# OPTIMISTIC GRADIENT LEARNING WITH HESSIAN CORRECTIONS FOR HIGH-DIMENSIONAL BLACK-BOX OPTIMIZATION

Anonymous authors

Paper under double-blind review

### Abstract

Black-box algorithms are designed to optimize functions without relying on their underlying analytical structure or gradient information, making them essential when gradients are inaccessible or difficult to compute. Traditional methods for solving black-box optimization (BBO) problems predominantly rely on non-parametric models and struggle to scale to large input spaces. Conversely, parametric methods that model the function with neural estimators and obtain gradient signals via backpropagation may suffer from significant gradient errors. A recent alternative, Explicit Gradient Learning (EGL), which directly learns the gradient using a firstorder Taylor approximation, has demonstrated superior performance over both parametric and non-parametric methods. In this work, we propose two novel gradient learning variants to address the robustness challenges posed by high-dimensional, complex, and highly non-linear problems. Optimistic Gradient Learning (OGL) introduces a bias toward lower regions in the function landscape, while Higher-order Gradient Learning (HGL) incorporates second-order Taylor corrections to improve gradient accuracy. We combine these approaches into the unified OHGL algorithm, achieving state-of-the-art (SOTA) performance on the synthetic COCO suite.

Additionally, we demonstrate OHGL's applicability to high-dimensional real-world machine learning (ML) tasks such as adversarial training and code generation. Our results highlight OHGL's ability to generate stronger candidates, offering a valuable tool for ML researchers and practitioners tackling high-dimensional, non-linear optimization challenges.

034 035

037

038

005

008 009 010

011 012 013

014

015

016

017

018

019

021

022

025

026

027

028

029

031

032

033

## 1 INTRODUCTION

Black-box optimization (BBO) is the process of searching for optimal solutions within a system's input domain without access to its internal structure or analytical properties Audet et al. (2017c). Unlike gradient-based optimization methods that rely on the calculation of analytical gradients, BBO algorithms query the system solely through input-output pairs, operating agnostically to the underlying function. This feature distinguishes BBO from traditional ML tasks, such as neural network training, where optimization typically involves backpropagation-based gradient computation.

Many real-world physical systems naturally fit into the BBO framework because their analytical behavior is difficult or impossible to model explicitly. In these cases, BBO algorithms have achieved remarkable success in diverse fields, such as ambulance deployment Zhen et al. (2014), robotic motor control Gehring et al. (2014); Prabhu et al. (2018), parameter tuning Olof (2018); Rimon et al. (2024), and signal processing Liu et al. (2020), among others Alarie et al. (2021). BBO applications extend beyond physical systems; many ML problems exhibit a black-box nature when the true gradient is absent. Examples include hyperparameter tuning Bischl et al. (2023), contextual bandit problems Bouneffouf et al. (2020), and large language model training with human feedback Bai et al. (2022), to name a few.

054 As the scale of data continues to grow, the dimensionality of the problem space increases in 055 tandem. This trend is particularly evident in ML, where model architectures, embedding and 056 latent representation sizes, and the number of hyperparameters are continually expanding. 057 High-dimensional optimization challenges traditional BBO algorithms, which often require 058 the number of collected samples at each step,  $n_s$ , to scale proportionally with the problem size, N. Parametric neural models, however, have shown that reusing past samples can 059 effectively reduce the required sample size such that  $n_s \ll N$  Sarafian et al. (2020); Lu et al. 060 (2023). Building on this, we propose a sampling profiler that further reduces the number of 061 samples needed at each step while maintaining performance. 062

063 While problem dimensions increase, the cost of evaluating intermediate solutions remains a 064 critical constraint, especially in real-world settings where interaction with the environment is expensive or in ML tasks where larger models counterbalance gains in computational power. 065 Therefore, modern BBO algorithms must not only reduce evaluation steps but also converge 066 more quickly Hansen et al. (2010). Achieving this requires algorithms capable of more 067 accurately predicting optimization directions, either through better gradient approximation 068 Anil et al. (2020); Lesage-Landry et al. (2020) or momentum-based strategies to handle non-069 convexity and noise. In this paper, we propose two key improvements to Explicit Gradient 070 Learning (EGL): (1) Optimistic Gradient Learning (OGL), a weighted gradient estimator 071 that biases toward promising solutions and (2) Higher-Order Gradient Learning (HGL), 072 which incorporates Hessian corrections to yield more accurate gradient approximations. 073

We combine the strengths of OGL and HGL to a unified algorithm termed OHGL which exhibits four key advantages:

- **Robustness**: OHGL consistently outperforms baseline algorithms across a diverse range of benchmark problems, including synthetic test suites and real-world ML applications. Its ability to handle noisy and non-convex environments.
- **Gradient Precision**: By integrating the second-order information via Hessian corrections, OHGL achieves significantly more accurate gradient approximations than standard EGL.
- 082

076

077

078

079

081

085

086

- Convergence Rate: OHGL demonstrates faster convergence rate.
- Utilizing the sampling profiler and the optimistic approach, OGL is able to solve **high-dimensional problems** with smaller budget and **converge faster** than baseline algorithms.

Related works: Black-box optimization (BBO) algorithms have a long history, with various approaches developed over the years. Some of the foundational techniques include grid search, coordinate search Audet et al. (2017e), simulated annealing Busetti (2003), and direct search methods like Generalized Pattern Search and Mesh Adaptive direct search Audet et al. (2017a), Gradient-less descent Golovin et al. (2019), and ZOO Chen et al. (2017). These approaches iteratively evaluate potential solutions and decide whether to continue in the same direction. However, they resample for every step and don't use the sampled budget from previous iterations, wasting a lot of budget.

095 Another prominent family of BBO algorithms is the genetic algorithm family Back (1996). 096 This includes methods such as Covariance Matrix Adaptation (CMA) Hansen (2016) and Particle Swarm Optimization (PSO) Clerc (2010). These algorithms simulate the process of 098 natural evolution, where a population of solutions evolves through mutation and selection Audet et al. (2017b). They are considered state-of-the-art (SOTA) in optimization due 099 to their effectiveness in tackling complex problems. However, they come with significant 100 drawbacks, particularly the need for extensive fine-tuning of parameters like generation 101 size and mutation rates. CMA, for example, struggles in higher-dimensional environments 102 and requires careful adjustment of hyperparameters and guidance to perform optimally 103 Loshchilov et al. (2013); Tang (2021). In this work, we propose a simpler method to enhance 104 the performance of CMA, particularly in high-dimensional settings. 105

Then there are model-based methods Audet et al. (2017d), which attempt to emulate the behavior of the function using a surrogate model. These models provide important analytical information, such as gradients Bertsekas (2015), to guide the optimization process and help

108 find a minimum. Within this class, we can further distinguish two sub-classes. To address 109 the issue of dimensionality, Explicit Gradient Learning (EGL) was proposed by Sarafian 110 et al. (2020). While many model-based methods focus on learning the function's structure to 111 derive analytical insights (e.g., Indirect Gradient Learning or IGL Lillicrap (2015); Sarafian 112 et al. (2020)), EGL directly learns the gradient information. EGL uses Taylor's theorem to estimate the gradient. The authors also emphasize the importance of utilizing a trust region 113 to handle black-box optimization problems. However, EGL has some drawbacks: it often 114 uses the available budget inefficiently, disregarding both the complexity and dimensionality 115 of the environment. Additionally, the datasets created by EGL can be naive, leading to 116 over-fitting or improper network learning. This work tackles these issues by showing the 117 importance of proper algorithm calibration and optimization. 118

Recent work also highlights the limitations of common assumptions in BBO algorithms, such as continuity or Gaussian distributions of functions, which can hinder optimization.
For instance, OPT-GAN Tang (2021), a generative model, seeks to bypass these assumptions by learning a function's distribution and generating better candidate solutions based on that knowledge.

The paper is organized as follows: Section 2 covers the algorithm's theoretical background and mathematical foundations. Sections 3 and 4 present our two enhanced variants
of the gradient learning algorithm: OGL and HGL, these are followed by section 5 where we
present the full algorithm OHGL. Section 6 provides experimental results on the synthetic
COCO test suite and Section 7 highlights 2 real-world high-dimensional applications and
potential uses. Finally, section 8 concludes and suggests future research directions. Our
code, experiments, and environment setup are available in the supplementary material.

131 132

133

136

137

143 144

# 2 BACKGROUND

The goal of black-box optimization (BBO) is to minimize a target function f(x) through a series of evaluations Audet et al. (2017c), over a predefined domain  $\Omega$ :

find: 
$$x^* = \arg\min_{x \in \Omega} f(x)$$
 (1)

The Explicit Gradient Learning method, as proposed by Sarafian et al. (2020) leverages the first-order Taylor's expansion:  $f(y) = f(x) + \nabla f(x)^{\top}(y-x) + R_1(x,y)$ . Here,  $R_1(x,y) = O(||y-x||^2)$  is a higher-order residual. By minimizing the residual term with a surrogate neural network model, EGL learns the *mean-gradient*: a smooth approximation of the function's gradient

$$g_{\varepsilon}^{EGL}(x) = \underset{g_{\theta}:\mathbb{R}^n \to \mathbb{R}^n}{\operatorname{arg\,min}} \int_{y \in B_{\varepsilon}(x)} \left( f(x) - f(y) + g_{\theta}(x)^{\top} (y - x) \right)^2 dy \tag{2}$$

145 146 where  $B_{\varepsilon}(x)$  is a ball around x. As  $\varepsilon \to 0$ , 147 the mean-gradient converges to  $\nabla f$ , this 148 property lets EGL explore for lower regions 149 in the function landscape when  $\varepsilon$  is suffi-149 ciently large and converge to a local mini-150 mum when  $\varepsilon$  converges to 0.

151 A key component of the EGL algorithm 152 is the *trust region* (TR), which restricts 153 the search space around the current es-154 This region standardizes inputtimate. 155 output statistics, enhancing the neural net-156 work's effectiveness. Although Sarafian 157 et al. (2020) suggested the TR framework as 158 part of the EGL algorithm, TR is not exclu-159 sive to EGL and can significantly improve other algorithms. In our work, we created 160 stronger baseline algorithms by adding TR 161



Figure 1: Comparing the effect of trust region. Trust region significantly improves both EGL and CMA algorithms.

to the classic Covariance Matrix Adaptation (CMA) algorithm (see Appendix: Algorithm

11). Notably, this seemingly minor change to CMA significantly improves its performance and outperforms vanilla EGL in the COCO optimization suite, as shown in Fig. 1.

