# Heuristic Quantum Advantage with Peaked Circuits

Hrant Gharibyan[1], Mohammed Zuhair Mullath[1], Nicholas E. Sherman[1], Vincent P. Su[1],
Hayk Tepanyan[1], and Yuxuan Zhang[2,3]

[1]BlueQubit
San Francisco, CA 94105, USA
[2]Department of Physics, University of Toronto
Toronto, ON M5S 1A7, Canada
[3]Vector Institute for Artificial Intelligence
Toronto, ON M5G 0C6, Canada

October 31, 2025

**Abstract**

We design and demonstrate heuristic quantum advantage with peaked circuits (HQAP circuits) on Quantinuum's System Model H2 quantum processor. Through extensive experimentation with state-of-the-art classical simulation strategies, we identify a clear gap between classical and quantum runtimes. Our largest instance involves all-to-all connectivity with 2000 two-qubit gates, which H2 can produce the target peaked bitstring directly in under 2 hours. Our extrapolations from leading classical simulation techniques such as tensor networks with belief propagation and Pauli path simulators indicate the same instance would take years on exascale systems (Frontier, Summit), suggesting a potentially exponential separation. This work marks an important milestone toward verifiable quantum advantage, as well as providing a useful benchmarking protocol for current utility-scale quantum hardware. We sketch our protocol for designing these circuits and provide extensive numerical results leading to our extrapolation estimates. Separate from our constructed HQAP circuits, we prove hardness on a decision problem involving generic peaked circuits. When both the input and output bitstrings of a peaked circuit are unknown, determining whether the circuit is peaked constitutes a QCMA-complete problem, meaning the problem remains hard even for a quantum polynomial-time machine under commonly accepted complexity assumptions. Inspired by this observation, we propose an application of the peaked circuits as a potentially quantum-safe encryption scheme [1–4]. We make our peaked circuits publicly available and invite the community to try additional methods to solve these circuits to see if this gap persists even with novel classical techniques.

# Contents

# 1 Introduction

Random circuit sampling (RCS) [5–7] has emerged as a primary benchmark for demonstrating quantum computational advantage, exemplified by Google's 2019 Sycamore experiment [8] and later followed by USTC Zuchongzhi works [9, 10]. However, recent large-scale demonstrations from Google's Willow [11] and Quantinuum's H-2 [12], while likely beyond classical simulation, expose a fundamental limitation of random circuit sampling: verification through cross-entropy benchmarking becomes exponentially expensive at large scales, rendering direct validation impossible.

We present a scalable method for demonstrating heuristic quantum advantage using peaked quantum circuits [13], circuits engineered to produce specific output bitstrings with anomalously high probability compared to an truly random circuit but otherwise look random. The verification protocol is straightforward: Alice constructs a peaked circuit knowing its peak bitstring, Bob executes the circuit with his quantum computer measuring in the computational basis, and Alice verifies by checking if Bob's output matches the known peak. No exponential classical computation is required.

Our contributions are threefold: (i) We develop scalable construction techniques for peaked circuits at larger qubit counts and depths using methods related to identity obfuscation [14–16]. (ii) Through extensive benchmarking, we demonstrate that our circuits resist state-of-the-art classical simulation methods including matrix product states [17–20], tensor networks with belief propagation [19, 21], and Pauli path simulators [22–29]. (iii) We motivate this construction further with hardness results for a decision problem involving general peaked circuits. In particular, we prove that determining whether an unknown circuit is peaked constitutes a QCMA-complete problem. Note that this is separate from our definition of "solving" our HQAP circuits, where the input string is fixed to the all zeros state.

Peaked circuits broadly have mounting evidence for their hardness. Simulating random peaked circuits is hard in the average case, analogous to RCS, particularly for circuits generated via post-selection [30] and error correction codes [31]. Our theoretical results in this work further add to the literature on hardness when the quantum circuit is peaked from a specific input bitstring to a particular output bitstring. However, our primary contribution lies in demonstrating the empirical difficulty of classically simulating circuits that not only admit efficient construction, but can also be implemented on existing quantum hardware. The relationship between our practical construction and the broader ensemble of peaked circuits remains an open question. In the absence of this theoretical connection, we refrain from asserting strong classical intractability results, instead characterizing the observed quantum-classical separation as a *heuristic* quantum advantage.

The rest of the paper is organized as follows. In Section 2, we summarize the empirical runtime gap between quantum and classical resources needed to solve our peaked circuits. In Section 3, we describe briefly the techniques used to construct these peaked circuits. In Section 4, we provide additional background and details on the classical methods we implement for comparison. In section 5, we study a potential practical application of peaked circuits as a post-quantum cryptographic primitive. We conclude in Section 6 and invite others to test other ways of finding the peak in Section 7.

## 2   A Gap in Classical and Quantum Runtime for Solving Peaked Circuits

We start with our main result, which is an empirical gap in the time it takes to solve peaked circuits that we have constructed. Though our methods can generalize to even larger circuits, we focus on what is achievable on existing hardware, in particular Quantinuum's H-2 processor [32] with 56 qubits and any-to-any connectivity. We create peaked circuits with 56 qubits and two-qubit gate counts ranging from 200 up to 2000. The quantum run time is approximately a few hours, while classical methods fail to extract the correct answer at around 700 RZZ gate circuits with reasonable computational resources. We estimate the largest circuits would require millions of GPU hours.

From Ref. [13], we take the following definition.

**Definition 1.1** (Peaked Circuit). *Given $\delta \in (0, 1]$, we call the unitary $C$ $\delta$-peaked if:*

$$\max_{s \in \{0,1\}^N} |\langle s|C|0^N\rangle|^2 \geq \delta$$

*with a corresponding peak weight $\delta_{\boldsymbol{s}} \equiv |\langle \boldsymbol{s}|C|0^N\rangle|^2$.* where **s** is the peaked bitstring. To "solve" a peaked circuit means given $C$, e.g. a circuit description in terms of gates, one correctly identifies the peaked bitstring **s**. For sufficiently small systems, one can compute the full statevector and compute the maximal amplitude squared but this clearly does not scale well with system size. With access to a quantum computer, one can simply run the circuit an $O(1/\delta_{\mathbf{s}})$ times to observe a signal. Of course, in the presence of noise, one may encounter an overhead in the required number of shots.
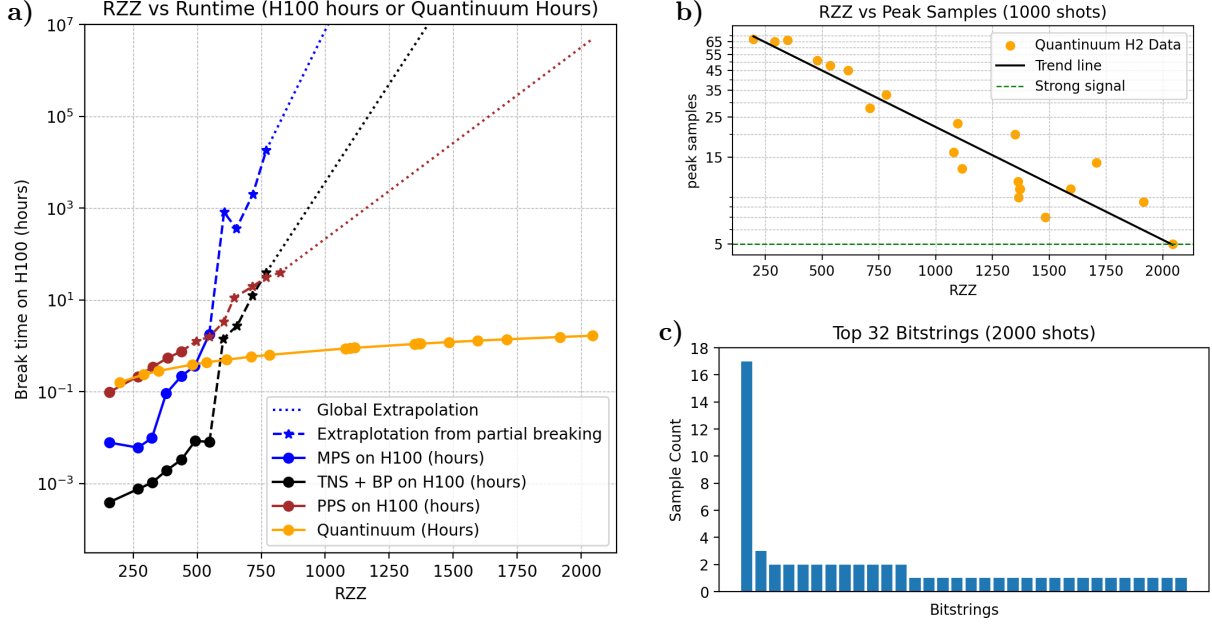
Figure 1: **Gap in Classical and Quantum Runtimes for Solving 56 qubit HQAP Circuits**. **a)** We compare classical techniques (MPS, PPS and TNS+BP) with quantum techniques for solving 56 qubit, 10% peaked circuits with all-to-all connectivity. As number of RZZs increases, classical methods require estimated millions of Nvidia H100 GPU hours, while the quantum device requires less than 2 hours. More details on the classical methods can be found in section 4 **b)** The performance of the quantum device drops with the growing number of RZZs in the HQAP circuit as expected, with still strong signal of 5/1000 (correct) peak samples for a 2044 RZZ circuit. **c)** Actual results from a single run on Quantinuum H2 for a 1917 RZZ peaked circuit, with 17/2000 correct peaked bitstrings. To our knowledge no other device in the world can find the peak for HQAP circuits of this scale and connectivity.

The main technical contribution of this paper is to create trainable circuits which go beyond the qubit count and depth of the results trained in Aaronson and Zhang [13]. Additionally, while the method they found had $\delta_{\mathbf{s}}$ which seemed nearly universal based on their training strategy, we have greater flexibility in tuning the peak weight. We call these HQAP (Heuristic Quantum Advantage Peaked) circuits as they demonstrate a large heuristic gap between classical and quantum runtimes for solving them. More on the protocol of how we generate HQAP circuits can be found in Section 3.

The size of the peak presents an interesting trade-off. Increasing $\delta_{\mathbf{s}}$ means fewer shots would be required, reducing the quantum runtime. On the other hand, one may worry that tuning it to be too large may lead to classically exploitable structure in the circuit description. This is an open question, and we invite curious researchers to participate in our open challenge (see Section 7). Without that exploitable structure, one is forced to deal with generic approximate classical simulation methods. Some of the leading ones include Matrix Product State (MPS) simulators [33], tensor network with belief propagation (TN + BP) [19, 21], and Pauli path simulators (PPS) [22–29]. For more details on our implementations and attack strategies, see Section 4.

With the explicit goal of constructing the largest possible separation in classical and quantum runtimes, we construct HQAP circuits that are designed to fit within the tolerance of existing quantum hardware and find an extreme gap. The largest such instance features 56 qubits, 2044 two-qubit gates, all-to-all connectivity and a peak weight $\delta_{\mathbf{s}} \approx 0.1$. With these parameters, Quantinuum's

H2 device can consistently detect the peak in under 2 hours with 1000 shots, assuming an average speed of 3000 HQC/hour.

On the classical side, we find that state-of-the-art simulation methods are only able to solve the circuits up to $\sim 700$ two-qubit gates within 10 hours. For HQAP circuits smaller than that, we find a nice extrapolation of runtime with two-qubit gate count. Assuming such larger computations were even feasible (e.g. with infinite memory and millions of CPUs or GPUs), our benchmarks estimate that solving the largest circuit we ran on Quantinuum would take years even on a Frontier scale supercomputer. More details on the classical methods used can be found in the section 4.

