

Quantum-inspired algorithms in practice

Juan Miguel Arrazola¹, Alain Delgado¹, Bhaskar Roy Bardhan¹, and Seth Lloyd^{2,1}

¹Xanadu, Toronto, Ontario, M5G 2C8, Canada

²Massachusetts Institute of Technology, Department of Mechanical Engineering, 77 Massachusetts Avenue, Cambridge, Massachusetts 02139, USA

We study the practical performance of quantum-inspired algorithms for recommendation systems and linear systems of equations. These algorithms were shown to have an exponential asymptotic speedup compared to previously known classical methods for problems involving low-rank matrices, but with complexity bounds that exhibit a hefty polynomial overhead compared to quantum algorithms. This raised the question of whether these methods were actually useful in practice. We conduct a theoretical analysis aimed at identifying their computational bottlenecks, then implement and benchmark the algorithms on a variety of problems, including applications to portfolio optimization and movie recommendations. On the one hand, our analysis reveals that the performance of these algorithms is better than the theoretical complexity bounds would suggest. On the other hand, their performance as seen in our implementation degrades noticeably as the rank and condition number of the input matrix are increased. Overall, our results indicate that **quantum-inspired algorithms can perform well in practice provided that stringent conditions are met: low rank, low condition number, and very large dimension of the input matrix. By contrast, practical datasets are often sparse and high-rank, precisely the type that can be handled by quantum algorithms.**

1 Introduction

A driving force for studying quantum computing is the conviction that quantum algorithms can solve some problems more efficiently than classical methods. But the boundary between classical and quantum computing changes constantly: new classical and quantum algorithms disrupt the landscape, continuously updating the border that separates these two paradigms. This interplay is fruitful; it allows classical and quantum algorithm developers to feed off each other's innovations and advance our understanding of the limits of computation. Indeed, several classical algorithms have been reported that rely at least partially on developments in quantum computing [1, 3, 5, 6, 12, 20, 26, 31, 32].

Quantum algorithms for linear algebra are a flagship application of quantum computing, particularly due to their relevance in machine learning [2, 14, 16, 17, 23]. These algorithms typically scale polylogarithmically with dimension, which, at the time they were reported, implied an asymptotic exponential speedup compared to state-of-the-art classical methods. For this reason, significant interest has been generated in the dequantization approach that led to breakthrough quantum-inspired classical algorithms for linear algebra problems with sublinear complexity [4, 10, 30, 31]. These dequantized algorithms work for general low-rank matrices, whereas quantum computers still exhibit an exponential speedup over all known classical algorithms for sparse, full-rank matrix problems, including the quantum Fourier transform, eigenvector and eigenvalue analysis, linear systems, and others. Dequantized algorithms for such problems – high-rank matrix inversion, for example – would imply that classical computers can efficiently simulate quantum computers, i.e.,

Juan Miguel Arrazola: juanmiguel@xanadu.ai

$\text{BQP} = \text{P}$, which is not currently considered to be likely.

The proof techniques used to determine the complexity of quantum-inspired algorithms lead to bounds that suggest that runtimes might be prohibitively large for practical applications: the proven complexity of the linear systems algorithm is $\tilde{O}(\kappa^{16}k^6\|A\|_F^6/\varepsilon^6)$ [10], while for recommendation systems it is $\tilde{O}(k^{12}/\varepsilon^{12})$ [31]¹. Here, A is the input matrix, $\|A\|_F$ is the Frobenius norm of A , κ is the condition number, k is the rank, and ε is the precision of the algorithm. These results raise a few immediate questions: are these complexity bounds the consequence of proof techniques, or do they reflect a fundamental limitation on the capabilities of these algorithms? How do these algorithms actually perform in practice?

We address these questions by analyzing, implementing, and testing quantum-inspired algorithms for linear algebra. First, we perform a theoretical analysis aimed at identifying potential bottlenecks in their practical implementation. This allows us to anticipate the regimes where quantum-inspired algorithms may outperform previously-known classical methods. We then test the algorithms on artificial examples based on random matrices, where it is possible to control the dimension, rank, and condition number of the input matrix. We conclude by testing the algorithms on problems of practical interest. Specifically, we run the quantum-inspired algorithm for linear systems of equations applied to portfolio optimization on stocks from the S&P 500, and analyze the algorithm for recommendation systems on the MovieLens dataset [13]. Based on our analysis and tests, we find that, provided that the input matrices have small rank and condition number, quantum-inspired algorithms can perform well in practice: they provide good approximations in reasonable times, even for very large-dimensional problems. This shows that quantum-inspired algorithms can have significantly smaller runtimes than their worst-case complexity bounds would suggest, although we find evidence that the dependence on the error ε may be tight for the linear systems algorithm. For problems whose input matrices have larger ranks and condition numbers, our implementation of the algorithms struggles noticeably to produce good results.

In the following sections, we begin by providing a brief description of quantum-inspired algorithms for linear algebra. This high-level summary may prove useful for those wishing to become more familiar with these techniques. We then analyze each of the main steps of the algorithm in further detail, focusing on identifying potential practical bottlenecks. We continue by implementing the algorithms and benchmarking their performance on artificially-generated problems and on real-world data. We conclude by summarizing the implications of our analysis.

2 Quantum-inspired algorithms for linear algebra

In this section, we give an overview of quantum-inspired algorithms for linear systems of equations and recommendation systems. More details can be found in Refs. [4, 10, 30, 31]. A unique feature of our description is that we view both algorithms as specific examples of a more general method to sample from vectors expressed in terms of the singular value decomposition (SVD) of an input matrix. Henceforth, capital letters are employed to denote matrices and bold letters to denote column vectors. For example, a linear system of equations is written as $A\mathbf{x} = \mathbf{b}$, where $A \in \mathbb{R}^{m \times n}$, $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$, and $\mathbf{b} = (b_1, b_2, \dots, b_m)^T$.

We consider problems of the following form: Given an $m \times n$ matrix $A \in \mathbb{R}^{m \times n}$ with SVD

$$A = \sum_{\ell=1}^k \sigma_\ell \mathbf{u}^{(\ell)} \mathbf{v}^{(\ell)\top}, \quad (1)$$

the goal is to sample entries of the n -dimensional vector

$$\mathbf{x} = \sum_{\ell=1}^k \lambda_\ell \mathbf{v}^{(\ell)}, \quad (2)$$

¹We have taken $\sigma = O(\|A\|_F/\sqrt{k})$ as stated in [31] and ignored scaling with respect to the error parameter η .

with respect to the *length-square* probability distribution $p_x(i) = x_i^2/\|\mathbf{x}\|^2$. The coefficients λ_ℓ depend on the matrix A and possibly on other inputs. For linear systems of equations, \mathbf{x} is the solution vector given by $\mathbf{x} = A^+ \mathbf{b}$, with A^+ the Moore-Penrose pseudoinverse. As discussed in Ref. [10], the coefficients λ_ℓ of Eq. (2) are given by the inner product

$$\lambda_\ell = \frac{1}{\sigma_\ell^2} \langle \mathbf{v}^{(\ell)}, A^T \mathbf{b} \rangle. \quad (3)$$

Similarly, for recommendation systems, A is the preference matrix, whose entries A_{ij} denote the rating that user i has given to product j . In this case, the solution vector \mathbf{x} is the i -th row of a low-rank approximation of A , which contains all the preferences of user i . The coefficients are given by

$$\lambda_\ell = \langle A_i^T, \mathbf{v}^{(\ell)} \rangle, \quad (4)$$

where A_i is the i -th row of A [31].

Quantum-inspired algorithms assume that the entries of A are given in a way that allows length-square sampling to be performed. This can be accomplished, for example, if the entries of A are stored in a data structure of the form proposed by [16], which stores both entries of vectors and their sub-norms. Length-square sampling allows us to preferentially sample the rows of A with large norm, and the large entries within each row. Given access to length-squared sampling, quantum-inspired algorithms solve low-rank linear algebra problems in three main steps:

1. *Approximate SVD*: The Frieze-Kannan-Vempala (FKV) algorithm [9] is used to compute approximate singular values $\tilde{\sigma}_\ell$ and approximate right singular vectors $\tilde{\mathbf{v}}^{(\ell)}$. The singular vectors are not calculated in their entirety; instead, a description is found allowing efficient computation of any given entry $v_i^{(\ell)}$.
2. *Coefficient estimation*: Based on the results of step 1, the coefficients $\tilde{\lambda}_\ell$ are approximated using Monte Carlo estimation techniques.
3. *Sampling solution vectors*: Using Eq. (2) and the results from steps 1 and 2, rejection sampling is employed to perform length-square sampling from the approximate vectors $\tilde{\mathbf{x}} = \sum_{\ell=1}^k \tilde{\lambda}_\ell \tilde{\mathbf{v}}^{(\ell)}$.

Before describing each of these steps in more detail, it is important to understand that the innovations of quantum-inspired algorithms are steps 2 and 3. In these steps, coefficient estimation and sampling from the solution vectors is performed in time $O(\text{poly}(k, \kappa, \varepsilon, \log n, \log m))$. By contrast, the strategy originally suggested in Ref. [9] is simply to compute coefficients and solution vectors directly from the approximate SVD of step 1, which requires time $O(kn)$. Asymptotically, quantum-inspired algorithms thus achieve an exponential asymptotic speedup in dimension, from $O(n, m)$ to $\text{polylog}(n, m)$, at the cost of a polynomial dependency on other parameters. In practice, the direct calculation can be done extremely fast since it scales only linearly with dimension, meaning that quantum-inspired algorithms require very large-dimensional problems before it becomes preferable to employ their sampling techniques.

2.1 Approximate SVD

The Frieze-Kannan-Vempala (FKV) algorithm [9] is a method to compute low-rank approximations of matrices. The FKV algorithm is a pioneering example of a randomized method for computing approximate matrix decompositions. These techniques have been studied extensively and several improvements have since been described [7, 8, 19, 24, 25, 27, 33, 34], see Ref. [11] for a survey of these techniques. As described in Ref. [11], these algorithms are variants of the same main strategy: (i) preprocess the matrix to calculate sampling probabilities, (ii) generate samples from the matrix, and (iii) use linear algebra techniques to postprocess the samples to compute a final approximation. It is an interesting question to understand to what extent other randomized methods for approximate matrix decompositions can be used to improve quantum-inspired algorithms. For concreteness and to allow a direct connection to Refs. [4, 10, 30, 31], in this work we focus on

$$\begin{aligned}
A &= \begin{pmatrix} -2.04 & 0.02 & 0.79 & 0.3 & 2.04 & 0.68 & -0.37 & 0.4 & -0.08 & 0.66 \\ i_1 & -3.44 & 1.5 & 2.68 & -0.93 & 4.04 & -0.02 & -0.78 & 0.96 & -0.58 & -1.99 \\ i_3 & 5.25 & -2.78 & -3.3 & -1.97 & -4.99 & -5.27 & 1.16 & 1.38 & -3.31 & 2.16 \\ i_4 & -0.03 & 3.14 & 2.49 & -0.33 & 1.04 & 0.81 & -0.21 & -0.84 & 1.8 & -3.86 \\ i_1 & -2.51 & 2.45 & 3.35 & -0.48 & 3.87 & 0.17 & -0.59 & 0.62 & 0.48 & -1.63 \\ i_2 & -0.89 & -0.79 & -1.5 & 1.41 & -0.97 & 3.14 & -0.19 & -1.58 & 1.37 & -0.68 \\ i_2 & 0.42 & 1.31 & 1.36 & -0.12 & 0.69 & -0.47 & 0.06 & 0.14 & 0.55 & -0.18 \\ i_2 & -0.22 & 2.63 & 1.94 & 0.78 & 0.85 & 2.08 & -0.18 & -1.3 & 2.58 & -2.31 \\ i_4 & 0.0 & 1.29 & 1.09 & 0.84 & 0.6 & 1.17 & -0.02 & -0.6 & 1.6 & 0.06 \\ i_1 & 3.69 & 1.48 & 1.27 & -1.99 & -1.2 & -4.43 & 0.69 & 1.24 & -1.13 & -0.39 \end{pmatrix} \\
R &= \begin{pmatrix} i_1 & -94.51 & 41.21 & 73.63 & -25.55 & 111. & -0.55 & -21.43 & 26.38 & -15.93 & -54.67 \\ i_2 & -7.33 & 87.64 & 64.64 & 25.99 & 28.32 & 69.31 & -6. & -43.32 & 85.97 & -76.97 \\ i_3 & 87.04 & -46.09 & -54.71 & -32.66 & -82.73 & -87.37 & 19.23 & 22.88 & -54.87 & 35.81 \\ i_4 & -35.23 & -31.27 & -59.38 & 55.82 & -38.4 & 124.31 & -7.52 & -62.55 & 54.24 & -26.92 \end{pmatrix} \\
R &= \begin{pmatrix} j_3 & & & j_1 & j_4 & & & & & j_2 \\ -94.51 & 41.21 & 73.63 & -25.55 & 111. & -0.55 & -21.43 & 26.38 & -15.93 & -54.67 \\ -7.33 & 87.64 & 64.64 & 25.99 & 28.32 & 69.31 & -6. & -43.32 & 85.97 & -76.97 \\ 87.04 & -46.09 & -54.71 & -32.66 & -82.73 & -87.37 & 19.23 & 22.88 & -54.87 & 35.81 \\ -35.23 & -31.27 & -59.38 & 55.82 & -38.4 & 124.31 & -7.52 & -62.55 & 54.24 & -26.92 \end{pmatrix} \\
C &= \begin{pmatrix} j_1 & j_2 & j_3 & j_4 \\ 106.02 & 57.6 & -129.5 & 106.02 \\ 93.08 & -94.6 & -10.04 & 93.08 \\ -78.78 & 49.96 & 119.26 & -78.78 \\ -85.51 & -136.6 & -48.28 & -85.51 \end{pmatrix}
\end{aligned}$$

Figure 1: Schematic representation of the FKV algorithm. The input is an $m \times n$ matrix A , which in in this example has ten rows and columns. After performing one pass through the matrix, it is possible to perform length-square sampling of rows and columns. A total of r rows are sampled from this length-square distribution, then renormalized and stacked together to construct a new $r \times n$ matrix R . In this example, we set $r = 4$ and obtain the row indices $(i_1, i_2, i_3, i_4) = (2, 8, 3, 6)$. Once the matrix R has been built, c column indices are drawn by repeatedly choosing rows of R uniformly at random, then sampling from the corresponding length-square distribution of column indices. In this example, we set $c = 4$ and obtain the column indices $(j_1, j_2, j_3, j_4) = (3, 8, 1, 3)$. Note that repeated indices are allowed in the algorithm. The final matrix C is built by renormalizing these columns of R and grouping them together. The singular values of C approximate those of A , while the singular vectors of C can be used to reconstruct approximations to the singular vectors of A .

analyzing the FKV algorithm and implement quantum-inspired algorithms that rely on FKV. The conclusions and performance of the implementation are therefore subject to the specific properties of this algorithm.