#### 3 **OPTIMISTIC GRADIENT LEARNING WITH WEIGHTED GRADIENTS**

While local search algorithms like EGL are based on the notion that the gradient descent 168 path is the optimal search path, following the true gradient path may not be the optimal 169 approach in terms of sampling budget utilization, avoiding shallow local minima and nav-170 igating through narrow ravines. Specifically, in case we obtain a batch D of sample pairs  $\{(x_i, f(x_i))\}_{i \in D}$  around our current solution x, a plausible and optimistic heuristic can be 172 to direct the search path towards low regions in the sampled function landscape regardless 173 the local curvature around x. In other words, to take a sensible guess and bias the search 174 path towards the lower  $(x_i, f(x_i))$  samples. To that end, we define the Optimistic Gradient Learning objective by adding an *importance sampling* weight to the integral of Eq. 2 175

$$g_{\varepsilon}^{OGL}(x) = \underset{g_{\theta}:\mathbb{R}^{n}\to\mathbb{R}^{n}}{\operatorname{arg\,min}} \int_{y\in B_{\varepsilon}(x)} W_{f}(x,y) \cdot \left(f(x) - f(y) + g_{\theta}(x)^{\top}(y-x)\right)^{2} dy \tag{3}$$

Here,  $W_f$  is a softmax-like weight function that normalizes the weight by the sum of exponents across the sampled batch

$$W_f(x,y) = \frac{e^{-\min(f(x),f(y))}}{\sum_{x_i \in D} e^{-f(x_i)}}$$
(4)

184 Notice that in practice, the theoretical objective in Eq. 3 is replaced by a sampled Monte-185 Carlo version (see Sec. 5) s.t. the sum of all weights across the sampled batch is smaller than 1. In the following theorem, we show that the *controllable accuracy* property of EGL, 186 which implies that the mean-gradient converges to the true gradient still holds for our biased 187 version, s.t. when  $\varepsilon \to 0$ ,  $g_{\varepsilon}^{OGL} \to \nabla f(x)$  this guarantees that the convergence properties 188 of the mean-gradient still hold 189

**Theorem 1** (Optimistic Controllable Accuracy) For any differentiable function f with a continuous gradient, there exists  $\kappa^{OGL} > 0$  such that for any  $\varepsilon > 0$ ,  $q_{\varepsilon}^{OGL}(x)$  satisfies

192 193 194

195

196

197

199

190

191

163

164 165

166 167

171

180 181

182 183

 $\|g_{\varepsilon}^{OGL}(x) - \nabla f(x)\| \le \kappa^{OGL} \varepsilon \quad for \ all \ x \in \Omega.$ 

In Fig. 2a we plot 4 typical trajectories of EGL and OGL which demonstrate why in practice and on average we can benefit from biasing the gradient. We find that the biased version avoids getting trapped in local minima and avoids long travels through narrow ravines which rapidly consume the sampling budget. In Fig. 2b we find statistically that OGL tends to progress faster and closer to the global minimum while EGL diverges from the global minimum in favor of local minima.

# 200 201 202

203

204

205

206 207 208

#### 4 GRADIENT LEARNING WITH HESSIAN CORRECTIONS

To learn the mean-gradient, EGL minimizes the first-order Taylor residual (see Sec. 2). Utilizing higher-order approximations has the potential of learning more accurate models for the mean-gradient. Specifically, the second-order Taylor expansion is

$$f(y) = f(x) + \nabla f(x)^{\top} (y - x) + \frac{1}{2} (y - x)^{\top} \nabla^2 f(x) (y - x) + R_2(x, y)$$
(5)

209 Here  $R_2(x,y) = O(||x-y||^3)$  is the second order residual. Like in EGL, we replace  $\nabla f$ 210 with a surrogate model  $g_{\theta}$  and minimize the surrogate residual to obtain our Higher-order 211 Gradient Learning (HGL) variant 212

213  
214  
215  

$$g_{\varepsilon}^{HGL}(x) = \underset{g_{\theta}:\mathbb{R}^{n}\to\mathbb{R}^{n}}{\arg\min} \int_{y\in B_{\varepsilon}(x)} \mathcal{R}_{HGL}^{2}(x,y) dy$$
(6)

$$\mathcal{R}_{g_{\theta}}^{HGL}(x,y) = f(x) - f(y) + g_{\theta}(x)^{\top}(y-x) + \frac{1}{2}(y-x)^{\top}J_{g_{\theta}}(x)(y-x)$$

243 244

257

258

259

260

261 262

263 264

265

266



(a) OGL vs EGL trajectories. 1st row: Gallagher's Gaussian 101-me. 2nd row: 21-hi.



(b) The average Euclidean distance at each step between the algorithm and the global minimum.



The new higher-order term  $J_{g_{\theta}}(x)$  is the Jacobean of  $g_{\theta}(x)$ , evaluated at x which approximates the function's Hessian matrix in the vicinity of our current solution, i.e.,  $J_{q_{\theta}}(x) \approx \nabla^2 f(x)$ . Next, we show theoretically that as expected, HGL converges faster to the true gradient which amounts to lower gradient error in practice.<sup>1</sup>

**Theorem 2** (Improved Controllable Accuracy): For any twice differentiable function  $f \in$  $\mathcal{C}^2$ , there exists  $\kappa_{HGL} > 0$  such that for any  $\varepsilon > 0$ , the second-order mean-gradient  $g_{\varepsilon}^{HGL}(x)$ satisfies

$$\|g_{\varepsilon}^{HGL}(x) - \nabla f(x)\| \le \kappa_{HGL} \varepsilon^2 \quad \text{for all } x \in \Omega.$$

245 In other words, in HGL the model error is in an order of magnitude of  $\varepsilon^2$  instead of  $\varepsilon$  in EGL and OGL. 246 In practice, we verified that this property of HGL 247 translates to more accurate gradient models. Figure 248 3 shows the gradient error of EGL and HGL for a 249 selected set of problems from the COCO test suite 250 where the analytical true gradient can be easily cal-251 culated and compared to our learned model. 252

Learning the gradient with the Jacobian corrections 253 introduces a computational challenge as double back-254 propagation can be expensive. This overhead can 255 hinder the scalability and practical application of the 256 method. A swift remedy is to detach the Jacobian matrix from the competition graph. While this step



Figure 3: Comparing the normalized MSE between the true gradient and EGL and HGL gradient models

slightly changes the objective's gradient (i.e. the gradient trough  $(R_{HGL})$  it removes the second-order derivative and in practice, we found that it achieves almost similar results compared to the full backpropagation through the residual  $R_{HGL}$ , see Fig. 4(b).

#### 5OHGL

In this section, we combine the optimistic approach from Sec. 3 and the Hessian corrections from Sec. 4. However, unlike previous sections, we will present the practical implementation

<sup>267</sup> <sup>1</sup>Notice that, while HGL incorporates Hessian corrections during the gradient learning phase, 268 we deliberately avoid modifying the gradient descent step to include inverse Hessian scaling, as is done in Newton's methods. In practice, inverting the approximated Hessian matrix can lead to 269 numerical challenges and instabilities and was found to be less effective in our experiments.

270 where integrals are replaced with Monte-Carlo sums of sampled pairs. In this case, the loss 271 function which is used to obtain the Optimistic Higher-order Gradient (OHGL) model is 272

274 275

276 277

281

$$\mathcal{L}_{\varepsilon}^{OHGL}(\theta) = \sum_{\|x_i - x_j\| \le \varepsilon} W_f(x_i, x_j) \times \left( f(x_i) - f(x_j) + g_{\theta}(x_i)^{\top} (x_j - x_i) + \frac{1}{2} (x_j - x_i)^{\top} J_{g_{\theta}}(x_i) (x_j - x_i) \right)^2$$
(7)

The summation is applied over sampled pairs which satisfy  $||x_i - x_j|| \leq \varepsilon$ . As explained in 278 Sec. 4, we detach the Jacobian of  $g_{\theta}$  from the computational graph to avoid second-order 279 derivatives. Our final algorithm is outlined in Alg. 1 and an extended version including additional technical details is found in the appendix (See Alg. 2).

Algorithm 1 OHGL (Optimistic Higher-Order Gradient Learning) 283 **Require:**  $x_0, \Omega, \alpha, \epsilon_0, \gamma_\alpha, \gamma_\epsilon, n_{\max}, \lambda$ , Budget 284  $k = 0, j = 0, \Omega_i \leftarrow \Omega$ 285 while Budget > 0 do 286 **Explore:** Generate and evaluate samples  $\mathcal{D}_k = \{\tilde{x}_i, y_i\}_{i=1}^m$ 287 **Create Dataset:** Select *m* tuples  $\mathcal{T}$  from  $\mathcal{D}_k$ 288 Weighted Output: Assign weights  $w_i$  and apply squashing function  $\tilde{y}_i = r_k(w_i y_i)$ 289 **Higher-Order Gradient Learning:** Minimize the loss of  $\theta_k$  (Eq. 7) **Gradient Descent:** Update solution:  $x_{k+1} \leftarrow x_k - \alpha \tilde{g}_{\theta_k}(x_k)$ 291 if  $f(h_j^{-1}(\tilde{x}_{k+1})) > f(h_j^{-1}(\tilde{x}_k))$  for  $n_{\max}$  times then Generate new trust region  $\Omega_{j+1}$  and update  $\epsilon_j$ 292 293 end if if  $f(h_j^{-1}(\tilde{x}_k)) < f(h_j^{-1}(\tilde{x}_{\text{best}}))$  then Update  $x_{\text{best}}$ 295 296 end if  $k \leftarrow k+1$ ; Budget = Budget - m 297 end while 298 return  $x_{\text{best}}$ 299

300 301

302 303

304

305

306

307

308

309

310

313

314

315

316

#### EXPERIMENTS IN THE COCO TEST SUITE 6

We evaluated the OGL, HGL and OHGL algorithms on the COCO framework Hansen et al. (2021) and compared them to EGL and other strong baselines: (1) CMA and its trust region variants L-CMA (linear trust region mapping function) and T-CMA (tanh trust region mapping function) and (2) Implicit Gradient Learning (IGL) Sarafian et al. (2020) where we follow the EGL protocol but train a model for the objective function and obtain the gradient estimation by backpropagation as in DDPG Lillicrap (2015). We also adjusted EGL hyper-parameters A.4 and improved the trust region A.5 to reduce the budget usage by our algorithms.

- 311 We use the following evaluation metrics: 312
  - Convergence Rate: Speed of reaching the global optimum.
  - **Success Rate**: Percentage of problems solved within a fixed budget.
  - **Robustness**: Performance stability across different hyperparameter settings.

