

Diagonal-Budgeted Trotterization for Efficient Quantum Hamiltonian Simulation

Srikar Chundury
schundu3@ncsu.edu
North Carolina State University
Raleigh, NC, USA

Blake Burgstahler
bburgst@ncsu.edu
North Carolina State University
Raleigh, NC, USA

Jiajia Li
jiajia.li@ncsu.edu
North Carolina State University
Raleigh, NC, USA

In-Saeng Suh
suh@ornl.gov
Oak Ridge National Laboratory
Oak Ridge, TN, USA

Frank Mueller
fmuelle@ncsu.edu
North Carolina State University
Raleigh, NC, USA

ABSTRACT

Efficient classical simulation of quantum Hamiltonian dynamics is often bottlenecked by exponential state growth and the overhead of generic sparse linear algebra. We introduce diagonal-budgeted Trotterization, a structure-aware strategy that decomposes Hamiltonians into factors preserving diagonal sparsity while tightly controlling fidelity loss.

Our implementation, HAMSIM, utilizes a compact diagonal-sparse data layout and specialized C++/CUDA kernels to bypass the overheads of generic formats like CSR. By leveraging SIMD vectorization, multithreading, and GPU acceleration, HAMSIM achieves high performance across heterogeneous architectures. Benchmarks on the HamLib suite show that HAMSIM significantly outperforms Qiskit-Aer. On CPUs, HAMSIM attains speedups of 182–1,269× on optimization instances (TSP, MaxCut) and 4.8–841× on physical models (TFIM, Heisenberg). On GPUs, it achieves up to 178× speedup for 12–16 qubit problems.

Unlike traditional Trotterization, HAMSIM maintains near-perfect fidelity without requiring exponential steps. This demonstrates that diagonal-aware numerical kernels provide a scalable foundation for high-fidelity classical Hamiltonian simulation.

CCS CONCEPTS

• **Computing methodologies** → **Quantum mechanic simulation; Simulation evaluation; Parallel algorithms.**

KEYWORDS

Quantum computing, Hamiltonian simulation, Sparse linear algebra, Diagonal sparsity, Trotterization

ACM Reference Format:

Srikar Chundury, Blake Burgstahler, Jiajia Li, In-Saeng Suh, and Frank Mueller. 2026. Diagonal-Budgeted Trotterization for Efficient Quantum Hamiltonian Simulation. In *2026 International Conference on Supercomputing (ICS '26)*, July 06–09, 2026, Belfast, United Kingdom. ACM, New York, NY, USA, 16 pages. <https://doi.org/10.1145/3797905.3807869>



This work is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License.

ICS '26, July 06–09, 2026, Belfast, United Kingdom

© 2026 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-2522-7/2026/07.

<https://doi.org/10.1145/3797905.3807869>

1 INTRODUCTION

Quantum computing has seen rapid progress recently, driven by foundational algorithms such as Shor’s factoring [37], Grover’s search [15], and HHL [17], as well as variational and hybrid approaches, including the variational quantum eigensolver (VQE) [34], the quantum approximate optimization algorithm (QAOA) [12, 14], and quantum neural networks (QNNs) [2]. These methods have enabled applications in cryptography, combinatorial optimization, physics and chemistry, machine learning, and finance [6, 13, 30, 38, 46]. Central to both algorithm development and hardware validation is the ability to accurately simulate quantum systems.

Hamiltonian simulation plays a particularly important role, as it models the continuous-time evolution of quantum systems governed by the Schrödinger equation [8]. Such simulations are critical for studying physical phenomena, including molecular dynamics and many-body systems, and are especially relevant to today’s quantum hardware, which is costly to access, provides limited coherence times, and is imperfect due to noise. Classical Hamiltonian simulation therefore remains a cornerstone for validating quantum algorithms, benchmarking [42], and exploring algorithmic behavior prior to hardware execution [39].

A common approach to Hamiltonian simulation is to approximate time evolution using Trotterization, decomposing the evolution operator into a sequence of discrete time steps. While higher-order Trotter schemes can improve accuracy [40], achieving high fidelity often requires a large number of time steps, resulting in deep quantum circuits that are difficult to execute on hardware. Consequently, HPC simulation remains indispensable, both for validation and exploration of regimes beyond current hardware [24–26, 28].

However, classical simulation of quantum systems is fundamentally limited by the exponential growth of the Hilbert space with the number of qubits. This “curse of dimensionality” leads to prohibitive computational and memory costs, even on modern HPC platforms. Although parallelism and acceleration help extend feasible problem sizes, performance is often constrained by the efficiency of sparse linear algebra kernels, particularly sparse matrix exponentiation and sparse matrix–vector multiplication [18].

A key observation motivating this work is that problem-Hamiltonians (that have Hermitian properties when represented using matrices) in many quantum applications exhibit strong and persistent structural sparsity. In contrast to generic sparse matrices with irregular non-zero patterns, such Hermitian operators

often concentrate their non-zero entries along a limited number of diagonals due to their local interactions and structured operator decompositions (e.g., Pauli strings). We refer to these as “active” diagonals. Fig. 1 empirically illustrates this phenomenon for Hamiltonians drawn from the HamLib benchmark suite [36]: Despite exponential growth in matrix dimension, both the overall matrix density and the number of active diagonals grow slowly and remain tightly bounded with qubit count.

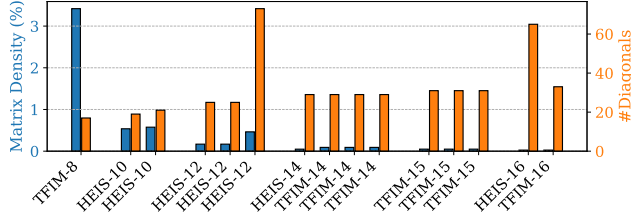


Figure 1: Matrix density (left y-axis) and number of active diagonals (right y-axis) for selected instances from HamLib [36]. Identical x-axis labels correspond to distinct Hamiltonians of the same category with the same qubit size.

Importantly, this diagonal sparsity is not incidental. As shown in Fig. 2, non-zero diagonals exhibit clear structure, symmetry, and widening gaps away from the main diagonal. These properties suggest that much of the computational cost incurred by existing simulation pipelines stems not from inherent complexity, but from a mismatch between Hamiltonian structure and the numerical representations traditionally used in classical simulation.

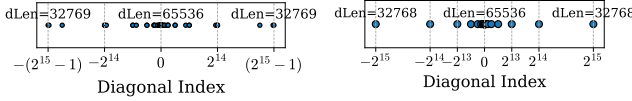


Figure 2: Distribution of diagonal indices in 16-qubit Hermitian matrices: Heisenberg (left) and TFIM (right).

We address this mismatch by designing a novel time-evolution strategy and numerical kernels around diagonal sparsity. We introduce **diagonal-budgeted Trotterization**, a sparsity-aware time-evolution technique that explicitly controls diagonal growth by selecting a carefully bounded number of evolution steps whose factors preserve diagonal structure while tightly controlling fidelity loss. We realize this technique efficiently in **HAMSIM**, a diagonal-aware Hamiltonian simulation approach that uses a compact diagonal-sparse data layout together with high-performance kernels for sparse matrix exponentiation, sparse matrix–matrix multiplication, and sparse matrix–vector multiplication. Implemented in a C++/CUDA framework with Python bindings, HAMSIM leverages SIMD vectorization, multithreading, and GPU acceleration to enable efficient classical Hamiltonian simulation at scales that are challenging for existing methods.

In summary, this paper makes the following contributions:

- We propose **diagonal-budgeted Trotterization**, a sparsity-aware time-evolution strategy that bounds the number of Trotter steps for high-fidelity Hamiltonian simulation.
- We develop a **diagonal-aware Hamiltonian simulation approach**, implemented using a diagonal-sparse data layout

and HPC-optimized kernels for sparse matrix exponentiation, sparse matrix–matrix multiplication, and sparse matrix–vector multiplication.

- We implement and evaluate our approach on multi-core CPUs and GPUs, demonstrating substantial speedups over today’s simulators and strong scaling to larger problem sizes. Across HamLib benchmarks, HAMSIM achieves 182–1,269× speedups on TSP and MaxCut and 4.8–841× on TFIM and Heisenberg on CPUs, while HAMSIM-GPU attains 7.8×–37.1× speedups on 8–10 qubits and about 178× on 12–16 qubit instances, all while maintaining near-perfect fidelity.

2 BACKGROUND

2.1 Hamiltonian Time Evolution

Hamiltonian simulation models the continuous-time evolution of quantum systems governed by the Schrödinger equation,

$$i\hbar \frac{d}{dt} |\psi(t)\rangle = H |\psi(t)\rangle \quad (1)$$

where H is a Hermitian operator encoding system interactions. The formal solution is

$$|\psi(t)\rangle = U(t) |\psi(0)\rangle, \quad U(t) = e^{-iHt}. \quad (2)$$

Hamiltonians in chemistry, condensed matter, and many-body physics are typically expressed as structured sums of non-commuting terms

$$H = \sum_{j=1}^m H_j. \quad (3)$$

A common classical approximation strategy is Trotterization, which factorizes the propagator into short-time evolutions,

$$e^{-iHt} \approx \left(\prod_{j=1}^m e^{-iH_j t/r} \right)^r, \quad (4)$$

where r is the number of Trotter steps. Increasing r improves accuracy but raises computational costs, while higher-order formulas reduce the step count at the expense of more complex per-step operators. Reducing the number of steps needed to reach a target fidelity is therefore critical for both classical simulation performance and near-term quantum execution.

2.2 Scalability and Sparsity

For an n -qubit system, both H and the propagator $U(t)$ are $N \times N$ matrices with $N = 2^n$. Dense storage requires $O(N^2)$ memory, and dense matrix multiplication incurs $O(N^3)$ time, making large-scale simulation infeasible since N grows exponentially with n . Applying a dense propagator to a state vector costs $O(N^2)$ per time step.

In practice, however, physically meaningful Hamiltonians exhibit strong structural sparsity. As shown in Fig. 1, the number of active diagonals grows slowly even as matrix dimension increases exponentially, reflecting locality, limited interaction range, and the structured nature of Pauli-operator decompositions. Preserving and exploiting this diagonal sparsity throughout time evolution is essential to reduce memory footprint and computational cost.

Unfortunately, existing simulation pipelines often fail to preserve this structure as they rely on sparsity-oblivious Trotterization and sparse formats that do not capture diagonal regularity.

2.3 Limitations of Existing Sparse Formats

Sparse matrix formats such as Compressed Sparse Row (CSR), Compressed Sparse Column (CSC), Coordinate (COO), Ellpack (ELL), and diagonal-based formats like SciPy’s DIA representation [27, 35, 43] are widely used in numerical computing. These formats are optimized for arbitrary sparsity patterns and deliver reasonable performance across general workloads.

Diagonal-based formats organize non-zero entries by offset from the main diagonal and are efficient when a matrix has *narrow bandwidth*, i.e., when all non-zero entries lie within a small range of diagonal offsets. However, when active diagonals lie far from the main diagonal, these formats must pad each diagonal to the full matrix width, producing large blocks of unused storage (Tab. 1). Additionally, common DIA implementations often fall back to CSR for matrix operations, reintroducing CSR overheads and limiting achievable performance. General-purpose formats are therefore

Table 1: Existing sparse matrix formats vs. HAMSIM’s layout.

