

---

# Memory-Constrained Algorithms for Convex Optimization

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We propose a family of recursive cutting-plane algorithms to solve feasibility  
2 problems with constrained memory, which can also be used for first-order convex  
3 optimization. Precisely, in order to find a point within a ball of radius  $\epsilon$  with a  
4 separation oracle in dimension  $d$ —or to minimize 1-Lipschitz convex functions to  
5 accuracy  $\epsilon$  over the unit ball—our algorithms use  $\mathcal{O}(\frac{d^2}{p} \ln \frac{1}{\epsilon})$  bits of memory, and  
6 make  $\mathcal{O}((C \frac{d}{p} \ln \frac{1}{\epsilon})^p)$  oracle calls, for some universal constant  $C \geq 1$ . The family  
7 is parametrized by  $p \in [d]$  and provides an oracle-complexity/memory trade-off in  
8 the sub-polynomial regime  $\ln \frac{1}{\epsilon} \gg \ln d$ . While several works gave lower-bound  
9 trade-offs (impossibility results) [29, 5]—we explicit here their dependence with  
10  $\ln \frac{1}{\epsilon}$ , showing that these also hold in any sub-polynomial regime—to the best of  
11 our knowledge this is the first class of algorithms that provides a positive trade-off  
12 between gradient descent and cutting-plane methods in any regime with  $\epsilon \leq 1/\sqrt{d}$ .  
13 The algorithms divide the  $d$  variables into  $p$  blocks and optimize over blocks  
14 sequentially, with approximate separation vectors constructed using a variant of  
15 Vaidya’s method. In the regime  $\epsilon \leq d^{-\Omega(d)}$ , our algorithm with  $p = d$  achieves the  
16 information-theoretic optimal memory usage and improves the oracle-complexity  
17 of gradient descent.

## 18 1 Introduction

19 Optimization algorithms are ubiquitous in machine learning, from solving simple regressions to  
20 training neural networks. Their essential roles have motivated numerous studies on their efficiencies,  
21 which are usually analyzed through the lens of oracle-complexity: given an oracle (such as function  
22 value, or subgradient oracle), how many calls to the oracle are needed for an algorithm to output an  
23 approximate optimal solution? [32]. However, ever-growing problem sizes have shown an inadequacy  
24 in considering only the oracle-complexity, and have motivated the study of the trade-off between  
25 oracle-complexity and other resources such as memory [49, 29, 5] and communication[23, 38, 40, 43,  
26 31, 50, 48, 47].

27 In this work, we study the oracle-complexity/memory trade-off for first-order non-smooth convex  
28 optimization, and the closely related feasibility problem, with a focus on developing memory efficient  
29 (deterministic) algorithms. Since [49] formally posed as open problem the question of characterizing  
30 this trade-off, there have been exciting results showing what is impossible: for convex optimization in  
31  $\mathbb{R}^d$ , [29] shows that any randomized algorithm with  $d^{1.25-\delta}$  bits of memory needs at least  $\tilde{\Omega}(d^{1+4\delta/3})$   
32 queries, and this has later been improved for deterministic algorithms to  $d^{1-\delta}$  bits of memory or  
33  $\tilde{\Omega}(d^{1+\delta/3})$  queries by [5]; in addition [5] shows that for the feasibility problem with a separation  
34 oracle, any algorithm which uses  $d^{2-\delta}$  bits of memory needs at least  $\tilde{\Omega}(d^{1+\delta})$  queries.

35 Despite these recent results on the lower bounds, all known first-order convex optimization algorithms  
 36 that output an  $\epsilon$ -suboptimal point fall into two categories: those that are quadratic in memory but  
 37 can potentially achieve the optimal  $\mathcal{O}(d \ln \frac{1}{\epsilon})$  query complexity, as represented by the center-of-mass  
 38 method, and those that have  $\mathcal{O}(\frac{d^2}{\epsilon})$  query complexity but only need the optimal  $\mathcal{O}(d \ln \frac{1}{\epsilon})$  bits of  
 39 memory, as represented by the classical gradient descent [49]. In addition, the above-mentioned  
 40 memory bounds apply only between queries, and in particular the center-of-mass method [49] is  
 41 allowed to use infinite memory during computations.

42 We propose a family of memory-constrained algorithms for the stronger feasibility problem in which  
 43 one aims to find a point within a set  $Q$  containing a ball of radius  $\epsilon$ , with access to a separation oracle.  
 44 In particular, this can be used for convex optimization since the subgradient information provides a  
 45 separation vector. Our algorithms use  $\mathcal{O}(\frac{d^2}{p} \ln \frac{1}{\epsilon})$  bits of memory (including during computations) and  
 46  $\mathcal{O}((C \frac{d}{p} \ln \frac{1}{\epsilon})^p)$  queries for some universal constant  $C \geq 1$ , and a parameter  $p \in [d]$  that can be chosen  
 47 by the user. Intuitively, in the context of convex optimization, the algorithms are based on the idea that  
 48 for any function  $f(\mathbf{x}, \mathbf{y})$  convex in the pair  $(\mathbf{x}, \mathbf{y})$ , the partial minimum  $\min_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$  as a function  
 49 of  $\mathbf{x}$  is still convex and, using a variant of Vaidya's method proposed in [25], our algorithm can  
 50 approximate subgradients for that function  $\min_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ , thereby turning an optimization problem  
 51 with variables  $(\mathbf{x}, \mathbf{y})$  to one with just  $\mathbf{x}$ . This idea, applied recursively with the variables divided into  
 52  $p$  blocks, gives our family of algorithms and the above-mentioned memory and query complexity.

53 When  $p = 1$ , our algorithm is a memory-constrained version of Vaidya's method [46, 25], and  
 54 improves over the center-of-mass [49] method by a factor of  $\ln \frac{1}{\epsilon}$  in terms of memory while having  
 55 optimal oracle-complexity. The improvements provided by our algorithms are more significant in  
 56 regimes when  $\epsilon$  is very small in the dimension  $d$ : increasing the parameter  $p$  can further reduce the  
 57 memory usage of Vaidya's method ( $p = 1$ ) by a factor  $\ln \frac{1}{\epsilon} / \ln d$ , while still improving over the  
 58 oracle-complexity of gradient descent. In particular, in a regime  $\ln \frac{1}{\epsilon} = \text{poly}(\ln d)$ , these memory  
 59 improvements are only in terms of  $\ln d$  factors. However, in sub-polynomial regimes with potentially  
 60  $\ln \frac{1}{\epsilon} = d^c$  for some constant  $c > 0$ , these provide polynomial improvements to the memory of  
 61 standard cutting-plane methods.

62 As a summary, this paper makes the following contributions.

- 63 • Our class of algorithms provides a trade-off between memory-usage and oracle-complexity  
 64 whenever  $\ln \frac{1}{\epsilon} \gg \ln d$ . Further, taking  $p = 1$  improves the memory-usage from center-of-  
 65 mass [49] by a factor  $\ln \frac{1}{\epsilon}$ , while preserving the optimal oracle-complexity.
- 66 • For  $\ln \frac{1}{\epsilon} \geq \Omega(d \ln d)$ , our algorithm with  $p = d$  is the first known algorithm that outperforms  
 67 gradient descent in terms of the oracle-complexity, but still maintains the optimal  $\mathcal{O}(d \ln \frac{1}{\epsilon})$   
 68 memory usage.
- 69 • We show how to obtain a  $\ln \frac{1}{\epsilon}$  dependence in the known lower-bound trade-offs [29, 5],  
 70 confirming that the oracle-complexity/memory trade-off is necessary for any regime  $\epsilon \lesssim \frac{1}{\sqrt{d}}$ .

## 71 2 Setup and Preliminaries

72 In this section, we precise the formal setup for our results. We follow the framework introduced in  
 73 [49], to define the memory constraint on algorithms with access to an oracle  $\mathcal{O} : \mathcal{S} \rightarrow \mathcal{R}$  which  
 74 takes as input a query  $q \in \mathcal{S}$  and outputs a response  $\mathcal{O}(q) \in \mathcal{R}$ . Here, the algorithm is constrained to  
 75 update an internal  $M$ -bit memory between queries to the oracle.

76 **Definition 2.1** ( $M$ -bit memory-constrained algorithm [49, 29, 5]). *Let  $\mathcal{O} : \mathcal{S} \rightarrow \mathcal{R}$  be an oracle. An*  
 77  *$M$ -bit memory-constrained algorithm is specified by a query function  $\psi_{\text{query}} : \{0, 1\}^M \rightarrow \mathcal{S}$  and*  
 78 *an update function  $\psi_{\text{update}} : \{0, 1\}^M \times \mathcal{S} \times \mathcal{R} \rightarrow \{0, 1\}^M$ . The algorithm starts with the memory*  
 79 *state  $\text{Memory}_0 = 0^M$  and iteratively makes queries to the oracle. At iteration  $t$ , it makes the query*  
 80  *$q_t = \psi_{\text{query}}(\text{Memory}_{t-1})$  to the oracle, receives the response  $r_t = \mathcal{O}(q_t)$  then updates its memory*  
 81  *$\text{Memory}_t = \psi_{\text{update}}(\text{Memory}_{t-1}, q_t, r_t)$ .*

82 The algorithm can stop at any iteration and the last query is its final output. Importantly, this model  
 83 does not enforce constraints on the memory usage during the computation of  $\psi_{\text{update}}$  and  $\psi_{\text{query}}$ .  
 84 This is ensured in the stronger notion of a memory-constrained algorithm with computations. These

85 are precisely algorithms that have constrained memory including for computations, with the only  
 86 specificity that they need a decoder function  $\phi$  to make queries to the oracle from their bit memory,  
 87 and a discretization function  $\psi$  to write a discretized response into the algorithm’s memory.

88 **Definition 2.2** (*M-bit memory-constrained algorithm with computations*). *Let  $\mathcal{O} : \mathcal{S} \rightarrow \mathcal{R}$  be an*  
 89 *oracle. We suppose that we are given a decoding function  $\phi : \{0, 1\}^* \rightarrow \mathcal{S}$  and a discretization*  
 90 *function  $\psi : \mathcal{R} \times \mathbb{N} \rightarrow \{0, 1\}^*$  such that  $\psi(r, n) \in \{0, 1\}^n$  for all  $r \in \mathcal{R}$ . An M-bit memory-*  
 91 *constrained algorithm with computations is only allowed to use an M-bit memory in  $\{0, 1\}^M$  even*  
 92 *during computations. The algorithm has three special memory placements  $Q, N, R$ . Say the contents*  
 93 *of  $Q$  and  $N$  are  $q$  and  $n$  respectively. To make a query,  $R$  must contain at least  $n$  bits. The algorithm*  
 94 *submits  $q$  to the encoder which then submits the query  $\phi(q)$  to the oracle. If  $r = \mathcal{O}(\phi(q))$  is the*  
 95 *oracle response, the discretization function then writes  $\psi(r, n)$  in the placement  $R$ .*

96 **Feasibility problem.** In this problem, the goal is to find a point  $\mathbf{x} \in Q$ , where  $Q \subset \mathcal{C}_d := [-1, 1]^d$   
 97 is a convex set. We choose the cube  $[-1, 1]^d$  as prior bound for convenience in our later algorithms,  
 98 but the choice of norm for this prior ball can be arbitrary and does not affect our results. The algorithm  
 99 has access to a *separation oracle*  $O_S : \mathcal{C}_d \rightarrow \{\text{Success}\} \cup \mathbb{R}^d$ , that for a query  $\mathbf{x} \in \mathbb{R}^d$  either returns  
 100 Success if  $\mathbf{x} \in Q$ , or a separating hyperplane  $\mathbf{g} \in \mathbb{R}^d$ , i.e., such that  $\mathbf{g}^\top \mathbf{x} < \mathbf{g}^\top \mathbf{x}'$  for any  $\mathbf{x}' \in Q$ .  
 101 We suppose that the separating hyperplanes are normalized,  $\|\mathbf{g}\|_2 = 1$ . An algorithm solves the  
 102 feasibility problem with accuracy  $\epsilon$  if the algorithm is successful for any feasibility problem such that  
 103  $Q$  contains an  $\epsilon$ -ball  $B_d(\mathbf{x}^*, \epsilon)$  for  $\mathbf{x}^* \in \mathcal{C}_d$ .

104 As an important remark, this formulation asks that the separation oracle is consistent over time:  
 105 when queried at the exact same point  $\mathbf{x}$ , the oracle always returns the same separation vector. In  
 106 this context, we can use the natural decoding function  $\phi$  which takes as input  $d$  sequences of bits  
 107 and outputs the vector with coordinates given by the sequences interpreted in base 2. Similarly, the  
 108 natural discretization function  $\psi$  takes as input the separation hyperplane  $\mathbf{g}$  and outputs a discretized  
 109 version up to the desired accuracy. From now, we can omit these implementation details and consider  
 110 that the algorithm can query the oracle for discretized queries  $\mathbf{x}$ , up to specified rounding errors.

111 **Remark 2.1.** *An algorithm for the feasibility problem with accuracy  $\epsilon/(2\sqrt{d})$  can be used for first-*  
 112 *order convex optimization. Suppose one aims to minimize a 1-Lipschitz convex function  $f$  over the unit*  
 113 *ball, and output an  $\epsilon$ -suboptimal solution, i.e., find a point  $\mathbf{x}$  such that  $f(\mathbf{x}) \leq \min_{\mathbf{y} \in B_d(0,1)} f(\mathbf{y}) + \epsilon$ .*  
 114 *A separation oracle for  $Q = \{\mathbf{x} : f(\mathbf{x}) \leq \min_{\mathbf{y} \in B_d(0,1)} f(\mathbf{y}) + \epsilon\}$  is given at a query  $\mathbf{x}$  by the*  
 115 *subgradient information from the first-order oracle:  $-\frac{\partial f(\mathbf{x})}{\|\partial f(\mathbf{x})\|}$ . Its computation can also be carried*  
 116 *memory-efficiently up to rounding errors since if  $\|\partial f(\mathbf{x})\| \leq \epsilon/(2\sqrt{d})$ , the algorithm can return  $\mathbf{x}$*   
 117 *and already has the guarantee that  $\mathbf{x}$  is an  $\epsilon$ -suboptimal solution ( $\mathcal{C}_d$  has diameter  $2\sqrt{d}$ ). Notice that*  
 118 *because  $f$  is 1-Lipschitz,  $Q$  contains a ball of radius  $\epsilon/(2\sqrt{d})$  (the factor  $1/(2\sqrt{d})$  is due to potential*  
 119 *boundary issues). Hence, it suffices to run the algorithm for the feasibility problem while keeping in*  
 120 *memory the queried point with best function value.*

## 121 2.1 Known trade-offs between oracle-complexity and memory

122 **Known lower-bound trade-offs.** All known lower bound apply to the more general class of  
 123 memory-constrained algorithms without computational constraints given in Definition 2.1. [32] first  
 124 showed that  $\mathcal{O}(d \ln \frac{1}{\epsilon})$  queries are needed for solving convex optimization to ensure that one finds an  
 125  $\epsilon$ -suboptimal solution. Further,  $\mathcal{O}(d \ln \frac{1}{\epsilon})$  bits of memory are needed even just to output a solution in  
 126 the unit ball with  $\epsilon$  accuracy [49]. These historical lower bounds apply in particular to the feasibility  
 127 problem and are represented in the pictures of Fig. 1 as the dashed pink region.

128 More recently, [29] showed that achieving both optimal oracle-complexity and optimal memory is  
 129 impossible for convex optimization. They show that a possibly randomized algorithm with  $d^{1.25-\delta}$   
 130 bits of memory makes at least  $\tilde{\Omega}(d^{1+4\delta/3})$  queries. This result was extended for deterministic  
 131 algorithms in [5] which shows that a deterministic algorithm with  $d^{1-\delta}$  bits of memory makes  
 132  $\tilde{\Omega}(d^{1+\delta/3})$  queries. For the feasibility problem, they give an improved trade-off: any deterministic  
 133 algorithm with  $d^{2-\delta}$  bits of memory makes  $\tilde{\Omega}(d^{1+\delta})$  queries. These trade-offs are represented in the  
 134 left picture of Fig. 1 as the pink, red, and purple solid region, respectively.

135 **Known upper-bound trade-offs.** Prior to this work, to the best of our knowledge only two  
 136 algorithms were known in the oracle-complexity/memory landscape. First, cutting-plane algorithms  
 137 achieve the optimal oracle-complexity  $\mathcal{O}(d \ln \frac{1}{\epsilon})$  but use quadratic memory. The memory-constrained  
 138 (MC) center-of-mass method analyzed in [49] uses in particular  $\mathcal{O}(d^2 \ln^2 \frac{1}{\epsilon})$  memory. Instead, if one  
 139 uses Vaidya’s method which only needs to store  $\mathcal{O}(d)$  cuts instead  $\mathcal{O}(d \ln \frac{1}{\epsilon})$ , we show that one can  
 140 achieve  $\mathcal{O}(d^2 \ln \frac{1}{\epsilon})$  memory. These algorithms only use the separation oracle and hence apply to  
 141 both convex optimization and the feasibility problem. On the other hand, the memory-constrained  
 142 gradient descent for convex optimization [49] uses the optimal  $\mathcal{O}(d \ln \frac{1}{\epsilon})$  memory but makes  $\mathcal{O}(\frac{1}{\epsilon^2})$   
 143 iterations. While the analysis in [49] is only carried for convex optimization, we can give a modified  
 144 proof showing that gradient descent can also be used for the feasibility problem.

## 145 2.2 Other related works

146 Vaidya’s method [46, 36, 1, 2] and the variant [25] that we use in our algorithms, belong to the  
 147 family of cutting-plane methods. Perhaps the simplest example of an algorithm in this family is the  
 148 center-of-mass method, which achieves the optimal  $\mathcal{O}(d \ln \frac{1}{\epsilon})$  oracle-complexity but is computa-  
 149 tionally intractable, and the only known random walk-based implementation [4] has computational  
 150 complexity  $\mathcal{O}(d^7 \ln \frac{1}{\epsilon})$ . Another example is the ellipsoid method, which has suboptimal  $\mathcal{O}(d^2 \ln \frac{1}{\epsilon})$   
 151 query complexity, but has an improved computational complexity  $\mathcal{O}(d^4 \ln \frac{1}{\epsilon})$ . [8] pointed out that  
 152 Vaidya’s method achieves the best of both worlds by sharing the  $\mathcal{O}(d \ln \frac{1}{\epsilon})$  optimal query complexity  
 153 of the center-of-mass, and achieving a computational complexity of  $\mathcal{O}(d^{1+\omega} \ln \frac{1}{\epsilon})^1$ . In a major  
 154 breakthrough, this computational complexity was improved to  $\mathcal{O}(d^3 \ln^3 \frac{1}{\epsilon})$  in [25], then to  $\mathcal{O}(d^3 \ln \frac{1}{\epsilon})$   
 155 in [20]. We refer to [8, 25, 20] for more detailed comparisons of these algorithms.

156 Another popular convex optimization algorithm that requires quadratic memory is the Broyden-  
 157 Fletcher- Goldfarb- Shanno (BFGS) algorithm [41, 7, 18, 19], which stores an approximated inverse  
 158 Hessian matrix as gradient preconditioner. Several works aimed to reduce the memory usage of  
 159 BFGS; in particular, the limited memory BFGS (L-BFGS) stores a few vectors instead of the entire  
 160 approximated inverse Hessian matrix [35, 28]. However, it is still an open question whether even the  
 161 original BFGS converges for non-smooth convex objectives [27].

162 Lying at the other extreme of the oracle-complexity/memory trade-off is gradient descent, which  
 163 achieves the optimal memory usage but requires significantly more queries than center-of-mass or  
 164 Vaidya’s method in the regime  $\epsilon \lesssim \frac{1}{\sqrt{d}}$ . There is a rich literature of works aiming to speed up  
 165 gradient descent, such as the optimized gradient method [15, 14], Nesterov’s Acceleration [33], the  
 166 triple momentum method [39], geometric descent [9], quadratic averaging [16], the information-  
 167 theoretic exact method [44], or Big-Step-Little-Step method [21]. Interested readers can find a  
 168 comprehensive survey on acceleration methods in [10]. However, these acceleration methods usually  
 169 require additional smoothness or strong convexity assumptions (or both) on the objective function,  
 170 due to the known  $\Omega(\frac{1}{\epsilon^2})$  query lower bound in the large-scale regime  $\epsilon \gtrsim \frac{1}{\sqrt{d}}$  for any first order  
 171 method where the query points lie in the span of the subgradients of previous query points [34].

172 Besides accelerating gradient descent, researchers have investigated more efficient ways to leverage  
 173 subgradients obtained in previous iterations. Of interest are bundle methods [3, 22, 26], that have  
 174 found a wide range of applications [45, 24]. In their simplest form, they minimize the sum of the  
 175 maximum of linear lower bounds constructed using past oracle queries, and a regularization term  
 176 penalizing the distance from the current iteration variable. Although the theoretical convergence rate  
 177 of the bundle method is the same as that of gradient descent, in practice, bundle methods can benefit  
 178 from previous information and substantially outperform gradient descent [3].

179 The increasing size of optimization problems has also motivated the development of communication-  
 180 efficient optimization algorithms in distributed settings such as [23, 38, 40, 43, 31, 50, 48,  
 181 47]. Moreover, recent works have explored the trade-off between sample complexity and mem-  
 182 ory/communication complexity for learning problems under the streaming model, with notable  
 183 contributions including [6, 11, 12, 37, 42, 30].

---

<sup>1</sup> $\omega < 2.373$  is the exponent of matrix multiplication

184 **3 Main results**

185 We first check that the memory-constrained gradient descent method solves feasibility problems. This  
 186 was known for convex optimization [49] and the same algorithm with a modified proof gives the  
 187 following result. For completeness, the proof is given in Appendix D.

188 **Proposition 3.1.** *The memory-constrained gradient descent algorithm solves the feasibility problem  
 189 with accuracy  $\epsilon \leq \frac{1}{\sqrt{d}}$  using  $\mathcal{O}(d \ln \frac{1}{\epsilon})$  bits of memory and  $\mathcal{O}(\frac{1}{\epsilon^2})$  separation oracle calls.*

190 Our main contribution is a class of algorithms based on Vaidya’s cutting-plane method that provide a  
 191 query-complexity / memory tradeoff. More precisely, we show the following.

192 **Theorem 3.2.** *For any  $1 \leq p \leq d$ , there is a deterministic first-order algorithm that solves  
 193 the feasibility problem for accuracy  $\epsilon \leq \frac{1}{\sqrt{d}}$ , using  $\mathcal{O}(\frac{d^2}{p} \ln \frac{1}{\epsilon})$  bits of memory (including during  
 194 computations), with  $\mathcal{O}((C \frac{d}{p} \ln \frac{1}{\epsilon})^p)$  calls to the separation oracle, and computational complexity  
 195  $\mathcal{O}((C \frac{d}{p})^{1+\omega} \ln \frac{1}{\epsilon})^p$ , where  $C \geq 1$  is a universal constant.*

196 For simplicity, in Section 4, we describe algorithms that achieve this trade-off without computation  
 197 concerns (Definition 2.1), which already provide the main elements of our method. The proof  
 198 of oracle-complexity and memory usage is given in Appendix A. In Appendix B, we consider  
 199 computational constraints and give corresponding algorithms using the cutting-plane method of [25].

200 To better understand the implications of Theorem 3.2, it is useful to compare the provided class of  
 201 algorithms to the two algorithms known in the oracle-complexity/memory tradeoff landscape: the  
 202 memory-constrained center-of-mass method and the memory-constrained gradient descent [49].

203 For  $p = 1$ , our resulting procedure, which is essentially a memory-constrained Vaidya’s algorithm,  
 204 has optimal oracle-complexity  $\mathcal{O}(d \ln \frac{1}{\epsilon})$  and uses  $\mathcal{O}(d^2 \ln \frac{1}{\epsilon})$  bits of memory. This improves by a  
 205  $\ln \frac{1}{\epsilon}$  factor the memory usage of the center-of-mass-based algorithm provided in [49], which used  
 206  $\mathcal{O}(d^2 \ln^2 \frac{1}{\epsilon})$  memory and had the same optimal oracle-complexity.

207 Next, we recall that the memory-constrained gradient descent method used the optimal number  
 208  $\mathcal{O}(d \ln \frac{1}{\epsilon})$  bits of memory (including for computations), and a sub-optimal  $\mathcal{O}(\frac{1}{\epsilon^2})$  oracle-complexity.  
 209 While the memory of our algorithms decreases with  $p$ , their oracle-complexity is exponential in  $p$ .  
 210 This significantly restricts the values of  $p$  for which the oracle-complexity is improved over that of  
 211 gradient descent. The range of application of Theorem 3.2 is given in the next result.