Again, though these circuits lack any theoretical guarantees, we claim that this is a form of *heuristic* quantum advantage, where there is an efficient verification protocol for the quantum device, and as far as we know, classical methods fail to produce the same answer.

# 3   HQAP Circuit Protocol

The protocol to design the peaked circuits in this work is based on the construction discussed in Ref. [13]. In that work, peaked circuits are composed of two main parts, a random quantum circuit $R$, and a variationally parmaterized circuit $P$ that serves as a *peaking layer*. Using gradient descent, the parameters of $P$ are trained so that the combined circuit, $R$ followed by $P$, results in a peak weight on a particular bitstring. When the depth of $P$ exceeds the depth of $R$, the peaking can be done trivially and to a peak weight of 1.[1] The intuition for using this structure is that $R$ mimics the structure of random quantum circuits, and thus may be hard to simulate classically, even if the final state is relatively simple to describe. This protocol has two main challenges when scaling to larger circuits. First, as the circuits become larger, the presence of barren plateaus [34] makes training difficult in general. Second, training the $P$ requires simulating the full circuit, which can become excessively expensive or even impossible in practice. These challenges limit the theoretical bound on the gap between classical and quantum runtimes possible.

Here, we improve upon this construction by extending both the depth and the width of such circuits, primarily enabled by tensor network methods. By utilizing generic tensor contraction functionality as in Quimb [35, 36], we can train $R$ and $P$ as in Ref. [13] for large qubit number as long as they are sufficiently shallow. Such shallow peaked circuits are known to be easy to solve [37]. To increase the depth, we introduce a large, obfuscated identity block between $R$ and $P$. In the worst case, it has been shown that checking whether a quantum circuit is equivalent to the identity is a QMA-complete problem [14], even when the circuit depth is shallow [38]. This indicates that identity obfuscation could be a feasible direction for generating peaked circuits that cannot be simplified by an attacker. In practice, to make the structure harder to reverse engineer, we employ a technique we call *Tensor Patch Optimization*, which serves two purposes, patch sweeping and patch masking, which will be explained shortly.

**A word on notation.** To enhance clarity in representing the sequential application of quantum operations, we introduce the left-to-right composition operator $\triangleright$, defined by:

$$U_1 \triangleright U_2 \triangleright U_3 := U_3 U_2 U_1$$

That is, $U_1$ is applied first, followed by $U_2$, then $U_3$. This reverses the standard right-to-left ordering of matrix multiplication to improve readability when reasoning about temporal circuit structure.

---

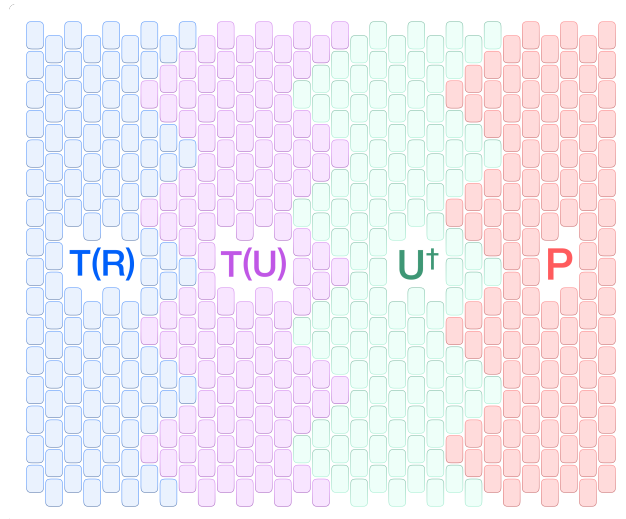[1] Assuming that the connectivity of $P$ matches that of $R^\dagger$

Figure 2: Visual schematic of HQAP building. We include details of the circuit manipulations used to scramble the circuit structure.

Note that the identity block $U \triangleright U^\dagger$ implemented naively as gate-by-gate inversion is easy to detect, so we obfuscate this exact structure. The full protocol to generate peaked circuits is illustrated in Figure 2, and is summarized as

- Train a shallow peaked circuit $R \triangleright P$ using the method in Ref. [13].

- Insert an identity of the form $U \triangleright U^\dagger$ between $R$ and $P$.

- Apply sweeping to variational circuit parameters.

- Apply masking to modify the connectivity and gate structure.

- Apply swap transformations throughout the circuit.

The sweeping and masking, whose definition we will return to momentarily, can be done modularly throughout the circuit. However, the swapping must be coordinated globally to preserve the action of the unitary on the initial state. Let us denote the collective action of swapping, sweeping and masking as $T$, then the final circuit has the form $T[R] \triangleright T[U] \triangleright U^\dagger \triangleright P$. Note that the last few circuit deformation techniques can be applied in any order. Further details on the obfuscation strategies in the protocol are discussed in the following sections.

For our construction there exist, in principle, simplifications that may reduce the simulation requirement. With knowledge of how they were constructed and subsequently modified, one can trace out the simplifications. However, we expect that the series of changes are likely hard to reverse. We have internally tested a number of strategies for attacking these circuits, which led us to the current set of techniques we describe below. Abstractly, these can be thought of as methods to make transpilation or finding an efficient contraction path of the quantum circuit more difficult. We liken our scrambling techniques (swaps, sweeps, masks) to how hash functions build one-way transformations. Our expectation is that while its non-trivial to prove how our scrambling techniques make the final circuit irreversible to the original circuit, a lack of successful attempts to crack them will eventually build conviction that our scrambling is in practice irreversible. That would be analogous to how hash functions like SHA-1 and SHA-2 build irreversible 1-way protocols based on

many small iterations. Similarly, SHA-1 and SHA-2 lack rigorous proof of irreversibility. However, the huge body of evidence shows that for all practical purposes they are hard to reverse. For this same reason, we invite everyone to try to unscramble and break our HQAP circuits, see Section 7.

## 3.1   Tensor Patch Optimization

The idea of tensor patch optimization is to take a circuit and train a second variational circuit which matches the unitary action of the first. More procedurally, one can isolate subcircuits, or subsets of gates, and by finding an approximate substitute, locally deform them. A similar idea to this type of obfuscation was discussed in Ref. [16], there called circuit unoptimization. These subcircuits are unitary and can generally be treated as tensors $\mathcal{T}_i$. The mathematical idea is to introduce a new trainable patch $\tilde{\mathcal{T}}(\theta)$ which is trained to mimic the original patch by taking the trace fidelity as a loss function

$$\mathcal{L}(\theta) = 1 - \frac{\left| \text{Tr}[\mathcal{T}_i^\dagger \tilde{\mathcal{T}}(\theta)] \right|}{d}, \tag{1}$$

where $d$ is the dimension of the full unitary $\mathcal{T}_i$. Since $\mathcal{T}_i$ and $\tilde{\mathcal{T}}(\theta)$ are unitary, this loss function is always in the interval $[0, 1]$, and is 0 if and only if $\tilde{\mathcal{T}}(\theta) = \mathcal{T}_i$ up to a global phase.

Though mathematically similar problems, there are two ways we employ this technique. The first we refer to as angle **sweeping**, where the variational patch has the same gates as the original patch, but the parameters are first kicked away from their original values and then trained to increase the fidelity. The second we call **masking**, where the trained patch actually has a different subcircuit structure, where the set of gates may be different or acting in a different order. Both offer a way to obfuscate structure that could potentially be used to simplify the circuit.

For the purposes of making the circuit harder to simulate, it can actually be beneficial to introduce a lossy approximation to the original patch. For example, when the loss function is small, but non-zero, one can no longer find an equivalence, thus making the transpiler's job harder.

## 3.2   Swap Transformations

The second obfuscation technique we employ is to apply a permutation enabled by a series of simple swap transformations between two qubit wires. Let $\sigma^{(i \leftrightarrow j)}$ denote the unitary operation that swaps the quantum states of qubits $i$ and $j$. We define the action of this transformation on a circuit $X$ as

$$\sigma^{(i \leftrightarrow j)}[X] := \sigma^{(i \leftrightarrow j)\dagger} X_\sigma \, \sigma^{(i \leftrightarrow j)}, \tag{2}$$

where $X_\sigma$ denotes the circuit $X$ with its qubit indices relabeled under the transposition $(i \leftrightarrow j)$. This can be interpreted either as a physical SWAP gate or as a virtual relabeling of the wire indices, depending on the context (e.g., compilation vs. simulation)

These swaps can be interspersed throughout the circuit as long as one is careful about appropriately implementing the swaps and relabeling throughout. Such techniques make it difficult to design tensor network ansatzes that retain small bond dimension. Additionally, it is useful to note that the initial state is invariant under swaps, allowing all swaps at the beginning of the final transformed circuit to be ignored.

# 4   Benchmarking Classical Approaches for Solving Peaked Circuits

In the context of quantum advantage, it is increasingly understood that compelling demonstrations must not only exploit quantum computational power but also move beyond the reach of classical

simulation. For sampling problems, this often means operating at a scale where even the most advanced classical techniques become infeasible due to computational or memory constraints. In particular, state-vector simulations, which scale exponentially in the number of qubits, are ruled out once the qubit count exceeds approximately 40. To test the classical hardness of our HQAP circuit constructions, we focus on scalable attack strategies that aim to approximately reproduce key features of the output distribution.

This section is broadly divided into three parts. We start by briefly mentioning our primary attack strategy, looking at marginal or single qubit expectation values. This, in general, may not always converge to the correct answer, but is a general purpose approach that is simple to understand and can work when there is sufficient weight in the peak bitstring. In the case of examples we have run on quantum hardware, where we want the peak to be large enough to be detectable, this is a reasonable strategy.

We subsequently explain the state-of-the-art classical methods used to implement this strategy. Details on simulation techniques are included. Because they fail to correctly output the target bitstring on moderately sized circuits with hundreds of two-qubit gates, we explain how our runtime extrapolations to the largest HQAP circuit sizes of $\sim 2000$ two-qubit gates are calculated.

Importantly, these benchmarks of scaling generally serve as a lower bound on the time needed to use these classical methods as an attacker. Because we built the circuits to encode a particular bitstring, we are able to search for the minimal computational requirements needed to reproduce that answer. However, from the perspective of the attacker, these methods do not readily come with convergence guarantees, apart from doing full exact simulation.

To give a sense of the classical computational resources used to reproduce these results, we spent thousands of A100 and H100 GPU hours and tens of thousands of CPU hours to solve the peaked circuits with the best known classical techniques. For statevector simulation we have been using qsim by Google on CPUs [39] and cuQuantum by Nvidia on GPUs [40]. For MPS, Tensor Network and Belief Propagation we used the quimb library on H100 GPUs [35, 36]. For Pauli path simulation we have been using a newly developed GPU implementation on H100 GPUs. All of the classical experiments have been orchestrated using the BlueQubit platform.

Finally, we explore solution strategies that exploit the structure of the circuit directly, bypassing the need for brute-force simulation. While our circuits do not yet guarantee security against all such attacks, we present evidence of baseline structural resistance.

## 4.1  Marginal attack strategies

Sampling from the output distribution of a general quantum circuit is widely believed to be classically intractable. In the average case, even approximate sampling to small total variation distance is believed to be classically intractable, and under plausible assumptions such as the non-collapse of the polynomial hierarchy, this forms the basis of most near-term quantum advantage proposals. For random peaked circuits, it has been shown that simulating the output distribution of peaked circuits to $1/\exp\{\mathrm{poly}(N)\}$ precision is average case #P hard; relaxing this precision to $1/\mathrm{poly}(N)$, in the worst case the problem is still BQP-complete, and for any sequential simulator it is average case BQP-complete under plausible complexity assumptions [30]. Similar hardness results have been observed in other types of peaked circuits [31]. Nevertheless, in this work there is a caveat one could leverage: our task is just to 'identify the peaked strings', leaving room for spoofing. Although exact evaluation of expectation values of arbitrary circuits is computationally prohibitive, estimating approximate expectation values, particularly of low-weight observables or marginal distributions, is often significantly more tractable for classical algorithms, especially when the circuit depth is moderate or its structure is constrained, such as in circuits with shallow magic depth [41].

