



CHORUS

This is the accepted manuscript made available via CHORUS. The article has been published as:

Quantum Principal Component Analysis Only Achieves an Exponential Speedup Because of Its State Preparation Assumptions

Ewin Tang

Phys. Rev. Lett. **127**, 060503 — Published 4 August 2021

DOI: [10.1103/PhysRevLett.127.060503](https://doi.org/10.1103/PhysRevLett.127.060503)

Quantum principal component analysis only achieves an exponential speedup because of its state preparation assumptions

Ewin Tang*
University of Washington
(Dated: July 9, 2021)

A central roadblock to analyzing quantum algorithms on quantum states is the lack of a comparable input model for classical algorithms. Inspired by recent work of the author [2], we introduce such a model, where we assume we can efficiently perform ℓ^2 -norm samples of input data, a natural analogue to quantum algorithms that assume efficient state preparation of classical data. Though this model produces less practical algorithms than the (stronger) standard model of classical computation, it captures versions of many of the features and nuances of quantum linear algebra algorithms. With this model, we describe classical analogues to Lloyd, Mohseni, and Rebentrost’s quantum algorithms for principal component analysis [3] and nearest-centroid clustering [4]. Since they are only polynomially slower, these algorithms suggest that the exponential speedups of their quantum counterparts is simply an artifact of state preparation assumptions.

INTRODUCTION

Quantum machine learning (QML) has shown great promise towards yielding new exponential quantum speedups in machine learning, ever since the pioneering linear systems algorithm of Harrow, Hassidim, and Lloyd [5]. Since ML routines often push real-world limits of computing power, an exponential improvement to algorithm speed would allow for machine learning systems with vastly greater capabilities. While we have found many fast QML subroutines for machine learning problems since HHL [6–10], researchers have not been able to prove that these subroutines can be used to achieve an exponentially faster algorithm for a classical machine learning problem, even in the strongest input and output models [1, 11]. A recent work of the author [2] suggests a surprising reason why: even our best QML algorithms, with issues with input and output models resolved, fail to achieve exponential speedups. This previous work constructs a classical algorithm matching, up to polynomial slowdown, a corresponding quantum algorithm for recommendation systems [12], which was previously believed to be one of the best candidates for an exponential speedup in machine learning [13]. In light of this result, we need to question our intuitions and reconsider one of the guiding questions of the field: when is quantum linear algebra exponentially faster than classical linear algebra?

The main challenge in answering this question is not in finding fast classical algorithms, as one might expect. Rather, most QML algorithms are *incomparable* to classical algorithms, since they take quantum states as input and output quantum states: we don’t even know an analogous classical model of computation where we can search for similar classical algorithms [1]. The quantum recommendation system is unique in that it has a classical input (a data structure implementing QRAM) and classical output (a sample from a vector in the computational basis), allowing for rigorous comparisons with classical algorithms.

In our previous work we suggest an idea for developing classical analogues to QML algorithms beyond this exceptional case [2]:

When QML algorithms are compared to classical ML algorithms in the context of finding speedups, any state preparation assumptions in the QML model should be matched with ℓ^2 -norm sampling assumptions in the classical ML model.

In this work, we implement this idea by introducing a new input model, *SQ access*, which is a form of ℓ^2 -norm sampling assumption. We can get SQ access to data under typical state preparation assumptions, so fast classical algorithms in this model are strong barriers to their QML counterparts admitting exponential speedups. To support that the resulting model is the right notion to consider, we use it to dequantize two seminal and well-known QML algorithms, quantum principal component analysis [3] and quantum supervised clustering [4]. That is, we give classical algorithms that, with classical SQ access assumptions replacing quantum state preparation assumptions, match the bounds and runtime of the corresponding quantum algorithms up to polynomial slowdown. Surprisingly, we do so using only the classical toolkit originally applied to the recommendation systems problem, demonstrating the power of this model in analyzing QML algorithms.

From this work, we conclude that the exponential speedups of the quantum algorithms that we consider arise from strong input assumptions, rather than from the quantum-ness of the algorithms, since they vanish when classical algorithms are given analogous assumptions. In other words, in a wide swathe of settings, on *classical data*, these algorithms do not give exponential speedups. Dequantized algorithms can still be useful for quantum data (say, states generated from a quantum system), though a priori it’s not clear if they give a speedup in that case, since the analogous “classical algorithm on

quantum data” isn’t well-defined.