212 **Corollary 3.1.** *The algorithms given in Theorem 3.2 effectively provide a tradeoff for  $p \leq \mathcal{O}(\frac{\ln \frac{1}{\epsilon}}{\ln d} \vee d)$ .  
 213 Precisely, this provides a tradeoff between*

- 214 • using  $\mathcal{O}(d^2 \ln \frac{1}{\epsilon})$  memory with optimal  $\mathcal{O}(d \ln \frac{1}{\epsilon})$  oracle-complexity, and
- 215 • using  $\mathcal{O}(d^2 \ln d \wedge d \ln \frac{1}{\epsilon})$  memory with  $\mathcal{O}(\frac{1}{\epsilon^2} \vee (C \ln \frac{1}{\epsilon})^d)$  oracle-complexity.

216 Importantly, for  $\epsilon \leq \frac{1}{d^{\Omega(d)}}$ , taking  $p = d$  yields an algorithm that uses the optimal memory  $\mathcal{O}(d \ln \frac{1}{\epsilon})$   
 217 and has an improved query complexity over gradient descent. In this regime of small (virtually  
 218 constant) dimension, for the same memory usage, gradient descent has a query complexity that is  
 219 polynomial in  $\epsilon$ ,  $\mathcal{O}(\frac{1}{\epsilon^2})$ , while our algorithm has poly-logarithmic dependence in  $\epsilon$ ,  $\mathcal{O}_d(\ln^d \frac{1}{\epsilon})$ , where  
 220  $\mathcal{O}_d$  hides an exponential constant in  $d$ . It remains open whether this  $\ln^d \frac{1}{\epsilon}$  dependence in the oracle-  
 221 complexity is necessary. To the best of our knowledge, this is the first example of an algorithm that  
 222 improves over gradient descent while keeping its optimal memory usage in any regime where  $\epsilon \leq \frac{1}{\sqrt{d}}$ .

223 While this improvement holds only in the exponential regime  $\epsilon \leq \frac{1}{d^{\Omega(d)}}$ , Theorem 3.2 still provides  
 224 a non-trivial trade-off whenever  $\ln \frac{1}{\epsilon} \gg \ln d$ , and improves over the known memory-constrained  
 225 center-of-mass in the standard regime  $\epsilon \leq \frac{1}{\sqrt{d}}$  [49]. Fig. 1 depicts the trade-offs in the two regimes  
 226 mentioned earlier.

227 Last, we note that the lower-bound trade-offs presented in [29, 5] do not show a dependence in the  
 228 accuracy  $\epsilon$ . Especially in the regime when  $\ln \frac{1}{\epsilon} \gg \ln d$ , this yields sub-optimal lower bounds (in  
 229 fact even in the regime  $\epsilon = 1/\text{poly}(d)$ , our more careful analysis improves the lower bound on the  
 230 memory by a  $\ln d$  factor). We show with simple arguments that one can extend their results to include  
 231 a  $\ln \frac{1}{\epsilon}$  factor for both memory and query complexity. Fig. 1 presented these improved lower bounds.

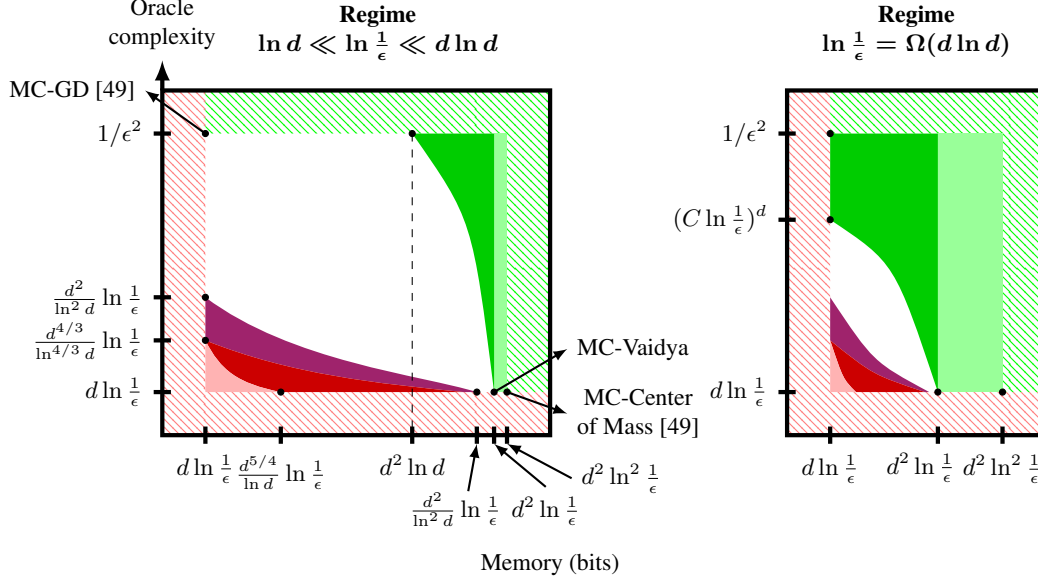


Figure 1: Trade-offs between available memory and first-order oracle-complexity for the feasibility problem over the unit ball. MC=Memory-constrained. GD=Gradient Descent. The left picture corresponds to the regime  $\epsilon \gg d^{-\Omega(d)}$  and  $\epsilon \leq 1/\text{poly}(d)$  and the right picture represents the regime  $\epsilon \leq d^{-\mathcal{O}(d)}$ . For both figures, the dashed pink "L" (resp. green inverted "L") region corresponds to historical lower (resp. upper) bounds for randomized algorithms. The solid pink (resp. red) lower bound tradeoff is due to [29] (resp. [5]) for randomized algorithms (resp. deterministic algorithms). The purple region is a lower bound tradeoff for the feasibility problem for accuracy  $\epsilon$  and deterministic algorithms [5]. All these lower-bound trade-offs are represented with their  $\ln \frac{1}{\epsilon}$  dependence (Theorem 3.3). We use memory-constrained Vaidya's method to gain a factor  $\ln \frac{1}{\epsilon}$  in memory compared to memory-constrained center-of-mass [49], which gives the light green region, and a class of algorithms represented in dark green, that allows trading query-complexity for an extra  $\ln \frac{1}{\epsilon} / \ln d$  factor saved in memory (Theorem 3.2). The dark green dashed region in the left figure emphasizes that the area covered by our class of algorithms depends highly on the regime for the accuracy  $\epsilon$ : the resulting improvement in memory is more significant as  $\epsilon$  is smaller. In the regime when  $\epsilon \leq d^{-\mathcal{O}(d)}$  (right figure), our class of algorithms improves over the oracle-complexity of gradient descent while keeping the optimal memory  $\mathcal{O}(d \ln \frac{1}{\epsilon})$ .

232 **Theorem 3.3.** For  $\epsilon \leq 1/\text{poly}(d)$  and any  $\delta \in [0, 1]$  (the notation  $\tilde{\Omega}$  hides  $\ln^{\mathcal{O}(1)} d$  factors),

- 233 1. any (randomized) algorithm guaranteed to minimize 1-Lipschitz convex functions over the  
234 unit ball with accuracy  $\epsilon$  uses  $d^{5/4-\delta} \ln \frac{1}{\epsilon}$  bits of memory or makes  $\tilde{\Omega}(d^{1+4\delta/3} \ln \frac{1}{\epsilon})$  queries,
- 235 2. any deterministic algorithm guaranteed to minimize 1-Lipschitz convex functions over the  
236 unit ball with accuracy  $\epsilon$  uses  $d^{2-\delta} \ln \frac{1}{\epsilon}$  bits of memory or makes  $\tilde{\Omega}(d^{1+\delta/3} \ln \frac{1}{\epsilon})$  queries,
- 237 3. any deterministic algorithm guaranteed to solve the feasibility problem over the unit ball  
238 with accuracy  $\epsilon$  uses  $d^{2-\delta} \ln \frac{1}{\epsilon}$  bits of memory or makes  $\tilde{\Omega}(d^{1+\delta} \ln \frac{1}{\epsilon})$  queries.

239 The proof is given in Appendix C and the arguments therein could readily be used to exhibit the  $\ln \frac{1}{\epsilon}$   
240 dependence of potential future works improving over these lower bounds trade-offs.

241 **Sketch of proof.** At a high level, [29, 5] use a barrier term  $\|\mathbf{Ax}\|_\infty$  where  $\mathbf{A}$  has  $\Theta(d)$  rows: if an  
242 algorithm does not have enough memory,  $\mathbf{A}$  cannot be fully stored which in turn incurs a sub-optimal  
243 oracle-complexity. To achieve a  $\ln \frac{1}{\epsilon}$  improvement in memory (Appendix C.1), we modify the  
244 sampling of rows of  $\mathbf{A}$ , from uniform on vertices of the hypercube to uniform in an  $\epsilon$ -net. The proof  
245 can then be adapted accordingly. Last, one can improve the oracle-complexity by a  $\ln \frac{1}{\epsilon} / \ln d$  factor  
246 (Appendix C.2) using a standard rescaling argument [32].

247 **4 Memory-constrained feasibility problem without computation**

248 In this section, we present a class of algorithms that are memory-constrained according to Defini-  
 249 tion 2.1 and achieve the desired memory and oracle-complexity bounds. We emphasize that the  
 250 memory constraint is only applied between calls to the oracle and as a result, the algorithm is allowed  
 251 infinite computation memory and computation power between calls to the oracle.

252 We start by defining discretization functions that will be used in our algorithms. For  $\xi > 0$  and  
 253  $x \in [-1, 1]$ , we pose  $\text{Discretize}_1(x, \xi) = \text{sign}(x) \cdot \xi \lfloor |x|/\xi \rfloor$ . Next, we define the discretization  
 254  $\text{Discretize}_d$  for general dimensions  $d \geq 1$ . For any  $\mathbf{x} \in C$  and  $\xi > 0$ ,

$$\text{Discretize}_d(\mathbf{x}, \xi) = \left( \text{Discretize}_1(x_1, \xi/\sqrt{d}), \dots, \text{Discretize}_1(x_d, \xi/\sqrt{d}) \right).$$

255 **4.1 Memory-constrained Vaidya's method**

256 Our algorithm recursively uses Vaidya's cutting-plane method [46] and subsequent works expanding  
 257 on this method. We briefly describe the method. Given a polyhedron  $\mathcal{P} = \{\mathbf{x} : \mathbf{A}\mathbf{x} \geq \mathbf{b}\}$ , we define  
 258  $s_i(\mathbf{x}) = \mathbf{a}_i^\top \mathbf{x} - b_i$  and  $\mathbf{S}_x = \text{diag}(s_i(x), i \in [d])$ . We will also use the shorthand  $\mathbf{A}_x = \mathbf{S}_x^{-1} \mathbf{A}$ .  
 259 The volumetric barrier is defined as

$$V_{\mathbf{A}, \mathbf{b}}(\mathbf{x}) = \frac{1}{2} \ln \det(\mathbf{A}_x^\top \mathbf{A}_x).$$

260 At each step, Vaidya's method queries the volumetric center of the polyhedron, which is the point  
 261 minimizing the volumetric barrier. For convenience, we denote by  $\text{VolumetricCenter}$  this function,  
 262 i.e., for any  $\mathbf{A} \in \mathbb{R}^{m \times d}$  and  $\mathbf{b} \in \mathbb{R}^d$  defining a non-empty polyhedron  $\mathcal{P} = \{\mathbf{x} : \mathbf{A}\mathbf{x} \geq \mathbf{b}\}$ ,

$$\text{VolumetricCenter}(\mathbf{A}, \mathbf{b}) = \arg \min_{\mathbf{x} : \mathbf{A}\mathbf{x} \geq \mathbf{b}} V_{\mathbf{A}, \mathbf{b}}(\mathbf{x}).$$

263 When the polyhedron is unbounded, we can for instance take  $\text{VolumetricCenter}(\mathbf{A}, \mathbf{b}) = \mathbf{0}$ . Vaidya's  
 264 method makes use of leverage scores for each constraint  $i$  of the polyhedron, defined as  $\sigma_i =$   
 265  $(\mathbf{A}_x \mathbf{H}^{-1} \mathbf{A}_x^\top)_{i,i}$ , where  $\mathbf{H} = \mathbf{A}_x^\top \mathbf{A}_x$ . We are now ready to define the update procedure for the  
 266 polyhedron considered by Vaidya's volumetric method. We denote by  $\mathcal{P}_t$  the polyhedron stored in  
 267 memory after making  $t$  queries. The method keeps in memory the constraints defining the current  
 268 polyhedron and the iteration index  $k$  when these constraints were added, which will be necessary for  
 269 our next procedures. Hence, the polyhedron will be stored in the form  $\mathcal{P}_t = \{(k_i, \mathbf{a}_i, b_i), i \in [m]\}$ ,  
 270 and the associated constraints are given via  $\{\mathbf{x} : \mathbf{A}\mathbf{x} \geq \mathbf{b}\}$  where  $\mathbf{A}^\top = [\mathbf{a}_1, \dots, \mathbf{a}_m]$  and  $\mathbf{b}^\top =$   
 271  $[b_1, \dots, b_m]$ . By abuse of notation, we will write  $\text{VolumetricCenter}(\mathcal{P})$  for the volumetric center of  
 272 the polyhedron  $\text{VolumetricCenter}(\mathbf{A}, \mathbf{b})$  where  $\mathbf{A}$  and  $\mathbf{b}$  define the constraints stored in  $\mathcal{P}$ .

273 Initially, the polyhedron is simply  $\mathcal{C}_d$ , these constraints are given  $-1$  index for convenience, and  
 274 they will not play a role in the next steps. At each iteration, if the constraint  $i \in [m]$  with minimum  
 275 leverage score  $\sigma_i$  falls below a given threshold  $\sigma_{\min}$ , it is removed from the polyhedron. Otherwise,  
 276 we query the volumetric center of the current polyhedron and add the separation hyperplane as a  
 277 constraint to the polyhedron. We bound the number of iterations of the procedure by

$$T(\delta, d) = \left\lceil c \cdot d \left( 1.4 \ln \frac{1}{\delta} + 2 \ln d + 2 \ln(1 + 1/\sigma_{\min}) \right) \right\rceil,$$

278 where  $\sigma_{\min}$  and  $c$  are parameters that will be fixed shortly. Instead of making a call directly to the  
 279 oracle  $O_S$ , we instead suppose that one has access to an oracle  $O : \mathcal{I}_d \rightarrow \mathbb{R}^d$  where  $\mathcal{I}_d = (\mathbb{Z} \times \mathbb{R}^{d+1})^*$   
 280 has exactly the shape of the memory storing the information from the polyhedron. This form of  
 281 oracle is used in our recursive calls to Vaidya's method. For example, such an oracle can simply be  
 282  $O : \mathcal{P} \in \mathcal{I}_d \mapsto O_S(\text{VolumetricCenter}(\mathcal{P}))$ . Last, in our recursive method, we will not assume that  
 283 oracle responses are normalized. As a result, we specify that if the norm of the response is too small,  
 284 we can stop the algorithm. We assume however that the oracle already returns discretized vectors,  
 285 which will be ensured in the following procedures. The cutting-plane algorithm is formally described  
 286 in Algorithm 1. With an appropriate choice of parameters, this procedure finds an approximate  
 287 solution of feasibility problems. We base the constants from [2].

288 **Lemma 4.1.** Fix  $\sigma_{\min} = 0.04$  and  $c = \frac{1}{0.0014} \approx 715$ . Let  $\delta, \xi \in (0, 1)$  and  $O : \mathcal{I}_d \rightarrow \mathbb{R}^d$ . Write  
 289  $\mathcal{P} = \{(k_i, \mathbf{a}_i, b_i), i \in [m]\}$  as the output of Algorithm 1 run with  $O, \delta$  and  $\xi$ . Then,

$$\min_{\substack{\lambda_i \geq 0, i \in [m], \\ \sum_{i \in [m]} \lambda_i = 1}} \max_{\mathbf{y} \in \mathcal{C}_d} \sum_{i=1}^m \lambda_i (\mathbf{a}_i^\top \mathbf{y} - b_i) = \max_{\mathbf{x} \in \mathcal{C}_d} \min_{i \in [m]} (\mathbf{a}_i^\top \mathbf{x} - b_i) \leq \delta.$$

**Input:**  $O : \mathcal{I}_d \rightarrow \mathbb{R}^d$ ,  $\delta, \xi \in (0, 1)$

- 1 Let  $T_{max} = T(\delta, d)$  and initialize  $\mathcal{P}_0 := \{(-1, \mathbf{e}_i, -1), (-1, -\mathbf{e}_i, -1), i \in [d]\}$
- 2 **for**  $t = 0, \dots, T_{max}$  **do**
- 3     **if**  $\{\mathbf{x} : \mathbf{A}\mathbf{x} \geq \mathbf{b}\} = \emptyset$  **then return**  $\mathcal{P}_t$ ;
- 4     **if**  $\min_{i \in [m]} \sigma_i < \sigma_{min}$  **then**
- 5          $\mathcal{P}_{t+1} = \mathcal{P}_t \setminus \{(k_j, \mathbf{a}_j, b_j)\}$  where  $j \in \arg \min_{i \in [m]} \sigma_i$
- 6     **else if**  $\boldsymbol{\omega} := \text{VolumetricCenter}(\mathcal{P}_t) \notin \mathcal{C}_d$  **then**
- 7          $\mathcal{P}_{t+1} = \mathcal{P}_t \cup \{(-1, -\text{sign}(\omega_j)\mathbf{e}_j, -1)\}$  where  $j \in [d]$  has  $|\omega_j| > 1$
- 8     **else**
- 9          $\mathbf{g} = O(\mathcal{P}_t)$  and  $\mathbf{b} = \xi \left[ \frac{\mathbf{g}^\top \boldsymbol{\omega}}{\xi} \right]$ , where  $\boldsymbol{\omega} = \text{VolumetricCenter}(\mathcal{P}_t)$
- 10          $\mathcal{P}_{t+1} = \mathcal{P}_t \cup \{(t, \mathbf{g}, \mathbf{b})\}$
- 11         **if**  $\|\mathbf{g}\| \leq \delta$  **then return**  $\mathcal{P}_{t+1}$ ;
- 12 **end**
- 13 **return**  $\mathcal{P}_{T_{max}+1}$ .

**Algorithm 1:** Memory-constrained Vaidya's volumetric method

290 From now, we use the parameters  $\sigma_{min} = 0.04$  and  $c = 1/0.0014$  as in Lemma 4.1. Since the  
 291 memory of both Vaidya's method and center-of-mass consists primarily of the constraints, we recall  
 292 an important feature of Vaidya's method that the number of constraints at any time is  $\mathcal{O}(d)$ .

293 **Lemma 4.2** ([46, 1, 2]). *At any time while running Algorithm 1, the number of constraints of the*  
 294 *current polyhedron is at most  $\frac{d}{\sigma_{min}} + 1$ .*

## 295 4.2 A recursive algorithm

296 We write  $\mathcal{C}_{m+n} = \mathcal{C}_m \times \mathcal{C}_n$  and aim to apply Vaidya's method to the first  $m$  coordinates. To do so,  
 297 we need to approximate a separation oracle on these  $m$  coordinates only, which corresponds to giving  
 298 separation hyperplanes with small values for the last  $n$  coordinates. This can be achieved using the  
 299 following auxiliary linear program. For  $\mathcal{P} \in \mathcal{I}_n$ , we define

$$\min_{\substack{\lambda_i \geq 0, i \in [m], \\ \sum_{i \in [m]} \lambda_i = 1}} \max_{\mathbf{y} \in \mathcal{C}_n} \sum_{i=1}^m \lambda_i (\mathbf{a}_i^\top \mathbf{y} - b_i), \quad m = |\mathcal{P}| \quad (\mathcal{P}_{aux}(\mathcal{P}))$$

300 where as before,  $\mathbf{A}$  and  $\mathbf{b}$  define the constraints stored in  $\mathcal{P}$ . The procedure to obtain an approximate  
 301 separation oracle on the first  $n$  coordinates  $\mathcal{C}_n$  is given in Algorithm 2 and using Lemma 4.1 we can  
 302 show that this procedure provides approximate separation vectors for the first  $n$  coordinates.

**Input:**  $\delta, \xi, O_x : \mathcal{I}_n \rightarrow \mathbb{R}^m$  and  $O_y : \mathcal{I}_n \rightarrow \mathbb{R}^n$

- 1 Run Algorithm 1 with  $\delta, \xi$  and  $O_y$  to obtain polyhedron  $\mathcal{P}^*$
- 2 Solve  $\mathcal{P}_{aux}(\mathcal{P}^*)$  to get a solution  $\boldsymbol{\lambda}^*$
- 3 Store  $\mathbf{k}^* = (k_i, i \in [m])$  where  $m = |\mathcal{P}^*|$ , and  $\boldsymbol{\lambda}^* \leftarrow \text{Discretize}(\boldsymbol{\lambda}^*, \xi)$
- 4 Initialize  $\mathcal{P}_0 := \{(-1, \mathbf{e}_i, -1), (-1, -\mathbf{e}_i, -1), i \in [d]\}$  and  $\mathbf{u} = \mathbf{0} \in \mathbb{R}^m$
- 5 **for**  $t = 0, 1, \dots, \max_i k_i$  **do**
- 6     **if**  $t = k_i^*$  for some  $i \in [m]$  **then**
- 7          $\mathbf{g}_x = O_x(\mathcal{P}_t)$
- 8          $\mathbf{u} \leftarrow \text{Discretize}_m(\mathbf{u} + \lambda_i^* \mathbf{g}_x, \xi)$
- 9         Update  $\mathcal{P}_t$  to get  $\mathcal{P}_{t+1}$  as in Algorithm 1
- 10 **end**
- 11 **return**  $\mathbf{u}$

**Algorithm 2:**  $\text{ApproxSeparationVector}_{\delta, \xi}(O_x, O_y)$

303 The next step involves using this approximation recursively. We write  $d = \sum_{i=1}^p k_i$ , and interpret  $\mathcal{C}_d$   
 304 as  $\mathcal{C}_{k_1} \times \dots \times \mathcal{C}_{k_p}$ . In particular, for  $\mathbf{x} \in \mathcal{C}_d$ , we write  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_p)$  where  $\mathbf{x}_i \in \mathcal{C}_{k_i}$  for  $i \in [p]$ .  
 305 Applying Algorithm 2 recursively, we can obtain an approximate separation oracle for the first  $i$   
 306 coordinates  $\mathcal{C}_{k_1} \times \dots \times \mathcal{C}_{k_i}$ . However, storing such separation vectors would be too memory-expensive,  
 307 e.g., for  $i = p$ , that would correspond to storing the separation hyperplanes from the oracle  $O_S$



308 directly. Instead, given  $j \in [i]$ , Algorithm 3 recursively computes the  $\mathbf{x}_j$  component of an approximate  
 309 separation oracle for the first  $i$  variables  $(\mathbf{x}_1, \dots, \mathbf{x}_i)$ , via the procedure  $\text{ApproxOracle}(i, j)$ .

**Input:**  $\delta, \xi, 1 \leq j \leq i \leq p, \mathcal{P}^{(r)} \in \mathcal{I}_{k_r}$  for  $r \in [i], O_S : \mathcal{C}_d \rightarrow \mathbb{R}^d$   
 1 **if**  $i = p$  **then**  
 2      $\mathbf{x}_r = \text{VolumetricCenter}(\mathbf{A}_r, \mathbf{b}_r)$  where  $(\mathbf{A}_r, \mathbf{b}_r)$  defines the constraints stored in  $\mathcal{P}^{(r)}$  for  
        $r \in [p]$   
 3      $(\mathbf{g}_1, \dots, \mathbf{g}_p) = O_S(\mathbf{x}_1, \dots, \mathbf{x}_p)$   
 4     **return**  $\text{Discretize}_{k_j}(\mathbf{g}_j, \xi)$   
 5 **end**  
 6 Define  $O_x : \mathcal{I}_{k_{i+1}} \rightarrow \mathbb{R}^{k_j}$  as  $\text{ApproxOracle}_{\delta, \xi, O_f}(i+1, j, \mathcal{P}^{(1)}, \dots, \mathcal{P}^{(i)}, \cdot)$   
 7 Define  $O_y : \mathcal{I}_{k_{i+1}} \rightarrow \mathbb{R}^{k_{i+1}}$  as  $\text{ApproxOracle}_{\delta, \xi, O_f}(i+1, i+1, \mathcal{P}^{(1)}, \dots, \mathcal{P}^{(i)}, \cdot)$   
 8 **return**  $\text{ApproxSeparationVector}_{\delta, \xi}(O_x, O_y)$

**Algorithm 3:**  $\text{ApproxOracle}_{\delta, \xi, O_S}(i, j, \mathcal{P}^{(1)}, \dots, \mathcal{P}^{(i)})$

310 We can then use  $\text{ApproxOracle}_{\delta, \xi, O_S}(1, 1, \cdot)$  to solve the original problem with the memory-  
 311 constrained Vaidya’s method. In Appendix A, we show that taking  $\delta = \frac{\epsilon}{4d}$  and  $\xi = \frac{\sigma_{\min} \epsilon}{32d^{5/2}}$  achieves  
 312 the desired oracle-complexity and memory usage. The final algorithm is given in Algorithm 4.