(a) Dense		
$\begin{pmatrix} a & 0 & 0 & b \\ 0 & c & 0 & 0 \\ 0 & 0 & d & 0 \\ e & 0 & 0 & f \end{pmatrix}$		
(b) CSR	(c) CSC	(d) COO
Val: [a, b, c, d, e, f] RPtr: [0, 2, 3, 4, 6] CIdx: [0, 3, 1, 2, 0, 3]	Val: [a, e, c, d, b, f] CPtr: [0, 2, 3, 4, 6] RIdx: [0, 3, 1, 2, 0, 3]	Val: [a, b, c, d, e, f] RIdx: [0, 0, 1, 2, 3, 3] CIdx: [0, 3, 1, 2, 0, 3]
(e) ELL	(f) DIA	(g) HAMSIM
Val: [a, b, c, *, d, *, e, f] CIdx: [0, 3, 1, *, 2, *, 0, 3]	Val: $\begin{bmatrix} a & c & d & f \\ * & * & * & e \\ b & * & * & * \end{bmatrix}$	Val: [a, c, d, f, e, b] Off: [0, -3, 3]

ill-suited for Hamiltonian simulation workloads where sparsity is highly structured, persistent across time steps, and concentrated along specific diagonal offsets. This mismatch motivates the development of approaches that explicitly preserve diagonal structure.

Fig. 2 illustrates the challenge for representative 16-qubit Hamiltonians. Although the number of active diagonals is small, these diagonals may lie far apart, creating large empty regions of just zeros. In diagonal-based formats such as DIA, this leads to substantial padding overhead and inefficient traversal during matrix operations. As system size increases, padding costs grow rapidly, and memory traffic begins to dominate runtime, even though the true number of non-zero entries remains small.

3 METHODOLOGY

This section presents the methodological components underlying HAMSIM. We proceed from the algorithmic strategy to the numerical representation choices and finally to the high-performance execution kernels that together enable scalable time evolution. We begin with diagonal-budgeted Trotterization, which controls sparsity growth during time evolution by selecting a timestep that preserves diagonal structure while meeting a target fidelity. We then describe the diagonal-aware sparse layout used internally by HAMSIM to compactly store and efficiently process the operators arising from this decomposition. Finally, we present the core diagonal-aware kernels for sparse matrix exponentiation, matrix-matrix multiplication, and matrix-vector multiplication that are tailored

for our novel trotterization strategy and data layout forming the computational backbone of HAMSIM.

All algorithms are presented in high-level pseudocode/walk-through form and are implemented in our C++/CUDA framework, with Python bindings for integration into existing simulation workflows. For clarity, the pseudocode/walk-through focuses on the main algorithmic ideas and omits low-level implementation details such as memory management, norm estimation, diagonal enumeration, and SIMD alignment logic. These aspects are handled explicitly in the implementation but are orthogonal to the core algorithmic ideas described here.

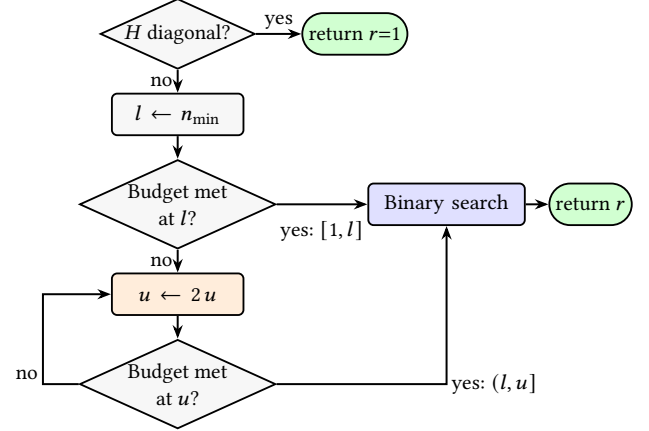


Figure 3: Adaptive timestep selection for diagonal budget D_{\max} . Given Hamiltonian H and evolution time t , find the smallest r satisfying Eq. (5). $r \rightarrow$ #steps, $t \rightarrow$ total time, $l, u \rightarrow$ bounds on r .

3.1 Diagonal-Budgeted Trotterization

Hamiltonian simulation via Trotterization approximates the time-evolution operator e^{-iHt} by decomposing the Hamiltonian $H = \sum_{j=1}^m H_j$ into a product of short-time exponentials. Classical implementations typically fix the number of Trotter steps r a priori, based on theoretical error bounds or rough heuristic choices. These approaches are agnostic to the sparsity structure of the intermediate operators and often lead to unnecessary fill-in during exponentiation, undermining the benefits of sparse representations.

We introduce diagonal-budgeted Trotterization, a sparsity-aware time-evolution strategy that explicitly links the Trotter step size to the diagonal structure of the short-time propagator. As illustrated in Fig. 3, instead of choosing a fixed number of steps, the method adaptively selects the smallest r such that the per-step operator $e^{-iH(t/r)}$ satisfies a user-defined diagonal budget D_{\max} . This constraint limits the growth of non-zero diagonals and ensures that each short-time propagator remains efficiently representable in a diagonal-oriented sparse layout.

A key empirical observation is that, for physically meaningful Hamiltonians, diagonal fill-in increases monotonically as the Trotter step size grows. Consequently, the feasibility of a given diagonal budget defines a monotone predicate over r , enabling efficient search procedures. Given total evolution time t , we determine the minimal number of steps r satisfying

$$\text{num_diags}\left(e^{-iH(t/r)}\right) \leq D_{\max}. \quad (5)$$

This adaptive selection balances sparsity preservation and accuracy: Larger steps reduce computation but risk diagonal growth, while smaller steps maintain sparsity and ensure efficient propagation within HAMSIM’s diagonal-centric kernels.

Diagonal-budgeted Trotterization fundamentally differs from traditional error-driven step selection. Rather than controlling truncation error directly, it constrains structural complexity, yielding a practical upper bound on the number of steps required for high-fidelity simulation while preserving the sparsity patterns common in real Hamiltonians. This structure-aware decomposition forms the algorithmic foundation upon which our diagonal-oriented representation and kernels operate. The overhead of this procedure is incurred as a one-time preprocessing step per Hamiltonian and chosen diagonal budget. In practice, we evaluate a small number of candidate values of r and estimate $\text{num_diags}(e^{-iH(t/r)})$ to identify the minimal feasible step count. This search is inexpensive compared to the overall simulation cost and is amortized across all subsequent time-evolution steps.

3.2 Diagonal Representation in HAMSIM

Many problem Hamiltonians and the short-time propagators generated by diagonal-budgeted Trotterization are sparse with nonzeros concentrated along a small set of diagonals. To exploit this structure, HAMSIM adopts a diagonal-oriented sparse representation that stores only the active diagonals needed for time evolution, without the padding overheads incurred by general-purpose diagonal formats such as SciPy’s DIA. Other aspects of the representation are rather similar to DIA with the key distinction coming from numerical kernels tailored to preserve diagonal structure during exponentiation and time evolution, as described in Sec. 3.3.

The representation follows a diagonal-major layout, where active diagonals are packed back-to-back into contiguous, SIMD-aligned buffers. Each diagonal occupies a fixed offset within these buffers, enabling predictable memory access and efficient vectorized kernels. Tab. 1 illustrates this layout relative to common sparse formats. Although only real values are shown for clarity, the representation fully supports complex-valued matrices.

Formally, a Hamiltonian or propagator in HAMSIM is stored as

$$\text{DiagRepr} = \left\{ \left(d_k, v_k^{\text{Re}}, v_k^{\text{Im}} \right) \mid d_k \in \mathbb{Z} \right\}, \quad (6)$$

where each diagonal index d_k corresponds to two contiguous arrays holding the real and imaginary parts of the diagonal’s entries. Following standard conventions, $d_k = 0$ denotes the main diagonal; positive (resp. negative) values denote superdiagonals (resp. subdiagonals). For an $N \times N$ operator, the logical length of diagonal d_k is $n_k = N - |d_k|$, and HAMSIM allocates exactly n_k elements, with minimal padding only when required for SIMD alignment.

To maximize performance, real and imaginary components are stored in a structure-of-arrays (SoA) layout. Although this is not ideal when every kernel performs coupled complex operations, the vast majority of HAMSIM’s inner loops operate on one component at a time (e.g., real accumulations in intermediate steps), and SoA yields more regular memory access and significantly better vectorization. This choice also eliminates the need to store unused imaginary parts for purely real Hamiltonians.

DIAGONAL-SET ESTIMATION AND PREALLOCATION. HAMSIM’s numerical kernels, including matrix exponentiation, matrix–matrix multiplication, and matrix–vector multiplication, are implemented out of place. Rather than generating diagonals dynamically, HAMSIM performs a lightweight structural dry run that propagates only diagonal offsets. This identifies the exact set of diagonals that will appear in the output, enabling a single preallocation of contiguous storage. Fixing the diagonal layout up-front avoids pointer chasing, dynamic growth, and allocation overhead, all of which are critical for predictable performance and high SIMD/GPU efficiency.

Storage implications. This diagonal-oriented layout becomes increasingly beneficial at scale. A $2^{16} \times 2^{16}$ Hamiltonian with 16 active diagonals at offsets $\pm 2^{15}$ contains roughly 16×2^{15} actual nonzeros. HAMSIM’s layout stores exactly these values, consuming about 8 MB for double-precision complex data. By contrast, SciPy’s DIA format pads each diagonal to length 2^{16} , adding another 8 MB of unused entries for the same matrix. Padding costs grow exponentially with system size and quickly dominate memory usage for realistic Hamiltonians. While memory savings are not the primary goal of HAMSIM, this compact layout offers useful secondary advantages on memory-constrained devices.

3.2.1 Contrast with SciPy’s DIA Format. SciPy’s DIA format stores all diagonals in a dense $n_{\text{diags}} \times N$ matrix, padding every diagonal to the full dimension N . This design works well for narrow-band matrices, where most diagonals lie close to the main diagonal. However, many Hamiltonians arising in physics and quantum algorithms contain active diagonals far from the main diagonal. As shown in Fig. 2, representative 16-qubit Heisenberg and TFIM Hamiltonians contain diagonals at offsets as large as $\pm 2^{15}$. Although these diagonals have logical length 2^{15} , DIA pads them to 2^{16} , wasting half their entries. At 16 bytes per double-precision complex entry, this corresponds to 512 KB of padding per diagonal, doubling with each added qubit. HAMSIM’s layout avoids this padding entirely by allocating storage proportional to the true number of nonzeros, independent of diagonal offset. Combined with its preallocated, SIMD-aligned layout, this provides both memory efficiency and regular access patterns well suited to the short-time propagators produced by diagonal-budgeted Trotterization.

Compact storage alone is not sufficient: Hamiltonian simulation also requires kernels that preserve diagonal structure during exponentiation and time evolution, which are described next.

3.3 HAMSIM Kernels

Building on the HAMSIM storage layout, we implement a suite of diagonal-aware numerical kernels that preserve sparsity and enable efficient Hamiltonian simulation. Unlike general-purpose sparse libraries that permit arbitrary fill-in during multiplication, all kernels in HAMSIM operate on fixed diagonal layouts determined in advance. This design avoids dynamic sparsity growth, enables contiguous preallocation of all output buffers, and yields predictable memory access patterns that map cleanly to SIMD on CPUs and to regular thread/block decompositions on GPUs.