Our dequantized algorithms in the SQ access model provide the first formal evidence supporting the crucial concern about strong input/output assumptions in QML. Based on these results, we recommend exercising care when analyzing quantum linear algebra algorithms, since some algorithms with poly-logarithmic runtimes only admit polynomial speedups. BQP-complete QML problems, such as sparse matrix inversion [5] and quantum Boltzmann machine training [14], still cannot be dequantized in full unless BQP=BPP. However, many QML problems that are not BQP-complete have strong input model assumptions (like QRAM) and low-rank-type assumptions (which makes sense for machine learning, where high-dimensional data often exhibits low-dimensional trends). This regime is precisely when the classical approaches we outline here work, so such problems are highly susceptible to dequantization. We believe continuing to explore the capabilities and limitations of this model is a fruitful direction for QML research.

Notation. $[n] := \{1, \dots, n\}$. Consider a vector $x \in \mathbb{C}^n$ and matrix $A \in \mathbb{C}^{m \times n}$. $A_{i,*}$ and $A_{*,i}$ will refer to A ’s i th row and column, respectively. $\|x\|$, $\|A\|_F$, and $\|A\|$ will refer to ℓ^2 , Frobenius, and spectral norm, respectively. $|x\rangle := \frac{1}{\|x\|} \sum_{i=1}^n x_i |i\rangle$ and $|A\rangle := \frac{1}{\|A\|_F} \sum_{i=1}^m \|A_{i,*}\| |i\rangle |A_{i,*}\rangle$ (where, by the previous definition, $|A_{i,*}\rangle = \frac{1}{\|A_{i,*}\|} \sum_{j=1}^n A_{i,j} |j\rangle$). $A = \sum_{i=1}^{\min m,n} \sigma_i u_i v_i^\dagger$ is A ’s singular value decomposition, where $u_i \in \mathbb{C}^m$, $v_i \in \mathbb{C}^n$, $\sigma_i \in \mathbb{R}$, $\{u_i\}$ and $\{v_i\}$ are sets of orthonormal vectors, and $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min m,n} \geq 0$. $A_\sigma := \sum_{\sigma_i \geq \sigma} \sigma_i u_i v_i^\dagger$ and $A_k := \sum_{i=1}^k \sigma_i u_i v_i^\dagger$ denote low-rank approximations to A . We assume basic arithmetic operations take unit time, and $\tilde{O}(f) := O(f \log f)$.

THE DEQUANTIZATION MODEL

A typical QML algorithm works in the model where state preparation of input is efficient and a quantum state is output for measurement and post-processing. (Here, we assume an ideal/fault-tolerant quantum computer.) In particular, given a data point $x \in \mathbb{C}^n$ as input, we assume we can prepare copies of $|x\rangle$. For m input data points as a matrix $A \in \mathbb{C}^{m \times n}$, we additionally assume efficient preparation of $|A\rangle$, to preserve relative scale. We wish to compare QML and classical ML on classical data, so state preparation usually requires access to this data and its normalization factors. This informs the classical input model for our quantum-inspired algorithms, where we assume such access, and instead of preparing states, we can prepare measurements of these states.

Definition. We have $O(T)$ -time *sample and query access* to $x \in \mathbb{C}^n$ (notated $\text{SQ}(x)$) if, in $O(T)$ time, we can query an index $i \in [n]$ for its entry x_i ; produce an independent measurement of $|x\rangle$ in the computational basis; and query

for $\|x\|$. If we can only query for an estimate of the squared norm $\bar{x} \in (1 \pm \nu)\|x\|^2$, then we denote this by $\text{SQ}^\nu(x)$. For $A \in \mathbb{C}^{m \times n}$, sample and query access to A (notated $\text{SQ}(A)$) is $\text{SQ}(A_{1,*}, \dots, A_{n,*})$ along with $\text{SQ}(\bar{A})$ where \bar{A} is the vector of row norms, i.e. $\bar{A}_i := \|A_{i,*}\|$.

Sample and query (SQ) access will be our classical analogue to quantum state preparation. As we noted previously [2], we should be able to assume that classical analogues can efficiently measure input states: QML algorithms shouldn’t rely on fast state preparation as the “source” of an exponential speedup. The algorithm itself should create the speedup.

For typical instantiations of state preparation oracles on classical input, we can get efficient SQ access to input. For example, given input in QRAM [15], a strong proposed generalization of classical RAM that supports state preparation, we can get log-dimension-time SQ access to input [2, Proposition 3.2]. Similarly, sparse and close-to-uniform vectors can be prepared efficiently, and correspondingly admit efficient SQ access [16].

So, in usual QML settings, SQ assumptions are easier to satisfy than state preparation assumptions.

