# Automated Design Optimization via Strategic Search with Large Language Models

**Anthony Carreon**[*]**, Vansh Sharma, and Venkat Raman**
Department of Aerospace Engineering
University of Michigan
Ann Arbor, MI 48109

December 1, 2025

## Abstract

Traditional optimization methods excel in well-defined search spaces but struggle with design problems where transformations and design parameters are difficult to define. Large language models (LLMs) offer a promising alternative by dynamically interpreting design spaces and leveraging encoded domain knowledge. To this end, we introduce AUTO, an LLM agent framework that treats design optimization as a gradient-free search problem guided by strategic LLM reasoning. The framework employs two collaborative agents: a Strategist that selects between exploration and exploitation strategies, and an Implementor that executes detailed designs. Applied to GPU code optimization – a domain critical to fields from machine learning to scientific computing – AUTO generates solutions competitive with expert implementations for chemical kinetics integration and dense matrix multiplication. The framework achieves 50-70% search efficiency relative to Bayesian optimization methodologies. It completes optimizations in approximately 8 hours at an estimated cost of up to $159 per run, compared to an estimated cost of up to $480 with median-wage software developers. These findings open the door to automating design optimization in ill-defined search spaces with limited prior information.

***Keywords*** Large Language Models · GPU Optimization · Agent Systems · Design Optimization · High-Performance Computing

## 1 Introduction

Traditional optimization methods have demonstrated notable success across several domains, from gradient descent in machine learning, to Bayesian optimization for materials discovery (Zhang et al., 2020) and chemical process design (Li et al., 2021), to evolutionary algorithms for antenna design (Hornby et al., 2006). However, such methods are unsuitable for navigating large, ill-defined search spaces with complex design constraints. Gradient-free methods like evolutionary algorithms and Bayesian optimization excel when the design space can be parameterized, but are incompatible with problems where defining design transformations is difficult or where valid design parametrizations are unclear.

Recent advances in large language models (LLMs) offer a promising alternative by dynamically interpreting design spaces and leveraging domain knowledge encoded in their training data or through semantic search. These advances have led to systems like FunSearch (Romera-Paredes et al., 2023), which discovered new mathematical solutions through evolutionary LLM-driven search, and AlphaEvolve (Novikov et al., 2025), which discovered new solutions to diverse computational problems. Building on these advances, we introduce AUTO, an LLM agent framework that treats design optimization as a strategic search problem. Complementary to hypothesis-driven approaches focused on high-level system design (Hamadanian et al., 2025), AUTO targets implementation-level optimization by exploring the design space. Our framework draws inspiration from evolutionary algorithms and Bayesian optimization by iteratively

---

[*]Corresponding author. Email: `acarreon@umich.edu`

refining designs through (a) improving existing candidates, (b) combining successful patterns, or (c) innovating when existing approaches plateau.

We demonstrate AUTO on GPU code optimization, a domain relevant to fields from machine learning to scientific computing. Recent work has explored LLM-based code optimization through various approaches. Du et al. (2025) demonstrated reinforcement learning methods for competitive programming, while Andrews and Witteveen (2025) targeted scenarios with limited hardware documentation, and Gupta et al. (2025) automated legacy code transformation using multi-agent systems. AUTO extends these efforts by framing optimization as a systematic search problem that aligns with established optimization methods.

The remainder of this paper is organized as follows: Section 2 describes the AUTO framework architecture and optimization workflow. Section 3 presents the GPU programming applications and implementation details. We analyze solution quality, search efficiency, and cost-effectiveness in Section 4, followed by conclusions and future directions in Section 5.

## 2   The LLM-Based Optimizer Framework (AUTO)

AUTO employs two LLMs as agents in a collaborative environment. Inspired by gradient-free optimization techniques such as Bayesian optimization and evolutionary algorithms, the framework iteratively improves designs by learning from a population of previously generated candidates. High-level planning by the Strategist agent is separate from low-level implementation by the Implementor agent. This separation breaks the optimization problem into manageable tasks while circumventing LLM context-window limitations. Figure 1 shows the complete workflow. The optimization loop runs for $N$ iterations, starting with context curation (step 0) and concluding with a performance report (step 5). After each iteration concludes, the chat histories and context windows for the agents are reset. The final output of AUTO is a detailed trace of every design iteration, with the best design having the highest score.
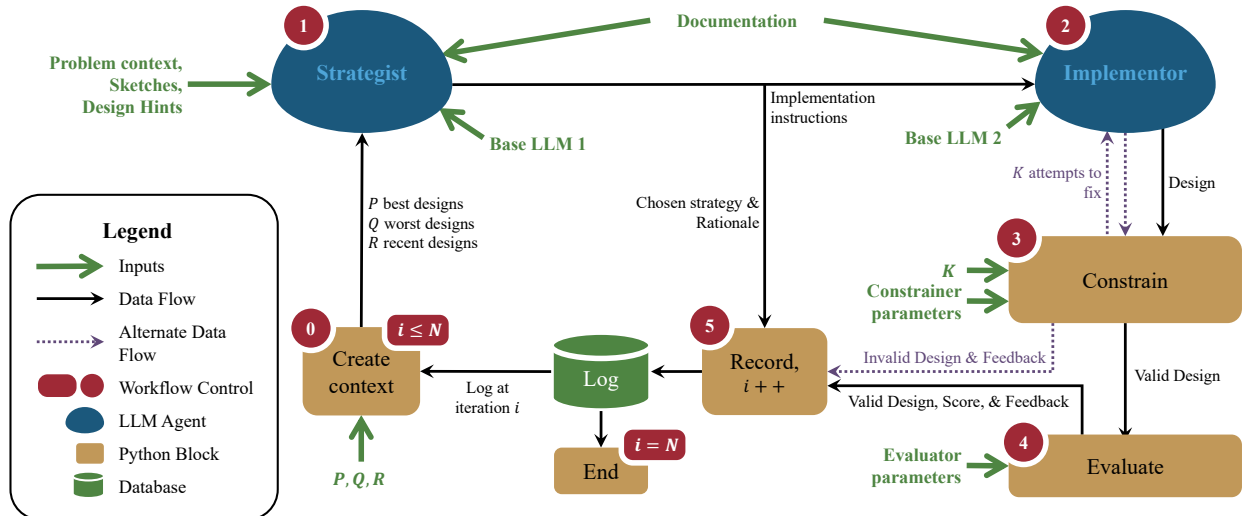


Figure 1: AUTO: An LLM-based design optimization framework. AUTO functions as follows: (0) Create context from historical design data to inform subsequent iterations. (1) The Strategist selects an optimization strategy. (2) The Implementor generates designs from strategic instructions. (3) Validate designs against constraints; (4) Evaluate and score the design; (5) Record the results. The cycle repeats for $N$ iterations. After each iteration concludes, the chat histories and context windows for the Strategist and Implementor are reset.

