

---

# Randomized Benchmarking of Local Zeroth-Order Optimizers for Variational Quantum Systems

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In the field of quantum information, classical optimizers play an important role.  
2 From experimentalists optimizing their physical devices to theorists exploring  
3 variational quantum algorithms, many aspects of quantum information require the  
4 use of a classical optimizer. For this reason, there are many papers that benchmark  
5 the effectiveness of different optimizers for specific quantum learning tasks and  
6 choices of parameterized algorithms. However, for researchers exploring new  
7 algorithms or physical devices, the insights from these studies don't necessarily  
8 translate. To address this concern, we compare the performance of classical  
9 optimizers across a series of partially-randomized tasks to more broadly sample the  
10 space of quantum learning problems. We focus on local zeroth-order optimizers  
11 due to their generally favorable performance and query-efficiency on quantum  
12 systems. We discuss insights from these experiments that can help motivate future  
13 works to improve these optimizers for use on quantum systems.

## 14 1 Introduction

15 Quantum computing has over time gathered more and more attention from researchers for the  
16 promise of significant computational speedups relative to classical computers. This has spurred many  
17 developments across all fronts in the field, from algorithms to building real quantum computers.  
18 However, many of these works still rely on the use of classical optimizers. For instance, variational  
19 quantum algorithms are a class of algorithms that have parameters that are then optimized by a  
20 classical optimizer [Cerezo et al., 2021]. These algorithms exist both as a way to do machine learning  
21 in a quantum system and to realize practically useful algorithms in noisy near-term devices. And  
22 beyond the scope of algorithm design, classical optimizers also play a part in assisting experimentalists  
23 working in quantum information. Beyond simply being used in practice to realize theorized algorithms  
24 on real devices [Ebadi et al., 2022], optimizers can tune the control of physical actions (such as laser  
25 pulses, injections of electrical current, etc.) that all need to be controlled precisely to produce desired  
26 quantum operation [Coopmans et al., 2021, Leng et al., 2023].

27 For these reasons, understanding how classical optimizers interact with quantum objects is both  
28 difficult and highly important for designing the best quantum devices / algorithms. As a result, there  
29 are many studies on exactly that. Some of these studies are general benchmarks that compare a wide  
30 variety of optimizers [Pellow-Jarman et al., 2021, Anand et al., 2021, Singh et al., 2023]; others are  
31 works that propose new optimizers for quantum circuits and show experimental evidence for the  
32 advantages of their optimizer [Sung et al., 2020, Gacon et al., 2021, Leng et al., 2023]. However,

33 most of these works benchmark for a fixed set of problems using a specific form of parameterized  
34 quantum ansatzes / models intended for each problem.

35 While these works are important for understanding specific use cases, it's not necessarily clear  
36 how much the insight from these works translates to new scenarios. And in the case of physics  
37 experimentalists, similar concerns arise when there may be noise or aspects of their system that shift  
38 over time [Proctor et al., 2020, Blume-Kohout et al., 2020]. So unless they constantly re-evaluate  
39 many optimizers on their own, what sort of confidence can they have in their choice of optimizer  
40 being wise?

41 These are the questions we take a shot at addressing in this work. We do this by benchmarking on  
42 tasks that are randomized. So in addition to the random parameter initialization common in other  
43 benchmarks, we also randomize the parameterized circuit / ansatz that is used, and in some cases  
44 randomize parts of the objective we are trying to minimize. While this is certainly not the most  
45 perfect way to answer the questions we posed before, we hope that by adding more variety via  
46 randomness in our benchmark we can begin to identify features of optimization algorithms that work  
47 more generically on many types of variational quantum learning problems.

48 However, when doing a study like this, you run into the risk of benchmarking something so generic  
49 that it's difficult to get any concrete insights from results. For this reason, we narrow our focus onto  
50 understanding how to improve a specific class of optimizers. First, we only consider **zeroth-order**  
51 methods. As mentioned previously, this is because methods that only sample and don't require  
52 gradients currently tend to be more easily realizable on quantum systems. Second, we use only  
53 **local** optimizers. This means that our optimizers sample the objective centered around a specific  
54 "candidate" point. Last, we only consider **sample-efficient** methods. This means that our optimizers  
55 make optimization decisions based on sampling as few points as possible per step. These choices  
56 mostly centered around us deciding to focus on studying the SPSA algorithm [Spall, 1998] and  
57 optimizers like it, because its generally favorable performance and runtime efficiency on quantum  
58 systems.

59 In this study, we benchmark randomized experiments for a variety of Hamiltonian minimization  
60 and generative modeling tasks. We side-by-side compare the performance of 7 optimizers: SPSA,  
61 AdamSPSA, 2-SPSA, QNSPSA, GES, xNES, and sNES. We produce plots illustrating both the  
62 average rate of convergence and statistics on the end-result performance of each optimizer, and  
63 discuss our thoughts on insights to be gained from these results. But in short, we believe there are two  
64 **main take-aways**. First, **more sophisticated optimizers are not generally better**. In our benchmarks  
65 SPSA tends to perform best overall, followed by the other simpler heuristic methods like AdamSPSA  
66 and GES. There is more nuance to this statement and it certainly isn't true in all circumstances, but  
67 under our randomized tasks simpler methods tended to be more reliably effective. Second,  
68 **there is a need for more robust or adaptive optimization heuristics**. While there are heuristics that  
69 can assist these optimizers in optimizing quickly in certain parts of the optimization process, at other  
70 points these heuristics can begin to hurt the optimizer's performance. As such, we argue that it would  
71 be beneficial to make these heuristics more robust to distribution shifts. But more broadly, we hope  
72 that this work helps stimulate thought into how to more generally compare and study optimizers of  
73 quantum systems.

## 74 **2 Prior Work**

75 There are a number of works that explicitly benchmark a variety of optimizers. Pellow-Jarman et al.  
76 [2021] compares a variety of both gradient and gradient-free optimizers on variational quantum linear  
77 solver problems, both in the presence and absence of noise. They show that while there's no clear  
78 best optimizer, SPSA tends to perform favorably in realistic noise scenarios. Anand et al. [2021]  
79 benchmarks natural evolutionary strategies (NES) on variational quantum eigensolver (VQE) and  
80 state preparation problems. They also empirically investigate and provide some justification for  
81 how NES could be used in a hybrid algorithm to assist gradient-based optimizers in barren plateau

82 regimes. Singh et al. [2023] benchmarks optimizers for a variety of quantum chemistry tasks. Like  
83 other studies, there’s no clear best algorithm, but SPSA tends to perform well in noisy conditions.

