

000 QUANTUM ALGORITHMS FOR PROJECTION-FREE 001 002 SPARSE CONVEX OPTIMIZATION 003 004

005 **Anonymous authors**

006 Paper under double-blind review

007 008 ABSTRACT 009

011 This paper considers the projection-free sparse convex optimization problem for
012 the vector domain and the matrix domain, which covers a large number of im-
013 portant applications in machine learning and data science. For the vector domain
014 $\mathcal{D} \subset \mathbb{R}^d$, we propose two quantum algorithms for sparse constraints that finds a
015 ε -optimal solution with the query complexity of $O(\sqrt{d}/\varepsilon)$ and $O(1/\varepsilon)$ by using
016 the function value oracle, reducing a factor of $O(\sqrt{d})$ and $O(d)$ over the best clas-
017 sical algorithm, respectively, where d is the dimension. For the matrix domain
018 $\mathcal{D} \subset \mathbb{R}^{d \times d}$, we propose two quantum algorithms for nuclear norm constraints that
019 improve the time complexity to $\tilde{O}(rd/\varepsilon^2)$ and $\tilde{O}(\sqrt{rd}/\varepsilon^3)$ for computing the up-
020 date step, reducing at least a factor of $O(\sqrt{d})$ over the best classical algorithm,
021 where r is the rank of the gradient matrix. Our algorithms show quantum advan-
022 tages in projection-free sparse convex optimization problems as they outperform
023 the optimal classical methods in dependence on the dimension d .

024 025 1 INTRODUCTION

027 In this paper, we consider the following *constrained* optimization problem of the form

$$029 \min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x}), \quad (1)$$

031 such objective covers many important application in operations research and machine learning. We
032 are interested in the case where 1) the objective function f is convex and continuously differentiable,
033 and 2) the domain $\mathcal{D} \subset \mathbb{R}^d$ is a feasible set that is convex, and the dimension d is high. Typical in-
034 stances of such high-dimensional optimization problems include multiclass classification, multitask
035 learning, matrix learning, network systems and many more Garber & Hazan (2016); Hazan & Kale
036 (2012); Hazan et al. (2012); Jaggi (2013); Dudik et al. (2012); Zhang et al. (2012); Harchaoui et al.
037 (2015); Hazan & Luo (2016). As an example, for matrix completion, the optimization problem is:

$$038 \min_{X \in \mathbb{R}^{m \times n}, \|X\|_{\text{tr}} \leq k} \sum_{(i,j) \in \Omega} (X_{i,j} - Y_{i,j})^2, \quad (2)$$

040 where X is the matrix to be recovered, Ω denotes the observed elements, $Y_{i,j}$ is the observed known
041 value at position (i, j) , and $\|X\|_{\text{tr}} \leq k$ represents the trace norm (nuclear norm) constraint.

043 Compared with unconstrained convex optimization problems, optimizing Equation (1) involves han-
044 dling constraints, which introduces new challenges. A straightforward method for optimizing Equa-
045 tion (1) is the projected gradient descent approach Levitin & Polyak (1966). This method first takes
046 a step in the gradient direction and then performs the projection to satisfy the constraint. However, in
047 practice, the dimensions of the feasible set can be very large, leading to prohibitively high computa-
048 tional complexity. For example, when solving Equation (2), the projection step involves performing
049 a singular value decomposition (SVD), whose time complexity is $O(mn \min\{m, n\})$ ($O(d^3)$ for
050 $X \in \mathbb{R}^{d \times d}$). Compared to the projected gradient descent approach, the Frank-Wolfe (FW) method
051 (also known as the conditional gradient method) is more efficient when dealing with structured con-
052 strained optimization problems. Rather than performing projections, it solves a computationally
053 efficient linear sub-problem to ensure that the solution lies within the feasible set \mathcal{D} . When solving
Equation (2), the time complexity of the Frank-Wolfe method is $O(mn)$ ($O(d^2)$ for $X \in \mathbb{R}^{d \times d}$),
which is significantly lower than the complexity of SVD-based projections. Since the Frank-Wolfe

method is efficient for optimizing many difficult machine learning problems, such as low-rank constrained problems and **sparsity-inducing** constrained problems, it has attracted significant attention and has been applied to solving Equation (3) and many of its variants.

Despite the efficiency of FW in handling structured constraints, it still incurs significant computational overhead when dealing with high-dimensional problems. The bottleneck of the computation is the linear subproblem over \mathcal{D} , which is either assumed to have efficient implementation or simply follows existing classical oracles, such as Dunn & Harshbarger (1978); Jaggi (2013); Garber & Hazan (2016). The overhead of these oracles, however, grows linearly or superlinearly in terms of dimension d .

Recently, quantum computing has emerged as a promising new paradigm to accelerate a large number of important optimization problems (see Appendix A.4). We aim to take a thorough investigation on whether quantum computing can accelerate FW algorithms, in particular the linear subproblem over structured constraints regarding dimension d . We aim to answer the following question:

Can one utilize quantum techniques to accelerate Frank-Wolfe algorithms in terms of dimension d ?

Chen et al. gave an initial answer to this question Chen & de Wolf (2023). They considered the linear regression problem with explicit functional form where the closed form of gradient is provided. Given the precomputed matrix factors of the closed-form objective function stored in specific data structures, they leveraged HHL-based algorithms to accelerate matrix multiplications in calculating the closed-form gradient, leading to a upper bound of $O\left(\sqrt{d}/\varepsilon^2\right)$. In this work, we consider a more general problem where the objective function is a smooth convex function accessible only through a function value oracle, and then we consider a more general constraint conditions (*the latent group norm ball*) to enhance the theoretical framework’s applicability. Besides, we also consider the case of matrix feasible set, under different assumptions. To our best knowledge, we are the first one to consider accelerating the matrix case of the FW algorithm by quantum computing.

Contributions. We give a systematic study on how to accelerate FW algorithms when \mathcal{D} is either a vector domain \mathbb{R}^d , or a matrix domain $\mathbb{R}^{d \times d}$ subject to various structured constraints. Note that our findings can be applied to non-square matrices, we express our results using square matrices for simplicity of presentation (Remark 1). We summarize our contributions as follows.

For the vector domain $\mathcal{D} \subset \mathbb{R}^d$:

- We propose the quantum Frank-Wolfe algorithm for the projection-free sparse convex optimization problem under ℓ_1 norm constraints (Theorem 1) and the d -dimensional simplex Δ_d (Theorem 2). We achieve a query complexity of $\tilde{O}(\sqrt{d}/\varepsilon)$ in finding an ε -optimal solution using the function value oracle, reducing a factor of $O(\sqrt{d})$ over the optimal classical algorithm. Furthermore, if the objective function is a Lipschitz continuous function, we prove that the query complexity can be reduced to $O(1/\varepsilon)$ by employing the bounded-error Jordan quantum gradient estimation algorithm, at the cost of more qubits and additional gates (Theorem 5). In addition, we consider the generalization to latent group norm constraints (Theorem 6) and achieve a query complexity of $\tilde{O}\left(\sqrt{|\mathcal{G}|}|\mathbf{g}|_{\max}\right)$, representing an $O\left(\sqrt{|\mathcal{G}|}\right)$ speedup over the classical algorithm. These results are presented in Section 3, Appendix A.1 and A.2. The comparison with the classical methods is shown in Table 1.
- Specifically, we develop a novel quantum subroutine for the Frank-Wolfe linear subproblem over latent group norm constraints, by computing dual norms coherently across all groups in quantum superposition and identifying the dominant group via quantum maximum finding. We establish a novel error propagation analysis for dual norm computation under gradient approximation, deriving bounds via Hölder’s inequality that enable precise control of linear subproblem accuracy throughout Frank-Wolfe iterations. The examples in the main text such as the ℓ_1 -norm constrained are special instances of the latent group constraints. In short, we develop quantum subroutines for dominant atom finding and show that the errors can be controlled by setting appropriate parameters.

For the matrix domain $\mathcal{D} \subset \mathbb{R}^{d \times d}$:

- For the projection-free sparse problem under nuclear norm constraints, we propose two complementary quantum Frank-Wolfe algorithms tailored to high-rank and low-rank gradient matrices, respectively (see Appendix A.5). For finding an ε -optimal solution, we achieve a time complexity of $\tilde{O}(rd/\varepsilon^2)$ (Theorem 3) and $\tilde{O}(\sqrt{rd}/\varepsilon^3)$ (Theorem 4) in computing the update direction, representing an at least $O(\sqrt{d})$ speedup over state-of-the-art classical algorithm, where r is the rank of the gradient matrix. These results are presented in Section 4 and the comparison with the classical methods is shown in Table 2.
- Specifically, in the first algorithm, we simplify the top- k singular vectors extraction method Bellante et al. (2022) by utilizing the quantum maximum finding algorithm, which avoids the overheads of repeated sampling to estimate the factor score ratio, and avoids the overheads of searching the threshold value. In the second algorithm, we introduce the quantum power method to extract the top singular vectors, which reduces the dependence on the rank of the gradient matrix, at the cost of higher sensitivity on solution precision.

Wide range of critical applications can be benefited from the acceleration of QFW, including sparse regression (Lasso), sparse signal recovery, matrix completion, boosting algorithms (e.g., AdaBoost), Support Vector Machines, and density estimation Jaggi (2013). Other applications include signal processing (sparsity constraints via ℓ_1 norm), game theory (zero sum games with simplex) and SDPs (nuclear norm optimization). The proposed top singular vectors extraction techniques also have a potential application for bi-quadratic programming Li et al. (2024). We discuss some of these applications in Appendix A.6

We notice an independent work on the quantum power method Chen et al. (2025a), whose second algorithm shares a conceptual similarity with our second approach: both iteratively apply quantum matrix-vector multiplication. They assume a sparse-query access to the matrix as input, and achieve a complexity of $\tilde{O}((d\sqrt{s}/\gamma\varepsilon)^{1+o(1)})$, where s is the sparsity, γ is the eigenvalue gap, whereas our method relies on the rank of the matrix, instead of the sparsity. In the case of dense full-rank matrix, their algorithm and ours are consistent in terms of dimensional dependence, which provides mutual validation of correctness.

The remainder of this paper is organized as follows. Section 2 introduces the basic concept of constrained optimization and the classical Frank-Wolfe algorithm. Appendix A.3 introduces the notations and assumptions of quantum computing. Section 3 and 4 presents our quantum FW methods for vector domain and matrix domain, respectively. Extension for the vector cases are presented in Appendix A.1 and A.2. Extended related works are presented in Appendix A.4, and we conclude with a discussion about the future work in Section 5. Proof details are given in Appendix B.

Table 1: Classical algorithms V.S. quantum algorithms of the vector case, where C_f is the curvature of the objective function f , ε is the precision of the solution, d is the dimension of the domain, G is the Lipschitz parameter of the objective function, p is the failure probability.

| Optimization Domain | Constraints | Algorithm | Iteration | Query complexity | Qubits | Gates |
|-----------------------------|---------------------------------|--------------------------|----------------------|--|--|--|
| Sparse Vectors | $\ \cdot\ _1$ -ball | FW Jaggi (2013) | $O(C_f/\varepsilon)$ | $O(d)$ | $O\left(d + \log\frac{1}{\varepsilon}\right)$ | $O(\sqrt{d})$ |
| | | QFW (Theorem 1) | $O(C_f/\varepsilon)$ | $O(\sqrt{d} \log(C_f/\varepsilon))$ | | |
| | | QFW (Theorem 5) | $O(1)$ | $O(1)$ | | |
| Sparse non-neg. vectors | Simplex Δ_d | Frank-Wolfe Jaggi (2013) | $O(C_f/\varepsilon)$ | $O(d)$ | $O\left(d + \log\frac{1}{\varepsilon}\right)$ | $O(\sqrt{d})$ |
| | | QFW (Theorem 2) | $O(C_f/\varepsilon)$ | $O(\sqrt{d} \log(C_f/\varepsilon))$ | | |
| | | QFW (Theorem 5) | $O(1)$ | $O(1)$ | | |
| Latent group sparse vectors | $\ \cdot\ _{\mathcal{G}}$ -ball | FW Jaggi (2013) | $O(C_f/\varepsilon)$ | $O(\sum_{g \in \mathcal{G}} g)$ | $O(d + \log \mathcal{G} + \mathcal{G} _{\max} \log(1/\varepsilon))$ | $\tilde{O}(\sqrt{ \mathcal{G} } \cdot \mathcal{G} _{\max})$ |
| | | QFW (Theorem 6) | $O(C_f/\varepsilon)$ | $O(\sqrt{ \mathcal{G} } \mathcal{G} _{\max} \log(C_f/\varepsilon))$ | | |

2 PRELIMINARIES

2.1 NOTATIONS AND ASSUMPTIONS FOR CONSTRAINED OPTIMIZATION PROBLEM

We consider constrained convex optimization problems of the form

$$\min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x}), \quad (3)$$

where $\mathbf{x} \in \mathbb{R}^d$, $f : \mathbb{R}^d \rightarrow \mathbb{R}$, and $\mathcal{D} \subseteq \mathbb{R}^d$ is the constraint set. In addition, as usually in constrained convex optimization, we also make the following assumptions:

162 Table 2: Classical algorithms V.S. quantum algorithms of the matrix case, where C_f is the curvature
 163 of the objective function f , ε is the precision of the solution, d is the dimension of the domain, T_∇
 164 is the times required to evaluate ∇f ; $\sigma_1(M)$ and $\sigma_2(M)$ are the largest and the second largest singular value,
 165 respectively; r is the rank of the gradient matrix; γ'_{\min} is a factor which depends on the relation of
 166 the singular value distribution of the gradient matrix and the direction of the initial vector.

| Domain | Constraints | Algorithm | Iteration | Complexity of the Update Computing |
|------------------------|------------------------|-------------------------------------|----------------------|--|
| Sparse Matrices 178 | $\ \cdot\ _{tr}$ -ball | FW with Power Method Jaggi (2013) | $O(C_f/\varepsilon)$ | $O\left(\frac{\sigma_1(M)d^2}{(\sigma_1(M)-\sigma_2(M))\varepsilon} + T_\nabla\right)$ |
| | | FW with Lanczos Method Jaggi (2013) | $O(C_f/\varepsilon)$ | $O\left(\frac{\sqrt{\sigma_1(M)d^2}}{\sqrt{(\sigma_1(M)-\sigma_2(M))\varepsilon}} + T_\nabla\right)$ |
| | | FW with QTSVE (Theorem 3) | $O(C_f/\varepsilon)$ | $\tilde{O}\left(\frac{r\sigma_1^3(M)d}{(\sigma_1(M)-\sigma_2(M))\varepsilon^2} + T_\nabla\right)$ |
| | | FW with QPM (Theorem 4) | $O(C_f/\varepsilon)$ | $\tilde{O}\left(\frac{\sqrt{r}\sigma_1^3(M)d}{(1-\sigma_1(M))\gamma'^3_{\min}\varepsilon^3} + T_\nabla\right)$ |

Algorithm 1 Classical Frank-Wolfe Algorithm with Approximate Linear Subproblems

1: **Input:** Solution precision ε , iterations T .
 2: **Output:** $\mathbf{x}^{(T)}$ such that $f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \varepsilon$.
 3: **Initialize:** Let $\mathbf{x}^{(1)} \in \mathcal{D}$.
 4: **for** $t = 1, \dots, T$ **do**
 5: Let $\gamma_t = \frac{2}{t+2}$.
 6: Find direction $\mathbf{s} \in \mathcal{D}$ such that
 182
$$\langle \mathbf{s}, \nabla f(\mathbf{x}^{(t)}) \rangle \leq \min_{\hat{\mathbf{s}} \in \mathcal{D}} \langle \hat{\mathbf{s}}, \nabla f(\mathbf{x}^{(t)}) \rangle + \frac{\delta}{2} \gamma_t C_f. \quad (5)$$

 184 7: Update $\mathbf{x}^{(t+1)} = (1 - \gamma_t) \mathbf{x}^{(t)} + \gamma_t \mathbf{s}$.
 185 8: **end for**

187
 188 **Assumption 1.** f is convex and L -smooth, i.e., the gradient of f satisfies $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq$
 189 $L\|\mathbf{x} - \mathbf{y}\|_2$ for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.

190 **Assumption 2.** \mathcal{D} is compact and convex, and the diameter of \mathcal{D} has an upper bound D , i.e.,
 191 $\forall \mathbf{x}, \mathbf{y} \in \mathcal{K}, \|\mathbf{x} - \mathbf{y}\|_2 \leq D$.

192 Typically, solving $\operatorname{argmin}_{\mathbf{x} \in \mathcal{D}} \mathbf{x}^\top \mathbf{y}$ for any $\mathbf{y} \in \mathbb{R}^d$, is much faster than the projection operation
 193 onto \mathcal{D} (i.e., solving $\operatorname{argmin}_{\mathbf{x} \in \mathcal{D}} \|\mathbf{x} - \mathbf{y}\|$). Examples of such domains include the set of sparse
 194 vectors, bounded norm matrices, flow polytope and many more Hazan & Kale (2012). Therefore,
 195 for such domains, the basic idea of the Frank-Wolfe algorithm is to replace the projection operation
 196 with a linear optimization problem.

197 In the design and analysis of the Frank-Wolfe algorithm, one important quantity is the curvature C_f ,
 198 which measures the “non-linearity” of f and is defined as follows,

$$C_f = \sup_{\mathbf{x}, \mathbf{s} \in \mathcal{D}, \beta \in [0, 1], \mathbf{y} = \mathbf{x} + \beta(\mathbf{s} - \mathbf{x})} \frac{2}{\beta^2} \times (f(\mathbf{y}) - f(\mathbf{x}) - \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle). \quad (4)$$

203 By Lemma 7 of Jaggi (2013), the curvature can be bounded as $C_f \leq LD^2$.

205 **2.2 CLASSICAL FRANK-WOLFE ALGORITHM**

207 The classical Frank-Wolfe algorithm is given in Algorithm 1. The key step is the linear subproblem
 208 of Equation (5) which seeks an approximate minimizer in \mathcal{D} of $\langle \mathbf{s}, \nabla f(\mathbf{x}^{(t)}) \rangle$. Classically, the per-
 209 step cost is $O(N)$ where N is the number of elements that need to be searched which introduces a
 210 large $O(N)$ cost. In this work, we will show that $O(\sqrt{N})$ quantum queries to solve this subproblem.

211 **Lemma 1.** [Jaggi (2013), Theorem 1] For each $t \geq 1$, the iterates of Algorithm 1 satisfy

$$f(\mathbf{x}^{(t)}) - f(\mathbf{x}^*) \leq \frac{2C_f}{t+2}(1 + \delta), \quad (6)$$

214 where \mathbf{x}^* is the optimal solution to Equation (3), and δ is the solution quality to which the internal
 215 linear subproblems are solved. That is, one can use $O(\frac{(1+\delta)C_f}{\varepsilon})$ iterations to have a ε -opt solution.

216 **Algorithm 2** Quantum Frank-Wolfe Algorithm for Sparsity/Simplex Constraint

217 1: **Input:** Solution precision ε , gradient precision $\{\sigma_t\}_{t=1}^T$.

218 2: **Output:** $\mathbf{x}^{(T)}$ such that $f(\mathbf{x}^T) - f(\mathbf{x}^*) \leq \varepsilon$.

219 3: **Initialize:** Let $\mathbf{x}^{(1)} \in \mathcal{D}$.

220 4: Let $T = \frac{4C_f}{\varepsilon} - 2$.

221 5: **for** $t = 1, \dots, T$ **do**

222 6: Let $\gamma_t = \frac{2}{t+2}$.

223 7: Prepare quantum state $\sum_{i=0}^{d-1} |i\rangle |\mathbf{x}^{(t)}\rangle |0\rangle$.

224 8: Perform quantum gradient circuit (Lemma 3) to get $\sum_{i=0}^{d-1} |i\rangle |\mathbf{x}^{(t)}\rangle \left| \frac{f(\mathbf{x}^{(t)} + \sigma_t \mathbf{e}_i) - f(\mathbf{x}^{(t)})}{\sigma_t} \right\rangle$.

225 9: Apply quantum maximum finding to the absolute value of the third register (to the third

226 register directly for the simplex constraint, respectively) (Lemma 4), and then measure the first

227 register to obtain measurement result i_t .

228 10: Set $\mathbf{s} = -\mathbf{e}_{i_t}$. Update $\mathbf{x}^{(t+1)} = (1 - \gamma_t)\mathbf{x}^{(t)} + \gamma_t \mathbf{s}$.

229 11: **end for**

233 **3 QUANTUM FRANK-WOLFE ALGORITHMS OVER VECTORS**

235 **3.1 QUANTUM FRANK-WOLFE WITH SPARSITY CONSTRAINTS**

237 We first consider the optimization problem

239
$$\min f(\mathbf{x}), \text{ s.t. } \mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\| \leq 1, \quad (7)$$

241 where the sparsity constraint $\mathcal{D} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_1 \leq 1\}$.

242 By Section 4 of Jaggi (2013), any linear function attains its minimum over a convex hull at a vertex.
243 Thus, for the ℓ_1 norm problem, the exact minimizer (i.e, corresponding to $\delta = 0$) of Equation (5) is
244 $\hat{\mathbf{s}} = -\mathbf{e}_{i_t}$ with

246
$$i_t \in \operatorname{argmax}_{i \in [d]} |\nabla_i f(\mathbf{x}^{(t)})|, \quad (8)$$

248 i.e., it is a coordinate corresponding to the largest absolute value of the gradient component.

250 Our approach will be to construct an approximate quantum maximum gradient component finding
251 algorithm to find such an i_t .

252 **Quantum access model U_f .** In this subsection, we assume that the value of the loss function is
253 accessed via a function value oracle as shown in Assumption 3.

