5 pp on A100). **In other words, within a few attempts, LASSI-EE is nearly guaranteed to generate passable code, while an LLM on its own might only get it right 70% of the time.**

*Cross-Device Consistency.* The difference of only 0.4 pp between AMD MI100 and NVIDIA A100 at $k = 3$ illustrates that LASSI-EE performs consistently at generating passable codes on both devices. At $k = 1$, the difference is marginally higher at 5.9 pp (78.5% vs. 84.4%), both still representing significant improvement over vanilla LLM code generation. Figure 2 illustrates these averaged cross-device pass@$k$ results, showing LASSI-EE's consistent improvement over vanilla LLM across both GPU architectures.

*Table 3.* **pass@$k$ Snapshot**: Functional equivalence screening success rates for representative applications. All values as percentages show LASSI-EE's performance at generating *passable* parallel code, with differences over vanilla LLM in parentheses ($\Delta$ pp). The cross-application results are averaged over all 22 benchmark source codes, showing similar performance on A100 and MI100 for each $k$ on average. Complete results for all applications are provided in Appendix E Table 14.

| App/GPU | $k=1$ LASSI-EE ($\Delta$) | $k=3$ LASSI-EE ($\Delta$) | $k=5$ LASSI-EE ($\Delta$) |
|---|---|---|---|
| **Cross-app average** | | | |
| MI100 | 78.5% (+35.8) | 97.1% (+24.3) | 99.2% (+14.9) |
| A100 | 84.4% (+45.2) | 97.5% (+28.5) | 98.9% (+18.1) |
| **Cross-device avg:** | 81.4% | 97.3% | 99.0% |
| **segment-reduce** | | | |
| MI100 | 86.7 (+6.7) | 99.9 (+0.4) | 100.0 (0.0) |
| A100 | 93.3 (+30.0) | 100.0 (+4.1) | 100.0 (+0.3) |
| **entropy** | | | |
| MI100 | 76.7 (+46.7) | 99.1 (+31.9) | 100.0 (+14.3) |
| A100 | 96.7 (+93.3) | 100.0 (+90.0) | 100.0 (+83.3) |
| **jacobi** | | | |
| MI100 | 90.0 (+46.7) | 100.0 (+16.7) | 100.0 (+4.3) |
| A100 | 96.7 (+66.7) | 100.0 (+32.8) | 100.0 (+14.3) |
| **lid-driven-cavity** | | | |
| MI100 | 93.3 (+66.7) | 100.0 (+37.9) | 100.0 (+18.5) |
| A100 | 86.7 (+73.3) | 99.9 (+63.9) | 100.0 (+46.2) |
| **XSBench** | | | |
| MI100 | 80.0 (+43.3) | 99.5 (+23.4) | 100.0 (+8.2) |
| A100 | 100.0 (+66.7) | 100.0 (+28.1) | 100.0 (+10.9) |
| **miniMDock** | | | |
| MI100 | 73.3 (−20.0) | 98.6 (−1.4) | 100.0 (0.0) |
| A100 | 73.3 (+30.0) | 98.6 (+15.4) | 100.0 (+4.3) |

*Diminishing Returns Beyond $k = 3$.* LASSI-EE achieves >97% average pass rates at $k = 3$ over both devices, with 20 of 22 applications reaching 99-100% pass@3, compared to 70.9% for vanilla LLM (Appendix E Table 14). The improvement in LASSI-EE from $k = 1$ to $k = 3$ is +15.9 pp, while improvement from $k = 3$ to $k = 5$ is only +1.7 pp despite 67% additional computational cost, estimated by the average pipeline runtime.