The kernel suite includes sparse matrix–vector multiplication (SpMV), sparse matrix–matrix multiplication (SpGEMM), matrix addition and scaling, Hermitian conjugation, transposition, and sparse

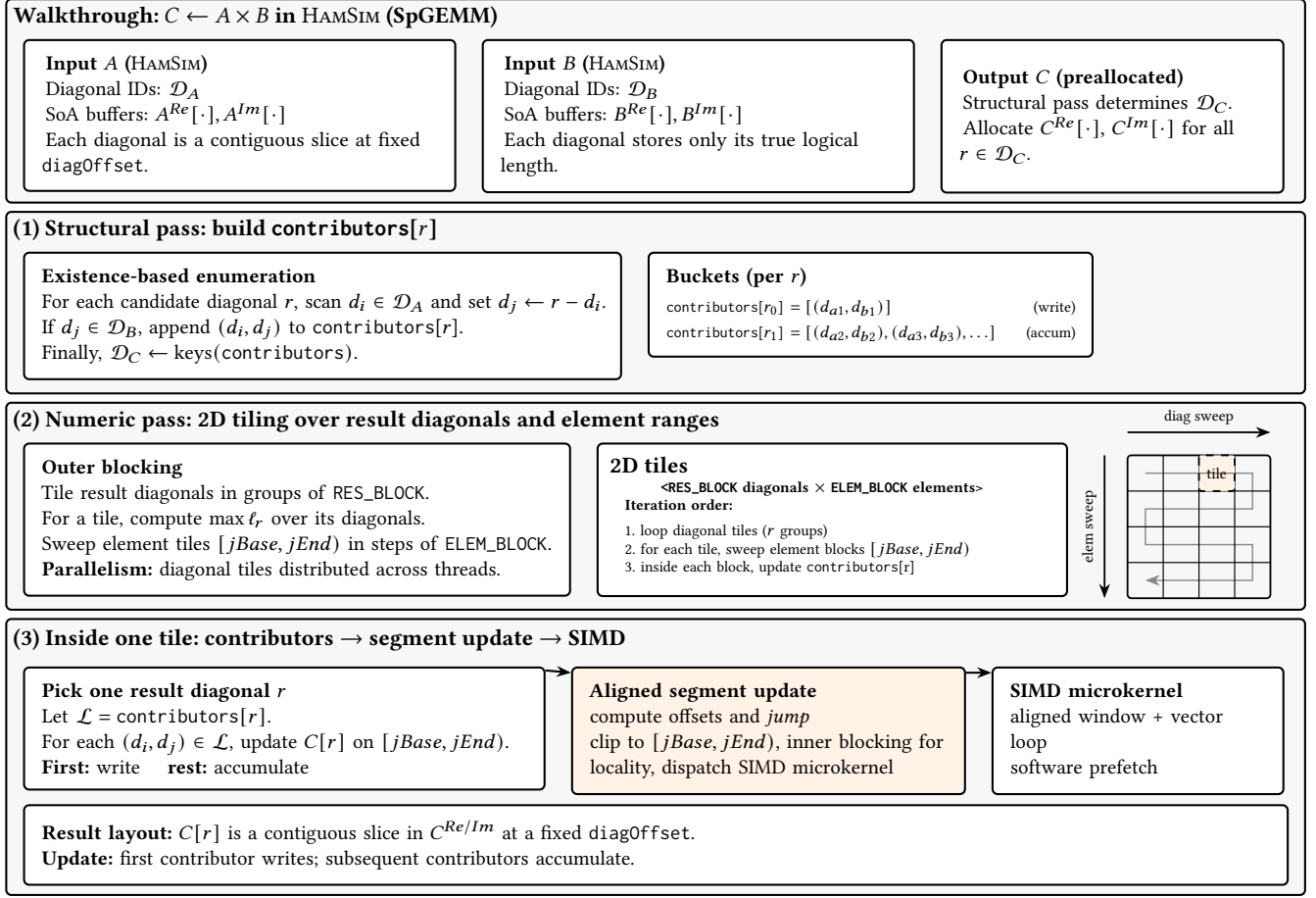


Figure 4: Walkthrough of HAMSIM SpGEMM ($C=AB$). (1) A structural pass builds `contributors[r]` via diagonal-index existence checks. (2) The numeric pass performs 2D tiling over result diagonals (`RES_BLOCK`) and element ranges (`ELEM_BLOCK`). (3) Inside a tile, each r iterates contributor pairs and applies an aligned segment update over $[jBase, jEnd]$, dispatching a SIMD microkernel to write/accumulate into the preallocated $C[r]$ slice.

matrix exponentiation (SpME). Together, these primitives provide the core linear-algebra building blocks required for Trotterized Hamiltonian evolution and unitary simulation, while preserving diagonal structure across intermediate computations.

3.3.1 SpGEMM. SpGEMM is the most performance-critical kernel in HAMSIM, supporting matrix exponentiation. Given two diagonally sparse inputs A and B , the product $C = AB$ is formed by enumerating all diagonal pairs (d_i, d_j) whose offsets satisfy $r = d_i + d_j$ for a candidate result diagonal r . This structural phase determines the exact set of diagonals in C prior to any numerical work.

Fig. 4 illustrates the full execution flow: The structural pass builds `contributors[r]` for every result diagonal; the numeric pass tiles both the diagonal index space and the element-range space (`RES_BLOCK` and `ELEM_BLOCK`); and each tile invokes the **MULTI-DIAGONALS** microkernel, which streams contiguous diagonal segments and performs complex fused multiply-add operations. The figure makes explicit how write-vs-accumulate behavior is handled and how SIMD parallelism is exposed inside each tile.

Numerical execution is parallelized across result diagonals. When a diagonal r has a single contributing pair, its output slice is written directly; when multiple contributors exist, partial results are accumulated explicitly. Each diagonal-level multiplication runs in time linear in the diagonal length.

If both inputs contain only the principal diagonal, SpGEMM reduces to $O(N)$. In general, the complexity is $O(d_A \cdot d_B \cdot N)$, where d_A and d_B denote the numbers of active (non-zero) diagonals, typically small constants for Hamiltonians exhibiting structured sparsity.

3.3.2 SpMV. SpMV is another dominant kernel during Trotterized time evolution, as each step applies a diagonally sparse propagator to a state vector. In HAMSIM, every diagonal d contributes an independent shifted element-wise multiply between the diagonal's contiguous buffer and the input vector. This yields $O(d \cdot N)$ work (where d is the number of active diagonals, typically a small constant) compared to $O(N^2)$ for dense multiplication. The structural phase (Alg. 1) partitions each diagonal into element-range segments

$(d, [start, end))$, enabling thread-level parallelism over tiles, while the numeric phase executes each tile using a streaming inner loop.

Fig. 5 provides a detailed view of the `DIAGONALTIMESVECTOR` microkernel that operates inside each segment. The figure highlights how the diagonal offset induces an affine mapping of the logical index i to input/output indices (j, k) , allowing uniform vectorized execution across all diagonal types. Each tile processes elements in units of `VEC_WIDTH` via SIMD loads/stores, applies fused complex FMA updates, and uses a masked tail to handle partial vectors. Software prefetching and SoA layout further improve memory locality, making the kernel bandwidth-efficient and highly predictable. The per-tile thread-local results are finally reduced into the global output vector, completing the SpMV in linear time.

Algorithm 1 SpMV: Sparse Matrix–Vector Product

Input: Matrix A in HAMSIM’s layout, vector x

Output: Result vector $y = Ax$

```

1:  $N \leftarrow |x|$ 
2:  $y \leftarrow \text{INITIALIZERESULTVECTOR}(N, 0)$ 
   Phase I (structural): tile diagonals into element-range tasks
3:  $\text{tasks} \leftarrow \text{partition each } (d, v^{\text{Re}}, v^{\text{Im}}) \in A.\text{diagonals}()$ 
   into segments  $(d, v^{\text{Re}}, v^{\text{Im}}, start, end)$ 
   Phase II (numeric): parallel over diagonal segments
4: allocate thread-local output buffers  $\{y_{\text{local}}\}$ 
5: #pragma omp parallel for schedule(static)
6: for  $(d, v^{\text{Re}}, v^{\text{Im}}, start, end) \in \text{tasks}$  do
7:   DIAGONALTIMESVECTOR $(d, v^{\text{Re}}, v^{\text{Im}}, x, y_{\text{local}}, start, end)$ 
   ▶ Fig. 5
8: end for
   Phase III (reduction): combine partial results
9:  $y \leftarrow \sum y_{\text{local}}$  ▶ thread-local reduction
10: return  $y$ 

```

3.3.3 SpME. Matrix exponentiation is required to form short-time propagators of the form $e^{-iH\Delta t}$. HAMSIM exploits diagonal layout to provide fast paths and scalable approximations. If the input matrix is purely diagonal, HAMSIM computes the exponential exactly via element-wise complex exponentiation in linear time. For general diagonally sparse Hermitian matrices, HAMSIM applies a truncated Taylor-series expansion (Alg. 2). In practice, we use the standard scaling-and-squaring variant for numerical stability; the squaring phase reduces to repeated SpGEMM calls and thus benefits directly from HAMSIM’s diagonal-aware multiplication. To prevent fill-in from overwhelming intermediate computations, we explicitly cap diagonal growth using the diagonal-budgeting strategy (Sec. 3.1).

Performance characteristics. Fig. 6 depicts the average execution time (y-axis in log-scale) of HAMSIM’s matrix exponential on a 12-qubit Hamiltonian for varying evolution times and diagonal counts. As expected, performance is strongest when the number of diagonals is small and degrades gracefully as intermediate matrices become denser. For small to moderate time steps and structured Hamiltonians (common in physical simulation workloads) HAMSIM maintains both numerical stability (using a scaled and squared

Algorithm 2 SpME: Matrix Exponential via Taylor Series

Input: Hermitian matrix A in HAMSIM’s layout, scalar time t , convergence tolerance (`=1e-6`), and maximum iterations

Output: Unitary $U = e^{-iAt}$ in HAMSIM’s layout

```

1:  $X \leftarrow A * (-i) * t$ 
2:  $U, T \leftarrow \text{identityMatrix}(A.\text{rows})$ 
3: for  $i = 1$  to max_iterations do
4:    $T \leftarrow \text{SpGEMM}(T, X), T \leftarrow T/i$  ▶ Fig. 4
5:   if  $T.\text{norm}() < \text{tolerance}$  then break
6:   end if
7:    $U.\text{add}(T)$ 
8: end for
9: return  $U$ 

```

approach to avoid division by small numbers) and substantial performance advantages over generic sparse and dense approaches.

3.3.4 Implementation details. All kernels are implemented in C++ with OpenMP multi-threading and SIMD vectorization (AVX-512, AVX2, or SSE selected at compile time) for CPUs. For GPUs, CUDA implementations follow the same preallocation strategy to ensure consistent sparsity layouts across backends. Memory allocation is alignment-aware. Optional loop blocking and software prefetching are supported for improved cache locality. Python bindings using `pybind11` [44] and `NumPy` [16] converters enable seamless integration with other Python-based simulation workflows.

Together, these diagonal-aware kernels enable the efficient realization of Trotterized Hamiltonian simulation and form the foundation for the sparse time-evolution methods described next.

3.4 Efficient Sparse Hamiltonian Simulation

The diagonal-aware kernels described in Sec. 3.3, together with the diagonal-budgeted trotterization strategy introduced in Sec. 3.1, enable an end-to-end Hamiltonian simulation pipeline that preserves sparsity throughout time evolution.

Given a Hermitian (problem-Hamiltonian) H and total evolution time T , the goal is to approximate $e^{-iHT}x$ for an initial state vector x without ever materializing dense matrices or allowing uncontrolled fill-in. HAMSIM’s layout stores the problem-Hamiltonian only with a small number of diagonals, reflecting the physical interactions.