254 **Assumption 3.** *There is a unitary U_f that, in time T_f , returns the function value, i.e., $U_f : |x\rangle |a\rangle \rightarrow |x\rangle |a + f(x)\rangle$, for any a , where $|x\rangle := |x_1\rangle |x_2\rangle \dots |x_d\rangle$.*

257 The preparation of the input state in Step 7 of Algorithm 2 is efficient. Initialize the algorithm at
258 $\mathbf{x}^{(0)} = 0$, each Frank-Wolfe step adds a single coordinate direction to the solution. Specifically, the
259 update rule $\mathbf{x}^{(t+1)} = (1 - \gamma_t)\mathbf{x}^{(t)} + \gamma_t \mathbf{s}_t$ —where \mathbf{s}_t is a standard basis vector—implies that the
260 solution $\mathbf{x}^{(t)}$ after t iterations is a sparse vector with at most t non-zero components. Consequently, the
261 quantum state $|\mathbf{x}^{(t)}\rangle$ is a sparse computational basis state. This state can be perform an incremental
262 update, setting at most one new coordinate to a non-zero value per iteration. The gate complexity
263 for this sparse update is $O(t)$. The total number of iterations T required for an ε -optimal solution
264 is $O(1/\varepsilon)$, which is independent of the dimension d . Therefore, the state preparation overhead per
265 iteration remains $O(1/\varepsilon)$, completely decoupled from the potentially large dimension d .

266 **Quantum gradient circuit.** Next, we present a general unitary U_g to approximate the gradient
267 $\nabla f(\mathbf{x}_t)$. Specifically, we use the forward difference $g_i(\mathbf{x}_t) = \frac{f(\mathbf{x}_t + \sigma \mathbf{e}_i) - f(\mathbf{x}_t)}{\sigma}$ to approximate each
268 item of $\nabla_i f(\mathbf{x}_t)$ with ℓ_∞ error ε_g , i.e., $\|\nabla f(\mathbf{x}_t) - g(\mathbf{x}_t)\|_\infty \leq \varepsilon_g$, where σ is the tunable parameter
269 for the desired accuracy.

270 **Lemma 2** (Theorem 3.1 Berahas et al. (2022)). *Under Assumption 1, let $g_i(\mathbf{x}) = \frac{f(\mathbf{x} + \sigma \mathbf{e}_i) - f(\mathbf{x})}{\sigma}$,
271 then for all $\mathbf{x} \in \mathbb{R}^d$,*

$$273 \quad \|g(\mathbf{x}) - \nabla f(\mathbf{x})\|_2 \leq \frac{\sqrt{d}L\sigma}{2}. \quad (9)$$

275 **Lemma 3.** *Given access to the quantum function value oracle \mathbf{U}_f , there exists a quantum circuit
276 to construct a quantum error bounded gradient oracle $\mathbf{U}_g : |i\rangle |\mathbf{x}\rangle |0\rangle \rightarrow |i\rangle |\mathbf{x}\rangle |g_i(\mathbf{x})\rangle$, where
277 $g_i(\mathbf{x}) = \frac{f(\mathbf{x} + \sigma \mathbf{e}_i) - f(\mathbf{x})}{\sigma}$ is the i -th component of the gradient and σ is the tunable parameter, with
278 two queries to the quantum function value oracle.*

280 The proof is given in Appendix B.1.

281 **Quantum maximum finding circuit.** Based on \mathbf{U}_g , leveraging the quantum minimum-finding al-
282 gorithm Durr & Hoyer (1996), we give an approximate search of the maximum gradient component
283 as shown in Lemma 4, with proof given in Appendix B.2. Note that Algorithm 3 in the matrix sec-
284 tion of this work also utilizes quantum maximum finding, but with a non-uniform input state. We
285 also provide a proof in Appendix B.2 that the quantum maximum finding procedure is applicable to
286 non-uniform input states.

287 **Lemma 4.** *(Approximate maximum gradient component finding) Given access to the quantum
288 error bounded gradient oracle $\mathbf{U}_g : |i\rangle |\mathbf{x}\rangle |0\rangle \rightarrow |i\rangle |\mathbf{x}\rangle |g_i(\mathbf{x})\rangle$ s.t. for each $i \in [d]$, after
289 measuring $|g_i(\mathbf{x})\rangle$, the measured outcome $g_i(\mathbf{x})$ satisfies $|g_i(\mathbf{x}) - \nabla f_i(\mathbf{x})| \leq \epsilon$. There exists a
290 quantum circuit \mathcal{A}_{\max} that finds the index i^* that satisfies $\nabla f_{i^*}(\mathbf{x}) \geq \max_{j \in [d]} \nabla f_j(\mathbf{x}) - 2\epsilon$ or
291 $|\nabla f_{i^*}(\mathbf{x})| \geq \max_{j \in [d]} |\nabla f_j(\mathbf{x})| - 2\epsilon$, using $O(\sqrt{d} \log(\frac{1}{\delta}))$ applications of \mathbf{U}_g , \mathbf{U}_g^\dagger and $O(\sqrt{d})$
292 elementary gates, with probability $1 - \delta$. For the non-uniform initial state, let p be the initial mea-
293 surement probability of the maximum component, then the algorithm finds the maximum with query
294 complexity of $O(\frac{1}{\sqrt{p}} \log(\frac{1}{\delta}))$.*

295 **Convergence Analysis.** Now we can conduct the convergence analysis with the help of approxi-
296 mate maximum finding sub-routine and show how to choose appropriate parameters, which gives
297 Theorem 1, with proof given in Appendix B.3.

298 **Theorem 1.** *(Quantum FW over the sparsity constraint) By setting $\sigma_t = \frac{C_f}{\sqrt{d}L(t+2)}$ for $t \in [T]$,
299 the quantum algorithm (Algorithm 2) solves the sparsity constraint optimization problem for any
300 precision ϵ such that $f(\mathbf{x}^T) - f(\mathbf{x}^*) \leq \epsilon$ in $T = \frac{4C_f}{\epsilon} - 2$ rounds, succeed with probability $1 - p$,
301 with $O\left(\sqrt{d} \log \frac{C_f}{p\epsilon}\right)$ calls to the function value oracle \mathbf{U}_f per round.*

302 If the objective function is a G -Lipschitz continuous function (i.e. $|f(\mathbf{x}) - f(\mathbf{y})| \leq G\|\mathbf{y} - \mathbf{x}\|$,
303 $\forall \mathbf{x}, \mathbf{y} \in \mathcal{D}$), an alternative approach for estimating the gradient of the objective function
304 involves employing the bounded-error Jordan algorithm to improve the query complexity of each
305 iteration to $O(1)$, at the cost of additional space complexity and extra gate operations. This result is
306 given in Appendix A.1.

310 3.2 EXTENSIONS: QUANTUM FRANK-WOLFE FOR ATOMIC SETS

311 Classically, the Frank-Wolfe algorithm has been shown to be well-suited to atomic sets Jaggi (2013),
312 i.e. where the constraint set is expressed as the convex hull of another (not-necessarily finite) set \mathcal{A} :
313 $\mathcal{D} = \text{conv}(\mathcal{A})$ In this case, the Frank-Wolfe update calculation requires a minimization only over
314 \mathcal{A} : $\min_{\hat{s} \in \mathcal{A}} \langle \hat{s}, \nabla f(x^{(t)}) \rangle$. The optimization over the ℓ_1 ball as studied above is a special case of
315 this, since $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_1\} = \text{conv}\{\pm \mathbf{e}_1, \pm \mathbf{e}_2, \dots, \pm \mathbf{e}_d\}$. Note also that quantum optimization
316 over the simplex $\Delta_d = \text{conv}\{\mathbf{e}_1, \dots, \mathbf{e}_d\}$ can be done by almost exactly the same method as for
317 the ℓ_1 case, with the only modification to account for the fact that only the unit vectors need to be
318 optimized over, which gives Theorem 2.

319 **Theorem 2.** *(Quantum FW over the simplex) By setting $\sigma_t = \frac{C_f}{\sqrt{d}L(t+2)}$ for $t \in [T]$, the quantum
320 algorithm (Algorithm 2) solves the simplex constraint optimization problem for any precision ϵ such
321 that $f(\mathbf{x}^T) - f(\mathbf{x}^*) \leq \epsilon$ in $T = \frac{4C_f}{\epsilon} - 2$ rounds, succeed with probability $1 - p$, with $O\left(\sqrt{d} \log \frac{C_f}{p\epsilon}\right)$
322 calls to the function value oracle \mathbf{U}_f per round.*

324 Two more extensions for atomic sets are given in Appendix A.2.
 325
 326

327 4 QUANTUM FRANK-WOLFE ALGORITHMS OVER MATRICES 328

329 In this section, we consider the matrix version of the constrained optimization problem in Equa-
 330 tion (1), specifically,
 331

$$332 \min f(X), \text{ s.t. } X \in \mathbb{R}^{d \times d}, \|X\|_{\text{tr}} \leq 1, \quad (10)$$

334 where the sparsity constraint is $\mathcal{D} = \{X \in \mathbb{R}^{d \times d} : \|X\|_{\text{tr}} \leq 1\}$. For simplicity of presentation, we
 335 first focus on square matrices, i.e., $X \in \mathbb{R}^{d \times d}$ (Remark 1).
 336

337 **Schatten matrix norm.** In contrast to the vector norm $\|\cdot\|$ on \mathbb{R}^d , the corresponding Schatten matrix
 338 norm $\|X\|$ is defined as $\|(\sigma_1, \dots, \sigma_d)\|$, where $\sigma_1, \dots, \sigma_d$ are singular values of X . It is known that
 339 the dual of the Schatten ℓ_p norm is the Schatten ℓ_q norm with $1/p + 1/q = 1$. The most prominent
 340 example is the trace norm $\|\cdot\|_{\text{tr}}$, also referred to as the nuclear norm or Schatten ℓ_1 norm, defined as
 341 the sum of the singular values $\|X\|_{\text{tr}} = \sum_{i=1}^d \sigma_i$.
 342

343 **Linear subproblem solver.** Following the classical Frank-Wolfe iteration framework, we aim to
 344 solve the linear optimization subproblem $\min_{S \in \mathcal{D}} \langle S, \nabla f(X_t) \rangle$ where X_t denotes the iterate matrix
 345 at step t , and $\langle X, Y \rangle = \text{tr } X^\top Y$ represents the Hilbert-Schmidt inner product. For convenience, let
 346 $M = \nabla f(X_t)$ in the rest of this section. To solve this subproblem, one can compute the singular
 347 value decomposition (SVD) $M = U \text{diag}(\sigma) V^\top$, where σ are singular values of M and $U, V \in$
 348 $\mathbb{R}^{d \times d}$ are orthogonal matrices. Since Schatten norms are invariant under orthogonal transformations,
 349 the optimal solution $S \in \mathcal{D}$ for the minimization problem $\min_{S \in \mathcal{D}} \langle S, M \rangle$ takes the forms of $S =$
 350 $U \text{diag}(s) V^\top$, where $\langle s, \sigma \rangle = \|\sigma\|_q$ with $\|s\|_p \leq 1$ and $1/p + 1/q = 1$. For the nuclear norm (i.e.,
 351 ℓ_1 Schatten norm), this reduces to $S = \mathbf{u} \mathbf{v}^\top$ where \mathbf{u}, \mathbf{v} are the left and right top singular vectors
 352 of M , corresponding to its largest singular value $\sigma_1(M)$. Thus, the core computational task is to
 353 efficiently approximate the top singular vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$, ensuring $|\mathbf{u}^\top M \mathbf{v} - \sigma_1(M)| \leq \varepsilon$.
 354

355 **Power method and Lanczos method.** Compared with the SVD that requires $O(d^3)$ computational
 356 cost per iteration to compute all d singular vectors, extracting only the top singular vector is much
 357 easier. Specifically, Kuczyński & Woźniakowski (1992) considers two iterative methods: the power
 358 method and the Lanczos method. The power method achieves $|\mathbf{u}^\top M \mathbf{v} - \sigma_1(M)| \leq \varepsilon'$ with the
 359 worst-case computation complexity of $O\left(\frac{\sigma_1(M)d^2 \ln d}{(\sigma_1(M) - \sigma_2(M))\varepsilon'}\right)$, while the Lanczos method achieves
 360 $|\mathbf{u}^\top M \mathbf{v} - \sigma_1(M)| \leq \varepsilon'$ with the worst-case computation complexity of $O\left(\frac{\sqrt{\sigma_1(M)}d^2 \ln d}{\sqrt{(\sigma_1(M) - \sigma_2(M))\varepsilon'}}\right)$,
 361 where ε' is the additive error. Similar to the convergence analysis in Section 3.1, setting $\varepsilon' =$
 362 $O(\varepsilon)$, the complexity of update computing are $O\left(\frac{\sigma_1(M)d^2 \ln d}{(\sigma_1(M) - \sigma_2(M))\varepsilon}\right)$ and $O\left(\frac{\sqrt{\sigma_1(M)}d^2 \ln d}{\sqrt{(\sigma_1(M) - \sigma_2(M))\varepsilon}}\right)$,
 363 respectively.
 364

365 **Quantum enhancement.** In the following, we propose two quantum subroutines to compute the top
 366 singular vector: the quantum top singular vector extraction method and the quantum power method.
 367 Note that for the matrix case, we could also assume the same function value oracle and naturally
 368 employ an improved Jordan’s algorithm to achieve a query complexity advantage in gradient esti-
 369 mation. However, in this work, we aimed to further investigate whether quantum algorithms can
 370 accelerate the computational complexity of the update step beyond just query counts. Therefore, the
 371 analysis focuses on the update direction computation and assumes that the gradient has been pre-
 372 computed and stored in the memory (Remark 3), following the classical convention of excluding
 373 gradient evaluation time Jaggi (2013).
 374

375 First, we assume the following gradient access model for matrix data. A detailed description of this
 376 data structure can be found in Section 1.A of Kerenidis & Prakash (2020b).
 377

378 **Assumption 4** (Quantum access to a matrix). *We assume that we have efficient quantum access to
 379 the matrix $M \in \mathbb{R}^{d \times d}$. That is, there exists a data structure that allows performing the mapping
 380 $|i\rangle |0\rangle \rightarrow |i\rangle |M_{i,\cdot}\rangle = |i\rangle \frac{1}{\|M_{i,\cdot}\|} \sum_j^d M_{ij} |j\rangle$ for all i , and $|0\rangle \rightarrow \frac{1}{\|M\|_F} \sum_i^d \|M_{i,\cdot}\| |i\rangle$ in time $\tilde{O}(1)$.*

378 4.1 QUANTUM FRANK-WOLFE WITH QUANTUM TOP SINGULAR VECTOR EXTRACTION
379380 Leveraging the quantum access defined in Assumption 4, quantum singular value estimation can be
381 performed efficiently.382 **Lemma 5** (Singular value estimation (Theorem 3 Bellante et al. (2022), Kerenidis & Prakash
383 (2020b)). *Let there be quantum access to $M \in \mathbb{R}^{d \times d}$, with singular value decomposition $M =$
384 $\sum_i^d \sigma_i \mathbf{u}_i \mathbf{v}_i^T$. Let $\epsilon > 0$ be a precision parameter. There exists a quantum circuit for performing the
385 mapping $\frac{1}{\|M\|_F} \sum_i^d \sum_j^d M_{ij} |i\rangle |j\rangle |0\rangle \rightarrow \frac{1}{\|M\|_F} \sum_i^d \sigma_i |\mathbf{u}_i\rangle |\mathbf{v}_i\rangle |\bar{\sigma}_i\rangle$ such that $|\sigma_i - \bar{\sigma}_i| \leq \epsilon$ with
386 probability at least $1 - 1/\text{poly}(d)$ in time $O\left(\frac{\|M\|_F \text{polylog } d}{\epsilon}\right)$.*388 To extract classical singular vectors corresponding to the largest singular value from a quantum state,
389 ℓ_2 norm quantum state tomography is required.391 **Lemma 6** (ℓ_2 state-vector tomography Kerenidis et al. (2020; 2019d)). *Given a unitary mapping
392 $U_x : |0\rangle \rightarrow |\mathbf{x}\rangle$ in time $T(U_x)$ and $\delta > 0$, there is an algorithm that produces an estimate
393 $\bar{\mathbf{x}} \in \mathbb{R}^d$ with $\|\bar{\mathbf{x}}\|_2 = 1$ such that $\|\mathbf{x} - \bar{\mathbf{x}}\|_2 \leq \delta$ with probability at least $1 - 1/\text{poly}(d)$ in time
394 $O\left(T(U_x) \frac{d \log d}{\delta^2}\right)$.*396 **Quantum top singular vector extraction (QTSVE).** The goal of the quantum subroutine in each
397 iteration is to find the top right / left singular vectors of the gradient matrix. First, we prepare the
398 gradient matrix state using the quantum access as stated in Assumption 4, then we perform QSVE to
399 this state. The quantum maximum finding is applied to obtain the quantum state corresponding to the
400 largest singular value. Prepare sufficient quantum states corresponding to the largest singular value
401 until satisfying the requirement of tomography, then perform quantum state tomography to extract
402 the corresponding right / left classical singular vectors. This procedure is shown in Lemma 7, with
403 the proof given in Appendix B.8. Note that the success probability of QTSVE can be improved by
404 repeating it logarithmic times and then taking the average.405 **Lemma 7.** *(Quantum top singular vector extraction) Let there be efficient quantum access to a
406 matrix $M \in \mathbb{R}^{d \times d}$, with singular value decomposition $M = \sum_i^d \sigma_i \mathbf{u}_i \mathbf{v}_i^T$. Define $p = \frac{\sigma_1^2(M)}{\sum_{i=1}^d \sigma_i^2}$.
407 There exist quantum algorithms that with time complexity $O\left(\frac{\|M\|_F d \text{polylog } d}{\sqrt{p} \delta^2}\right)$, give the estimated
408 top singular value $\bar{\sigma}_1$ of M to precision ϵ and the corresponding unit estimated singular vectors \mathbf{u}, \mathbf{v}
409 to precision δ such that $\|\mathbf{u} - \mathbf{u}_{top}\| \leq \delta, \|\mathbf{v} - \mathbf{v}_{top}\| \leq \delta$ with probability at least $1 - 1/\text{poly}(d)$.*410 **Convergence Analysis.** Our quantum Frank-Wolfe algorithm for nuclear norm constraint (Algorithm 3) then follows, with the analysis given in Appendix B.9.411 **Theorem 3.** *(Quantum FW with QTSVE) By setting $\delta_t = \frac{C_f}{2(t+2)\sigma_1(M_t)}$ and $\epsilon_t \leq (\sigma_1(M_t) -$
412 $\sigma_2(M_t))/2$ for $t \in [T]$, the quantum algorithm (Algorithm 3) solves the nuclear norm constraint
413 optimization problem for any precision ε such that $f(X^T) - f(X^*) \leq \varepsilon$ in $T = \frac{4C_f}{\varepsilon} - 2$ rounds,
414 with time complexity $\tilde{O}\left(\frac{r\sigma_1^3(M_t)d}{(\sigma_1(M_t) - \sigma_2(M_t))\varepsilon^2}\right)$ for computing the update direction per round, where
415 r is the rank of the gradient matrix.*416 In computing the update direction, Algorithm 3 reduces a $O(d\varepsilon/r\sigma_1^2(M))$ factor to the power
417 method and $O(d\varepsilon^{1.5}/r\sigma_1^{2.5}(M))$ to the Lanczos method, respectively. See Remark 2 for more in-
418 formation about parameter choosing.423 4.2 QUANTUM FRANK-WOLFE WITH QUANTUM POWER METHOD
424425 The second framework is to accelerate the power method directly with quantum matrix-vector
426 multiplication method and quantum tomography. The classical power method constructs a se-
427 quence $\mathbf{z}_0, \dots, \mathbf{z}_k$, where $\mathbf{z}_0 = \mathbf{b}$ is drawn uniformly random over a unit sphere $\mathbf{b} : \|\mathbf{b}\|_2 = 1$, and
428 $\mathbf{z}_{i+1} = M^\top M \mathbf{z}_i$ for $i \geq 1$, ($\mathbf{z}_{i+1} = MM^\top \mathbf{z}_i$ for the left singular vector, respectively). After
429 $k = \frac{C_0 \sigma_1(M) \ln d}{\varepsilon}$, we have $\left| \frac{\mathbf{z}_k^\top M \mathbf{z}_k}{\|\mathbf{z}_k\|_2^2} - \sigma_1(M) \right| \leq \varepsilon$, where C_0 is a constant.430 **Quantum power method (QPM).** Using the quantum access given in Assumption 4, the quantum
431 matrix-vector multiplication can be performed efficiently:

432 **Algorithm 3** Quantum Frank-Wolfe Algorithm for Nuclear Norm Constraint with QTSVE

433 1: **Input:** Solution precision ε , singular value precision $\{\epsilon_t\}_{t=1}^T$, tomography precision $\{\delta_t\}_{t=1}^T$.

434 2: **Output:** $X^{(T)}$ such that $f(X^T) - f(X^*) \leq \varepsilon$.

435 3: **Initialize:** Let $X^{(1)} \in \mathcal{D}$.

436 4: Let $T = \frac{4C_f}{\varepsilon} - 2$.

437 5: **for** $t = 1, \dots, T$ **do**

438 6: Let $\gamma_t = \frac{2}{t+2}$.

439 7: Prepare $\frac{1}{\|M\|_F} \sum_i^d \sum_j^d M_{ij} |i\rangle |j\rangle |0\rangle$.

440 8: Perform QSVE (Lemma 5) to get $\frac{1}{\sqrt{\sum_i^r \sigma_i^2}} \sum_i^r \sigma_i |\mathbf{u}_i\rangle |\mathbf{v}_i\rangle |\bar{\sigma}_i\rangle$, where $|\sigma_i - \bar{\sigma}_i| \leq \epsilon_t$.

441 9: Apply quantum maximum finding (Lemma 4) to the third register to get $|\mathbf{u}_{top}\rangle |\mathbf{v}_{top}\rangle |\bar{\sigma}_1\rangle$.

442 10: Perform ℓ_2 -norm tomography (Lemma 6), to obtain \mathbf{u}, \mathbf{v} , where $\|\mathbf{u} - \mathbf{u}_{top}\| \leq \delta_t$,

443 $\|\mathbf{v} - \mathbf{v}_{top}\| \leq \delta_t$.

444 11: Set $S = \mathbf{u}\mathbf{v}^\top$. Update $X^{(t+1)} = (1 - \gamma_t)X^{(t)} + \gamma_t S$.