The main idea behind FKV is the following: instead of performing a singular value decomposition of the large input matrix $A \in \mathbb{R}^{m \times n}$, a smaller matrix $C \in \mathbb{R}^{r \times c}$ is constructed by sampling r rows and c columns of A . The only restriction on r and c is that they have to be larger than the rank of the low-rank approximation of A : in particular, they can be independent of m and n , i.e., not even $O(\log m)$, $O(\log n)$. The matrix C captures the important features of A , making it possible to perform an SVD of the smaller matrix C to retrieve information about the singular values and vectors of A . Since the rank of C is bounded by its dimension, this approach is only possible when A is low-rank or when a low-rank approximation of A is sufficient.

In more detail, the first step of FKV is to perform a pass through the matrix A and compute the Frobenius norm of each row $\|A_i\|_F$, as well as the norm of the matrix $\|A\|_F$. These quantities define a length-square distribution over rows: $p(i) = \|A_i\|_F^2 / \|A\|_F^2$. For each row i , a length-square distribution over its entries is also computed, $q_i(j) = \|A_{ij}\|_F^2 / \|A_i\|_F^2$. To sample from these distributions, it is possible to employ a specialized data structure as suggested in Refs. [9, 16, 31]

that allows sampling in time that is logarithmic in dimension. The algorithm proceeds as follows:

1. Sample r row indices i_1, i_2, \dots, i_r from the row distribution $p(i)$. For each row index, select the row A_{i_s} and renormalize it as $R_{i_s} = \frac{\|A\|_F}{\sqrt{r}\|A_{i_s}\|} A_{i_s}$. This defines an $r \times n$ matrix R consisting of the rescaled rows R_{i_s} .
2. Select an index $s \in \{1, 2, \dots, r\}$ uniformly at random, then sample a column index j from the column distribution $q_{i_s}(j)$. Repeat this a total of c times to obtain column indices j_1, j_2, \dots, j_c . For each column index, select the column R_{\cdot, j_t} and renormalize it as $C_{\cdot, j_t} = \frac{\|A\|_F}{\sqrt{c}\|R_{\cdot, j_t}\|} R_{\cdot, j_t}$. This defines a matrix C consisting of the rescaled columns C_{\cdot, j_t} .
3. Compute the singular value decomposition of C .

These steps are illustrated in Fig. 1. The FKV algorithm performs a useful tradeoff: it is possible to perform SVDs of a significantly smaller matrix C at the cost of obtaining only approximations of the singular values and vectors of the original matrix A . We denote by $\tilde{\sigma}_\ell$ the singular values of C and by $\omega^{(\ell)}$ its left singular vectors. The FKV algorithm directly provides approximations of singular values, $\tilde{\sigma}_\ell \approx \sigma_\ell$ for $\ell = 1, 2, \dots, k$. Approximate right singular vectors $\tilde{v}^{(\ell)}$ of A are obtained using the expression

$$\tilde{v}^{(\ell)} = \frac{1}{\tilde{\sigma}_\ell} R^T \omega^{(\ell)}, \quad (5)$$

while approximate left singular vectors $\tilde{u}^{(\ell)}$ are given by

$$\tilde{u}^{(\ell)} = A \left(\frac{1}{\tilde{\sigma}_\ell^2} R^T \omega^{(\ell)} \right). \quad (6)$$

In existing classical methods, these vectors are computed directly, whereas in quantum-inspired algorithms it suffices to query their entries.

The quality of the FKV approximations depends on the number of sampled rows r and columns c . The challenge is to find values of r and c such that sufficiently good approximations are obtained, while also ensuring that the complexity of computing the SVD of matrix C is significantly less than for the original input matrix A , i.e., such that $r \ll m$ and $c \ll n$. Since r, c can be much smaller than m, n , a substantial runtime reduction can potentially be achieved.

FKV is the core of quantum-inspired algorithms for linear algebra. Without the use of randomized methods, computing the SVD of an $m \times n$ matrix A requires $O(\min\{m^2n, mn^2\})$ time using naive matrix multiplication. From that point onwards, computing coefficients λ and the solution vector x takes only linear time in m, n . In FKV, we instead calculate the SVD of an $r \times c$ matrix, which requires $O(\min\{r^2c, rc^2\})$ time. When r, c are much smaller than m, n , this is where the real savings of quantum-inspired algorithms happen. As explained previously, the subsequent steps of coefficient estimation and sampling from the solution vector are used to reduce a linear runtime to polylogarithmic runtime. As we discuss later, for practical problem sizes that are far from the asymptotic limit, it is preferable to compute coefficients and solution vectors explicitly by employing Eqs. (5) and (6) starting from the approximate SVD of matrix C in the FKV algorithm.

2.2 Coefficient estimation

The coefficients λ appearing in Eq. (2) are inner products between vectors, multiplied by a power of the singular values in the case of linear systems. Quantum-inspired algorithms compute approximations of these coefficients using Monte Carlo estimation. In general, a coefficient λ is given by

$$\lambda = \langle \mathbf{y}, \mathbf{z} \rangle = \sum_{i=1}^n y_i z_i, \quad (7)$$

for some appropriate vectors \mathbf{y}, \mathbf{z} , as in Eqs. (3) and (4). The strategy to estimate this inner product is to perform length-square sampling from one of the vectors. Without loss of generality,

we take this vector to be \mathbf{y} . Define the random variable χ that takes values $\chi_i = y_i z_i / p_x(i)$, where the indices i are sampled from the length-square distribution $p_y(i) = y_i^2 / \|\mathbf{y}\|^2$. The expectation value of the random variable satisfies

$$\mathbb{E}(\chi) = \sum_{i=1}^n \frac{y_i z_i}{p_y(i)} p_y(i) = \langle \mathbf{y}, \mathbf{z} \rangle = \lambda. \quad (8)$$

Similarly, the second moment is

$$\mathbb{E}(\chi^2) = \sum_{i=1}^n \left(\frac{y_i z_i}{p_y(i)} \right)^2 p_y(i) = \|\mathbf{y}\|^2 \|\mathbf{z}\|^2, \quad (9)$$

and $\sigma_\chi^2 = \mathbb{E}(\chi^2) - \mathbb{E}(\chi)^2$ is the variance of χ . The strategy to estimate coefficients λ is to draw N samples $\chi^{(1)}, \chi^{(2)}, \dots, \chi^{(N)}$ from χ and compute the unbiased estimator $\hat{\lambda} = \frac{1}{N} \sum_{j=1}^N \chi^{(j)} \approx \lambda$. This constitutes a form of importance sampling in Monte Carlo estimation: large entries of \mathbf{y} are preferably selected since they contribute more significantly to the inner product. The error in the estimation is quantified by the ratio between the standard deviation and the mean of the estimator: $\epsilon := \sqrt{\text{Var}(\hat{\lambda})}/|\lambda|$. The variance of the estimator is $\text{Var}(\hat{\lambda}) = \sigma_\chi^2/N$, leading to a precision $\epsilon = \sigma_\chi/(|\lambda|\sqrt{N})$ in the estimator. This implies that the number of samples N needed to, with high probability, achieve a precision ϵ is

$$N = \frac{\sigma_\chi^2}{\lambda^2 \epsilon^2} = \frac{1}{\epsilon^2} \left[\frac{\|\mathbf{y}\|^2 \|\mathbf{z}\|^2}{\langle \mathbf{y}, \mathbf{z} \rangle^2} - 1 \right] \quad (10)$$

$$= O\left(\frac{\|\mathbf{y}\|^2 \|\mathbf{z}\|^2}{\epsilon^2 \langle \mathbf{y}, \mathbf{z} \rangle^2}\right) = O\left(\frac{1}{\epsilon^2 \cos^2 \theta}\right), \quad (11)$$

where θ is the angle between the vectors \mathbf{y} and \mathbf{z} .

2.3 Complexity of coefficient estimation

Quantum-inspired algorithms differ from the FKV algorithm because they estimate the coefficients λ rather than calculating them directly from the approximate SVD of the input matrix. Therefore, it is worthwhile to examine the complexity of coefficient estimation step to determine the regime where quantum-inspired algorithms become beneficial. For n -dimensional random vectors, the expectation value of the angle between two vectors satisfies $\mathbb{E}[\frac{1}{\cos^2 \theta}] = n$, so we can anticipate that the number of samples needed may scale linearly with dimension. If $\frac{1}{\cos^2 \theta} = \text{poly}(m, n)$, the algorithm no longer has the desired polylogarithmic complexity. In fact, as we show in Appendix A, regardless of what coefficient estimation strategy is employed, in a worst-case setting, inner products of vectors cannot be approximated in sublinear time. However, for matrices of low rank and condition number, as we show below, it is possible to achieve good estimation of coefficients.

For linear systems of equations, the coefficients λ_ℓ are given by

$$\lambda_\ell = \frac{1}{\sigma_\ell^2} \langle \mathbf{v}^{(\ell)}, A^T \mathbf{b} \rangle = \frac{1}{\sigma_\ell} \langle \mathbf{u}^{(\ell)}, \mathbf{b} \rangle. \quad (12)$$

In the quantum-inspired algorithm, we do not have direct access to the length-square distributions of either $A^T \mathbf{b}$ or $\mathbf{v}^{(\ell)}$. Instead, the strategy in Ref. [10] is to express the coefficients as $\lambda_\ell = \frac{1}{\sigma_\ell^2} \text{Tr}(A^T \mathbf{b} \mathbf{v}^{(\ell)T})$ and sample the entries A_{ij} of A . In this case the the number of samples needed in the estimation is [10]:

$$N = \frac{1}{\epsilon^2} \left(\frac{\|\mathbf{v}^{(\ell)}\|^2 \|A\|_F^2 \|\mathbf{b}\|^2}{\langle \mathbf{v}^{(\ell)}, A^T \mathbf{b} \rangle^2} - 1 \right) = O\left(\frac{1}{\epsilon^2} \frac{\|\boldsymbol{\sigma}\|^2 \|\mathbf{b}\|^2}{\beta_\ell^2}\right), \quad (13)$$

where we have defined $\boldsymbol{\sigma} := (\sigma_1, \sigma_2, \dots, \sigma_k)$, $\beta_\ell := \langle \mathbf{u}^{(\ell)}, \mathbf{b} \rangle$ such that $\mathbf{b} = \sum_{\ell=1}^k \beta_\ell \mathbf{u}^{(\ell)}$, and we have used Eq. (12) as well as the identity $\|A\|_F^2 = \|\boldsymbol{\sigma}\|^2$. In the special cases of recommendation

systems and portfolio optimization, the coefficients are proportional to an inner product $\langle A_i^T, \mathbf{v}^{(\ell)} \rangle$. Writing $A_i = \sum_{\ell'=1}^k \sigma_{\ell'} u_i^{(\ell')} \mathbf{v}^{(\ell')}^T$ we have

$$\langle A_i^T, \mathbf{v}^{(\ell)} \rangle = \sum_{\ell'=1}^k \sigma_{\ell'} u_i^{(\ell')} \langle \mathbf{v}^{(\ell')}, \mathbf{v}^{(\ell')} \rangle = \sigma_\ell u_i^{(\ell)}. \quad (14)$$