*Application-Specific Results.* LASSI-EE delivers substantial improvements over vanilla LLM for several applications at $k = 1$: `entropy` on A100 improves +93 pp; `XSBench` by +67 pp on A100; `all-pairs-distance` by +83 and +73 on MI100 and A100, respectively; and `lid-driven-cavity` by +73 pp on A100. Three applications show nominal initial lower LASSI-EE pass rates than vanilla LLM on MI100 at $k = 1$ (`dense-embedding` by -3.3 pp, `layout` by -6.7 pp, and `jaccard` by -6.7 pp), with only `jaccard` slightly underperforming at $k \in \{3, 5\}$, while the others reach $\geq 99\%$ by $k = 3$. LASSI-EE struggled with `miniMDock` at $k = 1$ by -20 pp below vanilla LLM, but still reached $\geq 99\%$ by $k = 3$. Following our exemplar `jacobi` app, Appendix B Figure 7 illustrates LASSI-EE's increased capability over vanilla LLM in pass@$k$ for the A100 GPU.
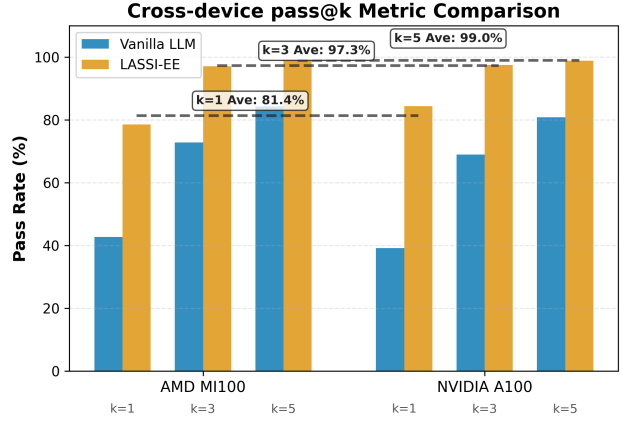


*Figure 2.* Cross-device pass@$k$ results averaged across all 22 applications. LASSI-EE (orange) achieves similar gains on AMD MI100 and NVIDIA A100 over vanilla LLM (blue) for each $k \in \{1, 3, 5\}$.

### 4.3 energy-reduction@$k$ Results

Table 4 lists a snapshot of the expected energy-reduction@$k$ metrics for vanilla LLM and LASSI-EE across representative applications on both GPU platforms, with full data in Appendix E Table 15. As described in §3.1, this metric combines the probability of `VALID` code screening by the LLM-as-a-Judge (pass@$k$) with measured energy savings to quantify the *expected* energy reduction when generating $k$ code candidates and selecting the most energy-efficient solution, given the possibility of code generation fails.

Across all applications and devices, LASSI-EE achieves $\sim$2.8$\times$ the expected energy reduction of vanilla LLM at $k = 1$ (29.1% vs. 10.5%), 1.8$\times$ at $k = 3$ (48.2% vs. 26.3%), and 1.6$\times$ at $k = 5$ (54.6% vs. 34.6%). **These improvements emerge from both higher pass rates and deeper energy reductions strictly due to the LASSI-EE pipeline.**
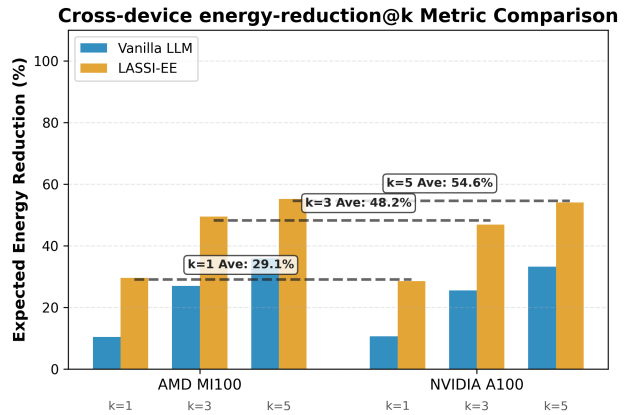


*Figure 3.* Cross-device expected energy-reduction@$k$ results averaged across all 22 applications. LASSI-EE (orange) achieves similar gains on AMD MI100 and NVIDIA A100 over vanilla LLM (blue) for each $k \in \{1, 3, 5\}$.