**Input:**  $\delta, \xi$ , and  $O_S : \mathcal{C}_d \rightarrow \mathbb{R}^d$  a separation oracle

**Check:** Throughout the algorithm, if  $O_S$  returned Success to a query  $\mathbf{x}$ , **return**  $\mathbf{x}$

1 Run Algorithm 1 with parameters  $\delta$  and  $\xi$  and oracle  $\text{ApproxOracle}_{\delta, \xi, O_S}(1, 1, \cdot)$

**Algorithm 4:** Memory-constrained algorithm for convex optimization

313 **Sketch of proof.** At the high level, the algorithm recursively runs Vaidya’s method Algorithm 1  
 314 for each level of computation  $i \in [p]$ . Since each run of Algorithm 4 requires  $\mathcal{O}(\frac{d}{p} \ln \frac{1}{\epsilon})$  queries, the  
 315 total number of calls to the oracle, which is exponential in the number of levels, is  $\mathcal{O}(\mathcal{O}(\frac{d}{p} \ln \frac{1}{\epsilon})^p)$ .  
 316 As for the memory usage, the algorithm mainly needs to keep in memory the constraints defining the  
 317 polyhedrons at each level  $i \in [p]$ . From Lemma 4.2, each polyhedron only requires  $\mathcal{O}(\frac{d}{p})$  constraints  
 318 that each require  $\mathcal{O}(\frac{d}{p} \ln \frac{1}{\epsilon})$  bits of memory. Hence, the total memory needed is  $\mathcal{O}(\frac{d^2}{p} \ln \frac{1}{\epsilon})$ . The main  
 319 difficulty lies in showing that the algorithm is successful. To do so, we need to show that the precision  
 320 in the successive approximated separation oracles Algorithm 2 is sufficient. To avoid an exponential  
 321 dependence of the approximation error in  $p$ —which would be prohibitive for the memory usage of  
 322 our method—each run of Vaidya’s method Algorithm 1 is run for more iterations than the precision of  
 323 the separation vectors would classically allow. To give intuition, if the separation oracle came from a  
 324 convex optimization subgradient oracle for a function  $f$ , the iterates at a level  $i$  do not converge to the  
 325 true “minimizer” of  $\min_{\mathbf{x}_i} f^{(i)}(\mathbf{x}_1, \dots, \mathbf{x}_i)$ , where  $f^{(i)}(\cdot) = \min_{\mathbf{x}_{i+1}, \dots, \mathbf{x}_p} f(\cdot, \mathbf{x}_{i+1}, \dots, \mathbf{x}_p)$ , but  
 326 instead converge to a close enough point while still providing meaningful approximate subgradients  
 327 at the higher level  $i - 1$  (in Algorithm 2).

## 328 5 Discussion and Conclusion

329 To the best of our knowledge, this work is the first to provide some positive trade-off between  
 330 oracle-complexity and memory-usage for convex optimization or the feasibility problem, as opposed  
 331 to lower-bound impossibility results [29, 5]. Our trade-offs are more significant in a high accuracy  
 332 regime: when  $\ln \frac{1}{\epsilon} \approx d^c$ , for  $c > 0$  our trade-offs are polynomial, while the improvements when  
 333  $\ln \frac{1}{\epsilon} = \text{poly}(\ln d)$  are only in  $\ln d$  factors. A natural open direction [49] is whether there exist  
 334 algorithms with polynomial trade-offs in that case. We also show that in the exponential regime  
 335  $\ln \frac{1}{\epsilon} \geq \Omega(d \ln d)$ , gradient descent is not Pareto-optimal. Instead, one can keep the optimal memory  
 336 and decrease the dependence in  $\epsilon$  of the oracle-complexity from  $\frac{1}{\epsilon^2}$  to  $(\ln \frac{1}{\epsilon})^d$ . The question of  
 337 whether the exponential dependence in  $d$  is necessary is left open. Last, our algorithms rely on  
 338 the consistency of the oracle, which allows re-computations. While this is a classical assumption,  
 339 gradient descent and classical cutting-plane methods do not need it; removing this assumption could  
 340 be an interesting research direction (potentially, this could also yield stronger lower bounds).

341 **References**

- 342 [1] Kurt M Anstreicher. “On Vaidya’s volumetric cutting plane method for convex programming”.  
 343 In: *Mathematics of Operations Research* 22.1 (1997), pp. 63–89.
- 344 [2] Kurt M Anstreicher. “Towards a practical volumetric cutting plane method for convex pro-  
 345 gramming”. In: *SIAM Journal on Optimization* 9.1 (1998), pp. 190–206.
- 346 [3] Aharon Ben-Tal and Arkadi Nemirovski. “Non-euclidean restricted memory level method for  
 347 large-scale convex optimization”. In: *Mathematical Programming* 102.3 (Jan. 2005), pp. 407–  
 348 456.
- 349 [4] Dimitris Bertsimas and Santosh Vempala. “Solving convex programs by random walks”. In:  
 350 *Journal of the ACM (JACM)* 51.4 (2004), pp. 540–556.
- 351 [5] Moise Blanchard, Junhui Zhang, and Patrick Jaillet. “Quadratic Memory is Necessary for  
 352 Optimal Query Complexity in Convex Optimization: Center-of-Mass is Pareto-Optimal”. In:  
 353 *arXiv preprint arXiv:2302.04963* (2023).
- 354 [6] Mark Braverman et al. “Communication Lower Bounds for Statistical Estimation Problems  
 355 via a Distributed Data Processing Inequality”. In: *Proceedings of the Forty-Eighth Annual  
 356 ACM Symposium on Theory of Computing*. STOC ’16. Cambridge, MA, USA: Association for  
 357 Computing Machinery, 2016, pp. 1011–1020.
- 358 [7] C. G. Broyden. “The Convergence of a Class of Double-rank Minimization Algorithms 1.  
 359 General Considerations”. In: *IMA Journal of Applied Mathematics* 6.1 (Mar. 1970), pp. 76–90.
- 360 [8] Sébastien Bubeck. “Convex Optimization: Algorithms and Complexity”. In: *Foundations and  
 361 Trends® in Machine Learning* 8.3-4 (2015), pp. 231–357.
- 362 [9] Sébastien Bubeck, Yin Tat Lee, and Mohit Singh. *A geometric alternative to Nesterov’s  
 363 accelerated gradient descent*. 2015.
- 364 [10] Alexandre d’Aspremont, Damien Scieur, and Adrien Taylor. “Acceleration Methods”. In:  
 365 *Foundations and Trends® in Optimization* 5.1-2 (2021), pp. 1–245.
- 366 [11] Yuval Dagan, Gil Kur, and Ohad Shamir. “Space lower bounds for linear prediction in the  
 367 streaming model”. In: *Proceedings of the Thirty-Second Conference on Learning Theory*. Ed.  
 368 by Alina Beygelzimer and Daniel Hsu. Vol. 99. Proceedings of Machine Learning Research.  
 369 PMLR, 25–28 Jun 2019, pp. 929–954.
- 370 [12] Yuval Dagan and Ohad Shamir. “Detecting Correlations with Little Memory and Communica-  
 371 tion”. In: *Proceedings of the 31st Conference On Learning Theory*. Ed. by Sébastien Bubeck,  
 372 Vianney Perchet, and Philippe Rigollet. Vol. 75. Proceedings of Machine Learning Research.  
 373 PMLR, June 2018, pp. 1145–1198.
- 374 [13] James Demmel, Ioana Dumitriu, and Olga Holtz. “Fast linear algebra is stable”. In: *Numerische  
 375 Mathematik* 108.1 (2007), pp. 59–91.
- 376 [14] Kim Donghwan and Jeffrey Fessler. “Optimized first-order methods for smooth convex mini-  
 377 mization”. In: *Mathematical Programming* 159 (June 2014).
- 378 [15] Yoel Drori and Marc Teboulle. “Performance of first-order methods for smooth convex mini-  
 379 mization: a novel approach”. In: *Mathematical Programming* 145.1 (June 2014), pp. 451–  
 380 482.
- 381 [16] Dmitriy Drusvyatskiy, Maryam Fazel, and Scott Roy. “An Optimal First Order Method Based  
 382 on Optimal Quadratic Averaging”. In: *SIAM Journal on Optimization* 28.1 (2018), pp. 251–  
 383 271.
- 384 [17] Uriel Feige and Gideon Schechtman. “On the optimality of the random hyperplane rounding  
 385 technique for MAX CUT”. In: *Random Structures & Algorithms* 20.3 (2002), pp. 403–440.
- 386 [18] R. Fletcher. “A new approach to variable metric algorithms”. In: *The Computer Journal* 13.3  
 387 (Jan. 1970), pp. 317–322.
- 388 [19] Donald Goldfarb. “A Family of Variable-Metric Methods Derived by Variational Means”. In:  
 389 *Mathematics of Computation* 24.109 (1970), pp. 23–26.
- 390 [20] Haotian Jiang et al. “An Improved Cutting Plane Method for Convex Optimization, Convex-  
 391 Concave Games, and Its Applications”. In: *Proceedings of the 52nd Annual ACM SIGACT  
 392 Symposium on Theory of Computing*. STOC 2020. Chicago, IL, USA: Association for Com-  
 393 puting Machinery, 2020, pp. 944–953.

- 394 [21] Jonathan Kelner et al. “Big-Step-Little-Step: Efficient Gradient Methods for Objectives with  
395 Multiple Scales”. In: *Proceedings of Thirty Fifth Conference on Learning Theory*. Ed. by  
396 Po-Ling Loh and Maxim Raginsky. Vol. 178. Proceedings of Machine Learning Research.  
397 PMLR, Feb. 2022, pp. 2431–2540.
- 398 [22] Guanghui Lan. “Bundle-level type methods uniformly optimal for smooth and nonsmooth  
399 convex optimization”. In: *Mathematical Programming* 149.1 (Feb. 2015), pp. 1–45.
- 400 [23] Guanghui Lan, Soomin Lee, and Yi Zhou. “Communication-efficient algorithms for decen-  
401 tralized and stochastic optimization”. In: *Mathematical Programming* 180.1 (Mar. 2020),  
402 pp. 237–284.
- 403 [24] Quoc Le, Alex Smola, and S.V.N. Vishwanathan. “Bundle Methods for Machine Learning”.  
404 In: *Advances in Neural Information Processing Systems*. Ed. by J. Platt et al. Vol. 20. Curran  
405 Associates, Inc., 2007.
- 406 [25] Yin Tat Lee, Aaron Sidford, and Sam Chiu-wai Wong. “A faster cutting plane method and  
407 its implications for combinatorial and convex optimization”. In: *2015 IEEE 56th Annual  
408 Symposium on Foundations of Computer Science*. IEEE, 2015, pp. 1049–1065.
- 409 [26] Claude Lemaréchal, Arkadi Nemirovski, and Yurii Nesterov. “New variants of bundle methods”.  
410 In: *Mathematical Programming* 69.1 (July 1995), pp. 111–147.
- 411 [27] Adrian S. Lewis and Michael L. Overton. “Nonsmooth optimization via quasi-Newton meth-  
412 ods”. In: *Mathematical Programming* 141.1 (Oct. 2013), pp. 135–163.
- 413 [28] Dong C. Liu and Jorge Nocedal. “On the limited memory BFGS method for large scale  
414 optimization”. In: *Mathematical Programming* 45.1 (Aug. 1989), pp. 503–528.
- 415 [29] Annie Marsden et al. “Efficient convex optimization requires superlinear memory”. In: *Confer-  
416 ence on Learning Theory*. PMLR, 2022, pp. 2390–2430.
- 417 [30] Dana Moshkovitz and Michal Moshkovitz. “Mixing Implies Lower Bounds for Space Bounded  
418 Learning”. In: *Proceedings of the 2017 Conference on Learning Theory*. PMLR, 2017,  
419 pp. 1516–1566.
- 420 [31] João F. C. Mota et al. “D-ADMM: A Communication-Efficient Distributed Algorithm for  
421 Separable Optimization”. In: *IEEE Transactions on Signal Processing* 61.10 (2013), pp. 2718–  
422 2723.
- 423 [32] Arkadi Nemirovski and David Borisovich Yudin. *Problem Complexity and Method Efficiency  
424 in Optimization*. A Wiley-Interscience publication. Wiley, 1983.
- 425 [33] Yurii Nesterov. “A method of solving a convex programming problem with convergence rate  
426  $O(1/k^2)$ ”. In: *Dokl. Akad. Nauk SSSR* 269 (3 1983), pp. 543–547.
- 427 [34] Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*. Vol. 87. Springer  
428 Science & Business Media, 2003.
- 429 [35] Jorge Nocedal. “Updating Quasi-Newton Matrices with Limited Storage”. In: *Mathematics of  
430 Computation* 35.151 (1980), pp. 773–782.
- 431 [36] Srinivasan Ramaswamy and John E Mitchell. *A long step cutting plane algorithm that uses the  
432 volumetric barrier*. Tech. rep. Citeseer, 1995.
- 433 [37] Ran Raz. “A Time-Space Lower Bound for a Large Class of Learning Problems”. In: *2017  
434 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*. 2017, pp. 732–  
435 742.
- 436 [38] Sashank J. Reddi et al. “AIDE: Fast and Communication Efficient Distributed Optimization”.  
437 In: *ArXiv abs/1608.06879* (2016).
- 438 [39] B. Van Scoy, R. A. Freeman, and K. M. Lynch. “The Fastest Known Globally Convergent  
439 First-Order Method for Minimizing Strongly Convex Functions”. In: *IEEE Control Systems  
440 Letters* PP.99 (2017), pp. 1–1.
- 441 [40] Ohad Shamir, Nati Srebro, and Tong Zhang. “Communication-Efficient Distributed Optimiza-  
442 tion using an Approximate Newton-type Method”. In: *Proceedings of the 31st International  
443 Conference on Machine Learning*. Ed. by Eric P. Xing and Tony Jebara. Vol. 32. Proceedings  
444 of Machine Learning Research 2. Beijing, China: PMLR, 22–24 Jun 2014, pp. 1000–1008.
- 445 [41] D. F. Shanno. “Conditioning of Quasi-Newton Methods for Function Minimization”. In:  
446 *Mathematics of Computation* 24.111 (1970), pp. 647–656.
- 447 [42] Vatsal Sharan, Aaron Sidford, and Gregory Valiant. “Memory-Sample Tradeoffs for Linear  
448 Regression with Small Error”. In: *Proceedings of the 51st Annual ACM SIGACT Symposium on  
449 Theory of Computing*. STOC 2019. Association for Computing Machinery, 2019, pp. 890–901.

- 450 [43] Virginia Smith et al. “CoCoA: A General Framework for Communication-Efficient Distributed  
451 Optimization”. In: *J. Mach. Learn. Res.* 18.1 (Jan. 2017), pp. 8590–8638.
- 452 [44] Adrien Taylor and Yoel Drori. “An optimal gradient method for smooth strongly convex  
453 minimization”. In: *Mathematical Programming* 199.1 (May 2023), pp. 557–594.
- 454 [45] Choon Hui Teo et al. “Bundle Methods for Regularized Risk Minimization”. In: *Journal of  
455 Machine Learning Research* 11.10 (2010), pp. 311–365.
- 456 [46] Pravin M Vaidya. “A new algorithm for minimizing convex functions over convex sets”. In:  
457 *Mathematical programming* 73.3 (1996), pp. 291–341.
- 458 [47] Jialei Wang, Weiran Wang, and Nathan Srebro. “Memory and Communication Efficient Dis-  
459 tributed Stochastic Optimization with Minibatch Prox”. In: *Proceedings of the 2017 Conference  
460 on Learning Theory*. Ed. by Satyen Kale and Ohad Shamir. Vol. 65. Proceedings of Machine  
461 Learning Research. PMLR, July 2017, pp. 1882–1919.
- 462 [48] Jianqiao Wang et al. “Gradient Sparsification for Communication-Efficient Distributed  
463 Optimization”. In: *Proceedings of the 32nd International Conference on Neural Information  
464 Processing Systems*. NIPS’18. Montréal, Canada: Curran Associates Inc., 2018, pp. 1306–  
465 1316.
- 466 [49] Blake Woodworth and Nathan Srebro. “Open problem: The oracle complexity of convex opti-  
467 mization with limited memory”. In: *Conference on Learning Theory*. PMLR. 2019, pp. 3202–  
468 3210.
- 469 [50] Yuchen Zhang, John C. Duchi, and Martin J. Wainwright. “Communication-efficient algorithms  
470 for statistical optimization”. In: *2012 IEEE 51st IEEE Conference on Decision and Control  
471 (CDC)*. 2012, pp. 6792–6792.

472 **A Proof of the query complexity and memory usage of Algorithm 4**

473 First, we give simple properties on the discretization functions. One can easily check that for any  
474  $\mathbf{x} \in C$ ,

$$\|\mathbf{x} - \text{Discretize}_d(\mathbf{x}, \xi)\| \leq \xi \quad \text{and} \quad \|\text{Discretize}_d(\mathbf{x}, \xi)\| \leq \|\mathbf{x}\|. \quad (1)$$

475 Further, one can easily check that to represent any output of  $\text{Discretize}_d(\cdot, \xi)$ , one needs at most  
476  $d \ln \frac{2\sqrt{d}}{\xi} = \mathcal{O}(d \ln \frac{d}{\xi})$  bits.

477 We next prove Lemma 4.1.

478 *Proof of Lemma 4.1.* We first consider the case when the algorithm terminates because of a query  
479  $\mathbf{g} = O(\mathcal{P}_t)$  such that  $\|\mathbf{g}\| \leq \delta/(2\sqrt{d})$ . Then, for any  $\mathbf{x} \in \mathcal{C}_d$ , one directly has

$$\mathbf{g}^\top \mathbf{x} - b \leq \mathbf{g}^\top (\mathbf{x} - \boldsymbol{\omega}) \leq 2\sqrt{d}\|\mathbf{g}\| \leq \delta.$$

480 where  $\boldsymbol{\omega}$  is the volumetric center of the resulting polyhedron. In the second inequality we used the  
481 fact that  $\boldsymbol{\omega} \in \mathcal{C}_d$ , otherwise the algorithm would not have terminated at that step.

482 We next turn to the other cases and start by showing that the output polyhedron does not contain a  
483 ball of radius  $\delta$ . This is immediate if the algorithm terminated because the polyhedron was empty.  
484 We then suppose this was not the case, and follow the same proof as given in [2]. Algorithm 1  
485 and the one provided in [2] coincide when removing a constraint of the polyhedron. Hence, it  
486 suffices to consider the case when we add a constraint. We use the notation  $\tilde{\mathbf{A}}^\top = [\mathbf{A}^\top, \mathbf{a}_{m+1}^\top]$ ,  
487  $\tilde{\mathbf{b}}^\top = [\mathbf{b}^\top, b_{m+1}]$  for the updated matrix  $\mathbf{A}$  and vector  $\mathbf{b}$  after adding the constraint. We also  
488 denote  $\boldsymbol{\omega} = \text{VolumetricCenter}(\mathbf{A}, \mathbf{b})$  (resp.  $\tilde{\boldsymbol{\omega}} = \text{VolumetricCenter}(\tilde{\mathbf{A}}, \tilde{\mathbf{b}})$ ) the volumetric center of  
489 the polyhedron before (resp. after) adding the constraint. Next, we consider the vector  $(\mathbf{b}')^\top =$   
490  $[\mathbf{b}^\top, \mathbf{a}_{m+1}^\top \boldsymbol{\omega}]$ , which would have been obtained if the cut was performed at  $\boldsymbol{\omega}$  exactly. We then denote  
491  $\boldsymbol{\omega}' = \text{VolumetricCenter}(\tilde{\mathbf{A}}, \mathbf{b}')$ . Then proof of [2] shows that

$$V_{\tilde{\mathbf{A}}, \mathbf{b}'}(\boldsymbol{\omega}') \geq V_{\mathbf{A}, \mathbf{b}}(\boldsymbol{\omega}) + 0.0340.$$

492 We now observe that by construction, we have  $\tilde{\mathbf{b}}_{m+1} \geq \mathbf{a}_{m+1}^\top \boldsymbol{\omega}$ , so that the polyhedron associated  
493 to  $(\tilde{\mathbf{A}}, \tilde{\mathbf{b}})$  is more constrained than the one associated to  $(\tilde{\mathbf{A}}, \mathbf{b}')$ . As a result, we have  $V_{\tilde{\mathbf{A}}, \tilde{\mathbf{b}}}(\mathbf{x}) \geq$   
494  $V_{\tilde{\mathbf{A}}, \mathbf{b}'}(\mathbf{x})$ , for any  $\mathbf{x} \in \mathbb{R}^d$  such that  $\tilde{\mathbf{A}}\mathbf{x} \geq \tilde{\mathbf{b}}$ . Therefore,

$$V_{\tilde{\mathbf{A}}, \tilde{\mathbf{b}}}(\tilde{\boldsymbol{\omega}}) \geq V_{\tilde{\mathbf{A}}, \mathbf{b}'}(\tilde{\boldsymbol{\omega}}) \geq V_{\tilde{\mathbf{A}}, \mathbf{b}'}(\boldsymbol{\omega}') \geq V_{\mathbf{A}, \mathbf{b}}(\boldsymbol{\omega}) + 0.0340.$$

495 This ends the modifications in the proof of [2]. With the notations of this paper, we still have  
496  $\Delta V^+ = 0.340$  and  $\Delta V^- = 0.326$ , so that  $\Delta V = 0.0014$ . Then, because  $c = \frac{1}{\Delta V}$ , the same  
497 proof shows that the procedure is successful for precision  $\delta$ : the final polyhedron  $(\mathbf{A}, \tilde{\mathbf{b}})$  returned by  
498 Algorithm 1 does not contains a ball of radius  $> \delta$ . As a result, whether the algorithm performed  
499 all  $T_{max}$  iterations or not,  $\{\mathbf{x} : \mathbf{A}\mathbf{x} \geq \tilde{\mathbf{b}}\}$  does not contain a ball of radius  $> \delta'$ , where  $\mathbf{A}$  and  
500  $\tilde{\mathbf{b}}$  define the constraints stored in the output  $\mathcal{P}$ . Now letting  $m$  be the objective value of the right  
501 optimization problem, there exists  $\mathbf{x} \in \mathcal{C}_d$  such that for all  $t \leq T$ ,  $\mathbf{g}_t^\top (\mathbf{x} - \mathbf{c}_t) \geq m$ . Therefore, for  
502 any  $\mathbf{x}' \in B_d(\mathbf{x}, m)$  one has

$$\forall i \in [m], \mathbf{a}_i^\top \mathbf{x}' - b_i \geq m + \mathbf{a}_i^\top (\mathbf{x}' - \mathbf{x}) \geq m - \|\mathbf{x}' - \mathbf{x}\| \geq 0.$$

503 In the last inequality we used  $\|\mathbf{a}_t\| \leq 1$ . This implies that the polyhedron contains  $B_d(\mathbf{x}, m)$ . Hence,  
504  $m \leq \delta$ .

505 This ends the proof of the right inequality. The left equality is a direct application of strong duality  
506 for linear programming.  $\square$

507 We now prove that Algorithm 4 has the desired oracle-complexity and memory usage.

508 We first describe the recursive calls of Algorithm 3 in more detail. To do so, consider running the pro-  
509 cedure  $\text{ApproxOracle}(i, j, \mathcal{P}^{(1)}, \dots, \mathcal{P}^{(i)})$  where  $i < p$ , which corresponds to running Algorithm 2  
510 for specific oracles. We say that this is a level- $i$  run. Then, the algorithm performs at most  $2T(\delta, k_{i+1})$   
511 calls to  $\text{ApproxOracle}(i+1, i+1, \mathcal{P}^{(1)}, \dots, \mathcal{P}^{(i)}, \cdot)$ , where the factor 2 comes from the fact that

512 Vaidya's method Algorithm 1 is effectively run twice in Algorithm 2. The solution to  $(\mathcal{P}_{aux}(\mathcal{P}))$  has  
 513 as many components as constraints in the last polyhedron, which is at most  $\frac{k_{i+1}}{\sigma_{min}} + 1$  by Lemma 4.2.  
 514 Hence, the number of calls to  $\text{ApproxOracle}(i+1, j, \mathcal{P}^{(1)}, \dots, \mathcal{P}^{(i)}, \cdot)$  is at most  $\frac{k_{i+1}}{\sigma_{min}} + 1$ . In total,  
 515 that is  $\mathcal{O}(k_{i+1} \ln \frac{1}{\delta})$  calls to the level  $i+1$  of the recursion.