The number of samples then satisfies

$$N = \frac{1}{\epsilon^2} \left(\frac{\|A_i^T\|^2 \|\mathbf{v}^{(\ell)}\|^2}{\langle A_i^T, \mathbf{v}^{(\ell)} \rangle^2} - 1 \right) = O \left(\frac{\sum_{\ell'=1}^k (\sigma_{\ell'} u_i^{(\ell')})^2}{\epsilon^2 \sigma_\ell^2 u_i^{(\ell)^2}} \right). \quad (15)$$

Define the k -dimensional vectors $\boldsymbol{\sigma} := (\sigma_1, \sigma_2, \dots, \sigma_k)^T$, $\boldsymbol{\mu} := (u_i^{(1)}, u_i^{(2)}, \dots, u_i^{(k)})^T$ and their Hadamard product $\boldsymbol{\nu} = \boldsymbol{\sigma} \circ \boldsymbol{\mu}$ with entries $\nu_\ell = \sigma_\ell u_i^{(\ell)}$. It follows that

$$N = O \left(\frac{1}{\epsilon^2} \frac{\|\boldsymbol{\nu}\|^2}{\nu_\ell^2} \right). \quad (16)$$

In all cases, the number of samples depends on the ratio between the norm squared of a k -dimensional vector and the square of the ℓ -th entry of that vector. We now determine the worst-case scaling for these ratios. We consider the ratio $\|\boldsymbol{\sigma}\|^2/\sigma_\ell^2$ for concreteness; the same will apply for \mathbf{b} and $\boldsymbol{\nu}$. The largest ratio occurs for the smallest singular value σ_{\min} and we have

$$\frac{\|\boldsymbol{\sigma}\|^2}{\sigma_{\min}^2} = \frac{\|\boldsymbol{\sigma}\|^2}{\sigma_{\max}^2} \frac{\sigma_{\max}^2}{\sigma_{\min}^2} = \frac{\|\boldsymbol{\sigma}\|^2}{\sigma_{\max}^2} \kappa^2. \quad (17)$$

Similarly, this ratio is largest when σ_{\max} is as small as possible. This occurs when the largest $k-1$ singular values are equal to each other. In that case, defining $\sigma'_\ell := \sigma_\ell/\|\boldsymbol{\sigma}\|$ such that $\|\boldsymbol{\sigma}'\| = 1$ we have

$$\|\boldsymbol{\sigma}'\|^2 = (\sigma'_{\max})^2 \left(k - 1 + \frac{1}{\kappa} \right) = 1 \quad (18)$$

$$\Rightarrow (\sigma'_{\max})^2 = \frac{1}{k - 1 + \frac{1}{\kappa}} = O(1/k), \quad (19)$$

and therefore

$$\frac{\|\boldsymbol{\sigma}\|^2}{\sigma_{\max}^2} \kappa^2 = \frac{1}{(\sigma'_{\max})^2} \kappa^2 = O(k\kappa^2). \quad (20)$$

Defining $\kappa_\beta := \beta_{\max}/\beta_{\min}$ and $\kappa_\nu := \nu_{\max}/\nu_{\min}$, where $\beta_{\max/\min}$ and $\nu_{\max/\min}$ are the largest and smallest entries of their corresponding vectors, we conclude that the number of samples required to estimate coefficients in quantum-inspired algorithms for linear systems scales as

$$N = O \left(\frac{k^2 \kappa^2 \kappa_\beta^2}{\epsilon^2} \right), \quad (21)$$

and for recommendation systems as

$$N = O \left(\frac{k \kappa_\nu^2}{\epsilon^2} \right). \quad (22)$$

These asymptotic formulas – or more precisely Eqs. (13) and (15) – can be used to estimate the problem sizes for which the complexity of coefficient estimation is smaller than a direct calculation. As an illustration, in linear systems, setting $k = \kappa = \kappa_\beta = 100$ means that an order of $N = 10^{16}$ samples are needed to estimate the coefficients with an error of $\epsilon = 10^{-2}$. Finally, it is important to understand that a large number of samples can improve the precision of the estimation, but not its accuracy: if the singular values and vectors obtained from FKV have large errors, the expectation value of the estimator will not coincide with the actual value of the coefficient.

2.4 Sampling solution vectors

The approximate SVD and coefficient estimation steps respectively provide approximate singular values $\tilde{\sigma}_\ell$ and approximate coefficients $\tilde{\lambda}_\ell$. The approximate right singular vectors of A are not constructed explicitly, but are implicitly defined by Eq. (5), namely $\tilde{\mathbf{v}}^{(\ell)} = \frac{1}{\tilde{\sigma}_\ell} R^T \mathbf{w}^{(\ell)}$. This allows an implicit calculation of the approximate solution vector $\tilde{\mathbf{x}} = R^T \mathbf{w}$, where

$$\mathbf{w} := \sum_{\ell=1}^k \frac{\tilde{\lambda}_\ell}{\tilde{\sigma}_\ell} \mathbf{w}^{(\ell)}. \quad (23)$$

The challenge is to sample from these vectors using only query access to the entries of \mathbf{w} and R . Quantum-inspired algorithms achieve this using a well-known technique known as rejection sampling. The steps are as follows [31]:

1. Sample a row index i uniformly at random.
2. Sample a column index j from the length-square distribution $q_i(j) = |R_{ij}|^2 / \|R_i\|^2$.
3. Output j with probability $\frac{|\langle \mathbf{w}, R_{\cdot,j} \rangle|^2}{\|R_{\cdot,j}\|^2 \|\mathbf{w}\|^2}$, or sample j again.

The expected number of repetitions before outputting a column index j is $r \|\mathbf{w}\|^2 / (\sum_{j=1}^n |\langle \mathbf{w}, R_{\cdot,j} \rangle|^2) = r \|\mathbf{w}\|^2 / \|\tilde{\mathbf{x}}\|^2$. The steps of rejection sampling can be performed quickly. The overhead arises from calculating the inner product $\langle \mathbf{w}, R_{\cdot,j} \rangle$, which requires $O(rk)$ time; significantly less than computing the approximate SVD. As we illustrate in the following sections, the number of expected repetitions is usually moderate in practice, so sampling from the approximate solution vector can be performed in comparatively less time than other steps in the algorithm.

3 Numerical benchmarking

In the previous section we performed a theoretical analysis of quantum-inspired algorithms. Here, we complement that analysis by benchmarking the algorithms on specific problems. The main goal of this numerical implementation is to individually examine the performance of each step of the algorithm and use this information to reveal which steps contribute most to the overall runtime as well as which ones lead to the most significant source of errors. A secondary goal is to gain insights on the actual scaling of the algorithm with respect to its various parameters.

As an initial benchmark, we study the algorithm for linear systems of equations on artificial problems of extremely large dimension. We design matrices such that performing length-square sampling from its rows and columns, as well as all other steps of the quantum-inspired algorithm, can be done using only query access to its entries. This allows to study the performance of the algorithms in their intended asymptotic regime, without handling large datasets directly. We then test the algorithm for linear systems on Gaussian random matrices. This choice is made to allow full control over properties of the input matrices: dimension, rank, and condition number, while collecting large amounts of statistics. Afterwards, we apply the algorithms to portfolio optimization and movie recommendations. All algorithms were implemented in Python and the source code is available at <https://github.com/XanaduAI/quantum-inspired-algorithms>.

3.1 Design principles

In conducting numerical experiments, choices must be made regarding how the algorithms are implemented. Here we describe the principles that guide the numerical experiments performed to benchmark quantum-inspired algorithms. The first choice is the programming language. To implement resource-intensive algorithms, it is common to employ compiler-based languages such as C++. However, our main goal is not to develop high-performance implementation of quantum-inspired algorithms. Rather, the goal is to identify practical bottlenecks. Crucially, it is also important to share code that is easily-reproducible by the quantum information community. For these reasons, the algorithm is in Python. The runtimes reported in this work can therefore be

likely improved by using a language more suitable for high-performance computation, and by employing specialized hardware.

In implementing the algorithms, whenever possible, we have been faithful to the descriptions in Refs. [10] and [31]. These algorithms require the use of a specialized data structure to ensure polylogarithmic scaling of length-square sampling. As shown in Fig. 2, for the matrix sizes that we study in this work, which have dimensions smaller than 10^5 , using built-in sampling functions is orders-of-magnitude faster than employing our implementation of the data structure. Therefore, to speed up the practical runtime of the algorithm for these problems, we employ direct sampling from the distribution using built-in functions.

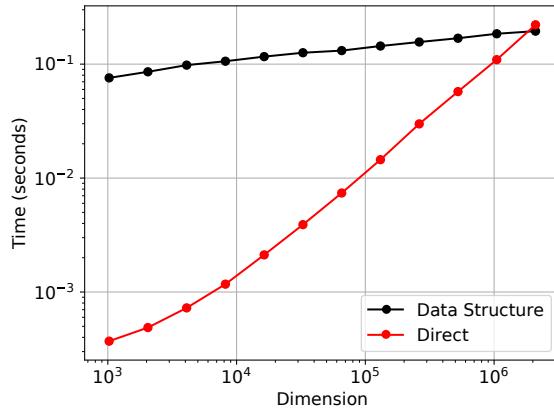


Figure 2: Runtimes for generating one thousand samples from a length-square distribution. We compare performance of the data structure of Ref. [16] and of direct sampling with built-in functions. Runtimes are plotted as a function of the dimension of the distribution. The polylogarithmic scaling is evident for the data structure, but highly-optimized direct sampling functions are orders-of-magnitude lower for dimensions smaller than 10^5 , as considered in this work.

Finally, although it is natural to consider the total variation distance as a way to quantify errors, we identify that the mean relative error, defined as $\eta = \sum_i \frac{|x_i - \tilde{x}_i|}{|x_i|}$ for two vectors x and \tilde{x} , is more informative in determining the quality of approximations. Moreover, it can be simultaneously used to measure errors in approximating singular values, coefficients, and solution vectors.

3.2 High-dimensional problems

Quantum-inspired algorithms for linear algebra are aimed at tackling extremely large datasets, but it is very challenging to test them in such scenarios. Instead, we design a class of problems where it is possible to run quantum-inspired algorithms using only black-box access to the entries of the input matrix and vector. We focus on the algorithm for linear systems and report the main formulas, whose derivations can be found in Appendix B.

Define the matrices

$$V^{(\ell)} = \frac{1}{2^n} \sum_{y,z \in \{0,1\}^n} (-1)^{\mathbf{x}^{(\ell)} \cdot (\mathbf{y} \oplus \mathbf{z})} e^{(y)} e^{(z)^T}, \quad (24)$$

where $\mathbf{x}, \mathbf{y}, \mathbf{z}$ are n -bit strings, the vectors $e^{(y)}$ form a basis of the 2^n -dimensional vector space, and the $\mathbf{x}^{(\ell)}$ are distinct n -bit strings for $\ell = 1, 2, \dots, k$. We define the input matrix as

$$A := \sum_{\ell=1}^k \sigma_\ell V^{(\ell)} = \frac{1}{2^n} \sum_{y,z \in \{0,1\}^n} a_{y,z} e^{(y)} e^{(z)^T}, \quad (25)$$

where $\sigma_1, \sigma_2, \dots, \sigma_k$ are the singular values and we have implicitly defined $a_{y,z} = \sum_{\ell=1}^k \sigma_\ell (-1)^{x^{(\ell)} \cdot (y \oplus z)}$. The right and left singular vectors of A are the same, given by $v^{(\ell)} = \frac{1}{\sqrt{2^n}} \sum_{y \in \{0,1\}^n} (-1)^{x^{(\ell)} \cdot y} e^{(y)}$. As we show in the Appendix B, by construction, all rows of A have the same norm, therefore length-square sampling can be done by choosing rows uniformly at random. For length-square sampling of columns, we use rejection sampling: given a row y , we select a column index z uniformly at random and accept it with probability

$$p(z) = \frac{|a_{y,z}|^2}{\max_{z'} |a_{y,z'}|^2}. \quad (26)$$

The optimization is trivial since $\max_{z'} |a_{y,z'}|^2 = \sum_{\ell=1}^k (\sigma_\ell)^2$. Once row indices i_1, \dots, i_r and column indices j_1, \dots, j_c have been sampled, the entries of the $r \times c$ matrix C in FKV can be expressed as

$$C_{s,t} = \frac{1}{\sqrt{2^k c}} \frac{\|\mathbf{a}\|}{\sqrt{\sum_{s'=1}^r (a_{i_{s'},j_t})^2}} a_{i_s,j_t}, \quad (27)$$

where $\mathbf{a} = (a_1, a_2, \dots, a_{2^k})^T$ is a vector containing all the possible values of the entries $a_{y,z}$. The approximate singular values and vectors of A can be reproduced from the SVD of C as detailed in previous sections. To estimate the coefficient $\lambda_\ell = \frac{1}{\sigma_\ell^2} \langle \mathbf{v}^{(\ell)}, A^T \mathbf{b} \rangle$, we sample the random variable

$$\chi_{y,z} = \frac{\|\mathbf{a}\|^2}{2^k} \frac{(-1)^{x^{(\ell)} \cdot z}}{\sum_{\ell=1}^k \sigma_\ell (-1)^{x^{(\ell)} \cdot (y \oplus z)}} \sum_{\ell=1}^k \beta_\ell (-1)^{x^{(\ell)} \cdot y}, \quad (28)$$

and use its sample mean as an estimator for the coefficients. Finally, the approximate solution vector $\tilde{\mathbf{x}}$ is given by

$$\tilde{\mathbf{x}} = \frac{1}{\sqrt{2^n}} \sum_{\ell=1}^k \tilde{\lambda}_\ell \tilde{\mathbf{v}}^{(\ell)}. \quad (29)$$

Any entry of the approximate solution vector can be computed efficiently from the entries of the approximate right singular vectors $\tilde{\mathbf{v}}^{(\ell)}$, which can be obtained from the approximate SVD of C .