The idea is relatively straightforward. If there is sufficient bias in the peak, then approximate expectation values, as long as they get the correct sign, can be used to attack or reconstruct what the peak bitstring is. Of course, as circuits get deeper, even getting approximate expectation values can be difficult. Indeed, for our largest circuits, our classical methods cannot solve them. Nonetheless, we introduce some of our notation that will be useful in describing their (lack of) success.

For sufficiently large peak weight $\delta_{\mathbf{s}}$, it is straightforward to observe that each of the bits will have marginal distributions $p(b_i)$ which are heavily biased towards the true value $s_i$. This is straightforwardly related to the $Z$ expectation value

$$p(b_i = 0) = \frac{1 + \langle Z_i \rangle}{2}. \tag{3}$$

Thus, by coming up with estimates for $\langle Z_i \rangle$, e.g. with approximate tensor network methods, one can construct a candidate bitstring $\hat{s}$ by looking at the sign of each $\langle Z_i \rangle$. We refer to this as the marginal attack. From an attacker's perspective, this is not guaranteed to work. Indeed, one can manipulate the peak weight and the distribution so that marginals, even if calculated correctly, would give the wrong result. However, for the purposes of the classical benchmarks we present here, we have constructed the circuits so that the marginals are correctly lined up with the target bitstring.

To keep track of the accuracy, one can consider the overlap (or inverse hamming distance) $O = N - |s - \hat{s}|_1$, or similarly, the fraction of correct bits $R = O/N$. Note that for a random guess one expects $R = 0.5$, and $R = 1$ is only achieved if $\mathbf{s} = \hat{\mathbf{s}}$.

Interestingly, this attack strategy can also be used when looking at empirical or approximately generated samples, and corresponds to majority vote on each qubit.

For the largest circuits, our experiments are dominated by noise with minimal or no signal at all $R \lesssim 0.5$. For this reason, we study circuits of different sizes and compute $R$. We then define $T_{\text{break}}$ as the minimum time needed to achieve $R = 1$, or a quantity such as $\chi_{\text{break}}$ defined as the smallest bond-dimension needed to achieve $R = 1$ for MPS and TNS simulations. Since the runtime and accuracy for simulating a circuit is fully characterized by the bond-dimension $\chi$, there is a direct relationship between $t_{\text{break}}$ and $\chi_{\text{break}}$.

One can play a similar game with access to multi-qubit observables, though these in general will not be any easier to calculate than single qubit observables. When dealing with samples, whether approximate or empirical, one can implement many classical post-processing schemes. They may consider correlations or do majority votes on subsets of bits rather than single bits. However, finding an optimal strategy when there are conflicts is also very difficult.

## 4.2 Matrix Product State Simulations

MPS has revolutionized the study of 1D quantum systems [33]. Its power for capturing and studying entanglement has thus become a standard tool in the kit of many quantum practitioners. MPS allows a dial between efficiency and expressibility through a simple parameter, the bond dimension $\chi$.

Our peaked circuits are explicitly designed with all-to-all connectivity. Thus, one would not expect MPS to be particularly well suited to simulate them. However, they serve as a baseline simulation method for understanding the classical difficulty. In particular because producing the correct the peak bitstring could be spoofed by a low rank MPS that is far from the true final state. One could worry that basic MPS simulation may crack the circuits. We find this naive method works up to a point. Beyond 600 two-qubit gates, however, MPS fails to predict the correct bitstring.

Because MPS works best for states with area law entanglement in 1D, the explicit transpilation of a circuit with long-range entangling gates to a chain layout can greatly affect the performance.
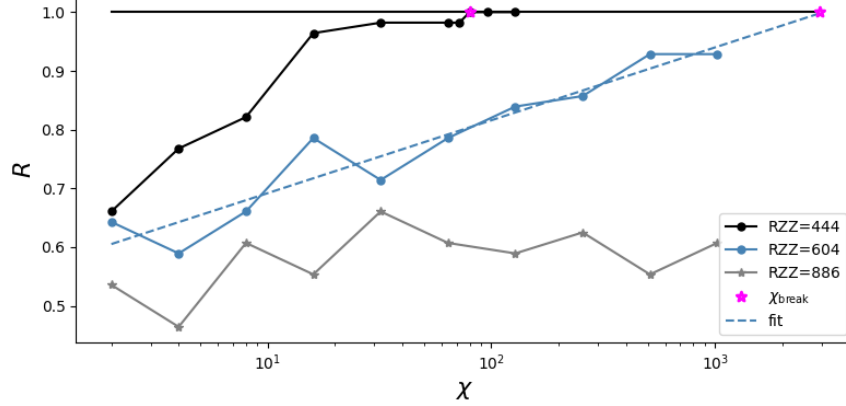
Figure 3: Estimating $\chi_{\text{break}}$. By tracking the accuracy $R$ as a function of $\chi$, we find 3 cases. The black curve is an example of a circuit where $\chi_{\text{break}}$ is found exactly. For the blue curve, the $\chi_{\text{break}}$ is not found exactly, and so we fit the data to find $\hat{R}(\chi)$ and extrapolate to estimate $\chi_{\text{break}}$, and show the fit with a dashed line. Lastly, the gray curve is an example of a circuit our fitting protocol rejects when estimating $\chi_{\text{break}}$.

In the worst case, a naive permutation of an otherwise local 1D chain can make MPS simulation struggle. To define the ordering of the qubits, we look at the adjacency matrix for our quantum circuit and use the Reverse Cuthill-McKee (RCM) algorithm [42] to define a global permutation to our circuit. The RCM algorithm defines a reordering of an adjacency matrix to minimize the bandwidth of the matrix and hence the maximum distance between qubits that interact.

To produce extrapolation estimates for our largest circuits, we study the relationships between bond dimension, circuit depth, simulation time, and accuracy of the predicted answer. We define $\chi_{\text{break}}$ as the minimum bond-dimension needed to yield the correct bitstring for a given circuit. For small circuits, this is found explicitly through simulation. For larger circuits, this can be extrapolated by looking at the trend of $R$, the correct fraction of bits, as a function of increasing $\chi$.

For a given circuit, we perform MPS simulation with bond dimension increasing in powers of 2. If the circuit is cracked, we perform binary search to find $\chi_{\text{break}}$. The maximum $\chi$ used in our experiments varies by circuit with $\chi = 5000$ the largest bond-dimension used. If the circuit remains uncracked for all $\chi$ studied, we perform an automated fitting procedure to determine $\hat{R}(\chi)$. We assume a linear fit in $\log \chi$, weighted more heavily towards large values of $\chi$ and reject fits that yield a correlation coefficient of $r^2 < 0.8$, or those that show a negative slope. We show examples of $R(\chi)$ for three circuits, one which was cracked exactly, one which we successfully extrapolate $R(\chi_{\text{break}})$, and one which we reject.

Our peaked circuit construction allows us to control the RZZ gate count. However, there is some variability in the resulting $\chi_{\text{break}}$ due to the stochastic nature of the training. Thus, for a given circuit depth, we apply this extrapolation procedure to five instantiations, reporting both the average $\chi_{\text{break}}$ as well as the standard deviation in the log. For accessible system sizes, we observe exponential scaling in the necessary bond dimension as a function of circuit depth, as might be expected for random circuits before the bond dimension saturates at its maximal value $\chi_{\text{sat}} = 2^{N/2} \simeq 10^8$ [33]. See Fig. 4. Interestingly, our extrapolated answer for $\chi_{\text{break}}$ on our deepest circuits reaches beyond this $\chi_{\text{max}}$, giving evidence for the difficulty of MPS to crack our system sizes.

Ultimately, one cares not only about the memory requirements determined by $\chi$ but also the total runtime of the classical method. To this end, we also study the time $T$ to simulate these circuits
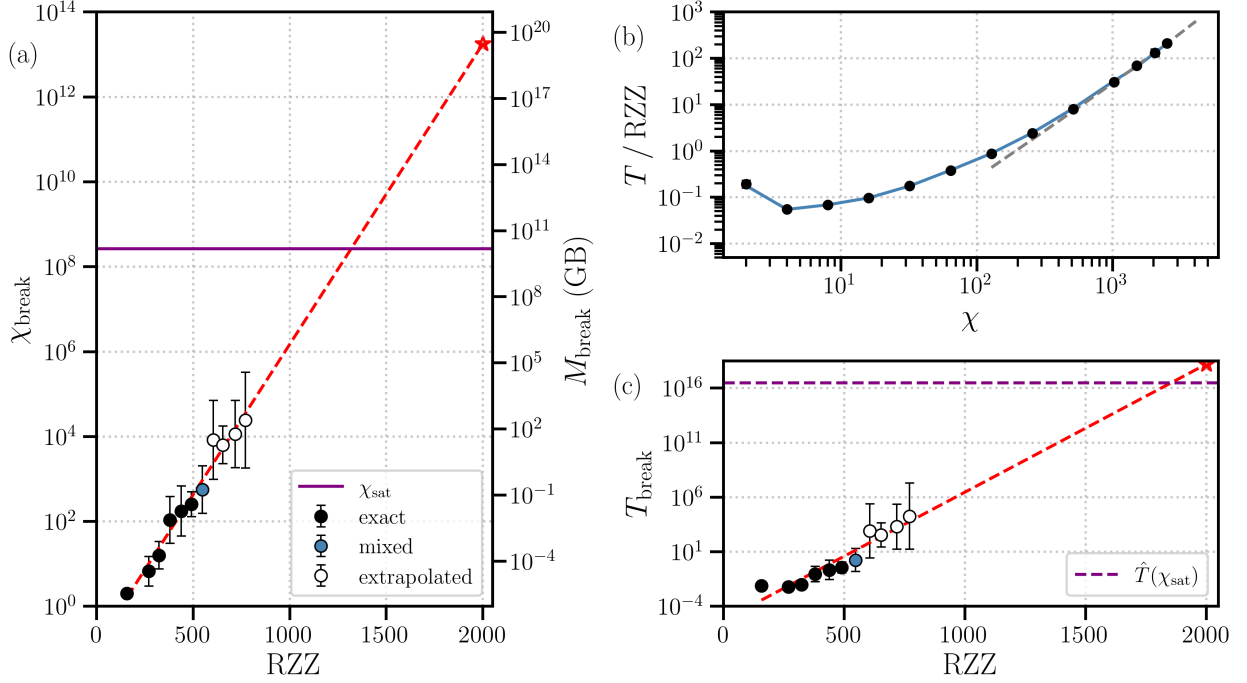
Figure 4: Benchmarking of bond dimension and time required to simulate peaked circuits of increasing depth with MPS. (a) Bond dimension vs two-qubit gate count. Black dots indicate MPS simulations that converged exactly, while white dots indicate circuit samples where the extrapolation procedure in Fig. 3 was used. Blue dots indicate the transition where some circuits of that size were fully solved. The red dashed line shows the exponential fit, while the solid purple line shows the $\chi$ corresponding to exact simulation. (b) Gate simulation cost for increasing $\chi$. We study how the simulation time grows with increasing bond dimension, normalized by the number of gates. The large $\chi$ behavior is modeled by by a cubic function of $\chi$. (c) We plot the solution time $T_{\text{break}}$ in hours as a function of circuit depth. Exact answers (black dots) can be read off, while white dots are extrapolated by using $\hat{\chi}_{\text{break}}$ and the time conversion in (b). The red dashed line is an exponential fit for $T_{\text{break}}$ compared with circuit depth. The purple dashed line is the estimate for the runtime of a depth 2000 circuit with $\chi_{\text{sat}}$.

as a function of $\chi$, normalized by the circuit depth. The runtime is dominated by the operations of applying a gate which scales as $\chi^2$, and the subsequent SVD which scales as $\chi^3$. We use this fact to fit our data to a functional form $T/RZZ \sim a\chi^2 + b\chi^3$, and show the fit in Figure 4. Combined, these give a way to extrapolate the total runtime needed for our deepest circuits. Assuming one even has access to enough memory to run the simulation, we can use this relation to estimate $T_{\mathrm{break}}$ based on $\chi_{\mathrm{sat}}$ for our largest circuits.