*Table 4.* **energy-reduction@$k$ Snapshot**: Expected energy reduction for representative applications. All values as percentages show the expected *energy savings* over the source code in the LASSI-EE-generated parallel code, with differences over vanilla LLM in parentheses ($\Delta$ pp). The cross-application results are averaged over all 22 benchmark source codes, showing similar performance on A100 and MI100 for each $k$ on average. Complete results for all applications are provided in Appendix E Table 15.

| App/GPU | $k = 1$<br>LASSI-EE ($\Delta$) | $k = 3$<br>LASSI-EE ($\Delta$) | $k = 5$<br>LASSI-EE ($\Delta$) |
|---|---|---|---|
| **Cross-app average** | | | |
| MI100 | 29.6% (+19.1) | 49.5% (+22.5) | 55.2% (+19.2) |
| A100 | 28.6% (+18.0) | 46.9% (+21.4) | 54.1% (+20.9) |
| **Cross-device avg:** | 29.1% | 48.2% | 54.6% |
| **segment-reduce** | | | |
| MI100 | 28.5 (+18.9) | 57.4 (+33.4) | 68.7 (+38.6) |
| A100 | 48.1 (+25.0) | 66.1 (+24.8) | 74.0 (+26.5) |
| **entropy** | | | |
| MI100 | 63.9 (+54.9) | 88.0 (+52.1) | 90.7 (+39.0) |
| A100 | 3.2 (+3.2) | 6.5 (+6.5) | 8.8 (+8.8) |
| **jacobi** | | | |
| MI100 | 22.7 (+20.5) | 46.5 (+36.3) | 52.5 (+34.6) |
| A100 | 44.8 (+38.8) | 66.4 (+37.7) | 75.2 (+28.3) |
| **lid-driven-cavity** | | | |
| MI100 | 7.8 (+7.7) | 16.4 (+15.9) | 22.9 (+21.9) |
| A100 | 16.6 (+13.9) | 41.3 (+24.7) | 51.0 (+17.7) |
| **XSBench** | | | |
| MI100 | 3.0 (+0.7) | 10.6 ($-$2.4) | 16.3 ($-$4.9) |
| A100 | 5.3 (+5.2) | 15.8 (+15.2) | 19.4 (+18.2) |
| **miniMDock** | | | |
| MI100 | 20.8 (+9.5) | 44.1 (+24.8) | 51.1 (+25.7) |
| A100 | 29.6 (+15.1) | 50.8 (+18.8) | 57.9 (+17.4) |

*Cross-Device Consistency.* As seen with pass@$k$, nominal differences between average energy-reduction@$k$ values for AMD MI100 and NVIDIA A100 show consistent LASSI-EE performance across GPUs. Figure 3 illustrates these cross-device results, showing LASSI-EE's improvement over vanilla LLM across both GPU architectures at $k \in \{1, 3, 5\}$, with the average savings of 29.1% at $k = 1$, 48.2% at $k = 3$, and 54.6% at $k = 5$, demonstrating consistent results with differences between devices of only 1 pp, 0.4 pp, and 1.1 pp, respectively.

*Diminishing Returns Beyond $k = 3$.* LASSI-EE increases its average expected energy savings by +19.1 pp from $k = 1$ to $k = 3$, corresponding to significant +22 pp average gains over vanilla LLM at $k = 3$. While improvement continues from $k = 3$ to $k = 5$, passing a notable threshold of >50% average expected energy savings, this is a smaller jump of +6.4 pp that requires additional overhead.