This leads to a model based on SQ access that we codify with the informal definition of “dequantization”. We say we *dequantize* a quantum protocol $\mathcal{S} : O(T)$ -time state preparation of $|\phi_1\rangle, \dots, |\phi_c\rangle \rightarrow |\psi\rangle$ if we describe a classical algorithm of the form $\mathcal{C}_\mathcal{S} : O(T)$ -time $\text{SQ}(\phi_1, \dots, \phi_c) \rightarrow \text{SQ}^\nu(\psi)$ with similar guarantees to \mathcal{S} up to polynomial slowdown. This is the sense in which we dequantized the quantum recommendation system in prior work [12]. In the rest of this article, we will dequantize two quantum algorithms, giving a detailed sketch of the algorithm and leaving proofs of correctness to the supplemental material [16]. These algorithms are applications of three protocols from our previous work [2] rephrased in our access model.

NEAREST-CENTROID CLASSIFICATION

Lloyd, Mohseni, and Rebentrost’s quantum algorithm for clustering estimates the distance of a data point to the centroid of a cluster of points [4]. The paper claims [17] that this quantum algorithm gives an exponential speedup over classical algorithms. We dequantize Lloyd et al’s quantum supervised clustering algorithm [4] with only quadratic slowdown. Though classical algorithms by Aaronson [1] and Wiebe et al. [18, Section 7] dequantize this algorithm for close-to-uniform input and sparse input, we are the first to give a general classical algorithm for this problem.

Problem 1 (Centroid distance). Suppose we are given access to $V \in \mathbb{C}^{n \times d}$ and $u \in \mathbb{C}^d$. Estimate $\|u - \frac{1}{n} \bar{I}V\|^2$ to ε additive error with probability $\geq 1 - \delta$.

Note that we are treating vectors as rows, with $\vec{1}$ the vector of ones. Let $\bar{u} := \frac{u}{\|u\|}$ and let \bar{V} be V , normalized so all rows have unit norm. Both classical and quantum algorithms argue about $M \in \mathbb{R}^{(n+1) \times d}$ and $w \in \mathbb{R}^{n+1}$ instead of u and V , where

$$M := \begin{bmatrix} \bar{u} \\ \frac{1}{\sqrt{n}} \bar{V} \end{bmatrix} \text{ and } w := \left[\|u\| \quad -\frac{1}{\sqrt{n}} \bar{V} \right].$$

Because $wM = u - \frac{1}{n} \vec{1}V$, we wish to estimate $\|wM\|^2 = wMM^\dagger w^\dagger$. Let $Z := \|w\|^2 = \|u\|^2 + \frac{1}{n} \|V\|_F^2$ be an ‘‘average norm’’ parameter appearing in our algorithms.

Theorem 2 (Quantum Nearest-Centroid [4]). *Suppose that, in $O(T)$ time, we can (1) determine $\|u\|$ and $\|V\|_F$; or (2) prepare a state $|u\rangle, |V_1\rangle, \dots, |V_n\rangle$, or $|\bar{V}\rangle$. Then we can solve Problem 1 in $O(T \frac{Z}{\varepsilon} \log \frac{1}{\delta})$ time.*

The quantum algorithm proceeds by constructing the states $|M\rangle$ and $|w\rangle$, then performing a swap test to get $|wM\rangle$. The swap test succeeds with probability $\frac{1}{2} wMM^\dagger w^\dagger$, so we can run amplitude amplification to get an estimate up to ε error with $O(\frac{1}{\varepsilon} \log \frac{1}{\delta})$ overhead.

Dequantizing this algorithm is simply a matter of dequantizing the swap test, which is done in Algorithm 1. Here, $Q(y)$ is *query access* to y , which supports querying y ’s entries in $O(1)$ time, but no sampling or norm queries.

Algorithm 1 Inner product estimation

Input: $O(T)$ -time $\text{SQ}^\nu(x) \in \mathbb{C}^n$, $Q(y) \in \mathbb{C}^n$

Output: an estimate of $\langle x|y\rangle$

- 1: Let $s = 54 \frac{1}{\varepsilon^2} \log \frac{2}{\delta}$
 - 2: Collect measurements i_1, \dots, i_s from $|x\rangle$
 - 3: Let $z_j = x_{i_j}^\dagger y_{i_j} \frac{\|x\|^2}{|x_{i_j}|^2}$ for all $j \in [s]$ $\triangleright \mathbb{E}[z_j] = \langle x|y\rangle$
 - 4: Separate the z_j ’s into $6 \log \frac{2}{\delta}$ buckets of size $\frac{\delta}{\varepsilon^2}$, and take the mean of each bucket
 - 5: Output the (component-wise) median of the means
-

From a simple analysis of the random variable z_i ’s, we get the following result.