445 12: **end for**

446 **Algorithm 4** Quantum Frank-Wolfe Algorithm for Nuclear Norm Constraint with QPM

447 1: **Input:** Solution precision ε , multiplication times $\{k_t\}_{t=1}^T$, multiplication precision $\{\delta_t\}_{t=1}^T$,
448 tomography precision $\{\delta'_t\}_{t=1}^T$.

449 2: **Output:** $X^{(T)}$ such that $f(X^T) - f(X^*) \leq \varepsilon$.

450 3: **Initialize:** Let $X^{(1)} \in \mathcal{D}$.

451 4: Let $T = \frac{4C_f}{\varepsilon} - 2$.

452 5: **for** $t = 1, \dots, T$ **do**

453 6: Let $\gamma_t = \frac{2}{t+2}$.

454 7: Prepare $\frac{1}{\|M\|_F} \sum_i^d \sum_j^d M_{ij} |i\rangle |j\rangle |\mathbf{b}\rangle |\mathbf{b}\rangle$, where \mathbf{b} is the uniform superposition state.

455 8: Apply quantum power method (Lemma 9) to get $\frac{1}{\|M\|_F} \sum_i^d \sum_j^d M_{ij} |i\rangle |j\rangle |\bar{\mathbf{z}}_u\rangle |\bar{\mathbf{z}}_v\rangle$, where
456 $\|\bar{\mathbf{z}}_u - (MM^\top)^k \mathbf{b}\| \leq \delta_t$, $\|\bar{\mathbf{z}}_v - (M^\top M)^k \mathbf{b}\| \leq \delta_t$.

457 9: Perform ℓ_2 -norm tomography (Lemma 6) to obtain \mathbf{u}, \mathbf{v} , where $\|\mathbf{u} - \bar{\mathbf{z}}_u\| \leq \delta'_t$,
458 $\|\mathbf{v} - \bar{\mathbf{z}}_v\| \leq \delta'_t$.

459 10: Set $S = \mathbf{u}\mathbf{v}^\top$. Update $X^{(t+1)} = (1 - \gamma_t)X^{(t)} + \gamma_t S$.

460 11: **end for**

461 **Lemma 8.** (Quantum matrix-vector multiplication (Theorem 4 Bellante et al. (2022)), Chakraborty
462 et al. (2019)) Let there be quantum access to the matrix $M \in R^{d \times d}$ with $\sigma_{\max} \leq 1$, and to a vector
463 $\mathbf{z} \in R^d$. Let $\|M\mathbf{z}\| \geq \gamma'$. There exists a quantum algorithm that creates a state $|\mathbf{y}\rangle$ such that
464 $\||\mathbf{y}\rangle - |M\mathbf{z}\rangle\| \leq \epsilon$ in time $\tilde{O}\left(\frac{1}{\gamma'} \|M\|_F \log(1/\epsilon)\right)$, with probability at least $1 - 1/\text{poly}(d)$.

465 Apply $2k$ times of quantum matrix-vector multiplication, we can get a quantum state corresponding
466 to \mathbf{z}_k , as shown in Lemma 9, with proof given in Appendix B.10. A similar process can be
467 constructed to compute $(MM^\top)^k \mathbf{b}$ (corresponding to the left singular vector) simultaneously.

468 **Lemma 9.** (Quantum power method) Let there be quantum access to the matrix $M \in R^{d \times d}$ with
469 $\sigma_{\max} \leq 1$, and to a vector $\mathbf{z} \in R^d$. Let γ'_{\min} be the lower bound of $\|(M^\top M)^i \mathbf{z}\|$ for all $i \in [k]$.
470 There exists a quantum algorithm that creates a state $|\mathbf{y}\rangle$ such that $\||\mathbf{y}\rangle - |(M^\top M)^k \mathbf{z}\rangle\| \leq \delta$ in
471 time $\tilde{O}\left(\frac{k}{\gamma'_{\min}} \|M\|_F \log(1/\delta)\right)$, with probability at least $1 - O(k/\text{poly}(d))$.

472 **Convergence Analysis.** After quantum state tomography, we can extract the classical top singular
473 vectors. Note that the success probability of QPM and tomography can be improved by repeating
474 the whole procedure logarithmic times and then taking the average. Our quantum Frank-Wolfe
475 algorithm (Algorithm 4) for nuclear norm constraint then follows, and the parameters choosing and
476 convergence analysis are given in Theorem 4, with the proof given in Appendix B.11.

477 **Theorem 4.** (Quantum FW with QPM) By setting $k_t = \frac{2C_0\sigma_1(M_t) \ln d}{\varepsilon}$, $\delta_t = \delta'_t = \frac{\varepsilon\gamma'_{\min}}{16\sigma_1(M_t)}$ for
478 $t \in [T]$, the quantum algorithm (Algorithm 4) solves the nuclear norm constraint optimization
479 problem for any precision ε such that $f(X^T) - f(X^*) \leq \varepsilon$ in $T = \frac{4C_f}{\varepsilon} - 2$ rounds, with time
480

486 complexity $\tilde{O}\left(\frac{\sqrt{r}\sigma_1^4(M_t)d}{(1-\sigma_1(M_t))\gamma'^3_{\min}\varepsilon^3}\right)$ for computing the update direction per round, where r is the
 487 rank of the gradient matrix, C_0 is a constant and γ'_{\min} is the lower bound of $\|(M_t^\top M_t)^i b\|$ for all
 488 $i \in [k]$.
 489

490 In computing the update direction, Algorithm 4 reduces a $O(d\varepsilon^2\gamma'^3_{\min}/\sqrt{r}\sigma_1^3(M))$ factor to the
 491 power method and $O(d\varepsilon^{2.5}\gamma'^3_{\min}/\sqrt{r}\sigma_1^{3.5}(M))$ to the Lanczos method. A discussion of this section
 492 is given in Appendix A.5.
 493

494 5 CONCLUSION AND FUTURE WORK

495 This paper addresses the projection-free sparse convex optimization problem. We propose several
 496 quantum Frank-Wolfe algorithms for both vector and matrix domains, demonstrating the quantum
 497 speedup over the classical methods with respect to the dimension of the feasible set.

498 For future work, we aim to extend quantum Frank-Wolfe methods to stochastic and online opti-
 499 mization frameworks, to characterize quantum advantages in projection-free regret minimization.
 500 Meanwhile, Jaggi (2013) highlights several interesting cases involving matrix norms, where classi-
 501 cal approaches often rely on computationally expensive singular value decomposition. A potential
 502 avenue of interest is determining whether quantum computing can yield greater speedups in such
 503 settings. Moreover, as mentioned in Appendix A.5, the gradient in the matrix completion is sparse,
 504 which might allow for further acceleration via quantum sparse matrix multiplication, constituting
 505 an interesting direction for future research. These investigations would collectively advance the
 506 understanding of quantum-enhanced projection-free optimization.
 507

508 REFERENCES

509
 510 Zeyuan Allen-Zhu, Elad Hazan, Wei Hu, and Yuanzhi Li. Linear convergence of a frank-wolfe type
 511 algorithm over trace-norm balls. *Advances in neural information processing systems*, 30, 2017.
 512
 513 A Ambainis. Variable time amplitude amplification and a faster quantum algorithm for solving
 514 systems of linear equations. In *Symp. Theoretical Aspects of Computer Science (STACS 2012)*,
 515 volume 14, pp. 636–47, 2012.
 516
 517 Andris Ambainis and Robert Špalek. Quantum algorithms for matching and network flows. In
 518 *Annual Symposium on Theoretical Aspects of Computer Science*, pp. 172–183, Berlin, 2006.
 519 Springer.
 520
 521 Simon Apers and Sander Gribling. Quantum speedups for linear programming via interior point
 522 methods. *arXiv preprint arXiv:2311.03215*, 2023.
 523
 524 Andreas Argyriou, Marco Signoretto, and Johan Suykens. Hybrid conditional gradient-smoothing
 525 algorithms with applications to sparse and low rank regularization. *Regularization, Optimization,*
 526 *Kernels, and Support Vector Machines*, pp. 53–82, 2014.
 527
 528 Michel Baes and Michael Bürgisser. Smoothing techniques for solving semi-definite programs with
 529 many constraints. *Optimization Online*, 2009.
 530
 531 Armando Bellante, Alessandro Luongo, and Stefano Zanero. Quantum algorithms for svd-based
 532 data representation and analysis. *Quantum Machine Intelligence*, 4(2):20, 2022.
 533
 534 Albert S Berahas, Liyuan Cao, Krzysztof Choromanski, and Katya Scheinberg. A theoretical and
 535 empirical comparison of gradient approximations in derivative-free optimization. *Foundations of*
 536 *Computational Mathematics*, 22(2):507–560, 2022.
 537
 538 Immanuel M Bomze, Francesco Rinaldi, and Damiano Zeffiro. Frank–wolfe and friends: a journey
 539 into projection-free first-order optimization methods. *4OR*, 19:313–345, 2021.
 540
 541 Fernando GSL Brandão and Krysta M Svore. Quantum speed-ups for solving semidefinite programs.
 542 In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science*, pp. 415–426, Pis-
 543 cataway, NJ, 2017. IEEE.

540 Fernando GSL Brandão, Amir Kalev, Tongyang Li, Cedric Yen-Yu Lin, Krysta M Svore, and Xiaodi
 541 Wu. Quantum SDP solvers: Large speed-ups, optimality, and applications to quantum learning.
 542 In *46th International Colloquium on Automata, Languages, and Programming*, pp. 27–1. Schloss
 543 Dagstuhl–Leibniz-Zentrum für Informatik, 2019.

544 Gilles Brassard, Peter Høyer, Michele Mosca, and Alain Tapp. Quantum amplitude amplification
 545 and estimation. *Contemporary Mathematics*, 305:53–74, 2002.

546 Michael D Canon and Clifton D Cullum. A tight upper bound on the rate of convergence of frank-
 547 wolfe algorithm. *SIAM Journal on Control*, 6(4):509–516, 1968.

548 Balthazar Casalé, Giuseppe Di Molfetta, Hachem Kadri, and Liva Ralaivola. Quantum bandits.
 549 *Quantum Machine Intelligence*, 2(1):1–7, 2020.

550 Shouvanik Chakrabarti, Andrew M Childs, Tongyang Li, and Xiaodi Wu. Quantum algorithms and
 551 lower bounds for convex optimization. *Quantum*, 4:221, 2020.

552 Shantanav Chakraborty, András Gilyén, and Stacey Jeffery. The power of block-encoded matrix
 553 powers: Improved regression techniques via faster hamiltonian simulation. In *46th International
 554 Colloquium on Automata, Languages, and Programming*, pp. 33–1. Schloss Dagstuhl–Leibniz-
 555 Zentrum für Informatik, 2019.

556 Lin Chen, Christopher Harshaw, Hamed Hassani, and Amin Karbasi. Projection-free online opti-
 557 mization with stochastic gradient: From convexity to submodularity. In *International Conference
 558 on Machine Learning*, pp. 814–823. PMLR, 2018.

559 Yanlin Chen and Ronald de Wolf. Quantum algorithms and lower bounds for linear regression with
 560 norm constraints. In *50th International Colloquium on Automata, Languages, and Programming*,
 561 pp. 38–1. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2023.

562 Yanlin Chen, András Gilyén, and Ronald de Wolf. A quantum speed-up for approximating the top
 563 eigenvectors of a matrix. In *Proceedings of the 2025 Annual ACM-SIAM Symposium on Discrete
 564 Algorithms (SODA)*, pp. 994–1036. SIAM, 2025a.

565 Zherui Chen, Yuchen Lu, Hao Wang, Yizhou Liu, and Tongyang Li. Quantum langevin dynamics
 566 for optimization. *Communications in Mathematical Physics*, 406(3):52, 2025b.

567 Andrew M Childs, Robin Kothari, and Rolando D Somma. Quantum algorithm for systems of linear
 568 equations with exponentially improved dependence on precision. *SIAM Journal on Computing*,
 569 46(6):1920–1950, 2017.

570 Alexandre d’Aspremont. Smooth optimization with approximate gradient. *SIAM Journal on Opti-
 571 mization*, 19(3):1171–1183, 2008.

572 Miroslav Dudik, Zaid Harchaoui, and Jérôme Malick. Lifted coordinate descent for learning with
 573 trace-norm regularization. In *Artificial intelligence and statistics*, pp. 327–336. PMLR, 2012.

574 Joseph C Dunn and S Harshbarger. Conditional gradient algorithms with open loop step size rules.
 575 *Journal of Mathematical Analysis and Applications*, 62(2):432–444, 1978.

576 Christoph Durr and Peter Hoyer. A quantum algorithm for finding the minimum. *arXiv preprint
 577 quant-ph/9607014*, 1996.

578 Christoph Dürr, Mark Heiligman, Peter Høyer, and Mehdi Mhalla. Quantum query complexity of
 579 some graph problems. *SIAM Journal on Computing*, 35(6):1310–1328, 2006.

580 Marguerite Frank et al. An algorithm for quadratic programming. *Naval research logistics quarterly*,
 581 3(1-2):95–110, 1956.

582 Dan Garber. Faster projection-free convex optimization over the spectrahedron. *Advances in Neural
 583 Information Processing Systems*, 29, 2016.

584 Dan Garber and Elad Hazan. A linearly convergent variant of the conditional gradient algorithm
 585 under strong convexity, with applications to online and stochastic optimization. *SIAM Journal on
 586 Optimization*, 26(3):1493–1528, 2016.

594 András Gilyén, Srinivasan Arunachalam, and Nathan Wiebe. Optimizing quantum optimization
 595 algorithms via faster quantum gradient computation. In *Proceedings of the Thirtieth Annual*
 596 *ACM-SIAM Symposium on Discrete Algorithms*, pp. 1425–1444, Philadelphia, PA, 2019. SIAM.
 597

598 Lov K Grover. A fast quantum mechanical algorithm for database search. In *Proceedings of the*
 599 *twenty-eighth annual ACM symposium on Theory of computing*, pp. 212–219, New York, NY,
 600 1996. ACM.

601 Zaid Harchaoui, Anatoli Juditsky, and Arkadi Nemirovski. Conditional gradient algorithms for
 602 norm-regularized smooth convex optimization. *Mathematical Programming*, 152(1):75–112,
 603 2015.

604 Aram W Harrow, Avinatan Hassidim, and Seth Lloyd. Quantum algorithm for linear systems of
 605 equations. *Physical review letters*, 103(15):150502, 2009.

606

607 Hamed Hassani, Amin Karbasi, Aryan Mokhtari, and Zebang Shen. Stochastic conditional gradi-
 608 ent++:(non) convex minimization and continuous submodular maximization. *SIAM Journal on*
 609 *Optimization*, 30(4):3315–3344, 2020.

610 Elad Hazan and Satyen Kale. Projection-free online learning. *arXiv preprint arXiv:1206.4657*,
 611 2012.

612 Elad Hazan and Haipeng Luo. Variance-reduced and projection-free stochastic optimization. In
 613 *International Conference on Machine Learning*, pp. 1263–1271. PMLR, 2016.

614

615 Elad Hazan, Satyen Kale, and Shai Shalev-Shwartz. Near-optimal algorithms for online matrix
 616 prediction. In *Conference on Learning Theory*, pp. 38–1. JMLR Workshop and Conference Pro-
 617 ceedings, 2012.

618 Jianhao He, Feidiao Yang, Jialin Zhang, and Lvzhou Li. Quantum algorithm for online convex
 619 optimization. *Quantum Science and Technology*, 7(2):025022, 2022.

620

621 Jianhao He, Chengchang Liu, Xutong Liu, Lvzhou Li, and John C.S. Lui. Quantum algorithm for
 622 online exp-concave optimization. In *Proceedings of the 41st International Conference on Machine*
 623 *Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 17946–17971. PMLR,
 624 21–27 Jul 2024.

625 Xiaoyu He, Jialin Zhang, and Xiaoming Sun. Quantum search with prior knowledge. *arXiv preprint*
 626 *arXiv:2009.08721*, 2020.

627

628 Martin Jaggi. Revisiting frank-wolfe: Projection-free sparse convex optimization. In *International*
 629 *Conference on Machine Learning*, pp. 427–435. PMLR, 2013.

630 Rodolphe Jenatton, Jean-Yves Audibert, and Francis Bach. Structured variable selection with
 631 sparsity-inducing norms. *The Journal of Machine Learning Research*, 12:2777–2824, 2011.

632

633 Stephen P Jordan. Fast quantum algorithm for numerical gradient estimation. *Physical review*
 634 *letters*, 95(5):050501, 2005.

635 Iordanis Kerenidis and Anupam Prakash. Quantum recommendation systems. In *8th Innovations in*
 636 *Theoretical Computer Science Conference (ITCS 2017)*. Schloss Dagstuhl-Leibniz-Zentrum fuer
 637 Informatik, 2017.

638 Iordanis Kerenidis and Anupam Prakash. A quantum interior point method for LPs and SDPs. *ACM*
 639 *Transactions on Quantum Computing*, 1(1):1–32, 2020a.

640

641 Iordanis Kerenidis and Anupam Prakash. Quantum gradient descent for linear systems and least
 642 squares. *Physical Review A*, 101(2):022316, 2020b.

643 Iordanis Kerenidis, Anupam Prakash, and Dániel Szilágyi. Quantum algorithms for portfolio opti-
 644 mization. In *Proceedings of the 1st ACM Conference on Advances in Financial Technologies*, pp.
 645 147–155, 2019a.

646 Iordanis Kerenidis, Anupam Prakash, and Dániel Szilágyi. Quantum algorithms for second-order
 647 cone programming and support vector machines. *arXiv preprint arXiv:1908.06720*, 2019b.

648 Iordanis Kerenidis, Anupam Prakash, and Dániel Szilágyi. A quantum interior-point method for
 649 second-order cone programming. *[Research Report] IRIF*. 2019. *ffhal-02138307*, 2019c.
 650

651 Iordanis Kerenidis, Anupam Prakash, and Dániel Szilágyi. Quantum algorithms for portfolio opti-
 652 mization. In *Proceedings of the 1st ACM Conference on Advances in Financial Technologies*, pp.
 653 147–155, 2019d.

654 Iordanis Kerenidis, Alessandro Luongo, and Anupam Prakash. Quantum expectation-maximization
 655 for gaussian mixture models. In *International Conference on Machine Learning*, pp. 5187–5197.
 656 PMLR, 2020.

657 Jacek Kuczyński and Henryk Woźniakowski. Estimating the largest eigenvalue by the power and
 658 lanczos algorithms with a random start. *SIAM journal on matrix analysis and applications*, 13
 659 (4):1094–1122, 1992.

660 Guanghui Lan. The complexity of large-scale convex programming under a linear optimization
 661 oracle. *arXiv preprint arXiv:1309.5550*, 2013.

663 Guanghui Lan and Yi Zhou. Conditional gradient sliding for convex optimization. *SIAM Journal on
 664 Optimization*, 26(2):1379–1409, 2016.

666 Evgeny S Levitin and Boris T Polyak. Constrained minimization methods. *USSR Computational
 667 mathematics and mathematical physics*, 6(5):1–50, 1966.

668 Kfir Levy and Andreas Krause. Projection free online learning over smooth sets. In *The 22nd
 669 international conference on artificial intelligence and statistics*, pp. 1458–1466. PMLR, 2019.

671 Shigui Li, Linzhang Lu, Xing Qiu, Zhen Chen, and Delu Zeng. Tighter bound estimation for efficient
 672 biquadratic optimization over unit spheres. *Journal of Global Optimization*, 90(2):323–353, 2024.

674 Tongyang Li and Ruizhe Zhang. Quantum speedups of optimizing approximately convex functions
 675 with applications to logarithmic regret stochastic convex bandits. *Advances in Neural Information
 676 Processing Systems*, 35:3152–3164, 2022.

677 Tongyang Li, Shouvanik Chakrabarti, and Xiaodi Wu. Sublinear quantum algorithms for training
 678 linear and kernel-based classifiers. In *International Conference on Machine Learning*, pp. 3815–
 679 3824, Palo Alto, CA, 2019. AAAI.

680 Debbie Lim and Patrick Rebentrost. A quantum online portfolio optimization algorithm. *arXiv
 681 preprint arXiv:2208.14749*, 2022.

683 Seth Lloyd, Masoud Mohseni, and Patrick Rebentrost. Quantum algorithms for supervised and
 684 unsupervised machine learning. *arXiv preprint arXiv:1307.0411*, 2013.

685 Seth Lloyd, Masoud Mohseni, and Patrick Rebentrost. Quantum principal component analysis.
 686 *Nature Physics*, 10(9):631–633, 2014.

688 Ari Mizel. Critically damped quantum search. *Physical review letters*, 102(15):150501, 2009.

689

690 Yurii Nesterov. Smoothing technique and its applications in semidefinite optimization. *Mathemati-
 691 cal Programming*, 110(2):245–259, 2007.

692 Federico Pierucci, Zaid Harchaoui, and Jérôme Malick. *A smoothing approach for composite con-
 693 ditional gradient with nonsmooth loss*. PhD thesis, INRIA Grenoble, 2014.

694

695 Sebastian Pokutta. The frank-wolfe algorithm: a short introduction. *Jahresbericht der Deutschen
 696 Mathematiker-Vereinigung*, pp. 1–33, 2023.

697 Patrick Rebentrost, Masoud Mohseni, and Seth Lloyd. Quantum support vector machine for big
 698 data classification. *Physical review letters*, 113(13):130503, 2014.

699

700 Patrick Rebentrost, Maria Schuld, Leonard Wossnig, Francesco Petruccione, and Seth Lloyd. Quan-
 701 tum gradient descent and newton’s method for constrained polynomial optimization. *New Journal
 of Physics*, 21(7):073023, 2019.

702 Przemysław Sadowski. Quantum search with prior knowledge. *arXiv preprint arXiv:1506.04030*,
 703 2015.

704

705 Aaron Sidford and Chenyi Zhang. Quantum speedups for stochastic optimization. *Advances in*
 706 *Neural Information Processing Systems*, 36:35300–35330, 2023.