*Application-Specific Results.* Expected energy reductions by LASSI-EE at $k = 1$ ranged from 1.0% (`keogh` on A100) to 88.8% (`colorwheel` on MI100) and at $k = 3$ from 3.7% (`marchingCubes` on MI100) to 99.5% (`pathfinder` on MI100). Fourteen of 22 applications achieved >20% expected energy reduction at $k = 3$ on both

devices. Applications with <10% reductions at $k = 3$ include `keogh` on A100 (5.9%), `entropy` on A100 (6.5%), `murmurhash3` on MI100 (9.6%), and `marchingCubes` on MI100 (3.7%). LASSI-EE also struggled with `jaccard` on MI100, and while delivering strong optimizations of up to 63.3% energy savings at $k = 5$, underperformed vanilla LLM by 11.3 pp. Following our exemplar `jacobi` app, Appendix B Figure 8 illustrates the clear optimization capability by LASSI-EE over vanilla LLM for the A100 GPU.

# 5. OPTIMIZATION ANALYSIS

LASSI-EE employed a variety of optimizations in its refactored versions that demonstrated improved energy efficiency. We characterize these strategies through a broad categorization of over 4,000 methods extracted from the successful code generations across all 44 application-device trials. These optimization methods were automatically identified during Stage 4 termination (§ 2.6), where the LLM generated a comparison report between the original source code and the final refactored version, explicitly listing the optimizations applied.

We sorted these optimizations into four high-level categories based on their optimization goal: Memory Hierarchy Optimization (MHO), Algorithmic & Computational Efficiency (ACE), Device-Specific Tuning (DST), and Parallelism & Thread Management (PTM). The relatively balanced distribution across these categories (Figure 4 a) demonstrates that LASSI-EE employs diverse optimization strategies rather than over-relying on any single approach, suggesting multiple dimensions of GPU energy efficiency are addressed from memory access patterns to thread organization to computational techniques.

This balanced approach is consistent across both AMD MI100 and NVIDIA A100 platforms (Figure 4 b), with each device showing similar proportional emphasis on all four category types, indicating strategic consistency at the conceptual level.

We also observe strong device specialization. Of the $\sim$3,500 unique optimization implementations identified across both platforms, only 124 (3.5%) were *syntactically identical* in the generated codes for AMD MI100 and NVIDIA A100, while 47% were AMD-specific and 49% were NVIDIA-specific This distribution suggests LASSI-EE appropriately translates optimization *concepts* into platform-specific APIs. For example, pinned host memory allocation appears frequently in both AMD and NVIDIA code but uses `hipHostMalloc` on AMD versus `cudaHostAlloc` on NVIDIA. Similarly, read-only data caching is implemented via compiler optimization hints on AMD but leverages NVIDIA's hardware-specific `__ldg` intrinsic on A100 GPUs. This pattern demonstrates that LASSI-EE adapts conceptual
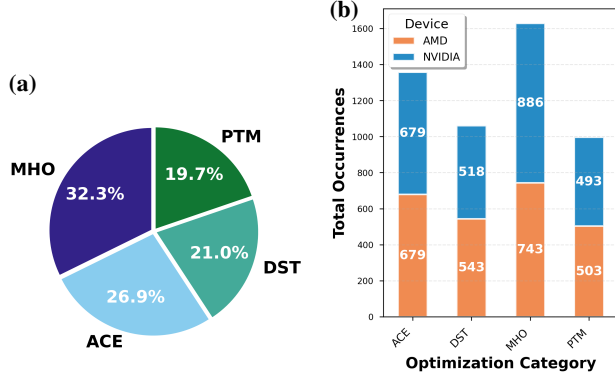
**(a)**

**(b)**



*Figure 4.* Optimization category distribution across all applications: (a) Overall distribution shows balanced utilization across four categories: Memory Hierarchy Optimization (MHO), Algorithmic & Computational Efficiency (ACE), Device-Specific Tuning (DST), and Parallelism & Thread Management (PTM). (b) Cross-device comparison demonstrates consistent categorical emphasis on both AMD MI100 and NVIDIA A100 platforms.

optimization strategies to device-appropriate implementations.