**Steps 0 and 5: Context Curation from Recorded Data**   Each design iteration is added to a growing database, including its strategies, validation status, detailed evaluation metrics, and agent interactions. This logging enables analysis of the optimization trajectory and is akin to the surrogate model updating in Bayesian optimization. The design space is represented by the set of all designs up to iteration $t$ (exclusive), $\mathcal{D}_t$. This set includes invalid designs (those that failed to meet constraints at step 3) and valid designs. During context curation at step 0 in Figure 1, the goal is to construct a subset of designs, $\mathcal{C}_t \subseteq \mathcal{D}_t$, that when presented to the Strategist, leads to an informed decision. This construction must balance presentation detail and richness with constraints on the context window size and LLM accuracy. We reason that, by composing $\mathcal{C}_t$ with a small mixture of $P$ high-scoring, $Q$ low-scoring, and $R$ recent

designs, the Strategist's decision will yield a design inspired by the top-scoring ones, improved from the low-scoring ones, and/or diversified from the recent ones. In highly nonlinear, non-smooth design surfaces, it is possible that low-scoring designs can be adjusted to achieve significant improvements; hence, tagging along $Q$ of the lowest-scoring designs may be beneficial.

**Step 1: The Strategist**    The Strategist is tasked with analyzing the problem context, sketches, design hints, documentation, and $\mathcal{C}_t$. The problem context refers to background information on the problem being optimized. The sketches may be pseudocodes, similar existing designs, or partial designs related to the application for starting design points. Design hints may include possible actions, design parameters, or guidelines for manipulating and exploring the design space. Next, the Strategist must decide on one of three predefined sampling strategies – refine, combine, or innovate – and provide detailed implementation instructions to the Implementor. Via its system prompt, the Strategist is instructed on when to select each option:

1. "**innovate**" is chosen when all or many implementations exhibit similar designs and scores, when a distinctly different approach could yield breakthrough improvements, or when recent implementations fail to satisfy constraints with no obvious fixes.

2. "**combine**" is selected when multiple implementations excel in different aspects, when some implementations have strong design patterns in one dimension while others excel in a different dimension, when hybrid approaches can capture benefits from multiple designs, or when different design approaches are complementary.

3. "**refine**" is chosen when an implementation is close to reaching a breakthrough score, when specific bottlenecks are identified and fixable, when constraint-related errors are minor and correctable, or when the core design is sound but needs tuning or fixing.

The system prompt additionally instructs the Strategist to operate under several key guidelines to ensure effective collaboration with the Implementor. For instance, the Strategist must provide specific instructions (e.g., "*use 32×32 tiles in shared memory with column-major layout* instead of *implement optimal memory access*"). It should also leverage the optimization hints when applicable. The Strategist must also consider the attainability of constraints and potential improvements in its recommendations. For the first 10 optimization iterations, the Strategist is instructed to prioritize innovating to fully explore and understand the design landscape.

**Steps 2 and 3: The Implementor and Constraints Enforcement**    The Implementor is tasked with implementing the Strategist's instructions and fixing constraint errors along the way. These constraints are domain-specific. For code optimization, this includes compilation, execution, and testing against ground-truth data points. For other design tasks, validation may involve meeting manufacturing constraints or regulatory requirements. As with any optimization algorithm, context on all possible design parameters and their values is needed to find an implementation that meets the design constraints; thus, the Implementor may be provided with domain-specific documentation. If any check fails, the Implementor is given feedback for correction in its persistent context window. After $K$ failed correction attempts, AUTO skips the "evaluate" step and records the most recent output of the failure at Step 5, along with the associated design and its metadata. The Strategist may leverage this design during Step 1 if it is among the $R$ most recent ones selected in Step 0 for future iterations.

**Step 4: Evaluation and Scoring**    Valid designs are evaluated according to problem-specific objectives. Ideally, the evaluation process measures relevant performance metrics across multiple operating conditions. When stochasticity is present, each operation condition should be evaluated multiple times to obtain a statistical measure. The evaluation step transforms all metrics across all operating conditions into a single design score. This score is used during context curation in Step 0; however, the Strategist is presented with full evaluation details for each design in the curated context to support informed decision-making.

## 3    Optimizer Applications

To evaluate AUTO's capabilities, we apply the framework to GPU code optimization. High-performance computing increasingly relies on GPU acceleration, from deep learning (Sharma et al., 2025) to fluid dynamics simulation (Sharma et al., 2024; Carreon et al., 2025). However, optimizing GPU kernels demands specialized expertise in hardware architecture, memory hierarchies, and parallel programming. To this end, we evaluate AUTO on two tasks with differing performance characteristics: chemical kinetics integration and dense matrix multiplication. These applications test AUTO's ability to discover optimizations competitive with expert-tuned implementations.

### 3.1 Chemical Kinetics

We deploy AUTO to optimize the performance of a chemical kinetics GPU code. Specifically, AUTO must optimize a CUDA program that integrates millions of chemical reaction systems in parallel. This problem is relevant to compressible combustion in computational fluid dynamics (CFD) and is a primary computational bottleneck in many combustion applications (Raman et al., 2023). For demonstration and tractability in this preliminary work, we consider the Robertson problem (Robertson, 1966), a three-species, three-reaction chemical system:

$$A \xrightarrow{0.04} B \tag{1}$$

$$B + B \xrightarrow{3 \times 10^7} C + B \tag{2}$$

$$B + C \xrightarrow{1 \times 10^4} A + C \tag{3}$$

This yields the following system of ordinary differential equations:

$$\dot{x} = -0.04x + 10^4 yz \tag{4}$$

$$\dot{y} = 0.04x - 10^4 yz - 3 \times 10^7 y^2 \tag{5}$$

$$\dot{z} = 3 \times 10^7 y^2 \tag{6}$$

The ODE system exhibits timescales spanning several orders of magnitude, leading to highly variable iteration counts during integration across different initial conditions. The LLM's task is to generate an optimal GPU implementation to time-integrate $N$ such systems with different initial conditions, mirroring the variation in thermochemical states across millions of computational cells in combustion CFD[2]. Conditional on the implementation approach, the GPU performance bottlenecks may include, but are not limited to, warp divergence, memory coalescence, and kernel launch overhead.

### 3.2 Dense Square Matrix Multiplication

We additionally deploy AUTO to optimize dense matrix multiplication performance on GPUs. This problem serves as a baseline benchmark and is ubiquitous across scientific computing, machine learning, and numerical methods. To demonstrate this preliminary work, we consider the dense, $N \times N$ matrix multiplication problem, $C = A \times B$, where all matrices are randomly generated with values between 0 and 1 and stored in column-major format in double precision. AUTO's task is to search for an optimal CUDA implementation across a range of problem sizes, from $N = 32$ to $N = 4{,}096$. The inherent data reuse in this operation makes it particularly amenable to GPU optimization through techniques such as tiling, memory hierarchy exploitation, and thread-level parallelism (Volkov and Demmel, 2008). We compare AUTO's solutions against the double-precision general matrix multiplication subroutine (GEMM), dgemm(), in NVIDIA's cuBLAS library (NVIDIA Corporation), which represents years of optimization efforts by teams of human experts and serves as the "ground truth" for dense matrix multiplication on NVIDIA hardware.