317 Performance was normalized against the best-known solutions to minimize bias: normalized\_value =  $\frac{y - y_{min}}{y_{max} - y_{min}}$ . A function was considered solved if the normalized value 318 319 was below 0.01.

320 321

322

SUCCESS AND CONVERGENCE RATE 6.1

Figure 4(c) illustrates the success rate of each algorithm relative to the distance from the 323 best point required for solving a function. The results show that both OGL and OHGL

4	Metric	T-CMA	L-CMA	EGL	OGL	HGL	OHGL
-	Budget to reach 0.01	$17,828 \pm 32648$	$6497 \pm 19731$	24,102	$8981 \pm 9177$	$22,146 \pm 29757$	$26,706 \pm 32676$
0	Mean	$0.01 \pm 0.05$	$0.01 \pm 0.05$	$0.01 \pm 0.03$	$0.003 \pm 0.02$	$0.006 \pm 0.03$	$0.002\pm0.02$
6	Solved Functions	$0.86 {\pm} 0.023$	$0.88 {\pm} 0.023$	$0.83 {\pm} 0.023$	$0.92 {\pm} 0.022$	$0.89 \pm 0.022$	$0.926 {\pm} 0.022$

Table 1: Comparison of different metrics: Budget used to reach 0.99 of the final score  $(\downarrow)$ , the mean normalized results  $(\downarrow)$ , std  $(\downarrow)$ ; and percentage of solved problems  $(\uparrow)$ .



Figure 4: Experiment results against the baseline: (a) Convergence for all our algorithms against baseline algorithms, (b) ablation test for EGL enhancements, (c) Success rate of algorithms as a function of the normalized distance from the best-known solution, (d) Percentage of solved algorithms when the distance from the best point is 0.01

consistently outperform all other algorithms. In particular where the error should be small (<0.01), with t-tests yielding p-values approaching 1, indicating strong statistical confidence (see the complete list of t-test p-values in Table 4 in the appendix).

Figure 4(a) presents the convergence rates for the seven algorithms, where OGL, HGL, and OHGL demonstrate superior convergence compared to the other methods. OHGL, in particular, consistently ranks among the top performers across most metrics, as seen in Table 1, which compares multiple performance indicators. Although OGL achieves better initial results, as it searches longer for the space with the optimal solution, making it a worse fit for problems with a low budget, it ultimately reaches superior results. Additionally, Figure 4(b) shows an ablation test, confirming that the core components of our work: the optimistic approach (OGL) and the Higher-order corrections (HGL) contributed substantially much more to the overall performance than other technical improvements to the algorithm (with respect to vanilla EGL), these technical improvements are denoted as EGL-OPTIMIZED and they amount to better trust-region management and more efficient sampling strategies which mainly helps in faster learning rate at the beginning of the process (as described in Sec. A.4 and A.5 in the appendix).

# 3786.2 Hyperparameter Tolerance

In this part, we evaluate the robustness of our algorithm to hyperparameter tuning. Our objective is to demonstrate that variations in key hyperparameters have minimal impact on the algorithm's overall performance. To assess this, we conducted systematic experiments, modifying several hyperparameters and analyzing their effects. Table 2 (and Table 6 in the appendix) reports the coefficient of variation (CV), defined as  $CV = \frac{\sigma}{\mu}$ , across different hyperparameter sweeps, highlighting the algorithm's stability under varying conditions.

386 Our findings indicate that certain hyperparameters—such as the epsilon factor, shrink fac-387 tor, and value normalizer learning rate (LR)—exhibit cumulative effects during training. 388 While small variations in these parameters may not have an immediate impact, their influence can accumulate over time, potentially leading to significant performance changes. In 389 A.11.2, we establish the relationship between the step size and the epsilon factor necessary 390 for ensuring progress toward a better optimal solution. When selecting these parameters, 391 this relationship must be considered. Additionally, the shrink factor for the trust region 392 should be chosen relative to the budget, enabling the algorithm to explore the maximum 393 number of sub-problems. This underscores the importance of fine-tuning these parameters 394 for optimal results. Conversely, we found that the structural configuration of the neural 395 network, including the number and size of layers, had minimal effect on performance. This suggests that the algorithm's reliance on Taylor loss enables effective learning even with 397 relatively simple network architectures, implying that increasing model complexity does not 398 necessarily yield substantial improvements.

Metric	Networks	Epsilon	Epsilon Factor	Training Bias	Perturbation
CV	0.0167	0.0339	0.2361	0.1292	0.1642
Metric	TR Shrink Method	Normalizer LR	Normalizer Outlier	TR Shrink Factor	Optimizer LR
CV	0.0333	0.1618	0.0333	0.4894	0.5625

Table 2: Coefficient of Variation  $(CV = \frac{\sigma}{\mu})$  over a Hyperparameter sweep experiment.

# 7 HIGH DIMENSIONAL APPLICATIONS

# 7.1 Adversarial Attacks

As powerful vision models like ResNet Targ et al. (2016) and Vision Transformers (ViT)
Han et al. (2022) grow in prominence, adversarial attacks have become a significant concern. These attacks subtly modify inputs, causing models to misclassify them, while the perturbation remains imperceptible to both human vision and other classifiers Tang et al. (2019). Formally, an adversarial attack is defined as:

415 416

404 405 406

407 408

409

 $x_a^* = \arg\min_x d(x, x_a) \quad \text{s.t.} f(x) \neq f(x_a) \tag{8}$ 

417 Where f is the classifier and d is some distance metric between elements.

418 Recent studies have extended adversarial attacks to domains like AI-text detection Sadasivan 419 et al. (2024) and automotive sensors Mahima et al. (2024). These attacks prevent tracking 420 and detection, posing risks to both users and pedestrians. Adversarial attacks are classified 421 into black-box and white-box methods. Black-box attacks only require query access, while 422 white-box methods use model gradients to craft perturbations Machado et al. (2021); Cao 423 et al. (2019). Despite some black-box methods relying on surrogate models Dong et al. (2018); Xiao et al. (2018); Madry et al. (2017); Goodfellow et al. (2014), approaches like 424 AutoZOOM Tu et al. (2019) generate random samples to approximate gradient estimation, 425 though they are computationally expensive. Other methods use GAN networks to search 426 latent spaces for adversarial examples Liu et al. (2021); Sarkar et al. (2017). Still, they 427 depend on existing GANs and their latent space diversity. 428

429 Our Enhanced Gradient Learning method offers a true black-box approach with precise
 430 perturbation control, avoiding gradient back-propagation. OGL directly optimizes perturbations to maintain low distortion while fooling the model, handling high-dimensional spaces with over 30,000 parameters.

441 442

443 444

454

456

478



Figure 5: Adversarial examples generated by OGL and CMA against the ImageNet model.

445 Methodology: We applied OGL to classifiers trained on MNIST, CIFAR-10, and Ima-446 geNet, aiming to minimize equation 8. To generate adversarial images with minimal distortion, we developed a penalty that jointly minimizes MSE and CE-loss A.9. This approach 447 successfully fooled the model, evading the top 5 classifications. 448

449 Results: We evaluated four different configurations: CMA, OGL, HGL, and a combination 450 of both (CMA+OGL) where the CMA run provides the initial guess for an OGL run. While 451 CMA alone was not able to converge to a satisfying adversarial example, the combined 452 CMA+OGL enjoyed the rapid start of CMA with the robustness of OGL s.t. it was able to 453 find a satisfying adversarial example in 50% of the time used by OGL and HGL.

#### 455 7.2CODE GENERATION

The development of large language models (LLMs) such as Transformers Vaswani (2017) 457 have advanced code generation Dehaerne et al. (2022). Despite these strides, fine-tuning 458 outputs based on parameters measured post-generation remains challenging. Recent algo-459 rithms have been developed to generate code tailored for specific tasks using LLMs. For 460 instance, FunSearch Romera-Paredes et al. (2024) generates new code solutions for complex 461 tasks, while Chain of Code Li et al. (2023) incorporates reasoning to detect and correct 462 errors in the output code. Similarly, our method uses black-box optimization to guide code 463 generation for runtime efficiency. Building on Zhang et al. (2024), which links LLM exper-464 tise to a small parameter set, we fine-tuned the embedding layer to reduce Python code 465 runtime. Using LoRA Hu et al. (2021), we optimized the generated code based on execution 466 time, scaling up to  $\sim 200$ k parameters.

467 **Fibonacci:** We tested this approach by having the model generate a Fibonacci function. 468 Initially incorrect, optimization guided the model to a correct and efficient solution. Figure 469 1 illustrates this progression, with the 25th step showing an optimized version. 470

**Line-Level Efficiency Enhancements:** We tested the model's ability to implement small 471 code efficiencies, such as replacing traditional for-loops with list comprehensions. The algo-472 rithm optimized the order of four functions—'initialize', 'start', 'activate', and 'stop'—each 473 with eight variants, minimizing overall runtime by optimizing the function order. 474

**Code Force Challenge:** For a more complex problem, we used the Count Triplets chal-475 lenge from Codeforces<sup>2</sup>. While the model initially struggled, once it found a correct solution, 476 the algorithm further optimized it for runtime performance A.3. 477

<sup>2</sup>https://codeforces.com/

Metric	$\mathbf{CMA}$	OGL	HGL	CMA+OGL
Accuracy	0.02	0.02	0.02	0.02
MSE	0.05	0.001	0.002	0.001
Time (Until Convergence)	20m	6H	7H	3H

Table 3: Comparison of methods on Accuracy, MSE, and Time.



Figure 6: Generated code samples by OGL algorithm

**Discussion:** Our method demonstrates the 500 ability to generate correct solutions while 501 applying micro-optimizations for efficiency. 502 In simple tasks like Fibonacci, OGL con-503 verged on an optimal solution 7, and in 504 more complex problems, it improved the 505 initial solutions. However, in more com-506 plex tasks, the LLM may generate code that 507 fails to solve the problem, hindering the 508 optimization of the run-time. To address this, either stronger models or methods focused on optimizing solution correctness are 510 needed. This would ensure that valid solu-511



Figure 7: Code generation experiments: penalty over time.

tions are generated first, which can then be further optimized for performance.

#### 514 8 CONCLUSION

515

513

497

498 499

In this work, we introduced several key enhancements to the EGL algorithm, resulting
in OHGL, a powerful black-box optimization tool capable of addressing various complex
problems. We demonstrated improvements such as second-order gradient approximation,
optimistic gradients, and trust region enhancements, which collectively led to faster convergence and more robust performance, particularly in high-dimensional and intricate tasks.