Proposition 3 ([2, Proposition 4.2]). For $x, y \in \mathbb{C}^n$, given $\text{SQ}^\nu(x)$ and $Q(y)$, Algorithm 1 outputs an estimate of $\langle x|y\rangle$ to $(\varepsilon + \nu + \varepsilon\nu) \|x\| \|y\|$ error with probability $\geq 1 - \delta$ in time $O(\frac{T}{\varepsilon^2} \log \frac{1}{\delta})$.

For this protocol, quantum algorithms can achieve a quadratic speedup via amplitude estimation (but no more, by unstructured search lower bounds [19]). To apply this to nearest-centroid, we write $wMM^\dagger w^\dagger$ as an inner product of tensors $\langle a|b\rangle$, where

$$a := \sum_{i=1}^d \sum_{j=1}^{n+1} \sum_{k=1}^{n+1} M_{ji} \|M_{k,*}\| |i\rangle |j\rangle |k\rangle = M \otimes \tilde{M};$$

$$b := \sum_{i=1}^d \sum_{j=1}^{n+1} \sum_{k=1}^{n+1} \frac{w_j w_k M_{ki}}{\|M_{k,*}\|} |i\rangle |j\rangle |k\rangle.$$

Then, we show we have SQ access to one of the tensors (a). With this, we see that the quadratic speedup from amplitude amplification is the only speedup that quantum nearest-centroid achieves:

Theorem 4 (Classical Nearest-Centroid). *Suppose we are given $O(T)$ -time $\text{SQ}(V) \in \mathbb{C}^{n \times d}$ and $\text{SQ}(u) \in \mathbb{C}^d$. Then one can output a solution to Problem 1 in $O(T \frac{Z^2}{\varepsilon^2} \log \frac{1}{\delta})$ time.*

PRINCIPAL COMPONENT ANALYSIS

We now dequantize Lloyd, Mohseni, and Rebentrost’s quantum principal component analysis (QPCA) algorithm [3], an influential early example of QML [20, 21]. While the paper describes a more general strategy for Hamiltonian simulation of density matrices, their central claim is an exponential speedup in an immediate application: producing quantum states corresponding to the top principal components of a low-rank dataset [3].

The setup for the problem is as follows: suppose we are given a matrix $A \in \mathbb{R}^{n \times d}$ whose rows correspond to data in a dataset. We will find the principal eigenvectors and eigenvalues of $A^\dagger A$; when A is a mean zero dataset, this corresponds to the top principal components.

Problem 5 (Principal component analysis). Suppose we are given access to $A \in \mathbb{C}^{n \times d}$ with singular values σ_i and right singular vectors v_i . Further suppose we are given σ , k , and η with the guarantee that, for all $i \in [k]$, $\sigma_i \geq \sigma$ and $\sigma_i^2 - \sigma_{i+1}^2 \geq \eta \|A\|_F^2$. With probability $\geq 1 - \delta$, output estimates $\hat{\sigma}_1^2, \dots, \hat{\sigma}_k^2$ and $\hat{v}_1, \dots, \hat{v}_k$ satisfying $|\hat{\sigma}_i^2 - \sigma_i^2| \leq \varepsilon_\sigma \|A\|_F^2$ and $\|\hat{v}_i - v_i\| \leq \varepsilon_v$ for all $i \in [k]$.

Denote $\|A\|_F^2 / \sigma^2$ by K . Lloyd et al. get the following:

Theorem 6. *Given $\|A\|_F$ and the ability to prepare copies of $|A\rangle$ in $O(T)$ time, a quantum algorithm can output the desired estimates for Problem 5 $\hat{\sigma}_1^2, \dots, \hat{\sigma}_k^2$ and $|\hat{v}_1\rangle, \dots, |\hat{v}_k\rangle$ in $\tilde{O}(TK \min(\varepsilon_\sigma, \delta)^{-3})$ time.*

Later results [12, 22, Theorems 5.2, 27] improve the runtime here to $\tilde{O}(TK \varepsilon_\sigma^{-1} \text{polylog}(nd/\delta))$ when A is given in QRAM. We will compare to the original QPCA result.

To dequantize QPCA, we use a similar high-level idea to that of the quantum-inspired recommendation system [2]. We begin by using a low-rank approximation algorithm, Algorithm 2, to output a description of approximate top singular values and vectors.