### 3.3 Implementation Details

The framework is written in Python using the Ollama client (Marcondes et al., 2025). GPT-OSS-20b (Agarwal et al., 2025) is used as the base LLM for both the Strategist and the Implementor. We ran AUTO with different configurations, as shown in Table 1, to study the impact of temperature-based sampling and the sketch provided to the Strategist. The chemical kinetics problem was optimized using two different sketches. Both sketches are Python codes implementing the same kinetics problem. Python code A performs data access by indexing using a column-major format flat index. In contrast, Python code B uses a 2D subscript for data access with the Pandas module, without any notion of the underlying data storage order. Data layout and access in GPU memory play a major role in the program's performance; therefore, using two different sketches that employ different data-access patterns could reveal different search paths during optimization. For matrix multiplication, we do not provide a sketch because the optimized algorithms are well-established and therefore often present in LLM training data. Figure 2 visualizes Steps 2 through 5 from Figure 1 as it applies to GPU code optimization.

In the "constrain" step, the CUDA code generated by the Implementor is inspected to ensure it compiles without errors (compilation check), executes without runtime errors (execution check), and produces the correct output for a given

---

[2]A note on terminology: we refer to a single ODE system as a "cell" for brevity and relevance to combustion CFD applications.
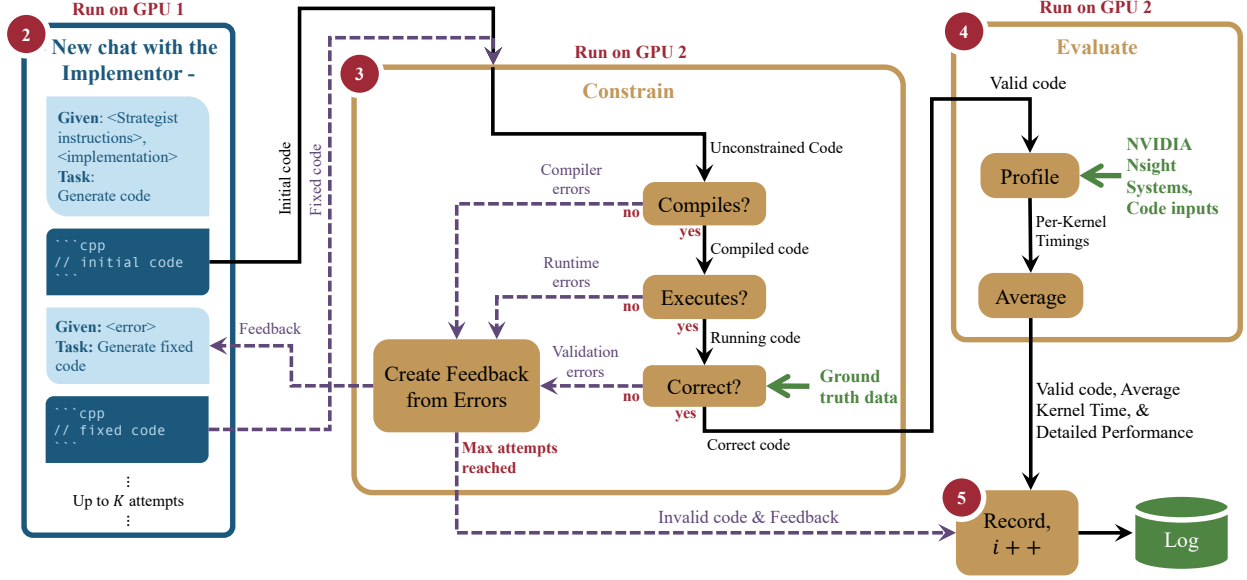
Figure 2: Details of Steps 2 through 5 from Figure 1 as it applies to GPU code optimization. The "Constrain" and "Evaluate" steps are Python blocks of code executed on one GPU, while interactions with the Implementor are LLM chats executed on a separate GPU.

input (correctness check). If any check fails, the Implementor receives feedback for up to $K = 4$ correction attempts in its persistent context window, which is reset after the optimization iteration completes. In the "evaluate" step, the code's GPU kernels are timed with NVIDIA Nsight Systems three times per input and averaged into a single performance value. Note that the detailed, per-input performance data is also retained. The scores are used during context curation in Step 0, while the Strategist uses the detailed design feedback for informed decision-making.

For a given application (kinetics or matrix multiplication), two different, fixed sets of "ground truth" input-output data pairs are generated prior to running AUTO. The first dataset contains input and ground-truth outputs and is used for correctness checking. The second dataset contains only inputs and is used for code profiling. For every generated code, the inputs are provided on the command line as raw values or as file locations. The generated codes are required to output data in a specific file format for compatibility with the constraint and evaluation steps. Notably, we do not provide the agents with access to relevant documentation or external knowledge bases, as the GPT-OSS-20b base model demonstrates sufficient coding proficiency to generate functional implementations solely from the problem context and performance feedback.

The framework runs on two NVIDIA H100 GPUs connected to one Intel Xeon Platinum 8468 CPU. The first GPU hosts the LLM, while the second GPU hosts the code execution environment for compiling, correctness checking, and performance profiling, to maintain exclusivity and avoid interference from the optimization process. The AUTO framework runs on the CPU, controls both GPU processes, and handles I/O.

## 4   Results and Discussion

The discussion and analysis of the results are split into three subsections. The first subsection analyzes AUTO's solutions for matrix multiplication and kinetics, along with relevant artifacts from select optimization runs. Next, we explore AUTO's performance as an optimization framework using the best runs from the first subsection. The final subsection discusses the costs of the optimization framework.

### 4.1   Solution Quality

Figure 3 compares human-optimized GPU code against AUTO's best solutions for (A) Robertson chemical kinetics and (B) square matrix multiplication. The framework generates high-performance code across different problem sizes with near-human quality.

Table 1: AUTO parameters for both optimization tasks.