The most frequently applied optimizations across both platforms included variations of pinned (page-locked) host memory allocation (Harris, 2021) (MHO), `__restrict__` pointer annotations (ACE) for aliasing disambiguation, and constant memory placement (MHO) for frequently accessed read-only data (the Top 10 for each device are shown in Appendix D Figure 9. Differences in optimization frequency is observed between devices: generated HIP codes for AMD showed higher utilization of precomputed lookup tables in constant memory (22.7% mean appearance rate), asynchronous device-to-host transfers (14.2%), and HIP-specific event timing with `hipEvent` (8.2%) (Advanced Micro Devices, Inc., 2024), while CUDA code for NVIDIA more frequently employed the platform-exclusive `__ldg` read-only cache intrinsic (16.0%) (NVIDIA Corporation, 2024), loop unrolling pragmas (15.9%), and CUDA-specific pinned memory allocation functions (14.0%). **These patterns demonstrate LASSI-EE's ability to select appropriate optimizations during code refactoring, incorporating both generic techniques and platform-specific API implementations to deliver energy-reducing variants from source code.**

## 6  RELATED WORK

Studies of LLMs and parallel programs started early (Chen et al., 2021; Godoy et al., 2023; Nichols et al.), including a critical evaluation by Nichols et al. (Nichols et al., 2024) finding that while LLMs (up to GPT-4) generated parallel codes with speedup over serial counterparts, they did not improve at maximizing these speedups. This suggested that LLMs can write parallel code but not performant par-

allel code that fully utilizes compute resources. LASSI-EE addresses this gap through iterative refinement guided by real-time runtime energy measurements rather than relying solely on base LLM capabilities.

Recent parallel code generation efforts include CodeRosetta (TehraniJamsaz et al., 2024) for HPC language translation (CUDA, Fortran, C++) using static similarity metrics. Ranasinghe et al.'s (Ranasinghe et al., 2025) work on code translation with FORTRAN → C++ migration and Zhu et al.'s MIRACLE (Zhu et al., 2024) using back-translation for fine-tuning, both evaluate on compilation and static output matching. Schmitz et al.'s (Schmitz et al., 2024) Parallel Pattern Language framework applies compile-time optimizations via a Domain Specific Language.

Muller and Lou (Lou & Muller, 2024) optimized CUDA kernels through static resource analysis (COpPER/RaCUDA), achieving 2-4% gains via *predicted* execution costs, instead of *measured* costs as in LASSI-EE. Finally, the PCEBench (Chen et al., 2025) benchmark evaluates parallel code generation across functional correctness, performance, and OpenMP/MPI support. LASSI-EE's compilation and execution pipeline is well-suited to incorporate PCEBench tasks in future scalable assessments.

Unlike prior studies focused on correctness, static analysis, or performance alone, LASSI-EE uniquely integrates runtime power profiling with iterative LLM-guided refinement to deliver energy-efficient parallel codes with double-digit gains. To the best of our knowledge, this is the first effort leveraging LLMs for automated parallel code refactoring aimed at improving energy efficiency on target systems.

## 7  CONCLUSION

We presented LASSI-EE, an automated LLM-based framework that generated energy-efficient parallel codes through iterative refinement guided by runtime power measurements and tailored to the target execution platform. Integrating system-aware prompting, real-time power profiling, self-correcting feedback loops, and LLM-as-a-Judge code screening, LASSI-EE achieves 29-48% average energy reduction on NVIDIA and AMD GPUs, as much as a 2.8× improvement over vanilla LLM prompting. Our new energy-reduction@$k$ metric enables systematic evaluation of expected energy savings in multi-attempt code generations, and our optimization analysis reveals diverse platform-specific approaches spanning memory hierarchy, algorithmic efficiency, and parallelism management.

We will release LASSI-EE, along with the associated datasets, as **open-source on GitHub** to support reproducibility and foster collaboration within the HPC and AI communities.