Our experiments showcased OHGL's superiority over state-of-the-art methods, especially
 in tasks where computational efficiency and precision are paramount, such as adversarial
 attacks on vision models and optimizing code generation. OHGL's ability to navigate com plex optimization spaces while minimizing computational cost demonstrates its utility across
 various applications.

Looking ahead, Future research should explore new applications of black-box optimization
algorithms like OHGL. While OHGL has shown promise in high-dimensional tasks, its efficiency could be further improved. Our research utilized dimension-reduction methods such
as Hu et al. (2021). Similar techniques (Zhao et al. (2024); Anil et al. (2020)) can help
calculate the Hessian more efficiently.

In conclusion, OHGL represents a significant advancement in BBO algorithm design, providing a robust framework for solving complex optimization problems. However, the true potential of BBO algorithms lies in their broader applicability, and future research should focus on unlocking new avenues for their use, especially in generative models and other cutting-edge areas.

- 536
- 537
- 538
- 539

# 540 REFERENCES

579

585

586

588

593

- 542 Stéphane Alarie, Charles Audet, Aïmen E Gheribi, Michael Kokkolaras, and Sébastien Le Di543 gabel. Two decades of blackbox optimization applications. *EURO Journal on Computa-*544 *tional Optimization*, 9:100011, 2021.
- Rohan Anil, Vineet Gupta, Tomer Koren, Kevin Regan, and Yoram Singer. Scalable second order optimization for deep learning. arXiv preprint arXiv:2002.09018, 2020.
- 548 Charles Audet, Warren Hare, Charles Audet, and Warren Hare. Direct search in derivative549 free and blackbox optimization. *Derivative-Free and Blackbox Optimization*, pp. 93–156, 2017a.
- 551
  552
  553
  554
  Charles Audet, Warren Hare, Charles Audet, and Warren Hare. Genetic methods in derivative-free and blackbox optimization. Derivative-Free and Blackbox Optimization, pp. 57–73, 2017b.
- Charles Audet, Warren Hare, Charles Audet, and Warren Hare. Chapter 1: The begining of dfo algorithms in derivative-free and blackbox optimization. *Derivative-Free and Blackbox Optimization*, pp. 11–15, 2017c.
- Charles Audet, Warren Hare, Charles Audet, and Warren Hare. Model-based methods in derivative-free and blackbox optimization. Derivative-Free and Blackbox Optimization, pp. 156–218, 2017d.
- 562 Charles Audet, Warren Hare, Charles Audet, and Warren Hare. Chapter 3: The begining of
   563 dfo algorithms in derivative-free and blackbox optimization. *Derivative-Free and Blackbox* 564 Optimization, pp. 33–47, 2017e.
- Thomas Back. Evolutionary algorithms in theory and practice: evolution strategies, evolutionary programming, genetic algorithms. Oxford university press, 1996.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma,
  Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful
  and harmless assistant with reinforcement learning from human feedback. arXiv preprint
  arXiv:2204.05862, 2022.
- 572 573 Dimitri Bertsekas. *Convex optimization algorithms*. Athena Scientific, 2015.
- Bernd Bischl, Martin Binder, Michel Lang, Tobias Pielok, Jakob Richter, Stefan Coors, Janek Thomas, Theresa Ullmann, Marc Becker, Anne-Laure Boulesteix, et al. Hyper-parameter optimization: Foundations, algorithms, best practices, and open challenges. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 13(2):e1484, 2023.
- Djallel Bouneffouf, Irina Rish, and Charu Aggarwal. Survey on applications of multi-armed and contextual bandits. In 2020 IEEE Congress on Evolutionary Computation (CEC), pp. 1–8. IEEE, 2020.
- Franco Busetti. Simulated annealing overview. World Wide Web URL www. geocities. com/francorbusetti/saweb. pdf, 4, 2003.
  - Yulong Cao, Chaowei Xiao, Dawei Yang, Jing Fang, Ruigang Yang, Mingyan Liu, and Bo Li. Adversarial objects against lidar-based autonomous driving systems. arXiv preprint arXiv:1907.05418, 2019.
- Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*, pp. 15–26, 2017.
  - Maurice Clerc. Particle swarm optimization, volume 93. John Wiley & Sons, 2010.

628

634

635

636

- Enrique Dehaerne, Bappaditya Dey, Sandip Halder, Stefan De Gendt, and Wannes Meert.
  Code generation using machine learning: A systematic review. *Ieee Access*, 10:82434–82455, 2022.
- Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo
  Li. Boosting adversarial attacks with momentum. In *Proceedings of the IEEE conference* on computer vision and pattern recognition, pp. 9185–9193, 2018.
- Christian Gehring, Stelian Coros, Marco Hutter, Michael Bloesch, Péter Fankhauser,
  Markus A Hoepflinger, and Roland Siegwart. Towards automatic discovery of agile gaits
  for quadrupedal robots. In 2014 IEEE international conference on robotics and automation (ICRA), pp. 4243–4248. IEEE, 2014.
- Daniel Golovin, John Karro, Greg Kochanski, Chansoo Lee, Xingyou Song, and Qiuyi
   Zhang. Gradientless descent: High-dimensional zeroth-order optimization. arXiv preprint arXiv:1911.06317, 2019.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing
   adversarial examples. arXiv preprint arXiv:1412.6572, 2014.
- Kai Han, Yunhe Wang, Hanting Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang, An Xiao, Chunjing Xu, Yixing Xu, et al. A survey on vision transformer. *IEEE transactions on pattern analysis and machine intelligence*, 45(1):87–110, 2022.
- N. Hansen, A. Auger, R. Ros, O. Mersmann, T. Tušar, and D. Brockhoff. COCO: A platform
  for comparing continuous optimizers in a black-box setting. *Optimization Methods and Software*, 36:114–144, 2021. doi: https://doi.org/10.1080/10556788.2020.1808977.
- <sup>618</sup> Nikolaus Hansen. The cma evolution strategy: A tutorial. arXiv preprint arXiv:1604.00772, 2016.
   <sup>620</sup> <sup>620</sup>
- Nikolaus Hansen, Anne Auger, Raymond Ros, Steffen Finck, and Petr Pošík. Comparing
   results of 31 algorithms from the black-box optimization benchmarking bbob-2009. In
   *Proceedings of the 12th annual conference companion on Genetic and evolutionary com- putation*, pp. 1689–1696, 2010.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang,
  Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685, 2021.
- Antoine Lesage-Landry, Joshua A Taylor, and Iman Shames. Second-order online nonconvex optimization. *IEEE Transactions on Automatic Control*, 66(10):4866–4872, 2020.
- <sup>631</sup> Chengshu Li, Jacky Liang, Andy Zeng, Xinyun Chen, Karol Hausman, Dorsa Sadigh, Sergey
  <sup>632</sup> Levine, Li Fei-Fei, Fei Xia, and Brian Ichter. Chain of code: Reasoning with a language
  <sup>633</sup> model-augmented code emulator. arXiv preprint arXiv:2312.04474, 2023.
  - TP Lillicrap. Continuous control with deep reinforcement learning. arXiv preprint arXiv:1509.02971, 2015.
- Bingyu Liu, Yuhong Guo, Jianan Jiang, Jian Tang, and Weihong Deng. Multi-view correlation based black-box adversarial attack for 3d object detection. In *Proceedings of the* 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining, pp. 1036–1044, 2021.
- Sijia Liu, Pin-Yu Chen, Bhavya Kailkhura, Gaoyuan Zhang, Alfred O. Hero III, and Pramod K. Varshney. A primer on zeroth-order optimization in signal processing and machine learning: Principals, recent advances, and applications. *IEEE Signal Processing Magazine*, 37(5):43–54, 2020. doi: 10.1109/MSP.2020.3003837.
- Ilya Loshchilov, Marc Schoenauer, and Michèle Sebag. Bi-population cma-es agorithms
   with surrogate models and line searches. In *Proceedings of the 15th annual conference companion on Genetic and evolutionary computation*, pp. 1177–1184, 2013.

683

684 685

686

687

688

700

- 648 Minfang Lu, Shuai Ning, Shuangrong Liu, Fengyang Sun, Bo Zhang, Bo Yang, and Lin 649 Wang. Opt-gan: a broad-spectrum global optimizer for black-box problems by learning 650 distribution. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 37, 651 pp. 12462–12472, 2023.
- 652 Gabriel Resende Machado, Eugênio Silva, and Ronaldo Ribeiro Goldschmidt. Adversarial 653 machine learning in image classification: A survey toward the defender's perspective. 654 ACM Computing Surveys (CSUR), 55(1):1–38, 2021. 655
- 656 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. arXiv preprint 658 arXiv:1706.06083, 2017.
- 659 K. T. Yasas Mahima, Asanka G. Perera, Sreenatha Anavatti, and Matt Garratt. Toward 660 robust 3d perception for autonomous vehicles: A review of adversarial attacks and coun-661 termeasures. *IEEE Transactions on Intelligent Transportation Systems*, pp. 1–27, 2024. 662 doi: 10.1109/TITS.2024.3456293. 663
- Skogby Steinholtz Olof. A comparative study of black-box optimization algorithms for 664 665 tuning of hyper-parameters in deep neural networks, 2018.
- 666 Shanker G Radhakrishna Prabhu, Richard C Seals, Peter J Kyberd, and Jodie C Wetherall. 667 A survey on evolutionary-aided design in robotics. *Robotica*, 36(12):1804–1821, 2018. 668
- 669 Zohar Rimon, Tom Jurgenson, Orr Krupnik, Gilad Adler, and Aviv Tamar. Mamba: 670 an effective world model approach for meta-reinforcement learning. arXiv preprint arXiv:2403.09859, 2024. 671
- 672 Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Ba-673 log, M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming 674 Wang, Omar Fawzi, et al. Mathematical discoveries from program search with large lan-675 guage models. *Nature*, 625(7995):468–475, 2024. 676
- Vinu Sankar Sadasivan, Aounon Kumar, Sriram Balasubramanian, Wenxiao Wang, and 677 Soheil Feizi. Can ai-generated text be reliably detected?, 2024. URL https://arxiv. 678 org/abs/2303.11156. 679
- 680 Elad Sarafian, Mor Sinay, Yoram Louzoun, Noa Agmon, and Sarit Kraus. Explicit gradient 681 learning for black-box optimization. In *ICML*, pp. 8480–8490, 2020. 682
  - Sayantan Sarkar, Ankan Bansal, Upal Mahbub, and Rama Chellappa. Upset and angri: Breaking high performance image classifiers. arXiv preprint arXiv:1707.01159, 2017.
  - Sanli Tang, Xiaolin Huang, Mingjian Chen, Chengjin Sun, and Jie Yang. Adversarial attack type i: Cheat classifiers by significant changes. *IEEE transactions on pattern analysis and* machine intelligence, 43(3):1100–1109, 2019.
- Yunhao Tang. Guiding evolutionary strategies with off-policy actor-critic. In AAMAS, pp. 689 1317-1325, 2021. 690
- 691 Sasha Targ, Diogo Almeida, and Kevin Lyman. Resnet in resnet: Generalizing residual 692 architectures. arXiv preprint arXiv:1603.08029, 2016. 693
- Chun-Chen Tu, Paishun Ting, Pin-Yu Chen, Sijia Liu, Huan Zhang, Jinfeng Yi, Cho-Jui 694 Hsieh, and Shin-Ming Cheng. Autozoom: Autoencoder-based zeroth order optimization 695 method for attacking black-box neural networks. In Proceedings of the AAAI conference 696 on artificial intelligence, volume 33, pp. 742–749, 2019. 697
- 698 A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 699 2017.
- Chaowei Xiao, Bo Li, Jun-Yan Zhu, Warren He, Mingyan Liu, and Dawn Song. Generating 701 adversarial examples with adversarial networks. arXiv preprint arXiv:1801.02610, 2018.