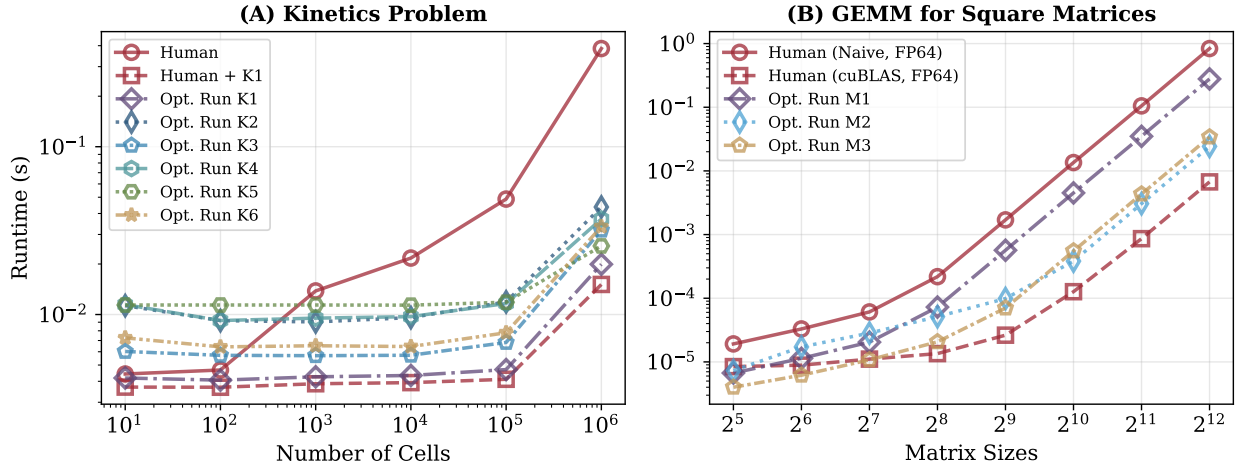| Parameter | Chemical Kinetics | Matrix Multiplication |
|---|---|---|
| Inference temperatures | 0.3, 0.7, and 1.0 | 0.3, 0.7, and 1.0 |
| Sketches | Python code A, Python code B | - |
| Max Iterations | 100 | 100 |
| Strategist Context (see Section 2) | $P = 4, Q = 3, R = 3$ | $P = 5, Q = 5, R = 5$ |
| Correctness Check Data | 100 different initial conditions and their solutions after $10^4$ time steps. | $N \times N$ matrices with values sampled from $x \sim U(0, 1)$ for $N = 10$, $10^3$, and $10^5$. |
| Performance Evaluation | Average GPU kernel runtimes from 10 to $10^6$ initital conditions out to $10^4$ time steps. | Average GPU kernel runtimes from $N = 32$ to $N = 4096$. |
| **Total optimization runs** | 6 (3 temperatures $\times$ 2 sketches) | 3 (3 temperatures $\times$ 1 sketch) |



Figure 3: Timing comparisons between the agent-optimized and human-optimized GPU code for (A) the kinetics problem by Robertson (1966) across different problem sizes and (B) matrix multiplication applied to two $N \times N$ matrices. See Table 2 for the parameters of each run.

In the kinetics problem, run K1 achieves the best solution across all optimization runs. From Table 2, this corresponds to a run with inference temperature of 0.3 and Python code A, which completed in 23 iterations. The resulting code from this iteration uses single-precision operations for all operations up to the final state and time updates of the chemical system. Operations are inlined, and loops are manually unrolled wherever possible. The agents further attempt to use intrinsic `fma`, `min`, `max`, and `load` instructions to achieve speedup. The largest difference between the first "Human" code and the solution from K1 is the use of register memory. After reviewing the solution code from K1, we updated the "Human" code to improve register usage and achieved better performance ("Human + K1"). It is expected that by introducing more of the optimization from code K1, further improvement would be achieved.

Runs M2 and M3 achieved the best matrix multiplication solutions; however, run M2 found the best solution after 1 iteration. Once the best solution is found on the first try, AUTO struggles to find anything better. The best solutions employ loop tiling, a well-established strategy to improve the performance of matrix multiplication. However, when AUTO attempts to optimize further using specialized libraries such as WMMA (Markidis et al., 2018), it fails to compile the code. For instance, in one iteration of M3, the Implementor incorrectly used WMMA's `load_matrix_sync` function. Compilation errors were not limited to matrix multiplication. After inspecting the optimization logs, it was found that the Implementor achieved 56.94% compilation success for matrix multiplication and 44.87% for kinetics. This is a result of hallucination and a lack of documentation of the CUDA API. WMMA (Warp Matrix Multiply Accumulate) provides a warp-level API for programming NVIDIA Tensor Cores to perform efficient mixed-precision

Table 2: Experimental results for different applications and parameter settings

| Application | Run ID | Inference Temperature | Sketch | Iters. to Best |
|---|---|---|---|---|
| **Kinetics Problem** | K1 | 0.3 | Python Code A | 23 |
| | K2 | | Python Code B | 1 |
| | K3 | 0.7 | Python Code A | 30 |
| | K4 | | Python Code B | 15 |
| | K5 | 1.0 | Python Code A | 86 |
| | K6 | | Python Code B | 82 |
| **Matrix Multiplication** | M1 | 0.3 | - | 28 |
| | M2 | 0.7 | - | 1 |
| | M3 | 1.0 | - | 68 |

matrix multiplication on GPUs. By providing AUTO with documentation for the WMMA API, we expect further improvements in the solutions generated and alignment with the performance of the cuBLAS-based GEMM subroutines.

Figure 4 shows a t-SNE visualization of the code vectors for run K1 from the kinetics problem and M3 from the matrix multiplication problem. To analyze the structural relationships among generated codes, we employed a bag-of-words approach to construct code vectors rather than using embedding models. This decision was made after preliminary experiments (not shown here) on similarity search using embedding models, as reported by (Sharma and Raman, 2024). Our preliminary findings highlighted the need to capture algorithmic and implementation-level details that may be obscured in high-dimensional similarity searches, as also noted by Galke and Scherp (2022). For each code, keywords and syntax were extracted with the Treesitter library (Brunsfeld and Tree-sitter contributors, 2024) to store position-aware data. The vectors were then formed based on the frequencies of keywords in each code. To identify stable code clusters, we applied consensus clustering using the evidence accumulation approach (Fred and Jain, 2005). We ran K-means clustering multiple times with varying cluster counts and constructed a co-occurrence matrix that tracks how frequently code vector pairs cluster together. Clusters were then formed by merging codes with co-occurrence rates exceeding 10%, yielding eight disjoint clusters for both applications. To validate the meaningfulness of the resulting clusters, we manually inspected random codes from each cluster. We confirmed that intra-cluster codes shared common implementation strategies, whereas inter-cluster codes exhibited different implementation characteristics. t-SNE embeddings were produced with perplexity values of 10 for the kinetics application and 15 for the matrix multiplication application. These perplexity settings were chosen by visually inspecting t-SNE plots at different perplexities to achieve close alignment with the clusters obtained via consensus K-means clustering.

The t-SNE visualizations reveal strong correlations between code structure and performance. For kinetics, Cluster A contains top performers, including the best solution (iteration 23). These codes used mixed-precision computation, double-precision accumulation, `__ldg` for read-only cache loads, and grid-stride loops. Cluster B explores synchronization strategies such as warp-level work stealing, which likely provided no benefits and incurred some overhead. Cluster C represents early explorations with pure double-precision and simpler approaches that achieved correctness without later optimizations. For matrix multiplication, optimal solutions span multiple clusters (B, C, D, F), with the best in Cluster C (iteration 68). Cluster B implements naive element-wise approaches. Cluster C contains tiled implementations and different padding strategies to avoid bank conflicts. Cluster D also shows tiling strategies without padding, while Cluster F explores kernels with loop unrolling. This diversity demonstrates that matrix multiplication admits multiple optimization pathways. Performance improvements followed irregular trajectories in both applications, which is expected. The search space is highly nonlinear, where small code changes could produce large performance variations.