707 Joran van Apeldoorn and András Gilyén. Improvements in quantum SDP-solving with applications.
 708 In *46th International Colloquium on Automata, Languages, and Programming*, pp. 99, Wadern,
 709 2019a. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik.

710

711 Joran van Apeldoorn and András Gilyén. Quantum algorithms for zero-sum games. *arXiv preprint*
 712 *arXiv:1904.03180*, 2019b.

713 Joran van Apeldoorn, András Gilyén, Sander Gribling, and Ronald de Wolf. Quantum SDP-solvers:
 714 Better upper and lower bounds. In *2017 IEEE 58th Annual Symposium on Foundations of Com-*
 715 *puter Science (FOCS)*, pp. 403–414, Piscataway, NJ, 2017. IEEE.

716

717 Joran van Apeldoorn, András Gilyén, Sander Gribling, and Ronald de Wolf. Convex optimization
 718 using quantum oracles. *Quantum*, 4:220, 2020.

719 Zongqi Wan, Zhijie Zhang, Tongyang Li, Jialin Zhang, and Xiaoming Sun. Quantum multi-armed
 720 bandits and stochastic linear bandits enjoy logarithmic regrets. In *Proceedings of the AAAI Con-*
 721 *ference on Artificial Intelligence*, volume 37, pp. 10087–10094, 2023.

722 Daochen Wang, Xuchen You, Tongyang Li, and Andrew M Childs. Quantum exploration algo-
 723 rithms for multi-armed bandits. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
 724 volume 35, pp. 10102–10110, 2021.

725

726 Philip Wolfe. Convergence theory in nonlinear programming. *Integer and nonlinear programming*,
 727 pp. 1–36, 1970.

728

729 Leonard Wossnig, Zhikuan Zhao, and Anupam Prakash. Quantum linear system algorithm for dense
 730 matrices. *Physical review letters*, 120(5):050502, 2018.

731 Jiahao Xie, Zebang Shen, Chao Zhang, Boyu Wang, and Hui Qian. Efficient projection-free online
 732 methods with stochastic recursive gradient. In *Proceedings of the AAAI Conference on Artificial*
 733 *Intelligence*, volume 34, pp. 6446–6453, 2020.

734

735 Theodore J Yoder, Guang Hao Low, and Isaac L Chuang. Fixed-point quantum search with an
 736 optimal number of queries. *Physical review letters*, 113(21):210501, 2014.

737

738 Alp Yurtsever, Suvrit Sra, and Volkan Cevher. Conditional gradient methods via stochastic path-
 739 integrated differential estimator. In *International Conference on Machine Learning*, pp. 7282–
 7291. PMLR, 2019.

740

741 Chenyi Zhang and Tongyang Li. Quantum lower bounds for finding stationary points of nonconvex
 742 functions. In *International Conference on Machine Learning*, pp. 41268–41299. PMLR, 2023.

743

744 Mingrui Zhang, Zebang Shen, Aryan Mokhtari, Hamed Hassani, and Amin Karbasi. One sample
 745 stochastic frank-wolfe. In *International Conference on Artificial Intelligence and Statistics*, pp.
 4012–4023. PMLR, 2020.

746

747 Xinhua Zhang, Dale Schuurmans, and Yao-liang Yu. Accelerated training for matrix-norm regular-
 748 ization: A boosting approach. *Advances in Neural Information Processing Systems*, 25, 2012.

749

750 Yixin Zhang, Chenyi Zhang, Cong Fang, Liwei Wang, and Tongyang Li. Quantum algorithms and
 751 lower bounds for finite-sum optimization. In *Proceedings of the 41st International Conference on*
 752 *Machine Learning*, pp. 60244–60270. PMLR, 2024.

753

754

755

756 **Appendix**
757758 **A EXTENSION AND DISCUSSION**
759760 **A.1 QUANTUM FRANK-WOLFE OVER VECTORS WITH BOUNDED-ERROR JORDAN**
761 **ALGORITHM**
762763 The quantum Frank-Wolfe Algorithm with Bounded-error Jordan's Algorithm is shown in Algo-
764 rithm 5. We reformulate the results of the bounded-error Jordan algorithm from He et al. (2024) in
765 terms of infinity norm error, with the proof detailed in the Appendix B.4.
766767 **Algorithm 5** Quantum Frank-Wolfe Algorithm with Bounded-error Jordan Algorithm
768

- 1: **Input:** Solution precision ε , gradient precision $\{\sigma_t\}_{t=1}^T$.
- 2: **Output:** $\mathbf{x}^{(T)}$ such that $f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \varepsilon$.
- 3: **Initialize:** Let $\mathbf{x}^{(1)} \in \mathcal{D}$.
- 4: Let $T = \frac{4C_f}{\varepsilon} - 2$.
- 5: **for** $t = 1, \dots, T$ **do**
- 6: Let $\gamma_t = \frac{2}{t+2}$.
- 7: Using Algorithm 7 to get the whole vector of estimated gradient $\tilde{\nabla} f_t(\mathbf{x}_t)$.
- 8: Scan all the component of $\tilde{\nabla} f_t(\mathbf{x}_t)$ to find the coordinate i_t corresponding to the largest
absolute value of the estimated gradient component.
- 9: Set $s = -e_{i_t}$. Update $\mathbf{x}^{(t+1)} = (1 - \gamma_t)\mathbf{x}^{(t)} + \gamma_t s$.
- 10: **end for**

780 **Lemma 10.** (Lemma 1 He et al. (2024)) If f is G -Lipschitz continuous and L -smooth convex function
781 and can be accessed by a quantum function value oracle, then there exists an quantum algorithm
782 that for any $r > 0$ and $1 \geq \rho > 0$, gives the estimated gradient $g(x)$, which satisfies

783
$$\Pr[\|g(x) - \nabla f(x)\|_\infty > 8\pi n^2(n/\rho + 1)Lr/\rho] < \rho, \quad (11)$$

784

785 using $O(1)$ applications of \mathbf{U}_f and $O(d \log d)$ elementary gates. The space complexity is
786 $O\left(d \log \frac{G\rho}{4\pi d^2 Lr}\right)$.
787788 The next step is to determine the quantum gradient estimated parameters r_t in each Frank-Wolfe
789 iteration through convergence analysis.
790791 **Theorem 5.** (Quantum FW with bounded-error Jordan algorithm) By setting $r_t =$
792 $\frac{\rho C_f}{16\pi d^2(d/\rho+1)L(t+2)}$ for $t \in [T]$, the quantum algorithm (Algorithm 5) solves the sparsity constraint
793 optimization problem for any precision ε such that $f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \varepsilon$ in $T = \frac{4C_f}{\varepsilon} - 2$ rounds,
794 with $O(1)$ calls to the function value oracle \mathbf{U}_f per round.
795796 The proof is given in Appendix B.5. Substituting the parameter r_t into the space complexity yields
797 the qubit requirement as $O\left(d \log \frac{Gd}{\rho\varepsilon}\right)$. Since each gradient estimation succeeds with probability
798 $1 - \rho$, the probability that all T iterations succeed is at least $1 - T\rho$. By setting $\rho = p/T$, we ensure
799 an overall success probability of at least $1 - p$.
800801 **A.2 MORE EXTENSIONS OVER VECTORS FOR ATOMIC SETS**
802803 In this appendix, we give two more extension for the vector case. The first extension is to consider
804 $|\mathcal{A}| = N$, with each $a_j \in \mathcal{A}$ being τ -sparse with non-zero (index, value) pairs $(i_k, (a_j)_k)$, i.e., each
805 $a_j \in \mathbb{R}^d$, but has only τ non-zero elements. Assume that the non-zero elements are accessed with a
806 quantum oracle V which implements the transformation $V |j\rangle |k\rangle |0\rangle |0\rangle \rightarrow |j\rangle |k\rangle |i_k\rangle |(a_j)_k\rangle$. One
807 can construct a coherent access to the non-zero elements
808

809
$$V^{\otimes \tau} |j\rangle \bigotimes_{k=1}^{\tau} |k\rangle |0\rangle |0\rangle = |j\rangle \bigotimes_{k=1}^{\tau} |k\rangle |i_k\rangle |(a_j)_k\rangle \quad (12)$$

810 using τ calls of V . Then, a slight modification of the method of Section 3.1 can compute the FW
 811 update using $O(\tau\sqrt{N}\log(1/\delta))$ queries to V and U_g .
 812

813 The second extension is to consider latent group norm constraints, which have found use in inducing
 814 sparsity in problems in machine learning Jenatton et al. (2011). The ℓ_1 norm, d -simplex, group lasso
 815 etc. are all special cases of this.

816 Following Jaggi (2013) we let $\mathcal{G} = \{\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_{|\mathcal{G}|}\}$, $\mathbf{g}_i \subseteq [d]$, $\bigcup_i \mathbf{g}_i = [d]$. Note that the \mathbf{g}_i need
 817 not be disjoint. For each $\mathbf{g} \in \mathcal{G}$, let $\|\cdot\|_{\mathbf{g}}$ be an arbitrary ℓ_p norm, and define the **latent group norm**
 818

$$\begin{aligned} 819 \quad \|x\|_{\mathcal{G}} &:= \min_{v_{(\mathbf{g})} \in \mathbb{R}^{|\mathcal{G}|}} \sum_{\mathbf{g} \in \mathcal{G}} \|v_{(\mathbf{g})}\|_{\mathbf{g}} \\ 820 \quad \text{s.t.} \quad x &= \sum_{\mathbf{g} \in \mathcal{G}} v_{[\mathbf{g}]} \end{aligned} \quad (13)$$

821 where $v_{(\mathbf{g})} \in \mathbb{R}^{\mathbf{g}}$ is the restriction of $v \in \mathbb{R}^d$ to coordinates in \mathbf{g} , and $v_{[\mathbf{g}]} \in \mathbb{R}^d$ has zeros outside
 822 the support of \mathbf{g} . In this case, the Frank-Wolfe update corresponds to finding the value $s : \|s\|_{\mathcal{G}} \leq 1$
 823 such that $s^\top \nabla f(x) = \|\nabla f(x)\|_{\mathcal{G}}^*$, where $\|\nabla f(x)\|_{\mathcal{G}}^* = \max_{s: \|s\|_{\mathcal{G}} \leq 1} s^\top \nabla f(x)$.
 824

825 By Section 4.1 in Jaggi (2013), this norm is an atomic norm, and the dual norm is given by
 826

$$827 \quad \|\nabla f(x)\|_{\mathcal{G}}^* = \max_{\mathbf{g} \in \mathcal{G}} \|\nabla f(x)_{(\mathbf{g})}\|_{\mathbf{g}}^*, \quad (14)$$

828 which implies that
 829

$$830 \quad \max_{s: \|s\|_{\mathcal{G}} \leq 1} (-s^\top \nabla f(x)) = \max_{\mathbf{g} \in \mathcal{G}} \max_{s: \|s\|_{\mathbf{g}} \leq 1} (-s^\top \nabla f(x)). \quad (15)$$

831 Therefore, it suffices to consider each $\|\cdot\|_{\mathbf{g}}$ ball separately, and then do quantum maximizing over
 832 all the $|\mathcal{G}|$ balls to find the one that has the largest value of $\|-\nabla f(x)_{\mathbf{g}_i}\|_{p_i}^*$. The quantum Frank-
 833 Wolfe algorithm over latent group norm ball is then given in Algorithm 6. Note that by the absolute
 834 homogeneity property of dual norms,
 835

$$836 \quad \|\nabla f(x)_{(\mathbf{g})}\|_{\mathbf{g}}^* = \|-\nabla f(x)_{(\mathbf{g})}\|_{\mathbf{g}}^*, \quad (16)$$

837 certain negative signs have been omitted in the algorithmic formulation.
 838

839 **Algorithm 6** Quantum Frank-Wolfe Algorithm over Latent Group Norm Ball

- 840 1: **Input:** Gap ε , accuracy $\{\sigma_t\}_{t=1}^T$, iterations T .
- 841 2: **Initialize:** Let $\mathbf{x}^{(1)} \in \mathcal{D}$.
- 842 3: **for** $t = 1, \dots, T$ **do**
- 843 4: Let $\gamma_t = \frac{2}{t+2}$, $\mathbf{x} = \mathbf{x}^{(t)}$.
- 844 5: Prepare state $\sum_{i=1}^n |i\rangle_A |\mathbf{x}\rangle \bigotimes_{j=1}^{|\mathbf{g}_i|} |\mathbf{g}_{i,j}\rangle |0\rangle |0\rangle |0\rangle |0\rangle$.
- 845 6: Perform quantum gradient circuit to get $\sum_{i=1}^n |i\rangle_A |\mathbf{x}\rangle \bigotimes_{j=1}^{|\mathbf{g}_i|} |\mathbf{g}_{i,j}\rangle |g_{\mathbf{g}_{i,j}}(\mathbf{x})\rangle |0\rangle |0\rangle |0\rangle$,
 846 where $g_{\mathbf{g}_{i,j}}(\mathbf{x}) = \frac{f(\mathbf{x} + \sigma_t \mathbf{e}_{\mathbf{g}_{i,j}}) - f(\mathbf{x})}{\sigma_t}$
- 847 7: Compute $\sum_{i=1}^n |i\rangle_A |\mathbf{x}\rangle \left(\bigotimes_{j=1}^{|\mathbf{g}_i|} |\mathbf{g}_{i,j}\rangle |g_{\mathbf{g}_{i,j}}(\mathbf{x})\rangle \left| \text{sgn}(g_{\mathbf{g}_{i,j}}(\mathbf{x})) |g_{\mathbf{g}_{i,j}}(\mathbf{x})|^{q_i-1} \right. \right)$
 848 $\left| \left\| g(\mathbf{x})_{(\mathbf{g}_i)} \right\|_{p_i} \right\rangle \left| \left\| g(\mathbf{x})_{(\mathbf{g}_i)} \right\|_{p_i}^* \right\rangle$.
- 849 8: Apply quantum maximum finding on the last register, and then measure the rest registers,
 850 denote the result of the first register as i_t .
- 851 9: Initial $\mathbf{s} = 0$, set $s_{\mathbf{g}_{i_t,j}} = \text{sgn}(g_{\mathbf{g}_{i_t,j}}(\mathbf{x})) |g_{\mathbf{g}_{i_t,j}}(\mathbf{x})|^{q_{i_t}-1}$ for $j = 1$ to $|\mathbf{g}_i|$, where $\frac{1}{p_{i_t}} + \frac{1}{q_{i_t}} = 1$.
- 852 10: Then normalize \mathbf{s} .
- 853 11: Update $\mathbf{x}^{(t+1)} = (1 - \gamma_t) \mathbf{x}^{(t)} + \gamma_t \mathbf{s}$.

861 To simplify the proof of the query complexity of the quantum FW update (Lemma 11), we first as-
 862 sume that the gradient estimation and the maximum-finding are exact, with proof given in Appendix
 863 B.6. Then we give the error analysis and show how to choose the parameters σ_t in Theorem 6, with
 864 proof given in Appendix B.7.

864 **Lemma 11.** [Quantum FW update over latent group norm ball] Let $\|\cdot\|_{\mathcal{G}}$ be a latent group norm
 865 corresponding to $\mathcal{G} = \{\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_{|\mathcal{G}|}\}$, and let $|\mathbf{g}|_{\max} = \max_j |\mathbf{g}_j|$. Then, there exists a quantum
 866 algorithm computing the Frank-Wolfe update $s^* := \operatorname{argmax}_{\hat{s} \in \|\cdot\|_{\mathcal{G}}\text{-Ball}} \langle \hat{s}^\top g(\mathbf{x}) \rangle$ in $O(\sqrt{|\mathcal{G}|} |\mathbf{g}|_{\max})$
 867 calls to \mathbf{U}_f .
 868

869 **Theorem 6.** [Quantum FW over latent group norm ball] By setting $\sigma_t = \frac{C_f}{\sqrt{d}L(t+2) \max_{i \in [|\mathcal{G}|]} |\mathbf{g}_i|^{1/p_i}}$
 870 for $t \in [T]$, the quantum algorithm (Algorithm 6) solves the latent group norm constraint optimiza-
 871 tion problem for any precision ε such that $f(\mathbf{x}^T) - f(\mathbf{x}^*) \leq \varepsilon$ in $T = \frac{4C_f}{\varepsilon} - 2$ rounds, succeed with
 872 probability $1 - p$, with $O\left(\sqrt{|\mathcal{G}|} |\mathbf{g}|_{\max} \log \frac{C_f}{p\varepsilon}\right)$ calls to the function value oracle \mathbf{U}_f per round.
 873
 874

875 A.3 NOTATIONS AND ASSUMPTIONS FOR QUANTUM COMPUTATION

876 **Basic Notions in Quantum Computing.** Quantum computing utilizes Dirac notation as its mathe-
 877 matical foundation. Let $\{|i\rangle\}_{i=0}^{d-1}$ denote the computational basis of \mathbb{C}^d as $\{|i\rangle\}_{i=0}^{d-1}$, where $|i\rangle$ is a
 878 d -dimensional unit vector with 1 at the i^{th} position and 0 elsewhere. A d -dimensional quantum state
 879 is represented as a unit vector $|v\rangle = (v_1, v_2, \dots, v_d)^T = \sum_i v_i |i\rangle \in \mathbb{C}^d$ with complex amplitudes
 880 v_i satisfying $\sum_i |v_i|^2 = 1$.
 881

882 **Composite Systems.** The joint state of two quantum systems $|v\rangle \in \mathbb{C}^{d_1}$ and $|u\rangle \in \mathbb{C}^{d_2}$ is described
 883 by the tensor product $|v\rangle \otimes |u\rangle = (v_1 u_1, v_1 u_2, \dots, v_2 u_1, \dots, v_{d_1} u_{d_2}) \in \mathbb{C}^{d_1 \times d_2}$. The \otimes symbol is
 884 omitted when context permits.
 885

886 **Quantum Dynamics.** Closed system evolution is described by unitary transformations. Quantum
 887 measurement in the computational basis probabilistically projects the state onto a basis vector $|i\rangle$
 888 with the probability of the square of the magnitude of its amplitude. For example, measuring $|v\rangle =$
 889 $\sum_i v_i |i\rangle$ yields outcome i with probability $|v_i|^2$, followed by post-measurement state $|i\rangle$.
 890

891 **Quantum Access Models.** In general, In quantum computing, access to the objective function is
 892 facilitated through quantum oracles Q_f , which is a unitary transformation that maps a quantum state
 893 $|x\rangle |q\rangle$ to the state $|x\rangle |q + f(x)\rangle$, where $|x\rangle$, $|q\rangle$ and $|q + f(x)\rangle$ are basis states corresponding to
 894 the floating-point representations of x , q and $q + f(x)$. Moreover, given the superposition input
 895 $\sum_{x,q} \alpha_{x,q} |x\rangle |q\rangle$, by linearity the quantum oracle will output the state $\sum_{x,q} \alpha_{x,q} |x\rangle |q + f(x)\rangle$.
 896

897 A.4 EXTENDED RELATED WORKS

898 The Frank-Wolfe (FW) algorithm, also known as the conditional gradient method, has evolved
 899 through several key theoretical and applied research phases. The original FW framework Frank
 900 et al. (1956) established a projection-free method for quadratic programming with optimal conver-
 901 gence rates when solutions lie on the feasible set boundary, a property later rigorously proven by
 902 Canon & Cullum (1968). Wolfe’s away-step modification Wolfe (1970) addressed boundary solu-
 903 tion limitations, while Dunn’s extension Dunn & Harshbarger (1978) generalized FW to smooth
 904 optimization over Banach spaces using linear minimization oracles.
 905

906 Modern convergence analyzes were unified by Jaggi (2013), who introduced duality gap certificates
 907 for primal-dual convergence in constrained convex optimization. For strongly convex objectives,
 908 Garber & Hazan (2016) demonstrated accelerated linear convergence rates. Projection-free optimi-
 909 zation on non-smooth objective functions was studied in Lan (2013); Argyriou et al. (2014);
 910 Pierucci et al. (2014). Data-dependent convergence bounds on spectahedrons were improved by
 911 Garber (2016) and Allen-Zhu et al. (2017).
 912

913 Note that the framework was extended to online and stochastic optimizations, inspiring a series
 914 of seminal contributions Hazan & Kale (2012); Garber & Hazan (2016); Levy & Krause (2019);
 915 Lan & Zhou (2016); Hazan & Luo (2016); Chen et al. (2018); Hassani et al. (2020); Xie et al.
 916 (2020); Yurtsever et al. (2019); Zhang et al. (2020). Our future research will explore quantum-
 917 enhanced acceleration for these online/stochastic settings. Meanwhile, in recent years, FW methods
 918 have gained attention for their effectiveness in dealing with structured constraint problem arising in
 919 machine learning and data science, such as LASSO, SVM training, matrix completion and clustering
 920 detection. Readers are referred to Bomze et al. (2021); Pokutta (2023) for more information.
 921

The algorithms we develop in the matrix domain belong to the quantum algorithmic family for linear systems. This family originated with the seminal HHL algorithm Harrow et al. (2009), which solves quantum linear systems and achieves exponential speedups over classical methods for well-conditioned sparse matrices. Subsequent improvements reduced dependency on condition number and sparsity Ambainis (2012); Childs et al. (2017); Wossnig et al. (2018). The HHL framework has been successfully adapted to machine learning tasks including support vector machines Rebentrost et al. (2014), supervised and unsupervised machine learning Lloyd et al. (2013), principal component analysis Lloyd et al. (2014) and recommendation systems Kerenidis & Prakash (2017). One can reduce the condition number by preprocessing the matrix itself, and QRAM can help to accelerate such preprocessing. Based on this, the quantum singular value estimation method was developed in Kerenidis & Prakash (2017) and was generalized in Kerenidis & Prakash (2020b). Furthermore, recent work integrates QSVE with state-vector tomography, amplitude amplification/estimation, and spectral norm analysis to enable top- k singular vector extraction Bellante et al. (2022).