702 703 704 705	Wei Zhang, Chaoqun Wan, Yonggang Zhang, Yiu-ming Cheung, Xinmei Tian, Xu Shen, and Jieping Ye. Interpreting and improving large language models in arithmetic calculation. arXiv preprint arXiv:2409.01659, 2024.
706 707 708 709	Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuan- dong Tian. Galore: Memory-efficient llm training by gradient low-rank projection. arXiv preprint arXiv:2403.03507, 2024.
710 711 712 713 714	Lu Zhen, Kai Wang, Hongtao Hu, and Daofang Chang. A simulation optimization framework for ambulance deployment and relocation problems. <i>Computers &amp; Industrial Engineering</i> , 72:12–23, 2014.
715 716 717	A Appendix
718 719	A.1 Full OHGL Algorithm
720 721	Algorithm 2 OHCL (Heggian Weighted Credient Learning)
722	Algorithm 2 Origi (nessian weighted Gradient Learning)
723	<b>Require:</b> $x_0, \Omega, \alpha, \epsilon_0, \gamma_\alpha < 1, \gamma_\epsilon < 1, n_{\max}, \lambda$
724	k = 0
725	j = 0
726	$M_j \leftarrow M$ $M_{rec} = h \to \Omega \to \mathbb{D}^n$
720	$\operatorname{Map} n_0: \Omega \to \mathbb{R}^n$
700	$ \begin{array}{c} \text{Map } W_0 : \mathbb{R} \to \mathbb{R} \\ \text{while Pudget } > 0 \ \text{de} \end{array} $
720	Fyplore:
729	Concrete samples $\mathcal{D}_{r} = \{\tilde{x}_{r}\}^{m}$ $\tilde{x}_{r} \in V$ $(\tilde{x}_{r})$
730	Evaluate samples $\nu_k = \{x_i\}_{i=1}, x_i \in V_{\epsilon_k}(x_k)$ Evaluate samples $u_i = f(h^{-1}(\tilde{x}_i)), i = 1$
731	Evaluate samples $y_i = f(n_0^{(i)}(x_i)), i = 1, \dots, m$
732	And tuples to the replay buller: $\mathcal{D} = \mathcal{D} \cup \mathcal{D}_k$
733	$\mathcal{T} \leftarrow \int (\tilde{x}, \tilde{x})   \ \tilde{x} - \tilde{x}\ _{2} < \epsilon  \forall i  i \end{bmatrix}$
734	$\mathcal{T} \leftarrow \{(x_i, x_j) \mid \ x_i - x_j\ _2 < \epsilon, \forall i, j\}$ Select the first <i>m</i> tuples from $\mathcal{T}$
735	Weighted Output Man:
736	Assign weights $w_i = W_i(u_i)$ to samples based on function values
737	Apply squashing function to the outputs: $\tilde{y}_i = r_i(w_i y_i)$ , $i = 1, \dots, m$
738	Higher-Order Gradient and Hessian Learning:
739	Calculate the Hessian from the Jacobian of the network: $H_k(x) = J(q_{\theta_k}(x))$
740	Optimize the network with GD using the formula in 7
741	Gradient Descent:
742	Update the current solution using second-order information:
743	$x_{k+1} \leftarrow x_k - \alpha \cdot \tilde{g}_{\theta_k}(x_k)$
744	if $f(h_i^{-1}(\tilde{x}_{k+1})) > f(h_i^{-1}(\tilde{x}_k))$ for $n_{\max}$ times in a row then
745	Generate a new trust region: $\Omega_{i+1} = \gamma_{\alpha}  \Omega_i $ with center at $x_{\text{best}}$
746	$\operatorname{Map} h_{i}: \Omega \to \mathbb{R}^{n}$
747	$\operatorname{Map}\check{W_j}:\mathbb{R}\to\mathbb{R}$
748	$j \leftarrow j + 1$
749	$\epsilon_j \leftarrow \gamma_\epsilon \epsilon_j$
750	end if $(1-1)(2-1)(2-1)(2-1)(2-1)(2-1)(2-1)(2-1)$
751	If $f(n_j^{-1}(x_k)) < f(n_j^{-1}(x_{\text{best}}))$ then
752	$x_{ ext{best}} = h_j^{-1}( ilde{x}_k)$
752	end if
753	$k \leftarrow k + 1; \text{Budget} \leftarrow \text{Budget} - m$
754	end while
100	roturn <i>x</i> ,

A.2 STATISTICAL ANALYSIS

756

757

758 759 Table 4: T-test of our algorithms against the benchmark in dimension 40 760 761 Metric OGL HOGL HGL 762 IGL 1.000000 1.000000 1.000000 763 OGL 0.5000000.278378 0.334161764 OHGL 0.7216220.5000000.473754765 HGL 0.6658390.5262460.500000766 EGL 0.999961 0.9999750.999098 767 CMA 1.000000 1.0000001.000000 768 0.999911T-CMA 0.9999340.999788 769 L-CMA 0.9998590.9998430.998842770 771 772 773 A.3 GENERATED CODE FROM CODE FORCE EXAMPLE 774 775 776 def count\_triplets(n: int, x: int) -> int: 777 778 This function calculates the number of triplets (a, b, c) of positive integers such that: 1. a\*b + a\*c + b\*c â% ¤ n 779 2.  $a + b + c \hat{a}_{a}^{\mu} x$ where the order of triplets matters (e.g., (1, 1, 2) and (1, 2, 1) are different triplets) and a, b, c must be strictly greater than 0. 781 ..... 782 pass def check\_function(): 783 assert count\_triplets(10, 5) == 3, 'Test case 1 failed' 784 assert count\_triplets(15, 10) == 6, 'Test case 2 failed' assert count\_triplets(20, 15) == 10, 'Test case 3 failed' 785 assert count\_triplets(100, 50) == 50, 'Test case 4 failed' 786 print('All test cases passed!') 787 788 Figure 8: Step 1 789 790 Figure 9: Generated code samples by the algorithm 791 792 793 794 795 def count\_triplets(n): 796 This function calculates the number of triplets (a, b, c) of positive integers such that: 797 1. a\*b + a\*c + b\*c <= n798 2.  $a + b + c \le x$ where the order of triplets matters (e.g., (1, 1, 2) and (1, 2, 1) are different triplets) 799 and a, b, c must be strictly greater than 0. 800 count = 0801 # Iterate through all possible values of a, b, c 802 for a in range(1, x + 1): for b in range(1, x + 1): 803 for c in range(1, x + 1): 804 if  $a * b + a * c + b * c \leq n$  and  $a + b + c \leq x$ : count += 1805 return count 806 807



Figure 10: Step 21

Figure 11: Generated code samples by the algorithm

```
811
812
813
814
815
816
817
818
819
820
```

```
822
823
824
```

825 826

831 832

837

843 844 845

821

```
This function calculates the number of triplets (a, b, c) of positive integers such that:
1. a*b + a*c + b*c <= n
2. a + b + c \le x
where the order of triplets matters (e.g., (1, 1, 2) and (1, 2, 1) are different triplets)
and a, b, c must be strictly greater than 0.
count = 0
# Iterate over a
for a in range(1, x + 1):
    # Iterate over b
   for b in range(1, x - a + 1):
        # Find the maximum c based on the condition a + b + c \le x
       max_c = min(x - a - b, (n - a * b) // (a + b)) # upper bound for c
        # Use list comprehension to count valid triplets
        count += sum(1 for c in range(1, max_c + 1) if a * b + a * c + b * c <= n)
```

return count

def count\_triplets(n):

#### Figure 12: Step 50

Figure 13: Generated code samples by the algorithm

#### A.4 Adaptive Sampling Size

Gradient learning algorithms collect samples to train the surrogate gradient model. While 833 more samples can potentially lead to more accurate models, it is important to curtail the 834 number of samples in each training iteration to be able to execute enough optimization steps 835 before consuming the budget. On the other hand, one advantage that gradient learning has 836 over direct objective learning is that we train the model with pairs of samples instead of single evaluation points, s.t. for a sample set of size  $n_s$  we can draw as many as  $n_p = n_s^2$ 838 pairs to train our model.

839 We empirically tested the optimal number of exploration size (i.e.  $n_s$ ) in each training step 840 and the optimal sample pair size  $(n_p)$  and found that a squared root profile is optimal both 841 hyperparameters. Therefore, we use the following equation to control these parameters: 842

$$n_s = 8 \cdot \left\lceil \sqrt{N} \right\rceil$$

$$n_p = 2000 \cdot \left\lceil \sqrt{N} \right\rceil$$
(9)

where N is the problem size.

850 851 852

853

854

#### TRUST REGION MANAGEMENT A.5

855 After finding the optimal solution inside the trust region, EGL shifts and scales down (i.e. 856 shrinks) the trust region so that the optimal solution is centered at its origin. In case the 857 decent path is long so that the minimum point is far away from the current trust region, 858 this shrinking has the potential to slow down the learning rate and impede convergence. 859 To prevent unnecessary shrinking of the trust region, we distinguish between two conver-860 gence types: **interior convergence**: in this case, we shrink the trust region while moving its center and **boundary convergence** where we only shift the center without applying 861 shrinking (When the algorithm reaches the edges of the TR). This adjustment prevents 862 fast convergence of the step size to zero before being able to sufficiently explore the input 863 domain.