Algorithm 2 finds the large singular vectors of A by reducing its dimension down to W , whose SVD we can compute quickly. Then, $S, \tilde{U}, \tilde{\Sigma}$ define approximate large singular vectors $\hat{V} := S^\dagger \tilde{U} \tilde{\Sigma}^{-1}$. The full set of guarantees on the output of Algorithm 2 are in the supplemental material [16], but in brief, for the right setting of parameters, the columns of \hat{V} and the diagonal entries

of $\hat{\Sigma}$ satisfy the desired constraints for our \hat{v}_i 's and $\hat{\sigma}_i$'s in Problem 5. The $\hat{\sigma}_i$'s are output explicitly, but the \hat{v}_i 's are described implicitly: $\hat{v}_i = S^\dagger \hat{U}_{*,i} / \hat{\sigma}_i$. We have $O(T)$ -time SQ(S) because all rows are normalized, and rows of S are simply rows of A . Thus, sampling from \hat{S} is a uniform sample from $[q]$ and sampling from $S_{i,*}$ is sampling from a row of A . $\hat{U}_{*,i}$ is an explicit vector, so in essence, we need SQ access to a linear combination of vectors, each of which we have SQ access to.

Algorithm 2 Low-rank approximation [23]

Input: $O(T)$ -time SQ(A) $\in \mathbb{R}^{m \times n}$, $\sigma, \varepsilon, \delta$

Output: SQ(S) $\in \mathbb{C}^{\ell \times n}$, $Q(\hat{U}) \in \mathbb{C}^{q \times \ell}$, $Q(\hat{\Sigma}) \in \mathbb{C}^{\ell \times \ell}$

- 1: Set $K = \|A\|_F^2 / \sigma^2$ and $q = \Theta(\frac{K^4}{\varepsilon^2} \log(\frac{1}{\delta}))$
 - 2: Sample rows i_1, \dots, i_q from \hat{A} and define $S \in \mathbb{R}^{q \times n}$ such that $S_{r,*} := A_{i_r,*} \frac{\|A\|_F}{\sqrt{q} \|A_{i_r,*}\|}$
 - 3: Sample columns j_1, \dots, j_q from \mathcal{F} , where \mathcal{F} denotes the distribution given by sampling a uniform $r \sim [q]$, then sampling c from S_r .
 - 4: Let $W \in \mathbb{C}^{q \times q}$ be the normalized submatrix $W_{*,c} := \frac{S_{*,j_c}}{q \mathcal{F}(j_c)}$
 - 5: Compute the left singular vectors of W $\hat{u}^{(1)}, \dots, \hat{u}^{(\ell)}$ that correspond to singular values $\hat{\sigma}^{(1)}, \dots, \hat{\sigma}^{(\ell)}$ larger than σ
 - 6: Output SQ(S), $\hat{U} \in \mathbb{R}^{q \times \ell}$ the matrix with columns $\hat{u}^{(i)}$, and $\hat{\Sigma} \in \mathbb{R}^{\ell \times \ell}$ the diagonal matrix with entries $\hat{\sigma}^{(i)}$.
-

Algorithm 3 does exactly this: it uses rejection sampling to dequantize the swap test over a subset of qubits (getting $|Vw\rangle$ via $\langle V | (|w\rangle \otimes I)$).

Algorithm 3 Matrix-vector SQ access

Input: $O(T)$ -time SQ(V^\dagger) $\in \mathbb{C}^{k \times n}$, $Q(w) \in \mathbb{C}^k$

Output: $\text{SQ}^\nu(Vw)$

- 1: **function** REJECTIONSAMPLE(SQ(V^\dagger), $Q(w)$)
 - 2: Sample $i \in [k]$ proportional to $|w_i|^2 \|V_{*,i}\|^2$ by manually calculating all k probabilities
 - 3: Sample $s \in [n]$ from $V_{*,i}$ using SQ(V^\dagger)
 - 4: Compute $r_s = (Vw)_s^2 / (k \sum_{j=1}^k (V_{s,j} w_j)^2)$ (after querying for w_j and $V_{s,j}$ for all $j \in [k]$)
 - 5: Output s with probability r_s (success); otherwise, output \emptyset (failure)
 - 6: **end function**
 - 7: QUERY: output $(Vw)_s$
 - 8: SAMPLE: run REJECTIONSAMPLE until success (outputting s) or $kC(V,w) \log \frac{1}{\delta}$ failures (outputting \emptyset)
 - 9: NORM(ν): Let p be the fraction of successes from running REJECTIONSAMPLE $\frac{k}{\nu^2} C(V,w) \log \frac{1}{\delta}$ times; output $pk \sum_{i=1}^k |w_i|^2 \|V_{*,i}\|^2$
-