516 We next aim to understand the output of running  $\text{ApproxOracle}(1, 1, \mathcal{P}^{(1)})$ . We denote by  $\lambda(\mathcal{P}^{(1)})$   
 517 the solution  $\mathcal{P}_{aux}(\mathcal{P}^*)$  computed at 1.2 of the first call to Algorithm 2, where  $\mathcal{P}^*$  is the output  
 518 polyhedron of the first call to Algorithm 1. Denote by  $\mathcal{S}(\mathcal{P}^{(1)})$  the set of indices of coordinates from  
 519  $\lambda(\mathcal{P}^{(1)})$  for which the procedure performed a call to  $\text{ApproxOracle}(2, 1, \mathcal{P}^{(1)}, \cdot)$ . In other words,  
 520  $\mathcal{S}(\mathcal{P}^{(1)})$  contains the indices of all coordinates of  $\lambda(\mathcal{P}^{(1)})$ , except those for which the corresponding  
 521 query lay outside of the unit cube, or the initial constraints of the cube. For any index  $l \in \mathcal{S}(\mathcal{P}^{(1)})$ ,  
 522 let  $\mathcal{P}_l^{(2)}$  denote the state of the current polyhedron ( $\mathcal{P}_l$  in 1.7 of Algorithm 2) when that call was  
 523 performed. Up to discretization issues, the output of the complete procedure is

$$\sum_{l \in \mathcal{S}(\mathcal{P}^{(1)})} \lambda_l(\mathcal{P}^{(1)}) \text{ApproxOracle}(2, 1, \mathcal{P}^{(1)}, \mathcal{P}_l^{(2)}).$$

524 We continue in the recursion, defining  $\lambda(\mathcal{P}^{(1)}, \mathcal{P}_l^{(2)})$  and  $\mathcal{S}(\mathcal{P}^{(1)}, \mathcal{P}_l^{(2)})$  for all  $l \in \mathcal{S}(\mathcal{P}^{(1)})$ , until we  
 525 define all vectors of the form  $\lambda(\mathcal{P}^{(1)}, \mathcal{P}_{l_2}^{(2)}, \dots, \mathcal{P}_{l_r}^{(r)})$  and sets of the form  $\mathcal{S}(\mathcal{P}^{(1)}, \mathcal{P}_{l_2}^{(2)}, \dots, \mathcal{P}_{l_r}^{(r)})$   
 526 for  $i+1 \leq r \leq p-1$ . To simplify the notation and emphasize that all these polyhedra depend on the  
 527 recursive computation path, we adopt the notation

$$\begin{aligned} \lambda^{l_2, \dots, l_r} &:= \lambda_{l_{r+1}}(\mathcal{P}^{(1)}, \mathcal{P}_{l_2}^{(2)}, \dots, \mathcal{P}_{l_r}^{(r)}) \\ \mathcal{S}^{l_2, \dots, l_r} &:= \mathcal{S}(\mathcal{P}^{(1)}, \mathcal{P}_{l_2}^{(2)}, \dots, \mathcal{P}_{l_r}^{(r)}) \end{aligned}$$

528 We recall that these polyhedron are kept in memory to query their volumetric center. For ease of  
 529 notation, we write  $\mathbf{x}_1 = \text{VolumetricCenter}(\mathcal{P}^{(1)})$ , and we write  $\mathbf{c}^{l_2, \dots, l_r} = \text{VolumetricCenter}(\mathcal{P}_{l_r}^{(r)})$   
 530 for  $2 \leq r \leq p$ , where  $l_2, \dots, l_{r-1}$  were the indices from the computation path leading up to  $\mathcal{P}_{l_r}^{(r)}$ .  
 531 Last, we write  $O_S = (O_{S,1}, \dots, O_{S,p})$ , where  $O_{S,i} : \mathcal{C}_d \rightarrow \mathbb{R}^{k_i}$  is the “ $\mathbf{x}_i$ ” component of  $O_S$ , for all  
 532  $i \in [p]$ .

533 With all these notations, we will show that the output of  $\text{ApproxOracle}(i, j, \mathcal{P}^{(1)}, \mathcal{P}_{l_2}^{(2)}, \dots, \mathcal{P}_{l_i}^{(i)})$  is  
 534 approximately equal to the vector

$$\begin{aligned} G(i, j, \mathbf{x}_1, \mathbf{c}^{l_2}, \dots, \mathbf{c}^{l_2, \dots, l_i}) \\ := \sum_{\substack{l_{i+1} \in \mathcal{S}, l_{i+2} \in \mathcal{S}^{l_{i+1}}, \\ \dots, l_p \in \mathcal{S}^{l_{i+1}, \dots, l_{p-1}}} \lambda^{l_{i+1}} \lambda^{l_{i+1}, l_{i+2}} \dots \lambda^{l_{i+1}, \dots, l_p} \cdot O_{S,j}(\mathbf{x}_1, \mathbf{c}^{l_2}, \dots, \mathbf{c}^{l_2, \dots, l_p}), \end{aligned}$$

535 with the convention that for  $i = p$ ,

$$G(p, j, \mathbf{x}_1, \mathbf{c}^{l_2}, \dots, \mathbf{c}^{l_2, \dots, l_p}) := O_{S,j}(\mathbf{x}_1, \mathbf{c}^{l_2}, \dots, \mathbf{c}^{l_2, \dots, l_p}).$$

536 The corresponding computation tree is represented in Fig. 2. For convenience, we omitted the term  
 537  $j = 1$ .

538 We start the analysis with a simple result showing that if the oracle  $O_S$  returns separation vectors of  
 539 norm bounded by one, then the responses from  $\text{ApproxOracle}$  also lie in the unit ball.

540 **Lemma A.1.** Fix  $\delta, \xi \in (0, 1)$ ,  $1 \leq j \leq i \leq p$  and an oracle  $O_S = (O_{S,1}, \dots, O_{S,p}) : \mathcal{C}_d \rightarrow \mathbb{R}^d$ .

541 Suppose that  $O_S$  takes values in the unit ball. For any  $s \in [i]$  let  $\mathcal{P}_{l_s}^{(s)} \in \mathcal{I}_{k_s}$  represent a bounded

542 polyhedrons with  $\text{VolumetricCenter}(\mathcal{P}_{l_s}^{(s)}) \in \mathcal{C}_{k_s}$ . Then, one has

$$\|\text{ApproxOracle}_{\delta, \xi, O_S}(i, j, \mathcal{P}_{l_1}^{(1)}, \dots, \mathcal{P}_{l_i}^{(i)})\| \leq 1.$$

543 *Proof.* We prove this by simple induction on  $i$ . For convenience, we define the point  $\mathbf{x}_k =$   
 544  $\text{VolumetricCenter}(\mathcal{P}_{l_k}^{(k)})$ . If  $i = p$ , we have

$$\begin{aligned} \|\text{ApproxOracle}_{\delta, \xi, O_S}(i, j, \mathcal{P}_{l_1}^{(1)}, \dots, \mathcal{P}_{l_i}^{(i)})\| &= \|\text{Discretize}_{k_j}(O_{S,j}(\mathbf{x}_1, \dots, \mathbf{x}_p), \xi)\| \\ &\leq \|O_{S,j}(\mathbf{x}_1, \dots, \mathbf{x}_p)\| \leq 1, \end{aligned}$$

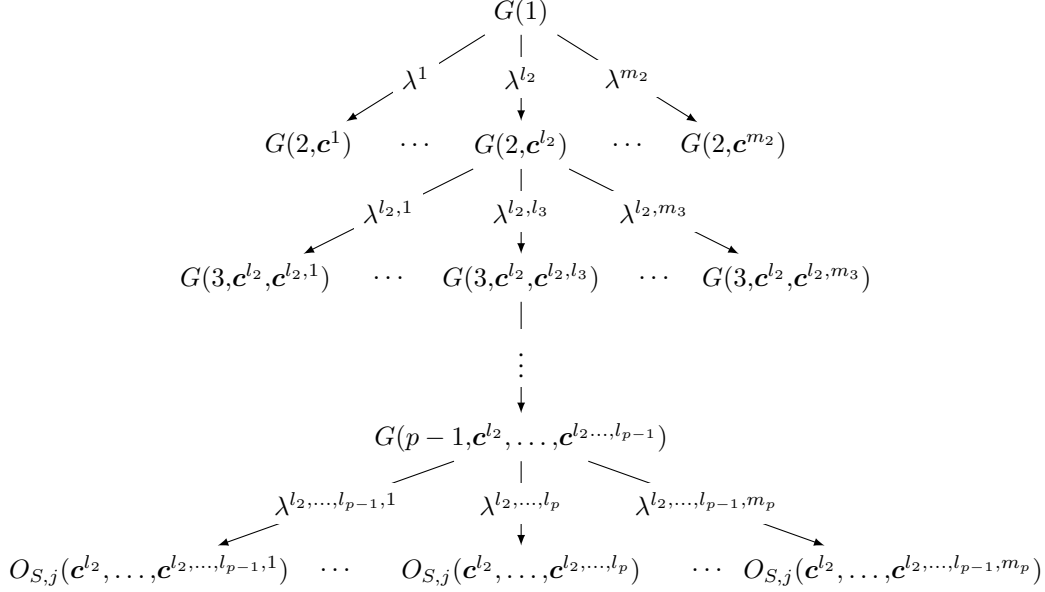


Figure 2: Computation tree representing the recursive calls to `ApproxOracle` starting from the calls to `ApproxOracle(1, 1, ·)` from Algorithm 4

545 where in the first inequality we used Eq (1) and in the second inequality we used the fact that  
 546  $O_S(\mathbf{x}_1, \dots, \mathbf{x}_p)$  has norm at most one. Now suppose that the result holds for  $i + 1 \leq p$ . Then by construction,  
 547 the output  $\text{ApproxOracle}_{\delta, \xi, O_S}(i, j, \mathcal{P}_{l_1}^{(1)}, \dots, \mathcal{P}_{l_i}^{(i)})$  is the result of iterative discretizations.  
 548 Using Eq (1) and the previously defined notations, we obtain

$$\begin{aligned} & \|\text{ApproxOracle}_{\delta, \xi, O_S}(i, j, \mathcal{P}_{l_1}^{(1)}, \dots, \mathcal{P}_{l_i}^{(i)})\| \\ & \leq \left\| \sum_{l_{i+1} \in \mathcal{S}^{l_1, \dots, l_i}} \lambda^{l_2, \dots, l_i} \text{ApproxOracle}_{\delta, \xi, O_S}(i+1, j, \mathcal{P}_{l_1}^{(1)}, \dots, \mathcal{P}_{l_i}^{(i)}, \mathcal{P}_{l_{i+1}}^{(i+1)}) \right\| \leq 1. \end{aligned}$$

549 In the last inequality, we used the induction hypothesis together with the fact that  $\sum_{l_{i+1}} \lambda^{l_2, \dots, l_{i+1}} \leq 1$   
 550 using Eq (1). This ends the induction and the proof.  $\square$

551 We are now ready to compare the output of Algorithm 3 to  $G(i, j, \mathbf{x}_1, \mathbf{c}^{l_2}, \dots, \mathbf{c}^{l_2, \dots, l_i})$ .

552 **Lemma A.2.** Fix  $\delta, \xi \in (0, 1)$ ,  $1 \leq j \leq i \leq p$  and an oracle  $O_S = (O_{S,1}, \dots, O_{S,p}) : \mathcal{C}_d \rightarrow \mathbb{R}^d$ .  
 553 Suppose that  $O_S$  takes values in the unit ball. For any  $s \in [i]$  let  $\mathcal{P}_{l_s}^{(s)} \in \mathcal{I}_{k_s}$  represent a bounded  
 554 polyhedron with  $\text{VolumetricCenter}(\mathcal{P}_{l_s}^{(s)}) \in \mathcal{C}_{k_s}$ . Denote  $\mathbf{x}_r = \mathbf{c}(\mathcal{P}_{l_r}^{(r)})$  for  $r \in [i]$ . Then,

$$\|\text{ApproxOracle}_{\delta, \xi, O_S}(i, j, \mathcal{P}_{l_1}^{(1)}, \dots, \mathcal{P}_{l_i}^{(i)}) - G(i, j, \mathbf{x}_1, \dots, \mathbf{x}_i)\| \leq \frac{4}{\sigma_{\min}} d\xi.$$

555 *Proof.* We prove by simple induction on  $i$  that

$$\begin{aligned} & \|\text{ApproxOracle}_{\delta, \xi, O_S}(i, j, \mathcal{P}_{l_1}^{(1)}, \dots, \mathcal{P}_{l_i}^{(i)}) - G(i, j, \mathbf{x}_1, \dots, \mathbf{x}_i)\| \\ & \leq \left( 1 + \frac{2}{\sigma_{\min}} (k_{i+1} + \dots + k_p) + 2(p-i) \right) \xi. \end{aligned}$$

556 First, for  $i = p$ , the result is immediate since the discretization is with precision  $\xi$  (1.4 of Algorithm 3).  
 557 Now suppose that this is the case for  $i \leq p$  and any valid values of other parameters. For conciseness,  
 558 we write  $\mathbf{G} = (\mathcal{P}_{l_1}^{(1)}, \dots, \mathcal{P}_{l_{i-1}}^{(i-1)})$ . Next, recall that by Lemma 4.2,  $|\mathcal{S}^{l_2, \dots, l_{i-1}}| \leq \frac{k_i}{\sigma_{\min}} + 1$ . Hence,

559 the discretizations due to 1.8 of Algorithm 2 can affect the estimate for at most that number of rounds.  
 560 Then, we have

$$\left\| \text{ApproxOracle}_{\delta, \xi, O_S}(i-1, j, \mathbf{G}) - \sum_{l_i \in \mathcal{S}^{l_2, \dots, l_{i-1}}} \tilde{\lambda}^{l_2, \dots, l_i} \text{ApproxOracle}_{\delta, \xi, O_S}(i, j, \mathbf{G}, \mathcal{P}_{l_i}^{(i)}) \right\| \leq \left( \frac{k_i}{\sigma_{\min}} + 1 \right) \xi,$$

561 where  $\tilde{\lambda}^{l_2, \dots, l_i}$  are the discretized coefficients that are used during the computation 1.8 of Algorithm 2.  
 562 Now using Lemma A.1, we have

$$\left\| \sum_{l_i \in \mathcal{S}^{l_2, \dots, l_{i-1}}} (\tilde{\lambda}^{l_2, \dots, l_i} - \lambda^{l_2, \dots, l_i}) \text{ApproxOracle}_{\delta, \xi, O_S}(i, j, \mathbf{G}, \mathcal{P}_{l_i}^{(i)}) \right\| \leq \|\tilde{\lambda}^{l_{i+1}, \dots, l_{i-1}} - \lambda^{l_{i+1}, \dots, l_{i-1}}\|_1 \leq \left( \frac{k_i}{\sigma_{\min}} + 1 \right) \xi.$$

563 In the last inequality we used the fact that  $\lambda$  has at most  $\frac{k_i}{\sigma_{\min}} + 1$  non-zero coefficients. As a result,  
 564 using the induction for each term of the sum, and the fact that  $\sum_{l_i} \lambda^{l_2, \dots, l_i} \leq 1$ , we obtain

$$\begin{aligned} & \|\text{ApproxOracle}_{\delta, \xi, O_S}(i-1, j, \mathbf{G}) - G(i-1, j, \mathbf{x}_1, \dots, \mathbf{x}_{i-1})\| \\ & \leq \left( 1 + \frac{2}{\sigma_{\min}} (k_{i+1} + \dots + k_p) + 2(p-i) \right) \xi + \left( \frac{2k_i}{\sigma_{\min}} + 2 \right) \xi, \end{aligned}$$

565 which completes the induction. Noting that  $k_{i+1} + \dots + k_p \leq k_1 + \dots + k_p \leq d$  and  $p-i \leq d-1$   
 566 ends the proof.  $\square$

567 Next, we show that the outputs of Algorithm 3 provide approximate separation hyperplanes for the  
 568 first  $i$  coordinates  $(\mathbf{x}_1, \dots, \mathbf{x}_i)$ .

569 **Lemma A.3.** Fix  $\delta, \xi \in (0, 1)$ ,  $1 \leq j \leq i \leq p$  and an oracle  $O_S = (O_{S,1}, \dots, O_{S,p}) : \mathcal{C}_d \rightarrow \mathbb{R}^d$  for  
 570 accuracy  $\epsilon > 0$ . Suppose that  $O_S$  takes values in the unit ball  $B_d(0, 1)$ . For any  $s \in [i]$  let  $\mathcal{P}_{l_s}^{(s)} \in \mathcal{I}_{k_s}$   
 571 represent a bounded polyhedron with  $\text{VolumetricCenter}(\mathcal{P}_{l_s}^{(s)}) \in \mathcal{C}_{k_s}$ . Denote  $\mathbf{x}_r = \mathbf{c}(\mathcal{P}_{l_r}^{(r)})$  for  
 572  $r \in [i]$ . Suppose that when running  $\text{ApproxOracle}_{\delta, \xi, O_S}(i, i, \mathcal{P}_{l_1}^{(1)}, \dots, \mathcal{P}_{l_i}^{(i)})$ , no successful vector  
 573 was queried. Then, any vector  $\mathbf{x}^* = (\mathbf{x}_1^*, \dots, \mathbf{x}_p^*) \in \mathcal{C}_d$  such that  $B_d(\mathbf{x}^*, \epsilon)$  is contained in the  
 574 successful set satisfies

$$\sum_{r \in [i]} \text{ApproxOracle}_{\delta, \xi, O_S}(i, r, \mathcal{P}_{l_1}^{(1)}, \dots, \mathcal{P}_{l_i}^{(i)})^\top (\mathbf{x}_r^* - \mathbf{x}_r) \geq \epsilon - \frac{8d^{5/2}}{\sigma_{\min}} \xi - d\delta.$$

575 *Proof.* For  $i \leq r \leq p$  and  $j \leq r$ , we use the notation

$$\mathbf{g}_j^{l_{i+1}, \dots, l_r} = \text{ApproxOracle}_{\delta, \xi, O_S}(r, j, \mathcal{P}_{l_1}^{(1)}, \dots, \mathcal{P}_{l_r}^{(r)}).$$

576 Using Lemma A.2, we always have for  $j \in [r]$ ,

$$\|\mathbf{g}_j^{l_{i+1}, \dots, l_r} - G(r, j, \mathbf{x}_1, \dots, \mathbf{x}_i, \mathbf{c}^{l_{i+1}}, \dots, \mathbf{c}^{l_{i+1}, \dots, l_r})\| \leq \frac{4d}{\sigma_{\min}} \xi. \quad (2)$$

577 Also, observe that by Lemma A.1 the recursive outputs of  $\text{ApproxOracle}$  always have norm bounded  
 578 by one.

579 Next, let  $\mathcal{T}^{l_{i+1}, \dots, l_{r-1}}$  be the set of indices corresponding to coordinates of  $\lambda^{l_{i+1}, \dots, l_{r-1}}$  for which the  
 580 procedure  $\text{ApproxOracle}$  did not call for a level- $r$  computation. These correspond to 1. constraints  
 581 from the initial cube  $\mathcal{P}_0$ , or 2. cases when the volumetric center was out of the unit cube (1.6-7 of  
 582 Algorithm 1) and as a result, the index of the added constraint was  $-1$  instead of the current iteration  
 583 index  $t$ . Similarly as above, for any  $t \in \mathcal{T}^{l_{i+1}, \dots, l_{r-1}}$ , we denote by  $\mathbf{g}_r^{l_{i+1}, \dots, l_{r-1}, t}$  the corresponding



584 vector  $\mathbf{a}_t$ . We recall that by construction, this vector is of the form  $\pm \mathbf{e}_j$  for some  $j \in [k_r]$ . Then, from  
 585 Lemma 4.1, since the responses of the oracle always have norm bounded by one, for all  $\mathbf{y}_r \in \mathcal{C}_{k_r}$ ,

$$\sum_{l_r \in \mathcal{S}^{l_{i+1}, \dots, l_{r-1}} \cup \mathcal{T}^{l_{i+1}, \dots, l_{r-1}}} \lambda^{l_{i+1}, \dots, l_r} (\mathbf{g}_r^{l_{i+1}, \dots, l_r})^\top (\mathbf{y}_r - \mathbf{c}^{l_{i+1}, \dots, l_r}) \leq \delta. \quad (3)$$

586 For conciseness, we use the shorthand  $(\mathcal{S} \cup \mathcal{T})^{l_{i+1}, \dots, l_{r-1}} := \mathcal{S}^{l_{i+1}, \dots, l_{r-1}} \cup \mathcal{T}^{l_{i+1}, \dots, l_{r-1}}$ , which  
 587 contains all indices from coordinates of  $\lambda^{l_{i+1}, \dots, l_{r-1}}$ . In particular,

$$\sum_{l_r \in (\mathcal{S} \cup \mathcal{T})^{l_{i+1}, \dots, l_{r-1}}} \lambda^{l_{i+1}, \dots, l_r} = 1. \quad (4)$$

588 We now proceed to estimate the precision of the vectors  $G(i, j, \mathbf{x}_1, \dots, \mathbf{x}_i)$  as approximate separation  
 589 hyperplanes for coordinates  $(\mathbf{x}_1, \dots, \mathbf{x}_i)$ . Let  $\mathbf{x}^* \in \mathcal{C}_d$  such that  $B_d(\mathbf{x}^*, \epsilon)$  is within the successful  
 590 set. Then, for any choice of  $l_{i+1} \in \mathcal{S}, \dots, l_p \in \mathcal{S}^{l_{i+1}, \dots, l_{p-1}}$ , since we did not query a successful  
 591 vector, we have for all  $\mathbf{z} \in B_d(\mathbf{x}^*, \epsilon)$ ,

$$O_S(\mathbf{x}_1, \dots, \mathbf{x}_i, \mathbf{c}^{l_{i+1}}, \dots, \mathbf{c}^{l_{i+1}, \dots, l_p})^\top (\mathbf{z} - (\mathbf{x}_1, \dots, \mathbf{x}_i, \mathbf{c}^{l_{i+1}}, \dots, \mathbf{c}^{l_{i+1}, \dots, l_p})) \geq 0.$$

592 As a result, because the responses from  $O_S$  have unit norm,

$$O_S(\mathbf{x}_1, \dots, \mathbf{x}_i, \mathbf{c}^{l_{i+1}}, \dots, \mathbf{c}^{l_{i+1}, \dots, l_p})^\top (\mathbf{x}^* - (\mathbf{x}_1, \dots, \mathbf{x}_i, \mathbf{c}^{l_{i+1}}, \dots, \mathbf{c}^{l_{i+1}, \dots, l_p})) \geq \epsilon. \quad (5)$$

593 Now write  $\mathbf{x}^* = (\mathbf{x}_1^*, \dots, \mathbf{x}_p^*)$ . In addition to the previous equation, for  $l_{i+1} \in \mathcal{S}, \dots, l_{r-1} \in$   
 594  $\mathcal{S}^{l_{i+1}, \dots, l_{r-2}}$  and any  $l_r \in \mathcal{T}^{l_{i+1}, \dots, l_{r-1}}$ , one has  $(\mathbf{g}_r^{l_{i+1}, \dots, l_r})^\top \mathbf{x}_r^* + 1 \geq \epsilon$ , because  $\mathbf{x}^*$  is within the  
 595 cube  $\mathcal{C}_d$  and at least at distance  $\epsilon$  from the constraints of the cube. Similarly as when  $l_r \in \mathcal{S}^{l_{i+1}, \dots, l_{r-1}}$ ,  
 596 for any  $l_r \in \mathcal{T}^{l_{i+1}, \dots, l_{r-1}}$  we denote by  $\mathbf{c}^{l_{i+1}, \dots, l_r}$  the volumetric center of the polyhedron  $\mathcal{P}_{l_r}^{(r)}$  along  
 597 the corresponding computation path, if  $l_r$  corresponded to an added constraints when  $\mathbf{c}^{l_{i+1}, \dots, l_r} \notin \mathcal{C}_{k_r}$ .  
 598 Otherwise, if  $l_r$  corresponded to the constraint  $\mathbf{a} = \pm \mathbf{e}_j$  of the initial cube, we pose  $\mathbf{c}^{l_{i+1}, \dots, l_r} = -\mathbf{a}$ .  
 599 Now by construction, in both cases one has  $(\mathbf{g}_r^{l_{i+1}, \dots, l_r})^\top \mathbf{c}^{l_{i+1}, \dots, l_r} \leq -1$  (I.7 of Algorithm 1). Thus,

$$(\mathbf{g}_r^{l_{i+1}, \dots, l_r})^\top (\mathbf{x}_r^* - \mathbf{c}^{l_{i+1}, \dots, l_r}) \geq \epsilon. \quad (6)$$