## A.6 Additional Empirical Results

#### 868 Parameter Value Description 869 870 $8 \cdot$ $\sqrt{\text{dimension}}$ The number of iterations the al-Exploration size 871 gorithm will run. 872 How much distance the algo-Maximum movement to shrink $0.2 \cdot \sqrt{\text{dimension}}$ 873 rithm needs to cover to prevent 874 shrinking when reaching conver-875 gence. 876 $0.4 \cdot \sqrt{\text{dimension}}$ Epsilon Epsilon size to sample data. 877 Epsilon Factor 0.97How much we shrink the epsilon 878 each iteration. 879 Minimum epsilon 0.0001The smallest epsilon we use. 880 √dimension Database size 20,000 · How many tuples we used to train for each iteration. 882 Gradient network dimension-10-15-10-dimension What kind of network we use. 883 Budget 150,000 How much budget the algorithm uses. 884

#### Table 5: EGL Best Parameters

#### A.7 BENCHMARK ALGORITHMS

```
\begin{array}{l} \textbf{Algorithm 3 CMA with a trust region} \\ \hline \textbf{Require: total_budget, budget = 0, startpointx, \delta = 1, \mu = 0, \gamma} \\ \textbf{while budget \leq total_budget do} \\ \textbf{new_generation, should_stop} = CMA(x) \\ \textbf{new_generation} = mu + \delta * \texttt{new_generation} \\ evalutaions = space(\texttt{new_generation}) \\ budget = budget + len(evalutaions) \\ UpdateCovarMatrix(evaluations) \\ \textbf{if should_stop then} \\ \delta = \delta * \gamma \\ \mu = argmin_{\texttt{x new_generation}}(space(x)) \\ \textbf{end if} \\ \textbf{end while} \end{array}
```

905

885 886 887

888 889

890

891

892

893

895

896

897

899

900

865 866

867

#### A.8 CMA VERSIONS

The trust region is key in the EGL algorithm (see Figures 1). We map between the search space  $(\Omega)$  and trust region  $(\Omega^*)$  using the tanh function, which provides smooth transitions but exhibits logarithmic behavior at the edges, slowing large steps. While helpful for problems with many local minima, this behavior restricts CMA's ability to take large steps, especially in high dimensions (see Figure 14). We developed two CMA variants: one with linear TR (L-CMA) and one with tanh TR(T-CMA). The linear TR allowed larger steps, improving exploration, while the tanh TR limited performance.

912 913 914

#### A.9 Adversarial Attack Implementation

915 Our adversarial attack leverages the flexibility of black-box optimization (BBO), which is 916 not constrained by differentiability requirements on the objective function. This allows us 917 to manipulate the search space freely. Specifically, our attack focuses on modifying the 918 positions of m fixed-size boxes within the image, along with the permutations applied to

Figure 14: Comparing CMA with different trust regions types across different dimensions

the pixels inside these boxes. This approach provides two key advantages: it enables us to limit the extent of the alterations we introduce to the image and simultaneously reduces the dimensionality of the search space  $\Omega$ .

The penalty function we use balances two competing goals: minimizing the cross-entropy loss to ensure misclassification and limiting the perturbation magnitude using the mean squared error (MSE) loss. To achieve this balance, we define the following loss function:

$$s(ce) = \sigma (b \cdot (ce - \epsilon))$$

$$Penalty_{mse}(ce, mse) (1 - s(ce)) \cdot mse^{n_1} + s(ce) \cdot (-mse^{n_2})$$

$$Penalty_{ce}(ce) (1 - s(ce))) \cdot e^{ce*n_1} + s(ce) \cdot (ln(ce^{n_2}))$$

$$Penalty(x) = Penalty_{mse}(ce(x), mse(x)) - Penalty_{ce}(ce(x))$$

In this formulation:

- mse  $\in [0, 1]$  and ce  $\in (-\infty, \infty)$ .
- The terms  $n_1 \leq n_2$  control the trade-off between the cross-entropy and MSE loss slopes.
- $\epsilon$  defines the minimum required cross-entropy loss to maintain misclassification.
- The parameter b governs the rate of transition between minimizing cross-entropy and MSE losses.
- The function  $\sigma$  is the sigmoid function, which controls how much focus is given to minimizing cross-entropy loss over MSE loss as ce increases.

This function balances the CR and the MSE, finding a minimum with CE to prevent detection by the classifier while minimizing the perturbation of the image.

A.10 Full Robustness Experiment

Table 6: Comparison of Algorithm Versions on coco benchmark, comparing the error, the std and the budget it takes to finish 99% of the progress

Version	Budget	Error	$\mathbf{Std}$	Solved Problems
		n	etworks	
$20 \ 40 \ 20$	58447.91	0.0060	0.0325	0.76
40 60 80 40	57960.44	0.0058	0.0325	0.76
60 80 60	58371.31	0.0060	0.0325	0.76
40 80 40	58159.34	0.0061	0.0325	0.76
$40 \ 15 \ 10 \ 40$	58667.00	0.0061	0.0325	0.76
		I	Epsilon	
0.1	39491.81	0.006	0.0325	0.9
				Continued on next pag

972	Version	Budget	Error	Std	Solved Problems		
973	0.4	40510.01	0.0057	0.0205	0.02		
974	0.4	48518.81	0.0057	0.0325	0.93		
975	0.5	49409.81	0.0000	0.0320	0.93		
976	0.0	54271.47 58667.00	0.0001	0.0520	0.9		
977	0.8	56007.00	0.0001	0.0525	0.9		
978			$\mathbf{Epsi}$	lon Facto	or		
979	0.8	39275.72	0.0105	0.0365	0.83		
980	0.9	47930.81	0.0068	0.0327	0.89		
981	0.95	53899.78	0.0060	0.0326	0.9		
982	0.97	58667.00	0.0061	0.0325	0.9		
983	0.99	63676.47	0.0065	0.0326	0.89		
984			Traini	ing Weig	$\mathbf{hts}$		
985	50%	58049.50	0.0059	0.0325	0.91		
986	60%	59477.47	0.0066	0.0333	0.89		
987	70%	57831.00	0.0057	0.0325	0.93		
988	83%	58667.00	0.0061	0.0325	0.9		
989	100%	59660.34	0.0061	0.0326	0.9		
990			Per	turbatio	n		
001	1	56807.97	0.0066	0.0333	0.89		
002	0.3	58935.38	0.0088	0.0420	0.87		
992	0.1	57699.78	0.0061	0.0325	0.89		
993	0.01	58075.03	0.0060	0.0325	0.89		
994	0	58667.00	0.0061	0.0325	0.89		
995	Euclidean	distance of	the algo	rithm be	efore shrinking trust region		
996	0.4	57172.53	0.0061	0.0325	0.9		
997	0.3	57324.53	0.0059	0.0325	0.91		
998	0.2	58667.00	0.0057	0.0325	0.93		
999	0.1	57327.50	0.0061	0.0325	0.89		
1000	0.01	57810.81	0.0062	0.0325	0.89		
1001	Value Normalizer I P						
1002	0.01	58211 59	0.0086	0.0420	0.83		
1003	0.01	58576.75	0.0061	0.0326	0.89		
1004	0.1	58667.00	0.0057	0.0325	0.93		
1005	0.2	59723.88	0.0059	0.0325	0.91		
1006	0.3	58194.97	0.0079	0.0370	0.87		
1007		<b>T</b> 7.	. 1 NT		Oct line		
1008	0.01	57999-41		nanzer 0.0225			
1009	0.01	57222.41 57776.38	0.0039	0.0325 0.0325	0.91		
1010	0.05	58667.00	0.0000 0.0057	0.0325 0.0325	0.93		
1011	0.2	58481 75	0.0061	0.0325	0.9		
1012	0.3	58629.59	0.0063	0.0325	0.89		
1013				0.001			
1014	0.00	102759.1C	ist Kegi	n Shrin	K Factor		
1015	0.99	103758.16	0.0177	0.0433	0.79		
1016	0.95	80000.09 58667.00	0.0093	0.0340 0.0225	0.02		
1017	0.9	20007.00 43503.22	0.0007	0.0320 0.0325	0.95		
1018	0.8	39790 50	0.0000 0.0075	0.0325	0.85		
1010	0.1	00100.00	0.0010	0.0020	0.00		
1013			Opti	imizer L	R		
1020	0.001	101550.59	0.0164	0.0364	0.76		
1021	0.005	64988.66	0.0066	0.0325	0.88		
1022	0.01	58667.00	0.0057	0.0325	0.93		
1023	0.02	62405.84	0.0056	0.0325	0.93		
1024					Continued on next page		
1025					. 0		

026	Version	Budget	Error	Std	Solved Problems
028	0.03	72041.81	0.0058	0.0325	0.93
029	Gradient LR				
030	0.0001	63283.41	0.0063	0.0325	0.89
)31	0.0005	59237.59	0.0086	0.0420	0.84
32	0.001	58667.00	0.0057	0.0325	0.93
33	0.002	58776.25	0.0061	0.0325	0.9
)34	0.003	57934.91	0.0061	0.0326	0.9
)35					

1038 A.11 THEORETICAL ANALYSIS

1040 A.11.1 MEAN GRADIENT ACCURACY

**Definition 1** The second-order mean gradient around x of radius  $\epsilon > 0$ :

$$\tilde{g}_{\varepsilon}(x) = \arg\min_{g \in \mathbb{R}^n} \int_{V_{\varepsilon}(x)} \left| g^{\top} \tau + \frac{1}{2} \tau^{\top} J(g) \tau - \left[ f(x+\tau) - f(x) \right] \right|^2 d\tau.$$

**Theorem 3** (Improved Controllable Accuracy): For any twice differentiable function  $f \in C^2$ , there exists  $\kappa_g(x) > 0$  such that for any  $\varepsilon > 0$ , the second-order mean-gradient  $\tilde{g}_{\varepsilon}(x)$  satisfies

$$\|\tilde{g}_{\varepsilon}(x) - \nabla f(x)\| \le \kappa_g(x)\varepsilon^2 \quad \text{for all } x \in \Omega$$