84 Additionally, while not explicitly a benchmark, a number of works compare optimizers against  
85 a variety of tasks. [Sung et al., 2020] introduces methods that use quadratic fitting of sampled  
86 points to evaluate the gradient and perform gradient descent and policy gradient descent. They  
87 additionally benchmark these methods against a variety of optimizers for three unique Hamiltonian-  
88 minimization problems with specific ansatzes. They also include some more practical considerations,  
89 such as the cost of evaluating different Hamiltonian measurements, the possibility of parallelizing  
90 multiple quantum circuit evaluations, and doing more robust hyperparameter tuning. In their results  
91 their method performs best, but SPSA can come close and often out-performs other methods in  
92 success rates. [Gacon et al., 2021], which proposes the QNSPSA algorithm we use in this study,  
93 also compares its performance to original SPSA on a variety of tasks and compares robustness to  
94 parameter initialization. [Leng et al., 2023] does the same for their proposed AdamSPSA to SPSA  
95 and similar finite-difference methods, but they instead compare on the task of tuning the performance  
96 of a qubit operation on quantum computer.

### 97 **3 Benchmarks**

98 In this section we outline each of the benchmarks we perform in this paper, motivate the reasoning  
99 behind each experiment choice, and provide the finer details of each. While we overall strive to  
100 include some aspect of randomness / broadness, we make a few distinct choice of fixing specific  
101 elements between different benchmarks. Some of these choices are just so we can help distinguish  
102 any differences between types of learning problems, and others are so we try to understand differences  
103 between different levels of difficulty within a type of learning problem.

104 To ensure fairness, for all experimental runs, all randomness (initialization parameters + random  
105 circuits / Hamiltonians / distributions) are controlled by the random key. So although each run is  
106 randomly sampled, because we use the same random keys across all optimizers, they all run the same  
107 variations of each problem. The statistics of each run only vary as a result of the differences between  
108 each optimizer.

109 We also want to ensure that our results aren’t biased by a poor choice of hyperparameters. However,  
110 especially because we are doing highly randomized tasks, it’s difficult to identify what’s a "good"  
111 choice of hyperparameter means. And even if they could be identified, it’s not always reasonable  
112 to assume the user of said optimizer would be able to properly find them. Our compromise is to  
113 do hyperparameter tuning only on a small subset of the possible problem space. We select the  
114 hyperparameters according to a random search run on 3 random keys. This means that as we try  
115 random hyperparameter combinations, they will be tried on 3 different random configurations of the  
116 learning problem. So while this is not as exhaustive as the 100 we test on, the tuning isn’t heavily  
117 biased to a single random problem sample. When there is a range of hyperparameters that all perform  
118 optimally, we bias our choice towards the default values typically used by the algorithm’s authors or  
119 commonly selected in the literature. Once we have the tuned hyperparameters, we benchmark each  
120 optimizer on each problem using 100 runs. All models in all experiments have their parameters  
121 initialized from a normal distribution of mean 0 and standard deviation  $\pi$ .

#### 122 **3.1 Hamiltonian Minimization Experiments**

123 First we run experiments on Hamiltonian minimization problems. This means that we choose a  
124 Hamiltonian as an observable, and the expected value of measuring this Hamiltonian becomes  
125 the "loss" with which we are aiming to produce a quantum state that minimizes this loss. These  
126 benchmarks are meant to encapsulate use-cases such as variational quantum eigensolvers and quantum  
127 optimization problems that map some problem to a specific Hamiltonian. We produce our candidate  
128 states by parameterizing a quantum circuit and optimizing it to map a simple state (usually  $|0\rangle$ ) to the  
129 state we measure with the Hamiltonian observable.

130 It is also important to note that we are simulating the noise-free version of this problem, as we assume  
 131 we have access to the exact expected value of the Hamiltonian to minimize. While this is certainly  
 132 not a realistic assumption, we wanted to first focus on how the aspects of each optimizer affect  
 133 performance on quantum systems first before considering varying levels of noise as a factor.

134 All of the experiments we run here use random circuits as their quantum circuit / ansatz. Specifically  
 135 we use the RandomLayers circuit by PennyLane [Bergholm et al., 2022], which are layers of randomly  
 136 placed parameterized single qubit X, Y, or Z rotation gates mixed with randomly placed CNOT gates.

137 **1D Ising Model:** The first set of experiments we run use the 1D Ising model as our Hamiltonian  
 138 observable. Specifically, we use

$$\mathcal{H} = - \sum_{i=1}^N Z_i \otimes Z_{i+1} - \frac{1}{2} \sum_{i=1}^N X_i.$$

139 As the 1D ising model is known to be an easily solvable problem, the intention of these experiments is  
 140 to provide a simple baseline to understand how our optimizers perform on easier quantum optimization  
 141 problems. These experiments are run on systems of 3 qubits, with circuits of 30 parameterized single  
 142 qubit gates and 10 2-qubit gates. This is relatively simpler and over-parameterized compared to the  
 143 other experiments in this section. Each run is executed for 500 update steps of the optimizer.

144 **2D Heisenberg Model:** Next we benchmark our optimizers on the 2D Heisenberg model. This  
 145 serves as our harder problem, as not only do we include additional interaction terms w.r.t. the 1D  
 146 Ising model, but now we also increase the dimensionality of the connectivity of our observable to a  
 147 2D lattice. We specifically use

$$\mathcal{H} = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \sum_{M \in \{X, Y, Z\}} M_{i,j} \otimes M_{i+1,j} + M_{i,j} \otimes M_{i,j+1} - \frac{1}{4} \sum_{i=1}^N \sum_{j=1}^N Z_{i,j}.$$

148 In contrast to our 1D Ising experiments, these experiments are intended to gain insight into how our  
 149 optimizers perform on a much more difficult problem. We use 9 qubit systems in these experiments,  
 150 with quantum circuits containing 162 parameterized single qubit gates and 49 2-qubit gates. Each  
 151 run is executed for 2000 update steps of the optimizer. One other import distinction is that these  
 152 runs were hyper-parameter tuned for 1000 update steps, but we increased the experiments to 2000  
 153 steps because a few runs seemed to not be fully converging. We believe this lead to some interesting  
 154 side-effects which we discuss in section 5.1.