The crucial property of this construction is that the only dependency on dimension comes from the sampling of rows and columns, which can be done efficiently even for extremely high-dimensional problems. The limitations arise in computing the matrix C , which takes $O(r^2 c)$ time, and from computing the 2^k different matrix elements. This limits the values of r, c and k that can be explored in this approach. We fix the input matrices to have an extremely large dimension of $2^{50} \approx 10^{15}$. We set $r = c = 150$, the largest that could be handled in a reasonable amount of time. This means that we are aiming to approximate the SVD of the input matrix by using a smaller matrix whose dimension differs by thirteen orders of magnitude. The bit strings $x^{(\ell)}$ that define each singular vector are chosen uniformly at random. We study three examples with $k = \kappa = \kappa_\beta = 3, 5, 10$, employing $N = 10^4$ samples in the estimation of the coefficients. We report the average and standard deviation of the errors $\eta_\sigma = \frac{1}{k} \sum_{\ell=1}^k \frac{|\tilde{\sigma}_\ell - \sigma_\ell|}{\sigma_\ell}$ and $\eta_\lambda = \frac{1}{k} \sum_{\ell=1}^k \frac{|\tilde{\lambda}_\ell - \lambda_\ell|}{|\lambda_\ell|}$ over ten repetitions of the algorithms. For vectors, we evaluate the first L entries and similarly report the average and standard deviation of the errors $\eta_v = \frac{1}{kL} \sum_{i=1}^L \sum_{\ell=1}^k \frac{|v_i^{(\ell)} - \tilde{v}_i^{(\ell)}|}{|v_i^{(\ell)}|}$ and $\eta_x = \frac{1}{kL} \sum_{i=1}^L \sum_{\ell=1}^k \frac{|x_i - \tilde{x}_i|}{|x_i|}$. The results are summarized in Table 1.

These results show that for problems with the appropriate properties – namely very low rank and condition number – quantum-inspired algorithms can provide good estimates even for extremely high-dimensional problems. Importantly, these results are obtained in considerably less time than would be expected from the theoretical complexity bound $\tilde{O}(\kappa^{16} k^6 \|A\|_F^6 / \varepsilon^6)$. Even setting $\|A\|_F = 1$ and $\varepsilon = 1$, for $k = \kappa = \kappa_\beta = 10$, the bound would suggest that an order of 10^{22} operations would be needed; this is not compatible with the runtimes we experience. This suggests that these theoretical bounds do not in fact reflect the practical runtime of the algorithms and therefore should not be used to evaluate their performance. Nevertheless, our findings do reflect that the quality of the estimates worsens considerably as rank and condition number are increased.

Table 1: Errors in estimating singular values, singular vectors, coefficients, and solution vectors for input matrices of dimension $n, m = 2^{50}$. The coefficients were estimated using $N = 10^4$ samples and the errors in the vectors were calculated by computing their first $L = 100$ entries. The error bars correspond to one standard deviation taken over ten repetitions of the algorithm. Runtimes correspond to a Python implementation of the algorithms running on 20 Intel Xeon CPUs operating at 2.4GHz with access to 252GB of shared memory.

Case study	Errors				Runtime (hours)
	η_σ	η_v	η_λ	η_x	
$k = \kappa = \kappa_\beta = 3$	0.011 ± 0.012	0.124 ± 0.071	0.285 ± 0.276	0.414 ± 0.154	4.3
$k = \kappa = \kappa_\beta = 5$	0.129 ± 0.016	0.212 ± 0.070	0.530 ± 0.616	1.235 ± 0.615	8.6
$k = \kappa = \kappa_\beta = 10$	0.626 ± 0.018	1.619 ± 0.264	1.193 ± 1.720	4.138 ± 2.321	29.0

3.3 Random matrices

We study the quantum-inspired algorithm for linear systems of equations $A\mathbf{x} = \mathbf{b}$ for randomly chosen A and \mathbf{b} . To generate a random matrix A of dimension $m \times n$, rank k , and condition number κ , instead of creating A directly, we build the components of its singular value decomposition in matrix form, $A = U\Sigma V$. Here, U is an $m \times k$ matrix whose columns are the left singular vectors of A . Similarly, V is a $k \times n$ matrix whose rows are the right singular vectors, and Σ is a $k \times k$ diagonal matrix whose entries are the corresponding singular values of A , which has rank k by construction. The details of numerical methods used to generate matrices U , V and Σ are explained in Appendix C. The vector \mathbf{b} has been chosen to have the form $\mathbf{b} = \sum_{\ell=1}^k \beta_\ell \mathbf{u}^{(\ell)}$, with coefficients β_ℓ drawn from the standard normal distribution $\mathcal{N}(0, 1)$.

We consider matrices consisting of $m = 40,000$ rows and $n = 20,000$ columns. These numbers are chosen to deal with the general case of rectangular matrices, while working at the limit of dimensions for which it would still be possible to calculate SVDs in a reasonable amount of time. This means, for simplicity, that in this setting we are not required to use specialized data structures to perform length-square sampling: we can instead use fast built-in numerical sampling algorithms in Python which are very fast in practice. More precisely, generating a single sample from a distribution of dimension 40,000 using built-in algorithms takes approximately 2×10^{-3} seconds on a standard desktop computer. This is even more efficient when several samples are taken: generating one million samples from such vectors requires roughly 0.2 seconds.

In Figs. 3(a)-3(c) we monitor, respectively, the mean relative error $\eta_\sigma = \frac{1}{k} \sum_{\ell=1}^k \frac{|\tilde{\sigma}_\ell - \sigma_\ell|}{\sigma_\ell}$ in estimating the singular values as well as the errors $\eta_A = \|\tilde{A} - A\|_F / \|A\|_F$ and $\eta_A^+ = \|\tilde{A}^+ - A^+\|_F / \|A^+\|_F$ of the reconstructed matrices A and A^+ . The approximation is accurate across all singular values even for a relatively small number of sampled rows and columns.

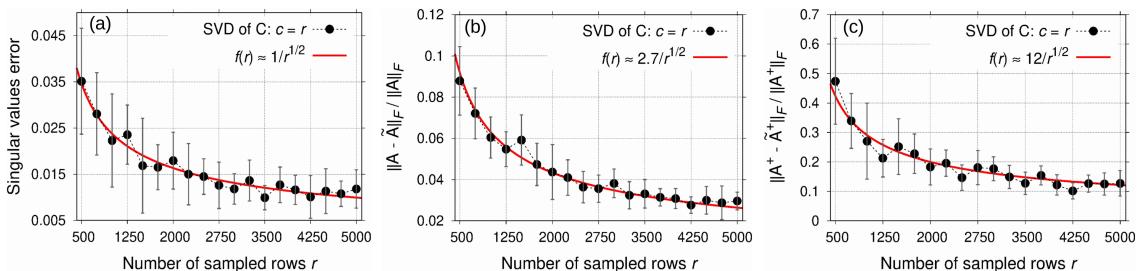


Figure 3: Errors in the approximate SVD of the random matrix A of dimension $40,000 \times 20,000$, rank $k = 5$ and condition number $\kappa = 5$ as the number of sampled rows r and columns $c = r$ is increased. The figures show the error of the (a) singular values $\eta_\sigma = \sum_{\ell=1}^k |\tilde{\sigma}_\ell - \sigma_\ell| / \sigma_\ell$ (b) reconstructed matrix $\eta_A = \|\tilde{A} - A\|_F / \|A\|_F$ with $\tilde{A} = \sum_{\ell=1}^k \tilde{\sigma}_\ell \tilde{\mathbf{u}}^{(\ell)} \tilde{\mathbf{v}}^{(\ell)T}$ and (c) reconstructed pseudo-inverse $\eta_A^+ = \|\tilde{A}^+ - A^+\|_F / \|A^+\|_F$ with $\tilde{A}^+ = \sum_{\ell=1}^k 1/\tilde{\sigma}_\ell \tilde{\mathbf{v}}^{(\ell)} \tilde{\mathbf{u}}^{(\ell)T}$. In all cases, error bars denote the standard deviation for 10 repetitions of the algorithm.

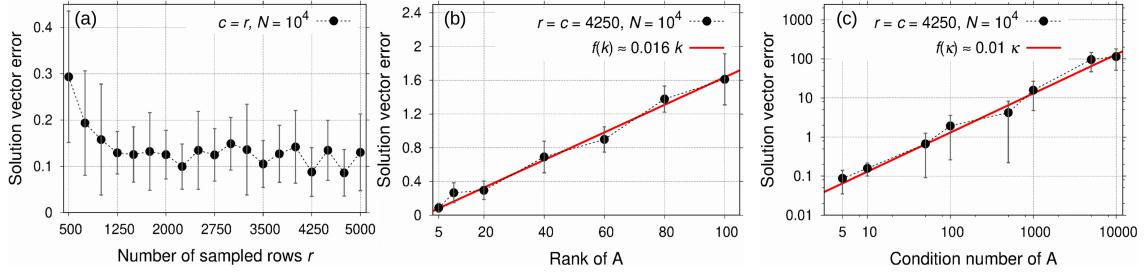


Figure 4: Error $\eta_{\mathbf{x}}$ of the solution vector for a $40,000 \times 20,000$ random matrix. (a) Error of the solution vector as a function of the number of sampled rows r and columns c entering the FKV algorithm, where we have fixed $c = r$. We set a rank $k = 5$ and condition number $\kappa = 5$. (b) Error in the solution vector as a function of matrix rank. There is a clear linear dependency as the matrix rank is increased. (c) Error in the solution vector as a function of the condition number of the input matrix. There is also a clear linear dependency, even for values spanning different orders of magnitude. In all cases, the coefficients λ_ℓ have been estimated by performing $N = 10^4$ samples. Error bars show the standard deviation over 10 independent repetitions of the FKV algorithm.

Empirically, the error in the approximation scales roughly as $1/\sqrt{r}$ for these random matrices, which is in agreement with the matrix Chernoff bound appearing in Theorem 3 of Ref. [10]. We also observe that the reconstructed matrix \tilde{A} has smaller errors compared to the approximated pseudo-inverse matrix A^+ . This is because in reconstructing matrix \tilde{A} , the largest singular values – which carry the dominant contributions – are approximated better by the FKV algorithm.

The main results are shown in Fig. 4, where we plot the error $\eta_{\mathbf{x}}$ of the approximated solution vector $\tilde{\mathbf{x}}$ calculated as the median of the relative errors $|\tilde{x}_i - x_i|/|x_i|$ with respect to all n entries of vector $\tilde{\mathbf{x}}$. First, we investigate in Fig. 4(a) how the error $\eta_{\mathbf{x}}$ changes as we increase the number of sampled rows r and columns $c = r$. All results have been obtained by taking $N = 10^4$ samples to estimate the coefficients λ_ℓ defined in Eq. (12). As expected, $\eta_{\mathbf{x}}$ decreases as the number of sampled rows and columns is increased. We have also superimposed error bars to show the standard deviation over 10 independent repetitions of the FKV algorithm.

In Fig. 4(b) we set $r = c = 4250$ and $\kappa = 5$ and investigate the error of the solution vector as a function of the rank k of matrix A . This plot reveals that the error $\eta_{\mathbf{x}}$ increases linearly with rank. For $k = 5$, a remarkably small relative error of 8.7% is found for the entries of the approximated vector. However, we observe that the accuracy of the algorithm deteriorates rapidly for larger values of the rank. A similar trend is observed in Fig. 4(c) for the solution vector error as a function of the condition number κ . Small errors of the order of 10% are only obtained for small condition numbers, while large values of $\eta_{\mathbf{x}} \geq 100\%$ already appear for $\kappa \geq 100$. Note that this is the behaviour we would expect from Eq. (21) if the dominant source of errors originates from coefficient estimation. Similar plots for the error in the approximate SVD of matrix A and in coefficients estimation are shown in Appendix C. Overall, we conclude that for random matrices, errors decrease as $1/\sqrt{r}$, and increase linearly with both rank and condition number.

Finally, Tables 2 and 3 respectively summarize the errors and running time of the quantum-inspired algorithm for the case where the matrix A has rank $k = 5$ and condition number $\kappa = 5$. In Table 2, we show the average relative errors of singular values η_σ and coefficients η_λ as defined in the previous section. We also report the relative errors $\eta_A = \|A - \tilde{A}\|_F/\|A\|_F$ and $\eta_{A^+} = \|A^+ - \tilde{A}^+\|_F/\|A^+\|_F$ of the reconstructed matrices \tilde{A} and \tilde{A}^+ . Sampling 4250 rows and columns from the original 40000×20000 matrix A is enough to have a good approximation full matrix. The largest error in Table 2 corresponds to the estimated coefficients.