## 4.3 Tensor Network State and Belief Propagation Simulations

Despite the success of MPS simulations in studying quantum dynamics in one-dimensional systems with local interactions, systems with higher-dimensions and long-range interactions have still remained a major hurdle in the field. A common approach has been to simply map such systems to quasi-one-dimensional systems at the cost of long-range interactions, as discussed in Section 4.2. Other tensor network methods such as projected entangled-pair states (PEPS) [18, 43] have been popular for studying the many-body dynamics of two-dimensional lattices, such as for a square lattice, and easily extend to arbitrary geometries. However, extracting useful quantities, such as local observables, from the final state remains difficult in general due to the complexity of the resulting tensor network contraction. The complexity of local measurements disappears in MPS systems by using the gauge freedom of the MPS ansatz to bring an MPS state into a canonical form [33]. Once in this form, local observables can be computed by a tensor contraction involving $\mathcal{O}(1)$ site tensors rather than a full network contraction.

Recently, belief propagation (BP) has been proposed as a method to gauge a generic Tensor Network State (TNS) to address this problem [19, 21]. BP is a set of equations for *message tensors* that can be solved self-consistently, whose fixed point defines a gauge transformation that brings the TNS into the Vidal gauge [44–47]. This procedure is known to be exact for tensor networks on a tree graph, and while approximate for loopy graphs, often leads to accurate results in practice.

Despite the benefits of TNS evolution and BP, this framework requires performing matrix decompositions on tensors, like QR or SVD, that scale as $\chi^z$ for geometries with a coordination number $z$. For two-dimensional systems with $z \simeq 4$, this scaling is manageable enough to lead to dramatic improvements in both accuracy and efficiency over traditional MPS-based simulations. However, for circuits with all-to-all connectivity, like the circuits studied here, the scaling becomes infeasible as $z \to N$ for a naive use of this method.

To remedy this, we define a simpler topology with a reduced coordination number, and then transpile the circuit to this simpler topology. This can be thought of as a generalization of what is done for MPS-based simulations, where the simpler topology is a chain. Once transpiled, we then perform the simulation using a generic TNS that matches the topology of the resulting circuit. With this increased flexibility comes the difficulty of picking a particular layout. There is a tradeoff between a simpler topology with a reduced $z$, allowing for more efficient tensor operations, and the increase in gate count from the added swap gates to perform the routing.

The limitless options for both topology and bond dimension of the ansatz are further complicated by the use of BP at the end to compute the one-qubit marginals. For example, when using a $k$-regular graph with $3 \leq k \leq 6$, we found surprisingly that increasing $\chi$ for TNS states can actually lead to *lower* accuracy with respect to the true expectation values. This is in sharp contrast with the MPS case where performance is expected to be monotonic in $\chi$. We believe this is due to the increase in long-range and loopy correlations throughout circuit simulation, which can spoil the accuracy of BP for computing the marginals [19, 21].

Given the stark effect of topology on the quality of the final answer, we found the most successful results using a tree tensor network. Notably, the lack of loops in the tree network satisfy a core
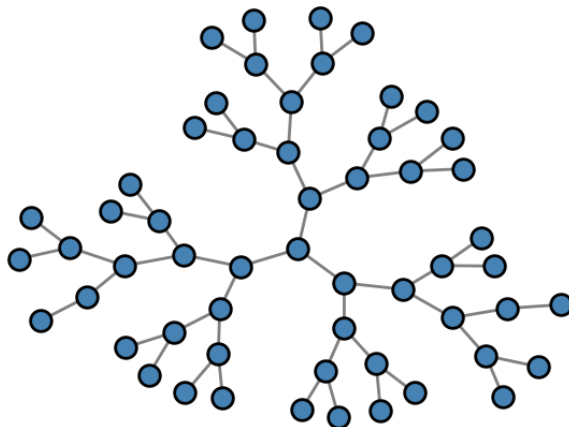
Figure 5: Example of the tree-like graph that we transpile our all-to-all circuits to before performing TNS simulations. The lack of loops lead to confidence in the convergence of BP for computing marginals in the tensor network state.

assumption of BP. This led to the desirable property that increasing bond dimension did lead to more accurate marginals. In the tradeoff of gate transpilation and topology, we found reasonable results using a graph with three binary trees stemming from the root. All reported TNS simulations we benchmark use this graph. See Figure 5 for an illustration.

We perform the transpilation using the Qiskit package [48], and the gate evolution was done using the simple update algorithm [49, 50] with the 'reduced tensor' method [21]. Once we have an approximate TNS representing the final state of the circuit, we perform BP to bring the TNS into an approximate canonical form, allowing us to efficiently approximate $\langle Z_q \rangle$ for each qubit $q$.

We then apply the marginal attack strategy discussed in Section 4.1. We also considered a greedy, iterative approach to generate a proposed bitstring $\hat{s}$. Rather than computing all marginals independently, one first selects the most biased marginal, and then fixes the tensor network based on the most likely bit. After fixing that bit, we re-run BP and compute the remaining marginals, conditioned on the outcome of the previous bits, and continue this process until all bits are fixed, defining $\hat{s}$. We then find $R$ of correct bits for both procedures, and use whichever method yielded a more accurate proposal.

Analogous to our MPS simulations, $\chi_{\text{break}}$ is determined by $R(\chi_{\text{break}}) = 1$. For sufficiently small circuits, we can compute $\chi_{\text{break}}$ exactly, while for larger circuits we rely on fitting $\hat{R}(\chi)$ and extrapolate to approximate $\chi_{\text{break}}$. In Figure 6, we show $\chi_{\text{break}}$ for peaked circuits with $N = 56$ and a range of RZZ gates.

Lastly, we create a conversion from $\chi_{\text{break}} \to T_{\text{break}}$ in the same way as we did for the MPS simulations described in Section 4.2. The runtime is dominated by the operations of applying a gate which scales as $\chi^2$, and the subsequent SVD which scales as $\chi^3$. We use this fact to fit our data to a functional form $T/RZZ \sim a\chi^2 + b\chi^3$. This fit allows us to extrapolate for larger $\chi$ to estimate the runtime for larger circuits. We note that from our fit of $T(\chi)$, the large $\chi$ scaling is only present for $\chi \gtrsim 256$. This leads to exponential scaling for $T_{\text{break}}(RZZ)$ to only be visible for $RZZ \gtrsim 550$. Due to limited data in this regime, the estimated $T_{\text{break}}$ should be treated as a rough estimate.
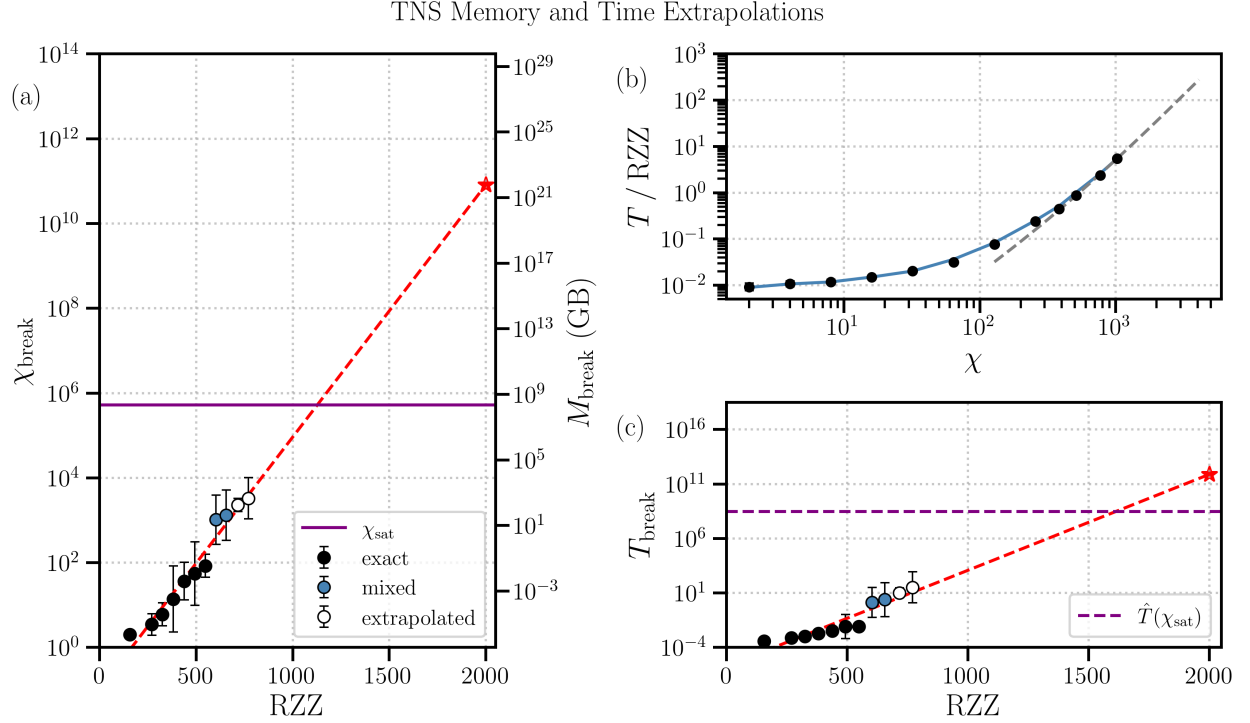
Figure 6: Memory and time extrapolations for TNS simulations. Notably, the TNS can crack circuits much more efficiently than in the case of MPS. However, we still find an exponential trend in the memory and time requirements. In the case of the network comprised of 3 binary trees, $\chi_{\mathrm{sat}}$ is determined by the size of the largest binary tree in Fig. 5. Though the $\chi_{\mathrm{sat}}$ is lower than in the MPS case, the required memory is comparable. Using the extrapolations with $\chi_{rmsat}$, we estimate that $T_{\mathrm{break}}$ for TNS simulations would take about $3 \times 10^8$ hours for a 2000 RZZ circuit, assuming the calculation would even fit in memory.

## 4.4 Pauli Path Simulations

Pauli Path Simulation (PPS) is a relatively new method to classicaly compute observable expectations through backward evolution of the circuit gates on the Pauli basis [22–29]. We leverage a state of the art GPU-accelerated version of this to implement the marginal, or $\langle Z \rangle$, attack.

A counterintuitive finding of Ref. [28], is that the convergence patterns of observables as a function of coefficient truncation parameter $\Delta$ is highly unpredictable. In the context of solving peaked circuits, this makes it difficult to predict, even for one $\langle Z_i \rangle$, what value of $\Delta$ will be sufficient to get an accurate expectation value. While intuitively one might expect smaller $\Delta$ should yield a more accurate answer, the deviations are non-monotonic, and thus a smaller $\Delta$ may actually be less accurate than a larger $\Delta$.