155 **Randomized Hamiltonians:** Lastly, we run experiments on randomized Hamiltonians. Specifically,  
 156 we generate each Hamiltonian by combining single qubit measurement terms (randomly sampled Pauli  
 157 X, Y, or Z gates on random qubits) with 2-qubit measurement terms (tensor products of randomly  
 158 sampled Pauli X, Y, or Z gates on random qubits). We define the Hamiltonian as

$$\mathcal{H} = \sum_{i=1}^{N_d} c_i (A_i \otimes B_i) + \sum_{i=1}^{N_s} s_i S_i$$

$$c_i, s_i \sim \mathcal{N}(0, \pi) \quad A_i, B_i, S_i \sim \mathcal{U} \left( \bigcup_{i \in [N]} \{X_i, Y_i, Z_i\} \right).$$

159 This benchmark exists for two purposes. First, it is intended to be a problem that is of "medium  
 160 hardness" that is in-between the 1D Ising and the 2D Heisenberg experiments. Second, by adding  
 161 randomness not only to the quantum circuit but also the objective, we hope to gain some additional  
 162 coverage of many possible hamiltonian minimization problems than we did through the prior experi-  
 163 ments to see if the insights from them have some evidence of generalization. These experiments are  
 164 run on systems of 10 qubits. For each run we randomly sample 10 single qubit terms and 20 2-qubit  
 165 terms to construct every Hamiltonian. Note that this means, unlike the prior two set of experiments,  
 166 the energy objective for each optimization procedure differs between run to run. For the random  
 167 circuit ansatz, we use 30 parameterized single qubit gates and 10 2-qubit gates. Each run is performed  
 168 for 500 update steps.

## 169 3.2 Quantum Generative Modelling Experiments

170 To study a larger variety of quantum learning problems, we also investigate quantum generative  
171 modeling. In this setting, instead of trying to produce a state that minimizes and observable, we  
172 desire to produce a state that when measured in the full computational basis, matches a provided  
173 target distribution. Like with Hamiltonian minimization, we produce our candidate states by using  
174 a parameterized circuit. For the same reasons mentioned in section 3.1, we focus on the noise-less  
175 setting where we assume we have direct access to the true loss function. In this case we use the  
176 negative log-likelihood (NLL) loss.

177 **Cardinality Constrained Distribution + QCBM:** Our first generative experiment is using the  
178 Quantum circuit Born machine [Benedetti et al., 2019] to learn a cardinality-constrained distribution.  
179 (So the only randomness in this benchmark is initialization parameters.) The purpose of this experiment  
180 is to serve as our baseline for quantum generative results. Because there aren't many optimizer  
181 benchmarks for quantum generative modelling, these results on a more standard test model and  
182 problem can help us interpret future more heavily randomized results. Specifically we run on a 10  
183 qubit system, with 10 layers of 1 and 2 qubit gates of the QCBM ansatz (illustrated as  $L = 20$  in  
184 figure 1 of Gili et al. [2023].) The cardinality we constrain to for our distribution is 5, meaning that  
185 our target distribution is the uniform distribution over any measurement of all 10 qubits that has 5 1's  
186 in the measurement result. Each run is executed for 5000 optimizer steps.

187 **Randomized Distribution + Random Circuits:** Our next generative experiment is a fully random-  
188 ized problem. We use the same random layers ansatz used in section 3.1, and our target distribution  
189 is fully random. Specifically we use the absolute value of a normal distribution with mean 0 and  
190 standard deviation  $\pi$ , and then divide by the sum to normalize it to a valid probability distribution. In  
191 contrast to the other generative modelling experiment, this experiment exists to try to broadly sample  
192 many possible generative models and target distributions. For these experiments we run on 5 qubit  
193 systems, with 100 parameterized single qubit gates and 30 2-qubit gates in the random layers ansatz.  
194 Each run is executed for 5000 optimizer steps.

## 195 4 Optimizers

196 In this section we briefly outline all of the optimizers we benchmarked. While this is certainly not  
197 a fully exhaustive comparison of all local zeroth-order optimizers, we chose this selection because  
198 they cover most methods and heuristics used in SPSA-like methods. Additionally, most of these  
199 optimizers have a history of being used for parameterized quantum circuit tasks. Table 1 in the  
200 appendix shows the detailed update rules of these optimizers.

201 **Simultaneous Perturbation Stochastic Approximation (SPSA)** [Spall, 1992] is a commonly used  
202 method, both inside and outside the context of optimizing quantum circuits. In a nutshell, SPSA is  
203 approximated gradient descent where we randomly sample directions in parameter space. Per step it  
204 samples a small random vector from a Rademacher distribution in parameter space, estimates the  
205 gradient along that vector with finite difference approximation, and then takes a step along said vector  
206 according to the sampled gradient to minimize loss. SPSA is also often used with learning rate and  
207 finite difference step size decay, which we also use here.

208 **AdamSPSA** [Leng et al., 2023] is the application of the Adam optimizer heuristic [Kingma and Ba,  
209 2017] on the SPSA algorithm. Specifically, it estimates via a running sum and updates according to  
210 momentum and variance normalization terms.

211 **2-SPSA** [Spall, 1997] is essentially an approximation of Newton's method, which is gradient descent  
212 where the gradient is multiplied by the inverse of the Hessian. To approximate the Hessian, it samples  
213 two random vectors (with a Rademacher distribution like in regular SPSA) and evaluates the 2nd  
214 order derivative along those two vectors. It then uses a weighted averaging of these samples to  
215 provide the Hessian used during optimization. Additionally, because the Hessian estimate can lead  
216 to more unstable updates, 2-SPSA also often blocks updates which increase the loss over a certain  
217 threshold from the prior value per step.

218 **Quantum Natural SPSA (QNPSA)** [Gacon et al., 2021] is a variation of SPSA that utilizes the  
 219 quantum natural gradient. Similar to classical natural information, using this metric can has a few  
 220 theoretical advantages to help accelerate and stabilize learning. In practice, this method functions  
 221 near-identically to 2-SPSA, except that when they sample the metric matrix, they compute the Hessian  
 222 of the Fubini-Study metric instead of the Hessian of the loss function.

223 **Guided evolutionary strategies (GES)** [Maheswaranathan et al., 2019] is an evolutionary method  
 224 with heuristic guiding. However, despite the different name, it is fundamentally very similar to  
 225 SPSA with only two major differences. First, GES instead samples its random parameter-space  
 226 vectors with a Gaussian instead of a Rademacher distribution. Second, GES biases the covariance of  
 227 the sampling Gaussian along the subspace of the recent prior gradients. The intuition behind this  
 228 choice is that, similar to momentum, that future gradients are more likely that not to be biased in the  
 229 direction of the most recent prior gradients. However, instead of just increasing the update size in  
 230 these directions, GES biases the sampling in this direction to increase information gain in this biased  
 231 direction-of-travel.