600 Recalling Eq (4), we then sum all equations of the form Eq (5) and Eq (6) along the computation  
 601 path, to obtain

$$\begin{aligned} (A) := & \sum_{\substack{l_{i+1} \in \mathcal{S}, \dots, \\ l_p \in \mathcal{S}^{l_{i+1}, \dots, l_{p-1}}}} \lambda^{l_{i+1}} \dots \lambda^{l_{i+1}, \dots, l_p} \\ & \cdot O_S(\mathbf{x}_1, \dots, \mathbf{x}_i, \mathbf{c}^{l_{i+1}}, \dots, \mathbf{c}^{l_{i+1}, \dots, l_p})^\top (\mathbf{x}^* - (\mathbf{x}_1, \dots, \mathbf{x}_i, \mathbf{c}^{l_{i+1}}, \dots, \mathbf{c}^{l_{i+1}, \dots, l_p})) \\ + & \sum_{i+1 \leq r \leq p} \sum_{\substack{l_{i+1} \in \mathcal{S}, \dots, l_{r-1} \in \mathcal{S}^{l_{i+1}, \dots, l_{r-2}}, \\ l_r \in \mathcal{T}^{l_{i+1}, \dots, l_{r-1}}}} \lambda^{l_{i+1}} \dots \lambda^{l_{i+1}, \dots, l_r} \cdot (\mathbf{g}_r^{l_{i+1}, \dots, l_r})^\top (\mathbf{x}_r^* - \mathbf{c}^{l_{i+1}, \dots, l_r}) \geq \epsilon. \end{aligned}$$

602 Now using the convention

$$G(r, r, \mathbf{x}_1, \dots, \mathbf{x}_i, \mathbf{c}^{l_{i+1}}, \dots, \mathbf{c}^{l_{i+1}, \dots, l_r}) := \mathbf{g}_r^{l_{i+1}, \dots, l_r}, \quad l_r \in \mathcal{T}^{l_{i+1}, \dots, l_{r-1}},$$

603 for any  $l_{i+1} \in \mathcal{S}, \dots, l_{r-1} \in \mathcal{S}^{l_{i+1}, \dots, l_{r-2}}$ , we can write

$$\begin{aligned} (A) = & \sum_{r \leq i} G(i, r, \mathbf{x}_1, \dots, \mathbf{x}_i)^\top (\mathbf{x}_r^* - \mathbf{x}_r) + \sum_{i+1 \leq r \leq p} \sum_{\substack{l_{i+1} \in \mathcal{S}, \dots, \\ l_{r-1} \in \mathcal{S}^{l_{i+1}, \dots, l_{r-2}}}} \lambda^{l_{i+1}} \dots \lambda^{l_{i+1}, \dots, l_{r-1}} \\ & \times \sum_{l_r \in (\mathcal{S} \cup \mathcal{T})^{l_{i+1}, \dots, l_{r-1}}} \lambda^{l_{i+1}, \dots, l_r} G(r, r, \mathbf{x}_1, \dots, \mathbf{x}_i, \mathbf{c}^{l_{i+1}}, \dots, \mathbf{c}^{l_{i+1}, \dots, l_r})^\top (\mathbf{x}_r^* - \mathbf{c}^{l_{i+1}, \dots, l_r}). \end{aligned}$$

604 We next relate the terms  $G$  to the output of ApproxOracle. For simplicity, let us write  $\mathbf{G} =$   
 605  $(\mathcal{P}_{l_1}^{(1)}, \dots, \mathcal{P}_{l_i}^{(i)})$ , which by abuse of notation was assimilated to  $(\mathbf{x}_1, \dots, \mathbf{x}_i)$ . Recall that by con-  
 606 struction and hypothesis, all points where the oracle was queried belong to  $\mathcal{C}_d$ , so that for instance

607  $\|\mathbf{x}_r^* - \mathbf{c}^{l_{i+1}, \dots, l_r}\| \leq 2\sqrt{k_r} \leq 2\sqrt{d}$  for any  $l_r \in \mathcal{S}^{l_{i+1}, \dots, l_{r-1}}$ . Using the above equations together  
 608 with Eq (2) and Lemma A.2 gives

$$\begin{aligned} \epsilon &\leq \sum_{r \leq i} \left[ \text{ApproxOracle}_{\delta, \xi, \mathcal{O}_f}(i, r, \mathbf{G})^\top (\mathbf{x}_r^* - \mathbf{x}_r) + \frac{8d^{3/2}}{\sigma_{\min}} \xi \right] + \sum_{i+1 \leq r \leq p} \sum_{\substack{l_{i+1} \in \mathcal{S}, \dots, \\ l_{r-1} \in \mathcal{S}^{l_{i+1}, \dots, l_{r-2}}} \\ \lambda^{l_{i+1}} \dots \lambda^{l_{i+1}, \dots, l_{r-1}} &\sum_{l_r \in (\mathcal{S} \cup \mathcal{T})^{l_{i+1}, \dots, l_{r-1}}} \lambda^{l_{i+1}, \dots, l_r} \left[ (\mathbf{g}_r^{l_{i+1}, \dots, l_r})^\top (\mathbf{x}_r^* - \mathbf{c}^{l_{i+1}, \dots, l_r}) + \frac{8d^{3/2}}{\sigma_{\min}} \xi \right] \\ &\leq \frac{8pd^{3/2}}{\sigma_{\min}} \xi + (p-i)\delta + \sum_{r \leq i} \text{ApproxOracle}_{\delta, \xi, \mathcal{O}_f}(i, r, \mathbf{G})^\top (\mathbf{x}_r^* - \mathbf{x}_r) \end{aligned}$$

609 where in the second inequality, we used Eq (3). Using  $p \leq d$ , this ends the proof of the lemma.  $\square$

610 We are now ready to show that Algorithm 4 is a valid algorithm for convex optimization.

611 **Theorem A.1.** *Let  $\epsilon \in (0, 1)$  and  $O_S : \mathcal{C}_d \rightarrow \mathbb{R}^d$  be a separation oracle such that the successful set*  
 612 *contains a ball of radius  $\epsilon$ . Pose  $\delta = \frac{\epsilon}{4d}$  and  $\xi = \frac{\sigma_{\min} \epsilon}{32d^{5/2}}$ . Next, let  $p \geq 1$  and  $k_1, \dots, k_p \leq \lceil \frac{d}{p} \rceil$  such*  
 613 *that  $k_1 + \dots + k_p = d$ . With these parameters, Algorithm 4 finds a successful vector with  $(C \frac{d}{p} \ln \frac{d}{\epsilon})^p$*   
 614 *queries and using memory  $\mathcal{O}(\frac{d^2}{p} \ln \frac{d}{\epsilon})$ , for some universal constant  $C > 0$ .*

615 *Proof.* Suppose by contradiction that Algorithm 4 never queried a successful point. Then, with the  
 616 chosen parameters, Lemma A.3 shows that, for any vector  $\mathbf{x}^* = (\mathbf{x}_1^*, \dots, \mathbf{x}_p^*)$  such that  $B_d(\mathbf{x}^*, \epsilon)$  is  
 617 within the successful set, with the same notations, one has

$$\sum_{r \leq i} \text{ApproxOracle}_{\delta, \xi, \mathcal{O}_S}(i, r, \mathcal{P}_{l_1}^{(1)}, \dots, \mathcal{P}_{l_i}^{(i)})^\top (\mathbf{x}_r^* - \mathbf{x}_r) \geq \epsilon - \frac{8d^{5/2}}{\sigma_{\min}} \xi - d\delta \geq \frac{\epsilon}{2}.$$

618 Now denote by  $(\mathbf{a}_t, b_t)$  the constraints that were added at any time during the run of Algorithm 1  
 619 when using the oracle ApproxOracle with  $i = j = 1$ . The previous equation shows that for all such  
 620 constraints,

$$\mathbf{a}_t^\top \mathbf{x}_1^* - b_t \geq \mathbf{a}_t^\top (\mathbf{x}_1^* - \omega_t) - \xi \geq \frac{\epsilon}{2} - \xi,$$

621 where  $\omega_t$  is the volumetric center of the polyhedron at time  $t$  during Vaidya's method Algorithm 1.  
 622 Now, since the algorithm terminated, by Lemma 4.1, we have that

$$\min_t (\mathbf{a}_t^\top \mathbf{x}_1^* - b_t) \leq \delta.$$

623 This is absurd since  $\delta + \xi < \frac{\epsilon}{2}$ . This ends the proof that Algorithm 4 finds a successful vector.

624 We now estimate its oracle-complexity and memory usage. First, recall that a run of ApproxOracle  
 625 of level  $i$  makes  $\mathcal{O}(k_{i+1} \ln \frac{1}{\delta})$  calls to level- $(i+1)$  runs of ApproxOracle. As a result, the oracle-  
 626 complexity  $Q_d(\epsilon; k_1, \dots, k_p)$  satisfies

$$Q_d(\epsilon; k_1, \dots, k_p) = \left( C k_1 \ln \frac{1}{\delta} \right) \times \dots \times \left( C k_p \ln \frac{1}{\delta} \right) \leq \left( C' \frac{d}{p} \log \frac{d}{\epsilon} \right)^p$$

627 for some universal constants  $C, C' \geq 2$ .

628 We now turn to the memory of the algorithm. For each level  $i \in [p]$  of runs for ApproxOracle, we  
 629 keep memory placements for

- 630 1. the value  $j^{(i)}$  of the corresponding call to ApproxOracle( $i, j^{(i)}, \cdot$ ) (for 1.6-7 of Algorithm 3):  
 631  $\mathcal{O}(\ln d)$  bits,
- 632 2. the iteration number  $t^{(i)}$  during the run of Algorithm 1 or within Algorithm 2:  $\mathcal{O}(\ln(k_i \ln \frac{1}{\delta}))$   
 633 bits
- 634 3. the polyhedron constraints contained in the state of  $\mathcal{P}^{(i)}$ :  $\mathcal{O}(k_i \times k_i \ln \frac{1}{\epsilon})$  bits,

Table 1: Memory structure for Algorithm 4

$i$	1	...	$p$
$j$	$j^{(1)}$		$j^{(p)}$
Iteration index	$t^{(1)}$		$t^{(p)}$
Polyhedron	$\mathcal{P}^{(1)} = \begin{pmatrix} k_1, \mathbf{a}_1, b_1 \\ k_2, \mathbf{a}_2, b_2 \\ \dots \\ k_m, \mathbf{a}_m, b_m \end{pmatrix}$		$\mathcal{P}^{(p)}$
Computed dual variables	$(\mathbf{k}^{*(1)}, \boldsymbol{\lambda}^{*(1)}) = \begin{pmatrix} k_1^*, \lambda_1^* \\ k_2^*, \lambda_2^* \\ \dots \end{pmatrix}$		$(\mathbf{k}^{*(p)}, \boldsymbol{\lambda}^{*(p)})$
Working separation vector	$\mathbf{u}^{(1)}$		$\mathbf{u}^{(p)}$

635 4. potentially, already computed dual variables  $\boldsymbol{\lambda}^*$  and their corresponding vector of constraint  
636 indices  $\mathbf{k}^*$  (1.3 of Algorithm 2):  $\mathcal{O}(k_i \times \ln \frac{1}{\xi})$  bits,

637 5. the working vector  $\mathbf{u}^{(i)}$  (updated 1.8 of Algorithm 2):  $\mathcal{O}(k_i \ln \frac{1}{\xi})$  bits.

638 The memory structure is summarized in Table 1.

639 We can then check that this memory is sufficient to run Algorithm 4. An important point is that for  
640 any run of  $\text{ApproxOracle}(i, j, \cdot)$ , in Algorithm 2, after running Vaidya's method Algorithm 1 and  
641 storing the dual variables  $\boldsymbol{\lambda}^*$  and corresponding indices  $\mathbf{k}^*$  within their placements  $(\mathbf{k}^{*(i)}, \boldsymbol{\lambda}^{*(i)})$   
642 (1.1-3 of Algorithm 2), the iteration index  $t^{(i)}$  and polyhedron  $\mathcal{P}^{(i)}$  memory placements are reset  
643 and can be used again for the second run of Vaidya's method (1.4-10 of Algorithm 2). During this  
644 second run, the vector  $\mathbf{u}$  is stored in its corresponding memory placement  $\mathbf{u}^{(i)}$  and updated along  
645 the algorithm. Once this run is finished, the output of  $\text{ApproxOracle}(i, j, \cdot)$  is readily available in  
646 the placement  $\mathbf{u}^{(i)}$ . For  $i = p$ , the algorithm does not need to wait for the output of a level- $(i + 1)$   
647 computation and can directly use the  $j^{(p)}$ -th component of the returned separation vector from the  
648 oracle  $O_S$ . As a result, the number of bits of memory used throughout the algorithm is at most

$$M = \sum_{i=1}^p \mathcal{O} \left( k_i^2 \ln \frac{1}{\xi} \right) = \mathcal{O} \left( \frac{d^2}{p} \ln \frac{d}{\epsilon} \right).$$

649 This ends the proof of the theorem.  $\square$

650 We can already give the useful range for  $p$  for our algorithms, which will also apply to the case with  
651 computational-memory constraints Appendix B.

652 *Proof of Corollary 3.1.* Suppose  $\epsilon \geq \frac{1}{d^d}$ . Then, for some  $p_{max} = \Theta(\frac{C \ln \frac{1}{\epsilon}}{2 \ln d}) \leq d$ , the algorithm  
653 from Theorem 3.2 yields a  $\mathcal{O}(\frac{1}{\epsilon^2})$  oracle-complexity. On the other hand, if  $\epsilon \leq \frac{1}{d^d}$ , we can take  
654  $p_{max} = d$ , which gives an oracle-complexity  $\mathcal{O}((C \ln \frac{1}{\epsilon})^d)$ .  $\square$

## 655 B Memory-constrained feasibility problem with computations

656 In the last section we gave the main ideas that allow reducing the storage memory. However,  
657 Algorithm 4 does not account for memory constraints in computations as per Definition 2.2. For  
658 instance, computing the volumetric center  $\text{VolumetricCenter}(\mathcal{P})$  already requires infinite memory for  
659 infinite precision. More importantly, even if one discretizes the queries, the necessary precision and  
660 computational power may be prohibitive with the classical Vaidya's method Algorithm 1. Even finding  
661 a feasible point in the polyhedron (let alone the volumetric center) using only the constraints is itself  
662 computationally intensive. There has been significant work to make Vaidya's method computationally  
663 tractable [46, 1, 2]. These works address the issue of computational tractability, but the memory issue

664 is still present. Indeed, the precision depends among other parameters on the condition number of the  
665 matrix  $\mathbf{H}$  in order to compute the leverage scores  $\sigma_i$  for  $i \in [m]$ , which may not be well-conditioned.  
666 Second, to avoid memory overflow, we also need to ensure that the points queried have bounded  
667 norm, which is again not a priori guaranteed in the original version Algorithm 1.

668 To solve these issues and also give a computationally-efficient algorithm, the cutting-plane subroutine  
669 Algorithm 1 needs to be modified. In particular, the volumetric barrier needs to include regularization  
670 terms. Fortunately, these have already been studied in [25]. In a major breakthrough, this paper gave  
671 a cutting-plane algorithm with  $\mathcal{O}(d^3 \ln^{(1)} \frac{d}{\epsilon})$  runtime complexity, improving over the seminal work  
672 from Vaidya and subsequent works which had  $\mathcal{O}(d^{1+\omega} \ln^{(1)} \frac{d}{\epsilon})$  runtime complexity, where  $\mathcal{O}(d^\omega)$  is  
673 the computational complexity of matrix multiplication. To achieve this result, they introduce various  
674 regularizing terms together with the logarithmic barrier. While the main motivation of [25] was  
675 computational complexity, as a side effect, these regularization terms also ensure that computations  
676 can be carried with efficient memory. We then use their method as a subroutine.

677 For the sake of exposition and conciseness, we describe a simplified version of their method, that  
678 is also deterministic. This comes at the expense of a suboptimal running time  $\mathcal{O}(d^{1+\omega} \ln^{(1)} \frac{1}{\epsilon})$ .  
679 We recall that our main concern is in memory usage rather than achieving the optimal runtime. The  
680 main technicality of this section is to show that their simplified method is numerically stable, and  
681 we emphasize that the original algorithm could also be shown to be numerically stable with similar  
682 techniques, leading to a time improvement from  $\tilde{\mathcal{O}}(d^{1+\omega})$  to  $\tilde{\mathcal{O}}(d^3)$ . The memory usage, however,  
683 would not be improved.

## 684 B.1 A memory-efficient Vaidya’s method for computations, via [25]

685 Fix a polyhedron  $\mathcal{P} = \{\mathbf{x} : \mathbf{A}\mathbf{x} \geq \mathbf{b}\}$ . Using the same notations as for Vaidya’s method in Section 4.1,  
686 we define the new leverage scores  $\psi(\mathbf{x})_i = (\mathbf{A}_x(\mathbf{A}_x^\top \mathbf{A}_x + \lambda \mathbf{I})^{-1} \mathbf{A}_x^\top)_{i,i}$  and  $\Psi(\mathbf{x}) = \text{diag}(\psi(\mathbf{x}))$ .  
687 Let  $\mu(\mathbf{x}) = \min_i \psi(\mathbf{x})_i$ . Last, let  $\mathbf{Q}(\mathbf{x}) = \mathbf{A}_x^\top (c_e \mathbf{I} + \Psi(\mathbf{x})) \mathbf{A}_x + \lambda \mathbf{I}$ , where  $c_e > 0$  is a constant  
688 parameter to be defined. In [25], they consider minimizing the volumetric-analytic hybrid barrier  
689 function

$$p(\mathbf{x}) = -c_e \sum_{i=1}^m \ln s_i(\mathbf{x}) + \frac{1}{2} \ln \det(\mathbf{A}_x^\top \mathbf{A}_x + \lambda \mathbf{I}) + \frac{\lambda}{2} \|\mathbf{x}\|_2^2.$$

690 We can check [25] that

$$\nabla p(\mathbf{x}) = -\mathbf{A}_x^\top (c_e \cdot \mathbf{1} + \psi(\mathbf{x})) + \lambda \mathbf{x},$$

691 where  $\mathbf{1}$  is the vector of ones. The following procedure gives a way to minimize this function  
692 efficiently given a good starting point.

**Input:** Initial point  $\mathbf{x}^{(0)} \in \mathcal{P} = \{\mathbf{x} : \mathbf{A}\mathbf{x} \geq \mathbf{b}\}$

**Input:** Number of iterations  $r > 0$

**Given :**  $\|\nabla p(\mathbf{x}^{(0)})\|_{\mathbf{Q}(\mathbf{x}^{(0)})^{-1}} \leq \frac{1}{100} \sqrt{c_e + \mu(\mathbf{x}^{(0)})} := \eta$ .

1 **for**  $k = 1$  **to**  $r$  **do**

2     **if**  $\|\nabla p(\mathbf{x}^{(k-1)})\|_{\mathbf{Q}(\mathbf{x}^{(0)})^{-1}} \leq 2(1 - \frac{1}{64})^r \eta$  **then Break;**

3      $\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} - \frac{1}{8} \mathbf{Q}(\mathbf{x}^{(0)})^{-1} \nabla p(\mathbf{x}^{(k-1)})$

4 **end**

**Output:**  $\mathbf{x}^{(k)}$

**Algorithm 5:**  $\mathbf{x}^{(r)} = \text{Centering}(\mathbf{x}^{(0)}, r)$

693 We then present their simplified cutting-plane method.

694 In both Algorithm 5 and Algorithm 6, notice that the updates require to compute in particular the  
695 leverage scores  $\psi(\mathbf{x})$ , which can be computed in  $\mathcal{O}(d^\omega)$  time using their formula. To achieve the  
696  $\mathcal{O}(d^3 \ln^{(1)} \frac{1}{\epsilon})$  computational complexity, an amortized computational cost  $\mathcal{O}(d^2)$  is needed. The  
697 algorithm from [25] achieves this through various careful techniques aiming to update estimates  
698 of these leverage scores. The above cutting-plane algorithm is exactly that of [25] when these  
699 estimates are always exact (i.e. recomputed at each iteration), which yields the  $d^{\omega-2}$  overhead time  
700 complexity. In particular, the original proof of convergence and correctness of [25] directly applies to  
701 this simplified algorithm.

**Input:**  $\epsilon, \delta > 0$  and a separation oracle  $O : \mathcal{C}_d \rightarrow \mathbb{R}^d$

**Check:** Throughout the algorithm, if  $s_i(\mathbf{x}^{(t)}) < 2\epsilon$  for some  $i$  then **return**  $(\mathcal{P}_t, \mathbf{x}^{(t)})$

```

1 Initialize  $\mathbf{x}^{(0)} = \mathbf{0}$  and  $\mathcal{P}_0 := \{(-1, \mathbf{e}_i, -1), (-1, -\mathbf{e}_i, -1), i \in [d]\}$ 
2 for  $t \geq 0$  do
3   if  $\min_{i \in [m]} \psi(\mathbf{x}^{(t)})_i \leq c_d$  then
4      $\mathcal{P}_{t+1} = \mathcal{P}_t \setminus \{(k_j, \mathbf{a}_j, b_j)\}$  where  $j \in \arg \min_{i \in [m]} \psi(\mathbf{x}^{(t)})_i$ 
5   else
6     if  $\mathbf{x}^{(t)} \notin \mathcal{C}_d$  then  $\mathbf{a} = -\text{sign}(x_i)\mathbf{e}_i$  where  $i \in \arg \min_{j \in [d]} |x_j^{(t)}|$ ;
7     else  $\mathbf{a} = O(\mathbf{x}^{(t)})$ ;
8     Let  $b = \mathbf{a}^\top \mathbf{x}^{(t)} - c_a^{-1/2} \sqrt{\mathbf{a}^\top (\mathbf{A}^\top \mathbf{S}_{\mathbf{x}^{(t)}}^{-2} \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{a}}$ 
9      $\mathcal{P}_{t+1} = \mathcal{P}_t \cup \{(t, \mathbf{a}, b)\}$ 
10     $\mathbf{x}^{(t+1)} = \text{Centering}(\mathbf{x}^{(t)}, 200, c_\Delta)$ 
11 end

```

**Algorithm 6:** An efficient cutting-plane method, simplified from [25]

702 It remains to check whether one can implement this algorithm with efficient memory, corresponding  
703 to checking this method's numerical stability.