Proposition 7 ([2, Proposition 4.3]). For $V \in \mathbb{C}^{n \times k}$, $w \in \mathbb{C}^k$, given SQ(V^\dagger) and $Q(w)$, Algorithm 3 simulates $\text{SQ}^\nu(Vw)$ where the time to query is $O(Tk)$, sample is $O(Tk^2 C(V,w) \log \frac{1}{\delta})$, and query norm is $O(Tk^2 C(V,w) \frac{1}{\nu^2} \log \frac{1}{\delta})$. Here, δ is the desired failure probability and $C(V,w) = \sum \|w_i V_{*,i}\|^2 / \|Vw\|^2$.

In general, $C(V,w)$ may be arbitrarily large, but in this application it is $O(K)$. Quantum algorithms achieve a speedup here when k is large and $C(V,w)$ is small, such as when V is a high-dimensional unitary, confirming our intuition that unitary operations are hard to simulate classically.

Altogether, we get our desired result.

Theorem 8. Given $O(T)$ -time SQ(A) $\in \mathbb{C}^{n \times d}$, with $\varepsilon_\sigma, \varepsilon_v, \delta \in (0, 0.01)$, there is an algorithm that output the desired estimates for Problem 5 $\hat{\sigma}_1, \dots, \hat{\sigma}_k$ and $O(T \frac{K^9}{\varepsilon^4} \log^3(\frac{k}{\delta}))$ -time SQ $^{0.01}(\hat{v}_1, \dots, \hat{v}_k)$ in $O(\frac{K^{12}}{\varepsilon^6} \log^3(\frac{k}{\delta}) + T \frac{K^8}{\varepsilon^4} \log^2(\frac{k}{\delta}))$ time, where $\varepsilon = \min(0.1 \varepsilon_\sigma K^{1.5}, \varepsilon_v^2 \eta, \frac{1}{4} K^{-1/2})$.

Under the non-degeneracy condition $\eta \leq \frac{1}{4} K^{-1/2}$, this runtime is $\tilde{O}(T \frac{K^{12}}{\varepsilon^6 \varepsilon_v^{12}} \log^3(\frac{1}{\delta}))$. While the classical runtime depends on ε_v , note that a quantum algorithm must also incur this error term to learn about v_i from copies of $|v_i\rangle$. For example, computing entries or expectations of observables of v_i given copies of $|v_i\rangle$ requires $\text{poly}(\frac{1}{\varepsilon_v})$ or $\text{poly}(n)$ time.

DISCUSSION

We have introduced the SQ access assumption as a classical analogue to the QML state preparation assumption and demonstrated two examples where, in this classical model, we can dequantize QML algorithms with ease. We now discuss the implications of this work with respect to related literature.

A natural question is of this work's relation to classical literature: does this work improve on classical algorithms for linear algebra in any regime? The answer may be no, for a subtle but fundamental reason: recall that our main idea is to introduce an input model *strong enough* to give classical versions of QML while being *weak enough* to extend to settings like QRAM, where classical computers can only access the input in very limited ways. In particular, the SQ access model that we study is *weaker* than the typical input model used for classical sketching algorithms [24–26]. $O(T)$ -time algorithms in the quantum-inspired access model are $\tilde{O}(\text{nnz} + T)$ -time algorithms in the usual RAM model (where nnz is the number of nonzero entries of the input), but not vice versa: typical sketching algorithms can exploit better data structures provided they only take $O(\text{nnz})$ time (e.g. oblivious sketches), whereas the quantum-inspired model can only use the QRAM data structure. The crucial insight of this work is that some algorithms (such as Algorithm 2 of Frieze et al. [23]) generalize to the weaker quantum-inspired model. Our algorithms give exponential speedups in the quantum-inspired setting, but since the model is weaker, one might expect that they perform worse in typical settings for classical computation (see

[27]). These model considerations also explain why we use Frieze et al. [23]: to our knowledge, this algorithm is the only one from the classical literature that naturally generalizes to the SQ input model.

The closest analogue to these results and techniques is a work by Van den Nest on probabilistic quantum simulation [28], which describes a notion of “computationally tractable” (CT) states that corresponds to our notion of SQ access for vectors. With this notion, the author describes special types of circuits on CT states where weak simulation is possible, using variants of Propositions 3 and 7. However, Van den Nest’s work does not have a version of Algorithm 2, since this technique only runs quickly on low-rank matrices, making it ineffective on generic quantum circuits. We exploit this low-rank structure for efficient quantum simulation of a small-but-practically-relevant class of circuits: quantum linear algebra on data with low-rank structure. So, our techniques used for supervised clustering are within the scope of Van den Nest’s work, whereas our techniques for PCA are new to this line of work.