## 4.2 Search Efficiency, Convergence, and Stability

This section aims to explore the Strategist's ability to choose a sensible strategy based on their observations of how different designs are performing. Given that a sensible strategy was selected, we also assess the Strategist's ability to communicate clear design instructions to the Implementor and the Implementor's ability to follow them. For a quantitative assessment of search efficiency, we first draw analogies with the well-established sampling techniques used in Bayesian optimization (Kushner, 1964). As with the AUTO framework, every iteration in Bayesian optimization begins with a set of existing points in the search space and their evaluations with respect to a predefined objective function (Frazier, 2018). Selection of the next point to sample and evaluate is based on an acquisition function that balances exploration and exploitation. Similarly, AUTO's Strategist must balance exploring novel design regions against exploiting promising areas. Thus, a measure of alignment between AUTO's and Bayesian optimization's sampling

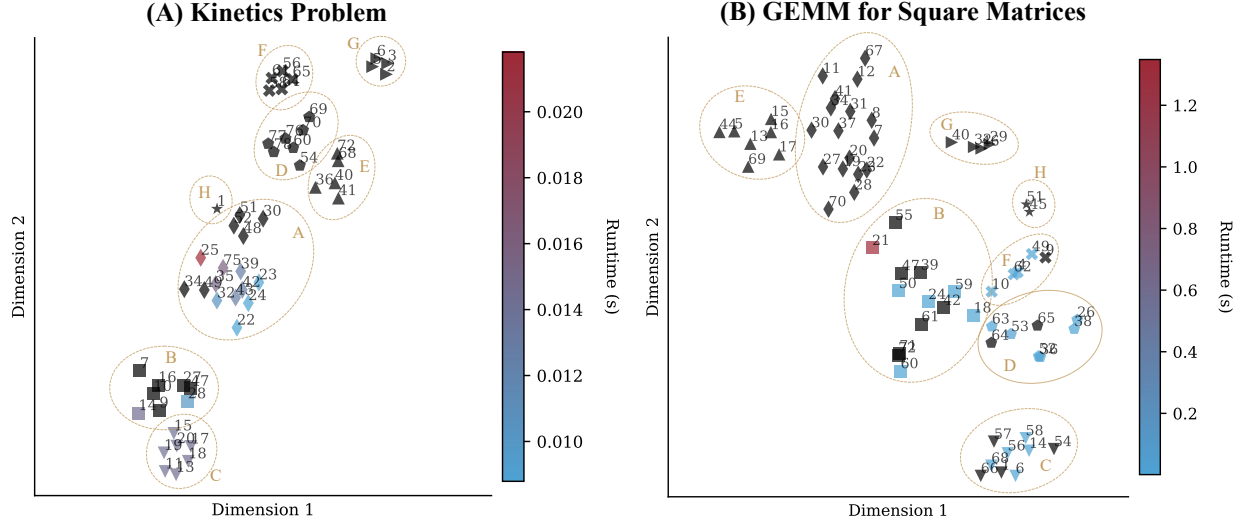**(A) Kinetics Problem**

**(B) GEMM for Square Matrices**



Figure 4: Code clusters and their correlations with the code runtimes for (A) kinetics run K1 and (B) matrix multiplication run M3. For each application, a "bag of words" approach produces the code vectors from the generated codes, which are then clustered using K-means. Visual embeddings are obtained using t-SNE with perplexity=10 for (A) and perplexity=15 for (B). The colors represent runtimes. The markers represent cluster memberships. The numbers next to each marker are the iteration at which the code was generated.

methods provides a rough estimate of search efficiency. To quantify these alignments, the first step is to classify the Strategist's decision at a given iteration, $t$, as exploratory or exploitive:

$$y_S^{(t)} = \begin{cases} \text{exploitation} & \text{if } \varsigma_t \in \{\text{refine, combine}\} \\ \text{exploration} & \text{if } \varsigma_t = \text{innovate} \end{cases} \tag{7}$$

$\varsigma_t$ is the Strategist's decision. The "refine" and "combine" strategies are considered exploitative, as they build on known design scores, whereas "innovate" is exploratory, as the Strategist ventures into uncharted regions of the design space. To compare the Strategist's choices against those of Bayesian optimization and to assess the Implementor's execution, we define the following mapping, $\Phi : \mathbb{R}^d \to \{\text{exploitation, exploration}\}$, for a consistent classification of any design point as exploratory or exploitative:

$$\Phi(\mathbf{x}; \mathcal{C}) = \begin{cases} \text{exploitation} & \text{if } D_{\min}(\mathbf{x}, \mathcal{C}) \leq 0.1 \times L(\mathcal{C}) \\ & \text{or } \mathbf{x} \in \text{ConvexHull}(\mathcal{C}) \\ \text{exploration} & \text{otherwise} \end{cases} \tag{8}$$

where $D_{\min}(\mathbf{x}, \mathcal{C})$ is the minimum distance from point $\mathbf{x} \in \mathbb{R}^d$ to any point in $\mathcal{C}$, and $L(\mathcal{C})$ is the max of max norms, $||\cdot||_\infty$, across all design points in $\mathcal{C}_t$. Intuitively, a sample point is exploitive if it is close to any design point in $\mathcal{C}$, or if it falls within the convex hull of $\mathcal{C}$. Otherwise, the point is considered exploratory. Using this classifier, we can determine the effective sample strategy of both the Bayesian optimizer and the Implementor:

$$y_B^{(t)} = \Phi(\mathbf{x}_{BO}; \mathcal{C}_t) \tag{9}$$

$$y_I^{(t)} = \Phi(\mathbf{x}_I; \mathcal{C}_t) \tag{10}$$

where $\mathbf{x}_{BO}(t)$ is the sample point selected by Bayesian optimization and $\mathbf{x}_I(t)$ is the design point produced by the Implementor. Search efficiency is quantified via two binary metrics:

$$A_{S,B}^{(t)} = \begin{cases} 1 & \text{if } y_S^{(t)} = y_B^{(t)} \\ 0 & \text{otherwise} \end{cases} \tag{11}$$

8

$$A_{S,I}^{(t)} = \begin{cases} 1 & \text{if } y_S^{(t)} = y_I^{(t)} \\ 0 & \text{otherwise} \end{cases} \tag{12}$$

Equation 11 measures whether the Strategist's sampling strategy matches that of a traditional Bayesian optimizer, while Equation 12 indicates whether the Implementor successfully executed the Strategist's intended strategy. Together, these metrics provide quantitative insight into the efficiency of the strategy-to-implementation pipeline. The final search efficiency marker is calculated as:

$$\text{search efficiency} = \left( \frac{1}{2n} \sum_t A_{S,B}^{(t)} + A_{S,I}^{(t)} \right) \times 100, \tag{13}$$

where $n$ is the number of iterations considered in the sum. This sum excludes optimization iterations that did not produce any code due to parsing errors in the LLM output or LLM inference timeouts. To obtain $\mathbf{x}_{BO}$, the upper confidence bound (UCB) acquisition function is used within the Bayesian Optimization Python library by Nogueira (2014). For details on the UCB function and its initial formulation, the reader is referred to Srinivas et al. (2010). The UCB function is parameterized by an exploration factor, $\xi$, which controls exploitation (low $\xi$) and exploration (high $\xi$).