704 **Lemma B.1.** *Suppose that each iterate of the centering Algorithm 5,  $\|\nabla p(\mathbf{x}^{(k-1)})\|_{\mathbf{Q}(\mathbf{x}^{(0)})^{-1}}$  is com-  
705 puted up to precision  $(1 - \frac{1}{64})^r \eta$  (l.2), and  $\mathbf{x}^{(k)}$  is computed up to an error  $\zeta^{(k)}$  with  $\|\zeta^{(k)}\|_{\mathbf{Q}(\mathbf{x}^{(0)})} \leq$   
706  $\frac{1}{2^{10r}}(1 - \frac{1}{64})^r \eta$  (l.3). Then, Algorithm 5 outputs  $\mathbf{x}^{(k)}$  such that  $\|\nabla p(\mathbf{x}^{(k)})\|_{\mathbf{Q}^{-1}(\mathbf{x}^{(k)})} \leq 3(1 - \frac{1}{64})^r \eta$   
707 and all iterates computed during the procedure satisfy  $\|\mathbf{S}_{\mathbf{x}^{(0)}}^{-1}(\mathbf{s}(\mathbf{x}^{(t)}) - \mathbf{s}(\mathbf{x}^{(0)}))\|_2 \leq \frac{1}{10}$ .*

708 *Proof.* As mentioned above, without computation errors, the result from [25] would apply directly.  
709 Here, we simply adapt the proof to the case with computational errors to show that it still applies.  
710 Denote  $\mathbf{Q} = \mathbf{Q}(\mathbf{x}^{(0)})$  for convenience. Let  $\eta = \frac{1}{100} \sqrt{c_e + \mu(\mathbf{x}^{(0)})}$ . We prove by induction that  
711  $\|\mathbf{x}^{(t)} - \mathbf{x}^{(0)}\|_{\mathbf{Q}} \leq 9\eta$ ,  $\|\nabla p(\mathbf{x}^{(t)})\|_{\mathbf{Q}^{-1}} \leq (1 - \frac{1}{64})^t \eta$  for all  $t \leq r$ . For a given iteration  $t$ , denote  
712  $\tilde{\mathbf{x}}^{(t+1)} = \mathbf{x}^{(k-1)} - \frac{1}{8} \mathbf{Q}^{-1} \nabla p(\mathbf{x}^{(k-1)})$  the result of the exact computation. The same arguments as  
713 in the original proof give  $\|\tilde{\mathbf{x}}^{(t+1)} - \mathbf{x}^{(0)}\|_{\mathbf{Q}} \leq 9\eta$ , and

$$\|\nabla p(\tilde{\mathbf{x}}^{(t+1)})\|_{\mathbf{Q}^{-1}} \leq \left(1 - \frac{1}{32}\right) \|\nabla p(\mathbf{x}^{(t)})\|_{\mathbf{Q}^{-1}}.$$

714 Now because  $\|\tilde{\mathbf{x}}^{(t+1)} - \mathbf{x}^{(t+1)}\|_{\mathbf{Q}} \leq \eta$ , we have  $\|\tilde{\mathbf{x}}^{(t+1)} - \mathbf{x}^{(0)}\|_{\mathbf{Q}}, \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(0)}\|_{\mathbf{Q}} \leq 10\eta$ , so that  
715 [25, Lemma 11] gives  $\nabla^2 p(\mathbf{y}(u)) \preceq 8\mathbf{Q}(\mathbf{y}(u)) \preceq 16\mathbf{Q}$ , where  $\mathbf{y}(u) = \mathbf{x}^{(t+1)} + u(\tilde{\mathbf{x}}^{(t+1)} - \mathbf{x}^{(t+1)})$   
716 for  $u \in [0, 1]$ . Thus,

$$\begin{aligned} \|\nabla p(\tilde{\mathbf{x}}^{(t+1)}) - \nabla p(\mathbf{x}^{(t+1)})\|_{\mathbf{Q}^{-1}} &\leq \left\| \int_0^1 \nabla^2 p(\mathbf{y}(u)) (\tilde{\mathbf{x}}^{(t+1)} - \mathbf{x}^{(t+1)}) \right\|_{\mathbf{Q}^{-1}} \\ &\leq 16 \|\tilde{\mathbf{x}}^{(t+1)} - \mathbf{x}^{(t+1)}\|_{\mathbf{Q}}. \end{aligned}$$

717 Now by construction of the procedure, if the algorithm performed iteration  $t + 1$ , we have  
718  $\|\nabla p(\mathbf{x}^{(t)})\|_{\mathbf{Q}^{-1}} \geq (1 - \frac{1}{64})^t \eta$ . Combining this with the fact that  $\|\tilde{\mathbf{x}}^{(t+1)} - \mathbf{x}^{(t+1)}\|_{\mathbf{Q}} \leq$   
719  $\frac{1}{2^{10r}}(1 - \frac{1}{64})^r \eta$ , obtain

$$\begin{aligned} \|\nabla p(\mathbf{x}^{(t+1)})\|_{\mathbf{Q}^{-1}} &\leq \|\nabla p(\tilde{\mathbf{x}}^{(t+1)}) - \nabla p(\mathbf{x}^{(t+1)})\|_{\mathbf{Q}^{-1}} + \|\nabla p(\tilde{\mathbf{x}}^{(t+1)})\|_{\mathbf{Q}^{-1}} \\ &\leq \left(1 - \frac{1}{64}\right) \|\nabla p(\mathbf{x}^{(t)})\|_{\mathbf{Q}^{-1}}. \end{aligned}$$

720 We now write

$$\begin{aligned} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(0)}\|_{\mathbf{Q}} &\leq \sum_{k=0}^t \|\tilde{\mathbf{x}}^{(k+1)} - \mathbf{x}^{(k+1)}\|_{\mathbf{Q}} + \frac{1}{8} \|\mathbf{Q}^{-1} \nabla p(\mathbf{x}^{(k)})\|_{\mathbf{Q}} \\ &\leq \eta + \frac{1}{8} \sum_{i=0}^{\infty} \left(1 - \frac{1}{64}\right)^i \eta \leq 9\eta. \end{aligned}$$

721 The induction is now complete. When the algorithm stops, either the  $r$  steps were performed, in which  
 722 case the induction already shows that  $\|\nabla p(\mathbf{x}^{(r)})\|_{\mathbf{Q}^{-1}} \leq (1 - \frac{1}{64})^r \eta$ . Otherwise, if the algorithm  
 723 terminates at iteration  $k$ , because  $\|\nabla p(\mathbf{x}^{(k)})\|_{\mathbf{Q}^{-1}}$  was computed to precision  $(1 - \frac{1}{64})^r \eta$ , we have  
 724 (see 1.2 of Algorithm 5)

$$\|\nabla p(\mathbf{x}^{(k)})\|_{\mathbf{Q}^{-1}} \leq 2 \left(1 - \frac{1}{64}\right)^r \eta + \left(1 - \frac{1}{64}\right)^r \eta = 3 \left(1 - \frac{1}{64}\right)^r \eta.$$

725 The same argument as in the original proof shows that at each iteration  $t$ ,

$$\|\mathbf{S}_{\mathbf{x}^{(0)}}^{-1}(\mathbf{s}(\mathbf{x}^{(t)}) - \mathbf{s}(\mathbf{x}^{(0)}))\|_2 = \|\mathbf{x}^{(t)} - \mathbf{x}^{(0)}\|_{\mathbf{A}^\top \mathbf{S}_{\mathbf{x}^{(0)}}^{-2} \mathbf{A}} \leq \frac{\|\mathbf{x}^{(t)} - \mathbf{x}^{(0)}\|_{\mathbf{Q}}}{\sqrt{\mu(\mathbf{x}^{(0)}) + c_e}} \leq \frac{1}{10}.$$

726 This ends the proof of the lemma.  $\square$

727 Because of rounding errors, Lemma B.1 has an extra factor 3 compared to the original guarantee in  
 728 [25, Lemma 14]. To achieve the same guarantee, it suffices to perform  $70 \geq \ln(3)/\ln(1/(1 - \frac{1}{64}))$   
 729 additional centering procedures at most. hence, instead of performing 200 centering procedures  
 730 during the cutting plane method, we perform 270 (1.10 of Algorithm 6). We next turn to the numerical  
 731 stability of the main Algorithm 6.

732 **Lemma B.2.** *Suppose that throughout the algorithm, when checking the stopping criterion*  
 733  *$\min_{i \in [m]} s_i(\mathbf{x}) < 2\epsilon$ , the quantities  $s_i(\mathbf{x})$  were computed with accuracy  $\epsilon$ . Suppose that at each*  
 734 *iteration of Algorithm 6, the leverage scores  $\psi(\mathbf{x}^{(t)})$  are computed up to multiplicative precision*  
 735  *$c_\Delta/4$  (1.3), that when a constraint is added, the response of the oracle  $\mathbf{a}$  (1.7) is stored perfectly but  $b$*   
 736 *(1.8) is computed up to precision  $\Omega(\frac{\epsilon}{\sqrt{n}})$ . Further suppose that the centering Algorithm 5 is run with*  
 737 *numerical approximations according to the assumptions in Lemma B.1. Then, all guarantees for the*  
 738 *original algorithm in [25] hold, up to a factor 3 for  $\epsilon$ .*

739 *Proof.* We start with the termination criterion. Given the requirement on the computational accuracy,  
 740 we know that the final output  $\mathbf{x}$  satisfies  $\min_{i \in [m]} s_i(\mathbf{x}) \leq 3\epsilon$ . Further, during the algorithm, if it  
 741 does not stop, then one has  $\min_{i \in [m]} s_i(\mathbf{x}) \geq \epsilon$ , which is precisely the guarantee of the original  
 742 algorithm in [25].

743 We next turn to the computation of the leverage scores in 1.4. In the original algorithm, only a  
 744  $c_\Delta$ -estimate is computed. Precisely, one computes a vector  $\mathbf{w}^{(t)}$  such that for all  $i \in [d]$ ,  $\psi(\mathbf{x}^{(t)})_i \leq$   
 745  $w_i \leq (1 + c_\Delta)\psi(\mathbf{x}^{(t)})_i$ , then deletes a constraint when  $\min_{i \in [m^{(t)}]} w_i^{(t)} \leq c_d$ . In the adapted  
 746 algorithm, let  $\tilde{\psi}(\mathbf{x}^{(t)})_i$  denote the computed leverage scores for  $i \in [d]$ . By assumption, we have

$$(1 - c_\Delta/4)\psi(\mathbf{x}^{(t)})_i \leq \tilde{\psi}(\mathbf{x}^{(t)})_i \leq (1 + c_\Delta/4)\psi(\mathbf{x}^{(t)})_i.$$

747 Up to re-defining the constant  $c_d$  as  $(1 - c_\Delta/4)c_d$ ,  $\tilde{\psi}(\mathbf{x}^{(t)})$  is precisely within the guarantee bounds  
 748 of the algorithm. For the accuracy on the separation oracle response and the second-term value  $b$ , [25]  
 749 emphasizes that the algorithm always changes constraints by a  $\delta$  amount where  $\delta = \Omega(\frac{\epsilon}{\sqrt{d}})$  so that  
 750 an inexact separation oracle with accuracy  $\Omega(\frac{\epsilon}{\sqrt{d}})$  suffices. Therefore, storing an  $\Omega(\frac{\epsilon}{\sqrt{d}})$  accuracy  
 751 of the second term keeps the guarantees of the algorithm. Last, we checked in Lemma B.1 that the  
 752 centering procedure Algorithm 5 satisfies all the requirements needed in the original proof [25].  $\square$

753 For our recursive method, we need an efficient cutting-plane method that also provides a proof  
 754 (certificate) of convergence. This is also provided by [25] that provide a proof that the feasible region  
 755 has small width in one of the directions  $\mathbf{a}_i$  of the returned polyhedron.

756 **Lemma B.3.** [25, Lemma 28] *Let  $(\mathcal{P}, \mathbf{x}, (\lambda_i)_i)$  be the output of Algorithm 7. Then,  $\mathbf{x}$  is feasible,*  
 757  *$\|\mathbf{x}\|_2 \leq 3\sqrt{d}$ ,  $\lambda_j \geq 0$  for all  $j$  and  $\sum_i \lambda_i = 1$ . Further,*

$$\left\| \sum_i \lambda_i \mathbf{a}_i \right\|_2 = \mathcal{O}\left(\epsilon \sqrt{d} \ln \frac{d}{\epsilon}\right), \quad \text{and} \quad \sum_i \lambda_i (\mathbf{a}_i^\top \mathbf{x} - b_j) \leq \mathcal{O}\left(d\epsilon \ln \frac{d}{\epsilon}\right).$$

758 We are now ready to show that Algorithm 6 can be implemented with efficient memory and also  
 759 provides a proof of the convergence of the algorithm.

**Input:**  $\epsilon > 0$  and a separation oracle  $O : \mathcal{C}_d \rightarrow \mathbb{R}^d$

1 Run Algorithm 6 to obtain a polyhedron  $\mathcal{P}$  and a feasible point  $\mathbf{x}$

2  $\mathbf{x}^* = \text{Centering}(\mathbf{x}, 64 \ln \frac{2}{\epsilon}, c_\Delta)$

3  $\lambda_i = \frac{c_\epsilon + \psi_i(\mathbf{x}^*)}{s_i(\mathbf{x}^*)} \left( \sum_j \frac{c_\epsilon + \psi_j(\mathbf{x}^*)}{s_j(\mathbf{x}^*)} \right)^{-1}$  for all  $i$

**Output:**  $(\mathcal{P}, \mathbf{x}^*, (\lambda_i)_i)$

**Algorithm 7:** Cutting-plane algorithm with certified optimality

760 **Proposition B.1.** *Provided that the output of the oracle are vectors discretized to precision  $\text{poly}(\frac{\epsilon}{d})$*   
 761 *and have norm at most 1, Algorithm 7 can be implemented with  $\mathcal{O}(d^2 \ln \frac{d}{\epsilon})$  bits of memory to output*  
 762 *a certified optimal point according to Lemma B.3. The algorithm performs  $\mathcal{O}(d \ln \frac{d}{\epsilon})$  calls to the*  
 763 *separation oracle and runs in  $\mathcal{O}(d^{1+\omega} \ln^{\mathcal{O}(1)} \frac{d}{\epsilon})$  time.*

764 *Proof.* We already checked the numerical stability of Algorithm 6 in Lemma B.2. It remains to check  
 765 the next steps of the algorithm. The centering procedure is stable again via Lemma B.1. It also  
 766 suffices to compute the coefficients  $\lambda_j$  up to accuracy  $\mathcal{O}(\epsilon/(\sqrt{d} \ln(d/\epsilon)))$  to keep the guarantees  
 767 desired since by construction all vectors  $\mathbf{a}_i$  have norm at most one.

768 It now remains to show that the algorithm can be implemented with efficient memory. We recall  
 769 that at any point during the algorithm, the polyhedron  $\mathcal{P}$  has at most  $\mathcal{O}(d)$  constraints [25, Lemma  
 770 22]. Hence, since we assumed that each vector  $\mathbf{a}_i$  composing a constraint is discretized to precision  
 771  $\text{poly}(\frac{\epsilon}{d})$ , we can store the polyhedron constraints with  $\mathcal{O}(d^2 \ln \frac{d}{\epsilon})$  bits of memory. The second  
 772 terms  $b$  are computed up to precision  $\Omega(\epsilon/\sqrt{d})$  hence only use  $\mathcal{O}(d \ln \frac{d}{\epsilon})$  bits of memory. The  
 773 algorithm also keeps the current iterate  $x^{(t)}$  in memory. These are all bounded throughout the  
 774 memory  $\|x^{(t)}\|_2 = \mathcal{O}(\sqrt{d})$  [25, Lemma 23], hence only require  $\mathcal{O}(d \ln \frac{d}{\epsilon})$  bits of memory for the  
 775 desired accuracy.

776 Next, the distances to the constraints are bounded at any step of the algorithm:  $s_i(x^{(t)}) \leq \mathcal{O}(\sqrt{d})$   
 777 [25, Lemma 24], hence computing  $s_i(x^{(t)})$  to the required accuracy is memory-efficient. Recall  
 778 that from the termination criterion, except for the last point, any point  $\mathbf{x}$  during the algorithm  
 779 satisfies  $s_i(\mathbf{x}) \geq \epsilon$  for all constraints  $i \in [m]$ . In particular, this bounds the eigenvalues of  $\mathbf{Q}$   
 780 since  $\lambda \mathbf{I} \preceq \mathbf{Q}(\mathbf{x}) \preceq (\lambda + m(c_\epsilon + 1)/\epsilon^2)\mathbf{I}$ . Thus, the matrix is sufficiently well-conditioned to  
 781 achieve the accuracy guarantees from Lemma B.1 using  $\mathcal{O}(d^2 \ln \frac{d}{\epsilon})$  memory during matrix inversions  
 782 (and matrix multiplications). Similarly, for the computation of leverage scores, we use  $\Psi(\mathbf{x}) =$   
 783  $\text{diag}(\mathbf{A}_x(\mathbf{A}_x^\top \mathbf{A}_x + \lambda \mathbf{I})^{-1} \mathbf{A}_x^\top)$ , where  $\lambda \mathbf{I} \preceq \mathbf{A}_x^\top \mathbf{A}_x + \lambda \mathbf{I} \preceq (\lambda + m\epsilon^{-2})\mathbf{I}$ . This same matrix  
 784 inversion appears when computing the second term of an added constraint. Overall, all linear algebra  
 785 operations are well conditioned and implementable with required accuracy with  $\mathcal{O}(d^2 \ln \frac{d}{\epsilon})$  memory.  
 786 Using fast matrix multiplication, all these operations can be performed in  $\tilde{\mathcal{O}}(d^\omega)$  time per iteration of  
 787 the cutting-plane algorithm since these methods are also known to be numerically stable [13]. Thus,  
 788 the total time complexity is  $\mathcal{O}(d^{1+\omega} \ln^{\mathcal{O}(1)} \frac{d}{\epsilon})$ . The oracle-complexity still has optimal  $\mathcal{O}(d \ln \frac{d}{\epsilon})$   
 789 oracle-complexity as in the original algorithm.  $\square$

790 Up to changing  $\epsilon$  to  $c \cdot \epsilon/(d \ln \frac{d}{\epsilon})$ , the described algorithm finds constraints given by  $\mathbf{a}_i$  and  $b_i$ ,  
 791  $i \in [m]$  returned by the normalized separation oracle, coefficients  $\lambda_i$ ,  $i \in [m]$ , and a feasible point  
 792  $\mathbf{x}^*$  such that for any vector in the unit cube,  $\mathbf{z} \in \mathcal{C}_d$ , one has

$$\min_{i \in [m]} \mathbf{a}_i^\top \mathbf{z} - b_i \leq \sum_{i \in [m]} \lambda_i (\mathbf{a}_i^\top \mathbf{z} - b_i) \leq \left( \sum_{i \in [m]} \lambda \mathbf{a}_i \right)^\top (\mathbf{x}^* - \mathbf{z}) + \sum_{i \in [m]} \lambda_i (\mathbf{a}_i^\top \mathbf{x}^* - b_i) \leq \epsilon.$$

793 This effectively replaces Lemma 4.1.

## 794 B.2 Merging Algorithm 7 within the recursive algorithm

795 Algorithms 2 to 4 from the recursive procedure need to be slightly adapted to the new format of  
 796 the cutting-plane method's output. In particular, the oracles do not take as input polyhedrons (and

797 eventually query their volumetric center as before), but directly take as input an point (which is an  
798 approximate volumetric center).

**Input:**  $\delta, \xi, O_x : \mathcal{C}_n \rightarrow \mathbb{R}^m$  and  $O_y : \mathcal{C}_n \rightarrow \mathbb{R}^n$   
1 Run Algorithm 7 with parameter  $c \cdot \delta / (d \ln \frac{d}{\delta})$ ,  $\xi$  and  $O_y$  to obtain  $(\mathcal{P}^*, \mathbf{x}^*, \boldsymbol{\lambda})$   
2 Store  $\mathbf{k}^* = (k_i, i \in [m])$  where  $m = |\mathcal{P}^*|$ , and  $\boldsymbol{\lambda}^* \leftarrow \text{Discretize}(\boldsymbol{\lambda}^*, \xi)$   
3 Initialize  $\mathcal{P}_0 := \{(-1, \mathbf{e}_i, -1), (-1 - \mathbf{e}_i, -1), i \in [d]\}$ ,  $\mathbf{x}^{(0)} = \mathbf{0}$  and let  $\mathbf{u} = \mathbf{0} \in \mathbb{R}^m$   
4 **for**  $t = 0, 1, \dots, \max_i k_i^*$  **do**  
5     **if**  $t = k_i^*$  for some  $i \in [m]$  **then**  
6          $\mathbf{g}_x = O_x(\mathbf{x}^{(t)})$   
7          $\mathbf{u} \leftarrow \text{Discretize}_m(\mathbf{u} + \lambda_i^* \mathbf{g}_x, \xi)$   
8     Update  $\mathcal{P}_t$  to get  $\mathcal{P}_{t+1}$ , and  $\mathbf{x}^{(t)}$  to get  $\mathbf{x}^{(t+1)}$  as in Algorithm 6  
9 **end**  
10 **return**  $\mathbf{u}$

**Algorithm 8:**  $\text{ApproxSeparationVector}_{\delta, \xi}(O_x, O_y)$

**Input:**  $\delta, \xi, 1 \leq j \leq i \leq p, \mathbf{x}^{(r)} \in \mathcal{C}_{k_r}$  for  $r \in [i], O_S : \mathcal{C}_d \rightarrow \mathbb{R}^d$   
1 **if**  $i = p$  **then**  
2      $(\mathbf{g}_1, \dots, \mathbf{g}_p) = O_S(\mathbf{x}_1, \dots, \mathbf{x}_p)$   
3     **return**  $\text{Discretize}_{k_j}(\mathbf{g}_j, \xi)$   
4 **end**  
5 Define  $O_x : \mathcal{C}_{k_{i+1}} \rightarrow \mathbb{R}^{k_j}$  as  $\text{ApproxOracle}_{\delta, \xi, O_f}(i+1, j, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(i)}, \cdot)$   
6 Define  $O_y : \mathcal{C}_{k_{i+1}} \rightarrow \mathbb{R}^{k_{i+1}}$  as  $\text{ApproxOracle}_{\delta, \xi, O_f}(i+1, i+1, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(i)}, \cdot)$   
7 **return**  $\text{ApproxSeparationVector}_{\delta, \xi}(O_x, O_y)$

**Algorithm 9:**  $\text{ApproxOracle}_{\delta, \xi, O_S}(i, j, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(i)})$

**Input:**  $\delta, \xi$ , and  $O_S : \mathcal{C}_d \rightarrow \mathbb{R}^d$  a separation oracle

**Check:** Throughout the algorithm, if  $O_S$  returned Success to a query  $\mathbf{x}$ , **return**  $\mathbf{x}$

1 Run Algorithm 6 with parameters  $\delta$  and  $\xi$  and oracle  $\text{ApproxOracle}_{\delta, \xi, O_S}(1, 1, \cdot)$

**Algorithm 10:** Memory-constrained algorithm for convex optimization

799 The same proof as for Algorithm 4 shows that Algorithm 10 run with the parameters in Theorem A.1  
800 also outputs a successful vector using the same oracle-complexity. We only need to analyze the  
801 memory usage in more detail.

802 *Proof of Theorem 3.2.* As mentioned above, we will check that Algorithm 10 with the same pa-  
803 rameters  $\delta = \frac{\epsilon}{4d}$  and  $\xi = \frac{\sigma_{\min} \epsilon}{32d^{5/2}}$  as in Theorem A.1 satisfies the desired requirements. We have  
804 already checked its correctness and oracle-complexity. Using the same arguments, the computational  
805 complexity is of the form  $\mathcal{O}(\mathcal{O}(\text{ComplexityCuttingPlanes})^p)$  where  $\text{ComplexityCuttingPlanes}$  is the  
806 computational complexity of the cutting-plane method used, i.e., here of Algorithm 7. Hence, the  
807 computational complexity is  $\mathcal{O}((C(d/p)^{1+\omega} \ln^{O(1)} \frac{d}{\epsilon})^p)$  for some universal constant  $C \geq 2$ . We  
808 now turn to the memory. In addition to the memory of Algorithm 4, described in Table 1, we need

- 809 1. a placement for all  $i \in [p]$  for the current iterate  $\mathbf{x}^{(i)}$ :  $\mathcal{O}(k_i \ln \frac{1}{\xi})$  bits,
- 810 2. a placement for computations, that is shared for all layers (used to compute leverage scores,  
811 centering procedures, etc. By Proposition B.1, since the vectors are always discretized to  
812 precision  $\xi$ , this requires  $\mathcal{O}(\max_{i \in [p]} k_i^2 \ln \frac{d}{\epsilon})$  bits,
- 813 3. the placement  $Q$  to perform queries is the concatenation of the placements  $(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(p)})$ :  
814 no additional bits needed.
- 815 4. a placement  $N$  to store the precision needed for the oracle responses:  $\mathcal{O}(\ln \frac{1}{\xi})$  bits



816 5. a placement  $R$  to receive the oracle responses:  $\mathcal{O}(d \ln \frac{1}{\epsilon})$  bits.

817 The new memory structure is summarized in Table 2.

818 With the same arguments as in the original proof of Theorem A.1, this memory is sufficient to run the  
 819 algorithm and perform computations, thanks to the computation placement. The total number of bits  
 820 used throughout the algorithm remains the same,  $\mathcal{O}(\frac{d^2}{p} \ln \frac{d}{\epsilon})$ . This ends the proof of the theorem.  $\square$

Table 2: Memory structure for Algorithm 10

$i$	1	...	$p$	Oracle response	Precision
$j$	$j^{(1)}$		$j^{(p)}$	$R = (R_1, \dots, R_p)$	$N$
Iteration index	$t^{(1)}$		$t^{(p)}$	Computation memory	
Polyhedron	$\mathcal{P}^{(1)} = \begin{pmatrix} k_1, \mathbf{a}_1, b_1 \\ k_2, \mathbf{a}_2, b_2 \\ \dots \\ k_m, \mathbf{a}_m, b_m \end{pmatrix}$		$\mathcal{P}^{(p)}$		
Current iterate	$\mathbf{x}^{(1)}$		$\mathbf{x}^{(p)}$		
Computed dual variables	$(\mathbf{k}^*, \boldsymbol{\lambda}^*) = \begin{pmatrix} k_1^*, \lambda_1^* \\ k_2^*, \lambda_2^* \\ \dots \end{pmatrix}$		$(\mathbf{k}^{*(p)}, \boldsymbol{\lambda}^{*(p)})$		
Working separation vector	$\mathbf{u}^{(1)}$		$\mathbf{u}^{(p)}$		

## 821 C Improved oracle-complexity/memory lower-bound trade-offs

822 We recall the three oracle-complexity/memory lower-bound trade-offs known in the literature.

- 823 1. First, [29] showed that any (including randomized) algorithm for convex optimization uses  
 824  $d^{1.25-\delta}$  memory or makes  $\tilde{\Omega}(d^{1+4\delta/3})$  queries.
- 825 2. Then, [5] showed that any deterministic algorithm for convex optimization uses  $d^{2-\delta}$   
 826 memory or makes  $\tilde{\Omega}(d^{1+\delta/3})$  queries.
- 827 3. Last, [5] show that any deterministic algorithm for the feasibility problem uses  $d^{2-\delta}$  memory  
 828 or makes  $\tilde{\Omega}(d^{1+\delta})$  queries.