# REFERENCES

Abadi, M., Keidar-Barner, S., Pidan, D., and Veksler, T. Verifying parallel code after refactoring using equivalence checking. 47(1):59–73, 2019. ISSN 0885-7458, 1573-7640. doi: 10.1007/s10766-017-0548-4. URL http://link.springer.com/10.1007/s10766-017-0548-4.

Advanced Micro Devices, Inc. HIP Programming Guide. ROCm Documentation, 2024. URL https://rocm.docs.amd.com/projects/HIP/. Accessed: 2025-10-29.

Blom, S. and Huisman, M. The vercors tool for verification of concurrent programs. In *FM*, volume 8442 of *Lecture Notes in Computer Science*, pp. 127–131. Springer, 2014. URL https://link.springer.com/chapter/10.1007/978-3-319-06410-9_9.

Chen, L., Ahmed, N., Capotă, M., Willke, T., Hasabnis, N., and Jannesari, A. PCEBench: A multi-dimensional benchmark for evaluating large language models in parallel code generation. In *2025 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pp. 546–557, 2025. doi: 10.1109/IPDPS64566.2025.00055.

Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. d. O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F. P., Cummings, D., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W. H., Nichol, A., Paino, A., Tezak, N., Tang, J., Babuschkin, I., Balaji, S., Jain, S., Saunders, W., Hesse, C., Carr, A. N., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., and Zaremba, W. Evaluating large language models trained on code, 2021. arXiv:2107.03374.

Dearing, M. T., Tao, Y., Wu, X., Lan, Z., and Taylor, V. Lassi: An llm-based automated self-correcting pipeline for translating parallel scientific codes. In *2024 IEEE International Conference on Cluster Computing Workshops (CLUSTER Workshops)*, pp. 136–143, 2024. doi: 10.1109/CLUSTERWorkshops61563.2024.00029.

Godoy, W., Valero-Lara, P., Teranishi, K., Balaprakash, P., and Vetter, J. Evaluation of OpenAI codex for HPC parallel programming models kernel generation. In *Proceedings of the 52nd International Conference on Parallel Processing Workshops*, pp. 136–144. ACM, 2023. ISBN 979-8-4007-0842-8. doi: 10.1145/3605731.3605886. URL https://dl.acm.org/doi/10.1145/3605731.3605886.

Harris, M. How to optimize data transfers in cuda c/c++, December 2021. URL https://developer.nvidia.com/blog/how-optimize-data-transfers-cuda-cc/.

Jin, Z. and Vetter, J. S. A benchmark suite for improving performance portability of the sycl programming model. In *2023 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, pp. 325–327, 2023. doi: 10.1109/ISPASS57527.2023.00041.

Kate Keahey, e. a. Lessons learned from the chameleon testbed. In *Proceedings of the 2020 USENIX Annual Technical Conference (USENIX ATC '20)*. USENIX Association, July 2020.

Lou, M. and Muller, S. K. Automatic static analysis-guided optimization of CUDA kernels. In *Proceedings of the 15th International Workshop on Programming Models and Applications for Multicores and Manycores*, pp. 11–21. ACM, 2024. ISBN 979-8-4007-0599-1. doi: 10.1145/3649169.3649249. URL https://dl.acm.org/doi/10.1145/3649169.3649249.

Nichols, D., Marathe, A., Menon, H., Gamblin, T., and Bhatele, A. HPC-coder: Modeling parallel programs using large language models. In *ISC High Performance 2024 Research Paper Proceedings (39th International Conference)*, pp. 1–12. IEEE. ISBN 978-3-9826336-0-2. doi: 10.23919/ISC.2024.10528929. URL https://ieeexplore.ieee.org/document/10528929/.