Table 3 shows that the time to compute the approximate SVD using the FKV algorithm is orders-of-magnitude smaller than a direct calculation. Most of the overhead in running the quantum-inspired algorithms is associated with constructing the length-square distributions and estimating the coefficients. However, it is important to keep in mind that length-square distributions can in principle be constructed on-the-fly as the data from the matrix is generated, so it is best to separate this from the runtime of the actual algorithm. For this particular case of a matrix A with low rank and small condition number, we find that the quantum-inspired algorithm outperforms the direct calculation.

We can place these runtimes into perspective by attempting to identify whether they are due to our implementation or are inherent to the algorithm. There are 8×10^8 entries in this matrix, so our implementation is taking roughly $\sim 1.8 \times 10^{-6}$ seconds per entry to compute the length-squared distribution. Similarly, to estimate the coefficients, we generate 10^4 samples, repeated 10 times to calculate the median. This is done for each of the five coefficients. For each such sample, the entries of the right singular vectors of the matrix R must be estimated, which requires summing over the $r = 4,250$ sampled rows. Multiplying these numbers gives a total of 2×10^9 operations, for roughly 2.6×10^{-7} seconds per operation. In both cases, runtimes per operation are reasonably low, showing that total runtimes are inherently large given the number of required operations.

Table 2: Relative errors with corresponding standard deviations computed over 10 repetitions of the FKV algorithm of the quantum-inspired algorithm for linear systems $A\mathbf{x} = \mathbf{b}$. A is a random matrix with dimension $m = 40000$, $n = 20000$, rank $k = 5$, and condition number $\kappa = 5$. Approximated singular values and vectors were obtained by decomposing matrix C with dimension $r = 4250$ and $c = 4250$. The coefficients $\tilde{\lambda}$ have been estimated by performing $N = 10^4$ samples. Here, η_σ is the error in approximating singular values $\tilde{\sigma}$, while η_A and η_{A+} are the errors of the reconstructed matrices \tilde{A} and \tilde{A}^+ , respectively. η_λ is the error in estimating coefficients, and $\eta_{\mathbf{x}}$ is the the median of the relative errors $|\tilde{x}_i - x_i|/|x_i|$ with respect to all $n = 20000$ entries of the approximate vector $\tilde{\mathbf{x}}$.

Case study	Parameters			Error				
	r	c	N	η_σ	η_A	η_{A+}	η_λ	$\eta_{\mathbf{x}}$
Random matrix	4250	4250	10^4	0.010 ± 0.005	0.028 ± 0.004	0.101 ± 0.027	0.387 ± 0.191	0.087 ± 0.053

Table 3: Running times (in seconds) for the quantum-inspired algorithm applied to solve the system of linear equations $A\mathbf{x} = \mathbf{b}$ with the same parameters as reported in Table 2. The parameter t_{LS} is the time required to construct length-squared (LS) probability distributions over rows and columns. t_{SVD}^C accounts for the time spent sampling rows and columns from matrix A and building and decomposing the sampled matrix C . t_λ is the running time of the estimation of the coefficients λ . Analogous quantities are defined for the direct calculation method. $t_{\mathbf{x}}$ is the elapsed time to sample 500 entries of the approximated solution vector in the quantum-inspired algorithm and the time spent to compute the exact solution vector in the direct calculation method. Runtimes correspond to a Python implementation of the algorithms running on two Intel Xeon CPUs operating at 2.4GHz with access to 252GB of shared memory.

Case study	Quantum-inspired algorithm					Direct calculation			
	t_{LS}	t_{SVD}^C	t_λ	$t_{\mathbf{x}}$	t_{total}	t_{SVD}^A	t_λ	$t_{\mathbf{x}}$	t_{total}
Random matrix	1488.8	83.9	554.7	343	2470.4	5191.1	1.4	0.0003	5192.5

3.4 Portfolio optimization

We study an application of the quantum-inspired algorithm for linear systems of equations to a canonical financial problem: portfolio optimization. In its simplest version, commonly known as the Markowitz mean-variance model [18], the goal is to optimally invest wealth across n possible assets that are modeled only by their expected returns and correlations. Consider a vector of returns $\mathbf{r}_j = (r_{1,j}, r_{2,j}, \dots, r_{n,j})^T$, where $r_{i,j}$ is the return of asset i on the j -th day. The vector of

expected returns \mathbf{r} and the correlation matrix Σ are defined as

$$\mathbf{r} = \frac{1}{n} \sum_{j=1}^n \mathbf{r}_j, \quad \Sigma = \frac{1}{n} \sum_{j=1}^n \mathbf{r}_j \mathbf{r}_j^T. \quad (30)$$

The vector \mathbf{r} captures the expected returns for each asset, while the correlation matrix Σ expresses fluctuations around mean values, which in the model is interpreted as risk. An optimal portfolio in this setting is one that minimizes risk for any given target expected return. This corresponds to the optimization problem

$$\underset{\mathbf{w}}{\text{Minimize:}} \quad \mathbf{w}^T \Sigma \mathbf{w} \quad (31)$$

$$\text{subject to: } \mathbf{r}^T \mathbf{w} = \mu, \quad (32)$$

where \mathbf{w} is a portfolio allocation vector that determines how much wealth is invested in each asset. We allow \mathbf{w} to have negative entries, which is interpreted as an indication to short-sell the corresponding asset. As discussed in Ref. [22], this problem can be equivalently cast as a linear system of equations:

$$\begin{bmatrix} 0 & \mathbf{r}^T \\ \mathbf{r} & \Sigma \end{bmatrix} \begin{bmatrix} \nu \\ \mathbf{w} \end{bmatrix} = \begin{bmatrix} \mu \\ \mathbf{0} \end{bmatrix}, \quad (33)$$

where ν is a Lagrange multiplier. In this case, the linear system $A\mathbf{x} = \mathbf{b}$ is given by $A = \begin{bmatrix} 0 & \mathbf{r}^T \\ \mathbf{r} & \Sigma \end{bmatrix}$, $\mathbf{b} = \begin{bmatrix} \mu \\ \mathbf{0} \end{bmatrix}$, and $\mathbf{x} = \begin{bmatrix} \nu \\ \mathbf{w} \end{bmatrix}$. To compute a solution to the linear system, we employ the Moore-Penrose pseudo-inverse A^+ , which can be expressed in terms of the singular values and vectors of A as $A^+ = \sum_{\ell=1}^k \frac{1}{\sigma_\ell} \mathbf{v}^{(\ell)} \mathbf{u}^{(\ell)T}$. The solution vector \mathbf{x} is then

$$\begin{aligned} \mathbf{x} = A^+ \mathbf{b} &= \sum_{\ell=1}^k \frac{1}{\sigma_\ell} \langle \mathbf{u}^{(\ell)}, \mathbf{b} \rangle \mathbf{v}^{(\ell)} \\ &= \sum_{\ell=1}^k \frac{1}{\sigma_\ell^2} \langle \mathbf{v}^{(\ell)}, A^T \mathbf{b} \rangle \mathbf{v}^{(\ell)} = \sum_{\ell=1}^k \lambda_\ell \mathbf{v}^{(\ell)}, \end{aligned} \quad (34)$$

where we have identified the coefficients $\lambda_\ell = \langle \mathbf{v}^{(\ell)}, A^T \mathbf{b} \rangle / \sigma_\ell^2$ as in Eq. (3). The vector \mathbf{b} has only one non-zero entry, so it holds that $A^T \mathbf{b} = \mu A_1$, where A_1 is the first row of A .

To test the algorithm in this setting, we employ publicly-available pricing data for the stocks comprising the S&P 500 stock index during the five-year period 2013-2018 [21]. Since not all stocks remain in the index during this time, we restrict to the 474 stocks that do, leading to a matrix A of dimension 475×475 . Returns are calculated on a daily basis based on opening price. The matrix has full rank of $k = 475$, its largest and smallest singular values are respectively $\sigma_{max} = 40.2$, $\sigma_{min} = 1.8 \times 10^{-4}$, and its condition number is $\kappa = 2.23 \times 10^4$. More details can be found in Appendix D. We set a target return μ equal to the average return over all 474 stocks in the index.

In Tables 4 we report the errors of the quantum-inspired algorithm. The SVD of matrix A was approximated by sampling 340 rows and columns from the full matrix. As in the previous example, we took $N = 10^4$ samples to estimate the coefficients λ_ℓ . The relative errors characterizing the approximate solution vector are considerable, with values of $\eta_x = 0.74$. In this example, we have calculated the approximate solution vector by using Eq. (34) within a low-rank approximation with $k = 10$. Notice from Table 4 that the large errors η_A^+ in the reconstructed pseudo-inverse matrix and estimated coefficients η_λ translate into large relative errors of the solution vector.

Finally, in Table 5 we report the running-times for the algorithm. Since the matrices involved in this example are small, exact and approximate SVDs can be computed quickly, but for the quantum-inspired algorithm, there is significant overhead due to coefficient estimation and sampling from the solution vector.

Table 4: Relative errors with corresponding standard deviations computed over 10 repetitions of the quantum-inspired algorithm applied to a portfolio optimization problem. Approximated singular values and vectors were obtained by decomposing matrix C with dimension $r = 340$ and $c = 340$. Coefficients λ have been estimated by performing $N = 10^4$ samples. Here, η_{σ} is the error in approximating singular values $\tilde{\sigma}$, while η_A and η_{A+} are the errors of the reconstructed matrices A and \tilde{A}^+ , respectively. η_{λ} is the error in estimating coefficients, and η_x is the median of the relative errors $|\tilde{x}_i - x_i|/|x_i|$ with respect to all $n = 475$ entries of the approximated vector \tilde{x} .

Case study	Parameters			Error				
	r	c	N	η_{σ}	η_A	η_A^+	η_{λ}	η_x
S&P 500 stock index	340	340	10^4	0.08 ± 0.02	0.16 ± 0.03	1.13 ± 0.12	1.58 ± 0.79	0.74 ± 0.19

Table 5: Running times (in seconds) for the quantum-inspired algorithm applied to the portfolio optimization problem. The algorithm was run using the parameters reported in Table 4. The parameter t_{LS} is the time required to construct length-squared (LS) probability distributions over rows and columns. t_{SVD}^C accounts for the time spent sampling rows and columns from matrix A and building and decomposing the sampled matrix C . t_{λ} is the running time of the estimation of the coefficients λ . Analogous quantities are defined for the direct calculation method. t_x is the elapsed time to sample 50 entries of the approximated solution vector in the quantum-inspired algorithm, and the time spent to compute the exact solution vector in the direct calculation method. Runtimes correspond to a python implementation of the algorithms running on two Intel Xeon CPUs operating at 2.4GHz with access to 252GB of shared memory.

Case study	Quantum-inspired algorithm					Direct calculation			
	t_{LS}	t_{SVD}^C	t_{λ}	t_x	t_{total}	t_{SVD}^A	t_{λ}	t_x	t_{total}
S&P 500 stock index	0.46	0.36	88.7	24.63	114.15	0.15	0.0005	0.0006	0.15

3.5 Movie recommendations

We analyze the performance of the algorithm for [recommendation systems on the MovieLens 100K database](#) [13], which consists of a preference matrix with 100,000 ratings from 611 users across 9,724 movies. Ratings are specified on a half-star scale in the range $[0.5, 5]$. For example, an entry of 4.5 corresponds to a rating of four-and-a-half stars. The lowest possible rating is 0.5 and the zero entries in the preference matrix correspond to movies that have not been rated by the user. Since most users watch only a small fraction of available movies, the matrix is sparse. The preference matrix has full rank of $k = 611$, its largest and smallest singular values are respectively $\sigma_{max} = 534.4$, $\sigma_{min} = 2.95$, and its condition number is $\kappa = 181.2$. More details can be found in Appendix D.

The goal of the algorithm is to predict missing ratings and subsequently recommend movies that have a high predicted rating. The approach is to consider a low-rank approximation A' of the preference matrix A , which can be written in terms of the SVD of A as $A' = \sum_{\ell=1}^k \sigma_{\ell} \mathbf{u}^{(\ell)} \mathbf{v}^{(\ell)T}$. Here k is a parameter of choice in the algorithm, not the actual rank of the original preference matrix. This low-rank approximation can be equivalently written as $A' = A \sum_{\ell=1}^k \mathbf{v}^{(\ell)} \mathbf{v}^{(\ell)T}$ and therefore the i -th row of A' is given by

$$A'_i = A_i \sum_{\ell=1}^k \mathbf{v}^{(\ell)} \mathbf{v}^{(\ell)T} = \sum_{\ell=1}^k \langle A_i^T, \mathbf{v}^{(\ell)} \rangle \mathbf{v}^{(\ell)T}. \quad (35)$$

Expressing A'_i as a column vector, we recognize the target vector as

$$\mathbf{x} := {A'_i}^T = \sum_{\ell=1}^k \lambda_{\ell} \mathbf{v}^{(\ell)}, \quad (36)$$

where $\lambda_{\ell} = \langle A_i^T, \mathbf{v}^{(\ell)} \rangle$ as in Eq. (4).