The way we then choose to extrapolate the difficulty is by looking for convergence of the expectation value of each $\langle Z_i \rangle$ *to the correct sign* and, hence, bit value. This assumption is then a *lower bound* on the time it would take to crack with this strategy since it equips the attacker with the most information possible. As in the MPS case, we are using our knowledge of the true answer for extrapolation purposes. A priori, without the ground truth, confidence in the convergence can be very difficult to establish [28].

To be slightly more precise, we compute $\langle \hat{Z}_i \rangle$ for each qubit with geometrically decreasing truncation parameters $\Delta_j$. We check that i) the sign has converged for 3 successive $\Delta_j$ and is the correct value and ii) that the expectation value exceeds a certain fraction of the operator weight that survives the truncation.

An interesting feature we notice when implementing this strategy is that there is generically a large gap between qubits that are easy to crack and those that are difficult. As can be seen from Figure 7 some qubits take more time to converge and crack, but the overall hardness grows exponentially with the number of 2-qubit gates in the circuit, as expected.
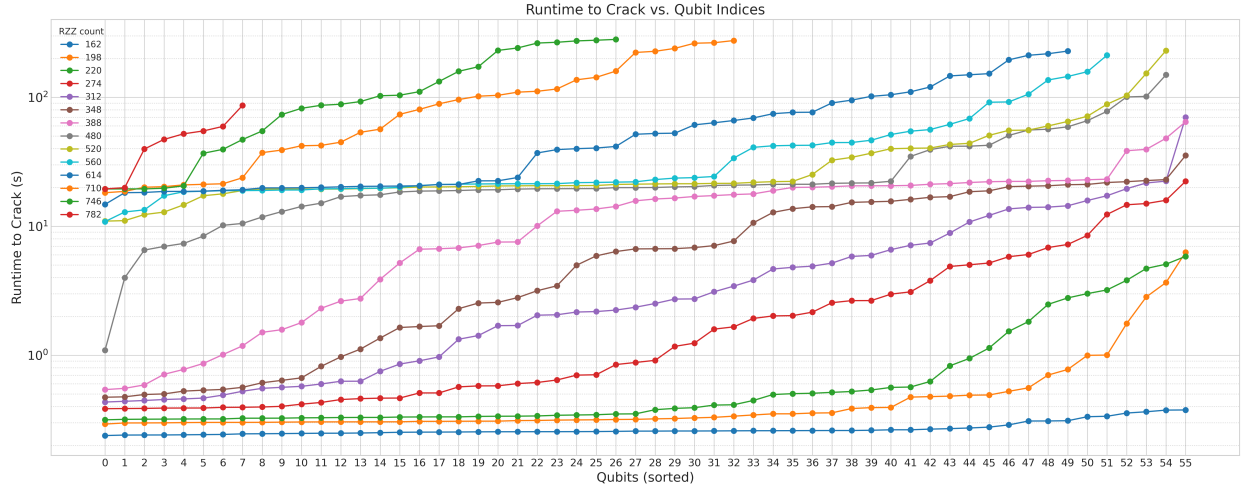


Figure 7: Required runtimes to crack individual bits of the peak bistring for 56-qubit HQAP circuit with varying RZZs going from 162 to 782 with PPS (standard coefficient truncation) on an H100 GPU. Qubit indices are ordered by difficulty (runtime to crack). We vary the truncation parameter $\Delta$ from $3.2 \times 10^{-3}$ to $2.5 \times 10^{-5}$ in powers of 2.
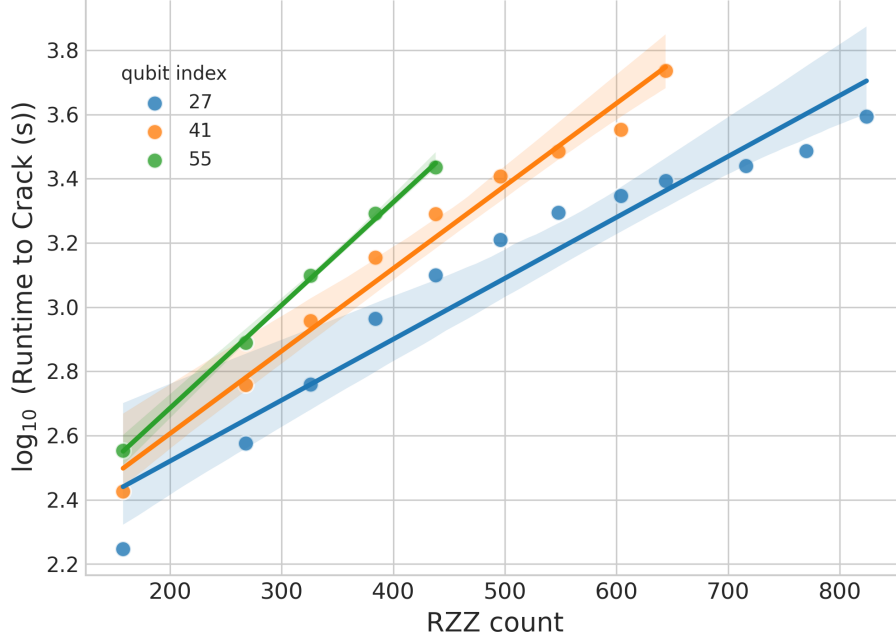
Figure 8: Plot showing exponential growth in runtime to crack bits of different difficulty levels for circuits with up to 824 RZZ gates. For this plot, we used the relative truncation method described in this section, and we tracked up to 1 billion Pauli terms, using $\sim$96 GB of memory on an H100 GPU. We note that solving qubits with index 55 amounts to fully solving the peaked circuit.

While PPS is very effective in shallower regimes and we can crack all 56 bits up to 388 RZZ HQAP circuit, memory very quickly becomes a constraint. Consequently, it was not feasible to explore truncation thresholds lower than $10^{-5}$ for these circuits, as the number of Paulis approached the memory limit ($\sim$1 billion Paulis). This is why the HQAP circuits after 388 RZZ have higher qubits missing from the plot in Figure 7.

One of the issues we encountered while running PPS with standard coefficient truncation for larger peaked circuits (with RZZ-count > 600) was that all the Pauli terms would get truncated midway through the simulation. The reason for this is the following: there are long stretches in the circuit where there is extensive branching and negligible merging of Paulis, resulting in rapid decay of all Pauli coefficients.

We briefly outline a modification to the coefficient truncation method devised to mitigate this issue and (potentially) yield a non-zero expectation value at the end of each PPS run. The modification, which we call *relative truncation*, is rather simple: after the $k$-th gate application, we drop only those Pauli strings $P$ whose coefficients $c_P$ satisfy

$$|c_P| \leq \|O_k\|\Delta,$$

where $\|O_k\|$ is the normalized Frobenius norm of the evolved Pauli-sum after $k$ steps of the simulation. This simple modification dynamically lowers the effective truncation threshold from $\Delta$ to $\|O_k\|\Delta$, and when more and more Pauli terms are dropped, resulting in a decline in $\|O_k\|$, we end up truncating fewer Pauli terms than we would otherwise have.

We also noticed an improvement in accuracy for the computed expectation values with this method, even after accounting for the fact that the new method utilizes more Pauli terms for a given threshold $\Delta$ compared to standard truncation.
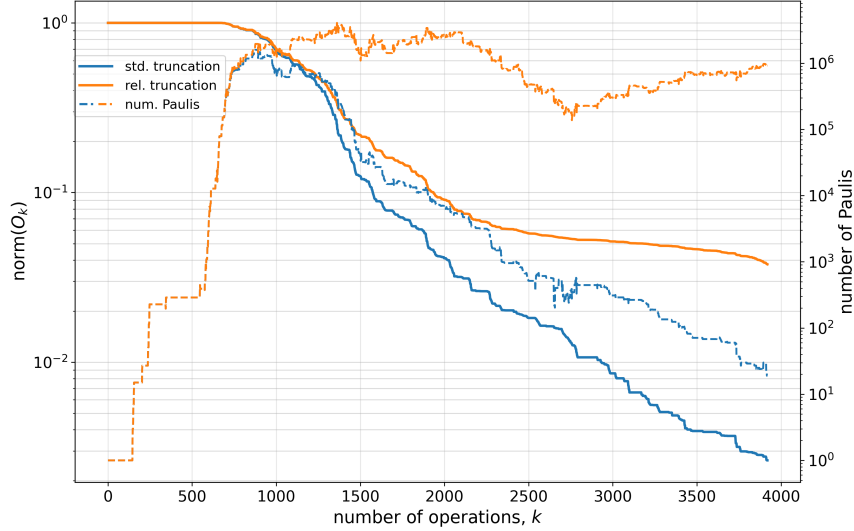
16

Figure 9: Comparison of change in observable norm $\|O_k\|$ and number of Pauli terms during the computation of $\langle Z_0 \rangle$ for both standard and relative coefficient truncation schemes for a circuit with 536 RZZ gates, and for threshold $\Delta = 2 \times 10^{-4}$. With standard coefficient truncation, we see that the number of Paulis decline rapidly in the latter half of the simulation, with only a couple of dozen terms remaining towards the end, whereas with relative truncation there are still around a million Pauli terms retained at the end, thereby improving the accuracy of the simulation.

## 4.5   Attacking the Circuit Structure

We have examined the computational requirements for direct classical simulation of the circuit. While this provides a useful baseline, we recognize that targeted attacks exploiting specific circuit structures could potentially reduce these computational requirements. We present an initial analysis of several such attack vectors and discuss how the techniques introduced in Section 2 - swapping, sweeping, and masking - may offer some protection against them.

To illustrate these concepts, we analyze a simplified test case: a circuit constructed as $U \triangleright U^\dagger$ (a random unitary circuit followed by its adjoint). This allows us to verify that our methods pass a baseline level of resistance to these attacks, though we emphasize that making stronger guarantees on the difficulty remains an open question.

**Correlating angles** For a random circuit that is comprised of multiple variational gates, the angles of gates in $U^\dagger$ will mirror those in the first half of the circuit. Thus, even though swaps were employed throughout, this would effectively allow someone to reverse engineer the resulting permutation. However, with tensor patch optimization used for the purpose of angle sweeping, we find that with a sufficient number of sweeps, one can hide the apparent correlations between different parts of the circuit. See Fig. 10.
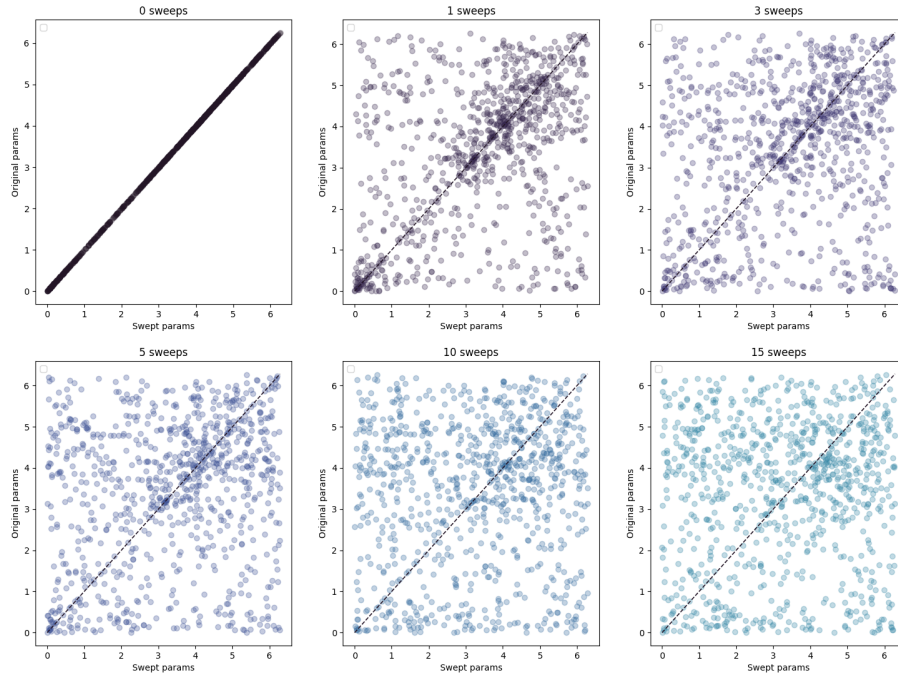
Figure 10: Angle sweeping. Taking a brickwork of two qubit gates and one qubit gates, we show the drift of angles from their original values after successive rounds of sweeping. The dashed indicates perfectly matching the original distribution.