829 Although these papers mainly focused on the regime  $\epsilon = 1/\text{poly}(d)$  and as a result  $\ln \frac{1}{\epsilon} = \mathcal{O}(\ln d)$ ,  
 830 neither of these lower bounds have an explicit dependence in  $\epsilon$ . This can lead to sub-optimal lower  
 831 bounds whenever  $\ln \frac{1}{\epsilon} \gg \ln d$ . Furthermore, in the exponential regime  $\epsilon \leq \frac{1}{2^{\mathcal{O}(d)}}$ , these results  
 832 do not effectively give useful lower bounds. Indeed, in this regime, one has  $d^2 = \mathcal{O}(d \ln \frac{1}{\epsilon})$  and  
 833 as a result, the lower bounds provided are weaker than the classical  $\Omega(d \ln \frac{1}{\epsilon})$  lower bounds for  
 834 oracle-complexity [32] and memory [49]. In particular, in this exponential regime, these results fail  
 835 to show that there is any trade-off between oracle-complexity and memory.

836 In this section, we aim to explicit the dependence in  $\epsilon$  of these lower-bounds. We show with simple  
 837 modifications and additional arguments that one can roughly multiply these oracle-complexity and  
 838 memory lower bounds by a factor  $\ln \frac{1}{\epsilon}$  each. We split the proofs in two. First we give arguments to  
 839 improve the memory dependence by a factor  $\ln \frac{1}{\epsilon}$ , which is achieved by modifying the sampling of  
 840 the rows of the matrix  $\mathbf{A}$  defining a wall term common to the functions considered in the lower bound  
 841 proofs [29, 5]. Then we show how to improve the oracle-complexity dependence by an additional  
 842  $\ln \frac{1}{\epsilon} / \ln d$  factor, via a standard rescaling argument.

### 843 C.1 Improving the memory lower bound

844 We start with some concentration results on random vectors. [29] gave the following result for random  
 845 vectors in the hypercube.

846 **Lemma C.1** ([29]). Let  $\mathbf{h} \sim \mathcal{U}(\{\pm 1\}^d)$ . Then, for any  $t \in (0, 1/2]$  and any matrix  $\mathbf{Z} =$   
 847  $[\mathbf{z}_1, \dots, \mathbf{z}_k] \in \mathbb{R}^{d \times k}$  with orthonormal columns,

$$\mathbb{P}(\|\mathbf{Z}^\top \mathbf{h}\|_\infty \leq t) \leq 2^{-c_H k}.$$

848 Instead, we will need a similar concentration result for random unit vectors in the unit sphere.

849 **Lemma C.2.** Let  $k \leq d$  and  $\mathbf{x}_1, \dots, \mathbf{x}_k$  be  $k$  orthonormal vectors, and  $\zeta \leq 1$ .

$$\mathbb{P}_{\mathbf{y} \sim \mathcal{U}(S^{d-1})} \left( |\mathbf{x}_i^\top \mathbf{y}| \leq \frac{\zeta}{\sqrt{d}}, i \in [k] \right) \leq \left( \frac{2}{\sqrt{\pi}} \zeta \right)^k \leq (\sqrt{2} \zeta)^k.$$

850 *Proof.* First, by isometry, we can suppose that the orthonormal vectors are simply  $\mathbf{e}_1, \dots, \mathbf{e}_k$ . We  
 851 now prove the result by induction on  $d$ . For  $d = 1$ , the result holds directly. Fix  $d \geq 2$ , and  $1 \leq k < d$ .  
 852 Then, if  $S_n$  is the surface area of  $S^n$  the  $n$ -dimensional sphere, then

$$\mathbb{P} \left( |y_1| \leq \frac{\zeta}{\sqrt{d}} \right) \leq \frac{S_{d-2}}{S_{d-1}} \frac{2\zeta}{\sqrt{d}} = \frac{2\zeta}{\sqrt{\pi d}} \frac{\Gamma(d/2)}{\Gamma(d/2 - 1/2)} \leq \frac{2}{\sqrt{\pi}} \zeta. \quad (7)$$

853 Conditionally on the value of  $y_1$ , the vector  $(y_2, \dots, y_d)$  follows a uniform distribution on the  
 854  $(d-2)$ -sphere of radius  $\sqrt{1 - y_1^2}$ . Then,

$$\mathbb{P} \left( |y_i| \leq \frac{\zeta}{\sqrt{d}}, 2 \leq i \leq k \mid y_1 \right) = \mathbb{P}_{\mathbf{z} \sim \mathcal{U}(S^{d-2})} \left( |z_i| \leq \frac{\zeta}{\sqrt{d(1 - y_1^2)}}, 2 \leq i \leq k \right)$$

855 Now recall that since  $|x_1| \leq 1/\sqrt{d}$ , we have  $d(1 - x_1^2) \geq d - 1$ . Therefore, using the induction,

$$\mathbb{P} \left( |y_i| \leq \frac{\zeta}{\sqrt{d}}, 2 \leq i \leq k \mid y_1 \right) \leq \mathbb{P}_{\mathbf{z} \sim \mathcal{U}(S^{d-2})} \left( |z_i| \leq \frac{\zeta}{\sqrt{d-1}}, 2 \leq i \leq k \right) \leq \left( \frac{2\zeta}{\sqrt{\pi}} \right)^{k-1}.$$

856 Combining this equation with Eq (7) ends the proof.  $\square$

857 We next use the following lemma to partition the unit sphere  $S^{d-1}$ .

858 **Lemma C.3** ([17] Lemma 21). For any  $0 < \delta < \pi/2$ , the sphere  $S^{d-1}$  can be partitioned into  
 859  $N(\delta) = (\mathcal{O}(1)/\delta)^d$  equal volume cells, each of diameter at most  $\delta$ .

860 Following the notation from [5], we denote by  $\mathcal{V}_\delta = \{V_i(\delta), i \in [N(\delta)]\}$  the corresponding partition,  
 861 and consider a set of representatives  $\mathcal{D}_\delta = \{\mathbf{b}_i(\delta), i \in [N(\delta)]\} \subset S^{d-1}$  such that for all  $i \in [N(\delta)]$ ,  
 862  $\mathbf{b}_i(\delta) \in V_i(\delta)$ . With these notations we can define the discretization function  $\phi_\delta$  as follows

$$\phi_\delta(\mathbf{x}) = \mathbf{b}_i(\delta), \quad \mathbf{x} \in V_i(\delta).$$

863 We then denote by  $\mathcal{U}_\delta$  the distribution of  $\phi_\delta(\mathbf{z})$  where  $\mathbf{z} \sim \mathcal{U}(S^{d-1})$  is sampled uniformly on the  
 864 sphere. Note that because the cells of  $\mathcal{V}_\delta$  have equal volume,  $\mathcal{U}_\delta$  is simply the uniform distribution on  
 865 the discretization  $\mathcal{D}_\delta$ .

866 We are now ready to give the modifications necessary to the proofs, to include a factor  $\ln \frac{1}{\epsilon}$  for  
 867 the necessary memory. For their lower bounds, [29] exhibit a distribution of convex functions  
 868 that are hard to optimize. Building upon their work [5] construct classes of convex functions that  
 869 are hard to optimize, but that also depend adaptively on the considered optimization algorithm.  
 870 For both, the functions considered a barrier term of the form  $\|\mathbf{A}\mathbf{x}\|_\infty$ , where  $\mathbf{A}$  is a matrix of  
 871  $\approx d/2$  rows that are independently drawn as uniform on the hypercube  $\mathcal{U}(\{\pm 1\}^d)$ . The argument  
 872 shows that memorizing  $\mathbf{A}$  is necessary to a certain extent. As a result, the lower bounds can only  
 873 apply for a memory of at most  $\mathcal{O}(d^2)$  bits, which is sufficient to memorize such a binary matrix.  
 874 Instead, we draw rows independently according to the distribution  $\mathcal{U}_\delta$ , where  $\delta \approx \epsilon$ . We explicit the  
 875 corresponding adaptations for each known trade-off. We start with the lower bounds from [5] for ease of  
 876 exposition; although these build upon those of [29], their parametrization makes the adaptation  
 877 more straightforward.

878 **C.1.1 Lower bound of [5] for convex optimization and deterministic algorithms**

879 For this lower bound, we use the exact same form of functions as they introduced,

$$\max \left\{ \|\mathbf{A}\mathbf{x}\|_\infty - \eta, \eta \mathbf{v}_0^\top \mathbf{x}, \eta \left( \max_{p \leq p_{\max}, l \leq l_p} \mathbf{v}_{p,l}^\top \mathbf{x} - p\gamma_1 - l\gamma_2 \right) \right\},$$

880 with the difference that rows of  $\mathbf{A}$  are take i.i.d. distributed according to  $\mathcal{U}_{\delta'}$  instead of  $\mathcal{U}(\{\pm 1\}^d)$ . As  
 881 a remark, they use  $n = \lceil d/4 \rceil$  rows for  $\mathbf{A}$ . Except for  $\eta$ , we keep all parameters  $\gamma_1, \gamma_2$ , etc as in the  
 882 original proof, and we will take  $\delta' = \epsilon$  and  $\eta = 2\sqrt{d}\epsilon$ . The reason why we introduced  $\delta'$  instead of  $\delta$   
 883 is that the original construction also needs the discretization  $\phi_\delta$ . This is used during the optimization  
 884 procedure which constructs adaptively this class of functions, and only needs  $\delta = \text{poly}(1/d)$  instead  
 885 of  $\delta$  of order  $\epsilon$ .

886 **Theorem C.1.** *For  $\epsilon \leq 1/(2d^{4.5})$  and any  $\delta \in [0, 1]$ , a deterministic first-order algorithm guaranteed  
 887 to minimize 1-Lipschitz convex functions over the unit ball with  $\epsilon$  accuracy uses at least  $d^{2-\delta} \ln \frac{1}{\epsilon}$   
 888 bits of memory or makes  $\tilde{\Omega}(d^{1+\delta/3})$  queries.*

889 With the changes defined above, we can easily check that all results from [5] which reduce convex  
 890 optimization to the optimization procedure, then the optimization procedure to their Orthogonal Vector  
 891 Game with Hints (OVGH) [5, Game 2], are not affected by our changes. The only modifications to  
 892 perform are to the proof of query lower bound for the OVGH [5, Proposition 14]. We emphasize that  
 893 the distribution of  $\mathbf{A}$  is changed in the optimization procedure but also in OVGH as a result.

894 **Proposition C.2.** *Let  $k \geq 20 \frac{M+3d \log(2d)+1}{n \log_2(\sqrt{2}(\zeta+\delta'\sqrt{d}))^{-1}}$ . And let  $0 < \alpha, \beta \leq 1$  such that  $\alpha(\sqrt{d}/\beta)^{5/4} \leq$   
 895  $\zeta/\sqrt{d}$  where  $\zeta \leq 1$ . If the Player wins the adapted OVGH with probability at least  $1/2$ , then  
 896  $m \geq \frac{1}{8} \left(1 + \frac{30 \log_2 d}{\log_2(\sqrt{2}(\zeta+\delta'\sqrt{d}))^{-1}}\right)^{-1} d$ .*

897 *Proof.* We use the same proof and only highlight the modifications. The proof is unchanged until the  
 898 step when the concentration result Lemma C.1 is used. Instead, we use Lemma C.2. With the same  
 899 notations as in the original proof, we constructed  $\lceil k/5 \rceil$  orthonormal vectors  $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_{\lceil k/5 \rceil}]$   
 900 such that all rows  $\mathbf{a}$  of  $\mathbf{A}'$  (which is  $\mathbf{A}$  up to some observed and unimportant rows) one has

$$\|\mathbf{Z}^\top \mathbf{a}\|_\infty \leq \frac{\zeta}{\sqrt{d}}.$$

901 Next, by Lemma C.2, we have

$$\begin{aligned} \left| \left\{ \mathbf{a} \in \mathcal{D}_{\delta'} : \|\mathbf{Z}^\top \mathbf{a}\|_\infty \leq \frac{\zeta}{\sqrt{d}} \right\} \right| &\leq |\mathcal{D}_{\delta'}| \cdot \mathbb{P}_{\mathbf{a} \sim \mathcal{U}_{\delta'}} \left( \|\mathbf{Z}^\top \mathbf{a}\|_\infty \leq \frac{\zeta}{\sqrt{d}} \right) \\ &\leq |\mathcal{D}_{\delta'}| \cdot \mathbb{P}_{\mathbf{z} \sim \mathcal{U}(S^{d-1})} \left( \|\mathbf{Z}^\top \mathbf{z}\|_\infty \leq \frac{\zeta}{\sqrt{d}} + \delta' \right) \\ &\leq |\mathcal{D}_{\delta'}| \cdot \left( \sqrt{2}(\zeta + \delta'\sqrt{d}) \right)^{\lceil k/5 \rceil}. \end{aligned}$$

902 Hence, using the same arguments as in the original proof, we obtain

$$H(\mathbf{A}' | \mathbf{Y}) \leq (n - m) \left( \log_2 |\mathcal{D}_{\delta'}| + \mathbb{P}(\mathcal{E}) \cdot \frac{k}{5} \log_2 \left( \sqrt{2}(\zeta + \delta'\sqrt{d}) \right) \right),$$

903 where  $\mathcal{E}$  is the event when the algorithm succeeds at the OVGH game. In the next step, we need  
 904 to bound  $H(\mathbf{A} | \mathbf{V}) - H(\mathbf{G}, \mathbf{j}, \mathbf{c})$  where  $\mathbf{V}$  stores hints received throughout the game,  $\mathbf{G}$  stores  
 905 observed rows of  $\mathbf{A}$  during the game, and  $\mathbf{j}, \mathbf{c}$  are auxiliary variables. The latter can be treated as in  
 906 the original proof. We obtain

$$\begin{aligned} H(\mathbf{A} | \mathbf{V}) - H(\mathbf{G}, \mathbf{j}, \mathbf{c}) &\geq H(\mathbf{A}) - H(\mathbf{G}) - I(\mathbf{A}; \mathbf{V}) - 3m \log_2(2d) \\ &\geq (n - m) \log_2 |\mathcal{D}_{\delta'}| - 3m \log_2(2d) - I(\mathbf{A}, \mathbf{V}). \end{aligned}$$

907 Now the same arguments as in the original proof show that we still have  $I(\mathbf{A}, \mathbf{V}) \leq 3km \log_2 d + 1$ ,  
 908 and that as a result, if  $M$  is the number of bits stored in memory,

$$M \geq \frac{k}{10} \log_2 \left( \frac{1}{\sqrt{2}(\zeta + \delta'\sqrt{d})} \right) (n - m) - 3km \log_2 d - 1 - 3d \log_2(2d).$$

909 Then, with the same arguments as in the original proof, we can conclude.  $\square$

910 We are now ready to prove Theorem C.1. With the parameter  $k = \lceil 20 \frac{M+3d \log(2d)+1}{n \log_2(\sqrt{2}(\epsilon d^4/2+\delta'\sqrt{d}))^{-1}} \rceil$  and  
 911 the same arguments, we show that an algorithm solving the convex optimization up to precision  
 912  $\eta/(2\sqrt{d}) = \epsilon$  yields an algorithm solving the OVGH where the parameters  $\alpha = \frac{2\eta}{\gamma_1}$  and  $\beta = \frac{\gamma_2}{4}$   
 913 satisfy

$$\alpha \left( \frac{\sqrt{d}}{\beta} \right)^{5/4} \leq \frac{\eta d^3}{4} = \frac{d^{3.5} \epsilon}{2}.$$

914 We can then apply Proposition C.2 with  $\zeta = d^4 \epsilon/2$ . Hence, if  $Q$  is the maximum number of queries  
 915 of the convex optimization algorithm, we obtain

$$\lceil Q/p_{max} \rceil + 1 \geq \frac{1}{8} \left( 1 + \frac{30 \log_2 d}{\log_2 \frac{1}{d^4 \epsilon} - 1/2} \right)^{-1} d \geq \frac{d}{8 \cdot 61},$$

916 where in the last inequality we used  $\epsilon \leq 1/(2d^{4.5})$ . As a result, with the same arguments, we obtain

$$Q = \Omega \left( \frac{d^{5/3} \ln^{1/3} \frac{1}{\epsilon}}{(M + \ln d)^{1/3} \ln^{2/3} d} \right).$$

917 This ends the proof of Theorem C.1.

### 918 C.1.2 Lower bound of [5] for feasibility problems and deterministic algorithms

919 We improve the memory dependence by showing the following result.

920 **Theorem C.3.** *For  $\epsilon = 1/(48d^3)$  and any  $\delta \in [0, 1]$ , a deterministic algorithm guaranteed to solve  
 921 the feasibility problem over the unit ball with  $\epsilon$  accuracy uses at least  $d^{2-\delta} \ln \frac{1}{\epsilon}$  bits of memory or  
 922 makes at least  $\tilde{\Omega}(d^{1+\delta})$  queries.*

923 We use the exact same class of feasibility problems and only change the parameter  $\eta_0$  which  
 924 constrained successful points to satisfy  $\|\mathbf{A}\mathbf{x}\|_\infty \leq \eta_0$ , as well as the rows of  $\mathbf{A}$  that are sampled i.i.d.  
 925 from  $\mathcal{U}_\delta$ . The other parameter  $\eta_1 = 1/(2\sqrt{d})$  is unchanged. We also take  $\delta' = \epsilon$ . Because the rows of  
 926  $\mathbf{A}$  are already normalized, we can take  $\eta_0 = \epsilon$  directly. Then, the same proof as in [5] shows that if an  
 927 algorithm solves feasibility problems with accuracy  $\epsilon$ , there is an algorithm for OVGH for parameters  
 928  $\alpha = \eta/\eta_1$  and  $\beta = \eta_1/2$ . Then, we have  $\alpha(\sqrt{d}/\beta)^{5/4} \leq 12d^2 \eta_0$  and we can apply Proposition C.2  
 929 with  $\zeta = 12d^{2.5} \eta_0 = 12d^{2.5} \epsilon$ . Similar computations as above then show that  $m \geq d/(8 \cdot 61)$ , with  
 930  $k = \Theta(\frac{M+\ln d}{d \ln \frac{1}{\epsilon}})$ , so that the query lower bound finally becomes

$$Q \geq \Omega \left( \frac{d^3 \ln \frac{1}{\epsilon}}{(M + \ln d) \ln^2 d} \right).$$

931 **Remark C.1.** *The more careful analysis—involving the discretization  $\mathcal{D}_\delta$  of the unit sphere at scale  
 932  $\delta$  instead of the hypercube  $\{\pm 1\}^d$ —allowed to add a  $\ln \frac{1}{\epsilon}$  factor to the final query lower bound but  
 933 also an additional  $\ln d$  factor for both convex-optimization and feasibility-problem results. Indeed,  
 934 the improved Proposition C.2 shows that the OVGH with adequate parameters requires  $\mathcal{O}(d)$  queries,  
 935 instead of  $\mathcal{O}(d/\ln d)$  in [5, Proposition 14]. At a high level, each hint queried brings information  
 936  $\mathcal{O}(d \ln d)$  but memorizing a binary matrix  $\mathbf{A} \in \{\pm 1\}^{\lceil d/4 \rceil \times d}$  only requires  $d^2$  bits of memory: hence  
 937 the query lower bound is limited to  $\mathcal{O}(d/\ln d)$ . Instead, memorizing the matrix  $\mathbf{A}$  where each row  
 938 lies in  $\mathcal{D}_\delta$  requires  $\Theta(d^2 \ln \frac{1}{\epsilon})$  memory, hence querying  $d$  hints (total information  $\mathcal{O}(d^2 \ln d)$ ) is not  
 939 prohibitive for the lower bound.*

### 940 C.1.3 Lower bound of [29] for convex optimization and randomized algorithms

941 We aim to improve the result to obtain the following.

942 **Theorem C.4.** *For  $\epsilon \leq 1/d^4$  and any  $\delta \in [0, 1]$ , any (potentially randomized) algorithm guaranteed  
 943 to minimize 1-Lipschitz convex functions over the unit ball with  $\epsilon$  accuracy uses at least  $d^{1.25-\delta} \ln \frac{1}{\epsilon}$   
 944 bits of memory or makes  $\tilde{\Omega}(d^{1+4\delta/3})$  queries.*

945 The distribution considered in [29] is given by the functions

$$\frac{1}{d^6} \max \left\{ d^5 \|\mathbf{A}\mathbf{x}\|_\infty - 1, \max_{i \in [N]} (\mathbf{v}_i^\top \mathbf{x} - i\gamma) \right\},$$

946 where  $N \leq d$  is a parameter,  $\mathbf{A}$  has  $\lfloor d/2 \rfloor$  rows drawn i.i.d. from  $\mathcal{U}(\{\pm 1\}^d)$ , and the vectors  $\mathbf{v}_i$  are  
 947 drawn i.i.d. from the rescaled hypercube  $\mathbf{v}_i \sim \mathcal{U}(d^{-1/2}\{\pm 1\}^d)$ . We adapt the class of functions by  
 948 simply changing pre-factors as follows

$$\mu \max \left\{ \frac{1}{\mu} \|\mathbf{A}\mathbf{x}\|_\infty - 1, \max_{i \in [N]} (\mathbf{v}_i^\top \mathbf{x} - i\gamma) \right\}, \quad (8)$$

949 where  $\mathbf{A}$  has the same number of rows but they are draw i.i.d. from  $\mathcal{U}_\delta$ , and  $\delta, \mu > 0$  are parameters  
 950 to specify. We use the notation  $\mu$  instead of  $\eta$  as in the previous sections because [29] already use  
 951 a parameter  $\eta$  which in our context can be interpreted as  $\eta = 1/(\mu\sqrt{d})$ . We choose the parameters  
 952  $\mu = 16\sqrt{d}\epsilon$  and  $\delta' = \epsilon$ .

953 Again, as for the previous sections, the original proof can be directly used to show that if an algorithm  
 954 is guaranteed to find a  $\frac{\mu}{16\sqrt{N}}$  ( $\geq \epsilon$ )-suboptimal point for the above function class, there is an algorithm  
 955 that wins at their Orthogonal Vector Game (OVG) [29, Game 1], with the only difference that the  
 956 parameter  $d^{-4}$  (1.8 of OVG) is replaced by  $\sqrt{d}\mu$ . OVG requires the output to be *robustly-independent*  
 957 (defined in [29]) and effectively corresponds to  $\beta = 1/d^2$  in OVGH. As a result, there is a successful  
 958 algorithm for the OVGH with parameters  $\alpha = \sqrt{d}\mu$  and  $\beta = 1/d^2$  and that even completely ignores  
 959 the hints. Hence, we can now directly use Proposition C.2 with  $\zeta = d^{1+25/16}\mu$  (from the assumption  
 960  $\epsilon \leq d^{-4}$  we have  $\zeta \leq 1/\sqrt{d}$ ). This shows that with the adequate choice of  $k = \Theta(\frac{M+d\ln d}{d\ln \frac{1}{\epsilon}})$ , the  
 961 query lower bound is  $\Omega(d)$ .

962 Putting things together, a potentially randomized algorithm for convex optimization that uses  $M$   
 963 memory makes at least the following number of queries

$$Q \geq \Omega\left(\frac{Nd}{k}\right) = \Omega\left(\frac{d^{4/3}}{\ln^{1/3} d} \left(\frac{d \ln \frac{1}{\epsilon}}{M + d \ln d}\right)^{4/3}\right).$$

## 964 C.2 Proof sketch for improving the query-complexity lower bound

965 We now turn to improving the query-complexity lower bound by a factor  $\frac{\ln \frac{1}{\epsilon}}{\ln d}$ . At the high level, the  
 966 idea is to replicate these constructed “difficult” class of functions at  $\frac{\ln \frac{1}{\epsilon}}{\ln d}$  different scales or levels,  
 967 similarly to the manner that the historical  $\Omega(d \ln \frac{1}{\epsilon})$  lower bound is obtained for convex optimization  
 968 [32]. This argument is relatively standard and we only give details in the context of improving the  
 969 bound from [29] for randomized algorithms in convex optimization for conciseness. This result uses  
 970 a simpler class of functions, which greatly eases the exposition. We first present the construction with  
 971 2 levels, then present the generalization to  $p = \Theta(\frac{\ln \frac{1}{\epsilon}}{\ln d})$  levels. For convenience, we write

$$Q(\epsilon; M, d) = \Omega\left(\frac{d^{4/3}}{\ln^{1/3} d} \left(\frac{d \ln \frac{1}{\epsilon}}{M + d \ln d}\right)^{4/3}\right).$$

972 This is the query lower bound given in Theorem C.5 for convex optimization algorithms with memory  
 973  $M$  that optimize the defined class of functions (Eq (8)) to accuracy  $\epsilon$ .