Recently, quantum computing has emerged as a promising new paradigm to accelerate a large number of important optimization problems, e.g., combinatorial optimization Grover (1996); Ambainis & Špalek (2006); Dürr et al. (2006); Durr & Hoyer (1996); Mizel (2009); Yoder et al. (2014); Sadowski (2015); He et al. (2020), linear programming Kerenidis & Prakash (2020a); Li et al. (2019); van Apeldoorn & Gilyén (2019b); Apers & Gribling (2023), second-order cone programming Kerenidis et al. (2019c;b;a), quadratic programming Kerenidis & Prakash (2020b), polynomial optimization Rebentrost et al. (2019), semi-definite optimization Kerenidis & Prakash (2020a); van Apeldoorn & Gilyén (2019a); Brandão & Svore (2017); Brandão et al. (2019); van Apeldoorn et al. (2017), convex optimization van Apeldoorn et al. (2020); Chakrabarti et al. (2020); Zhang et al. (2024), nonconvex optimization Zhang & Li (2023); Chen et al. (2025b), stochastic optimization Sidford & Zhang (2023) online optimization He et al. (2022; 2024); Lim & Rebentrost (2022), multi-arm bandit Casalé et al. (2020); Wang et al. (2021); Li & Zhang (2022); Wan et al. (2023). The quantum community is actively pursuing further accelerations of quantum computing in the field of optimization.

A.5 DISCUSSION OF THE TWO QUANTUM FRANK-WOLFE ALGORITHMS FOR THE MATRIX CASE

We essentially developed two complementary algorithms tailored to high-rank and low-rank gradient matrices, respectively. For Algorithm 3, quantum advantage exists when $d > r/\sqrt{\sigma_1 - \sigma_2}\epsilon$. For Algorithm 4, quantum advantage holds when $d > \sqrt{r}\sqrt{\sigma_1 - \sigma_2}/\epsilon^2(1 - \sigma_1)$. Since the quantum subroutines in the matrix section effectively process the gradient matrix normalized by its Frobenius norm, when this matrix has very low rank, $1 - \sigma_1$ tends to be small (approaching 0 when the rank is 1). In such cases, Algorithm 3 delivers better performance, whereas Algorithm 4 is more suitable otherwise. These two complementary algorithms deliver a quantum speedup of at least $\mathcal{O}(\sqrt{d})$.

Furthermore, the repetition steps required for quantum state tomography can be parallelized in the quantum computing cluster. By utilizing $\mathcal{O}(d)$ quantum computers simultaneously, the dependence of d in time complexity can be eliminated, giving a parallel time complexity of $\tilde{\mathcal{O}}\left(\frac{r\sigma_1^3(M)}{(\sigma_1(M) - \sigma_2(M))\epsilon^2}\right)$ and $\tilde{\mathcal{O}}\left(\frac{\sqrt{r}\sigma_1^4(M)}{(1 - \sigma_1(M))\gamma'^3\min\epsilon^3}\right)$.

Remark 1. Note that in Section 4, for simplicity of presentation, we focus on square matrices. However, all of the quantum techniques mentioned above can also be applied to non-square matrices, since the quantum singular value estimation can be applied to non-square matrices Kerenidis & Prakash (2020b).

Remark 2. All parameters can be determined during preprocessing. Since tomography constitutes the dominant part of the computational overhead, this preprocessing will not affect the final asymptotic complexity. The choice of δ_t relates to the maximum singular value of the current gradient matrix. Its range can be determined by running Quantum Singular Value Estimation (QSVE) followed by a maximum-value search algorithm. The purpose of ϵ_t is to ensure that the ordering of the largest and second-largest singular values does not become misordered during QSVE execution. This parameter can be determined via two methods: 1. During preprocessing, run QSVE-quantum maximum search and perform a binary search to find the critical point where two measurement outcomes appear. Then perform another binary search on ϵ_t to locate the critical point that distinguishes between these two outcomes. 2. Use the results of amplitude estimation as an indicator to

972 identify the critical point where a sudden jump in amplitude occurs. Since tomography remains the
 973 primary source of algorithmic overhead, the computational cost of this process will not impact the
 974 final asymptotic complexity.

975 **Remark 3.** Note that both the classical and quantum algorithms in this section assume that the
 976 gradients are pre-stored at the memory. In some applications, obtaining the gradients may not be
 977 easy, and even directly loading them into the memory would scale linearly with the size of the matrix.
 978 This work focuses only on the computation of the update direction, but the gradient calculation
 979 time, which is also ignored in classical algorithms Jaggi (2013), is explicitly included in the result
 980 Table 2. This is because in quantum computing, there exist several well-established algorithms for
 981 gradient estimation Jordan (2005); Gilyén et al. (2019). Moreover, in some applications (such as
 982 the matrix completion problem, which we will clarify below), the gradient matrix is sparse. In such
 983 applications, the construction of the corresponding quantum memory depends on the sparsity rather
 984 than the dimension. The potential acceleration in the gradient calculation and state preparation are
 985 left for future exploration.

986 To show that solving Equation (2) is a special case of solving Equation (10), let $Z = X/k$. Then, the
 987 constraint $\|X\|_{\text{tr}} \leq k$ becomes $\|Z\|_{\text{tr}} = \|X/k\|_{\text{tr}} = \|X\|_{\text{tr}}/k \leq 1$. Substituting into the objective
 988 function of Equation (2):

$$990 \sum_{(i,j) \in \Omega} (X_{i,j} - Y_{i,j})^2 = \sum_{(i,j) \in \Omega} (kZ_{i,j} - Y_{i,j})^2. \quad (17)$$

992 Define the function $f(Z) = \sum_{(i,j) \in \Omega} (kZ_{i,j} - Y_{i,j})^2$. Then, Equation (2) is equivalent to:

$$994 \min_{\|Z\|_{\text{tr}} \leq 1} f(Z). \quad (18)$$

996 This matches the form of Equation (10).

998 **Satisfaction of Assumption 2.** The trace norm $\|\cdot\|_{\text{tr}}$ is a convex function, so the set $\{Z : \|Z\|_{\text{tr}} \leq 1\}$ is convex.
 999 In the finite-dimensional space $\mathbb{R}^{d \times d}$, the set $\{Z : \|Z\|_{\text{tr}} \leq 1\}$ is closed (because the
 1000 trace norm is continuous) and bounded (since $\|Z\|_F \leq \|Z\|_{\text{tr}} \leq 1$), hence it is compact. For any
 1001 $Z_1, Z_2 \in D$, we have $\|Z_1\|_F \leq 1$ and $\|Z_2\|_F \leq 1$, so:

$$1002 \|Z_1 - Z_2\|_F \leq \|Z_1\|_F + \|Z_2\|_F \leq 2. \quad (19)$$

1003 Thus, the diameter $D \leq 2$. Therefore, Assumption 2 is satisfied.

1005 **Satisfaction of Assumption 1.** The function $f(Z) = \sum_{(i,j) \in \Omega} (kZ_{i,j} - Y_{i,j})^2$ is a sum of squares,
 1006 hence it is convex. For $(i, j) \in \Omega$, the partial derivative is $2(kZ_{i,j} - Y_{i,j})$; for $(i, j) \notin \Omega$, it is 0.
 1007 Therefore, the gradient $\nabla f(Z) = 2P_{\Omega}(kZ - Y)$, where P_{Ω} is the projection operator that preserves
 1008 elements in Ω and sets others to zero. For any Z_1, Z_2 ,

$$1009 \nabla f(Z_1) - \nabla f(Z_2) = 2P_{\Omega}(kZ_1 - kZ_2). \quad (20)$$

1011 Since P_{Ω} is a linear operator and does not increase the Frobenius norm, we have

$$1012 \|\nabla f(Z_1) - \nabla f(Z_2)\|_F = \|2P_{\Omega}(Z_1 - Z_2)\|_F \leq 2k\|Z_1 - Z_2\|_F. \quad (21)$$

1014 Thus, ∇f is Lipschitz continuous with constant $L = 2k$. Therefore, Assumption 1 is satisfied.

1015 In conclusion, we can apply the algorithms from Section 4 to solve the matrix completion problem.
 1016 Furthermore, since the gradient of the matrix completion problem is sparse (with only $|\Omega|$ non-zero
 1017 entries and zeros elsewhere), the construction of quantum memory depends solely on $|\Omega|$ rather
 1018 than the dimension d . Moreover, the computation of the update rule can be further accelerated by
 1019 leveraging quantum multiplication for sparse matrix. This aspect is left for future investigation.

1021 A.6 POTENTIAL APPLICATIONS

1023 Our proposed quantum Frank-Wolfe algorithms are applicable to a broad class of convex optimization
 1024 problems with structured constraints. This section elaborates on the applications of our algo-
 1025 rithms in three key domains: sparsity constraints in signal processing, zero-sum games in game
 theory, and semidefinite programming.

1026
 1027 **Signal Processing: Sparsity Constraints via ℓ_1 Norm.** In signal processing, a common problem is
 1028 recovering sparse signals from noisy observations, typically achieved through ℓ_1 norm regularization
 1029 to promote sparsity in the solution. Consider the basis pursuit denoising problem:
 1030

$$1031 \min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 \quad \text{subject to} \quad \|\mathbf{x}\|_1 \leq \tau, \quad (22)$$

1034 where $\mathbf{A} \in \mathbb{R}^{m \times d}$ is the measurement matrix, $\mathbf{b} \in \mathbb{R}^m$ is the observation vector, and $\tau > 0$ is the
 1035 constraint radius. The feasible domain $\mathcal{D} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_1 \leq \tau\}$ is an ℓ_1 -norm ball. As discussed
 1036 in Section 3, the core of the Frank-Wolfe update step under this constraint involves solving the linear
 1037 subproblem $\min_{\mathbf{s} \in \mathcal{D}} \langle \mathbf{s}, \nabla f(\mathbf{x}^{(t)}) \rangle$, whose exact solution is given by the coordinate with the largest
 1038 absolute gradient component (i.e., $\hat{\mathbf{s}} = -\tau \cdot \text{sign}(\nabla_i f(\mathbf{x}^{(t)})) \cdot \mathbf{e}_i$, where $i = \text{argmax}_j |\nabla_j f(\mathbf{x}^{(t)})|$).
 1039 Our quantum Frank-Wolfe algorithm (Theorem 1) can be used to reduce the per-iteration query
 1040 complexity from the classical $\mathcal{O}(d)$ to $\mathcal{O}(\sqrt{d})$.
 1041

1042 **Game Theory: Zero-Sum Games with Simplex Constraints.** In game theory, Nash equilibria
 1043 for two-player zero-sum games can be found by solving a linear programming problem over the
 1044 simplex. Consider a game with payoff matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$. The row player's mixed strategy is
 1045 a vector $\mathbf{x} \in \Delta_m$ (m -dimensional simplex), and the column player's mixed strategy is a vector
 1046 $\mathbf{y} \in \Delta_n$. The row player aims to minimize the expected loss $\mathbf{x}^\top \mathbf{A} \mathbf{y}$. Finding the Nash equilibrium
 1047 can be formulated as:
 1048

$$1049 \min_{\mathbf{x} \in \Delta_m} \max_{\mathbf{y} \in \Delta_n} \mathbf{x}^\top \mathbf{A} \mathbf{y}. \quad (23)$$

1052 Through linear programming duality or its variants, this problem can be transformed into an optimi-
 1053 zation problem over the simplex. The feasible domain is the simplex $\mathcal{D} = \Delta_d$. The solution to the
 1054 Frank-Wolfe linear subproblem under this constraint corresponds to the unit vector with the largest
 1055 gradient component (i.e., $\hat{\mathbf{s}} = \mathbf{e}_i$, where $i = \text{argmin}_j |\nabla_j f(\mathbf{x}^{(t)})|$). Our quantum Frank-Wolfe
 1056 algorithm (Theorem 2) similarly accelerates this step, achieving quantum speedup with respect to
 1057 dimension.
 1058

1059 **Semidefinite Programming.** Our quantum algorithms for computing top singular vectors have
 1060 potential applications in semidefinite programming (SDP). Many SDP solvers, particularly those
 1061 based on first-order methods, require repeatedly solving linear minimization oracles over the spec-
 1062 trahedron. The solution to this subproblem is given by the outer product of the eigenvector corre-
 1063 sponding to the smallest eigenvalue of a symmetric matrix A Nesterov (2007); d'Aspremont (2008);
 1064 Baes & Bürgisser (2009). Computing this vector is equivalent to finding the top eigenvector of the
 1065 shifted matrix $-A$. This computational bottleneck is structurally analogous to the top singular vector
 1066 extraction problem addressed by our quantum subroutines in Section 4. Therefore, our quantum
 1067 top singular vector extraction (QTSVE) and quantum power method (QPM) algorithms can be inte-
 1068 grated into SDP solvers to accelerate this subroutine, providing quantum speedup for a wide class
 1069 of SDP problems.
 1070

1071 B PROOF DETAIL

1072 B.1 PROOF OF LEMMA 3

1073 **Lemma 3.** *Given access to the quantum function value oracle \mathbf{U}_f , there exists a quantum circuit
 1074 to construct a quantum error bounded gradient oracle $\mathbf{U}_g : |i\rangle |\mathbf{x}\rangle |0\rangle \rightarrow |i\rangle |\mathbf{x}\rangle |g_i(\mathbf{x})\rangle$, where
 1075 $g_i(\mathbf{x}) = \frac{f(\mathbf{x} + \sigma \mathbf{e}_i) - f(\mathbf{x})}{\sigma}$ is the i -th component of the gradient and σ is the tunable parameter, with
 1076 two queries to the quantum function value oracle.*

1080 **Proof.** By choosing appropriate σ , we now construct a gradient unitary $U_g : |i\rangle|\mathbf{x}\rangle|0\rangle \rightarrow |i\rangle|\mathbf{x}\rangle|g_i(\mathbf{x})\rangle$ as follows:

$$1083 \quad |i\rangle|\mathbf{x}\rangle|0\rangle|0\rangle|0\rangle|0\rangle \quad (24)$$

$$1084 \quad \rightarrow |i\rangle|\mathbf{x}\rangle|\mathbf{x} + \sigma\mathbf{e}_i\rangle|0\rangle|0\rangle|0\rangle \quad (25)$$

$$1085 \quad \rightarrow |i\rangle|\mathbf{x}\rangle|\mathbf{x} + \sigma\mathbf{e}_i\rangle|f(\mathbf{x} + \sigma\mathbf{e}_i)\rangle|f(\mathbf{x})\rangle|0\rangle \quad (26)$$

$$1086 \quad \rightarrow |i\rangle|\mathbf{x}\rangle|\mathbf{x} + \sigma\mathbf{e}_i\rangle|f(\mathbf{x} + \sigma\mathbf{e}_i)\rangle|f(\mathbf{x})\rangle \left| \frac{f(\mathbf{x} + \sigma\mathbf{e}_i) - f(\mathbf{x})}{\sigma} \right\rangle \quad (26)$$

$$1087 \quad \rightarrow |i\rangle|\mathbf{x}\rangle \left| \frac{f(\mathbf{x} + \sigma\mathbf{e}_i) - f(\mathbf{x})}{\sigma} \right\rangle \quad (27)$$

$$1088 \quad = |i\rangle|\mathbf{x}\rangle|g_i(\mathbf{x})\rangle, \quad (28)$$

1089 where Equation (24) is by adding σ at the i -th entry of the third register, Equation (25) is by applying
1090 U_f based on the second and the third register, Equation (26) is by applying addition and division
1091 based on the fourth and the fifth register, Equation (27) is by uncomputing the third, fourth and fifth
1092 register. For the complexity, this U_g takes two queries of U_f and $O(1)$ elementary gates to get the
1093 approximate gradient. ■

1100 B.2 PROOF OF LEMMA 4

1101 **Lemma 4.** (Approximate maximum gradient component finding) *Given access to the quantum*
1102 *error bounded gradient oracle $U_g : |i\rangle|\mathbf{x}\rangle|0\rangle \rightarrow |i\rangle|\mathbf{x}\rangle|g_i(\mathbf{x})\rangle$ s.t. for each $i \in [d]$, after*
1103 *measuring $|g_i(\mathbf{x})\rangle$, the measured outcome $g_i(\mathbf{x})$ satisfies $|g_i(\mathbf{x}) - \nabla f_i(\mathbf{x})| \leq \epsilon$. There exists a*
1104 *quantum circuit \mathcal{A}_{\max} that finds the index i^* that satisfies $\nabla f_{i^*}(\mathbf{x}) \geq \max_{j \in [d]} \nabla f_j(\mathbf{x}) - 2\epsilon$ or*
1105 *$|\nabla f_{i^*}(\mathbf{x})| \geq \max_{j \in [d]} |\nabla f_j(\mathbf{x})| - 2\epsilon$, using $O(\sqrt{d} \log(\frac{1}{\delta}))$ applications of U_g , U_g^\dagger and $O(\sqrt{d})$*
1106 *elementary gates, with probability $1 - \delta$. For the non-uniform initial state, let p be the initial mea-*
1107 *surement probability of the maximum component, then the algorithm finds the maximum with query*
1108 *complexity of $O(\frac{1}{\sqrt{p}} \log(\frac{1}{\delta}))$.*

1109 **Proof.** We restate the quantum minimum finding algorithm here for reader benefits Durr & Hoyer
1110 (1996): Choose threshold index $0 \leq j \leq d - 1$ uniformly at random. Repeat the following and
1111 return j when the total running time is more than $22.5\sqrt{d} + 1.4 \log(d)$:

- 1114 1. Prepare the state $\sum_i^d |i\rangle|\mathbf{x}\rangle|g_i(\mathbf{x})\rangle|0\rangle$.
- 1115 2. Set the third register to $|1\rangle$ conditioned on the value of the second register smaller than $g_j(\mathbf{x})$
- 1116 3. Apply the quantum exponential Grover search algorithm for the third register being 1.
- 1117 4. Measure the first and the third registers in computation basis, if the measurement result
1118 of the third register is smaller than $g_j(\mathbf{x})$, set j to be the measurement result of the first
1119 register.

1120 By Theorem 1 of Durr & Hoyer (1996), the algorithm finds the minimum $g_i(\mathbf{x})$ with probability
1121 $1/2$, $O(\sqrt{d})$ applications of U_g , U_g^\dagger and $O(\sqrt{d})$ elementary gates. The probability can be boost to
1122 $1 - \delta$ with $O(\log(1/\delta))$ repeats and taking the minimum of the outputs.

1123 This algorithm can be modified into the quantum maximum absolute value finding algorithm by
1124 setting the third register to $|1\rangle$ conditioned on the value of the second register greater than $|g_j(\mathbf{x})|$ in
1125 Step 2, and set j to be the measurement result that is greater than $|g_j(\mathbf{x})|$ in Step 4.

1126 However, with the estimated error, the greatest estimated gradient component $g_{\max}(\mathbf{x})$ may not have
1127 the same index of $\nabla f_{\max}(\mathbf{x})$. As $|g_i(\mathbf{x}) - \nabla f_i(\mathbf{x})| \leq \epsilon$ for each i , in the worst case, there exists
1128 i such that $|g_i(\mathbf{x})| = |\nabla f_i(\mathbf{x})| + \epsilon \geq |g_{i^*}(\mathbf{x})| = \max_{j \in [d]} |\nabla f_j(\mathbf{x})| - \epsilon$, the maximum finding
1129 algorithm will give such $g_i(\mathbf{x})$ as outcome, which is greater than $\max_{j \in [d]} |\nabla f_j(\mathbf{x})| - 2\epsilon$.

1134 Similarly, As $|g_i(\mathbf{x}) - \nabla f_i(\mathbf{x})| \leq \epsilon$ for each i , in the worst case, there exists i such that
 1135 $|g_i(\mathbf{x})| = |\nabla f_i(\mathbf{x})| - \epsilon \leq |g_{i^*}(\mathbf{x})| = \min_{j \in [d]} |\nabla f_j(\mathbf{x})| + \epsilon$, the minimum finding algorithm
 1136 will give such $g_i(\mathbf{x})$ as outcome, which is less than $\min_{j \in [d]} |\nabla f_j(\mathbf{x})| + 2\epsilon$. Similar proof processes
 1137 can be employed to derive the error bounds for the minimum/maximum search.

1138 Note that in the matrix case of this work, the state prepared to apply quantum maximum finding
 1139 is not a uniform superposition, but the algorithm in Durr & Hoyer (1996) is only for the uniform
 1140 superposition input. For the non-uniform input, in the third step, the Grover operator should be
 1141 replaced with the amplitude amplification operator. We now prove the complexity of the algorithm
 1142 for the non-uniform initial state. For the analysis of the probability of success, assume that there
 1143 is no time-out, that is, the algorithm runs long enough to find the minimum. Then we analyze the
 1144 probability that an element of a given rank becomes the threshold during the algorithm (Lemma 12)
 1145 and then bound the expected number of iterations (Lemma 13), which extend Lemma 1 and 2 in
 1146 Durr & Hoyer (1996).

1147 Then, by Lemma 13, the expected running time of finding the maximum is $O\left(\frac{1}{\sqrt{p_1}}\right)$. By Markov's
 1148 inequality, after running the algorithm for twice the expected time, the probability of success is at
 1149 least $1/2$. The probability can be boost to $1 - \delta$ with $O(\log(1/\delta))$ repeats and taking the maximum
 1150 of the outputs. This extends the Dürr-Høyer minimum finding algorithm to the weighted case and
 1151 provides a complexity analysis tailored to singular value distributions for the matrix case of this
 1152 work. ■
 1153

1154
 1155

1156 **Lemma 12** (Probability of Selecting Threshold of Rank r). *Let $p(t, r)$ be the probability that the
 1157 element of rank r (where rank 1 is the maximum) will ever be chosen when the infinite algorithm is
 1158 searching among t elements. Then, for $r \leq t$, $p(t, r) = P_r = \frac{p_r}{\sum_{j=1}^{t+1} p_j}$, and for $r > t$, $p(t, r) = 0$.*
 1159

1160

1161 **Proof.** The case $r > t$ is trivial. For $r \leq t$, we proceed by induction on t for fixed r .
 1162

1163 **Base step:** When $t = r$, the algorithm starts by measuring the initial state, which yields the element
 1164 of rank r with probability P_r . Since the relative amplitudes of the basis states constituting the
 1165 marked state remain invariant throughout the amplification process, the probability of selecting rank
 1166 r as the threshold is exactly P_r .