When H is purely diagonal, the time evolution operator can be applied exactly in $O(N)$ time via element-wise complex exponentiation. For more general diagonally sparse Hamiltonians, we rely on a Trotterized time evolution scheme in which the total time T is split into n steps of size $\Delta t = T/n$. The key challenge is to choose n large enough to control fill-in during exponentiation, but not so large as to introduce unnecessary computational overhead.

Rather than fixing n a priori, HAMSIM determines it adaptively using the diagonal-budgeted search procedure introduced in Sec. 3.1. This routine selects the smallest n such that the short-time propagator $e^{-iH\Delta t}$ remains within a user-specified diagonal budget D_{max} .

Alg. 3 summarizes the resulting simulation workflow. Once n is determined, a sparse unitary $U \approx e^{-iH\Delta t}$ is constructed using the SpME kernel. The state is then evolved by repeated application of U using the SpMV kernel. Because both kernels operate directly on the

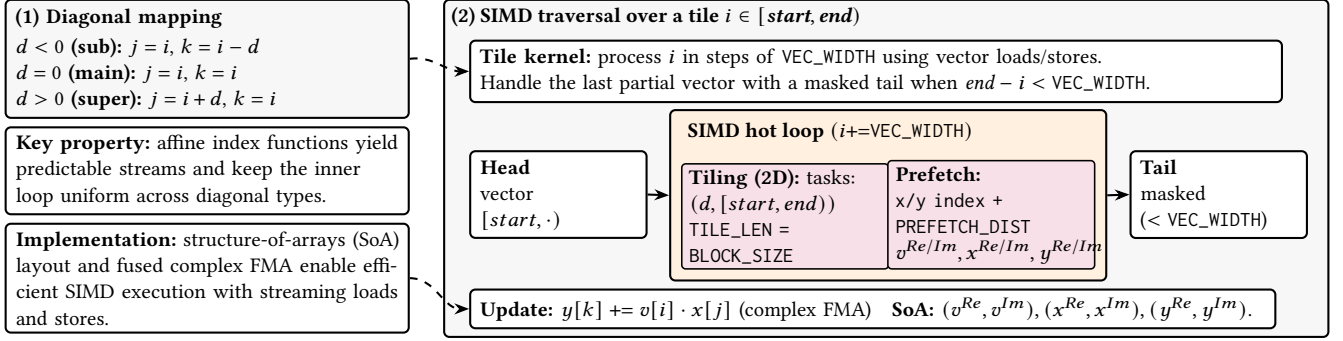


Figure 5: Schematic of the DIAGONALTIMESVECTOR kernel. A diagonal offset induces an affine mapping (j, k) from logical position i to input/output indices, enabling a uniform streaming update. In the implementation, each diagonal is partitioned into $(d, [start, end])$ tiles (2D tasking), and each tile executes a SIMD hot loop with vector loads/stores and a masked tail to handle partial vectors; software prefetching overlaps memory latency with fused complex FMA updates. GPU: we support a single-diagonal element-parallel kernel and a multi-diagonal row-parallel kernel.

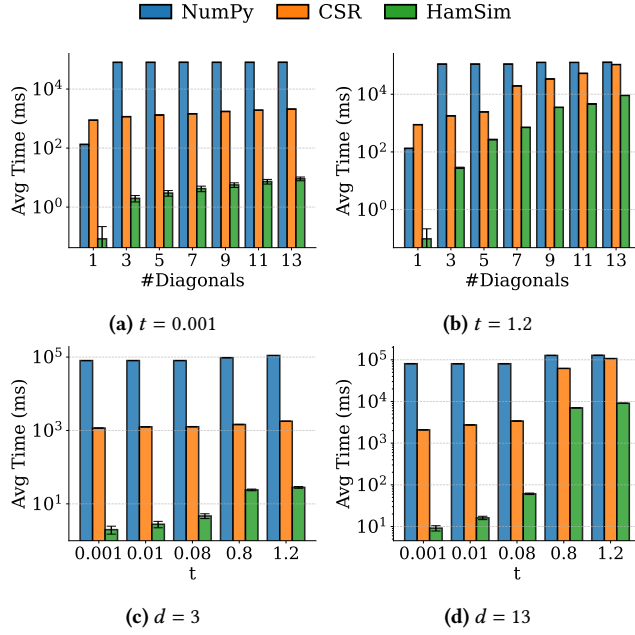


Figure 6: Performance of HAMSIM's matrix exponential on a 12-qubit system for varying t and diagonal counts d .

HAMSIM layout and respect the diagonal budget, all intermediate operators remain sparse and efficiently computable.

This diagonal-budgeted approach achieves two objectives simultaneously, which are difficult to satisfy jointly using conventional error-driven Trotterization. First, it bounds memory usage and arithmetic cost by explicitly limiting the number of active diagonals in each operator. Second, it avoids overly conservative timestep choices dictated by worst-case dense approximations. As shown in Sec. 5, this combination allows HAMSIM to achieve effectively unit

Algorithm 3 Diagonal-Budgeted Hamiltonian Simulation

Input: Hermitian matrix A in HAMSIM layout, initial state x , total evolution time T , diagonal budget D_{\max}

Output: Final state $x_T \approx e^{-iAT}x$

- 1: $n \leftarrow \text{ESTIMATE_TIMESTEPS}(A, T, D_{\max})$ ▷ Fig. 3
- 2: $\Delta t \leftarrow T/n$
- 3: $U \leftarrow \text{SPME}(A, \Delta t)$ ▷ Scaled sparse exponential, Alg. 2
- 4: **for** $i = 1$ **to** n **do**
- 5: $x \leftarrow \text{SPMV}(U, x)$ ▷ Alg. 1
- 6: **end for**
- 7: **return** x

fidelity relative to dense baselines while delivering substantial performance gains for structured Hamiltonians common in quantum simulation workloads.

4 EVALUATIONS

4.1 Setup

Experiments were conducted on a high-performance CPU-based platform, an AMD EPYC 8124P 16-Core Processor with 192 GB of DDR5 4800 MHz ECC memory. Programs are compiled using GCC 12.3 with the `-O3` optimization flag. To leverage parallelism, OpenMP is enabled across all runs, using all 32 hardware threads with simultaneous multithreading (SMT) on this platform.

GPU experiments were performed on a single NVIDIA H100 NVL PCIe Gen5 GPU (94 GB HBM3, SM 9.0). All GPU kernels are implemented in CUDA and compiled with `nvcc` using `-O3` and the default device-side optimizations. The kernels operate directly on the diagonal-sparse HAMSIM layout: for sparse matrix–vector multiplication we use a specialized single-diagonal kernel when the matrix contains only one diagonal, and a row-parallel kernel for the general multi-diagonal case; for sparse matrix–matrix multiplication, each result diagonal is assigned to a thread block, and threads within the block iterate over diagonal positions and accumulate contributions from all contributing diagonal pairs.

Unless otherwise stated, kernels are launched with 256 threads per block and a one-dimensional grid sized to cover all active rows (SpMV) or result diagonals (SpGEMM). Complex-valued data are stored in a structure-of-arrays layout (separate real and imaginary arrays) to improve global-memory coalescing. Reported GPU times correspond to kernel execution time after explicit synchronization.

Table 2: Comparison of time-evolution methods.

Method	Compute	Memory	Adaptive
NumPy	$O(N^3)$	$O(N^2)$	✗
CSR	$O(T \cdot \text{nnz}(H))$	$O(\text{nnz}(H))$	✗
expm_multiply	$O(T \cdot \text{nnz}(H))$	$O(\text{nnz}(H))$	✗
Qiskit-Aer	$O(T \cdot L \cdot N)$	$O(N)$	✓
HAMSIM	$O(T \cdot D \cdot N)$	$O(DN)$	✓

n : number of qubits; $N=2^n$: Hilbert-space dimension; H : Hermitian Hamiltonian; $\text{nnz}(H)$: nonzeros in H ; L : number of terms in H ; D : active diagonals; t : final evolution time; T : timesteps; $\Delta t = t/T$: step size. “Adaptive” indicates adaptive timestep selection.

4.2 Baselines

We compare HAMSIM against state-of-the-art approaches under identical inputs and precision. A high-level compute/memory summary is presented in Tab. 2. All baselines use OpenMP/SIMD (CPU) or CUDA (GPU) parallelism with a consistent number of threads to ensure fair comparison.

NumPy: Constructs the full $N \times N$ propagator or state explicitly in double precision. It is useful as a fidelity reference but quickly becomes infeasible as memory grows quadratically with Hilbert-space size.

SciPy-CSR: Operates on matrices in Compressed Sparse Row format without diagonal awareness. It supports general sparse kernels but does not exploit Hamiltonian structure.

SciPy-DIA: Uses SciPy’s built-in DIAgonal sparse matrix format and kernels. It is optimized for diagonal-sparse matrices but has a few differences from HAMSIM’s layout as discussed in Sec. 3.2.1. To add, while SciPy’s DIA format is designed for diagonal-sparse matrices, it does not maintain a fully diagonal-aware execution path. For operations such as matrix exponentiation, it typically performs format conversions (e.g., to CSR), after which computation proceeds using general-purpose sparse kernels. SciPy does not exploit diagonal-wise blocking, SIMD/vectorized fused kernels, or structure-aware accumulation across diagonals.

SciPy-expm_multiply: Applies e^{-iHt} directly to a vector without forming the full propagator. It is efficient for moderately sparse matrices, but not tuned for diagonal sparsity.

Qiskit-Aer (1st-order Trotter): Simulates time evolution by sequentially applying Hamiltonian terms to a full state-vector, whose dimension doubles with each qubit. The number of steps is chosen to match HAMSIM’s slicing, i.e., runtime scales with both the number of terms and the statevector length, but this often incurs noticeable fidelity loss at the same step budget. Our GPU results compare against state-of-the-art GPU simulation provided by **Qiskit-Aer-GPU** using the cuQuantum [1] (specifically, cuStateVec) library.

Hamsim: Our novel diagonal-budgeted sparse Hamiltonian simulation technique on CPU (HAMSIM) and GPU (HAMSIM-GPU).

Table 3: HamLib benchmarks used in this study; Labels match the figures.

Application	Label	Qubits	#Diagonals	Sparsity (%)
TSP	TSP-1	8	1	99.6094
	TSP-2	8	1	99.6094
	TSP-3	15	1	99.9969
	TSP-4	15	1	99.9969
	TSP-5	16	1	99.9985
	TSP-6	16	1	99.9985
Heisenberg	HEIS-1	6	13	93.7500
	HEIS-2	6	13	93.7500
	HEIS-3	8	25	97.3145
	HEIS-4	8	25	97.3267
	HEIS-5	10	19	99.4629
	HEIS-6	10	19	99.4629
	HEIS-7	14	27	99.9548
MaxCut	MAXCUT-1	10	1	99.9390
	MAXCUT-2	10	1	99.9386
	MAXCUT-3	12	1	99.9840
	MAXCUT-4	12	1	99.9834
	MAXCUT-5	14	1	99.9957
	MAXCUT-6	14	1	99.9958
	MAXCUT-7	16	1	99.9989
	MAXCUT-8	16	1	99.9989
TFIM	TFIM-1	7	15	93.7500
	TFIM-2	8	17	96.5820
	TFIM-3	9	19	98.0469
	TFIM-4	10	21	98.9258
	TFIM-5	14	29	99.9084
	TFIM-6	16	33	99.9741

4.3 Metrics

We evaluate our approach using the following metrics:

Simulation Time: Total wall-clock time is measured, averaged over 10 runs, with whiskers of bar charts providing the standard deviation. All methods are executed under identical input and software configurations.