**Correlating gate connectivity** Beyond the angles themselves, one may hope to unravel the circuit by looking at the gate connectivity in different parts of the circuit. The swaps and permutations make immediate cancellations less obvious, but one could try to reveal the permutations by identifying the relabeling. An attacker could potentially use the connectivity of the gates themselves, e.g. by looking at pairs of qubits that receive two-qubit interactions. Doing this identification would essentially be solving graph isomorphism, which is in NP, though not NP-complete. However, with the addition of masking, which replaces a subset of gates with a different set of gates, connectivity changes can be introduced, leading to an even more difficult version of the problem.

**Transpiler simplifications** A primary sanity check of these techniques is to make transpilers, such as the one in Qiskit [48], struggle to identify and simplify the identity block. To illustrate the efficacy of the compiler, we perform a simple test in which an increasing number of swap transformations are enacted in a densely packed all-to-all random circuit and its inverse. Here, every layer contains $N/2$ two-qubit gates where neighbors are determined by taking a random permutation and pairing sequential qubits. Indeed, before any swap transformations are introduced, the transpiler can find the full simplification, eliminating all two-qubit gates.

While one may expect that the swaps have a gradual effect on the size of remaining gates which were not transpiled away, we find that actually the difficulty jumps relatively quickly. One can think of this swap interaction as being a bottleneck in the simplification of the full circuit. Thus, even a few bottlenecks may lead to a large fraction of gates surviving in the lightcone of those swaps.
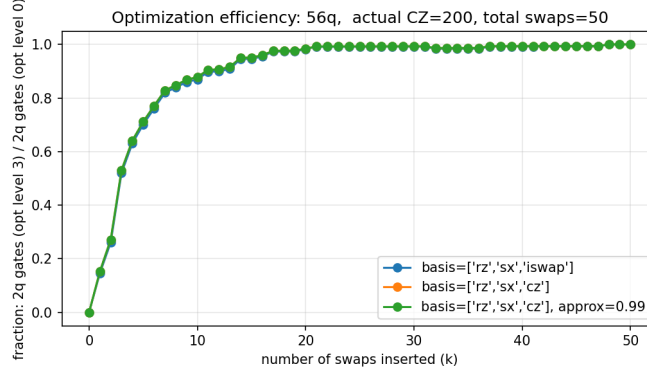
18

Figure 11: Transpiling away the identity. We show the difficulty of transpilers in simplifying these circuits. We start with an identity circuit comprised of $U$ and $U^\dagger$. The qiskit transpiler easily identifies it and reduces it. We then implement a series of swap transformations gradually and show the fraction of two qubits that survive the transpilation.

**Middle MPO Attack** To test the resiliency of our HQAP circuits, we employ a MPO compression technique we call the *middle MPO attack*. To build intuition, consider how time evolution is applied to MPS. A gate gets applied as a tensor operation before performing compression back to the MPS form, possibly with higher bond dimension. Our MPO attack is based on a similar principle, except it can be in the middle of the circuit (with gates before and after), as opposed to an MPS with all gates in the future.

This attack is particularly well suited for circuits which build up the identity from a sequence of random gates followed by the adjoint sequence. By starting with a rank 1 MPO object initialized to the identity in the middle of the circuit, joint evolution of gates on either side results in an MPO which is again representable with low rank since it remains the identity. Repeating this reveals that the whole circuit indeed has a much simpler representation and does not require simulation with exponentially large dimension.

The combination of sweeping, masking, and applying permutations has thus far prevented this attack from successfully compressing the circuit. This shows that HQAP circuits may remain hard even when the boundary between where $T(U)$ ends and $U^\dagger$ begins is known. Though it is known that worst-case identity checking is QMA complete [14], an open direction of importance is understanding how our version of identity obfuscation compares. In the meantime, we invite practictioners to try this and other strategies in our online portal. See Section 7.

# 5 HQAP Circuits as a potential candidate for post-quantum encryption

Throughout the work, we have studied the hardness of finding the peak when the input state is known and set to the all-0 bitstring. In this section, we discuss the possibility of using HQAP circuits as a post-quantum cryptographic protocol [1]. We first prove that, for generic quantum circuits, deciding whether a given input circuit is peaked is QCMA-complete. Note this decision problem of distinguishing peaked circuits differs from the task we consider in the previous sections, which is to produce the peaked bitstring given a circuit that is peaked. Therefore, unless BQP = QCMA, which is unlikely to be true under standard complexity assumptions, even a quantum computer cannot distinguish a peaked circuit from a non-peaked one. As a corollary, finding the input computational basis on which a circuit is peaked is hard even for a quantum computer. This provides a necessary

condition that HQAP-based encryption can be secure: an efficient adversary without the trapdoor cannot even decide peakedness with non-negligible advantage, let alone recover the planted input string and obtaining the encrypted output string or weight.

On the other hand, the obfuscation protocol we build in this work permits efficient generation and would generate a non-negligible signal for a designated input-output string pair. Therefore, *if there exists an HQAP algorithm such that "finding the input string on which the circuit is peaked" is sufficiently hard on average, we can realize a post-quantum encryption scheme by treating the input state as a secret trapdoor: the encryption party publishes an obfuscated, peaked circuit for the target string, while the decryption party applies the trapdoor to recover it. In this case, correctness follows from peakedness, and security reduces to the hardness of learning the peak without the trapdoor.*

## 5.1 Deciding whether a quantum circuit is peaked is QCMA-hard

To this end, we cast the task of distinguishing a peaked circuit from a non-peaked one as the following decision problem:

**Definition 5.1** (Peakedness check on basis states (PCBS)). *Let $U$ be an $n$-qubit unitary and fix thresholds $1 \geq \delta_{\text{yes}} > \delta_{\text{no}} \geq 0$ with $\delta_{\text{yes}} - \delta_{\text{no}} \geq 1/\text{poly}(N)$. The promise problem PCBS asks to decide whether*

$$\textsf{YES: } \exists y \in \{0,1\}^N \text{ with } |\langle y|U|y\rangle| \geq \delta_{\text{yes}} \quad vs \quad \textsf{NO: } \forall y \in \{0,1\}^N, \ |\langle y|U|y\rangle| \leq \delta_{\text{no}}\,.$$

Next, we introduce a complexity class called QCMA. Intuitively, QCMA is the probabilistic version of NP (more formally, MA) with a quantum verifier and a classical witness.

**Definition 5.2** (QCMA). *Fix completeness–soundness thresholds $1 \geq c > s \geq 0$ with gap $c - s \geq 1/\text{poly}(N)$. A language $L \subseteq \{0,1\}^\star$ is in QCMA if there is a classical, polynomial-time uniform generator that, on input $x$, outputs a polynomial-size quantum verifier circuit $W_x$ acting on registers*

$$Y \text{ (witness, } N_x \text{ qubits)}, \quad A \text{ (ancillas, } M_x \text{ qubits initialized to } |0^{M_x}\rangle), \quad B \text{ (one-qubit output)},$$

*such that*

$$\textsf{YES: } \exists y \in \{0,1\}^{N_x} \quad \Pr\big[B{=}1 \text{ after } W_x \text{ on } |y\rangle_Y |0^{M_x}\rangle_A |0\rangle_B\big] \ \geq \ c,$$
$$\textsf{NO: } \forall y \in \{0,1\}^{N_x} \quad \Pr\big[B{=}1 \text{ after } W_x \text{ on } |y\rangle_Y |0^{M_x}\rangle_A |0\rangle_B\big] \ \leq \ s.$$

*Equivalently, writing $\Pi_B := |1\rangle\langle 1|_B$ and $\rho_y := |y\rangle\langle y|_Y \otimes |0^{M_x}\rangle\langle 0^{M_x}|_A \otimes |0\rangle\langle 0|_B$, the conditions are $\text{tr}(W_x \rho_y W_x^\dagger \Pi_B) \geq c$ vs. $\leq s$.*

**Theorem 5.3** (PCBS is QCMA-complete). *PCBS is QCMA-complete under a polynomial-time reduction.*

From the definition it is clear that PCBS has the right flavor for QCMA: in QCMA the verifier's acceptance probability on some classical witness separates YES and NO instances; in PCBS we ask whether some diagonal entry of $U$ is large. First, it's easy to see the containment of PCBS $\in$ QCMA: a QCMA verifier uses $y$ as the witness and estimates $\Pr[\text{meas. } y \text{ after } U|y\rangle] = |\langle y|U|y\rangle|^2$ by sampling; since $\delta_{\text{yes}}^2 - \delta_{\text{no}}^2 \geq 1/\text{poly}(N)^2$, this is distinguishable in poly time.

Next, we show QCMA-hardness via a gadget that maps acceptance probabilities of any QCMA witness circuit $U$ to solving PCBS problem on an embedded larger circuit $Z$. This proof is inspired by [51].

**Theorem 5.4.** *Let $U$ be a circuit on registers $(Y, A, B)$ and let $P := \Pi_{B=1}$. For a basis witness $y$, define*

$$p(y) := \langle y, 0^a, 0 | U^\dagger P U | y, 0^a, 0 \rangle \in [0, 1].$$

*Fix $0 < \phi \le \frac{\pi}{2}$ and the one-qubit rotation $R(\phi) = \begin{pmatrix} \cos\phi & -\sin\phi \\ \sin\phi & \cos\phi \end{pmatrix}$ on a marker qubit $C$. Define*

$R_1 :=$ *apply* $R(\phi)$ *to* $C$ *controlled on* $(A{=}0^a)$,     $R_2 :=$ *apply* $R(-\phi)$ *to* $C$ *controlled on* $(B{=}1)$,

*and introduce a fresh guard qubit $D$ (initialized to $|0\rangle$) with the unitary*

$$G := X_D \otimes \left(I_{YABC} - |0^a 0\rangle\langle 0^a 0|_{AB}\right) + I_D \otimes |0^a 0\rangle\langle 0^a 0|_{AB}.$$

*Next, we set*

$$Z := U^\dagger R_2 U R_1 G.$$

*For the canonical basis state $|z\rangle := |0\rangle_C |0\rangle_D \otimes |y, 0^a, 0\rangle_{YAB}$ we have*

$$\left|\langle z|Z|z\rangle\right|^2 = \left(p(y) + (1 - p(y))\cos\phi\right)^2. \tag{4}$$

*Moreover, for any basis state $|z'\rangle$ with $(A, B) \ne (0^a, 0)$, that is: when the ancillas are not in the canonical input staes*

$$\langle z' | Z | z' \rangle = 0 \qquad \text{(since $G$ flips $D$ to an orthogonal basis state).} \tag{5}$$

**YES bound.** *If there exists $y$ with $p(y) \ge 1 - \varepsilon$, then for the corresponding canonical $|z\rangle$,*

$$\left|\langle z|Z|z\rangle\right|^2 \ge \left((1-\varepsilon) + \varepsilon\cos\phi\right)^2 = \left(1 - \varepsilon(1 - \cos\phi)\right)^2. \tag{6}$$

**NO bound.** *If for all $y$ we have $p(y) \le \varepsilon$, then for every canonical $|z\rangle$,*

$$\left|\langle z|Z|z\rangle\right|^2 \le \left(\varepsilon + (1-\varepsilon)\cos\phi\right)^2 \le \left(\cos\phi + \varepsilon\right)^2. \tag{7}$$

*Proof. Step 0 (apply $G$).* Recall $Z = U^\dagger R_2 U R_1 G$, so $G$ acts first. If $(A, B) = (0^a, 0)$, $G$ is the identity and the analysis below applies unchanged. If $(A, B) \ne (0^a, 0)$, i.e., in the non-canonical input case, then $G$ flips $D$; since the rest of the circuit acts trivially on $D$, $\langle z'|Z|z'\rangle = 0$, proving (5).