Figure 5 shows how the search efficiency varies with the exploration factor for the best optimization run per application. The graph shows search efficiency values ranging from $\sim 50\%$ to $\sim 70\%$ across exploration factors spanning four orders of magnitude. This indicates that, for the majority of optimization iterations, the Strategist's search strategies are aligned with those of the Bayesian optimizer, with the Implementor, or with both. The higher search efficiency obtained by the kinetics application at $\xi = 10$ indicates that AUTO exhibits more exploratory behavior. This is consistent with the Strategist's decisions being mostly to "innovate", as reported in Table 3. On the other hand, AUTO exhibits more exploitative behavior at $\xi = 0.01$ when optimizing matrix multiplication, which correlates with the Strategist's decisions being mostly to "combine" codes, as seen in Table 3.

Table 3: Strategist decisions distribution for each application.

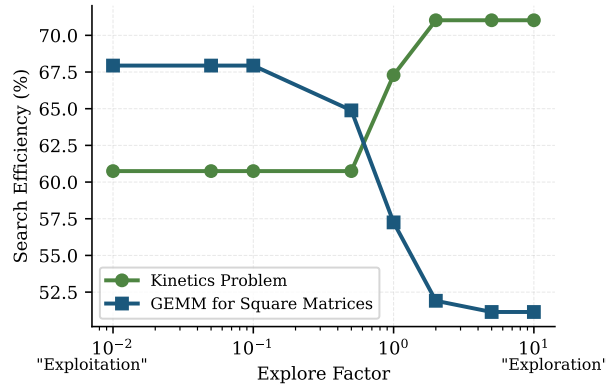| Strategist Decision | Kinetics Run K1 | | matrix multiplication Run M3 | |
| --- | --- | --- | --- | --- |
| | Count | Percentage | Count | Percentage |
| Combine | 10 | 12.82% | 42 | 58.33% |
| Innovate | 45 | 57.69% | 16 | 22.22% |
| Refine | 3 | 3.85% | 9 | 12.50% |
| N/A | 20 | 25.64% | 5 | 6.94% |
| **Total** | **78** | **100.00%** | **72** | **100.00%** |



Figure 5: Search efficiency (calculated from Equation 13) as a function of the exploration factor, $\xi$, from the upper confidence bound (UCB) acquisition function for kinetics run K1 (circles) and matrix multiplication run M3 (squares).

Figure 6 shows convergence patterns for both applications at runs K1 and M3 from Table 2. The top row shows the relative distance between code vectors across successive iterations. The relative distance is computed as:

$$\text{relative distance at } t \equiv \frac{||v_t - v_{t-1}||_2}{||v_{t-1}||_2}, \tag{14}$$

where $v_t$ is the code vector at iteration $t$ and $||\cdot||$ indicates the L2 norm. A large spike indicates significant code changes and possible exploratory behavior, while small values suggest minor changes and possible exploitive behavior. The bottom row shows the improvement (as a percentage relative to the best overall runtime) of the best runtime so far up to iteration $t$. This value is computed as:

$$\text{best solution at } t \equiv \left( 2 - \frac{\min_{i \leq t} r^{(i)}}{\min_{i \leq n} r^{(i)}} \right) \times 100 \tag{15}$$
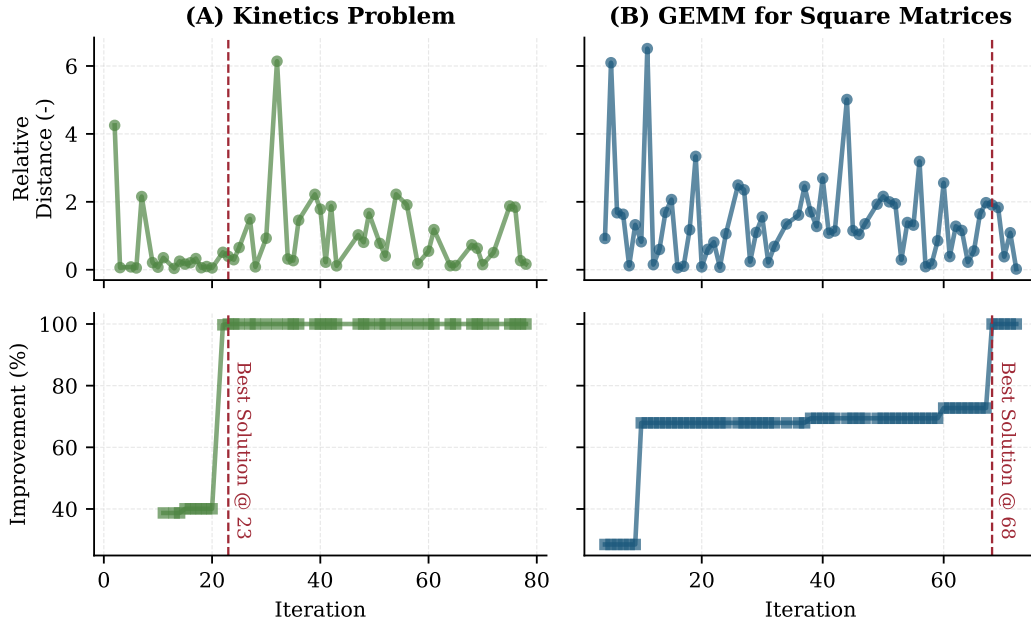


Figure 6: Relative distance (see Equation 14) between code vectors of successive iterations (top row) and code runtime improvement as a percentage (bottom row) versus each iteration for (A) kinetics run K1 and (B) matrix multiplication run M3.

In the kinetics problem, the relative distance shows low variability as it approaches the best solution at 23 iterations, suggesting exploitative behavior as a local optimum is reached. After 23 iterations, there is significantly greater variability, indicating that AUTO continues to explore the design space for other optima. The matrix multiplication application shows more frequent variations. Combined with observations from Figure 5, these seemingly large variations could be a result of jumping from one known similar set of design points to a completely different safe set of design points.

In the improvement subplots of Figure 6, initial codes may not have passed compiler and correctness checks, and therefore some values are missing. The kinetics problem shows a sharp performance jump at iteration 23, reaching its optimal solution for the run. The matrix multiplication optimization exhibits a steadier convergence pattern, reaching its optimum by iteration 68. Both applications show that breakthroughs in optimization can occur after extended periods of stagnation. However, this behavior poses challenges in determining the optimal final iteration, as premature termination can lead to suboptimal solutions. However, note that in this preliminary study, no CUDA documentation was provided, potentially resulting in wasted code generations (and LLM tokens) due to very minor, fixable syntax issues. With access to a knowledge base, stagnant performance segments are likely to lessen.

### 4.3   Token Usage and Cost Analysis

In this section, we demonstrate the economic viability of the AUTO framework by using optimization run K1 from Table 2 as an example. Run K1 contains 53 out of 78 iterations with successfully generated code, with the smallest code having 228 lines (iteration 23, coincidentally the best solution) and the largest having 393 lines (iteration 9). The average code size is $\sim 285 \pm 31$ lines, which indicates relatively consistent code generation. Table 4 reports the token usage per iteration, based on the GPT-OSS-20b tokenizer, and the context usage is based on the max context size of 128k tokens. The Implementor consumed an average of 10,054 input and 8,601 output tokens per iteration, while the Strategist processed significantly larger contexts with an average of 42,632 input and 3,684 output tokens per iteration. This is expected given the amount of code that the Strategist must analyze when provided the curated context (see Section 2). The average "per line of code" token usage was computed by dividing the total tokens by the total lines across all generated code.