State Fidelity: To assess numerical accuracy, we compute the fidelity between the quantum state produced by HAMSIM and a reference solution obtained using a dense simulation (using NumPy, when feasible).

Scalability and Memory Efficiency: We assess scalability by increasing the number of qubits in each benchmark. Existing methods can achieve high fidelity only for small problem sizes, but become infeasible as the system size grows due to memory or runtime constraints. In contrast, methods that scale to larger qubit counts typically suffer significant fidelity loss. HAMSIM uniquely combines high fidelity with scalability, enabling accurate simulation of large quantum systems beyond the reach of other approaches.

4.4 Benchmarks

We evaluate HAMSIM on a diverse set of Hamiltonian simulation problems drawn from the HamLib [36] benchmark suite. We select four representative problems spanning multiple quantum domains: the **Transverse Field Ising Model (TFIM)** and the **Heisenberg**

model from condensed matter physics, and the **Traveling Salesman Problem (TSP)** and **MaxCut** from discrete and binary optimization. HamLib provides each Hamiltonian as Pauli strings with complex coefficients, which we load into whatever representation each simulator expects: dense matrices for NumPy, CSR/DIA for SciPy, SparsePauliOp for Qiskit-Aer, and directly into HAMSIM’s diagonal-sparse target format without any conversion overhead.

Tab. 3 summarizes our benchmarks’ categories, applications, and structural properties (sparsity and diagonal count).

5 RESULTS AND DISCUSSION

5.1 Performance Comparison against Baselines

Fig. 7 reports simulation times (log-scale y) over HamLib benchmarks [36] along x . We compare HAMSIM (green; our work) against Qiskit-Aer (red), NumPy (blue; dense, 100%-fidelity baseline), CSR (orange), DIA (brown), and `expm_multiply` (purple). Below, r denotes the number of timesteps used to apply the propagator (i.e., the number of time-slicing steps).

TSP (Fig. 7a): On the 8-qubit `enc-stdbinary` benchmark, HAMSIM achieves **637–652** \times speedups over Aer and **90–125** \times over `expm_multiply`, while matching the dense reference exactly (fidelity ≈ 1.000000). At 15 qubits, HAMSIM further widens the gap to **315–323** \times (vs. Aer) and **442–487** \times (vs. `expm_multiply`), with fidelity again at dense level. On the 16-qubit `enc-unary` instances, HAMSIM remains sub-millisecond as baselines diverge, delivering **(3.6–3.9)** $\times 10^3$ over `expm_multiply`, **182–193** \times over Aer, and **(1.4–1.5)** $\times 10^6$ over CSR; in this regime, baselines do not provide fidelity (or stable accuracy), whereas HAMSIM matches dense references whenever fidelity is reported.

MaxCut (Fig. 7d): Across 10–14 qubits, HAMSIM converts second-scale sparse simulation into sub-millisecond execution, yielding **543–1,269** \times speedups over Aer and **(7.8** $\times 10^3$ **)–(1.2** $\times 10^5$ **)** over CSR, while also outperforming `expm_multiply` by **16–43** \times . When fidelity is available, HAMSIM reproduces the dense solver to numerical precision. At 16 qubits, HAMSIM sustains strong advantage with **70–74** \times over `expm_multiply` and **187–201** \times over Aer; baselines provide no dense-level fidelity guarantee at this scale.

Heisenberg and TFIM (Figs. 7b, 7c): For 6–8 qubits, HAMSIM delivers dramatic improvements: **628–739** \times (Heisenberg) and **506–841** \times (TFIM) over Aer, and **34–98** \times (Heisenberg) and **91–124** \times (TFIM) over `expm_multiply`. Crucially, HAMSIM maintains dense-level fidelity (≈ 1.000000), while Aer’s accuracy collapses at the same r (Heisenberg: as low as 0.18; TFIM: 0.03–0.37). As propagators densify and r increases, speedups tighten but HAMSIM preserves accuracy: for the 10-qubit TFIM case ($r=1120$), HAMSIM achieves **3.2–4.8** \times over `expm_multiply` and **1.4–1.6** \times over Aer while maintaining fidelity ≈ 1.000000 . For 10-qubit Heisenberg ($r=71–107$), HAMSIM continues to provide **32–68** \times speedups over dense baselines with strict fidelity preservation. At the largest sizes (14–16 qubits), HAMSIM completes all simulations reliably, while multiple baselines either underperform severely or fail outright and provide no fidelity assurances. For HEIS-7, Aer may appear competitive in runtime at the chosen r ; however, this occurs at substantially lower fidelity, whereas HAMSIM maintains near-exact accuracy. Achieving comparable fidelity with Aer would require increasing r , leading to significantly higher runtime.

5.2 Fidelity

As seen in Fig. 8, techniques other than Qiskit-Aer reach $\geq 99.99\%$ fidelity if/when they complete. HAMSIM completes with far fewer steps (often a single one) because ESTIMATE_TIMESTEPS (Fig. 3) selects the smallest Δt that preserves diagonal sparsity of $U(\Delta t)$ while meeting a fidelity target. This diagonal-aware budgeting is key: It keeps the sparse Taylor propagator cheap (SpGEMM over a banded pattern). In contrast, there exist Hamiltonians where $U(\Delta t)$ densifies rapidly. Exact diagonal methods can slow in this regime. Reducing Δt (more, smaller steps) restores diagonal sparsity and performance without sacrificing accuracy. Aer also reduces Δt but does not maintain diagonal structure per slice, so both per-step cost (due to fill-in) and step count grow. Dense and CSR baselines require full fill-in throughout.

5.3 Sensitivity to Final Time

HAMSIM maintains high fidelity and speed across $T \in \{0.2, 0.4, 0.8, 1.8\}$ (Fig. 9) but requires significantly fewer steps than Aer. NumPy and CSR remain accurate but slow, and degrade more steeply as T increases. For instance, on a fixed 12-qubit Hamiltonian (Fig. 6), performance is flat and stable for small t and degrades gracefully as t grows and the number of populated diagonals increases (as seen from 3 to 13 diagonals). The benefit of diagonal sparsity is largest when few diagonals are active. As the band widens, the gap to dense-style kernels narrows, which is precisely the behavior observed in TFIM/Heisenberg at larger sizes.

5.4 Multithreading and SIMD

Fig. 10 shows that HAMSIM strong-scales up to 16 threads on our platform, delivering 8–11 \times over the single-threaded baseline on multi-diagonal Hamiltonians. Scaling tapers beyond 16 because the node switches to SMT (32 logical threads on 16 cores) as sibling threads share core execution units and caches, so additional threads add little throughput or sometimes even slightly regress.

Our OpenMP design parallelizes across nonzero diagonals. Single-diagonal inputs (e.g. TSP) expose little outer-loop parallelism due to the simplicity of a single diagonal.

We also prototyped threading within a diagonal using per-thread accumulators and a final reduction. In our bandwidth-sensitive SpMV/SpGEMM kernels this added synchronization and memory traffic and was slower end-to-end, so HAMSIM adopts diagonal-level threading with SIMD inside the inner loop, orthogonal to OpenMP, and accelerates each diagonal’s inner loop, adding a further 1.4–2 \times on larger instances (≥ 10 qubits), with smaller gains on small ones where the loop trip count approaches the vector width.

5.5 GPU Performance

Fig. 11 compares HAMSIM’s GPU kernels against Qiskit-Aer-GPU (which uses NVIDIA’s cuQuantum library under the hood). For smaller problems (8–10 qubits, across TSP, MaxCut, TFIM, Heisenberg), HAMSIM’s GPU kernels deliver substantial speedups over Aer-GPU, typically in the range of 7.8 \times –37.1 \times . For larger problems (12–16 qubits), HAMSIM-GPU continues to outperform Aer-GPU, achieving 45 \times –178 \times speedups across benchmark families, with peak gains of 178.6 \times on several 16-qubit TSP and MaxCut instances.

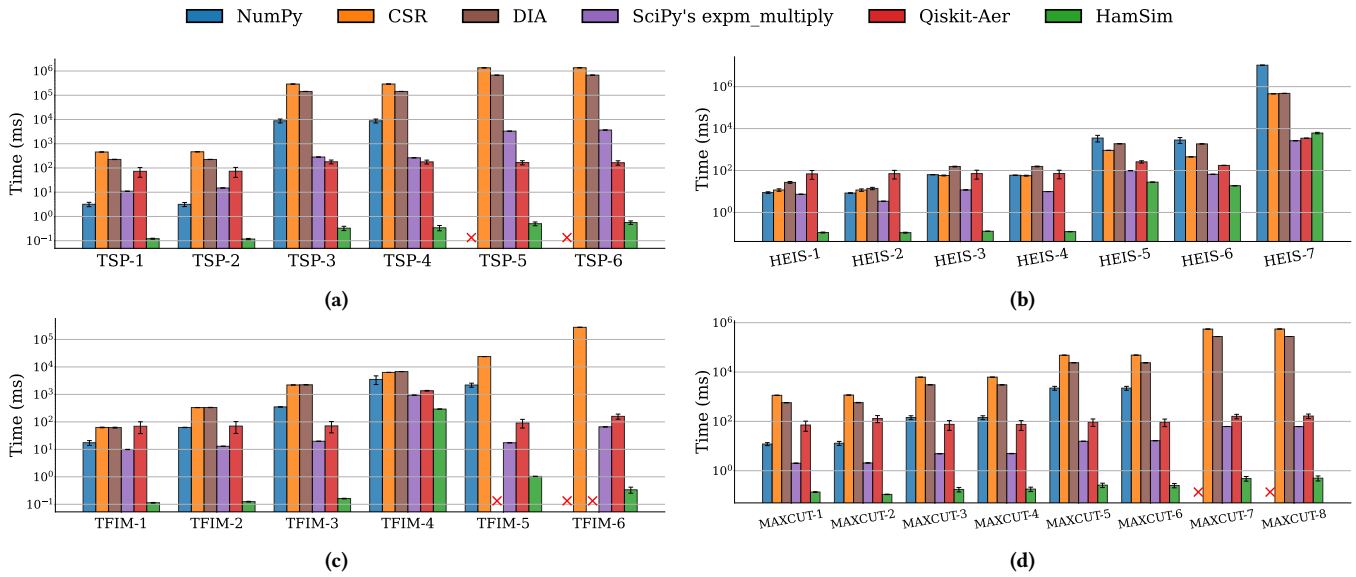


Figure 7: Benchmark runtimes across methods for (a) TSP, (b) Heisenberg, (c) TFIM, and (d) MaxCut.

Note: For TFIM and Heisenberg, the Qiskit-Aer simulator results are shown without fidelity constraints for a fair runtime comparison. Their fidelities are significantly lower, as the number of steps is chosen to match HAMSIM's settings. A red X indicates that the method did not produce a result due to memory-allocation failure (out of memory).

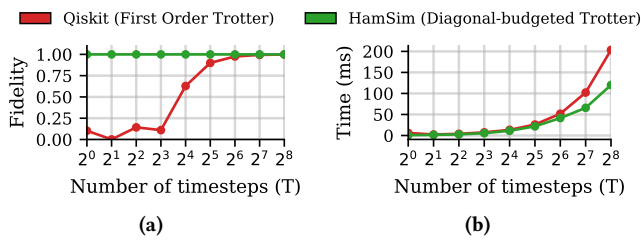


Figure 8: Heisenberg (12-qubit): (a) Fidelity vs. number of timesteps; (b) Simulation time vs. number of timesteps.