1167 **Inductive step:** Assume that for all $k \in [r, t]$, $p(k, r) = P_r$. Now consider $t + 1$ elements. The
 1168 initial threshold is chosen with probability p_r for rank r . If the initial threshold has rank greater than
 1169 r , then the algorithm will update the threshold only if it finds an element with rank between r and
 1170 the current threshold. By the induction hypothesis, the probability that rank r is eventually selected
 1171 when starting from a threshold of rank k (where $r < k \leq t + 1$) is $p(k - 1, r) = P_r$. Therefore,

$$\begin{aligned} p(t + 1, r) &= \frac{p_r}{\sum_{j=1}^{t+1} p_j} + \sum_{k=r+1}^{t+1} \frac{p_k}{\sum_{j=1}^{t+1} p_j} \cdot p(k - 1, r) \\ &= \frac{1}{\sum_{j=1}^{t+1} p_j} \left(p_r + \sum_{k=r+1}^{t+1} p_k \cdot p(k - 1, r) \right). \end{aligned} \quad (29)$$

1178 By the inductive hypothesis,

$$p(t + 1, r) = \frac{1}{\sum_{j=1}^{t+1} p_j} \left(p_r + \sum_{k=r+1}^{t+1} p_k \cdot P_r \right). \quad (30)$$

1184 Substitute P_r into the equation, we have

$$p(t + 1, r) = \frac{1}{\sum_{j=1}^{t+1} p_j} \left(p_r + \sum_{k=r+1}^{t+1} p_k \cdot \frac{p_r}{\sum_{j=1}^r p_j} \right). \quad (31)$$

1188 Then, after some simple equivalent transformations, we have
 1189

$$\begin{aligned}
 1190 \quad p(t+1, r) &= \frac{p_r}{\sum_{j=1}^{t+1} p_j} \left(1 + \frac{\sum_{k=r+1}^{t+1} p_k}{\sum_{j=1}^r p_j} \right) = \frac{p_r}{\sum_{j=1}^{t+1} p_j} \left(\frac{\sum_{j=1}^r p_j + \sum_{k=r+1}^{t+1} p_k}{\sum_{j=1}^r p_j} \right) \\
 1191 \\
 1192 \quad &= \frac{p_r}{\sum_{j=1}^{t+1} p_j} \frac{\sum_{j=1}^{t+1} p_j}{\sum_{j=1}^r p_j} = \frac{p_r}{\sum_{j=1}^r p_j} = P_r
 \end{aligned} \tag{32}$$

1193 This completes the induction. Therefore, the lemma follows. \blacksquare
 1194

1195 **Lemma 13** (Expected Running Time). *The expected number of iterations of the quantum maximum
 1196 finding algorithm for non-uniform initial state is $O\left(\frac{1}{\sqrt{p_1}}\right)$.*

1200 **Proof.** Let E be the expected number of iterations to find the maximum (rank 1). By Lemma 12,
 1201 the probability that the initial threshold has rank r is $P_r = \frac{p_r}{\sum_{j=1}^r p_j}$. When the current threshold
 1202 has rank r , the quantum search algorithm finds a better element (with rank less than r) in expected
 1203 $O(1/\sqrt{S_{r-1}})$ iterations, where $S_{r-1} = \sum_{j=1}^{r-1} p_j$.

1204 Since the threshold rank decreases monotonically, each rank r is visited as a threshold at most once,
 1205 with probability P_r . Thus,

$$E = \sum_{r=1}^N P_r \cdot O\left(\frac{1}{\sqrt{S_{r-1}}}\right) = O\left(\sum_{r=2}^N \frac{p_r}{S_r} \frac{1}{\sqrt{S_{r-1}}}\right), \tag{33}$$

1206 where $S_r = \sum_{j=1}^r p_j$, and for $r = 1$, $S_0 = 0$ and the search time is 0. We have
 1207

$$\sum_{r=2}^N \frac{p_r}{S_r \sqrt{S_{r-1}}} \leq \sum_{r=2}^N \int_{S_{r-1}}^{S_r} x^{-3/2} dx = \int_{p_1}^1 x^{-3/2} dx = 2\left(\frac{1}{\sqrt{p_1}} - 1\right) = O\left(\frac{1}{\sqrt{p_1}}\right), \tag{34}$$

1208 where the first inequality holds because
 1209

$$\begin{aligned}
 1210 \quad \left(1 - \sqrt{\frac{S_{r-1}}{S_r}}\right)^2 \geq 0 &\implies 1 - \frac{S_{r-1}}{S_r} \leq 2\left(1 - \sqrt{\frac{S_{r-1}}{S_r}}\right) \\
 1211 \\
 1212 \quad &\implies \frac{p_r}{S_r \sqrt{S_{r-1}}} \leq 2\left(\frac{1}{\sqrt{S_{r-1}}} - \frac{1}{\sqrt{S_r}}\right) = \int_{S_{r-1}}^{S_r} x^{-3/2} dx.
 \end{aligned} \tag{35}$$

1213 Therefore, $E = O(1/\sqrt{p_1})$, which gives the lemma. \blacksquare
 1214

1215 B.3 PROOF OF THEOREM 1

1216 **Theorem 1.** (Quantum FW over the sparsity constraint) *By setting $\sigma_t = \frac{C_f}{\sqrt{d}L(t+2)}$ for $t \in [T]$,
 1217 the quantum algorithm (Algorithm 2) solves the sparsity constraint optimization problem for any
 1218 precision ε such that $f(\mathbf{x}^T) - f(\mathbf{x}^*) \leq \varepsilon$ in $T = \frac{4C_f}{\varepsilon} - 2$ rounds, succeed with probability $1 - p$,
 1219 with $O\left(\sqrt{d} \log \frac{C_f}{p\varepsilon}\right)$ calls to the function value oracle \mathbf{U}_f per round.*

1220 **Proof.** By Lemma 2 and the inequality between ℓ_2 norm and ℓ_∞ norm, we have
 1221

$$|g_i(\mathbf{x}) - \nabla f_i(\mathbf{x})| \leq \|g(\mathbf{x}) - \nabla f(\mathbf{x})\|_\infty \leq \|g(\mathbf{x}) - \nabla f(\mathbf{x})\|_2 \leq \frac{\sqrt{d}L\sigma}{2}. \tag{36}$$

1222 By Lemma 4, after the quantum approximate maximum absolute value finding, we have an estimated
 1223 maximum gradient component which satisfied
 1224

$$|\nabla f_{i^*}(\mathbf{x})| \geq \max_{j \in [d]} |\nabla f_j(\mathbf{x})| - \sqrt{d}L\sigma \tag{37}$$

1242 Set $s = -e_{i^*}$, we have
 1243

$$\begin{aligned}
 1244 \quad \langle s, \nabla f(\mathbf{x}^{(t)}) \rangle &= -|\nabla f_{i^*}(\mathbf{x}^{(t)})| \\
 1245 &\leq -\max_{j \in [d]} |\nabla f_j(\mathbf{x}^{(t)})| + \sqrt{d}L\sigma_t \\
 1246 &= -\langle e_{\arg\max_{i \in [d]} |\nabla_i f(\mathbf{x}^{(t)})|}, \nabla f(\mathbf{x}^{(t)}) \rangle + \sqrt{d}L\sigma_t \\
 1247 &= \min_{\hat{s} \in \mathcal{D}} \langle \hat{s}, \nabla f(\mathbf{x}^{(t)}) \rangle + \sqrt{d}L\sigma_t. \tag{38}
 \end{aligned}$$

1251 By the update rule and the definition of the curvature, we have
 1252

$$f(\mathbf{x}^{(t+1)}) = f((1 - \gamma_t)\mathbf{x}^{(t)} + \gamma_t s) \leq f(\mathbf{x}^{(t)}) + \gamma_t \langle s - \mathbf{x}^{(t)}, \nabla f(\mathbf{x}^{(t)}) \rangle + \frac{\gamma_t^2}{2} C_f \tag{39}$$

1256 Combining Inequality 38 and 39, we have
 1257

$$f(\mathbf{x}^{(t+1)}) \leq f(\mathbf{x}^{(t)}) + \gamma_t \left(\min_{\hat{s} \in \mathcal{D}} \langle \hat{s}, \nabla f(\mathbf{x}) \rangle - \langle \mathbf{x}^{(t)}, \nabla f(\mathbf{x}^{(t)}) \rangle \right) + \sqrt{d}\gamma_t L\sigma_t + \frac{\gamma_t^2}{2} C_f. \tag{40}$$

1261 Let $h(\mathbf{x}^{(t)}) := f(\mathbf{x}^{(t)}) - f(x^*)$, we have
 1262

$$\begin{aligned}
 1263 \quad h(\mathbf{x}^{(t+1)}) &\leq h(\mathbf{x}^{(t)}) + \gamma_t \left(\min_{\hat{s} \in \mathcal{D}} \langle \hat{s}, \nabla f(\mathbf{x}) \rangle - \langle \mathbf{x}^{(t)}, \nabla f(\mathbf{x}^{(t)}) \rangle \right) + \sqrt{d}\gamma_t L\sigma_t + \frac{\gamma_t^2}{2} C_f \\
 1264 &\leq h(\mathbf{x}^{(t)}) - \gamma_t h(\mathbf{x}^{(t)}) + \sqrt{d}\gamma_t L\sigma_t + \frac{\gamma_t^2}{2} C_f \\
 1265 &= (1 - \gamma_t)h(\mathbf{x}^{(t)}) + \sqrt{d}\gamma_t L\sigma_t + \frac{\gamma_t^2}{2} C_f. \tag{41}
 \end{aligned}$$

1270 Set $\gamma_t = \frac{2}{t+2}$, $\sigma_t = \frac{\gamma_t C_f}{2\sqrt{d}L}$, we have
 1271

$$h(\mathbf{x}^{(t+1)}) \leq \left(1 - \frac{2}{t+2}\right) h(\mathbf{x}^{(t)}) + \left(\frac{2}{t+2}\right)^2 C_f. \tag{42}$$

1275 Using a similar induction as shown in Jaggi (2013) over t , we have
 1276

$$h(\mathbf{x}^{(t)}) \leq \frac{4C_f}{t+2}. \tag{43}$$

1279 We will restate this induction in Lemma 14 for reader benefit.
 1280

1281 Thus, set $\gamma_t = \frac{2}{t+2}$, $\sigma_t = \frac{C_f}{\sqrt{d}L(t+2)}$ for all $t \in [T]$, after $T = \frac{4C_f}{\varepsilon} - 2$ rounds, we have
 1282

$$f(\mathbf{x}^{(T)}) - f(x^*) \leq \varepsilon, \tag{44}$$

1284 for any $\varepsilon > 0$.
 1285

1286 In each round, by Lemma 3, two queries to the quantum function value oracle are needed to
 1287 construct the quantum gradient oracle. Then by lemma 4, $O(\sqrt{d} \log \frac{1}{\delta})$ queries to the quantum gradient
 1288 oracle are needed to find the index of the estimated maximum gradient component with successful
 1289 probability of $1 - \delta$. Since each maximum finding succeeds with probability $1 - \delta$, the probability
 1290 that all T iterations succeed is at least $1 - T\delta$. By setting $\delta = p/T$, we ensure an overall success
 1291 probability of at least $1 - p$. Therefore, $O\left(\sqrt{d} \log \frac{C_f}{p\varepsilon}\right)$ queries to the quantum function value oracle
 1292 are needed in each iteration. Then the theorem follows. ■
 1293

1294 We restate the proof of the induction we use in Theorem 1 for reader benefit.
 1295

1296

Lemma 14. (Jaggi (2013)) If for any $t \in [N]$,

1297

1298

1299

1300

then

1301

1302

1303

$$h(\mathbf{x}^{(t+1)}) \leq \left(1 - \frac{2}{t+2}\right) h(\mathbf{x}^{(t)}) + \left(\frac{2}{t+2}\right)^2 C_f, \quad (45)$$

Proof. For $t = 0$, we have

1304

1305

1306

1307

Assume that $h(\mathbf{x}^{(t)}) \leq \frac{4C_f}{t+2}$, we have

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

$$h(\mathbf{x}^{(1)}) \leq \left(1 - \frac{2}{0+2}\right) h(\mathbf{x}^{(0)}) + \left(\frac{2}{0+2}\right)^2 C_f = C_f. \quad (47)$$

$$\begin{aligned} h(\mathbf{x}^{(t+1)}) &\leq \left(1 - \frac{2}{t+2}\right) h(\mathbf{x}^{(t)}) + \left(\frac{2}{t+2}\right)^2 C_f \\ &\leq \left(1 - \frac{2}{t+2}\right) \frac{4C_f}{t+2} + \left(\frac{2}{t+2}\right)^2 C_f \\ &= \left(1 - \frac{1}{t+2}\right) \frac{4C_f}{t+2} + \left(\frac{2}{t+2}\right)^2 C_f \\ &= \frac{t+1}{t+2} \frac{4C_f}{t+2} \leq \frac{t+2}{t+3} \frac{4C_f}{t+2} = \frac{4C_f}{t+3}, \end{aligned} \quad (48)$$

which gives the lemma. ■

1318

1319

1320

B.4 PROOF OF LEMMA 10

1321

The framework of quantum gradient estimator originates from Jordan quantum gradient estimation method Jordan (2005), but Jordan algorithm did not give any error bound because the analysis of it was given by omitting the high-order terms of Taylor expansion of the function directly. In 2019, the quantum gradient estimation method with error analysis was given in Gilyén et al. (2019), and was applied to the general convex optimization problem van Apeldoorn et al. (2020); Chakrabarti et al. (2020). In those case, however, $O(\log n)$ repetitions were needed to estimate the gradient within an acceptable error. The query complexity was then improved to $O(1)$ in He et al. (2022; 2024). Here we use the version of He et al. (2024) (Algorithm 7).

1328

Algorithm 7 Bounded-error Jordan quantum gradient estimation He et al. (2024)

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

1: **Input:** point x , parameters r, ρ, ϵ .2: **Output:** $g(x)$ 3: Prepare the initial state: d b -qubit registers $|0^{\otimes b}, 0^{\otimes b}, \dots, 0^{\otimes b}\rangle$ where $b = \log_2 \frac{G\rho}{4\pi d^2 \beta r}$.Prepare 1 c -qubit register $|0^{\otimes c}\rangle$ where $c = \log_2 \frac{16\pi d}{\rho} - 1$. And prepare $|y_0\rangle = \frac{1}{\sqrt{2^d}} \sum_{a \in \{0,1,\dots,2^d-1\}} e^{\frac{2\pi i a}{2^d}} |a\rangle$.4: Apply Hadamard transform to the first d registers.5: Perform the quantum query oracle Q_F to the first $d+1$ registers, where $F(u) = \frac{2^b}{2Gr} \left[f\left(x + \frac{r}{2^b} \left(u - \frac{2^b}{2} \mathbb{1}\right)\right) - f(x) \right]$, and the result is stored in the $(d+1)$ th register.6: Perform the addition modulo 2^c operation to the last two registers.7: Apply the inverse evaluating oracle Q_F^{-1} to the first $d+1$ registers.8: Perform quantum inverse Fourier transformations to the first d registers separately.9: Measure the first d registers in computation bases respectively to get m_1, m_2, \dots, m_n .10: $g(x) = \tilde{\nabla} f(x) = \frac{2G}{2^b} \left(m_1 - \frac{2^b}{2}, m_2 - \frac{2^b}{2}, \dots, m_n - \frac{2^b}{2} \right)^T$.

1350
 1351 **Lemma 10.** (Lemma 1 He et al. (2024)) If f is G -Lipschitz continuous and L -smooth convex function
 1352 and can be accessed by a quantum function value oracle, then there exists an quantum algorithm
 1353 that for any $r > 0$ and $1 \geq \rho > 0$, gives the estimated gradient $g(x)$, which satisfies

$$1354 \quad \Pr[\|g(x) - \nabla f(x)\|_\infty > 8\pi n^2(n/\rho + 1)Lr/\rho] < \rho, \quad (11)$$

1355 using $O(1)$ applications of U_f and $O(d \log d)$ elementary gates. The space complexity is
 1356 $O\left(d \log \frac{G\rho}{4\pi d^2 Lr}\right)$.
 1357

1358 **Proof.** The primary additional gate overhead originates from the quantum Fourier transformation
 1359 (QFT). Each QFT requires $O(\log d)$ elementary gates, and for d such operations, the total additional
 1360 elementary gate overhead is $O(d \log d)$. Consequently, the additional elementary gate overhead is
 1361 $O(d \log d)$.
 1362

1363 The states after Step 3 will be:

$$1364 \quad \frac{1}{\sqrt{2^n}} \sum_{a \in \{0,1,\dots,2^n-1\}} e^{\frac{2\pi i a}{2^n}} |0^{\otimes b}, 0^{\otimes b}, \dots, 0^{\otimes b}\rangle |0^{\otimes c}\rangle |a\rangle. \quad (49)$$

1367 After Step 4:

$$1368 \quad \frac{1}{\sqrt{2^{bn+c}}} \sum_{u_1, u_2, \dots, u_n \in \{0,1,\dots,2^b-1\}} \sum_{a \in \{0,1,\dots,2^c-1\}} e^{\frac{2\pi i a}{2^n}} |u_1, u_2, \dots, u_n\rangle |0^{\otimes c}\rangle |a\rangle. \quad (50)$$

1371 After Step 5:

$$1373 \quad \frac{1}{\sqrt{2^{bn+c}}} \sum_{u_1, u_2, \dots, u_n \in \{0,1,\dots,2^b-1\}} \sum_{a \in \{0,1,\dots,2^c-1\}} e^{\frac{2\pi i a}{2^n}} |u_1, u_2, \dots, u_n\rangle |F(u)\rangle |a\rangle. \quad (51)$$

1375 After Step 6:

$$1377 \quad \frac{1}{\sqrt{2^{bn+c}}} \sum_{u_1, u_2, \dots, u_n \in \{0,1,\dots,2^b-1\}} \sum_{a \in \{0,1,\dots,2^c-1\}} e^{2\pi i F(u)} e^{\frac{2\pi i a}{2^n}} |u_1, u_2, \dots, u_n\rangle |F(u)\rangle |a\rangle. \quad (52)$$

1380 After Step 7:

$$1381 \quad \frac{1}{\sqrt{2^{bn+c}}} \sum_{u_1, u_2, \dots, u_n \in \{0,1,\dots,2^b-1\}} \sum_{a \in \{0,1,\dots,2^c-1\}} e^{2\pi i F(u)} e^{\frac{2\pi i a}{2^n}} |u_1, u_2, \dots, u_n\rangle |0^{\otimes c}\rangle |a\rangle. \quad (53)$$

1384 In the following, the last two registers will be omitted:

$$1385 \quad \frac{1}{\sqrt{2^{bn}}} \sum_{u_1, u_2, \dots, u_n \in \{0,1,\dots,2^b-1\}} e^{2\pi i F(u)} |u_1, u_2, \dots, u_n\rangle. \quad (54)$$

1388 And then we simply relabel the state by changing $u \rightarrow v = u - \frac{2^b}{2}$:

$$1390 \quad \frac{1}{\sqrt{2^{bn}}} \sum_{v_1, v_2, \dots, v_n \in \{-2^{b-1}, -2^{b-1}+1, \dots, 2^{b-1}\}} e^{2\pi i F(v)} |v\rangle. \quad (55)$$

1393 We denote Formula (55) as $|\phi\rangle$. Let $g = \nabla f(x)$, and consider the idealized state

$$1394 \quad |\psi\rangle = \frac{1}{\sqrt{2^{bn}}} \sum_{v_1, v_2, \dots, v_n \in \{-2^{b-1}, -2^{b-1}+1, \dots, 2^{b-1}\}} e^{\frac{2\pi i g \cdot v}{2G}} |v\rangle. \quad (56)$$

1397 After Step 9, from the analysis of phase estimation Brassard et al. (2002):

$$1399 \quad \Pr\left[\left|\frac{Ng_i}{2G} - m_i\right| > e\right] < \frac{1}{2(e-1)}, \forall i \in [n]. \quad (57)$$

1401 Let $e = n/\rho + 1$, where $1 \geq \rho > 0$. We have

$$1403 \quad \Pr\left[\left|\frac{Ng_i}{2G} - m_i\right| > n/\rho + 1\right] < \frac{\rho}{2n}, \forall i \in [n]. \quad (58)$$

Note that the difference in the probabilities of measurement on $|\phi\rangle$ and $|\psi\rangle$ can be bounded by the trace distance between the two density matrices:

$$\begin{aligned} F(v) &\leq \frac{2^b}{2Gr} [f(x + \frac{rv}{N}) - f(x)] + \frac{1}{2^{c+1}} \\ &\leq \frac{2^b}{2Gr} [\frac{r}{2^b} g \cdot v + \frac{L(rv)^2}{2^{2b}}] + \frac{1}{2^{c+1}} \\ &\leq \frac{g \cdot v}{2G} + \frac{2^b Lrn}{4G} + \frac{1}{2^{c+1}}. \end{aligned} \quad (60)$$

Then,

$$\begin{aligned} \|\phi\rangle - |\psi\rangle\|^2 &= \frac{1}{2^{bn}} \sum_v |e^{2\pi i F(v)} - e^{\frac{2\pi i g \cdot v}{2G}}|^2 \\ &\leq \frac{1}{2^{bn}} \sum_v |2\pi i F(v) - \frac{2\pi i g \cdot v}{2G}|^2 \\ &\leq \frac{1}{2^{bn}} \sum_v 4\pi^2 (\frac{2^b Lrn}{4G} + \frac{1}{2^{c+1}})^2. \end{aligned} \quad (61)$$

Set $b = \log_2 \frac{G\rho}{4\pi n^2 Lr}$, $c = \log_2 \frac{4G}{2^b n Lr} - 1$. We have

$$\|\phi\rangle - |\psi\rangle\|^2 \leq \frac{\rho^2}{16n^2}, \quad (62)$$

which implies $\|\phi\rangle\langle\phi| - |\psi\rangle\langle\psi\|_1 \leq \frac{\rho}{2n}$. Therefore, by the union bound,

$$\Pr \left[\left| \frac{2^b g_i}{2G} - m_i \right| > n/\rho + 1 \right] < \frac{\rho}{n}, \forall i \in [n]. \quad (63)$$

Furthermore, there is

$$\Pr \left[\left| g_i - \tilde{\nabla}_i f(x) \right| > \frac{2G(n/\rho + 1)}{2^b} \right] < \frac{\rho}{n}, \forall i \in [n], \quad (64)$$

as $b = \log_2 \frac{G\rho}{4\pi n^2 Lr}$, we have

$$\Pr \left[\left| g_i - \tilde{\nabla}_i f(x) \right| > 8\pi n^2 (n/\rho + 1) Lr/\rho \right] < \frac{\rho}{n}, \forall i \in [n]. \quad (65)$$

By the union bound, we have

$$\Pr \left[\|g - \tilde{\nabla} f(x)\|_\infty > 8\pi n^2 (n/\rho + 1) Lr/\rho \right] < \rho, \quad (66)$$

which gives the lemma. ■

B.5 PROOF OF THEOREM 5

Theorem 5. (Quantum FW with bounded-error Jordan algorithm) By setting $r_t = \frac{\rho C_f}{16\pi d^2(d/\rho+1)L(t+2)}$ for $t \in [T]$, the quantum algorithm (Algorithm 5) solves the sparsity constraint optimization problem for any precision ε such that $f(\mathbf{x}^T) - f(\mathbf{x}^*) \leq \varepsilon$ in $T = \frac{4C_f}{\varepsilon} - 2$ rounds, with $O(1)$ calls to the function value oracle \mathbf{U}_f per round.