In the canonical case, we expand $U |y, 0^a, 0\rangle = c_1 |1\rangle_B |\psi_1\rangle + c_0 |0\rangle_B |\psi_0\rangle$, with $|c_1|^2 = p(y)$ and $|c_0|^2 = 1 - p(y)$.

*Step 1 (apply $R_1$).* Since $A = 0^a$ on the canonical input, $R_1$ rotates the marker: $|0\rangle_C \mapsto |c_0\rangle_C := R(\phi)|0\rangle = \cos\phi\,|0\rangle + \sin\phi\,|1\rangle$.

*Step 2 (apply $U$).* We get

$$|c_0\rangle_C \otimes U|y, 0^a, 0\rangle = c_1 |c_0\rangle_C \otimes |1\rangle_B |\psi_1\rangle + c_0 |c_0\rangle_C \otimes |0\rangle_B |\psi_0\rangle.$$

*Step 3 (apply $R_2$ controlled on $B = 1$).* On the $B{=}1$ branch the marker sees $R(-\phi)$, so $R(-\phi)R(\phi)|0\rangle = |0\rangle$; on the $B{=}0$ branch nothing happens. Hence the joint state becomes

$$c_1 |0\rangle_C |1\rangle_B |\psi_1\rangle + c_0 |c_0\rangle_C |0\rangle_B |\psi_0\rangle.$$

*Step 4 (apply $U^\dagger$ and overlap).* Applying $U^\dagger$ on $(YAB)$ and projecting onto $\langle 0|_C \otimes \langle y, 0^a, 0|$ gives

$$\langle z|U^\dagger R_2 U R_1|z\rangle = |c_1|^2 + \cos\phi\,|c_0|^2 = p(y) + (1 - p(y))\cos\phi,$$

which proves (4). Notice that the quantity is real and nonnegative for $\phi \in (0, \frac{\pi}{2})$. The bounds (6)-(7) follow from the monotonicity of $f(p) := p + (1 - p)\cos\phi$. $\qquad\square$

Given a QCMA instance $x$ with verifier $W_x$ and thresholds $c > s$, build $Z$ from $U := W_x$ as in Theorem 5.4 using any fixed constant $\phi \in (0, \frac{\pi}{2})$ (e.g. $\phi = \pi/3$). Set PCBS thresholds

$$\delta_{\text{yes}} := f(c) = c + (1 - c)\cos\phi, \qquad \delta_{\text{no}} := f(s) = s + (1 - s)\cos\phi.$$

Then $\delta_{\text{yes}} - \delta_{\text{no}} = (1 - \cos\phi)(c - s) \geq 1/\text{poly}(N)$, and:

$$x \in L \Longleftrightarrow \exists y : \ |\langle z|Z|z\rangle| \geq \delta_{\text{yes}} \quad \text{vs} \quad x \notin L \Longleftrightarrow \forall y : \ |\langle z|Z|z\rangle| \leq \delta_{\text{no}},$$

establishing QCMA-hardness. Together with our previous containment argument, this proves Theorem 5.3. An immediate corollary of this theorem is:

**Corollary 5.5.** *There exists no quantum polynomial-time algorithm that, given a generic $\text{poly}(N)$-sized quantum circuit $U$ (described classically), outputs the peaked input string with non-negligible advantage than random guess, unless* $\text{BQP} = \text{QCMA}$.

This can be proved by contradiction: assume that a quantum polynomial-time (QPT) algorithm can return the secret peaked input, then running the quantum circuit would generate the desired peakedness, resolving PCBS.

## 5.2   Symmetric encryption with HQAP circuits

The QCMA-hardness in the previous section suggests that finding the peaked input–output pair of a peaked circuit is hard even for quantum algorithms. This motivates using HQAP-generated circuits as ciphertexts: in our obfuscation protocols, a quantum adversary is not expected to efficiently identify the designated peaked input or output. Intuitively, hiding a single input–output pair among exponentially many basis labels makes it difficult to locate without the trapdoor. We conjecture:

**Conjecture 5.6** (Peak-Search Hardness (PSH))**.** *There exist an efficient HQAP protocol which, given $(x^\star, y^\star) \in \{0,1\}^N \times \{0,1\}^N$ and fresh coins $\rho$, the procedure $\text{HQAP}(x^\star, y^\star; \rho)$ outputs a circuit $E$ such that measuring $E|x^\star\rangle$ in the computational basis returns $y^\star$ with probability at least $\delta \in (0,1)$. For all $x \neq x^\star$, the output distribution of $E|x\rangle$ is not comparably concentrated on $y^\star$.*
*Moreover, the security of this protocol is guaranteed in the adversarial setting: give the adversary the classical description $\text{desc}(E)$. Then no QPT outputs $x^\star$ (or the pair $(x^\star, y^\star)$) with probability better than $2^{-n} + negl(n)$.*

Then with a slight modification to that HQAP circuit protocol can also be used for encrypting messages with quantum-safe property. Here is a sketch of the protocol: Let $k \in \{0,1\}^d$ be a pre-shared key and $F_k$ a pseudo random function (PRF) to derive per-block designated inputs. To encrypt $M$, parse it into $n$-bit blocks $s_1, \ldots, s_t$. For each block $i$:

1. Publish a nonce $\text{ctr}_i$ and set $x_i^\star := F_k(\text{ctr}_i)$.

2. Set $y_i^\star := s_i$ (optionally, mask as $s_i \oplus G_k(\text{ctr}_i)$ with an independent PRF $G$).

3. Sample $\rho_i$ and compute $E_i \leftarrow \text{HQAP}(x_i^\star, y_i^\star; \rho_i)$.

4. Output ciphertext $C_i := \big(\text{desc}(E_i), \text{ctr}_i\big)$.

The receiver decrypts by computing $x_i^\star = F_k(\text{ctr}_i)$, running $E_i$ on $|x_i^\star\rangle$ for $N_s = O(\delta^{-2})$ shots, taking the empirical mode $\hat{y}_i$, and outputting $\hat{s}_i := \hat{y}_i$. As a check, the receiver accepts only if the measured peak weight $\hat{p}^{\text{peak}} = \frac{1}{N_s} \sum_{j=1}^{N} \mathbf{1}\{y_{i,j} = \hat{y}_i\}$ exceeds a threshold $\delta_{\min}$.

It's not hard to see that the protocol remains secure under any QPT adversary against one-wayness under chosen-plaintext attacks (OW-CPA): Assume PRF security for $F$ and PSH for the HQAP distribution. Given a ciphertext block $C_i = (\mathrm{desc}(E_i), \mathsf{ctr}_i)$, recovering $s_i$ requires either (a) computing $x_i^\star = F_k(\mathsf{ctr}_i)$ (which breaks the PRF), or (b) finding $x_i^\star$ by probing $E_i$ on inputs and detecting the heavy output (which breaks PSH), then reading off $y_i^\star$ and removing the mask via $G_k$ (which again requires the PRF). Thus, under PRF security and PSH, the scheme achieves OW-CPA.

On the receiver side, checking if the output peakedness have reached a designated threshold provides a practical sanity layer against malformed or dishonest ciphertexts. Beyond this, this HQAP protocol offers three unique properties other PQC schemes may lack: (i) validity here is semantic, an inherent property of how the ciphertext was generated and, crucially, receiver-local; (ii) the same machinery naturally doubles as a proof-of-quantumness-flavored check; and (iii) unlike standardized post-quantum cryptography (e.g., ML-KEM/ML-DSA/SLH-DSA) [1], which is deliberately classical and deployable on today's networks, our protocol is inherently quantum because decryption or validation requires executing the public circuit on the designated input. Of course, these properties hinge on an average-case PSH conjecture holding for our instance generation: this is precisely why establishing rigorously or empirically testing on the hardness properties of HQAP circuits remain a crucial open problem of broad applications.

## 6   Discussion

In this work, we have presented a protocol for constructing peaked quantum circuits that is both scalable and verifiable. Our construction builds on the original peaked circuit framework by introducing pre-processing steps and architectural refinements that allow these circuits to scale well beyond previous instantiations, in both qubit number and circuit depth. Importantly, we have designed specific circuits that stretch the computational limits of existing quantum processors while fitting within the error budget. At the same time, we have implemented a suite of classical simulation attack strategies, ranging from matrix product states to tensor network contraction methods to Pauli path simulators, that represent the current state of the art in classical circuit simulation. Many of these methods have been used in past work to challenge or reinterpret quantum advantage claims, including those involving IBM's utility-scale benchmarking efforts.

These results offer a timely contribution to the evolving landscape of quantum advantage. It is now broadly accepted that certain quantum processors can outperform classical computers on well-posed computational tasks, at least in the absence of strong classical shortcuts. However, in the case of RCS, these demonstrations are often based on output distributions that are difficult to interpret or verify at scale. Our work addresses this gap by providing a construction in which the quantum output is both easy to produce and easy to verify, yet remains classically intractable to simulate by all known methods. This puts us in what might be described as the *heuristic* phase of quantum advantage: a period in which empirical results, rather than unconditional complexity-theoretic proofs, serve as the main form of evidence.

In this sense, our peaked circuits occupy a role not unlike that of RSA in classical cryptography. Just as the hardness of factoring large integers is widely believed, despite the lack of a formal proof, because of decades of failed attempts to break it, the intractability of simulating peaked circuits may come to be seen as an empirical fact.

On the other hand, there may be continual progress on the theoretical side, bridging the gap between what is provably hard and our constructions. In this work, we build on the existing literature on the hardness of peaked circuits. In general, a peaked circuit may lead to concentration on a particular output bitstring from an input state other than the all zeros case. When the input

bitstring is unknown, determining whether a given circuit is peaked is QCMA-hard. Establishing a concrete connection between these complexity results and our existing construction would greatly add to the strength of claims of quantum advantage.

We do not claim that classical simulation of our peaked circuits is impossible. Indeed, further analysis may reveal weaknesses in specific instances, or new classical algorithms may improve upon current simulation strategies. But in the absence of such breakthroughs, we view these circuits as offering a valuable benchmark for the quantum community. They provide a target that is both meaningful and falsifiable, anyone who can simulate them classically is encouraged to do so (see Section 7). In this way, our work is as much a challenge as it is a result. By continuing to scale these constructions and openly testing their classical boundaries, we hope to sharpen the contours of what practical quantum advantage looks like in the NISQ era.

## 7 Open Peaked Circuits challenge

We have observed a big gap between quantum and all known classical tenchiques, however we do not exclude the possibility that a smarter classical technique exists. Hence we are announcing a peaked circuits challenge, where we will be releasing increasingly difficult peaked circuits to cement the runtime gap between quantum and classical solutions. These circuits can be found on BlueQubit Peak Portal where anyone is welcome to submit their classical solutions as well. This is similar to the RSA challenge announced in the 1990s and expanded in the 2000s [52], where the challenge was to factor a large product of primes. The hardness of the challenge was never formally proved, but it was heuristically agreed upon based on many tries from the community.