Table 4: Token statistics (tokens/iter) from run K1 (see Table 2)

| Token Type | Min | Max | Avg. | Std (%) | Context Usage (%) |
|---|---|---|---|---|---|
| Implementor (input) | 4,085 | 25,135 | 10,054 | 67 | 7.85 |
| Implementor (output) | 2,543 | 20,461 | 8,601 | 63 | 6.72 |
| Strategist (input) | 2,998 | 60,212 | 42,632 | 38 | 33.31 |
| Strategist (output) | 2,135 | 8,921 | 3,684 | 38 | 2.88 |
| Per Line of Code | 9 | 15 | 11 | 11 | 0.01 |

Figure 7 presents token estimated costs per optimization iteration across various commercial LLM providers. These estimates are based on the average input and output tokens per iteration from Table 4 and pricing data from OpenAI and Anthropic in October 2025. Note that tokenization may vary across commercial LLMs, potentially affecting the estimated costs. For example, Claude models typically require 16-30% more tokens than OpenAI models for the same content (Gupta, 2024), with code requiring the highest overhead. For the complete K1 optimization run (78 total iterations), the total cost ranges from approximately $0.78 for GPT-5-nano to $159.12 for o3-pro, with mid-tier models like GPT-5 and Claude Sonnet-4.5 costing $14.82 and $26.52, respectively. These mid-tier models are likely a good balance between affordability and accuracy for the AUTO framework.

According to The U.S. Bureau of Labor Statistics (2024), the median hourly wage for software developers was $63.98 in May 2024, with the bottom 10% earning $38.39 per hour. To estimate human optimization costs, we account for the iterative nature of performance optimization. Unlike AUTO, we assume that human developers make steadier, more focused modifications of $\sim$10 lines per iteration, with each iteration requiring $\sim$45 minutes for profiling, modification, compilation, and testing. Expert developers may achieve optimal solutions within 10 iterations, compared to AUTO's 78 total iterations (23 to the best solution). Furthermore, AUTO's broader exploration included 44-57% compilation failures, which human developers typically avoid through domain expertise and modern AI pair programming tools like GitHub Copilot (Peng et al., 2023).

Table 5 compares optimization costs. For human developers, we estimate 0.75 hours per iteration at median wage ($48 per iteration), with 10 iterations to reach competitive performance ($480 total). Entry-level developers may require more iterations at lower hourly rates, yielding similar total costs. AUTO with mid-tier models (GPT-5, Sonnet-4.5) costs $15-$27 versus $435-$480 for human optimization, which is at least a $16\times$ cost reduction. However, this comparison is nuanced, and a deeper analysis is needed as the model version affects the number of iterations required to reach the optimal solution. AUTO explored 53 codes in 78 iterations and 8 hours, while human developers may devote 12 hours over several days with deep expertise. AUTO's advantage lies in systematic exploration without requiring GPU-optimization expertise, though it incurs costs due to compilation failures. For organizations seeking GPU optimization but lacking GPU programming talent, AUTO's capabilities and potential cost savings may be a viable option. On the other hand, organizations with deep GPU expertise could augment their capabilities with AUTO's.

## 5   Conclusions

We introduced AUTO, an LLM-based framework for automated design optimization that treats the process as a gradient-free search problem. The framework separates strategic planning from implementation through a Strategist-Implementor agent architecture, managing context limitations while learning from historical design performance. To demonstrate the framework's capabilities, we applied AUTO to two GPU programming tasks: chemical kinetics integration and

Table 5: Cost comparison estimates for GPU code optimization

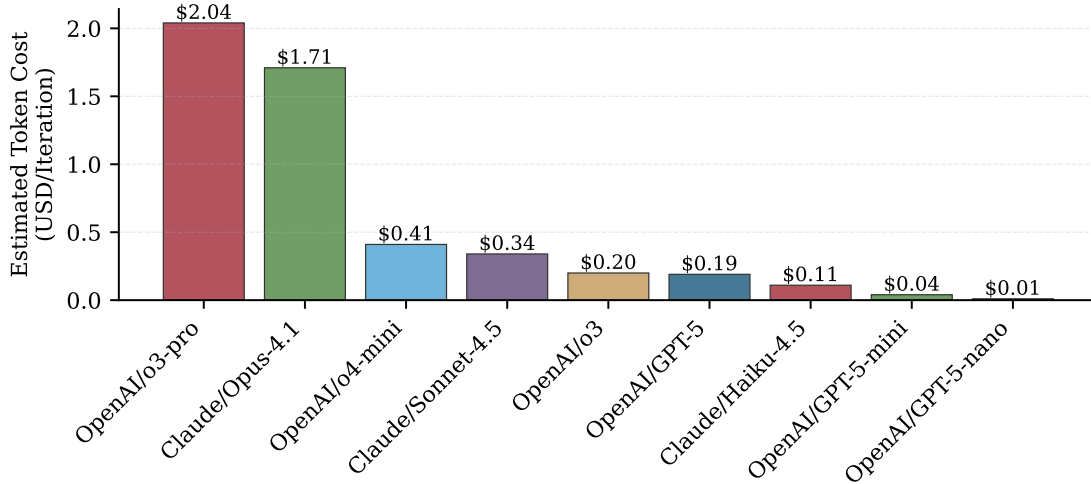| Approach | Cost per Iteration | Total Cost |
|---|---|---|
| *LLM-based (AUTO)* - see Figure 7 | | |
| GPT-5 | $0.19 | $14.82 (78 iter.) |
| Sonnet-4.5 | $0.34 | $26.52 (78 iter.) |
| Opus-4.1 | $1.71 | $133.38 (78 iter.) |
| o3-pro | $2.04 | $159.12 (78 iter.) |
| *Human Developer* - see discussion in Section 4.3 | | |
| Entry-level ($38.39/hr) | $29 | $435 (15 iter.) |
| Median ($63.98/hr) | $48 | $480 (10 iter.) |



Figure 7: Estimated costs per AUTO iteration across various commercial LLM providers

dense matrix multiplication. On these demonstration problems, the AUTO framework generated solutions competitive with human-optimized implementations, achieving search efficiency of 50-70% when compared against Bayesian optimization baselines. The kinetics optimization matched and, in some cases, exceeded expert performance through interesting optimization strategies. On the other hand, the matrix multiplication application implemented practical approaches but lacked familiarity with niche tensor-core libraries, hindering the Implementor's ability to realize working designs.