*Proof.* Since  $f \in C^2$ , the Taylor expansion around x is

$$f(x+\tau) = f(x) + \nabla f(x)^{\top} \tau + \frac{1}{2} \tau^{\top} H(x) \tau + R_x(\tau),$$

where H(x) is the Hessian matrix at x, and the remainder  $R_x(\tau)$  satisfies

 $|R_x(\tau)| \leq \frac{1}{6}k_q ||\tau||^3$ 

1059 By the Definition of  $g_{\epsilon}$ , an upper bound of  $\mathcal{L}(g_{\epsilon}(x))$  is:

$$\begin{aligned} & \underset{1061}{\overset{1061}{1062}} \qquad \mathcal{L}(g_{\epsilon}(x)) \leq \mathcal{L}(\nabla f(x)) = \int_{V_{\varepsilon}(x)} \left| \nabla f(x)^{\top} \tau + \frac{1}{2} \tau^{\top} H(x) \tau - \left( f(x+\tau) - f(x) \right) \right|^2 d\tau \\ & \underset{1063}{\overset{1064}{1064}} \end{aligned}$$

$$\leq \left(\frac{1}{6}k_g\right)^2 \int_{V_{\varepsilon}(x)} \|\tau\|^6 d\tau = \left(\frac{1}{6}k_g\right)^2 |V_1(x)|\varepsilon^{n+6}.$$

1068 We now find a lower bound:

lets define for convenience  $\delta_g = (g_{\varepsilon}(x) - \nabla f(x))$ ,  $\delta_H = J(g_{\varepsilon}(x)) - H(x)$  $\mathcal{L}(g_{\varepsilon}(x)) = \int_{V_{\varepsilon}(x)} |g_{\varepsilon}(x)\tau - \nabla f(x)\tau + \nabla f(x)\tau + \frac{1}{2}\tau J(g_{\varepsilon}(x))\tau - \nabla f(x)\tau + \frac{1}{2}\tau J(g_{\varepsilon}(x))\tau + \frac{1}{2}\tau J(g_{\varepsilon}(x)$  $\frac{1}{2}\tau H(x)\tau + \frac{1}{2}\tau H(x)\tau - f(x+\tau) + f(x)|^2 d\tau$  $= \int_{V_r(x)} (\delta_g^\top \tau + \frac{1}{2} \tau^\top \delta_H \tau - R_x(\tau))^2 d\tau$  $= \int_{V_{\tau}(x)} ((\delta_g^{\top} \tau)^2 + \frac{1}{4} (\tau^{\top} \delta_H \tau)^2 + R_x(\tau)^2 - 2\delta_g^{\top} \tau R_x(\tau) - \tau^{\top} \delta_H \tau R_x(\tau) + \tau^{\top} \delta_g \tau^{\top} \delta_H \tau) d\tau$ 



Figure 15: The search trajectory for an adversarial image generation with OGL

1103  
1104 Since 
$$a^2 + b^2 \ge 2ab$$
 we conclude that  $\frac{1}{4}(\tau^{\top}\delta_H\tau)^2 + R_x(\tau)^2 - \tau^{\top}\delta_H\tau R_x(\tau) \ge 0$   
1105  
1106  $\mathcal{L}(g_{\varepsilon}(x)) = \int_{V_{\varepsilon}(x)} ((\delta_g^{\top}\tau)^2 + \frac{1}{4}(\tau^{\top}\delta_H\tau)^2 + R_x(\tau)^2 - 2\delta_g^{\top}\tau R_x(\tau) - \tau^{\top}\delta_H\tau R_x(\tau) + \tau^{\top}\delta_g\tau^{\top}\delta_H\tau)d\tau$   
1108  $\ge \int_{V_{\varepsilon}(x)} ((\delta_g^{\top}\tau)^2 - 2\delta_g^{\top}\tau R_x(\tau) + \tau^{\top}\delta_g\tau^{\top}\delta_H\tau)d\tau$   
1109  $= \int_{V_{\varepsilon}(x)} (\delta_g^{\top}\tau)^2d\tau - 2\int_{V_{\varepsilon}(x)} \delta_g^{\top}\tau R_x(\tau)d\tau + \int_{V_{\varepsilon}(x)} \tau^{\top}\delta_g\tau^{\top}\delta_H\tau d\tau$   
1113  $= \int_{V_{\varepsilon}(x)} (\delta_g^{\top}\tau)^2d\tau - 2\int_{V_{\varepsilon}(x)} \delta_g^{\top}\tau R_x(\tau)d\tau + \int_{V_{\varepsilon}(x)} \tau^{\top}\delta_g\tau^{\top}\delta_H\tau d\tau$   
1114 115  
1116  
1117 We will colit the equation into 3 components: First  $A = \int_{V_{\varepsilon}(x)} (\delta_g^{\top}\tau)^2d\tau$  is we open this

We will split the equation into 3 components: First  $A = \int_{V_{\varepsilon}(x)} (\delta_g^{\top} \tau)^2 d\tau$ , is we open this expression we see that 

$$(\delta_g^{\top}\tau)^2 = \sum_{i=1}^n \sum_{j=1}^n (\delta_g)_i (\delta_g)_j \tau_i \tau_j.$$

Since  $V_{\varepsilon}(x)$  is symmetric around x odd moments of  $\tau$  integrate to zero, therefore 

$$\begin{array}{l} \mathbf{1124} \\ \mathbf{1125} \\ \mathbf{1126} \\ \mathbf{1126} \\ \mathbf{1127} \end{array} A = \int_{V_{\varepsilon}(x)} (\delta_g^{\top} \tau)^2 d\tau = \int_{V_{\varepsilon}(x)} (\sum_{i=1}^n \sum_{j=1}^n (\delta_g)_i (\delta_g)_j \tau_i \tau_j) d\tau = \int_{V_{\varepsilon}(x)} (\sum_{i=1}^n (\delta_g)_i^2 \tau_i^2) d\tau = ||\delta_g||^2 \varepsilon^{n+2} |V_1(x)| \\ \mathbf{1127} \end{array}$$

$$B = 2 \int_{V_{\varepsilon}(x)} \delta_{g}^{\top} \tau R_{x}(\tau) d\tau \leq 2 \int_{V_{\varepsilon}(x)} ||\delta_{g}|| \cdot ||\tau|| \cdot ||R_{x}(\tau)||d\tau \leq 2||\delta_{g}|| \int_{V_{\varepsilon}(x)} ||\tau||^{\frac{1}{6}} k_{g} \tau^{3} d\tau$$
$$\leq \frac{1}{3} k_{g} ||\delta_{g}|| \int ||\tau||^{4} d\tau = \frac{1}{3} k_{g} ||\delta_{g}|| \varepsilon^{n+4} |V_{1}(x)|$$

1132  
1133 
$$\leq \frac{1}{3}k_g ||\delta_g|| \int_{V_{\varepsilon}(x)} ||\tau||^4 d\tau = \frac{1}{3}k_g ||\delta_g||\varepsilon^{n+4}|V_1(x)|$$

Let  $B = 2 \int_{V_{\varepsilon}(x)} \delta_g^{\top} \tau R_x(\tau) d\tau$ 

Finally  $C = \int_{V_{\varepsilon}(x)} \tau^{\top} \delta_g \tau^{\top} \delta_H \tau d\tau$ , we should remember the property derived from  $V_{\varepsilon}(x)$ symmetry. If we open the vector multiplication we get

$$\tau^{\top} \delta_{g} \tau^{\top} \delta_{H} \tau = \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} \tau_{i}(\delta_{g})_{i}(\delta_{H})_{jk} \tau_{j} \tau_{k} = \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} (\delta_{g})_{i}(\delta_{H})_{jk} \tau_{i} \tau_{j} \tau_{k}.$$

Here there is no way for an even moment of  $\tau$  to exist so C = 0 To sum this up, we get

$$\mathcal{L}(g_{\varepsilon}(x)) \ge A - B + C \ge ||\delta_g||^2 \varepsilon^{n+2} |V_1(x)| - \frac{1}{3}k_g||\delta_g||\varepsilon^{n+4}|V_1(x)|$$

We combine the lower and upper-bound

1137 1138 1139

1141 1142

 $\begin{aligned} ||\delta_g||^2 \varepsilon^{n+2} |V_1(x)| &- \frac{1}{3} k_g ||\delta_g||\varepsilon^{n+4} |V_1(x)| \le (\frac{1}{6} k_g)^2 |V_1(x)|\varepsilon^{n+6} \\ ||\delta_g||^2 &- \frac{1}{3} k_g ||\delta_g||\varepsilon^2 - (\frac{1}{6} k_g)^2 \varepsilon^4 \le 0 \\ ||\delta_g||^2 &\le \frac{1}{3} k_g \varepsilon^2 + \sqrt{\frac{1}{9} k_g^2 \varepsilon^4 + 4(\frac{1}{6} k_g)^2 \varepsilon^4}}{2} = \frac{\frac{1}{3} k_g \varepsilon^2 + \sqrt{2\frac{1}{9} k_g^2 \varepsilon^4}}{2} = \frac{1}{3\sqrt{2}} k_g \varepsilon^2 = k_g(x) \varepsilon^2 \end{aligned}$ 

1147 1148 1149

1156 1157 1158

1159

1162 1163

1165 1166

1170 1171

1182

# 1150 A.11.2 CONVERGENCE ANALYSIS

**Theorem 4** Let  $f: \Omega \to \mathbb{R}$  be a convex function with Lipschitz continuous gradient, i.e.  $f \in \mathcal{C}^{+1}$  and a Lipschitz constant  $\kappa_f$  and let  $f(x^*)$  be its optimal value. Suppose a controllable mean-gradient model  $g_{\varepsilon}$  with error constant  $\kappa_g$ , the gradient descent iteration  $x_{k+1} = x_k - \alpha g_{\varepsilon}(x_k)$  with a sufficiently small  $\alpha$  s.t.  $\alpha \leq \min(\frac{1}{\kappa_g}, \frac{1}{\kappa_f})$  guarantees:

1. For 
$$\varepsilon \leq \frac{\|\nabla f(x)\|}{5\alpha}$$
, monotonically decreasing steps s.t.  $f(x_{k+1}) \leq f(x_k) - 2.25 \frac{\varepsilon^2}{\alpha}$ .

2. After a finite number of iterations, the descent process yields  $x^*$  s.t.  $\|\nabla f(x^*)\| \leq \frac{5\varepsilon}{\alpha}$ .