The matrix A' is typically no longer sparse, so the previously missing ratings can be predicted based on the non-zero entries of this low-rank approximation. In this sense, the algorithm can be interpreted as an approximate reconstruction of a preference matrix P that contains ratings for all users and movies, of which we are only given a few sample ratings in the form of the sparse matrix A . If P is well approximated by a low-rank matrix and A has sufficiently many entries of P , then a low-rank approximation of A is also close to P .

As in the previous section, we summarize the results of applying the quantum-inspired algorithm for recommendation systems in Tables 6 and 7. We observe from Table 6 that sampling 450 rows and 4500 columns from the full preference matrix is enough to describe the first ten singular values with a small relative error of $\eta_\sigma = 0.06$. On the other hand, the errors of the reconstructed matrices A , A^+ , and estimated coefficients show larger values of 0.32, 0.66 and 0.58, respectively. Finally, we find that the entries of the approximated solution vector, computed using Eq. (36) within the low-rank approximation with $k = 10$ show typical relative errors of $\eta_x = 0.71$. In Table 7 we compare the running times between the quantum-inspired algorithm and the direct calculation method. As in the portfolio optimization example, the matrices are small enough that computing their SVD can be done in a short time. The runtime of the quantum-inspired algorithm is again dominated by the coefficient estimation step, even when the resulting errors in estimating these coefficients is relatively large.

The numerical examples we have studied indicate that when applied to moderately-sized real-life data sets, because they rely on a more intricate procedure, the quantum-inspired algorithms take more time than exact diagonalization, and because they rely on sampling for coefficient estimation, lead to higher inaccuracies. These results suggest that in order to provide a speedup over preexisting classical algorithms, the quantum-inspired algorithms must be applied to very large data sets where exact diagonalization is impossible and where even the linear scaling of FKV prevents its direct application.

Table 6: Relative errors with corresponding standard deviations computed over 10 repetitions of the FKV algorithm in the quantum-inspired algorithm for recommendation systems. Here, matrix A is a preference matrix with dimension $m = 611$ and $n = 9274$. Approximated singular values and vectors were obtained by decomposing matrix C with dimension $r = 450$ and $c = 4500$. The coefficients λ have been estimated by performing $N = 10^4$ samples. Here, η_σ is the error in approximating singular values $\tilde{\sigma}$, while η_A and η_{A^+} are the errors of the reconstructed matrices \tilde{A} and \tilde{A}^+ , respectively. η_λ is the error in estimating coefficients, and η_x is the median of the relative errors $|\hat{x}_i - x_i|/|x_i|$ with respect to all $n = 9274$ entries of the approximated solution vector x .

Case study	Parameters			Error				
	r	c	N	η_σ	η_A	η_{A^+}	η_λ	η_x
MovieLens 100K	450	4500	10^4	0.06 ± 0.01	0.32 ± 0.02	0.66 ± 0.05	0.58 ± 0.18	0.71 ± 0.13

Table 7: Running times (in seconds) for the quantum-inspired algorithm for recommendation systems. The algorithm was run using the parameters reported in Table 6. The parameter t_{LS} is the time required to construct length-squared (LS) probability distributions over rows and columns. t_{SVD}^C accounts for the time spent sampling rows and columns from matrix A and building and decomposing the sampled matrix C . t_λ is the running time of the estimation of the coefficients λ . Analogous quantities are defined for the direct calculation method. t_x is the elapsed time to sample 500 entries of the approximated solution vector in the quantum-inspired algorithm and the time spent to compute the exact solution vector in the direct calculation method. Runtimes correspond to a python implementation of the algorithms running on two Intel Xeon CPUs operating at 2.4GHz with access to 252GB of shared memory.

Case study	Quantum-inspired algorithm					Direct calculation			
	t_{LS}	t_{SVD}^C	t_λ	t_x	t_{total}	t_{SVD}^A	t_λ	t_x	t_{total}
MovieLens 100K	11.05	6.22	124.2	14.2	155.7	2.00	0.0003	0.001	2.00

3.6 Practical complexity of quantum-inspired algorithms

In this section, we comment on the complexity of quantum-inspired algorithms based on the numerical tests and theoretical analysis performed thus far. The matrix Chernoff bound of Ref. [10] indicates that the error in the FKV algorithm scales as $\varepsilon = O(1/\sqrt{r})$ with the dimension r of the matrix C . This is consistent with the results of empirical tests on Gaussian random matrices as reported in Sec. 3.3. Since the complexity of computing the SVD of an $r \times c$ matrix with $c = O(r)$ scales as $O(r^3)$, this shows that the complexity of approximate SVD calculations with the FKV algorithm grows as $O(1/\varepsilon^6)$. This is compatible with the complexity bound from Ref. [10], suggesting that complexity with respect to error may be tight.

It was shown in Ref. [9] that the FKV algorithm can produce an approximate SVD with error ε by asymptotically setting $r = c = O(\max\{k^4/\varepsilon^2, k^2/\varepsilon^4\})$. This was improved in Ref. [10], where it was shown that an error ε could be obtained when $r \neq c$ by choosing $r = O(c) = O(k^2/\varepsilon^2)$. This implies a complexity of $O(k^6/\varepsilon^6)$ when all other parameters are fixed. Importantly, such a choice of r only makes sense if the resulting $r \times c$ matrix is smaller than the original input matrix, i.e., if $r \leq m$ and $c \leq n$. In practice, it is always possible to set $r \leq m$ and $c \leq n$, resulting in a guaranteed reduced runtime in computing the SVD at the expense of an error in the approximation. This is the strategy adopted in our experiments. It is important to understand that, when implemented properly, the runtime of quantum-inspired algorithms is never larger than a direct calculation. The crucial point is whether the resulting errors are sufficiently low, since it becomes prohibitively expensive to reduce them to arbitrarily small values due to the $O(1/\varepsilon^6)$ scaling.

In regimes where the largest source of errors originates from coefficient estimation, it is possible that the runtime of quantum-inspired algorithms is dominated by the complexity of this step. This occurred in several examples considered in this work, including high-dimensional problems with low-rank matrices. The complexity of coefficient estimation is captured in Eqs. (13) and (15), which have smaller exponents than the complexity bounds of Refs. [10] and [31]. Numerical tests support the theoretical calculations. The linear dependency shown in Fig. 4 between error, rank, and condition number for a fixed number of samples N indicates that N must be a polynomial in $(k\kappa/\varepsilon)$. Our bound states that $N = O(k^2\kappa^2/\varepsilon^2)$ in accordance with this behavior. Depending on the properties of the input problem, different complexity regimes are possible, each determined by the steps of the algorithms that dominate the runtime.

4 Conclusion

In terms of asymptotic complexity, quantum-inspired algorithms constitute a breakthrough in our understanding of the boundary between classical and quantum algorithms: they imply that certain linear algebra problems can be performed in sublinear classical time. Our results show that the proven complexity bounds for these algorithms do not actually reflect their practical runtimes. The proven complexity of the linear systems algorithm is $\tilde{O}(\kappa^{16}k^6\|A\|_F^6/\varepsilon^6)$ [10], while for recommendation systems it is $\tilde{O}(k^{12}/(\varepsilon^{12}\eta^6))$ [31], where η is an error parameter. In the implementation of the algorithm, we observe a significantly faster runtime than these bounds would suggest. This indicates that care must be taken when employing these bounds to make statements about the performance of the algorithms. Our results are also encouraging for complexity theorists aiming to improve these bounds.

In our analysis, we showed that quantum-inspired algorithms can provide reasonably low errors in relatively short times even in the regime of extremely large-dimensional problems. However, in our implementation, the performance requires matrices of very low rank and condition number; the errors in the outputs grow noticeably when rank and condition number are increased. Compared to previously-known methods such as the Frieze-Kannan-Vempala (FKV) algorithm, quantum-inspired algorithms differ in their use of sampling techniques to estimate coefficients and sample from the solution vectors. On the one hand, a direct calculation, as done in FKV, requires linear time in the dimension of the input matrix, which can be done extremely fast even for problems of large size. On the other hand, sampling methods scale polylogarithmically with dimension, but

they incur additional polynomial costs in the rank, condition number, and error in the estimation. Therefore, one conclusion from our work is that despite their asymptotic scaling, sampling techniques for coefficient estimation do not lead to practical improvements compared to the direct computation used in previously-existing classical algorithms: **quantum-inspired techniques only become advantageous for problems of extremely large dimension.**

Additionally, we have shown that when employing the FKV algorithm for approximate SVD, a scaling of $O(k^6/\varepsilon^6)$ with respect to the error ε and rank k follows directly from the matrix Chernoff bound and the error bounds from Ref. [10]. Therefore, these exponents are likely to be fundamental and not the result of proof techniques. Overall, **our results indicate that quantum-inspired algorithms perform well in practice, but only under the restrictive conditions of large-dimensional input matrices with very low rank and condition number.** It remains unclear whether datasets with these properties actually appear in practice.

By contrast, the complexity of quantum algorithms does not depend on matrix rank, and therefore they work properly even for full-rank problems. Matrices originating in practical applications are often sparse and they typically have large effective ranks; these are the problems that can be tackled by quantum algorithms. Dealing with matrices of large condition number remains challenging for all techniques.

Several open questions remain. The numerical tests on random matrices performed in Section 3.3 indicate that when all other parameters are fixed, errors in the approximation of the output vector increase linearly with rank and condition number. It is important to understand whether this is a general feature of the algorithm. Additionally, it is still of interest to determine the empirical scaling of quantum-inspired algorithms with respect to all relevant parameters. Our results indicate that for linear systems, the dependency on rank and error is likely tight, but the findings of section 3.2 show that not all exponents are tight. Finally, it would be of great interest to implement highly-optimized versions of our code on powerful supercomputers capable of operating much larger datasets than those tested in this work.

Acknowledgements

We thank Christian Weedbrook, Nathan Killoran, Nicolás Quesada, and Iordanis Kerenidis for useful discussions. We are grateful to András Gilyén for providing valuable feedback on an early draft of this manuscript. S. Lloyd was funded by AFOSR under a MURI on Optimal Quantum Measurements and State Verification, by IARPA under the QEO program, by ARO, and by NSF.

Appendix

A Worst-case hardness of coefficient estimation

Here we show that the coefficient estimation cannot generally be done in sublinear time. The main strategy is to encode an NP-Hard problem – approximating the partition function of Ising models – into the coefficients to be estimated. The exponential time hypothesis then implies that this coefficients cannot be estimated in sublinear time.

We begin by noting that, for any sum of the form $\mathcal{Z} = \sum_{i=1}^n w_i$, it is possible to find vectors \mathbf{y}, \mathbf{z} such that $w_i = y_i z_i$ for all $i = 1, 2, \dots, n$. This means that any sum \mathcal{Z} can be expressed as an inner product $\mathcal{Z} = \langle \mathbf{y}, \mathbf{z} \rangle$. If estimating \mathcal{Z} is NP-Hard, then so is calculating the coefficient $\lambda = \langle \mathbf{x}, \mathbf{y} \rangle$. As a concrete example, consider the partition function of the Ising model

$$\mathcal{Z}_I = \sum_s e^{-\beta E(s)}, \quad (37)$$

where $s = (s_1, \dots, s_m) \in \{-1, 1\}^m$ is the configuration of m spins, β is the inverse temperature, and $E(s) = -\sum_{i < j} J_{ij} s_i s_j - \mu \sum_j h_j s_j$ is the energy of the configuration. It is known that

approximating this partition function is NP-Hard in a general setting [29]. Since there are $n := 2^m$ different configurations s , each can be labeled by a number $k = 1, \dots, n$, also expressing energies as a function of k . As discussed above, we can also define n -dimensional vectors \mathbf{x}, \mathbf{y} such that $\mathcal{Z}_I = \sum_{k=1}^n x_k y_k = \langle \mathbf{x}, \mathbf{y} \rangle$, for instance by setting $x_k = y_k = e^{-\beta E(k)/2}$. Therefore, there exist vectors such that computing their inner product is as difficult as computing the NP-Hard partition function \mathcal{Z}_I . The exponential time hypothesis [15] states that the best possible algorithms for computing NP-complete problems require $2^{\delta m}$ time for some $\delta > 0$. Thus, unless this hypothesis is false, estimating coefficients requires $2^{\delta m} = O(n)$ time in a worst-case setting.