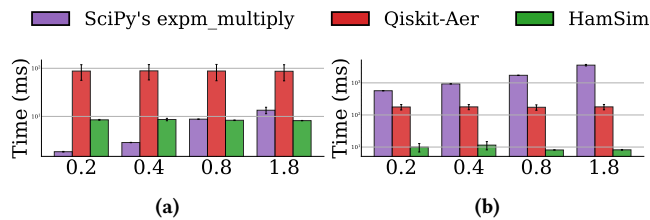


Figure 9: Simulation time vs. final time T for TSP for (a) 8-qubits and (b) 16-qubits.

In terms of accuracy, like our CPU results, HAMSIM-GPU also consistently attains near-perfect fidelity across all benchmarks, including the largest 14-16 qubit cases. Aer-GPU, in contrast, displays substantial fidelity variation as the system size and Trotter structure grow. For example, on TFIM-1 (7 qubits), Aer-GPU reports a

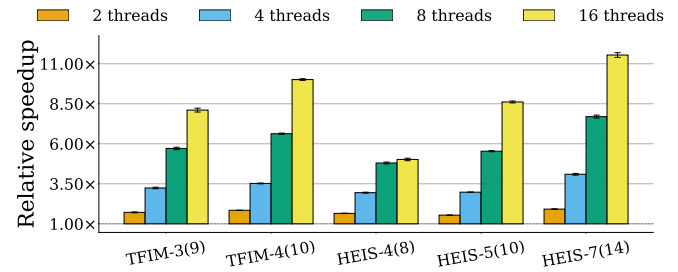


Figure 10: OpenMP strong-scaling of HAMSIM on our Platform (16-core), baseline is single-thread.

fidelity of just **0.0066** despite taking comparable simulation time, whereas HAMSIM-GPU remains effectively at ≈ 1.0 . Even in the one instance where Aer-GPU's runtime approaches ours, namely HEIS-5 (10 qubits, 14.99 s vs. 15.75 s), HAMSIM-GPU maintains its high accuracy (fidelity 1.0 vs. 0.8243).

On GPU, HAMSIM operates directly on the diagonal-sparse HAMSIM layout and avoids materializing dense intermediates. For SpMV we use a specialized single-diagonal kernel when $D=1$ and a row-parallel kernel for the general multi-diagonal case. For SpGEMM, each result diagonal is assigned to a thread block whose threads iterate over diagonal positions and accumulate contributions from all contributing diagonal pairs.

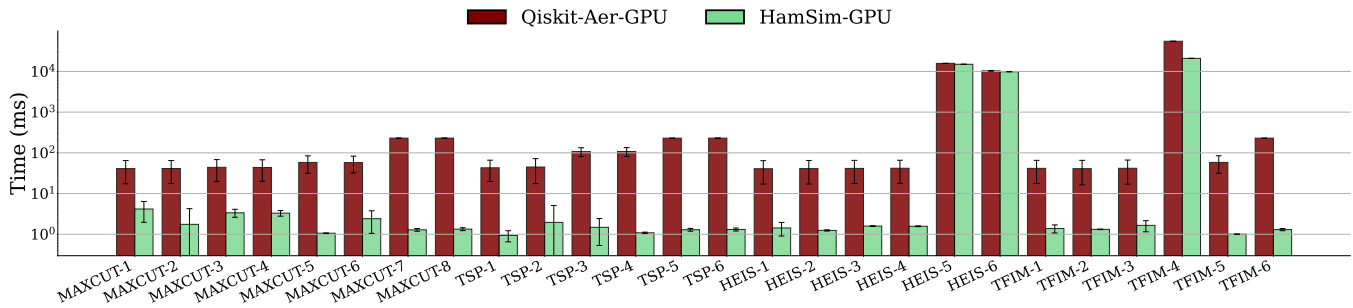


Figure 11: GPU performance of HAMSIM-GPU against Qiskit-Aer-GPU (cuQuantum underneath)

6 RELATED WORK

Several quantum simulators, such as SparQSim [29] and GraFeyn [45], accelerate simulation by tracking only non-zero components of quantum state vectors. While effective for shallow or structured circuits, their performance degrades as circuits deepen or state vectors densify — limitations that are particularly evident in Hamiltonian simulation. In contrast, HAMSIM focuses on the Hamiltonian itself, decoupling performance from state-vector sparsity and enabling efficient time evolution.

Other approaches, including Azure Sparse Simulator [22] and related work [21], compress state vectors with dictionary-based or CSR-like representations. Diagonal formats have likewise been applied to gate-based state-vector simulators [10], yielding performance gains on large circuits. There exist sparse libraries like cuSPARSE [32], MKL Sparse BLAS [20], and other block formats [11] but none preserve the long-range diagonal offsets, leading to padding or scattered accesses. Hence, these simulators/methods often incur overhead from format conversions and are not tuned for Hamiltonian simulations. In contrast, HAMSIM maintains diagonal locality and directly supports matrix exponentiation and matrix-vector multiplication natively for Hamiltonian simulation.

Theoretical advances in sparse Hamiltonian simulation, such as Taylor-series-based methods [4, 5, 9] and quantum walk techniques [3, 7], provide asymptotic complexity bounds on the cost of simulating e^{-iHt} in terms of t , ϵ , and $\|H\|$, with complexities polylogarithmic in $1/\epsilon$ and near-linear in $t\|H\|$, but are formulated in abstract algorithmic models and do not directly apply to high-performance classical simulation. In particular, they assume oracle access to sparse Hamiltonians and do not account for concrete data layouts, memory movement, or practical kernel-level implementations. HAMSIM draws on these insights but addresses practical challenges through diagonal-budgeted Trotterization.

Adaptive-step product formulas have been explored as an alternative to fixed-step Trotter schemes. Zhao et al. propose ADA-Trotter [47], which uses a feedback loop to select the largest Trotter step size consistent with bounded deviations in conserved quantities, effectively performing a search over feasible step sizes. Ostmeier [33] surveys optimized Suzuki–Trotter decompositions [33, 40] and highlights adaptive stepping and embedded error estimators as promising directions. Ikeda and Fujii [19] develop a

precision-guaranteed method that estimates Trotter error via embedded decompositions and adjusts the step size accordingly. HAMSIM differs from these techniques by targeting structure-preserving time evolution under a diagonal-sparsity budget, using step selection to maintain diagonal locality and fidelity rather than only controlling global approximation error.

General-purpose libraries like Quantuloop [41], QuTiP [23], and QSW_MPI [31] support sparse simulation using CSR, CSC, or bit-wise formats. However, they do not exploit structured diagonal sparsity, and often rely on dense matrix exponentiation or fixed-format sparse kernels. Thus, comparisons against dense (NumPy, Qiskit-Aer) and various sparse baselines are more representative for evaluating HAMSIM.

7 CONCLUSION

This work demonstrates that the efficiency of classical Hamiltonian simulation depends on the alignment between numerical kernels and physical operator structure. While many quantum Hamiltonians exhibit persistent diagonal sparsity, conventional simulators use sparsity-oblivious representations that fail to exploit this property.

We introduced diagonal-budgeted Trotterization, a strategy that explicitly bounds diagonal growth to maintain high simulation fidelity. Our implementation, HAMSIM, realizes this via a domain-specific sparse layout and HPC-optimized kernels providing efficient execution on multicore CPUs and GPUs.

Evaluating HamLib benchmarks, HAMSIM achieves CPU speedups of 182–1,269 \times for optimization problems and 4.8–841 \times for physical models. On GPUs, HAMSIM attains up to 178 \times speedup. Crucially, HAMSIM maintains **near-perfect fidelity** (≈ 1.0) in regimes where existing simulators suffer from accuracy loss or prohibitive latency.

This work underscores the value of structure-driven design in scientific computing. HAMSIM provides a foundation for accelerating near-diagonal operators, offering a scalable path for exploiting problem-specific structure in large-scale numerical workloads.

ACKNOWLEDGMENTS

This work was supported in part by NSF CISE-2217020, CISE-2316201, OMA-2120757, PHY-1818914, PHY-2325080 and DOE DE-SC0025384. This research used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department

of Energy under Contract No. DE-AC05-00OR22725. *Notice:* This manuscript has been authored by UT-Battelle, LLC, under contract DE-AC05-00OR22725 with the US Department of Energy (DOE). The US government retains and the publisher, by accepting the article for publication, acknowledges that the US government retains a nonexclusive, paid-up, irrevocable, worldwide license to publish or reproduce the published form of this manuscript, or allow others to do so, for US government purposes. DOE will provide public access to these results of federally sponsored research in accordance with the DOE Public Access Plan (<https://www.energy.gov/doi-public-access-plan>).

Appendix

8 FUTURE WORK

Fused Matrix-Exponential SpMV: Right now, each Trotter step computes $e^{-iH\Delta t}$ as an explicit DiaQ matrix and then multiplies it into the state — two passes over memory. The obvious next step is to fuse these into a single kernel that accumulates the Taylor terms directly into $|\psi\rangle$ (Eq. 7):

$$|\psi^{(k)}\rangle = \frac{-iH\Delta t}{k} |\psi^{(k-1)}\rangle, \quad |\psi(t + \Delta t)\rangle \approx \sum_{k=0}^K |\psi^{(k)}\rangle. \quad (7)$$

This is essentially what `scipy.sparse.linalg.expm_multiply` does, but for the diagonal-sparse case the per-iteration cost is much lower and the intermediate matrix never needs to be stored at all. At large qubit counts where memory is the bottleneck, this matters.

Time-Dependent Hamiltonians: All benchmarks in this paper use time-independent H , but driven systems, adiabatic schedules, and pulse-level control all require $H(t)$. For piecewise-constant protocols this is straightforward — re-run the diagonal-budget search at each segment. The more interesting case is a parameterized family $H(t) = H_0 + f(t)H_1$, where the budget search on the templates H_0 and H_1 carries over to every instance, and only the scalar $f(t)$ changes per step. We plan to expose this through a `-schedule` argument in *Artifact: HamSim*.

Diagonal Sparsity Beyond Quantum: The DiaQ format has no quantum physics in it — it is just a compact representation for matrices whose nonzeros fall on a bounded number of diagonals. That structure shows up in many places outside quantum simulation: 1D and 2D lattice force matrices in molecular dynamics, circulant operators in signal processing, and graph Laplacians for regular-degree graphs. DiaQ’s kernels’ performance carries over to those settings and is worth exploring.

Higher-Order Trotter: We use first-order Trotterization throughout. Fourth-order Suzuki-Trotter cuts the per-step error significantly, which should reduce the step count r needed to hit a fidelity target — but it adds more sub-steps per iteration, each of which is a DiaQ SpME on a fraction of H . Whether the reduced r outweighs the extra sub-step overhead, and how the diagonal budget interacts with higher-order splittings, is something the current benchmark suite is well-positioned to answer.

Open Quantum Systems: The Lindblad master equation $\dot{\rho} = -i[H, \rho] + \sum_k (L_k \rho L_k^\dagger - \frac{1}{2}\{L_k^\dagger L_k, \rho\})$ is the standard model for dissipative quantum dynamics. For local noise channels — dephasing, amplitude damping — the jump operators L_k are diagonal-sparse, so every term in the superoperator is a Kronecker product of DiaQ matrices. Extending *Artifact: libdiaq* to density-matrix simulation would let us study decoherence effects without leaving the diagonal-sparse framework.