1160 *Proof.* For a convex function with Lipschitz continuous gradient, the following inequality 1161 holds for all  $x_k$ 

$$f(x) \le f(x_k) + (x - x_k) \cdot \nabla f(x_k) + \frac{1}{2} \kappa_f ||x - x_k||^2$$
(10)

Plugging in the iteration update  $x_{k+1} = x_k - \alpha g_{\varepsilon}(x_k)$  we get

$$f(x_{k+1}) \le f(x_k) - \alpha g_{\varepsilon}(x_k) \cdot \nabla f(x_k) + \alpha^2 \frac{1}{2} \kappa_f \|g_{\varepsilon}(x_k)\|^2$$
(11)

1167 For a controllable mean-gradient we can write  $||g_{\varepsilon}(x) - \nabla f(x)|| \leq \varepsilon^2 \kappa_g$ , therefore we can 1168 write  $g_{\varepsilon}(x) = \nabla f(x) + \varepsilon^2 \kappa_g \xi(x)$  s.t.  $||\xi(x)|| \leq 1$  so the inequality is

$$f(x_{k+1}) \le f(x_k) - \alpha \|\nabla f(x_k)\|^2 - \alpha \varepsilon^2 \kappa_g \xi(x_k) \cdot \nabla f(x_k) + \alpha^2 \frac{1}{2} \kappa_f \|\nabla f(x) + \varepsilon^2 \kappa_g \xi(x)\|^2$$
(12)

Using the equality  $||a + b||^2 = ||a||^2 + 2a \cdot b + ||b||^2$  and the Cauchy-Schwartz inequality inequality  $a \cdot b \le ||a|| ||b||$  we can write

$$\begin{aligned}
f(x_{k+1}) &\leq f(x_k) - \alpha \|\nabla f(x_k)\|^2 - \alpha \varepsilon^2 \kappa_g \|\xi(x_k)\| \cdot \|\nabla f(x_k)\| \\
+ \alpha^2 \frac{1}{2} \kappa_f \left( \|\nabla f(x)\|^2 + 2\varepsilon^2 \kappa_g \|\xi(x)\| \cdot \|\nabla f(x_k)\| + \varepsilon^4 \kappa_g^2 \|\xi(x)\|^2 \right) \\
&\leq f(x_k) - \alpha \|\nabla f(x_k)\|^2 - \alpha \varepsilon^2 \kappa_g \|\nabla f(x_k)\| \\
+ \frac{\alpha^2}{2} \kappa_f \|\nabla f(x)\|^2 + \alpha^2 \kappa_f \varepsilon^2 \kappa_g \|\nabla f(x_k)\| + \frac{\alpha^2 \varepsilon^4}{2} \kappa_f \kappa_g^2
\end{aligned}$$
(13)
$$\begin{aligned}
+ \frac{\alpha^2}{2} \kappa_f \|\nabla f(x)\|^2 + \alpha^2 \kappa_f \varepsilon^2 \kappa_g \|\nabla f(x_k)\| + \frac{\alpha^2 \varepsilon^4}{2} \kappa_f \kappa_g^2
\end{aligned}$$

$$= f(x_k) - \left(\alpha - \frac{k_f}{2}\alpha^2\right) \|\nabla f(x_k)\|^2 + \left(-\alpha + k_f\alpha^2\right) \|\nabla f(x_k)\|\varepsilon^2 \kappa_g + \frac{k_f}{2}\alpha^2 \varepsilon^4 \kappa_g^2$$

Using the requirement  $\alpha \leq \min(\frac{1}{\kappa_g}, \frac{1}{\kappa_f})$  it follows that  $\alpha \kappa_g \leq 1$  and  $\alpha \kappa_f \leq 1$  so

1185  
1186 
$$f(x_{k+1}) \le f(x_k) - (\alpha - \frac{k_f}{2}\alpha^2) \|\nabla f(x_k)\|^2 + (-\alpha + k_f\alpha^2) \|\nabla f(x_k)\|\varepsilon^2 \kappa_g + \frac{k_f}{2}\alpha^2 \varepsilon^4 \kappa_g^2$$
1187 (14)

Now, for x s.t.  $\|\nabla f(x)\| \geq \frac{5\varepsilon^2}{\alpha}$  then  $\varepsilon^2 \leq \|\nabla f(x)\|_{\frac{\alpha}{5}}^{\alpha}$ . Plugging it into our inequality we obtain 

$$f(x_{k+1}) \leq f(x_k) - \frac{\alpha}{2} \|\nabla f(x_k)\|^2 + \frac{\alpha^2}{50} \kappa_g \|\nabla f(x_k)\|^2$$
  
$$\leq f(x_k) - \frac{\alpha}{2} \|\nabla f(x_k)\|^2 + \frac{\alpha}{50} \|\nabla f(x_k)\|^2$$
  
$$= f(x_k) - 0.48\alpha \|\nabla f(x_k)\|^2$$
  
$$\leq f(x_k) - 12 \frac{\varepsilon^4}{\alpha}$$
(15)

Therefore, for all x s.t.,  $\|\nabla f(x)\| \ge \frac{5\varepsilon^2}{\alpha}$  we have a monotonically decreasing step with finite size improvement, after a finite number of steps we obtain  $x^*$  for which  $\|\nabla f(x^*)\| \leq \frac{5\varepsilon^2}{\alpha}$ . 

A.11.3 Optimistic Gradient Accuracy 

**Definition 2** The optimistic gradient around x of radius  $\epsilon > 0$ : 

$$\tilde{g}_{\varepsilon}(x) = \arg\min_{g \in \mathbb{R}^n} \int_{V_{\varepsilon}(x)} w(\tau) \left| g^{\top} \tau - \left[ f(x+\tau) - f(x) \right] \right|^2 d\tau.$$

**Theorem 5** (Optimistic gradient controllable Accuracy): For any twice differentiable func-tion  $f \in \mathcal{C}^2$ , there exists  $\kappa_a(x) > 0$  such that for any  $\varepsilon > 0$ , the second-order mean-gradient  $\tilde{g}_{\varepsilon}(x)$  satisfies 

 $\|\tilde{g}_{\varepsilon}(x) - \nabla f(x)\| \leq \kappa_q(x)\varepsilon \text{ for all } x \in \Omega.$ 

*Proof.* Since  $f \in C^2$  the Taylor expansion around x is 

$$f(x+\tau) = f(x) + \nabla f(x)^{\top} \tau + R_x(\tau)$$

where H(x) is the Hessian matrix at x, and the remainder  $R_x(\tau)$  satisfies 

$$|R_x(\tau)| \le \frac{1}{2}k_g \|\tau\|^2$$

By definition  $g_{\epsilon}$ , an upper bound  $\mathcal{L}(g_{\epsilon}(x))$  is: 

$$\mathcal{L}(g_{\epsilon}(x)) \leq \mathcal{L}(\nabla f(x)) = \int_{V_{\varepsilon}(x)} \left| \nabla f(x)^{\top} \tau - \left( f(x+\tau) - f(x) \right) \right|^{2} d\tau = \int_{V_{\varepsilon}(x)} |R_{x}(\tau)|^{2} d\tau$$
$$\leq \left( \frac{1}{2} k_{g} \right)^{2} \int_{V_{\varepsilon}(x)} \|\tau\|^{4} d\tau = \left( \frac{1}{2} k_{g} \right)^{2} |V_{1}(x)| \varepsilon^{n+4}.$$

*Proof.* The optimistic mean gradient around x of radius  $\epsilon > 0$ : 

$$\tilde{g}_{\varepsilon}(x) = \arg\min_{g \in \mathbb{R}^n} \int_{V_{\varepsilon}(x)} w(\tau) \left| g^{\top} \tau - \left[ f(x+\tau) - f(x) \right] \right|^2 d\tau.$$

Since  $V_{\varepsilon}(x)$  is symmetric around x odd moments of  $\tau$  integrate to zero, therefore

$$A = \int_{V_{\varepsilon}(x)} w(\tau) (\delta_g^{\top} \tau)^2 d\tau = \int_{V_{\varepsilon}(x)} w(\tau) (\sum_{i=1}^n \sum_{j=1}^n (\delta_g)_i (\delta_g)_j \tau_i \tau_j) d\tau$$

1237  
1238  
1239 
$$= \int_{V_{\varepsilon}(x)} w(\tau) (\sum_{i=1}^{n} (\delta_g)_i^2 \tau_i^2) d\tau = \int_{V_{\varepsilon}(x)} w(\tau) \|\delta_g\|^2 \cdot \|\tau\|^2 d\tau$$

1240  
1241 = 
$$\|\delta_{\varepsilon}\|^{2} \varepsilon^{n+2} |V_{1}(x)| \int w(\tau) < \|\delta_{\varepsilon}\|^{2} \varepsilon^{r}$$

1240  
1241 
$$= \|\delta_g\|^2 \varepsilon^{n+2} |V_1(x)| \int_{V_{\varepsilon}(x)} w(\tau) \le \|\delta_g\|^2 \varepsilon^{n+2} |V_1(x)| W_u$$

Let 
$$B = 2 \int_{V_{\tau}(x)} w(\tau) \delta_{0}^{T} \tau R_{x}(\tau) d\tau$$
  
 $B = 2 \int_{V_{\tau}(x)} w(\tau) \delta_{0}^{T} \tau R_{x}(\tau) d\tau \leq 2 \int_{V_{\tau}(x)} w(\tau) ||\delta_{y}|| \cdot ||\tau|| \cdot R_{x}(\tau) d\tau \leq 2 ||\delta_{y}|| \int_{V_{\tau}(x)} w(\tau) ||\tau|| \frac{1}{2} k_{y} \tau^{2} d\tau$   
 $\leq \frac{1}{2} k_{y} ||\delta_{y}|| \int_{V_{\tau}(x)} w(\tau) ||\tau||^{3} d\tau = k_{y} ||\delta_{y}||^{2} t^{n+3} |V_{1}(x)| W_{u}$   
To sum this up, we get  
 $\mathcal{L}(g_{x}(x)) \geq A - B \geq ||\delta_{y}||^{2} t^{n+2} |V_{1}(x)| W_{u} - k_{y} ||\delta_{y}|| t^{n+3} |V_{1}(x)| W_{u}$   
We combine the lower and upper-bound  
 $||\delta_{y}||^{2} t^{n+2} |V_{1}(x)| W_{u} - k_{y} ||\delta_{y}|| t^{n+3} |V_{1}(x)| W_{u} \leq (\frac{1}{2} k_{y})^{2} |V_{1}(x)| t^{n+4} ||\delta_{y}|^{2} t^{n+2} ||\delta_{y}||^{2} t^{n+2} \leq 0$   
 $||\delta_{y}|| \leq \frac{k_{y} t^{n} W_{u} + \sqrt{k_{y}^{2} t^{n} W_{u} + (\frac{1}{2} k_{y})^{2} t^{2}}{2W_{u}} = k_{y} t^{n} \frac{W_{u} + \sqrt{W_{u}^{2} + W_{u}}}{2W_{u}}$