B Formulas for high-dimensional problems

We begin by showing that all rows of the high-dimensional matrices considered in this paper have equal norm. Denote by A_y the y -th row of A . This can be written as

$$\begin{aligned} A_y &= \frac{1}{2^n} \sum_{\mathbf{z}} \sum_{\ell=1}^k \sigma_\ell(-1)^{\mathbf{x}^{(\ell)} \cdot (\mathbf{z} \oplus \mathbf{y})} e^{(\mathbf{z})^T} \\ &= \frac{1}{2^n} \sum_{\mathbf{z}'} \left[\sum_{\ell=1}^k \sigma_\ell(-1)^{\mathbf{x}^{(\ell)} \cdot \mathbf{z}'} \right] e^{(\mathbf{z}' \oplus \mathbf{y})^T} \\ &= \frac{1}{2^n} \sum_{\mathbf{z}'} \left[\sum_{\ell=1}^k \sigma_\ell(-1)^{\mathbf{x}^{(\ell)} \cdot \mathbf{z}'} \right] \tilde{e}^{(\mathbf{z}')^T}, \end{aligned} \quad (38)$$

where $\mathbf{z}' = \mathbf{z} \oplus \mathbf{y}$ and we have defined a new basis $\{\tilde{e}^{(\mathbf{z}')}\} = \{e^{(\mathbf{z}' \oplus \mathbf{y})}\}$. Therefore the norm of the row is

$$\|A_y\|^2 = \frac{1}{2^n} \sum_{\mathbf{z}'} \left[\sum_{\ell=1}^k \sigma_\ell(-1)^{\mathbf{x}^{(\ell)} \cdot \mathbf{z}'} \right]^2, \quad (39)$$

which does not depend on \mathbf{y} , meaning it is the same for all rows. Indeed, the rows of A have the same elements; they simply appear in different positions depending on the corresponding row label \mathbf{y} .

The entries of matrix C can be expressed in terms of known quantities as

$$\begin{aligned} C_{s,t} &= \frac{\|A\|_F^2}{\sqrt{rc} \|A_{i_s}\| \|R_{\cdot,j_t}\|} A_{i_s,j_t} \\ &= \frac{2^n \|A_{i_s}\|}{\sqrt{rc} \|R_{\cdot,j_t}\|} A_{i_s,j_t} \\ &= \frac{2^n \|\mathbf{a}\|}{\sqrt{2^{n+k}} \sqrt{rc} \|R_{\cdot,j_t}\|} A_{i_s,j_t} \\ &= \frac{1}{\sqrt{2^k c}} \frac{\|\mathbf{a}\|}{\sqrt{\sum_{s'=1}^r A_{i_{s'},j_t}^2}} A_{i_s,j_t} \\ &= \frac{1}{\sqrt{2^k c}} \frac{\|\mathbf{a}\|}{\sqrt{\sum_{s'=1}^r a_{i_{s'},j_t}^2}} a_{i_s,j_t}, \end{aligned} \quad (40)$$

where $\mathbf{a} = (a_1, a_2, \dots, a_{2^k})^T$ is the vector containing all the possible values of the entries $a_{\mathbf{y},\mathbf{z}}$ and we used $\|R_{\cdot,j_t}\| = \sqrt{\frac{2^n}{r}} \sqrt{\sum_{s'=1}^r A_{i_{s'},j_t}^2}$.

Finally, to estimate the coefficient $\lambda_\ell = \frac{1}{\sigma_\ell^2} \langle \mathbf{v}^{(\ell)}, A^T \mathbf{b} \rangle$, we sample the random variable $X_{\mathbf{y},\mathbf{z}} =$

$\frac{\|A\|_F^2}{A_{y,z}} b_y v_z^{(\ell)}$ whose average is equal to $\langle \mathbf{v}^{(\ell)}, A^T \mathbf{b} \rangle$. This random variable is given by

$$\begin{aligned} X_{y,z} &= \frac{\|A\|_F^2}{A_{ij}} v_z^{(\ell)} b_y \\ &= \frac{2^{n-k} \|\mathbf{a}\|^2}{a_{y,z}} v_z^{(\ell)} b_y \\ &= \frac{2^{n-k} \|\mathbf{a}\|^2}{a_{y,z}} \frac{(-1)^{\mathbf{x}^{(\ell)} \cdot z}}{\sqrt{2^n}} \frac{1}{\sqrt{2^n}} \sum_{\ell=1}^k \beta_\ell (-1)^{\mathbf{x}^{(\ell)} \cdot y} \\ &= \frac{\|\mathbf{a}\|^2}{2^k} \frac{(-1)^{\mathbf{x}^{(\ell)} \cdot z}}{\sum_{\ell=1}^k \sigma_\ell (-1)^{\mathbf{x}^{(\ell)} \cdot (y \oplus z)}} \sum_{\ell=1}^k \beta_\ell (-1)^{\mathbf{x}^{(\ell)} \cdot y}. \end{aligned}$$

C Complementary results for random matrices

In this section we explain the methodology to generate the low-rank Gaussian random matrices used in Section 3.3 and show further numerical results to benchmark the accuracy of the algorithms used to approximate the SVD of matrix A and to estimate the coefficients λ . To generate U , we first sample an $m \times k$ Gaussian random matrix U' with entries drawn independently from the standard normal distribution $\mathcal{N}(0, 1)$. The columns of this matrix are generally not orthogonal, so we perform a QR decomposition $U' = QR$, where Q is an $m \times k$ orthogonal matrix and R is a $k \times k$ upper triangular matrix. We then simply set $U := Q$. An analogous method is used to generate V . Finally, given a target condition number κ , we select the largest singular value σ_{\max} uniformly at random in the interval $[1, 500]$. This fixes the smallest singular value as $\sigma_{\min} = \sigma_{\max}/\kappa$. Other singular values are sampled from the interval $(\sigma_{\min}, \sigma_{\max})$ by using the quadrant law for singular values [28].

In Figs. 5(a)-5(c), we benchmark the approximate SVD of random matrices with the same dimensions as the one in Fig. 3, but with increasing values of the rank k . In this case the condition number of these matrices has been fixed to $\kappa = 5$. Although the errors in the estimation of singular values does not appear to be strongly influenced by the rank, errors in the approximated singular vector become apparent in the relative errors η_A and η_{A^+} for the reconstructed matrices as k increases. In particular, we notice that the error of the pseudo-inverse matrix, which is the one required to compute the solution vector, scales linearly with rank showing relatively large errors larger than 45 % for $k = 100$.

In Figs. 6(a)-6(c) we benchmark the approximate SVD of low-rank random matrices with $k = 5$ while increasing their condition number. We observe from Figs 6(a)-6(b) almost no dependence of the errors η_σ and η_A with condition number. On the contrary, Fig. 6(c) demonstrates that the relative error η_{A^+} of the reconstructed pseudo-inverse matrix grows linearly with condition number, exhibiting values greater than 100% already for $\kappa > 100$. The latter is somehow expected since the error of the approximated singular vectors are propagated to matrix A^+ roughly as $1/\sigma_\ell$, i.e., the smaller the singular values, the larger the error η_{A^+} due to the approximate SVD of matrix A .

In Figs. 7(a)-(c) we investigate the error $\eta_\lambda = \sum_{\ell=1}^k |\tilde{\lambda}_\ell - \lambda_\ell|/|\lambda_\ell|$ of the estimated coefficients λ_ℓ . In all cases, the estimation of the coefficients was performed by taking $N = 10^4$ samples of the corresponding random variable. In Fig. 7(a) we show that the values of η_λ display a weak dependence on the number of sampled rows and columns. We notice, however, that this is not the case as we increase the rank and condition number of the matrix. For example, large values of $\eta_\lambda > 100$ % are observed for matrices with rank $k > 20$ as shown in Fig. 7(b). Furthermore, Fig. 7(c) shows clearly that even for a matrix with rank as low as $k = 5$, when the condition number is increased, the error of the estimated coefficients ramps up rapidly to large values.

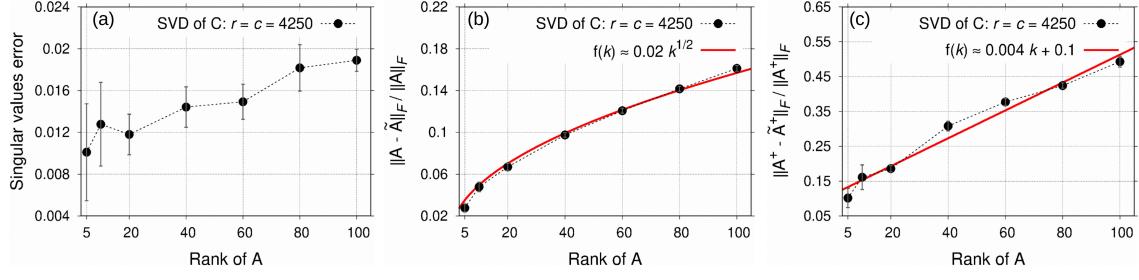


Figure 5: Errors in the approximate SVD of random matrices A of dimension $40,000 \times 20,000$ and condition number $\kappa = 5$ as the rank k is increased. The figures show the error of the (a) singular values $\eta_\sigma = \sum_{\ell=1}^k |\tilde{\sigma}_\ell - \sigma_\ell|/\sigma_\ell$, (b) reconstructed matrix $\eta_A = \|\tilde{A} - A\|_F/\|A\|_F$ with $\tilde{A} = \sum_{\ell=1}^k \tilde{\sigma}_\ell \tilde{\mathbf{u}}^{(\ell)} \tilde{\mathbf{v}}^{(\ell)T}$, and (c) reconstructed pseudo-inverse $\eta_A^+ = \|\tilde{A}^+ - A^+\|_F/\|A^+\|_F$ with $\tilde{A}^+ = \sum_{\ell=1}^k 1/\tilde{\sigma}_\ell \tilde{\mathbf{v}}^{(\ell)} \tilde{\mathbf{u}}^{(\ell)T}$. In all cases, error bars denote the standard deviation for 10 repetitions of the algorithm.

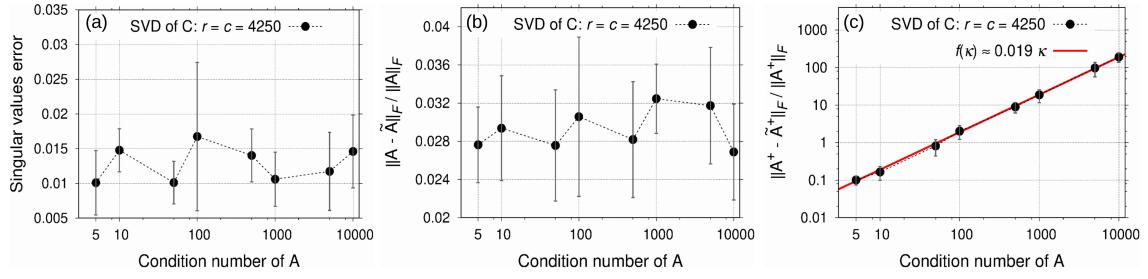


Figure 6: Errors in the approximate SVD of random matrices A of dimension $40,000 \times 20,000$ and rank $k = 5$ as the condition number κ is increased from 5 to 10000. Error of the (a) singular values $\eta_\sigma = \sum_{\ell=1}^k |\tilde{\sigma}_\ell - \sigma_\ell|/\sigma_\ell$ (b) reconstructed matrix $\eta_A = \|\tilde{A} - A\|_F/\|A\|_F$ with $\tilde{A} = \sum_{\ell=1}^k \tilde{\sigma}_\ell \tilde{\mathbf{u}}^{(\ell)} \tilde{\mathbf{v}}^{(\ell)T}$ and (c) reconstructed pseudo-inverse $\eta_A^+ = \|\tilde{A}^+ - A^+\|_F/\|A^+\|_F$ with $\tilde{A}^+ = \sum_{\ell=1}^k 1/\tilde{\sigma}_\ell \tilde{\mathbf{v}}^{(\ell)} \tilde{\mathbf{u}}^{(\ell)T}$. In all cases, error bars denote the standard deviation for 10 repetitions of the algorithm.

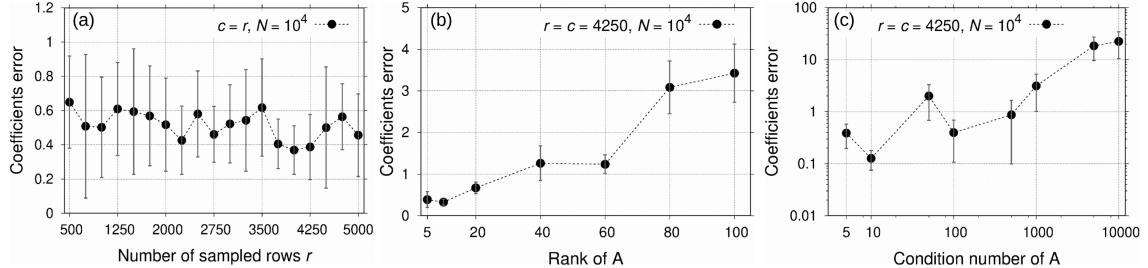


Figure 7: Mean relative error $\eta_\lambda = \sum_{\ell=1}^k |\tilde{\lambda}_\ell - \lambda_\ell|/|\lambda_\ell|$ of the estimated coefficients $\lambda_\ell = \langle \mathbf{v}^{(\ell)}, \mathbf{A}^T \mathbf{b} \rangle$ for random matrices matrices A of dimension $40,000 \times 20,000$ as a function of (a) the number of rows r and columns c sampled from matrix A with $k = \kappa = 5$, (b) the rank k and (c) condition number κ of matrix A . In all cases, coefficients have been estimated by performing $N = 10^4$ samples of the random variable. Error bars denote the standard deviation for 10 repetitions of the algorithm.