Proof. By Lemma 10, with probability greater than ρ , we have

$$|g_i(\mathbf{x}) - \nabla f_i(\mathbf{x})| \leq \|g(\mathbf{x}) - \nabla f(\mathbf{x})\|_\infty \leq 8\pi d^2 (d/\rho + 1) Lr/\rho. \quad (67)$$

1458 Then the maximum component's coordinate of the estimated gradient $i^* = \operatorname{argmax}_{i \in [d]} |g_i(\mathbf{x}^{(t)})|$
 1459 satisfies

$$1460 \quad |\nabla f_{i^*}(\mathbf{x})| \geq \max_{j \in [d]} |\nabla f_j(\mathbf{x})| - 16\pi d^2(d/\rho + 1)Lr/\rho \quad (68)$$

1463 Set $\mathbf{s} = -\mathbf{e}_{i^*}$, we have

$$\begin{aligned} 1464 \quad \langle \mathbf{s}, \nabla f(\mathbf{x}^{(t)}) \rangle &= -|\nabla f_{i^*}(\mathbf{x}^{(t)})| \\ 1465 \quad &\leq -\max_{j \in [d]} |\nabla f_j(\mathbf{x}^{(t)})| + 16\pi d^2(d/\rho + 1)Lr/\rho \\ 1466 \quad &= -\langle \mathbf{e}_{\operatorname{argmax}_{i \in [d]} |\nabla_i f(\mathbf{x}^{(t)})|}, \nabla f(\mathbf{x}^{(t)}) \rangle + 16\pi d^2(d/\rho + 1)Lr/\rho \\ 1467 \quad &= \min_{\hat{\mathbf{s}} \in \mathcal{D}} \langle \hat{\mathbf{s}}, \nabla f(\mathbf{x}^{(t)}) \rangle + 16\pi d^2(d/\rho + 1)Lr/\rho. \end{aligned} \quad (69)$$

1472 By the update rule and the definition of the curvature, we have

$$1474 \quad f(\mathbf{x}^{(t+1)}) = f((1 - \gamma_t)\mathbf{x}^{(t)} + \gamma_t \mathbf{s}) \leq f(\mathbf{x}^{(t)}) + \gamma_t \langle \mathbf{s} - \mathbf{x}^{(t)}, \nabla f(\mathbf{x}^{(t)}) \rangle + \frac{\gamma_t^2}{2} C_f \quad (70)$$

1477 Combining Inequality 69 and 70, we have

$$1479 \quad f(\mathbf{x}^{(t+1)}) \leq f(\mathbf{x}^{(t)}) + \gamma_t \left(\min_{\hat{\mathbf{s}} \in \mathcal{D}} \langle \hat{\mathbf{s}}, \nabla f(\mathbf{x}) \rangle - \langle \mathbf{x}^{(t)}, \nabla f(\mathbf{x}^{(t)}) \rangle \right) + 16\pi d^2(d/\rho + 1)L\gamma_t r/\rho + \frac{\gamma_t^2}{2} C_f. \quad (71)$$

1482 Let $h(\mathbf{x}^{(t)}) := f(\mathbf{x}^{(t)}) - f(x^*)$, we have

$$\begin{aligned} 1484 \quad h(\mathbf{x}^{(t+1)}) &\leq h(\mathbf{x}^{(t)}) + \gamma_t \left(\min_{\hat{\mathbf{s}} \in \mathcal{D}} \langle \hat{\mathbf{s}}, \nabla f(\mathbf{x}) \rangle - \langle \mathbf{x}^{(t)}, \nabla f(\mathbf{x}^{(t)}) \rangle \right) + 16\pi d^2(d/\rho + 1)L\gamma_t r/\rho + \frac{\gamma_t^2}{2} C_f \\ 1485 \quad &\leq h(\mathbf{x}^{(t)}) - \gamma_t h(\mathbf{x}^{(t)}) + 16\pi d^2(d/\rho + 1)L\gamma_t r/\rho + \frac{\gamma_t^2}{2} C_f \\ 1486 \quad &= (1 - \gamma_t)h(\mathbf{x}^{(t)}) + 16\pi d^2(d/\rho + 1)L\gamma_t r/\rho + \frac{\gamma_t^2}{2} C_f. \end{aligned} \quad (72)$$

1491 Set $\gamma_t = \frac{2}{t+2}$, $r_t = \frac{\rho\gamma_t C_f}{32\pi d^2(d/\rho+1)L}$, we have

$$1494 \quad h(\mathbf{x}^{(t+1)}) \leq \left(1 - \frac{2}{t+2}\right) h(\mathbf{x}^{(t)}) + \left(\frac{2}{t+2}\right)^2 C_f. \quad (73)$$

1497 Using a similar induction as shown in Jaggi (2013) over t , we have

$$1499 \quad h(\mathbf{x}^{(t)}) \leq \frac{4C_f}{t+2}. \quad (74)$$

1502 Thus, set $\gamma_t = \frac{2}{t+2}$, $r_t = \frac{\rho C_f}{16\pi d^2(d/\rho+1)L(t+2)}$ for all $t \in [T]$, after $T = \frac{4C_f}{\varepsilon} - 2$ rounds, we have

$$1504 \quad f(\mathbf{x}^{(T)}) - f(x^*) \leq \varepsilon, \quad (75)$$

1506 for any $\varepsilon > 0$.

1508 In each round, by Lemma 10, $O(1)$ queries to the quantum function value oracle are needed to get
 1509 the estimated gradient vector. Subsequent steps no longer require queries to the oracle. Therefore,
 1510 in each round, $O(1)$ queries to the quantum function value oracle are needed. Then the theorem
 1511 follows. ■

1512 B.6 PROOF OF LEMMA 11
1513

1514 **Lemma 11.** [Quantum FW update over latent group norm ball] Let $\|\cdot\|_{\mathcal{G}}$ be a latent group norm
1515 corresponding to $\mathcal{G} = \{\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_{|\mathcal{G}|}\}$, and let $|\mathbf{g}|_{\max} = \max_j |\mathbf{g}_j|$. Then, there exists a quantum
1516 algorithm computing the Frank-Wolfe update $s^* := \operatorname{argmax}_{\hat{s} \in \|\cdot\|_{\mathcal{G}}\text{-Ball}} \langle \hat{s}^\top g(\mathbf{x}) \rangle$ in $O(\sqrt{|\mathcal{G}|} |\mathbf{g}|_{\max})$
1517 calls to U_f .
1518

1519 **Proof.** Assume that all $\|\cdot\|_{\mathbf{g}}$ are ℓ_p -norms, i.e. $\|\cdot\|_{\mathbf{g}_i} = \|\cdot\|_{p_i}$ for some $(p_i \in [1, \infty])$, and have
1520 quantum access to each $\mathbf{g}_i = \{\mathbf{g}_{i,1}, \mathbf{g}_{i,2}, \dots, \mathbf{g}_{i,|\mathbf{g}_i|}\} \subseteq [d]$ that load \mathbf{g}_i into quantum registers via
1521

$$1523 U_{\mathcal{G}} |i\rangle_A |0\rangle \rightarrow |i\rangle_A |\mathbf{g}_{i,1}\rangle |\mathbf{g}_{i,2}\rangle \dots |\mathbf{g}_{i,|\mathbf{g}_i|}\rangle \quad (76)$$

1524 where A is a log $|\mathcal{G}|$ qubit register. For each $|\mathbf{g}_{i,j}\rangle$ one can compute an approximation $|g_{\mathbf{g}_{i,j}}(\mathbf{x})\rangle$ to
1525 the $\mathbf{g}_{i,j}$ -th component of the gradient at \mathbf{x} by the method in Sec. 3.1.
1526

1527 Noting that $\operatorname{max}_{\hat{s} \in \|\cdot\|_p\text{-Ball}} \hat{s}^\top \mathbf{y} := \|\mathbf{y}\|_p^*$ and that
1528

$$1529 \mathbf{s}^* := \operatorname{argmax}_{\hat{s} \in \|\cdot\|_p\text{-ball}} \hat{s}^\top \mathbf{y} \quad (77)$$

1531 has components
1532

$$1533 \mathbf{s}_i^* \propto \operatorname{sgn}(\mathbf{y}_i) |\mathbf{y}_i|^{q-1} \quad (78)$$

1535 where $\frac{1}{p} + \frac{1}{q} = 1$, one can compute
1536

$$1537 \begin{aligned} & |i\rangle_A \bigotimes_{j=1}^{|\mathbf{g}_i|} |\mathbf{g}_{i,j}\rangle |0\rangle |0\rangle |0\rangle |0\rangle |0\rangle \\ & \rightarrow |i\rangle_A \bigotimes_{j=1}^{|\mathbf{g}_i|} |\mathbf{g}_{i,j}\rangle |g_{\mathbf{g}_{i,j}}(\mathbf{x})\rangle |0\rangle |0\rangle |0\rangle |0\rangle \\ & \rightarrow |i\rangle_A \bigotimes_{j=1}^{|\mathbf{g}_i|} |\mathbf{g}_{i,j}\rangle |g_{\mathbf{g}_{i,j}}(\mathbf{x})\rangle \left| \operatorname{sgn}(g_{\mathbf{g}_{i,j}}(\mathbf{x})) |g_{\mathbf{g}_{i,j}}(\mathbf{x})|^{q_i-1} \right\rangle |0\rangle |0\rangle |0\rangle \\ & \rightarrow |i\rangle_A \left(\bigotimes_{j=1}^{|\mathbf{g}_i|} |\mathbf{g}_{i,j}\rangle |g_{\mathbf{g}_{i,j}}(\mathbf{x})\rangle \left| \operatorname{sgn}(g_{\mathbf{g}_{i,j}}(\mathbf{x})) |g_{\mathbf{g}_{i,j}}(\mathbf{x})|^{q_i-1} \right\rangle \right) \left| \|\mathbf{g}(\mathbf{x})_{(\mathbf{g}_i)}\|_{p_i} \right\rangle |0\rangle |0\rangle |0\rangle |0\rangle |0\rangle |0\rangle \\ & \rightarrow |i\rangle_A \left(\bigotimes_{j=1}^{|\mathbf{g}_i|} |\mathbf{g}_{i,j}\rangle |g_{\mathbf{g}_{i,j}}(\mathbf{x})\rangle \left| \operatorname{sgn}(g_{\mathbf{g}_{i,j}}(\mathbf{x})) |g_{\mathbf{g}_{i,j}}(\mathbf{x})|^{q_i-1} \right\rangle \right) \left| \|\mathbf{g}(\mathbf{x})_{(\mathbf{g}_i)}\|_{p_i} \right\rangle \left| \|\mathbf{g}(\mathbf{x})_{(\mathbf{g}_i)}\|_{p_i}^* \right\rangle |0\rangle |0\rangle |0\rangle |0\rangle |0\rangle |0\rangle \end{aligned} \quad (79)$$

1554 Apply quantum maximum finding to the last register can then be used to find s^* in $O(\sqrt{|\mathcal{G}|})$ iterations.
1555 Each $g_{\mathbf{g}_{i,j}}(\mathbf{x})$ requires 2 queries to U_f , totally $O(|\mathbf{g}_i|)$ queries for a fixed i . In the above the
1556 index i ranges over $i = 1, 2, \dots, |\mathcal{G}|$. The query complexity is therefore $O(\sqrt{|\mathcal{G}|} |\mathbf{g}|_{\max})$, compared
1557 with the classical $\sum_{\mathbf{g} \in \mathcal{G}} |\mathbf{g}|$. Then the lemma follows. \blacksquare
1558

1559 B.7 PROOF OF THEOREM 6
1560

1561 **Theorem 6.** [Quantum FW over latent group norm ball] By setting $\sigma_t = \frac{C_f}{\sqrt{d}L(t+2)^{\max_{i \in [|\mathcal{G}|]} |\mathbf{g}_i|^{1/p_i}}}$
1562 for $t \in [T]$, the quantum algorithm (Algorithm 6) solves the latent group norm constraint optimiza-
1563 tion problem for any precision ε such that $f(\mathbf{x}^T) - f(\mathbf{x}^*) \leq \varepsilon$ in $T = \frac{4C_f}{\varepsilon} - 2$ rounds, succeed with
1564 probability $1 - p$, with $O\left(\sqrt{|\mathcal{G}|} |\mathbf{g}|_{\max} \log \frac{C_f}{p\varepsilon}\right)$ calls to the function value oracle U_f per round.
1565

1566 **Proof.** Let the true gradient component be $g_{\mathbf{g}_{i,j}}(\mathbf{x})$, and its estimated value be $\tilde{g}_{\mathbf{g}_{i,j}}(\mathbf{x})$ such that
 1567 $|\tilde{g}_{\mathbf{g}_{i,j}}(\mathbf{x}) - g_{\mathbf{g}_{i,j}}(\mathbf{x})| \leq \frac{\sqrt{d}L\sigma}{2}$. According to Step 7 of the algorithm, the dual norm computation
 1568 involves:
 1569

$$1570 \quad 1571 \quad \|g(\mathbf{x})_{(\mathbf{g}_i)}\|_{p_i}^* = \max_{\mathbf{s} \in \mathbb{R}^{|\mathbf{g}_i|}} \left\{ \sum_{j=1}^{|\mathbf{g}_i|} s_j g_{\mathbf{g}_{i,j}}(\mathbf{x}) \mid \|\mathbf{s}\|_{q_i} \leq 1 \right\}, \quad (80)$$

1573 where $\frac{1}{p_i} + \frac{1}{q_i} = 1$. The estimated dual norm is:
 1574

$$1575 \quad 1576 \quad \|\tilde{g}(\mathbf{x})_{(\mathbf{g}_i)}\|_{p_i}^* = \max_{\mathbf{s} \in \mathbb{R}^{|\mathbf{g}_i|}} \left\{ \sum_{j=1}^{|\mathbf{g}_i|} s_j \tilde{g}_{\mathbf{g}_{i,j}}(\mathbf{x}) \mid \|\mathbf{s}\|_{q_i} \leq 1 \right\}. \quad (81)$$

1578 The dual norm error can be decomposed as
 1579

$$1580 \quad 1581 \quad \left| \|\tilde{g}(\mathbf{x})_{(\mathbf{g}_i)}\|_{p_i}^* - \|g(\mathbf{x})_{(\mathbf{g}_i)}\|_{p_i}^* \right| \leq \max_{\|\mathbf{s}\|_{q_i} \leq 1} \left| \sum_{j=1}^{|\mathbf{g}_i|} \frac{s_j \sqrt{d}L\sigma}{2} \right|. \quad (82)$$

1583 By Hölder's inequality, for any \mathbf{s} satisfying $\|\mathbf{s}\|_{q_i} \leq 1$, let $\delta_{(\mathbf{g}_i)}$ be the vector in $\mathbb{R}^{|\mathbf{g}_i|}$ with all the
 1584 component being $\frac{\sqrt{d}L\sigma}{2}$ we have
 1585

$$1586 \quad 1587 \quad \left| \sum_{j=1}^{|\mathbf{g}_i|} \frac{s_j \sqrt{d}L\sigma}{2} \right| \leq \|\mathbf{s}\|_{q_i} \cdot \|\delta_{(\mathbf{g}_i)}\|_{p_i} \leq \|\delta_{(\mathbf{g}_i)}\|_{p_i}. \quad (83)$$

1590 Since $\|\delta_{(\mathbf{g}_i)}\|_{p_i} \leq \frac{\sqrt{d}L\sigma|\mathbf{g}_i|^{1/p_i}}{2}$, it follows that
 1591

$$1592 \quad 1593 \quad \left| \|\tilde{g}(\mathbf{x})_{(\mathbf{g}_i)}\|_{p_i}^* - \|g(\mathbf{x})_{(\mathbf{g}_i)}\|_{p_i}^* \right| \leq \|\delta_{(\mathbf{g}_i)}\|_{p_i} \leq \frac{\sqrt{d}L\sigma|\mathbf{g}_i|^{1/p_i}}{2}. \quad (84)$$

1594 Then, by Lemma 4 and 11, after Step 8, we have
 1595

$$1596 \quad 1597 \quad \|\tilde{g}(\mathbf{x})_{(\mathbf{g}_{i_t})}\|_{p_{i_t}}^* \geq \max_{i \in |\mathcal{G}|} \|g(\mathbf{x})_{(\mathbf{g}_i)}\|_{p_i}^* - \sqrt{d}L\sigma \max_{i \in |\mathcal{G}|} |\mathbf{g}_i|^{1/p_i}, \quad (85)$$

1598 succeed with probability at least $1 - \delta$, with query complexity of $O\left(\sqrt{|\mathcal{G}|}|\mathbf{g}|_{\max} \log \frac{1}{\delta}\right)$. Set $\delta = \frac{p}{T}$
 1599 to ensure that this procedure succeeds for all T iterations.
 1600

1601 The rest parallels the proof of Theorem 1. Set $\sigma_t = \frac{C_f}{\sqrt{d}L(t+2) \max_{i \in |\mathcal{G}|} |\mathbf{g}_i|^{1/p_i}}$ for all $t \in [T]$, after
 1602 $T = \frac{4C_f}{\varepsilon} - 2$ rounds, we have
 1603

$$1604 \quad f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \varepsilon, \quad (86)$$

1605 for any $\varepsilon > 0$. Then the theorem follows. \blacksquare
 1606

1607 B.8 PROOF OF LEMMA 7

1609 **Lemma 7. (Quantum top singular vector extraction)** Let there be efficient quantum access to a
 1610 matrix $M \in \mathbb{R}^{d \times d}$, with singular value decomposition $M = \sum_i^d \sigma_i \mathbf{u}_i \mathbf{v}_i^T$. Define $p = \frac{\sigma_1^2(M)}{\sum_{i=1}^d \sigma_i^2}$.
 1611 There exist quantum algorithms that with time complexity $O\left(\frac{\|M\|_F d \text{poly} \log d}{\sqrt{p} \epsilon \delta^2}\right)$, give the estimated
 1612 top singular value $\bar{\sigma}_1$ of M to precision ϵ and the corresponding unit estimated singular vectors \mathbf{u}, \mathbf{v} to precision δ such that $\|\mathbf{u} - \mathbf{u}_{top}\| \leq \delta$, $\|\mathbf{v} - \mathbf{v}_{top}\| \leq \delta$ with probability at least $1 - 1/\text{poly}(d)$.
 1613

1615 **Proof.** Initialize the quantum registers to the uniform superposition state by using Hadamard gates,
 1616 we have
 1617

$$1618 \quad 1619 \quad H^{\otimes d} |0\rangle |0\rangle |0\rangle \rightarrow \sum_i^d |i\rangle |0\rangle |0\rangle. \quad (87)$$

1620 By Assumption 4, we can perform the mapping
 1621

$$1622 \sum_i^d |i\rangle |0\rangle |0\rangle \rightarrow \frac{1}{\|M\|_F} \sum_i^d \sum_j^d M_{ij} |i\rangle |j\rangle |0\rangle, \quad (88)$$

1623
 1624

1625 in time $\tilde{O}(1)$. Note that
 1626

$$1627 \frac{1}{\|M\|_F} \sum_i^d \sum_j^d M_{ij} |i\rangle |j\rangle |0\rangle = \frac{1}{\|M\|_F} \sum_i^k \sigma_i |\mathbf{u}_i\rangle |\mathbf{v}_i\rangle |\bar{\sigma}_i\rangle. \quad (89)$$

1628
 1629

1630 Then by the quantum singular estimation algorithm (QSVE, Lemma 5), we have
 1631

$$1632 \frac{1}{\|M\|_F} \sum_i^d \sum_j^d M_{ij} |i\rangle |j\rangle |0\rangle \rightarrow \frac{1}{\|\nabla\|_F} \sum_i^k \sigma_i |\mathbf{u}_i\rangle |\mathbf{v}_i\rangle |\bar{\sigma}_i\rangle, \quad (90)$$

1633
 1634

1635 with the cost of $O\left(\frac{\|M\|_F \text{poly log } d}{\epsilon}\right)$. This process of generating such a state is treated as an oracle
 1636 which will be invoked multiple times in the quantum maximum finding. This requires that the
 1637 errors in the estimates of the singular values should be consistent across multiple runs. Note that
 1638 the randomness of QSVE comes from the quantum phase estimation algorithm, and the QSVE
 1639 algorithm of Lemma 5 uses a consistent version of phase estimation. This consistency in phase
 1640 estimation guarantees that the error patterns are reproducible, thereby maintaining uniform errors
 1641 over repeated oracle calls.