232 **Exponential Natural Evolutionary Strategies (xNES)** [Wierstra et al., 2011] is an extension of  
 233 evolutionary strategies to improve trainability. In this algorithm, it is assumed we have a multi-  
 234 variate Gaussian in our space of model parameters, and our goal is to optimize this Gaussian to, in  
 235 expectation, sample parameters that produce the lowest loss on the underlying problem. This is done  
 236 by performing stochastic gradient descent on the parameters of the multi-variate Gaussian. xNES  
 237 then augments this by instead using natural gradient descent to improve convergence guarantees, and  
 238 utilizes an exponential matrix mapping to make the algorithm more computationally efficient.

239 **Seperable Natural Evolutionary Strategies (sNES)** [Wierstra et al., 2011] is simply xNES that  
 240 assumes independence between parameters in order to be even more computationally efficient. It  
 241 is functionally equivalent to xNES except where the covariance matrix only allows terms along the  
 242 diagonal.

## 243 5 Results

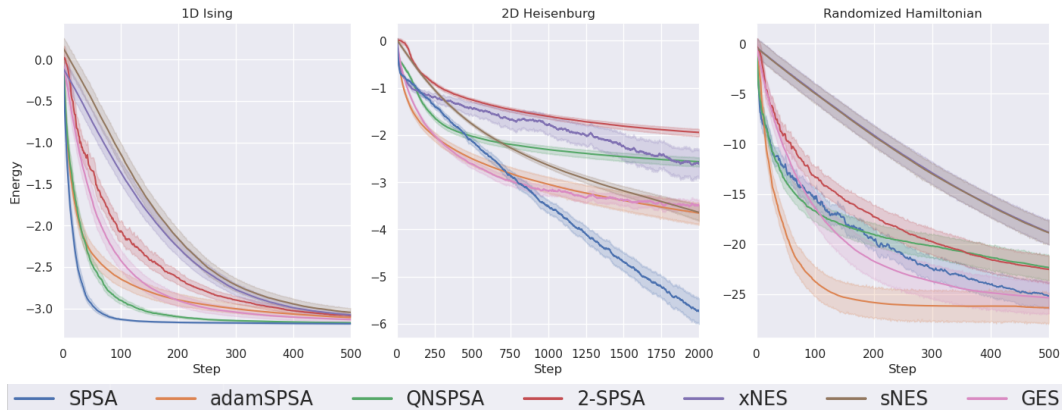


Figure 1: Hamiltonian minimization experiment convergence plots w.r.t. number of optimizer steps. Plots mean value of all runs with the 95% confidence interval. Experiment details can be found in section 3.

244 We produce two types of plots for all experiments in section 3: convergence plots and box plots. The  
 245 convergence plots show the average loss of the optimization during each stage of the process, where  
 246 the colored error area is the 95% confidence interval of the mean.

247 The box plots illustrate the statistics of the end-result of each optimization run. In these plots, the  
 248 center line of the colored region is the median loss value. The box region is the interquartile range  
 249 (the range centered around the median that contains 50% of the samples). The plot whiskers contain

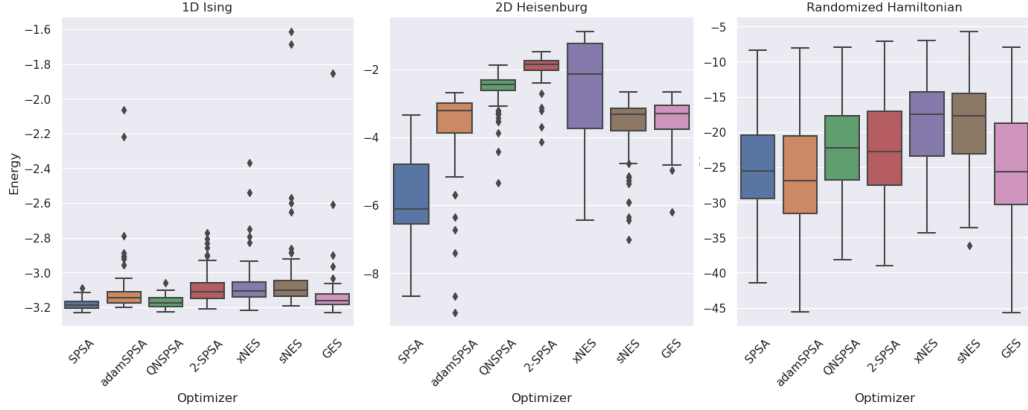


Figure 2: Hamiltonian minimization experiment box plots. Plots the statistics of the final loss value from each run. Experiment details can be found in section 3.

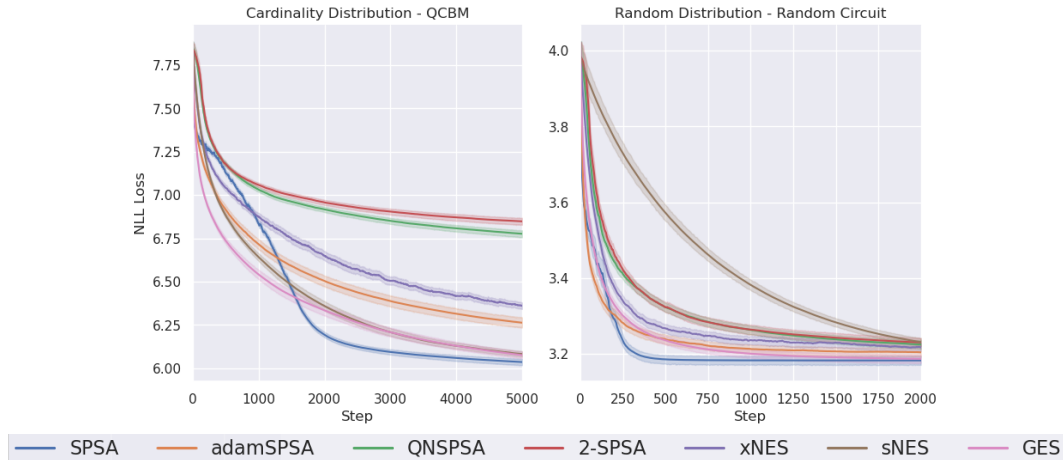


Figure 3: Generative modelling experiments convergence plots w.r.t. number of optimizer steps. Plots mean value of all runs with the 95% confidence interval. Experiment details can be found in section 3.

250 all points that are within 1.5 times the size of the interquartile range from the median. All other points  
 251 are considered outliers and are plotted individually.

252 The convergence plots are figures 1 and 3 for the Hamiltonian minimization and the generative  
 253 modelling experiments respectively. The box plots are figures 2 and 4, likewise for the hamiltonian  
 254 minimization and the generative modelling experiments respectively.

## 255 5.1 Insights

256 There are a few main take-aways from these results that we believe these results illustrate. While we  
 257 don't claim these results are concrete truths, we believe each of them warrant further study.