These techniques are not new to quantum simulation in general. Others have considered applying randomized numerical linear algebra to quantum simulation [29], but does not make the connection towards dequantizing quantum algorithms, especially in large generality. Low-rank approximation is crucial for tensor network simulations of quantum systems [30, 31], where simulation can be done efficiently provided the input is, say, a matrix product state with low tensor rank. In this context, low-rank ap-

proximation is often performed exactly and only on a subset of the space, instead of approximately done on the full state, as is done here. This reflects the fact that tensor network algorithms assume that the system is reasonably approximated by a tensor network and aims to work well in practice, whereas our “dequantized” algorithms must work on a broader class of input and prioritizes provable guarantees in an abstract computational model over real-world performance. Nevertheless, some of these dequantized algorithms might be able to be matched by tensor network contraction techniques, when the input has low *tensor* rank. See the supplemental material for further discussion of this comparison [16].

Since this work, numerous follow-ups have cemented the significance of the SQ access model introduced here [32–35]. In particular, a recent work [34] essentially dequantizes the singular value transformation framework of Gilyen et al. [36] when input is given in QRAM. These works use fundamentally the same techniques to dequantize a wide swathe of low-rank quantum machine learning—an exciting step forward in understanding QML.

Thanks to Ronald de Wolf for giving the initial idea to look at QPCA. Thanks to Nathan Wiebe for helpful comments on this document. Thanks to Daniel Liang and Patrick Rall for their help fleshing out these ideas and reviewing a draft of this document. Thanks to Scott Aaronson for helpful discussions. This material is based upon work supported by the National Science Foundation Graduate Research Fellowship Program under Grant No. DGE-1762114.

* ewint@cs.washington.edu; ewintang.com

- [1] S. Aaronson, *Nature Physics* **11**, 291 (2015).
- [2] E. Tang, in *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing - STOC 2019* (ACM Press, 2019) arXiv:1807.04271 [cs.IR].
- [3] S. Lloyd, M. Mohseni, and P. Rebentrost, *Nature Physics* **10**, 631 (2014).
- [4] S. Lloyd, M. Mohseni, and P. Rebentrost, arXiv (2013), arXiv:1307.0411 [quant-ph].
- [5] A. W. Harrow, A. Hassidim, and S. Lloyd, *Physical review letters* **103**, 150502 (2009).
- [6] P. Rebentrost, M. Mohseni, and S. Lloyd, *Physical review letters* **113**, 130503 (2014).
- [7] N. Wiebe, D. Braun, and S. Lloyd, *Physical review letters* **109**, 050505 (2012).
- [8] S. Lloyd, S. Garnerone, and P. Zanardi, *Nature Communications* **7**, 10138 (2016).
- [9] Z. Zhao, J. K. Fitzsimons, and J. F. Fitzsimons, *Physical Review A* **99**, 052331 (2019), arXiv:1512.03929 [quant-ph].
- [10] F. G. Brandao and K. M. Svore, in *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)* (IEEE, 2017).
- [11] A. M. Childs, *Nature Physics* **5**, 861 (2009).
- [12] I. Kerenidis and A. Prakash, in *8th Innovations in Theoretical Computer Science Conference (ITCS 2017)*, LIPIcs, Vol. 67 (Schloss Dagstuhl, 2017) pp. 49:1–49:21.
- [13] J. Preskill, *Quantum* **2**, 79 (2018).
- [14] M. Kieferová and N. Wiebe, *Phys. Rev. A* **96**, 062327 (2017).
- [15] V. Giovannetti, S. Lloyd, and L. Maccone, *Physical review letters* **100**, 160501 (2008).
- [16] See supplemental material for details on when sample and query access is possible; discussion on the relation of this work to the tensor networks literature; and full proofs for the results stated here. It includes Refs. [37–47].
- [17] We were not able to verify the quantum algorithm (namely, the Hamiltonian simulation for preparing $|\phi\rangle$) as stated. For our purposes, we can make the minor additional assumption of efficient state preparation access to $|\phi\rangle$, which makes correctness obvious. When we refer to the quantum algorithm in this letter, we mean this version of it.