1642 Set $\epsilon \leq (\sigma_1 - \sigma_2)/2$ to ensure that even with the error of singular value estimation, the estimated
 1643 largest singular value is still larger than the estimated second largest singular value, which can
 1644 ensure that when we use the quantum maximum finding algorithm, if succeed, we will always get
 1645 the superposition state corresponding to the largest singular value. By Lemma 4, the cost of finding
 1646 the largest singular value is $O\left(\frac{1}{\sqrt{p}}\right)$. By Lemma 6, $O\left(\frac{d \log d}{\delta^2}\right)$ repeats are needed to tomography the
 1647 corresponding singular vectors of the largest singular value.
 1648

1649 Therefore, the overall complexity is $O\left(\frac{\|M\|_F d \text{poly log } d}{\sqrt{p} \epsilon \delta^2}\right)$. ■
 1650

1651 B.9 PROOF OF THEOREM 3

1652 **Theorem 3.** (Quantum FW with QTSVE) By setting $\delta_t = \frac{C_f}{2(t+2)\sigma_1(M_t)}$ and $\epsilon_t \leq (\sigma_1(M_t) -$
 1653 $\sigma_2(M_t))/2$ for $t \in [T]$, the quantum algorithm (Algorithm 3) solves the nuclear norm constraint
 1654 optimization problem for any precision ϵ such that $f(X^T) - f(X^*) \leq \epsilon$ in $T = \frac{4C_f}{\epsilon} - 2$ rounds,
 1655 with time complexity $\tilde{O}\left(\frac{r\sigma_1^3(M_t)d}{(\sigma_1(M_t) - \sigma_2(M_t))\epsilon^2}\right)$ for computing the update direction per round, where
 1656 r is the rank of the gradient matrix.
 1657

1658 **Proof.** By Lemma 7, set $\epsilon_t \leq (\sigma_1(M) - \sigma_2(M))/2$ to ensure that the quantum maximum finding
 1659 algorithm, if succeed, will always get the superposition state of the largest singular value. As the
 1660 QSVE algorithm from Lemma 5 use a consistent version of phase estimation, the estimated error
 1661 of the singular value will keep unchanged. Thus, we can measure the register of singular value
 1662 in the computational basic, to check whether the quantum maximum finding succeed, to boost up
 1663 the success probability. By Lemma 7, we obtain the estimated singular vectors \mathbf{u}, \mathbf{v} , which satisfy
 1664 $\|\mathbf{u} - \mathbf{u}_{top}\| \leq \delta_t, \|\mathbf{v} - \mathbf{v}_{top}\| \leq \delta_t$, with time complexity $O\left(\frac{\|M\|_F d \text{poly log } d}{\sqrt{p} \epsilon \delta^2}\right)$.
 1665

1666 Note that in the matrix case, the linear optimization subproblem of the Frank-Wolfe framework
 1667

$$1668 \min_{\hat{S} \in \mathcal{D}} \langle \hat{S}, M_t \rangle \quad \text{s.t.} \quad \text{tr}\{\hat{S}\} \leq 1 \quad (91)$$

1669
 1670

1671 is equivalent to the following problem
 1672

$$1673 \min_{\mathbf{x}, \mathbf{y} \in \mathbb{R}^d} \mathbf{x}^\top M_t \mathbf{y} \quad \text{s.t.} \quad \|\mathbf{x}\|, \|\mathbf{y}\| \leq 1. \quad (92)$$

1674 Therefore, since the update direction $S = \mathbf{u}^\top \mathbf{v}$, the solution quality of the linear subproblem can be
 1675 bounded with the solution quality of the equivalent problem, that is

$$\begin{aligned} 1676 \langle S, M_t \rangle - \min_{\hat{S} \in \mathcal{D}} \langle \hat{S}, M_t \rangle &= \langle \mathbf{u}^\top \mathbf{v}, M_t \rangle - \min_{\hat{S} \in \mathcal{D}} \langle \hat{S}, M_t \rangle \\ 1677 &= \mathbf{u}^\top M_t \mathbf{v} - \min_{\mathbf{x}, \mathbf{y} \in \mathbb{R}^d} \mathbf{x}^\top M_t \mathbf{y} \\ 1678 &= \mathbf{u}^\top M_t \mathbf{v} - \mathbf{u}_{top}^\top M_t \mathbf{v}_{top}. \end{aligned} \quad (93)$$

1682 Then by Lemma 7 and 15, we have

$$1683 |\mathbf{u}^\top M_t \mathbf{v} - \mathbf{u}_{top}^\top M_t \mathbf{v}_{top}| \leq 2\sigma_1(M_t)\delta_t. \quad (94)$$

1685 By the update rule and the definition of the curvature, for each round t , we have

$$1686 f(X^{(t+1)}) = f((1 - \gamma_t)X^{(t)} + \gamma_t S) \leq f(X^{(t)}) + \gamma_t \langle S - X^{(t)}, M_t \rangle + \frac{\gamma_t^2}{2} C_f \quad (95)$$

1688 Combining Inequality 93, 94 and 95, we have

$$1689 f(X^{(t+1)}) \leq f(X^{(t)}) + \gamma_t (\min_{\hat{S} \in \mathcal{D}} \langle \hat{S}, M_t \rangle - \langle X^{(t)}, M_t \rangle) + 2\gamma_t \sigma_1(M_t)\delta_t + \frac{\gamma_t^2}{2} C_f. \quad (96)$$

1692 Let $h(X^{(t)}) := f(X^{(t)}) - f(X^*)$, we have

$$\begin{aligned} 1693 h(X^{(t+1)}) &\leq h(X^{(t)}) + \gamma_t (\min_{\hat{S} \in \mathcal{D}} \langle \hat{S}, M_t \rangle - \langle X^{(t)}, M_t \rangle) + 2\gamma_t \sigma_1(M_t)\delta_t + \frac{\gamma_t^2}{2} C_f \\ 1694 &\leq h(X^{(t)}) - \gamma_t h(X^{(t)}) + 2\gamma_t \sigma_1(M_t)\delta_t + \frac{\gamma_t^2}{2} C_f \\ 1695 &= (1 - \gamma_t)h(X^{(t)}) + 2\gamma_t \sigma_1(M_t)\delta_t + \frac{\gamma_t^2}{2} C_f. \end{aligned} \quad (97)$$

1700 Set $\gamma_t = \frac{2}{t+2}$, $\delta_t = \frac{\gamma_t C_f}{4\sigma_1(M_t)}$, we have

$$1703 h(X^{(t+1)}) \leq \left(1 - \frac{2}{t+2}\right) h(X^{(t)}) + \left(\frac{2}{t+2}\right)^2 C_f. \quad (98)$$

1705 Using a similar induction as shown in Jaggi (2013) over t , we have

$$1707 h(X^{(t)}) \leq \frac{4C_f}{t+2}. \quad (99)$$

1709 In summary, set $\gamma_t = \frac{2}{t+2}$, $\delta_t = \frac{C_f}{2(t+2)\sigma_1(M_t)}$, after $T = \frac{4C_f}{\varepsilon} - 2$ rounds, we have

$$1711 f(\mathbf{x}^{(T)}) - f(\mathbf{x}^*) \leq \varepsilon, \quad (100)$$

1712 for any $\varepsilon > 0$. Since $\delta_t = \frac{C_f}{2(t+2)\sigma_1(M_t)} \geq \frac{C_f}{2(T+2)\sigma_1(M_t)} = \frac{\varepsilon}{2\sigma_1(M)}$, in each round, the time
 1713 complexity of update computing is $O\left(\frac{\|M\|_F \sigma_1^2(M) d \cdot \text{poly log } d}{\sqrt{p}(\sigma_1(M) - \sigma_2(M))\varepsilon^2}\right)$. Since $\|M\|_F \leq \sqrt{r}\sigma_1(M)$, $p \geq \frac{1}{r}$,
 1714 the time complexity is upper bounded by $O\left(\frac{r\sigma_1^3(M) d \cdot \text{poly log } d}{(\sigma_1(M) - \sigma_2(M))\varepsilon^2}\right)$, where r is the rank of the gradient
 1715 matrix. \blacksquare

1718 **Lemma 15.** For any $\|\mathbf{x} - \mathbf{x}'\|_2, \|\mathbf{y} - \mathbf{y}'\|_2 \leq \delta < 1$, $\|\mathbf{x}\|, \|\mathbf{y}\| \leq 1$, we have

$$1719 \left| \mathbf{x}^\top M \mathbf{y} - \mathbf{x}'^\top M \mathbf{y}' \right| \leq 2\sigma_1(M)\delta. \quad (101)$$

1722 **Proof.** Since $\mathbf{x}^\top M \mathbf{y} - \mathbf{x}'^\top M \mathbf{y}' = (x - x')^\top M \mathbf{y} + \mathbf{x}'^\top M(\mathbf{y} - \mathbf{y}')$, we have

$$1723 \left| \mathbf{x}^\top M \mathbf{y} - \mathbf{x}'^\top M \mathbf{y}' \right| \leq \sigma_1(M)\|\mathbf{x} - \mathbf{x}'\|_2 \|\mathbf{y}\|_2 + \sigma_1(M)\|\mathbf{x}'\|_2 \|\mathbf{y} - \mathbf{y}'\|_2. \quad (102)$$

1725 Thus, for any $\|\mathbf{x} - \mathbf{x}'\|_2, \|\mathbf{y} - \mathbf{y}'\|_2 \leq \delta < 1$, $\|\mathbf{x}\|, \|\mathbf{y}\| \leq 1$, we have

$$1726 \left| \mathbf{x}^\top M \mathbf{y} - \mathbf{x}'^\top M \mathbf{y}' \right| \leq 2\sigma_1(M)\delta. \quad (103)$$

1728 B.10 PROOF OF LEMMA 9
1729

1730 **Lemma 9.** (Quantum power method) Let there be quantum access to the matrix $M \in R^{d \times d}$ with
1731 $\sigma_{\max} \leq 1$, and to a vector $\mathbf{z} \in R^d$. Let γ'_{\min} be the lower bound of $\|(M^\top M)^i \mathbf{z}\|$ for all $i \in [k]$.
1732 There exists a quantum algorithm that creates a state $|\mathbf{y}\rangle$ such that $\|\langle \mathbf{y} | (M^\top M)^k \mathbf{z} \rangle\| \leq \delta$ in
1733 time $\tilde{O}\left(\frac{k}{\gamma'_{\min}} \|M\|_F \log(1/\delta)\right)$, with probability at least $1 - O(k/\text{poly}(d))$.
1734

1735 **Proof.** Suppose $\|z_l - Mz_{l-1}\| \leq \epsilon$ with $z_l = Mz_{l-1}$ for $l \in [L]$, and $z_0 = x$, we have
1736

$$\begin{aligned}
 \|z_1 - Mx\| &\leq \epsilon \\
 \|z_2 - M^2x\| &\leq \|z_2 - Mz_1 + Mz_1 - M^2x\| \\
 &\leq \|z_2 - Mz_1\| + \|Mz_1 - M^2x\| \\
 &\leq \epsilon + \|M(z_1 - Mx)\| \\
 &\leq \epsilon + \sigma_{\max} \|z_1 - Mx\| \\
 &\leq (1 + \sigma_{\max})\epsilon \\
 \|z_3 - M^3x\| &\leq \|z_3 - Mz_2 + Mz_2 - M^3x\| \\
 &\leq \|z_3 - Mz_2\| + \|Mz_2 - M^3x\| \\
 &\leq \epsilon + \|M(z_2 - M^2x)\| \\
 &\leq \epsilon + \sigma_{\max} \|z_2 - M^2x\| \\
 &\leq \epsilon + \sigma_{\max}(1 + \sigma_{\max})\epsilon \\
 &\leq (1 + \sigma_{\max} + \sigma_{\max}^2)\epsilon. \tag{104}
 \end{aligned}$$

1752 We use $\sigma_{\min}\|x\| \leq \|Mx\| \leq \sigma_{\max}\|x\|$, where $\sigma_{\max} = \max_{x \neq 0} x^\top Mx/\|x\|^2$. By induction, we
1753 have
1754

$$\|z_L - M^Lx\| \leq \sum_{i \in [L]} \sigma_{\max}^{i-1} \epsilon = \frac{\sigma_{\max}^L - 1}{\sigma_{\max} - 1} \epsilon. \tag{105}$$

1757 Let γ'_{\min} be the lower bound of $\|(M^\top M)^i \mathbf{z}\|$ for all $i \in [k]$. As each multiplication requires time
1758 complexity of $\tilde{O}\left(\frac{1}{\gamma'} \|M\|_F \log(1/\epsilon)\right)$ (Lemma 8), k steps of multiplication require time complexity
1759 of $\tilde{O}\left(\frac{k}{\gamma'_{\min}} \|M\|_F \log(1/\epsilon)\right)$. Furthermore, since
1760

$$\log \frac{1 - \sigma_1^k(M)}{1 - \sigma_1(M)} \leq -\log(1 - \sigma_1(M)) \leq \frac{1}{1 - \sigma_1(M)}, \tag{106}$$

1764 if we want $\|z_k - M^kx\| \leq \delta$, the time complexity will be $\tilde{O}\left(\frac{k \|M\|_F}{(1 - \sigma_1(M)) \gamma'_{\min}} \log(1/\delta)\right)$.
1765 ■
1766

1768 B.11 PROOF OF THEOREM 4
1769

1770 **Theorem 4.** (Quantum FW with QPM) By setting $k_t = \frac{2C_0\sigma_1(M_t) \ln d}{\varepsilon}$, $\delta_t = \delta'_t = \frac{\varepsilon\gamma'_{\min}}{16\sigma_1(M_t)}$ for
1771 $t \in [T]$, the quantum algorithm (Algorithm 4) solves the nuclear norm constraint optimization
1772 problem for any precision ε such that $f(X^T) - f(X^*) \leq \varepsilon$ in $T = \frac{4C_f}{\varepsilon} - 2$ rounds, with time
1773 complexity $\tilde{O}\left(\frac{\sqrt{r}\sigma_1^4(M_t)d}{(1 - \sigma_1(M_t))\gamma'_{\min}^3 \varepsilon^3}\right)$ for computing the update direction per round, where r is the
1774 rank of the gradient matrix, C_0 is a constant and γ'_{\min} is the lower bound of $\|(M_t^\top M_t)^i \mathbf{b}\|$ for all
1775 $i \in [k]$.
1776

1777 **Proof.** Denote $(MM^\top)^k \mathbf{b}$ as \mathbf{z}_u , $(M^\top M)^k \mathbf{b}$ as \mathbf{z}_v . For the quantum power method, we first use
1778 the Lemma 8 to construct a unitary U_1 which computes k steps of multiplication: $U_1 : |\mathbf{b}\rangle |\mathbf{b}\rangle \rightarrow$
1779 $|\bar{\mathbf{z}}_u\rangle |\bar{\mathbf{z}}_v\rangle$ with $\|\bar{\mathbf{z}}_u - \mathbf{z}_u\|_2 \leq \delta$ and $\|\bar{\mathbf{z}}_v - \mathbf{z}_v\|_2 \leq \delta$ (Lemma 9). Then we tomography $|\bar{\mathbf{z}}_u\rangle |\bar{\mathbf{z}}_v\rangle$
1780 to get \mathbf{u}, \mathbf{v} . Simalar to the proof of Theorem 3, our goal is to ensure $\left| \frac{\mathbf{u}^\top M \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} - \sigma_1(M) \right| \leq \varepsilon$.
1781

1782 First, we settle down $k = \frac{2C_0\sigma_1(M)\ln d}{\varepsilon}$ so that we have
 1783

$$1784 \quad \left| \frac{\mathbf{z}_u^\top M \mathbf{z}_v}{\|\mathbf{z}_u\| \|\mathbf{z}_v\|} - \sigma_1(M) \right| \leq \varepsilon/2. \quad (107)$$

1787 Suppose $\sigma_1(M) < 1$ and $\|(M^\top M)^i \mathbf{b}\| \in [\gamma'_{\min}, 1]$ for $i = 1, \dots, k$, after applying k times of
 1788 quantum matrix-vector multiplication (U_1) as described by Lemma 8, we obtain $|\bar{\mathbf{z}}_u\rangle |\bar{\mathbf{z}}_v\rangle$ with
 1789 $\|\bar{\mathbf{z}}_u - \mathbf{z}_u\|_2 \leq \delta$ and $\|\bar{\mathbf{z}}_v - \mathbf{z}_v\|_2 \leq \delta$ in time $T(U_1) = \tilde{O}\left(\frac{k\|M\|_F}{(1-\sigma_1(M))\gamma'_{\min}} \log(1/\delta)\right)$. Using
 1790 U_1 , we can tomography $|\bar{\mathbf{z}}_u\rangle |\bar{\mathbf{z}}_v\rangle$ and obtain \mathbf{u}, \mathbf{v} with $\|\mathbf{u} - \bar{\mathbf{z}}_u\| \leq \delta', \|\mathbf{v} - \bar{\mathbf{z}}_v\| \leq \delta'$ in time
 1791 $O\left(\frac{T(U_1)d \log d}{(\delta')^2}\right)$. By the triangle inequality, we have
 1792

$$1794 \quad \|\mathbf{u} - \mathbf{z}_u\|_2 \leq \delta + \delta' \leq 1, \|\mathbf{v} - \mathbf{z}_v\|_2 \leq \delta + \delta' \leq 1 \quad (108)$$

1795 Notice that

$$1797 \quad \left\| \frac{\mathbf{u}}{\|\mathbf{u}\|} - \frac{\mathbf{z}_u}{\|\mathbf{z}_u\|} \right\| = \left\| \frac{\mathbf{u}}{\|\mathbf{u}\|} - \frac{\mathbf{u}}{\|\mathbf{z}_u\|} + \frac{\mathbf{u}}{\|\mathbf{z}_u\|} - \frac{\mathbf{z}_u}{\|\mathbf{z}_u\|} \right\| \\ 1798 \quad \leq \left\| \frac{\mathbf{u}}{\|\mathbf{u}\|} - \frac{\mathbf{u}}{\|\mathbf{z}_u\|} \right\| + \left\| \frac{\mathbf{u}}{\|\mathbf{z}_u\|} - \frac{\mathbf{z}_u}{\|\mathbf{z}_u\|} \right\| \\ 1799 \quad \leq 2 \frac{\|\mathbf{u} - \mathbf{z}_u\|}{\|\mathbf{z}_u\|}, \quad (109)$$

1800 we have

$$1805 \quad \left\| \frac{\mathbf{u}}{\|\mathbf{u}\|} - \frac{\mathbf{z}_u}{\|\mathbf{z}_u\|} \right\| \leq 2 \frac{\delta + \delta'}{\gamma'_{\min}}. \quad (110)$$

1806 Similarly, we have

$$1808 \quad \left\| \frac{\mathbf{v}}{\|\mathbf{v}\|} - \frac{\mathbf{z}_v}{\|\mathbf{z}_v\|} \right\| \leq 2 \frac{\delta + \delta'}{\gamma'_{\min}}. \quad (111)$$

1809 Thus, we have

$$1812 \quad \left| \frac{\mathbf{u}^\top M \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} - \frac{\mathbf{z}_u^\top M \mathbf{z}_v}{\|\mathbf{z}_u\| \|\mathbf{z}_v\|} \right| \leq \left| \frac{\mathbf{u}^\top M \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} - \frac{\mathbf{u}^\top M \mathbf{z}_v}{\|\mathbf{u}\| \|\mathbf{z}_v\|} \right| + \left| \frac{\mathbf{u}^\top M \mathbf{z}_v}{\|\mathbf{u}\| \|\mathbf{z}_v\|} - \frac{\mathbf{z}_u^\top M \mathbf{z}_v}{\|\mathbf{z}_u\| \|\mathbf{z}_v\|} \right| \\ 1813 \quad \leq \frac{\|M\| \|\mathbf{v} - \mathbf{z}_v\|}{\|\mathbf{z}_v\|} + \frac{\|M\| \|\mathbf{u} - \mathbf{z}_u\|}{\|\mathbf{z}_u\|} \\ 1814 \quad \leq 4 \frac{(\delta + \delta')\sigma_1(M)}{\gamma'_{\min}}. \quad (112)$$

1815 The remaining proof is similar to that of Theorem 3. Now we set $\delta = \delta' = \frac{\varepsilon\gamma'_{\min}}{16\sigma_1(M)}$,
 1816 $\left| \frac{\mathbf{u}^\top M \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} - \frac{\mathbf{z}_u^\top M \mathbf{z}_v}{\|\mathbf{z}_u\| \|\mathbf{z}_v\|} \right| \leq \varepsilon/2$. Therefore, $\left| \frac{\mathbf{u}^\top M \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} - \sigma_1(M) \right| \leq \varepsilon$. The time complexity is
 1817 $\tilde{O}(T(U_1)d/(\delta')^2) = \tilde{O}\left(\frac{\sqrt{r}\sigma_1^4(M)d}{(1-\sigma_1(M))\gamma'_{\min}^3 \varepsilon^3}\right)$, where r is the rank of the gradient matrix. \blacksquare
 1818

1824 C ETHICS STATEMENT

1825 This work is a theoretical study. As such, we do not foresee any immediate specific ethical issues
 1826 arising from this work. We have conducted our research with integrity and in accordance with the
 1827 academic standards of our community.

1831 D REPRODUCIBILITY STATEMENT

1832 As a theoretical paper, all claims and results are supported by detailed mathematical proofs provided
 1833 in the main text and appendices. Therefore, the results can be reproduced by verifying the logical
 1834 steps of the proofs. We have endeavored to make our proofs as clear and self-contained as possible
 1835 to facilitate verification by the reader.

1836 **E LLM USAGE**
18371838 Large Language Models (LLMs) were used solely to aid or polish writing. This includes polishing
1839 sentences, improving grammar, and enhancing the readability and fluency of the text.
1840

1841

1842

1843

1844

1845

1846

1847

1848

1849

1850

1851

1852

1853

1854

1855

1856

1857

1858

1859

1860

1861

1862

1863

1864

1865

1866

1867

1868

1869

1870

1871

1872

1873

1874

1875

1876

1877

1878

1879

1880

1881

1882

1883

1884

1885

1886

1887

1888

1889