258 **1) Hyperparameter tuning for optimizer generalization is extremely important.** Generally this is  
 259 illustrated by our results showing that most optimizer variants don't show clear benefits compared to  
 260 their original versions when some hyperparameter tuning is done, which often contrasts the results  
 261 shown in the original papers. However, this is more specifically illustrated in the convergence plot  
 262 results of the 2D Heisenberg model shown in figure 1. For this experiment, we hyperparameter  
 263 tuned to optimizers taking 1000 steps but ran our experimental results out to 2000 steps. If you cut

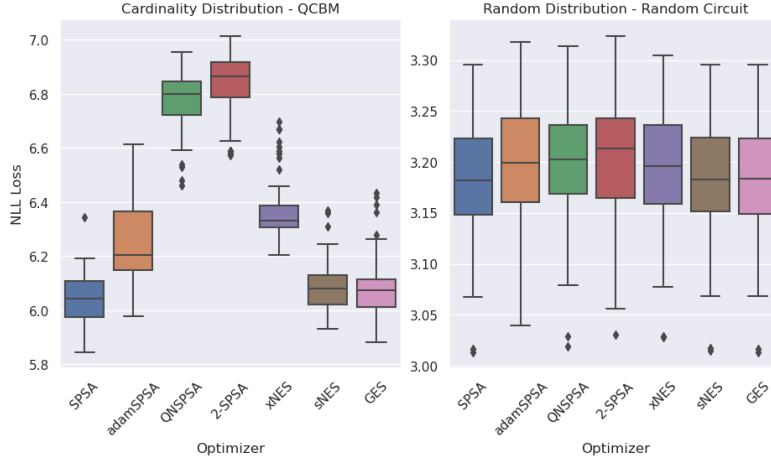


Figure 4: Generative modelling experiments box plots. Plots the statistics of the final loss value from each run. Experiment details can be found in section 3.

264 this plot off at 1000 steps, it would look very close to the random Hamiltonian experiments where  
 265 SPSA, GES, and AdamSPSA perform similarly aside from the latter two converging more quickly.  
 266 However, by choosing to optimize for longer SPSA is suddenly able to do significantly better than all  
 267 optimizers. While this does raise the question of what a truly realistic hyperparameter tuning scenario  
 268 is in quantum systems, we think it’s likely more fruitful to sidestep this concern altogether and work  
 269 to design optimizers that are adaptive or more robust to hyperparameter choice.

270 **2) More elaborate optimization strategies aren’t generally better.** In all of these benchmarks, no  
 271 optimizer clearly out-performs SPSA at the end of the optimization procedure. And the ones that do  
 272 are often ones like AdamSPSA and GES that rely on relatively simple and cheap heuristics. So while  
 273 there is certainly something to be said about the theoretical benefits of using a method like QNSPSA,  
 274 these results indicate that it may be possible to practically achieve better performance for cheaper by  
 275 using simpler methods like step size decay and guiding heuristics.

276 **3) When it comes to the convergence speed in the initial phase, you may benefit from acceleration**  
 277 **strategies.** It varies depending on the benchmark you look at, but often a few methods are able to  
 278 converge to lower loss early on in the optimization procedure before being met or overtaken by SPSA.  
 279 So while it may not be clearly better to use one of these methods as-is, these results indicate it might  
 280 be possible to develop new adaptive methods to improve the convergence speed of methods like  
 281 SPSA while not sacrificing overall performance.

## 282 6 Conclusion

283 In this work, we benchmark SPSA-like optimizers on a variety of parameterized quantum learning  
 284 tasks with randomized quantum circuits and randomized objectives. These results provide evidence  
 285 to suggest that certain heuristics can help accelerate optimization, they often do not perform better  
 286 than the simpler methods in general. However, we believe that in a broader sense this study helps  
 287 illustrate the need not only for methods that are adaptive / more robust to hyperparameter choice,  
 288 but also for broader thought on how we can effectively compare optimizers in quantum systems that  
 289 aligns with the realistic scenarios they will be used in. While the take-aways from this work are  
 290 intuitions and ultimately only serve to inform directions of future study, our hope is that this work  
 291 will inspire more thought into how to best categorize and compare optimizers in quantum learning  
 292 problems as a whole. (We also include some suggested areas of future studies in section A of the  
 293 appendix.)