### 974 C.2.1 Construction of a bi-level class of functions $F_{\mathbf{A}, \mathbf{v}_1, \mathbf{v}_2}$ to optimize

975 In the lower-bound proof, [29] introduce the point

$$\bar{\mathbf{x}} = -\frac{1}{2\sqrt{N}} \sum_{i \in [N]} P_{\mathbf{A}^\perp}(\mathbf{v}_i),$$

976 where  $P_{\mathbf{A}^\perp}$  is the projection onto the orthogonal space to the rows of  $\mathbf{A}$ . They show that with failure  
 977 probability at most  $2/d$ ,  $\bar{\mathbf{x}}$  has good function value

$$F_{\mathbf{A}, \mathbf{v}}(\bar{\mathbf{x}}) := \mu \max \left\{ \frac{1}{\mu} \|\mathbf{A}\bar{\mathbf{x}}\|_\infty - 1, \max_{i \in [N]} (\mathbf{v}_i^\top \bar{\mathbf{x}} - i\gamma) \right\} \leq -\frac{\mu}{8\sqrt{N}}.$$

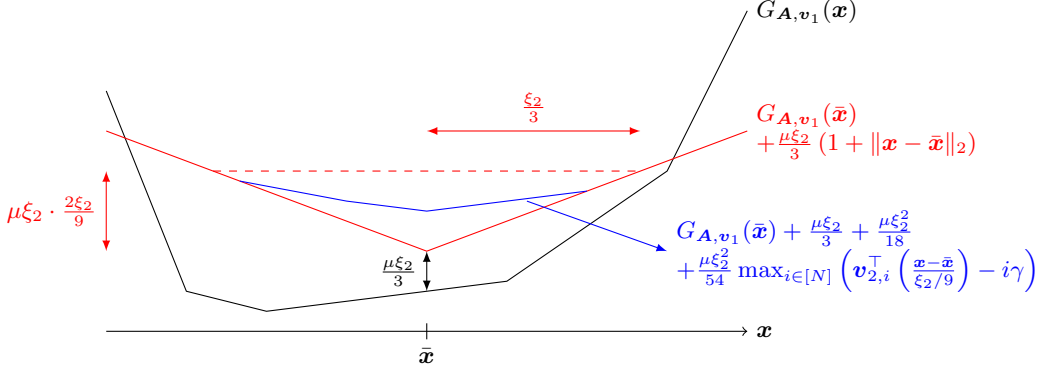


Figure 3: Representation of the procedure to rescale the optimization function.

978 This is shown in [29, Lemma 25]. On the other hand, from Theorem C.4, during the first

$$Q_1 = Q(\epsilon; M, d)$$

979 queries of any algorithm, with probability at least  $1/3$ , all queries are at least  $\mu/(16\sqrt{N})$ -suboptimal  
 980 compared to  $\bar{\mathbf{x}}$  in function value [29, Theorem 28, Lemma 14 and Theorem 16]. Precisely, if  $F_{A,v}$  is  
 981 the sampled function to optimize, with probability at least  $1/3$ ,

$$F_{A,v}(\mathbf{x}_t) \geq F_{A,v}(\bar{\mathbf{x}}) + \frac{\mu}{16\sqrt{N}} \geq F_{A,v}(\bar{\mathbf{x}}) + \frac{\mu}{16\sqrt{d}}, \quad \forall t \leq Q_1.$$

982 As a result, we can replicate the term  $\max_{i \in [N]} (\mathbf{v}_i^\top \mathbf{x} - i\gamma)$  at a smaller scale within the ball  
 983  $B_d(\bar{\mathbf{x}}, 1/(16\sqrt{d}))$ . For convenience, we introduce  $\xi_2 = 1/(16\sqrt{d})$  which will be the scale of the  
 984 duplicate function. We separate the wall term  $\|\mathbf{A}\mathbf{x}\|_\infty - \mu$  for convenience. Hence, we define

$$G_{A,v_1}(\mathbf{x}) := \mu \max_{i \in [N]} (\mathbf{v}_{1,i}^\top \mathbf{x} - i\gamma)$$

$$G_{A,v_1,v_2}(\mathbf{x}) := \max \left\{ G_{A,v_1}(\mathbf{x}), G_{A,v_1}(\bar{\mathbf{x}}) + \frac{\mu\xi_2}{3}, \right.$$

$$\left. \max \left\{ 1 + \|\mathbf{x} - \bar{\mathbf{x}}\|_2, 1 + \frac{\xi_2}{6} + \frac{\xi_2}{18} \max_{i \in [N]} \left( \mathbf{v}_{2,i}^\top \left( \frac{\mathbf{x} - \bar{\mathbf{x}}}{\xi_2/9} \right) - i\gamma \right) \right\} \right\}$$

985 An illustration of the construction is given in Fig. 3. The resulting optimization functions are given  
 986 by adding the wall term:

$$F_{A,v_1}(\mathbf{x}) = \max \{ \|\mathbf{A}\mathbf{x}\|_\infty - \mu, G_{A,v_1}(\mathbf{x}) \}$$

$$F_{A,v_1,v_2}(\mathbf{x}) = \max \{ \|\mathbf{A}\mathbf{x}\|_\infty - \mu, G_{A,v_1,v_2}(\mathbf{x}) \}$$

987 We first explain the choice of parameters. First observe that since  $\|\mathbf{A}\bar{\mathbf{x}}\| = 0$ , we have  $G_{A,v_1}(\bar{\mathbf{x}}) =$   
 988  $F_{A,v_1}(\bar{\mathbf{x}})$ . We can then check that for all  $\mathbf{x} \in B_d(0, 1)$ ,

$$G_{A,v_1,v_2}(\mathbf{x}) \leq \max \left\{ G_{A,v_1}(\mathbf{x}), G_{A,v_1}(\bar{\mathbf{x}}) + \frac{2}{3}\mu\xi_2 \right\}. \quad (9)$$

989 Further, for any  $\mathbf{x} \in B_d(\bar{\mathbf{x}}, \xi_2/3)$ , since  $F_{A,v_1}$  is 1-Lipschitz, we can easily check that

$$G_{A,v_1,v_2}(\mathbf{x}) - G_{A,v_1}(\bar{\mathbf{x}})$$

$$= \frac{\mu\xi_2}{3} \max \left\{ 1 + \|\mathbf{x} - \bar{\mathbf{x}}\|_2, 1 + \frac{\xi_2}{6} + \frac{\xi_2}{18} \max_{i \in [N]} \left( \mathbf{v}_{2,i}^\top \left( \frac{\mathbf{x} - \bar{\mathbf{x}}}{\xi_2/9} \right) - i\gamma \right) \right\} \leq \frac{2}{3}\mu\xi_2.$$

990 Thus,  $G_{A,v_1,v_2}(\mathbf{x})$  does not coincide with  $G_{A,v_1}(\mathbf{x})$  on  $B_d(\bar{\mathbf{x}}, \xi_2/3)$ . Then, the  $\|\mathbf{x} - \bar{\mathbf{x}}\|_2$  term  
 991 ensures that any minimizer of  $G_{A,v_1,v_2}$  is contained within the closed ball  $B_d(\bar{\mathbf{x}}, \xi_2/3)$ . Also,  
 992 to obtain a  $\mu\xi_2/3$ -suboptimal solution of  $F_{A,v_1,v_2}$ , the algorithm needs to find what would be a  
 993  $\mu\xi_2$ -suboptimal solution of  $F_{A,v_1}$ , while receiving the same response as when optimizing the latter.

994 Next, for any  $\mathbf{x} \in B_d(\bar{\mathbf{x}}, \xi_2/9)$ , the term  $\max_{i \in [N]} \left( \mathbf{v}_{2,i}^\top \left( \frac{\mathbf{x} - \bar{\mathbf{x}}}{\xi_2/9} \right) - i\gamma \right)$  lies in  $[-1, 1]$ . Hence, we  
 995 can check that for  $\mathbf{x} \in B_d(\bar{\mathbf{x}}, \xi_2/9)$ ,

$$G_{\mathbf{A}, v_1, v_2}(\mathbf{x}) = G_{\mathbf{A}, v_1}(\bar{\mathbf{x}}) + \frac{\mu \xi_2}{3} + \frac{\mu \xi_2^2}{18} + \frac{\mu \xi_2^2}{54} \max_{i \in [N]} \left( \mathbf{v}_{2,i}^\top \left( \frac{\mathbf{x} - \bar{\mathbf{x}}}{\xi_2/9} \right) - i\gamma \right). \quad (10)$$

996 We now argue that  $F_{\mathbf{A}, v_1, v_2}$  acts as a duplicate function. Until the algorithm reaches a point with  
 997 function value at most  $G_{\mathbf{A}, v_1}(\bar{\mathbf{x}}) + \mu \xi_2$ , the optimization algorithm only receives responses consistent  
 998 with the function  $F_{\mathbf{A}, v_1}$  by Eq (9). Next, all minimizers of  $F_{\mathbf{A}, v_1, v_2}$  are contained in  $B_d(\bar{\mathbf{x}}, \xi_2/3)$ ,  
 999 which was the goal of introducing the term in  $\|\mathbf{x} - \bar{\mathbf{x}}\|_2$ . As a result, optimizing  $F_{\mathbf{A}, v_1, v_2}$  on this  
 1000 ball is equivalent to minimizing

$$\tilde{F}_{\mathbf{A}, v_2}(\mathbf{y}) = \max \left\{ \|\mathbf{A}\mathbf{y}\|_\infty - \mu_2, c_2 + \nu_2 \max_{i \in [N]} (\mathbf{v}_{2,i}^\top \mathbf{y} - i\gamma), c'_2 + \nu'_2 \|\mathbf{y}\| \right\}, \quad \mathbf{y} \in B_d(0, 3),$$

1001 where  $\mathbf{y} = \frac{\mathbf{x} - \bar{\mathbf{x}}}{\xi_2/9}$ . The function has been rescaled by a factor  $\xi_2/9$  compared to  $F_{\mathbf{A}, v_1, v_2}$  so that  
 1002  $\mu_2 = \frac{9\mu}{\xi_2}$ ,  $\nu_2 = \frac{\mu \xi_2}{6}$ ,  $\nu'_2 = 6\mu$ ,  $c_2 = \frac{9}{\xi_2} G_{\mathbf{A}, v_1}(\bar{\mathbf{x}}) + 3\mu + \frac{\mu \xi_2}{2}$ , and  $c'_2 = \frac{9}{\xi_2} G_{\mathbf{A}, v_1}(\bar{\mathbf{x}}) + 3\mu$ . By  
 1003 Eq (10), the two first terms of  $\tilde{F}_{\mathbf{A}, v_1}$  are preponderant for  $\mathbf{y} \in B_d(0, 1)$ .

1004 The form of  $\tilde{F}_{\mathbf{A}, v_2}$  is very similar to the original form of functions

$$F_{\mathbf{A}, v_2} = \max \left\{ \|\mathbf{A}\mathbf{y}\|_\infty - \mu'_1, \mu'_2 \max_{i \in [N]} (\mathbf{v}_{2,i}^\top \mathbf{y} - i\gamma) \right\},$$

1005 In fact, the same proof structure for the query-complexity/memory lower-bound can be applied in  
 1006 this case. The main difference is that originally one had  $\mu'_1 = \mu'_2$ ; here we would instead have  
 1007  $\mu'_1 = \mu_2 + c_2 = \Theta(\mu/\xi_2)$  and  $\mu'_2 = \nu_2 = \Theta(\mu \xi_2)$ . Intuitively, this corresponds to increasing the  
 1008 accuracy to  $\Theta(\epsilon \xi_2^2)$ —a factor  $\xi_2$  is due to the fact that  $\tilde{F}_{\mathbf{A}, v_2}$  was rescaled by a factor  $\xi_2/9$  compared  
 1009 to  $F_{\mathbf{A}, v_1, v_2}$ , and a second factor  $\xi_2$  is due to the fact that within  $\tilde{F}_{\mathbf{A}, v_2}$ , we have  $\mu'_2 = \Theta(\mu \xi_2)$ —while  
 1010 the query lower bound is similar to that obtained for  $\Theta(\epsilon/\xi_2)$ . As a result, during the first

$$Q_2 = Q \left( \Theta \left( \frac{\epsilon}{\xi_2} \right); M, d \right)$$

1011 queries of any algorithm optimizing  $\tilde{F}_{\mathbf{A}, v_2}$ , with probability at least  $1/3$  on the sample of  $\mathbf{A}$  and  $\mathbf{v}_2$ ,  
 1012 all queries are at least  $\Theta(\epsilon \xi_2)$ -suboptimal compared to

$$\bar{\mathbf{y}} = -\frac{1}{2\sqrt{N}} \sum_{i \in [N]} P_{\mathbf{A}^\perp}(\mathbf{v}_{2,i}).$$

1013 We are now ready to give lower bounds on the queries of an algorithm minimizing  $F_{\mathbf{A}, v_1, v_2}$  to  
 1014 accuracy  $\Theta(\epsilon \xi_2^2)$ . Let  $T_2$  be the index of the first query with function value at most  $G_{\mathbf{A}, v_1}(\bar{\mathbf{x}}) + \mu \xi_2$ .  
 1015 We already checked that before that query, all responses of the oracle are consistent with minimizing  
 1016  $F_{\mathbf{A}, v_1}$ , hence on an event  $\mathcal{E}_1$  of probability at least  $1/3$ , one has  $T_2 \geq Q_1$ . Next, consider the  
 1017 hypothetical case when at time  $T_2$ , the algorithm is also given the information of  $\bar{\mathbf{x}}$  and is allowed to  
 1018 store this vector. Given this information, optimizing  $F_{\mathbf{A}, v_1, v_2}$  reduces to optimizing  $\tilde{F}_{\mathbf{A}, v_2}$  since we  
 1019 already know that the minimum is achieved within  $B_d(\bar{\mathbf{x}}, \xi_2/3)$ . Further, any query outside of this  
 1020 ball either

- 1021 • returns a vector  $\mathbf{v}_{1,i}$  which does not give any useful information for the minimization ( $\mathbf{v}_1$   
 1022 and  $\mathbf{v}_2$  are sampled independently and  $\bar{\mathbf{x}}$  is given),
- 1023 • or returns a row from  $\mathbf{A}$ , as covered by the original proof.

1024 Hence, on an event  $\mathcal{E}_2$  of probability at least  $1/3$ , even with the extra information of  $\bar{\mathbf{x}}$ , during the  
 1025 next  $Q_2$  queries starting from  $T_2$ , the algorithm does not query a  $\Theta(\mu \xi_2^3)$ -suboptimal solution to  
 1026  $F_{\mathbf{A}, v_1, v_2}$ . This holds a fortiori for the model when the algorithm is not given  $\bar{\mathbf{x}}$  at time  $T_2$ .

1027 **C.2.2 Recursive construction of a  $p$ -level class of functions  $F_{\mathbf{A}, \mathbf{v}_1, \dots, \mathbf{v}_p}$**

1028 Similarly as in the last section, one can inductively construct the sequence of functions  $F_{\mathbf{A}, \mathbf{v}_1}$ ,  
 1029  $F_{\mathbf{A}, \mathbf{v}_1, \mathbf{v}_2}$ ,  $F_{\mathbf{A}, \mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3}$ , etc. Formally, the induction is constructed as follows: let  $(\mathbf{v}_p)_{p \geq 1}$  be an i.i.d.  
 1030 sequence of  $N$  i.i.d. vectors  $(\mathbf{v}_{k,i})_{i \in [N]}$  sampled from the rescaled hypercube  $d^{-1/2} \{\pm 1\}^d$ . Next,  
 1031 we pose

$$G_{\mathbf{A}, \mathbf{v}_1}(\mathbf{x}) = \mu^{(1)} \max_{i \in [N]} \left( \mathbf{v}_{1,i}^\top \left( \frac{\mathbf{x} - \bar{\mathbf{x}}^{(1)}}{s^{(1)}} \right) - i\gamma \right),$$

1032 where  $\mu^{(1)} = \mu$ ,  $\bar{\mathbf{x}}^{(1)} = \mathbf{0}$  and  $s^{(1)} = 1$ . For  $k \geq 1$ , we pose

$$\bar{\mathbf{x}}^{(k+1)} = \bar{\mathbf{x}}^{(k)} - \frac{s^{(k)}}{2\sqrt{N}} \sum_{i \in [N]} P_{\mathbf{A}^\perp}(\mathbf{v}_{k,i}), \quad \text{and} \quad F^{(k)} := G_{\mathbf{A}, \mathbf{v}_1, \dots, \mathbf{v}_k}(\bar{\mathbf{x}}^{(k)}) + \mu^{(k)} \xi_{k+1},$$

1033 for a certain parameter  $\xi_{k+1}$  to be specified. We then define the next level as

$$G_{\mathbf{A}, \mathbf{v}_1, \dots, \mathbf{v}_{k+1}}(\mathbf{x}) := \max \left\{ G_{\mathbf{A}, \mathbf{v}_1, \dots, \mathbf{v}_k}(\mathbf{x}), G_{\mathbf{A}, \mathbf{v}_1, \dots, \mathbf{v}_k}(\bar{\mathbf{x}}^{(k+1)}) + \frac{\mu^{(k)} \xi_{k+1}}{3}, \right. \\ \left. \max \left\{ 1 + \frac{\|\mathbf{x} - \bar{\mathbf{x}}^{(k+1)}\|_2}{s^{(k)}}, 1 + \frac{\xi_{k+1}}{6} + \frac{\xi_{k+1}}{18} \max_{i \in [N]} \left( \mathbf{v}_{k+1,i}^\top \left( \frac{\mathbf{x} - \bar{\mathbf{x}}^{(k+1)}}{s^{(k)} \xi_{k+1}/9} \right) - i\gamma \right) \right\} \right\}.$$

1034 We then pose  $\mu^{(k+1)} := \mu^{(k)} \xi_{k+1}^2 / 54$  and  $s^{(k+1)} := s^{(k)} \xi_{k+1} / 9$ , which closes the induction. The  
 1035 optimization functions are defined simply as

$$F_{\mathbf{A}, \mathbf{v}_1, \dots, \mathbf{v}_{k+1}}(\mathbf{x}) = \max \{ \|\mathbf{A}\mathbf{x}\|_\infty - \mu, G_{\mathbf{A}, \mathbf{v}_1, \dots, \mathbf{v}_{k+1}}(\mathbf{x}) \}.$$

1036 We checked before that we can use  $\xi_2 = 1/(16\sqrt{d})$ . For general  $k \geq 0$ , given that the form of  
 1037 the function slightly changes to incorporate the absolute term (see  $\tilde{F}_{\mathbf{A}, \mathbf{v}_2}$ ), this constant may differ  
 1038 slightly. In any case, one has  $\xi_k = \Theta(1/\sqrt{d})$ . Now fix a construction level  $p \geq 1$  and for any  $k \in [p]$ ,  
 1039 let  $T_k$  be the first time that a point with function value at most  $F^{(k)}$  is queried. For convenience  
 1040 let  $T_0 = 0$ . Using the same arguments as above recursively, we can show that on an event  $\mathcal{E}_k$  with  
 1041 probability at least  $1/3$ ,

$$T_k - T_{k-1} \geq Q_k = Q \left( \Theta \left( \frac{\mu}{s^{(k)}} \right); M, d \right)$$

1042 Next note that the sequence  $F^{(k)}$  is decreasing and by construction, if one finds a  $\mu^{(p)} \xi_{p+1}$ -suboptimal  
 1043 point of  $F_{\mathbf{A}, \mathbf{v}_1, \dots, \mathbf{v}_p}$ , then this point has value at most  $F^{(p)}$ . As a result, for an algorithm that finds a  
 1044  $\mu^{(p)} \xi_{p+1}$ -suboptimal point, the times  $T_0, \dots, T_p$  are all well defined and non-decreasing. We recall  
 1045 that  $\mu = \Theta(\sqrt{d}\epsilon)$ . Therefore, we can still have  $\mu/s^{(p)} \leq \sqrt{\epsilon}$  and  $\mu^{(p)} \xi_{p+1} \geq \epsilon^2$  for  $p = \Theta(\frac{\ln \frac{1}{\epsilon}}{\ln d})$ .  
 1046 Combining these observations, we showed that when optimizing the functions  $F_{\mathbf{A}, \mathbf{v}_1, \dots, \mathbf{v}_p}$  to accuracy  
 1047  $\Theta(\mu^{(p)} \xi_{p+1}) = \Omega(\epsilon^2)$ , the total number of queries  $Q$  satisfies

$$\mathbb{E}[Q] \geq \frac{1}{3} \sum_{k \in [p]} Q_k \geq \frac{p}{3} Q(\sqrt{\epsilon}; M, d) = \Theta \left( \frac{d^{4/3} \ln \frac{1}{\epsilon}}{\ln^{4/3} d} \left( \frac{d \ln \frac{1}{\epsilon}}{M + d \ln d} \right)^{4/3} \right).$$

1048 Changing  $\epsilon$  to  $\epsilon^2$  proves the desired result.

1049 **Theorem C.5.** For  $\epsilon \leq 1/d^8$  and any  $\delta \in [0, 1]$ , any (potentially randomized) algorithm guaranteed  
 1050 to minimize 1-Lipschitz convex functions over the unit ball with  $\epsilon$  accuracy uses at least  $d^{1.25-\delta} \ln \frac{1}{\epsilon}$   
 1051 bits of memory or makes  $\tilde{\Omega}(d^{1+4\delta/3} \ln \frac{1}{\epsilon})$  queries.

1052 The same recursive construction can be applied to the results from Theorems C.1 and C.3 to improve  
 1053 their oracle-complexity lower bounds by a factor  $\frac{\ln \frac{1}{\epsilon}}{\ln d}$ , albeit with added technicalities due to the  
 1054 adaptivity of their class of functions. This yields Theorem 3.3.



**Input:** Number of iterations  $T$ , computation accuracy  $\eta \leq 1$ , target accuracy  $\epsilon \leq 1$   
Initialize:  $\mathbf{x} = \mathbf{0}$ ;  
**for**  $t = 0, \dots, T$  **do**  
    Query the oracle at  $\mathbf{x}$   
    **if**  $\mathbf{x}$  *successful* **then return**  $\mathbf{x}$ ;  
    Receive a separation vector  $\mathbf{g}$  with accuracy  $\eta$   
    Update  $\mathbf{x}$  as  $\mathbf{x} - \epsilon\mathbf{g}$  up to accuracy  $\eta$   
**end**  
**return**  $\mathbf{x}$

**Algorithm 11:** Memory-constrained gradient descent

## 1055 D Memory-constrained gradient descent for the feasibility problem

1056 In this section, we prove a simple result showing that memory-constrained gradient descent applies  
1057 to the feasibility problem. We adapt the algorithm described in [49].

1058 We now prove that this memory-constrained gradient descent gives the desired result of Proposi-  
1059 tion 3.1.

1060 *Proof of Proposition 3.1.* Denote by  $\mathbf{x}_t$  the state of  $\mathbf{x}$  at iteration  $t$ , and  $\mathbf{g}_t$  (resp.  $\tilde{\mathbf{g}}_t$ ) the separation  
1061 oracle without rounding errors (resp. with rounding errors) at  $\mathbf{x}_t$ . By construction,

$$\|\mathbf{x}_{t+1} - (\mathbf{x}_t + \epsilon\tilde{\mathbf{g}}_t)\| \leq \eta \quad \text{and} \quad \|\tilde{\mathbf{g}}_t - \mathbf{g}_t\| \leq \eta. \quad (11)$$

1062 As a result, recalling that  $\|\mathbf{g}_t\| = 1$ ,

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq (\|\mathbf{x}_t + \epsilon\tilde{\mathbf{g}}_t - \mathbf{x}^*\| + \eta)^2 \leq (\|\mathbf{x}_t + \epsilon\mathbf{g}_t - \mathbf{x}^*\| + (1 + \epsilon)\eta)^2 \leq \|\mathbf{x}_t + \epsilon\mathbf{g}_t - \mathbf{x}^*\|^2 + 20\eta.$$

1063 By assumption,  $Q$  contains a ball  $B_d(\mathbf{x}^*, \epsilon)$  for  $\mathbf{x}^* \in B_d(0, 1)$ . Then, because  $\mathbf{g}_t$  separates  $\mathbf{x}_t$  from  
1064  $B_d(\mathbf{x}^*, \epsilon)$ , one has  $\mathbf{g}_t^\top(\mathbf{x}^* - \mathbf{x}_t) \geq \epsilon$ . Therefore,

$$\begin{aligned} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 &\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 + 2\epsilon\mathbf{g}_t^\top(\mathbf{x}_t - \mathbf{x}^*) + \epsilon^2\|\mathbf{g}_t\|^2 + 20\eta \\ &\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \epsilon^2 + 20\eta. \end{aligned}$$

1065 Then, take  $\eta = \epsilon^2/40$  and  $T = \frac{8}{\epsilon^2}$ . If iteration  $T$  was performed, we have using the previous equation

$$\|\mathbf{x}_T - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \frac{\epsilon^2}{2}T \leq 4 - \frac{\epsilon^2}{2}T \leq 0.$$

1066 Hence,  $\mathbf{x}_T$  is an  $\epsilon$ -suboptimal solution.

1067 We now turn to the memory usage of gradient descent. It only needs to store  $\mathbf{x}$  and  $\mathbf{g}$  up to the desired  
1068 accuracy  $\eta = \mathcal{O}(\epsilon^2)$ . Hence, this storage and the internal computations can be done with  $\mathcal{O}(d \ln \frac{d}{\epsilon})$   
1069 memory. Because we suppose that  $\epsilon \leq \frac{1}{\sqrt{d}}$ , this gives the desired result.  $\square$