D Singular values and coefficients for portfolio optimization and recommendation systems

In Fig. 8, we show the first ten singular values of the portfolio optimization matrix as well as the exact and approximate coefficients. Similarly, Fig. 9 shows analogous plots for the case of movie recommendations.

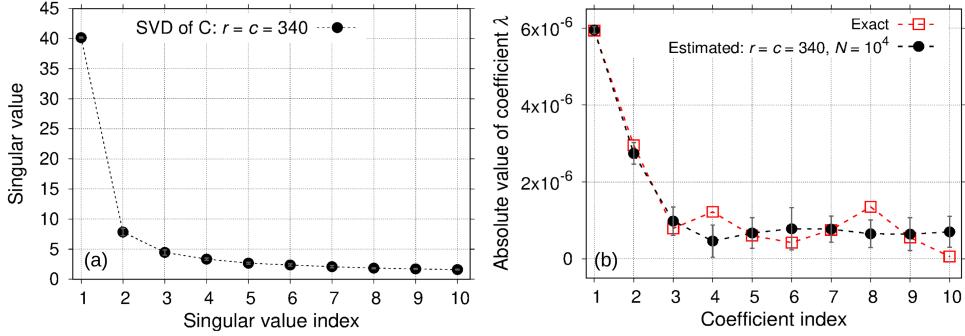


Figure 8: (a) First ten singular values σ_ℓ calculated with FKV by sampling $r = c = 340$ rows and columns of the full-rank input matrix A with dimension 475×475 for the portfolio optimization problem on S&P 500 data (see Sec. 3.4). (b) Exact and estimated coefficients λ_ℓ for indices $\ell = 1, \dots, 10$ by taking $N = 10^4$ samples of the corresponding random variable. The coefficient corresponding to the largest singular value is estimated very accurately with a relative error of 0.1 %. Error bars denote the standard deviation for 10 repetitions of the algorithm.

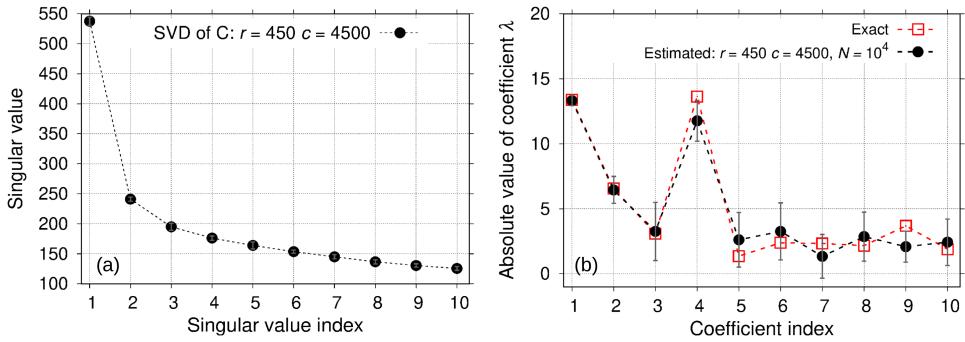


Figure 9: (a) First ten singular values σ_ℓ calculated with FKV by sampling $r = 450$ rows and $c = 4500$ columns of the full-rank preference matrix A of dimension 611×9724 for the recommendation system on the MovieLens 100K database (see Sec. 3.5). (b) Exact and estimated coefficients λ_ℓ for indices $\ell = 1, \dots, 10$ by taking $N = 10^4$ samples of the corresponding random variable. The coefficients corresponding to the two largest singular values are estimated very accurately with relative errors of 0.4 % and 1.5 %, respectively. Error bars denote the standard deviation for 10 repetitions of the algorithm.

References

- [1] Boaz Barak, Ankur Moitra, Ryan O'Donnell, Prasad Raghavendra, Oded Regev, David Steurer, Luca Trevisan, Aravindan Vijayaraghavan, David Witmer, and John Wright. Beating the random assignment on constraint satisfaction problems of bounded degree. *arXiv:1505.03424*, 2015. URL <https://arxiv.org/abs/1505.03424>.
- [2] Jacob Biamonte, Peter Wittek, Nicola Pancotti, Patrick Rebentrost, Nathan Wiebe, and Seth Lloyd. Quantum machine learning. *Nature*, 549(7671):195, 2017. doi:[10.1038/nature23474](https://doi.org/10.1038/nature23474).
- [3] L Chakhmakhchyan, NJ Cerf, and R Garcia-Patron. Quantum-inspired algorithm for estimating the permanent of positive semidefinite matrices. *Physical Review A*, 96(2):022329, 2017. doi:[10.1103/PhysRevA.96.022329](https://doi.org/10.1103/PhysRevA.96.022329).
- [4] Nai-Hui Chia, Han-Hsuan Lin, and Chunhao Wang. Quantum-inspired sublinear classical algorithms for solving low-rank linear systems. *arXiv:1811.04852*, 2018. URL <https://arxiv.org/abs/1811.04852>.
- [5] Nai-Hui Chia, Tongyang Li, Han-Hsuan Lin, and Chunhao Wang. Quantum-inspired classical sublinear-time algorithm for solving low-rank semidefinite programming via sampling approaches. *arXiv:1901.03254*, 2019. URL <https://arxiv.org/abs/1901.03254>.
- [6] AV Abs da Cruz, Marley Maria Bernardes Rebuzzi Vellasco, and Marco Aurélio Cavalcanti

- Pacheco. Quantum-inspired evolutionary algorithm for numerical optimization. In *Hybrid evolutionary algorithms*, pages 19–37. Springer, 2007. doi:[10.1007/978-3-540-73297-6_2](https://doi.org/10.1007/978-3-540-73297-6_2).
- [7] Yogesh Dahiya, Dimitris Konomis, and David P Woodruff. An empirical evaluation of sketching for numerical linear algebra. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1292–1300. ACM, 2018. doi:[10.1145/3219819.3220098](https://doi.org/10.1145/3219819.3220098).
 - [8] Petros Drineas, Ravi Kannan, and Michael W Mahoney. Fast monte carlo algorithms for matrices ii: Computing a low-rank approximation to a matrix. *SIAM Journal on computing*, 36(1):158–183, 2006. doi:[10.1137/S0095439704442696](https://doi.org/10.1137/S0095439704442696).
 - [9] Alan Frieze, Ravi Kannan, and Santosh Vempala. Fast monte-carlo algorithms for finding low-rank approximations. *Journal of the ACM (JACM)*, 51(6):1025–1041, 2004. doi:[10.1145/1039488.1039494](https://doi.org/10.1145/1039488.1039494).
 - [10] András Gilyén, Seth Lloyd, and Ewin Tang. Quantum-inspired low-rank stochastic regression with logarithmic dependence on the dimension. *arXiv:1811.04909*, 2018. URL <https://arxiv.org/abs/1811.04909>.
 - [11] Nathan Halko, Per-Gunnar Martinsson, and Joel A Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM review*, 53(2):217–288, 2011. doi:[10.1137/090771806](https://doi.org/10.1137/090771806).
 - [12] Kuk-Hyun Han and Jong-Hwan Kim. Quantum-inspired evolutionary algorithm for a class of combinatorial optimization. *IEEE transactions on evolutionary computation*, 6(6):580–593, 2002. doi:[10.1109/TEVC.2002.804320](https://doi.org/10.1109/TEVC.2002.804320).
 - [13] F Maxwell Harper and Joseph A Konstan. The movieLens datasets: History and context. *ACM Transactions on Interactive Intelligent Systems (TIIS)*, 5(4):19, 2016. doi:[10.1145/2827872](https://doi.org/10.1145/2827872).
 - [14] Aram W Harrow, Avinatan Hassidim, and Seth Lloyd. Quantum algorithm for linear systems of equations. *Physical Review Letters*, 103(15):150502, 2009. doi:[10.1103/PhysRevLett.103.150502](https://doi.org/10.1103/PhysRevLett.103.150502).
 - [15] Russell Impagliazzo and Ramamohan Paturi. On the complexity of k-SAT. *Journal of Computer and System Sciences*, 62(2):367–375, 2001. doi:[10.1006/jcss.2000.1727](https://doi.org/10.1006/jcss.2000.1727).
 - [16] Iordanis Kerenidis and Anupam Prakash. Quantum recommendation systems. *arXiv:1603.08675*, 2016. URL <https://arxiv.org/abs/1603.08675>.
 - [17] Seth Lloyd, Masoud Mohseni, and Patrick Rebentrost. Quantum principal component analysis. *Nature Physics*, 10(9):631, 2014. doi:[10.1038/nphys3029](https://doi.org/10.1038/nphys3029).
 - [18] Harry Markowitz. Portfolio selection. *The Journal of Finance*, 7(1):77–91, 1952.
 - [19] Per-Gunnar Martinsson, Vladimir Rokhlin, and Mark Tygert. A randomized algorithm for the decomposition of matrices. *Applied and Computational Harmonic Analysis*, 30(1):47–68, 2011. doi:[10.1016/j.acha.2007.12.002](https://doi.org/10.1016/j.acha.2007.12.002).
 - [20] Ajit Narayanan and Mark Moore. Quantum-inspired genetic algorithms. In *Evolutionary Computation, 1996., Proceedings of IEEE International Conference on*, pages 61–66. IEEE, 1996. doi:[10.1109/ICEC.1996.542334](https://doi.org/10.1109/ICEC.1996.542334).
 - [21] Cam Nugent. S and P 500 stock data. 2018. URL <https://www.kaggle.com/camnugent/sandp500>.
 - [22] Patrick Rebentrost and Seth Lloyd. Quantum computational finance: quantum algorithm for portfolio optimization. *arXiv:1811.03975*, 2018. URL <https://arxiv.org/abs/1811.03975>.
 - [23] Patrick Rebentrost, Masoud Mohseni, and Seth Lloyd. Quantum support vector machine for big data classification. *Physical Review Letters*, 113(13):130503, 2014. doi:[10.1103/PhysRevLett.113.130503](https://doi.org/10.1103/PhysRevLett.113.130503).
 - [24] Vladimir Rokhlin and Mark Tygert. A fast randomized algorithm for overdetermined linear least-squares regression. *Proceedings of the National Academy of Sciences*, 105(36):13212–13217, 2008. doi:[10.1073/pnas.0804869105](https://doi.org/10.1073/pnas.0804869105).
 - [25] Vladimir Rokhlin, Arthur Szlam, and Mark Tygert. A randomized algorithm for principal component analysis. *SIAM Journal on Matrix Analysis and Applications*, 31(3):1100–1124, 2009. doi:[10.1137/080736417](https://doi.org/10.1137/080736417).
 - [26] Troels F Rønnow, Zhihui Wang, Joshua Job, Sergio Boixo, Sergei V Isakov, David Wecker, John M Martinis, Daniel A Lidar, and Matthias Troyer. Defining and detecting quantum speedup. *Science*, 345(6195):420–424, 2014. doi:[10.1126/science.1252319](https://doi.org/10.1126/science.1252319).

- [27] Tamas Sarlos. Improved approximation algorithms for large matrices via random projections. In *Foundations of Computer Science, 2006. FOCS'06. 47th Annual IEEE Symposium on*, pages 143–152. IEEE, 2006. doi:[10.1109/FOCS.2006.37](https://doi.org/10.1109/FOCS.2006.37).
- [28] Jianhong Shen. On the singular values of Gaussian random matrices. *Linear Algebra and its Applications*, 326(1):1 – 14, 2001. ISSN 0024-3795. doi:[10.1016/S0024-3795\(00\)00322-0](https://doi.org/10.1016/S0024-3795(00)00322-0).
- [29] Allan Sly and Nike Sun. The computational hardness of counting in two-spin models on d-regular graphs. In *Foundations of Computer Science (FOCS), 2012 IEEE 53rd Annual Symposium on*, pages 361–369. IEEE, 2012. doi:[10.1109/FOCS.2012.56](https://doi.org/10.1109/FOCS.2012.56).
- [30] Ewin Tang. Quantum-inspired classical algorithms for principal component analysis and supervised clustering. *arXiv:1811.00414*, 2018. URL <https://arxiv.org/abs/1811.00414>.
- [31] Ewin Tang. A quantum-inspired classical algorithm for recommendation systems. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pages 217–228, 2019. doi:[10.1145/3313276.3316310](https://doi.org/10.1145/3313276.3316310).
- [32] Leslie G Valiant. Holographic algorithms. *SIAM Journal on Computing*, 37(5):1565–1594, 2008. doi:[10.1137/070682575](https://doi.org/10.1137/070682575).
- [33] David P Woodruff et al. Sketching as a tool for numerical linear algebra. *Foundations and Trends® in Theoretical Computer Science*, 10(1–2):1–157, 2014. doi:[10.1561/0400000060](https://doi.org/10.1561/0400000060).
- [34] Franco Woolfe, Edo Liberty, Vladimir Rokhlin, and Mark Tygert. A fast randomized algorithm for the approximation of matrices. *Applied and Computational Harmonic Analysis*, 25(3):335–366, 2008. doi:[10.1016/j.acha.2007.12.002](https://doi.org/10.1016/j.acha.2007.12.002).