## 294 References

- 295 Abhinav Anand, Matthias Degroote, and Alán Aspuru-Guzik. Natural evolutionary strategies for  
296 variational quantum computation. *Machine Learning: Science and Technology*, 2(4):045012,  
297 July 2021. ISSN 2632-2153. doi: 10.1088/2632-2153/abf3ac. URL [https://dx.doi.org/10.](https://dx.doi.org/10.1088/2632-2153/abf3ac)  
298 [1088/2632-2153/abf3ac](https://dx.doi.org/10.1088/2632-2153/abf3ac). Publisher: IOP Publishing.
- 299 Marcello Benedetti, Delfina Garcia-Pintos, Oscar Perdomo, Vicente Leyton-Ortega, Yunseong Nam,  
300 and Alejandro Perdomo-Ortiz. A generative modeling approach for benchmarking and training  
301 shallow quantum circuits. *npj Quantum Information*, 5(1):1–9, May 2019. ISSN 2056-6387. doi: 10.  
302 [1038/s41534-019-0157-8](https://www.nature.com/articles/s41534-019-0157-8). URL <https://www.nature.com/articles/s41534-019-0157-8>.  
303 Number: 1 Publisher: Nature Publishing Group.
- 304 Ville Bergholm, Josh Izaac, Maria Schuld, Christian Gogolin, Shah Nawaz Ahmed, Vishnu Ajith,  
305 M. Sohaib Alam, Guillermo Alonso-Linaje, B. Akash Narayanan, Ali Asadi, Juan Miguel Arrazola,  
306 Utkarsh Azad, Sam Banning, Carsten Blank, Thomas R. Bromley, Benjamin A. Cordier, Jack  
307 Ceroni, Alain Delgado, Olivia Di Matteo, Amintor Dusko, Tanya Garg, Diego Guala, Anthony  
308 Hayes, Ryan Hill, Aroosa Ijaz, Theodor Isaacsson, David Ittah, Soran Jahangiri, Prateek Jain,  
309 Edward Jiang, Ankit Khandelwal, Korbinian Kottmann, Robert A. Lang, Christina Lee, Thomas  
310 Loke, Angus Lowe, Keri McKiernan, Johannes Jakob Meyer, J. A. Montañez-Barrera, Romain  
311 Moyard, Zeyue Niu, Lee James O’Riordan, Steven Oud, Ashish Panigrahi, Chae-Yeun Park,  
312 Daniel Polatajko, Nicolás Quesada, Chase Roberts, Nahum Sá, Isidor Schoch, Borun Shi, Shuli  
313 Shu, Sukin Sim, Arshpreet Singh, Ingrid Strandberg, Jay Soni, Antal Száva, Slimane Thabet,  
314 Rodrigo A. Vargas-Hernández, Trevor Vincent, Nicola Vitucci, Maurice Weber, David Wierichs,  
315 Roeland Wiersema, Moritz Willmann, Vincent Wong, Shaoming Zhang, and Nathan Killoran.  
316 PennyLane: Automatic differentiation of hybrid quantum-classical computations, July 2022. URL  
317 <http://arxiv.org/abs/1811.04968>. arXiv:1811.04968 [physics, physics:quant-ph].
- 318 Robin Blume-Kohout, Kenneth Rudinger, Erik Nielsen, Timothy Proctor, and Kevin Young. Wildcard  
319 error: Quantifying unmodeled errors in quantum processors, December 2020. URL [http://](http://arxiv.org/abs/2012.12231)  
320 [arxiv.org/abs/2012.12231](http://arxiv.org/abs/2012.12231). arXiv:2012.12231 [quant-ph].
- 321 M. Cerezo, Andrew Arrasmith, Ryan Babbush, Simon C. Benjamin, Suguru Endo, Keisuke Fujii,  
322 Jarrod R. McClean, Kosuke Mitarai, Xiao Yuan, Lukasz Cincio, and Patrick J. Coles. Variational  
323 Quantum Algorithms. *Nature Reviews Physics*, 3(9):625–644, August 2021. ISSN 2522-5820. doi:  
324 [10.1038/s42254-021-00348-9](https://doi.org/10.1038/s42254-021-00348-9). URL <http://arxiv.org/abs/2012.09265>. arXiv:2012.09265  
325 [quant-ph, stat].
- 326 Luuk Coopmans, Di Luo, Graham Kells, Bryan K. Clark, and Juan Carrasquilla. Protocol Dis-  
327 covery for the Quantum Control of Majoranas by Differentiable Programming and Natural Evo-  
328 lution Strategies. *PRX Quantum*, 2(2):020332, June 2021. ISSN 2691-3399. doi: 10.1103/  
329 [PRXQuantum.2.020332](https://doi.org/10.1103/PRXQuantum.2.020332). URL <http://arxiv.org/abs/2008.09128>. arXiv:2008.09128 [cond-  
330 mat, physics:physics, physics:quant-ph].
- 331 Sepehr Ebadi, Alexander Keesling, Madelyn Cain, Tout T. Wang, Harry Levine, Dolev Bluvstein,  
332 Giulia Semeghini, Ahmed Omran, Jinguo Liu, Rhine Samajdar, Xiu-Zhe Luo, Beatrice Nash,  
333 Xun Gao, Boaz Barak, Edward Farhi, Subir Sachdev, Nathan Gemelke, Leo Zhou, Soonwon  
334 Choi, Hannes Pichler, Shengtao Wang, Markus Greiner, Vladan Vuletic, and Mikhail D. Lukin.  
335 Quantum Optimization of Maximum Independent Set using Rydberg Atom Arrays. *Science*, 376  
336 (6598):1209–1215, June 2022. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.abo6587.  
337 URL <http://arxiv.org/abs/2202.09372>. arXiv:2202.09372 [cond-mat, physics:physics,  
338 physics:quant-ph].
- 339 Enrico Fontana, Manuel S. Rudolph, Ross Duncan, Ivan Rungger, and Cristina Cîrstoiu. Classical  
340 simulations of noisy variational quantum circuits, June 2023. URL [http://arxiv.org/abs/](http://arxiv.org/abs/2306.05400)  
341 [2306.05400](http://arxiv.org/abs/2306.05400). arXiv:2306.05400 [quant-ph].

- 342 Julien Gacon, Christa Zoufal, Giuseppe Carleo, and Stefan Woerner. Simultaneous Perturbation  
343 Stochastic Approximation of the Quantum Fisher Information. *Quantum*, 5:567, October 2021.  
344 ISSN 2521-327X. doi: 10.22331/q-2021-10-20-567. URL <http://arxiv.org/abs/2103.09232>. arXiv:2103.09232 [quant-ph].  
345
- 346 Kaitlin Gili, Mohamed Hibat-Allah, Marta Mauri, Chris Ballance, and Alejandro Perdomo-Ortiz.  
347 Do quantum circuit Born machines generalize? *Quantum Science and Technology*, 8(3):035021,  
348 May 2023. ISSN 2058-9565. doi: 10.1088/2058-9565/acd578. URL <https://dx.doi.org/10.1088/2058-9565/acd578>. Publisher: IOP Publishing.  
349
- 350 Andi Gu, Angus Lowe, Pavel A. Dub, Patrick J. Coles, and Andrew Arrasmith. Adaptive shot  
351 allocation for fast convergence in variational quantum algorithms. 2021. doi: 10.48550/arXiv.  
352 2108.10434. URL <http://arxiv.org/abs/2108.10434>. arXiv:2108.10434 [quant-ph].
- 353 Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization, January 2017.  
354 URL <http://arxiv.org/abs/1412.6980>. arXiv:1412.6980 [cs].
- 355 Zhaoqi Leng, Pranav Mundada, Saeed Ghadimi, and Andrew Houck. Efficient Algorithms for  
356 High-Dimensional Quantum Optimal Control of a Transmon Qubit. *Physical Review Applied*, 19  
357 (4):044034, April 2023. doi: 10.1103/PhysRevApplied.19.044034. URL <https://link.aps.org/doi/10.1103/PhysRevApplied.19.044034>. Publisher: American Physical Society.  
358
- 359 Di Luo, Jiayu Shen, Rumen Dangovski, and Marin Soljačić. Koopman Operator learning for  
360 Accelerating Quantum Optimization and Machine Learning, November 2022. URL <https://arxiv.org/abs/2211.01365v1>.  
361
- 362 Niru Maheswaranathan, Luke Metz, George Tucker, Dami Choi, and Jascha Sohl-Dickstein. Guided  
363 evolutionary strategies: Augmenting random search with surrogate gradients, June 2019. URL  
364 <http://arxiv.org/abs/1806.10230>. arXiv:1806.10230 [cs, stat].
- 365 Aidan Pellow-Jarman, Ilya Sinayskiy, Anban Pillay, and Francesco Petruccione. A comparison  
366 of various classical optimizers for a variational quantum linear solver. *Quantum Information  
367 Processing*, 20(6):202, June 2021. ISSN 1573-1332. doi: 10.1007/s11128-021-03140-x. URL  
368 <https://doi.org/10.1007/s11128-021-03140-x>.
- 369 Timothy Proctor, Melissa Revella, Erik Nielsen, Kenneth Rudinger, Daniel Lobser, Peter Maunz,  
370 Robin Blume-Kohout, and Kevin Young. Detecting and tracking drift in quantum information  
371 processors. *Nature Communications*, 11(1):5396, October 2020. ISSN 2041-1723. doi: 10.1038/  
372 s41467-020-19074-4. URL <https://www.nature.com/articles/s41467-020-19074-4>.
- 373 Maria Schuld, Ville Bergholm, Christian Gogolin, Josh Izaac, and Nathan Killoran. Evaluating  
374 analytic gradients on quantum hardware. *Physical Review A*, 99(3):032331, March 2019. ISSN  
375 2469-9926, 2469-9934. doi: 10.1103/PhysRevA.99.032331. URL <http://arxiv.org/abs/1811.11184>. arXiv:1811.11184 [quant-ph].  
376
- 377 Harshdeep Singh, Sabyashachi Mishra, and Sonjoy Majumder. Benchmarking of Different Optimizers  
378 in the Variational Quantum Algorithms for Applications in Quantum Chemistry, February 2023.  
379 URL <http://arxiv.org/abs/2208.10285>. arXiv:2208.10285 [quant-ph].
- 380 J.C. Spall. Multivariate stochastic approximation using a simultaneous perturbation gradient approxi-  
381 mation. *IEEE Transactions on Automatic Control*, 37(3):332–341, March 1992. ISSN 1558-2523.  
382 doi: 10.1109/9.119632. Conference Name: IEEE Transactions on Automatic Control.
- 383 J.C. Spall. Accelerated second-order stochastic optimization using only function measurements. In  
384 *Proceedings of the 36th IEEE Conference on Decision and Control*, volume 2, pages 1417–1424  
385 vol.2, December 1997. doi: 10.1109/CDC.1997.657661. ISSN: 0191-2216.