Distributed Simulation: The current code is single-node. Getting past roughly 20 qubits on CPU requires distributing the state vector across nodes with MPI. HAMSIM’s diagonal structure actually helps here: diagonal d maps element i to $i + d$, so inter-node communication at each SpMV step is a fixed-offset boundary exchange rather than an all-to-all. This is a cleaner communication pattern

than general sparse formats, and it should be possible to pipeline the boundary exchange with the interior compute.

Automatic Diagonal-Budget Selection: In *Artifact: HamSim*, D_{\max} is currently user-specified as the `PERCENT_DIAGONALS` argument to `benchmark_apps/app2_hamlib_latest.py`. There is no principled reason a user should have to guess this. A short profiling pass that evaluates diagonal fill-in of $e^{-iH(t/r)}$ at a few values of r would be enough to find the budget that minimizes total flops for a given fidelity. Alternatively, the diagonal count curve as a function of t/r is smooth and predictable enough that a simple model fit at reduced qubit count should extrapolate well.

Sparse Tensor Network Contraction: In DMRG and MPS algorithms, the site-local operators being contracted are often diagonal-sparse in the on-site basis. The *Artifact: libdiaq* `tensor_dot` and `reshape` primitives already support the necessary contractions — we demonstrated this in the Quick Start of the *Artifact: libdiaq* repository. Whether DiaQ can serve as a useful backend for tensor network codes in practice — and what the performance ceiling looks like against dense or CSR contractions on current hardware — is a natural next experiment.

BENCHMARK DETAILS

Tab. 4 extends Tab. 3 with the full **HDF5File/Key** paths used as input to `benchmark_apps/app2_hamlib_latest.py`.

ARTIFACT: LIBDIAQ

The `libdiaq` C++ library is available at <https://github.com/srikarchundry/diaq>. Prerequisites: `cmake` ≥ 3.20 , a C++17-capable compiler (`gcc` ≥ 12 or `Apple Clang` ≥ 14), and optionally `CUDA` ≥ 11 for GPU kernels.

Build and test:

```
# Clone with submodules (pybind11, parallel-
  hashmap)
$ git clone --recurse-submodules \
  https://github.com/srikarchundry/diaq.git
$ cd diaq

# Build CPU library and test binary
$ ./build.sh cpu -DMAKE_TESTS=ON
```

Python quick start (SpMV):

```
import numpy as np
import diaq as dq

H_np = np.array([[1, 0, 0, 0], [0, -1, 2, 0],
                 [0, 2, -1, 0], [0, 0, 0, 1]],
                 dtype=np.complex128)
psi = np.array([1, 0, 0, 0], dtype=np.complex128)

H      = dq.from_numpy(H_np)
psi_dq = dq.from_numpy_vector(psi)
y_dq   = dq.spMV(H, psi_dq)
y      = dq.to_numpy_vector(y_dq) # y == H @ psi
```

Table 4: HamLib benchmarks used in this study; Labels match the figures.

Application	Label	HDF5File/Key	Qubits	#Diagonals	Sparsity (%)
TSP	TSP-1	TSP_Ncity-4/tsp_rand-008_Ncity-4_enc-stdbinary	8	1	99.6094
	TSP-2	TSP_Ncity-4/tsp_rand-007_Ncity-4_enc-stdbinary	8	1	99.6094
	TSP-3	TSP_Ncity-5/tsp_rand-007_Ncity-5_enc-stdbinary	15	1	99.9969
	TSP-4	TSP_Ncity-5/tsp_rand-006_Ncity-5_enc-stdbinary	15	1	99.9969
	TSP-5	TSP_Ncity-4/tsp_rand-008_Ncity-4_enc-unary	16	1	99.9985
	TSP-6	TSP_Ncity-4/tsp_rand-006_Ncity-4_enc-unary	16	1	99.9985
Heisenberg	HEIS-1	heis/graph-1D-grid-pbc-qubitnodes_Lx-6_h-5	6	13	93.7500
	HEIS-2	heis/graph-1D-grid-pbc-qubitnodes_Lx-6_h-3	6	13	93.7500
	HEIS-3	heis/graph-3D-grid-pbc-qubitnodes_Lx-2_Ly-2_Lz-2_h-5	8	25	97.3145
	HEIS-4	heis/graph-3D-grid-pbc-qubitnodes_Lx-2_Ly-2_Lz-2_h-3	8	25	97.3267
	HEIS-5	heis/graph-1D-grid-nonpbc-qubitnodes_Lx-10_h-5	10	19	99.4629
	HEIS-6	heis/graph-1D-grid-nonpbc-qubitnodes_Lx-10_h-3	10	19	99.4629
	HEIS-7	heis/graph-1D-grid-nonpbc-qubitnodes_Lx-14_h-0.5	14	27	99.9548
MaxCut	MAXCUT-1	ham-graph-regular_reg-4/reg-4_n-10_rinst-16	10	1	99.9390
	MAXCUT-2	ham-graph-regular_reg-4/reg-4_n-10_rinst-17	10	1	99.9386
	MAXCUT-3	ham-graph-regular_reg-4/reg-4_n-12_rinst-18	12	1	99.9840
	MAXCUT-4	ham-graph-regular_reg-4/reg-4_n-12_rinst-16	12	1	99.9834
	MAXCUT-5	ham-graph-regular_reg-4/reg-4_n-14_rinst-17	14	1	99.9957
	MAXCUT-6	ham-graph-regular_reg-4/reg-4_n-14_rinst-18	14	1	99.9958
	MAXCUT-7	ham-graph-regular_reg-4/reg-4_n-16_rinst-16	16	1	99.9989
	MAXCUT-8	ham-graph-regular_reg-4/reg-4_n-16_rinst-19	16	1	99.9989
TFIM	TFIM-1	tfim/graph-1D-grid-pbc-qubitnodes_Lx-7_h-6	7	15	93.7500
	TFIM-2	tfim/graph-3D-grid-pbc-qubitnodes_Lx-2_Ly-2_Lz-2_h-6	8	17	96.5820
	TFIM-3	tfim/graph-2D-triag-pbc-qubitnodes_Lx-3_Ly-5_h-6	9	19	98.0469
	TFIM-4	tfim/graph-1D-grid-pbc-qubitnodes_Lx-10_h-6	10	21	98.9258
	TFIM-5	tfim/graph-1D-grid-nonpbc-qubitnodes_Lx-14_h-0.5	14	29	99.9084
	TFIM-6	tfim/graph-1D-grid-nonpbc-qubitnodes_Lx-16_h-0.1	16	33	99.9741

ARTIFACT: HAMSIM

The HamSim benchmarking suite is available at https://github.com/srikarchundury/diaq_for_hamsim. It implements diagonal-budgeted Trotterization on top of libdiaq and benchmarks it against Qiskit-Aer and SciPy baselines across all Hamiltonians in Tab. 3.

Setup:

```
$ git clone --recurse-submodules \
  https://github.com/srikarchundury/
  diaq_for_hamsim.git
$ cd diaq_for_hamsim
$ pip install -r requirements.txt
```

Download HamLib datasets (TFIM example):

```
$ wget -r -np -nH --cut-dirs=4 -R "index.html*" \
  https://portal.nersc.gov/cfs/m888/dcams/
  hamlib/condensedmatter/tfim/
```

Run end-to-end benchmark:

```
$ python benchmark_apps/app2_hamlib_latest.py \
  10 tfim.hdf5 \
  "graph-1D-grid-pbc-qubitnodes_Lx-10_h-6" \
  diaq 1.8 5 True
```

The seven positional arguments are: PERCENT_DIAGONALS (diagonal budget as a percentage of the maximum possible diagonal

count, e.g., 10 for 10%), HDF5_FILE (path to a HamLib HDF5 file), HDF5_KEY (dataset key identifying the Hamiltonian instance, matching the **HDF5File/Key** column of Tab. 4), METHOD (one of diaq, diaq-gpu, csr, expm_multiply, or qiskit-aer), FINAL_TIME (total evolution time t), ITRS (timing repetitions for averaging), and TEST_FIDELITY (True to compute fidelity against a dense NumPy reference).

REFERENCES

- [1] Harun Bayraktar, Ali Charara, David Clark, Saul Cohen, Timothy Costa, Yao-Lung L. Fang, Yang Gao, Jack Guan, John Gunnels, Azzam Haidar, Andreas Hehn, Markus Hohnerbach, Matthew Jones, Tom Lubowe, Dmitry Lyakh, Shinya Morino, Paul Springer, Sam Stanwyck, Igor Terentyev, Satya Varadhan, Jonathan Wong, and Takuma Yamaguchi. 2023. cuQuantum SDK: A High-Performance Library for Accelerating Quantum Science. arXiv:2308.01999 [quant-ph]
- [2] Ville Bergholm, Josh Izaac, Maria Schuld, Christian Gogolin, Shah Nawaz Ahmed, Vishnu Ajith, M Sohaib Alam, Guillermo Alonso-Linaje, B AkashNarayanan, Ali Asadi, et al. 2018. Pennylane: Automatic differentiation of hybrid quantum-classical computations. arXiv preprint arXiv:1811.04968 (2018).
- [3] Dominic W. Berry, Graeme Ahokas, Richard Cleve, and Barry C. Sanders. 2007. Efficient Quantum Algorithms for Simulating Sparse Hamiltonians. *Communications in Mathematical Physics* 270 (2007), 359–371.
- [4] Dominic W. Berry, Andrew M. Childs, Richard Cleve, Robin Kothari, and Rolando D. Somma. 2014. Exponential Improvement in Precision for Simulating Sparse Hamiltonians. *Proceedings of the 46th Annual ACM Symposium on Theory of Computing (STOC)* (2014), 283–292.
- [5] Dominic W. Berry, Andrew M. Childs, Richard Cleve, Robin Kothari, and Rolando D. Somma. 2015. Simulating Hamiltonian Dynamics with a Truncated Taylor Series. *Physical Review Letters* 114, 9 (March 2015). <https://doi.org/10.1103/physrevlett.114.090502>
- [6] Jacob Biamonte, Peter Wittek, Nicola Pancotti, Patrick Rebentrost, Nathan Wiebe, and Seth Lloyd. 2017. Quantum machine learning. *Nature* 549, 7671 (sep 2017),