- [18] N. Wiebe, A. Kapoor, and K. M. Svore, *Quantum Information and Computation* **15**, 316–356 (2015).
- [19] C. H. Bennett, E. Bernstein, G. Brassard, and U. Vazirani, *SIAM Journal on Computing* **26**, 1510 (1997).
- [20] J. Biamonte, P. Wittek, N. Pancotti, P. Rebentrost, N. Wiebe, and S. Lloyd, *Nature* **549**, 195 (2017).
- [21] C. Ciliberto, M. Herbster, A. D. Ialongo, M. Pontil, A. Rocchetto, S. Severini, and L. Wossnig, *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* **474**, 20170551 (2018).
- [22] S. Chakraborty, A. Gilyén, and S. Jeffery, in *46th International Colloquium on Automata, Languages, and Programming (ICALP 2019)*, LIPIcs (Schloss Dagstuhl, 2019) arXiv:1804.01973 [quant-ph].
- [23] A. Frieze, R. Kannan, and S. Vempala, *Journal of the ACM (JACM)* **51**, 1025 (2004).
- [24] M. W. Mahoney, *Foundations and Trends® in Machine Learning* **3**, 123 (2011).
- [25] D. P. Woodruff, *Foundations and Trends® in Theoretical Computer Science* **10**, 10.1561/0400000060 (2014).
- [26] R. Kannan and S. Vempala, *Acta Numerica* **26**, 95 (2017).
- [27] J. M. Arrazola, A. Delgado, B. R. Bardhan, and S. Lloyd, *Quantum* 10.22331/q-2020-08-13-307 (2020), arXiv:1905.10415 [quant-ph].
- [28] M. Van Den Nest, *Quantum Info. Comput.* **11** (2011).
- [29] A. Rudi, L. Wossnig, C. Ciliberto, A. Rocchetto, M. Pontil, and S. Severini, *Quantum* **4**, 234 (2020), arXiv:1804.02484 [quant-ph].
- [30] U. Schöllwöck, *Annals of Physics* **326**, 96 (2011).
- [31] R. Orús, *Annals of Physics* **349**, 117 (2014).
- [32] N.-H. Chia, A. Gilyén, H.-H. Lin, S. Lloyd, E. Tang, and C. Wang, in *31st International Symposium on Algorithms and Computation (ISAAC 2020)*, LIPIcs (Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2020).
- [33] N.-H. Chia, T. Li, H.-H. Lin, and C. Wang, in *45th International Symposium on Mathematical Foundations of Computer Science (MFCS 2020)*, LIPIcs (Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2020) arXiv:1901.03254 [cs.DS].
- [34] N.-H. Chia, A. Gilyén, T. Li, H.-H. Lin, E. Tang, and C. Wang, in *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing - STOC 2020* (ACM Press, 2020) arXiv:1910.06151 [cs.DS].
- [35] D. Jethwani, F. L. Gall, and S. K. Singh, in *45th International Symposium on Mathematical Foundations of Computer Science (MFCS 2020)*, LIPIcs (Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2020) arXiv:1910.05699 [cs.DS].
- [36] A. Gilyén, Y. Su, G. H. Low, and N. Wiebe, in *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing - STOC 2019* (ACM Press, 2019) arXiv:1806.01838 [quant-ph].
- [37] A. Prakash, *Quantum algorithms for linear algebra and machine learning.*, Ph.D. thesis, UC Berkeley (2014).
- [38] L. Grover and T. Rudolph, arXiv (2002), arXiv:0208112 [quant-ph].
- [39] F. Verstraete and J. I. Cirac, *Phys. Rev. B* **73**, 094423 (2006).
- [40] S. R. White, *Phys. Rev. Lett.* **69**, 2863 (1992).
- [41] G. Vidal, *Phys. Rev. Lett.* **91**, 147902 (2003).
- [42] G. Vidal, *Phys. Rev. Lett.* **93**, 040502 (2004).
- [43] A. Cichocki, N. Lee, I. Oseledets, A.-H. Phan, Q. Zhao, and D. P. Mandic, *Found. Trends Mach. Learn.* **9**, 249–429 (2016).
- [44] J. C. Bridgeman and C. T. Chubb, *Journal of Physics A: Mathematical and Theoretical* **50**, 223001 (2017).
- [45] J. Eisert, *Phys. Rev. Lett.* **97**, 260501 (2006).
- [46] Z. Landau, U. Vazirani, and T. Vidick, *Nature Physics* **11**, 566 (2015).
- [47] W. Hoeffding, *Probability inequalities for sums of bounded random variables*, in *The Collected Works of Wassily Hoeffding*, edited by N. I. Fisher and P. K. Sen (Springer New York, New York, NY, 1994) pp. 409–426.