- 386 J.C. Spall. Implementation of the simultaneous perturbation algorithm for stochastic optimization.  
387 *IEEE Transactions on Aerospace and Electronic Systems*, 34(3):817–823, July 1998. ISSN 1557-  
388 9603. doi: 10.1109/7.705889. Conference Name: IEEE Transactions on Aerospace and Electronic  
389 Systems.
- 390 Kevin J. Sung, Jiahao Yao, Matthew P. Harrigan, Nicholas C. Rubin, Zhang Jiang, Lin Lin, Ryan  
391 Babbush, and Jarrod R. McClean. Using models to improve optimizers for variational quantum  
392 algorithms. *Quantum Science and Technology*, 5(4):044008, October 2020. ISSN 2058-9565. doi:  
393 10.1088/2058-9565/abb6d9. URL <http://arxiv.org/abs/2005.11011>. arXiv:2005.11011  
394 [quant-ph].
- 395 Daan Wierstra, Tom Schaul, Tobias Glasmachers, Yi Sun, and Jürgen Schmidhuber. Natural Evolution  
396 Strategies, June 2011. URL <http://arxiv.org/abs/1106.4487>. arXiv:1106.4487 [cs, stat].

397 **A Future Works**

398 Beyond the straightforward expansions on this work (larger studies, more optimizers, more relevant  
399 randomization, etc.) and works related to the insights in section 5.1, we feel the following would be  
400 particularly fruitful follow up studies:

401 **1) Noise** - This benchmark only compared the noiseless setting, but in practice noise is almost always  
402 involved with a quantum computer. Although there are papers that look at shot noise and how to  
403 select the right amount of circuit evaluations to still optimize efficiently [Gu et al., 2021], we would  
404 like to see a study similar to the one we’ve done here that adds noise as an additional dimension to  
405 study where the transitions of an optimizer performing well v.s. poorly occurs.

406 **2) Sampling Amount** - For this study, we restricted ourselves to considering only sampling from  
407 a minimal number of parameter-difference vectors to get the information needed to take an update  
408 step. However, it’s possible we could do better by relaxing this constraint. (It’s especially surprising  
409 methods like GES worked as well as they did in this study, as they were intended to be used with  
410 multiple samples per step.) And on the other side of this, can we re-design some aspects of methods  
411 that typically benefit from additional samples per step to instead work well in a small-sample setting  
412 by aggregating information between steps? Both considerations would be critical for understanding  
413 the truly best strategies in this class of optimizer.

414 **3) Sampling Distribution** - The optimizers we study here use either a Gaussian or a Rademacher  
415 distribution to sample parameter vectors. While they both seem to be able to produce effective opti-  
416 mizers, studying in a more principled way the effects of different choices of probability distributions  
417 would be interesting. This could especially become more relevant if we look at optimizing parameters  
418 that add discrete constraints, as would be the case in some experimental setups, error-corrected  
419 quantum computation, coordinate descent parameter-shift rule based optimizers [Schuld et al., 2019],  
420 or other more exotic forms of parameterization.

421 **4) Adaptive Methods** - Along the lines of what was mentioned in section 5.1, studying methods that  
422 adaptively change during the optimization process could be fruitful for a number of reasons. First, it  
423 could help combine the benefits of multiple strategies. Beyond the optimizers we covered here, there  
424 are works that find other ways to accelerate learning. (For instance, Luo et al. [2022] use machine  
425 learning to predict optimization trajectories, and Fontana et al. [2023] classically models the loss  
426 landscape of certain parameterized circuits.) When it’s possible to combine the information from  
427 quantum computer queries into multiple methods, having a strategy that can learn to rely either more  
428 or less on a specific strategy during parts of the optimization process could allow us to have desirable  
429 properties of multiple methods (speed of convergence, ability to optimize well in difficult landscapes,  
430 flexibility of a method to work without prior assumptions, robustness to distribution shift, etc.) with  
431 a minimal cost-regret overhead. Second, such a method could provide insights into the limits of  
432 each of the above methods. By studying in which parts of an optimization process one method  
433 begins to be unable to optimize as well as another, it could provide insight to researchers looking to  
434 mathematically understand and characterize optimizers and loss landscapes. Lastly, such methods  
435 would make future benchmarking studies much simpler. Instead of having to be concerned about  
436 what reasonable hyperparameter tuning is and expending the resources to perform it, benchmarks  
437 could just compare adaptive versions of the methods in question. (And if said method has a regret  
438 bound, they can have precise confidence in the robustness of their results.)