- 195–202. <https://doi.org/10.1038/nature23474>
- [7] Andrew M. Childs and Robin Kothari. 2010. Simulating Sparse Hamiltonians with Star Decompositions. In *Theory of Quantum Computation, Communication, and Cryptography (TQC)*.
- [8] Andrew M. Childs, Dmitri Maslov, Yunseong Nam, Neil J. Ross, and Yuan Su. 2018. Toward the first quantum simulation with quantum speedup. *Proceedings of the National Academy of Sciences* 115, 38 (Sept. 2018), 9456–9461. <https://doi.org/10.1073/pnas.1801723115>
- [9] Andrew M. Childs and Nathan Wiebe. 2012. Hamiltonian Simulation Using Linear Combinations of Unitary Operations. *Quantum Information and Computation* 12, 11&12 (Nov. 2012), 901–924. <https://doi.org/10.26421/QIC12.11-12>
- [10] Srikar Chundury, Jiajia Li, In-Saeng Suh, and Frank Mueller. 2024. DiaQ: Efficient State-Vector Quantum Simulation. [arXiv preprint arXiv:2405.01250](https://arxiv.org/abs/2405.01250) (2024).
- [11] I.S. Duff. 1977. A survey of sparse matrix research. *Proc. IEEE* 65, 4 (1977), 500–535. <https://doi.org/10.1109/PROC.1977.10514>
- [12] Edward Farhi, Jeffrey Goldstone, and Sam Gutmann. 2014. A Quantum Approximate Optimization Algorithm. [arXiv:1411.4028 \[quant-ph\]](https://arxiv.org/abs/1411.4028) <https://arxiv.org/abs/1411.4028>
- [13] Edward Farhi, Jeffrey Goldstone, and Sam Gutmann. 2014. A Quantum Approximate Optimization Algorithm. [arXiv:1411.4028 \[quant-ph\]](https://arxiv.org/abs/1411.4028)
- [14] Edward Farhi, Jeffrey Goldstone, Sam Gutmann, and Leo Zhou. 2022. The Quantum Approximate Optimization Algorithm and the Sherrington-Kirkpatrick Model at Infinite Size. *Quantum* 6 (July 2022), 759. <https://doi.org/10.22331/q-2022-07-07-759>
- [15] Lov K. Grover. 1996. A fast quantum mechanical algorithm for database search. In *Proceedings of the Twenty-Eighth Annual ACM Symposium on Theory of Computing (Philadelphia, Pennsylvania, USA) (STOC '96)*. Association for Computing Machinery, New York, NY, USA, 212–219. <https://doi.org/10.1145/237814.237866>
- [16] Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. 2020. Array programming with NumPy. *Nature* 585, 7825 (Sept. 2020), 357–362. <https://doi.org/10.1038/s41586-020-2649-2>
- [17] Aram W. Harrow, Avinandan Hassidim, and Seth Lloyd. 2009. Quantum Algorithm for Linear Systems of Equations. *Phys. Rev. Lett.* 103 (Oct 2009), 150502. Issue 15. <https://doi.org/10.1103/PhysRevLett.103.150502>
- [18] Guyue Huang, Guohao Dai, Yu Wang, and Huazhong Yang. 2020. GE-SPMM: General-Purpose Sparse Matrix-Matrix Multiplication on GPUs for Graph Neural Networks. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*. 1–12. <https://doi.org/10.1109/SC41405.2020.00076>
- [19] Tatsuhiko N. Ikeda, Hideki Kono, and Keisuke Fujii. 2024. Measuring Trotter error and its application to precision-guaranteed Hamiltonian simulations. *Physical Review Research* 6, 3 (Sept. 2024). <https://doi.org/10.1103/physrevresearch.6.033285>
- [20] Intel Corporation. 2024. Intel® oneAPI Math Kernel Library (oneMKL) Developer Reference. <https://www.intel.com/Version/2024.2>
- [21] Samuel Jaques and Thomas Häner. 2021. Leveraging State Sparsity for More Efficient Quantum Simulations. [arXiv preprint arXiv:2105.01533](https://arxiv.org/abs/2105.01533) (2021).
- [22] Samuel Jaques and Thomas Häner. 2025. Azure Quantum Sparse Simulator. <https://learn.microsoft.com/en-us/azure/quantum/user-guide/machines/sparse/>.
- [23] J. R. Johansson, P. D. Nation, and F. Nori. 2012. QuTiP: An open-source Python framework for the dynamics of open quantum systems. *Computer Physics Communications* 183, 8 (2012), 1760–1772.
- [24] Tyson Jones, Anna Brown, Ian Bush, and Simon C Benjamin. 2019. QuEST and high performance simulation of quantum computers. *Scientific reports* 9, 1 (2019), 10736.
- [25] Ang Li, Bo Fang, Christopher Granade, Guen Prawiroatmodjo, Bettina Hein, Martin Rotteler, and Sriram Krishnamoorthy. 2021. SV-Sim: Scalable PGAS-based State Vector Simulation of Quantum Circuits. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*.
- [26] Ang Li, Omer Subasi, Xiu Yang, and Sriram Krishnamoorthy. 2020. Density Matrix Quantum Circuit Simulation via the BSP Machine on Modern GPU Clusters. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*.
- [27] Jiajia Li, Guangming Tan, Mingyu Chen, and Ninghui Sun. 2013. SMAT: An input adaptive auto-tuner for sparse matrix-vector multiplication. In *Proceedings of the 34th ACM SIGPLAN conference on Programming language design and implementation*. 117–126.
- [28] Shuang Liang, Yuncheng Lu, Ce Guo, Paul Kelly, Wayne Luk, and Hongxiang Fan. 2026. Advancing full-stack acceleration for schrödinger-style quantum simulation. In *32nd IEEE International Symposium on High Performance Computer Architecture (HPCA) 2026*.
- [29] Huan-Yu Liu et al. 2025. SparQSim: Simulating Scalable Quantum Algorithms via Sparse Quantum State Representations. [arXiv preprint arXiv:2503.15118](https://arxiv.org/abs/2503.15118) (2025).
- [30] Guang Hao Low and Isaac L. Chuang. 2019. Hamiltonian Simulation by Qubitization. *Quantum* 3 (jul 2019), 163. <https://doi.org/10.22331/q-2019-07-12-163>
- [31] Oliver Mülken, Maxim Dolgushev, and Alexander Blumen. 2017. QSW_MPI: Parallel Simulation of Quantum Stochastic Walks. *New Journal of Physics* 19, 8 (2017), 083027.
- [32] NVIDIA Corporation. 2024. cuSPARSE. [https://docs.nvidia.com/cuda/cusparse/Version 12.6](https://docs.nvidia.com/cuda/cusparse/Version%2012.6).
- [33] Johann Ostmeier. 2023. Optimised Trotter decompositions for classical and quantum computing. *Journal of Physics A: Mathematical and Theoretical* 56, 28 (June 2023), 285303. <https://doi.org/10.1088/1751-8121/acde7a>
- [34] Alberto Peruzzo, Jarrod McClean, Peter Shadbolt, Man-Hong Yung, Xiao-Qi Zhou, Peter J. Love, Alán Aspuru-Guzik, and Jeremy L. O’Brien. 2014. A variational eigenvalue solver on a photonic quantum processor. *Nature Communications* 5, 1 (jul 2014). <https://doi.org/10.1038/ncomms5213>
- [35] Yousef Saad. 1990. A Basic Tool Kit for Sparse Matrix Computations. <https://api.semanticscholar.org/CorpusID:207974787>
- [36] Nicolas PD Sawaya, Daniel Marti-Dafcik, Yang Ho, Daniel P Tabor, David E Bernal Neira, Alicia B Magann, Shavindra Premaratne, Pradeep Dubey, Anne Matsuura, Nathan Bishop, Wibe A de Jong, Simon Benjamin, Ojas Parekh, Norm Tubman, Katherine Klymko, and Daan Camps. 2024. HamLib: A library of Hamiltonians for benchmarking quantum algorithms and hardware. *Quantum* 8 (Dec. 2024), 1559. <https://doi.org/10.22331/q-2024-12-11-1559>
- [37] P.W. Shor. 1994. Algorithms for quantum computation: discrete logarithms and factoring. In *Proceedings 35th Annual Symposium on Foundations of Computer Science*. 124–134. <https://doi.org/10.1109/SFCS.1994.365700>
- [38] Peter W. Shor. 1997. Polynomial-Time Algorithms for Prime Factorization and Discrete Logarithms on a Quantum Computer. *SIAM J. Comput.* 26, 5 (oct 1997), 1484–1509. <https://doi.org/10.1137/s0097539795293172>
- [39] Mikhail Smelyanskiy, Nicolas P. D. Sawaya, and Alán Aspuru-Guzik. 2016. qHiP-STER: The Quantum High Performance Software Testing Environment. (1 2016). [arXiv:1601.07195 \[quant-ph\]](https://arxiv.org/abs/1601.07195)
- [40] Masuo Suzuki. 1992. General theory of higher-order decomposition of exponential operators and symplectic integrators. *Physics Letters A* 165, 5 (1992), 387–395. [https://doi.org/10.1016/0375-9601\(92\)90335-J](https://doi.org/10.1016/0375-9601(92)90335-J)
- [41] Quantuloop Team. 2024. Quantuloop HPC: Sparse Simulation of Quantum Circuits. <https://quantuloop.com>.
- [42] Benjamin Villalonga, Dmitry Lyakh, Sergio Boixo, Hartmut Neven, Travis S Humble, Rupak Biswas, Eleanor G Rieffel, Alan Ho, and Salvatore Mandrà. 2020. Establishing the quantum supremacy frontier with a 281 Pflöp/s simulation. *Quantum Science and Technology* 5, 3 (apr 2020), 034003. <https://doi.org/10.1088/2058-9565/ab7eeb>
- [43] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C. J. Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, Aditya Vijaykumar, Alessandro Pietro Bardelli, Alex Rothberg, Andreas Hilboll, Andreas Kloeckner, Anthony Scopatz, Antony Lee, Ariel Rokem, C. Nathan Woods, Chad Fulton, Charles Masson, Christian Häggström, Clark Fitzgerald, David A. Nicholson, David R. Hagen, Dmitrii V. Pasechnik, Emanuele Olivetti, Eric Martin, Eric Wieser, Fabrice Silva, Felix Lenders, Florian Wilhelm, G. Young, Gavin A. Price, Gert-Ludwig Ingold, Gregory E. Allen, Gregory R. Lee, Hervé Audren, Irvin Probst, Jörg P. Dietrich, Jacob Silterra, James T. Webber, Janko Slavič, Joel Nothman, Johannes Buchner, Johannes Kulick, Johannes L. Schönberger, JoséVinicius de Miranda Cardoso, Joscha Reimer, Joseph Harrington, Juan Luis Cano Rodríguez, Juan Nunez-Iglesias, Justin Kuczynski, Kevin Tritz, Martin Thoma, Matthew Newville, Matthias Kümmerer, Maximilian Bolingbroke, Michael Tartre, Mikhail Pak, Nathaniel J. Smith, Nikolai Nowaczyk, Nikolay Shebanov, Oleksandr Pavlyk, Per A. Brodtkorb, Perry Lee, Robert T. McGibbon, Roman Feldbauer, Sam Lewis, Sam Tygier, Scott Sievert, Sebastiano Vigna, Stefan Peterson, Surhud More, Tadeusz Pudlik, Takuya Oshima, Thomas J. Pingel, Thomas P. Robitaille, Thomas Spura, Thouis R. Jones, Tim Cera, Tim Leslie, Tiziano Zito, Tom Krauss, Utkarsh Upadhyay, Yaroslav O. Halchenko, Yoshiki Vázquez-Baeza, and SciPy 1.0 Contributors. 2020. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature Methods* 17, 3 (2020), 261–272. <https://doi.org/10.1038/s41592-019-0686-2>
- [44] W. Jakob, J. Rhineland, D. Moldovan and others. 2017. pybind11 – Seamless operability between C++11 and Python.
- [45] Sam Westrick, Pengyu Liu, et al. 2024. GraFeyn: Efficient Parallel Sparse Simulation of Quantum Circuits. In *IEEE Quantum Week (QCE)*.
- [46] Stefan Woerner and Daniel J. Egger. 2019. Quantum risk analysis. *npj Quantum Information* 5, 1 (feb 2019). <https://doi.org/10.1038/s41534-019-0130-6>

[47] Hongzheng Zhao, Marin Bukov, Markus Heyl, and Roderich Moessner. 2023. Making Trotterization adaptive and energy-self-correcting for NISQ devices and

beyond. arXiv:2209.12653 [quant-ph] <https://arxiv.org/abs/2209.12653>