## Acknowledgements

## References

[1] Lily Chen et al. *Report on post-quantum cryptography*. Vol. 12. US Department of Commerce, National Institute of Standards and Technology . . ., 2016.

[2] Manoj Kumar and Pratap Pattnaik. "Post quantum cryptography (pqc)-an overview." In: *2020 IEEE High Performance Extreme Computing Conference (HPEC)*. IEEE. 2020, pp. 1–9.

[3] David Joseph et al. "Transitioning organizations to post-quantum cryptography." In: *Nature* 605.7909 (2022), pp. 237–243.

[4] Duc-Thuan Dam et al. "A survey of post-quantum cryptography: Start of a new race." In: *Cryptography* 7.3 (2023), p. 40.

[5] Scott Aaronson and Lijie Chen. "Complexity-Theoretic Foundations of Quantum Supremacy Experiments." In: (Dec. 2016). arXiv: 1612.05903 [quant-ph].

[6] Adam Bouland et al. "On the complexity and verification of quantum random circuit sampling." In: *Nature Phys.* 15.2 (2018), pp. 159–163. DOI: 10.1038/s41567-018-0318-2.

[7] Dominik Hangleiter and Jens Eisert. "Computational advantage of quantum random sampling." In: *Rev. Mod. Phys.* 95.3 (2023), p. 035001. DOI: `10.1103/RevModPhys.95.035001`.

[8] F. Arute, K. Arya, R. Babbush, et al. "Quantum supremacy using a programmable superconducting processor." In: *Nature* 574 (2019), pp. 505–510. DOI: `10.1038/s41586-019-1666-5`.

[9] Yulin Wu et al. "Strong quantum computational advantage using a superconducting quantum processor." In: *Physical review letters* 127.18 (2021), p. 180501.

[10] Qingling Zhu et al. "Quantum computational advantage via 60-qubit 24-cycle random circuit sampling." In: *Science bulletin* 67.3 (2022), pp. 240–245.

[11] A. Morvan, B. Villalonga, X. Mi, et al. "Phase transitions in random circuit sampling." In: *Nature* 634 (2024), pp. 328–333. DOI: `10.1038/s41586-024-07998-6`.

[12] M. DeCross et al. "Computational Power of Random Quantum Circuits in Arbitrary Geometries." In: *Phys. Rev. X* 15 (2 May 2025), p. 021052. DOI: `10.1103/PhysRevX.15.021052`.

[13] Scott Aaronson and Yuxuan Zhang. *On verifiable quantum advantage with peaked circuit sampling*. 2024. arXiv: `2404.14493 [quant-ph]`. URL: `https://arxiv.org/abs/2404.14493`.

[14] Dominik Janzing, Pawel Wocjan, and Thomas Beth. ""Non-Identity-Check" is QMA-Complete." In: *International Journal of Quantum Information* 03.03 (2005), pp. 463–473. DOI: `10.1142/S0219749905001067`.

[15] Zhengfeng Ji and Xiaodi Wu. *Non-Identity Check Remains QMA-Complete for Short Circuits*. 2009. arXiv: `0906.5416 [quant-ph]`. URL: `https://arxiv.org/abs/0906.5416`.

[16] Yusei Mori et al. "Quantum circuit unoptimization." In: *Phys. Rev. Res.* 7.2 (2025), p. 023139. DOI: `10.1103/PhysRevResearch.7.023139`.

[17] Román Orús. "A practical introduction to tensor networks: Matrix product states and projected entangled pair states." In: *Annals of Physics* 349 (2014), pp. 117–158. ISSN: 0003-4916. DOI: `https://doi.org/10.1016/j.aop.2014.06.013`.

[18] J. Ignacio Cirac et al. "Matrix product states and projected entangled pair states: Concepts, symmetries, theorems." In: *Rev. Mod. Phys.* 93 (4 Dec. 2021), p. 045003. DOI: `10.1103/RevModPhys.93.045003`.

[19] R. Alkabetz and I. Arad. "Tensor networks contraction and the belief propagation algorithm." In: *Phys. Rev. Res.* 3 (2 Apr. 2021), p. 023073. DOI: `10.1103/PhysRevResearch.3.023073`.

[20] Jacob Biamonte et al. "Quantum machine learning." In: *Nature* 549.7671 (Sept. 2017), pp. 195–202. ISSN: 1476-4687. DOI: `10.1038/nature23474`.

[21] Joseph Tindall and Matt Fishman. "Gauging tensor networks with belief propagation." In: *SciPost Phys.* 15 (2023), p. 222. DOI: `10.21468/SciPostPhys.15.6.222`.

[22] Patrick Rall et al. "Simulation of qubit quantum circuits via Pauli propagation." In: *Phys. Rev. A* 99 (2019), p. 062337. DOI: `10.1103/PhysRevA.99.062337`.

[23] Dorit Aharonov et al. "A Polynomial-Time Classical Algorithm for Noisy Random Circuit Sampling." In: *Proceedings of the 55th Annual ACM Symposium on Theory of Computing (STOC '23)*. See also arXiv:2211.03999 (2022). 2023. DOI: `10.1145/3564246.3585234`.

[24] Thomas Schuster et al. "A polynomial-time classical algorithm for noisy quantum circuits." In: (July 2024). arXiv: `2407.12768 [quant-ph]`.

[25] Tomislav Begušić, Johnnie Gray, and Garnet Kin-Lic Chan. "Fast and converged classical simulations of evidence for the utility of quantum computing before fault tolerance." In: *Science Advances* 10.3 (2024), eadk4321. DOI: `10.1126/sciadv.adk4321`.

[26] Tomislav Begušić, Kasra Hejazi, and Garnet Kin-Lic Chan. "Simulating quantum circuit expectation values by Clifford perturbation theory." In: *J. Chem. Phys.* 162.15 (2025), p. 154110. DOI: 10.1063/5.0269149.

[27] Tomislav Begušić and Garnet Kin-Lic Chan. "Real-Time Operator Evolution in Two and Three Dimensions via Sparse Pauli Dynamics." In: *PRX Quantum* 6 (2 Apr. 2025), p. 020302. DOI: 10.1103/PRXQuantum.6.020302.

[28] Hrant Gharibyan et al. "A Practical Guide to using Pauli Path Simulators for Utility-Scale Quantum Experiments." In: (July 2025). arXiv: 2507.10771 [quant-ph].

[29] Manuel S. Rudolph et al. "Pauli Propagation: A Computational Framework for Simulating Quantum Systems." In: (May 2025). arXiv: 2505.21606 [quant-ph].

[30] Yuxuan Zhang. "Complexity and hardness of random peaked circuits." In: (Sept. 2025). arXiv: 2510.00132 [quant-ph].

[31] Abhinav Deshpande et al. "Peaked quantum advantage using error correction." In: *arXiv preprint arXiv:2510.05262* (2025).

[32] S. A. Moses et al. "A Race-Track Trapped-Ion Quantum Processor." In: *Physical Review X* 13.4 (Dec. 2023). ISSN: 2160-3308. DOI: 10.1103/physrevx.13.041052.

[33] Ulrich Schollwock. "The density-matrix renormalization group in the age of matrix product states." In: *Annals of Physics* 326.1 (2011). January 2011 Special Issue, pp. 96–192. ISSN: 0003-4916. DOI: https://doi.org/10.1016/j.aop.2010.09.012.

[34] Jarrod R. McClean et al. "Barren plateaus in quantum neural network training landscapes." In: *Nature Communications* 9.1 (Nov. 2018). ISSN: 2041-1723. DOI: 10.1038/s41467-018-07090-4.

[35] Johnnie Gray. "quimb: A python package for quantum information and many-body calculations." In: *Journal of Open Source Software* 3.29 (2018), p. 819. DOI: 10.21105/joss.00819.

[36] Johnnie Gray and Stefanos Kourtis. "Hyper-optimized tensor network contraction." In: *Quantum* 5 (2021), p. 410. DOI: 10.22331/q-2021-03-15-410.

[37] Sergey Bravyi, David Gosset, and Yinchen Liu. "Classical simulation of peaked shallow quantum circuits." In: *Proceedings of the 56th Annual ACM Symposium on Theory of Computing*. 2024, pp. 561–572.

[38] Zhengfeng Ji and Xiaodi Wu. "Non-identity check remains QMA-complete for short circuits." In: *arXiv preprint arXiv:0906.5416* (2009).

[39] Quantum AI team and collaborators. *qsim*. Sept. 2020. DOI: 10.5281/zenodo.4023103.

[40] Harun Bayraktar et al. "cuQuantum SDK: A High-Performance Library for Accelerating Quantum Science." In: *arXiv preprint arXiv:2308.01999* (2023). DOI: 10.48550/arXiv.2308.01999.

[41] Yifan Zhang and Yuxuan Zhang. "Classical simulability of quantum circuits with shallow magic depth." In: *PRX Quantum* 6.1 (2025), p. 010337.

[42] E. Cuthill and J. McKee. "Reducing the bandwidth of sparse symmetric matrices." In: *Proceedings of the 1969 24th National Conference*. ACM '69. New York, NY, USA: Association for Computing Machinery, 1969, pp. 157–172. ISBN: 9781450374934. DOI: 10.1145/800195.805928.

[43] F. Verstraete and J. I. Cirac. *Renormalization algorithms for Quantum-Many Body Systems in two and higher dimensions*. 2004. arXiv: cond-mat/0407066 [cond-mat.str-el]. URL: https://arxiv.org/abs/cond-mat/0407066.

[44]    Guifré Vidal. "Efficient Classical Simulation of Slightly Entangled Quantum Computations."
        In: *Phys. Rev. Lett.* 91 (14 Oct. 2003), p. 147902. DOI: `10.1103/PhysRevLett.91.147902`.

[45]    Guifré Vidal. "Efficient Simulation of One-Dimensional Quantum Many-Body Systems." In:
        *Phys. Rev. Lett.* 93 (4 July 2004), p. 040502. DOI: `10.1103/PhysRevLett.93.040502`.

[46]    G. Vidal. "Classical Simulation of Infinite-Size Quantum Lattice Systems in One Spatial
        Dimension." In: *Phys. Rev. Lett.* 98 (7 Feb. 2007), p. 070201. DOI: `10.1103/PhysRevLett.98.`
        `070201`.

[47]    R. Orús and G. Vidal. "Infinite time-evolving block decimation algorithm beyond unitary
        evolution." In: *Phys. Rev. B* 78 (15 Oct. 2008), p. 155117. DOI: `10.1103/PhysRevB.78.155117`.

[48]    Ali Javadi-Abhari et al. *Quantum computing with Qiskit*. 2024. arXiv: `2405.08810` `[quant-ph]`.
        URL: `https://arxiv.org/abs/2405.08810`.

[49]    H. C. Jiang, Z. Y. Weng, and T. Xiang. "Accurate Determination of Tensor Network State
        of Quantum Lattice Models in Two Dimensions." In: *Phys. Rev. Lett.* 101 (9 Aug. 2008),
        p. 090603. DOI: `10.1103/PhysRevLett.101.090603`.

[50]    Saeed S. Jahromi and Román Orús. "Universal tensor-network algorithm for any infinite lattice."
        In: *Phys. Rev. B* 99 (19 May 2019), p. 195105. DOI: `10.1103/PhysRevB.99.195105`.

[51]    Pawel Wocjan, Dominik Janzing, and Thomas Beth. "Two QCMA-complete problems." In:
        *arXiv preprint quant-ph/0305090* (2003).

[52]    Burt Kaliski. "RSA factoring challenge." In: *Encyclopedia of Cryptography and Security*. Ed. by
        Henk C. A. van Tilborg. Boston, MA: Springer US, 2005, pp. 531–532. ISBN: 978-0-387-23483-0.
        DOI: `10.1007/0-387-23483-7_362`.