The framework completed optimizations in approximately 8 hours at costs ranging from $0.78 to $159.12, depending on the LLM provider, and compared favorably with estimated manual optimization costs of $435-$480. However, the observed weak correlation between the framework's hyperparameter and search convergence reveals challenges in navigating discrete, highly nonlinear design landscapes with ill-defined search spaces. Future work should scale AUTO to larger, more complex tasks and expand to other design domains beyond code optimization, including hardware-software co-design and multi-objective engineering problems. This expansion introduces challenges that may require (a) surrogate models for design quality prediction, (b) improving the stopping criteria to balance monetary costs and improvement potential, and (c) integrating domain-specific knowledge bases to reduce constraint violations (e.g., compilation errors in the matrix multiplication problem). Automating design optimization in ill-defined search spaces with limited information remains an open challenge. However, AUTO's ability to explore diverse optimization strategies opens avenues towards automated scientific and engineering design and discovery.

## Acknowledgments

## References

Yichi Zhang, Daniel W. Apley, and Wei Chen. Bayesian optimization for materials design with mixed quantitative and qualitative variables. *Scientific Reports*, 10(1):4924, 2020.

Yifan Li, Phillip M. Maffettone, Dionisios G. Vlachos, and Stavros Caratzoulas. Nextorch: A design and bayesian optimization toolkit for chemical sciences and engineering. *Journal of Chemical Information and Modeling*, 61(11): 5312–5319, 2021.

Gregory Hornby, Al Globus, Derek Linden, and Jason Lohn. Automated antenna design with evolutionary algorithms. In *Space 2006*, page 7242. 2006.

Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog, M. Pawan Kumar, Emilien Dupont, Francisco J. R. Ruiz, Jordan Ellenberg, Pengming Wang, Omar Fawzi, Pushmeet Kohli, and Alhussein Fawzi. Mathematical discoveries from program search with large language models. *Nature*, 2023. doi: 10.1038/s41586-023-06924-6.

Alexander Novikov, Ngân Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco JR Ruiz, Abbas Mehrabian, et al. Alphaevolve: A coding agent for scientific and algorithmic discovery. *arXiv preprint arXiv:2506.13131*, 2025.

Pouya Hamadanian, Pantea Karimi, Arash Nasr-Esfahany, Kimia Noorbakhsh, Joseph Chandler, Ali ParandehGheibi, Mohammad Alizadeh, and Hari Balakrishnan. Glia: A human-inspired ai for automated systems design and optimization. *arXiv preprint arXiv:2510.27176*, 2025.

Mingzhe Du, Luu Anh Tuan, Yue Liu, Yuhao Qing, Dong Huang, Xinyi He, Qian Liu, Zejun Ma, and See-kiong Ng. Afterburner: Reinforcement learning facilitates self-improving code efficiency optimization. *arXiv preprint arXiv:2505.23387*, 2025.

Martin Andrews and Sam Witteveen. Gpu kernel scientist: An llm-driven framework for iterative kernel optimization. *arXiv preprint arXiv:2506.20807*, 2025.

Sparsh Gupta, Kamalavasan Kamalakkannan, Maxim Moraru, Galen Shipman, and Patrick Diehl. From legacy fortran to portable kokkos: An autonomous agentic ai workflow. *arXiv preprint arXiv:2509.12443*, 2025.

Vansh Sharma, Andreas H Rauch, and Venkatramanan Raman. Accelerating cfd simulations with super-resolution feedback-informed adaptive mesh refinement. In *AIAA SCITECH 2025 Forum*, page 1467, 2025.

Shivank Sharma, Ral Bielawski, Oliver Gibson, Shuzhi Zhang, Vansh Sharma, Andreas H Rauch, Jagmohan Singh, Sebastian Abisleiman, Michael Ullman, Shivam Barwey, et al. An amrex-based compressible reacting flow solver for high-speed reacting flows relevant to hypersonic propulsion. *arXiv preprint arXiv:2412.00900*, 2024.

Anthony Carreon, Jagmohan Singh, Shivank Sharma, Shuzhi Zhang, and Venkat Raman. A gpu-based compressible combustion solver for applications exhibiting disparate space and time scales. *arXiv preprint arXiv:2510.23993*, 2025.

Venkat Raman, Supraj Prakash, and Mirko Gamba. Nonidealities in rotating detonation engines. *Annual Review of Fluid Mechanics*, 55(1):639–674, 2023.

H. H. Robertson. The solution of a set of reaction rate equations. In *Numerical Analysis: An Introduction*. 1966.

Vasily Volkov and James W Demmel. Benchmarking gpus to tune dense linear algebra. In *SC'08: Proceedings of the 2008 ACM/IEEE conference on Supercomputing*, pages 1–11. IEEE, 2008.

NVIDIA Corporation. cuBLAS Library. https://developer.nvidia.com/cublas. Accessed: 2025.

Francisco S Marcondes, Adelino Gala, Renata Magalhães, Fernando Perez de Britto, Dalila Durães, and Paulo Novais. Using ollama. In *Natural Language Analytics with Generative Large-Language Models: A Practical Approach with Ollama and Open-Source LLMs*, pages 23–35. Springer, 2025.

Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K Arora, Yu Bai, Bowen Baker, Haiming Bao, et al. gpt-oss-120b & gpt-oss-20b model card. *arXiv preprint arXiv:2508.10925*, 2025.

Stefano Markidis, Steven Wei Der Chien, Erwin Laure, Ivy Bo Peng, and Jeffrey S. Vetter. NVIDIA tensor core programmability, performance & precision. In *2018 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, pages 522–531. IEEE, 2018. doi: 10.1109/IPDPSW.2018.00091.

Vansh Sharma and Venkat Raman. A reliable knowledge processing framework for combustion science using foundation models. *Energy and AI*, 16:100365, 2024.

Lukas Galke and Ansgar Scherp. Bag-of-words vs. graph vs. sequence in text classification: Questioning the necessity of text-graphs and the surprising strength of a wide mlp, 2022. URL https://arxiv.org/abs/2109.03777.

Max Brunsfeld and Tree-sitter contributors. Tree-sitter: An incremental parsing system for programming tools. https://github.com/tree-sitter/tree-sitter, 2024. Accessed: 2024-11-24.

Ana LN Fred and Anil K Jain. Combining multiple clusterings using evidence accumulation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(6):835–850, 2005.

Harold J Kushner. A new method of locating the maximum point of an arbitrary multipeak curve in the presence of noise. 1964.

Peter I Frazier. A tutorial on bayesian optimization. *arXiv preprint arXiv:1807.02811*, 2018.

Fernando Nogueira. Bayesian Optimization: Open source constrained global optimization tool for Python, 2014. URL https://github.com/bayesian-optimization/BayesianOptimization.

Niranjan Srinivas, Andreas Krause, Sham M Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In *International Conference on Machine Learning (ICML)*, pages 1015–1022, 2010.

Lavanya Gupta. Hidden costs in AI deployment: Why Claude models may be 20-30% more expensive than GPT in enterprise settings, August 2024.

The U.S. Bureau of Labor Statistics. Software developers, quality assurance analysts, and testers. Occupational Outlook Handbook, May 2024. URL https://www.bls.gov/ooh/computer-and-information-technology/software-developers.htm. Accessed November 25, 2025.

Sida Peng, Eirini Kalliamvakou, Peter Cihon, and Mert Demirer. The impact of ai on developer productivity: Evidence from github copilot. *arXiv preprint arXiv:2302.06590*, 2023.