AdamSPSA	SPSA
$\Delta_i \sim \mathcal{U}(\{-1, 1\}^d)$ $\epsilon_i = \epsilon_0/i^\gamma, \eta_i = \eta_0/(c+i)^\alpha, \beta_i = \beta_0/i^\lambda$ $\nabla_{\text{est}} f(\theta_i) = \frac{f(\theta_i + \epsilon_i \Delta_i) - f(\theta_i - \epsilon_i \Delta_i)}{2\epsilon_i}$ $(m_1 = \nabla_{\text{est}} f(\theta_1), v_1 = (\nabla_{\text{est}} f(\theta_1))^2)$ $m_i = \beta_i m_{i-1} + (1 - \beta_i) \nabla_{\text{est}} f(\theta_i)$ $v_i = \gamma v_{i-1} + (1 - \gamma) (\nabla_{\text{est}} f(\theta_i))^2$ $\theta_{i+1} = \theta_i - \frac{\eta_i}{\sqrt{v_i + \delta}} m_i$	$\Delta_i \sim \mathcal{U}(\{-1, 1\}^d)$ $\epsilon_i = \epsilon_0/i^\gamma, \eta_i = \eta_0/(c+i)^\alpha$ $\nabla_{\text{est}} f(\theta_i) = \frac{f(\theta_i + \epsilon_i \Delta_i) - f(\theta_i - \epsilon_i \Delta_i)}{2\epsilon_i}$ $\theta_{i+1} = \theta_i - \eta_i \nabla_{\text{est}} f(\theta_i)$
2-SPSA	GES
$\Delta_i, \Delta'_i \sim \mathcal{U}(\{-1, 1\}^d)$ $\nabla_{\text{est}} f(\theta_i) = \frac{f(\theta_i + \epsilon_i \Delta_i) - f(\theta_i - \epsilon_i \Delta_i)}{2\epsilon_i}$ $\delta f = f(\theta_i + \epsilon \Delta_i + \epsilon \Delta'_i) - f(\theta_i + \epsilon \Delta_i) - f(\theta_i - \epsilon \Delta_i + \epsilon \Delta'_i) + f(\theta_i - \epsilon \Delta_i)$ $\hat{H}_i = \frac{\delta f}{2\epsilon^2} \frac{\Delta_i (\Delta'_i)^T + \Delta'_i (\Delta_i)^T}{2}$ $H_i = \frac{i}{i+1} H_{i-1} + \frac{1}{i+1} \hat{H}_i$ $\theta_{i+1} = \theta_i - \eta H_i^{-1} \nabla_{\text{est}} f(\theta_i)$	$U_i = \text{orthonormal basis of span of } \{\nabla_{\text{est}} f(\theta_{i-k}), \dots, \nabla_{\text{est}} f(\theta_i)\}$ $\Sigma_i = \frac{\alpha}{n} I + \frac{1-\alpha}{k} U U^T, \Sigma_{0, \dots, k} = \frac{1}{n} I$ $\Delta_i \sim \mathcal{N}(0, \sigma^2 \Sigma_i)$ $\nabla_{\text{est}} f(\theta_i) = \beta \frac{f(\theta_i + \Delta_i) - f(\theta_i - \Delta_i)}{2\sigma^2}$ $\theta_{i+1} = \theta_i - \eta \nabla_{\text{est}} f(\theta_i)$
QNPSA	xNES
$\Delta_i, \Delta'_i \sim \mathcal{U}(\{-1, 1\}^d)$ $F(\theta, \theta') =  \langle \psi(\theta)   \psi(\theta') \rangle ^2$ $\nabla_{\text{est}} f(\theta_i) = \frac{f(\theta_i + \epsilon_i \Delta_i) - f(\theta_i - \epsilon_i \Delta_i)}{2\epsilon_i}$ $\delta F = F(\theta_i, \theta_i + \epsilon \Delta_i + \epsilon \Delta'_i) - F(\theta_i, \theta_i + \epsilon \Delta_i) - F(\theta_i, \theta_i - \epsilon \Delta_i + \epsilon \Delta'_i) + F(\theta_i, \theta_i - \epsilon \Delta_i)$ $\hat{H}_i = \frac{-\delta F}{4\epsilon^2} \frac{\Delta_i (\Delta'_i)^T + \Delta'_i (\Delta_i)^T}{2}$ $H_i = \frac{i}{i+1} H_{i-1} + \frac{1}{i+1} \hat{H}_i$ $\theta_{i+1} = \theta_i - \eta H_i^{-1} \nabla_{\text{est}} f(\theta_i)$	$\Delta_i, \Delta'_i \sim \mathcal{N}(0, I), B_1 = I$ $z = \theta_i + \sigma_i B_i^T \Delta_i, z' = \theta_i + \sigma_i B_i^T \Delta'_i$ $u = 0.5 \text{ if } f(z) < f(z') \text{ else } -0.5$ $u' = 0.5 \text{ if } f(z') < f(z) \text{ else } -0.5$ $\nabla_\mu J = us + u' s'$ $\nabla_M J = u(ss^T - I) + u'(s' s'^T - I)$ $\nabla_\sigma J = \text{tr}(\nabla_M J) / d$ $\nabla_B J = \nabla_M J - \nabla_\sigma J \cdot I$ $\sigma_{i+1} = \sigma_i \exp(\eta_\sigma / 2 \cdot \nabla_\sigma J)$ $B_{i+1} = B_i \exp(\eta_B / 2 \cdot \nabla_B J)$ $\theta_{i+1} = \theta_i + \eta_\mu \sigma_i B_i \nabla_\mu J$
	sNES
	$\Delta_i, \Delta'_i \sim \mathcal{N}(0, I)$ $z = \theta_i + \sigma_i \Delta_i, z' = \theta_i + \sigma_i \Delta'_i$ $u = 0.5 \text{ if } f(z) < f(z') \text{ else } -0.5$ $u' = 0.5 \text{ if } f(z') < f(z) \text{ else } -0.5$ $\nabla_\mu J = us + u' s'$ $\nabla_\sigma J = u(s^2 - 1) + u'(s'^2 - 1)$ $\sigma_{i+1} = \sigma_i \exp(\eta_\sigma / 2 \cdot \nabla_\sigma J)$ $\theta_{i+1} = \theta_i + \eta_\mu \sigma_i B_i \nabla_\mu J$

Table 1: Optimizer algorithms. Illustrates a step of each optimizer in equation form. Note that these algorithms may slightly differ from the original works due to simplifying choices / constraints we made in this study.  $f$  is the loss function and  $\theta$  are the parameters per step.  $\psi$  is the parameterized quantum circuit model used by  $f$ .  $\mathcal{U}$  is the uniform distribution and  $\mathcal{N}$  is the Gaussian distribution.  $(\cdot)^2, \sqrt{(\cdot)}$  are element-wise on vectors. All other un-defined variables are hyperparameters.