Nichols, D., Davis, J. H., Xie, Z., Rajaram, A., and Bhatele, A. Can large language models write parallel code? In *Proceedings of the 33rd International Symposium on High-Performance Parallel and Distributed Computing*, pp. 281–294. ACM, 2024. ISBN 979-8-4007-0413-0. doi: 10.1145/3625549.3658689. URL https://dl.acm.org/doi/10.1145/3625549.3658689.

NVIDIA Corporation. CUDA C Programming Guide. NVIDIA Developer Documentation, 2024. URL https://docs.nvidia.com/cuda/cuda-c-programming-guide/. Accessed: 2025-10-29.

Ranasinghe, N. R., Jones, S. M., Kucer, M., Biswas, A., O'Malley, D., Most, A., Wanna, S. L., and Sreekumar, A. LLM-assisted translation of legacy FORTRAN codes to c++: A cross-platform study. In Jansen, P., Dalvi Mishra, B., Trivedi, H., Prasad Majumder, B., Hope, T., Khot, T., Downey, D., and Horvitz, E. (eds.), *Proceedings of the 1st Workshop on AI and Scientific Discovery: Directions and Opportunities*, pp. 58–69. Association for Computational Linguistics, 2025. ISBN 979-8-89176-224-

4. doi: 10.18653/v1/2025.aisd-main.6. URL https://aclanthology.org/2025.aisd-main.6/.

Schmitz, A., Miller, J., Burak, S., and Müller, M. S. Parallel pattern language code generation. In *Proceedings of the 15th International Workshop on Programming Models and Applications for Multicores and Manycores*, pp. 32–41. ACM, 2024. ISBN 979-8-4007-0599-1. doi: 10.1145/3649169.3649245. URL https://dl.acm.org/doi/10.1145/3649169.3649245.

Siegel, S. F., Mironova, A., Avrunin, G. S., and Clarke, L. A. Using model checking with symbolic execution to verify parallel numerical programs. In *Proceedings of the ACM SIGSOFT International Symposium on Software Testing and Analysis*, pp. 157–168, 2006.

Siegel, S. F., Zheng, M., Luo, Z., Zirkel, T. K., Marianiello, A. V., Edenhofner, J. G., Dwyer, M. B., and Rogers, M. S. CIVL: the concurrency intermediate verification language. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, SC '15, pp. 1–12. Association for Computing Machinery, 2015. ISBN 978-1-4503-3723-6. doi: 10.1145/2807591.2807635. URL https://doi.org/10.1145/2807591.2807635.

Strohmaier, E., Dongarra, J., Simon, H., and Meuer, M. Top500, June 2025. URL http://www.top500.org/.

TehraniJamsaz, A., Bhattacharjee, A., Chen, L., Ahmed, N. K., Yazdanbakhsh, A., and Jannesari, A. CodeRosetta: Pushing the boundaries of unsupervised code translation for parallel programming. In *NeurIPS*, 2024.

Thavappiragasam, M., Scheinberg, A., Elwasif, W., Hernandez, O., and Sedova, A. Performance portability of molecular docking miniapp on leadership computing platforms. In *2020 IEEE/ACM International Workshop on Performance, Portability and Productivity in HPC (P3HPC)*, pp. 36–44, GA, USA, 2020. IEEE. doi: 10.1109/P3HPC51967.2020.00009.

Tramm, J. R., Siegel, A. R., Islam, T., and Schulz, M. XS-Bench - the development and verification of a performance abstraction for Monte Carlo reactor analysis. In *PHYSOR 2014 - The Role of Reactor Physics toward a Sustainable Future*, Kyoto, 2014. URL https://www.mcs.anl.gov/papers/P5064-0114.pdf.

Zhu, M., Karim, M., Lourentzou, I., and Yao, D. Semi-supervised code translation overcoming the scarcity of parallel code data. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*, pp. 1545–1556. ACM, 2024. ISBN 979-8-4007-1248-7. doi: 10.1145/3691620.3695524. URL https://dl.acm.org/doi/10.1145/3691620.3695524.