

# EXPLICIT GRADIENT LEARNING FOR BLACK-BOX OPTIMIZATION

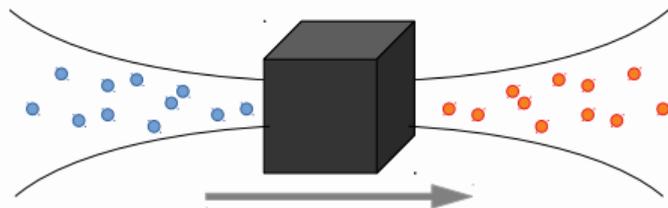
ELAD SARAFIAN\*, MOR SINAY\*,  
YORAM LOUZOUN, NOA AGMON AND SARIT KRAUS

BAR-ILAN UNIVERSITY



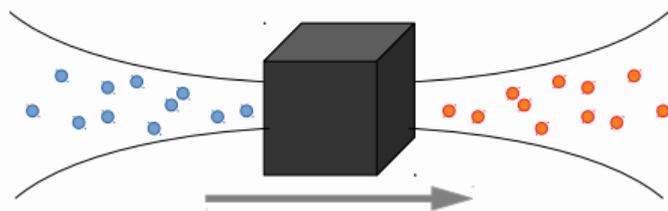
## Definition: A Black-Box Function

$f : \Omega \rightarrow \mathbb{R}$ ,  $\Omega \subseteq \mathbb{R}^n$  is a Black-Box function if one can sample  $y = f(x)$  at  $x \in \Omega$ , but has no prior knowledge of its analytical form.



## Definition: A Black-Box Function

$f : \Omega \rightarrow \mathbb{R}$ ,  $\Omega \subseteq \mathbb{R}^n$  is a Black-Box function if one can sample  $y = f(x)$  at  $x \in \Omega$ , but has no prior knowledge of its analytical form.



## Black-Box Optimization (BBO)

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

For  $C$  number of evaluation points, search for  $x^*$  and return the best candidate

$$\hat{x} = \arg \min_{x_i, i=1, \dots, C} f(x_i) \quad (2)$$

Machine Learning:

- **Hyperparameter tuning:** all the non-differential parameters of a learning algorithm.

Machine Learning:

- **Hyperparameter tuning:** all the non-differential parameters of a learning algorithm.
- **Reinforcement Learning:** Find the optimal parameters of a policy  $\pi_\theta : s \rightarrow a$  s.t. the expected utility function is maximized

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)]$$

In the real-world everything is a Black-Box...

# BLACK-BOX OPTIMIZATION: APPLICATIONS

Machine Learning:

- **Hyperparameter tuning:** all the non-differential parameters of a learning algorithm.
- **Reinforcement Learning:** Find the optimal parameters of a policy  $\pi_\theta : s \rightarrow a$  s.t. the expected utility function is maximized

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)]$$

In the real-world everything is a Black-Box...



# BLACK-BOX OPTIMIZATION: APPLICATIONS

Machine Learning:

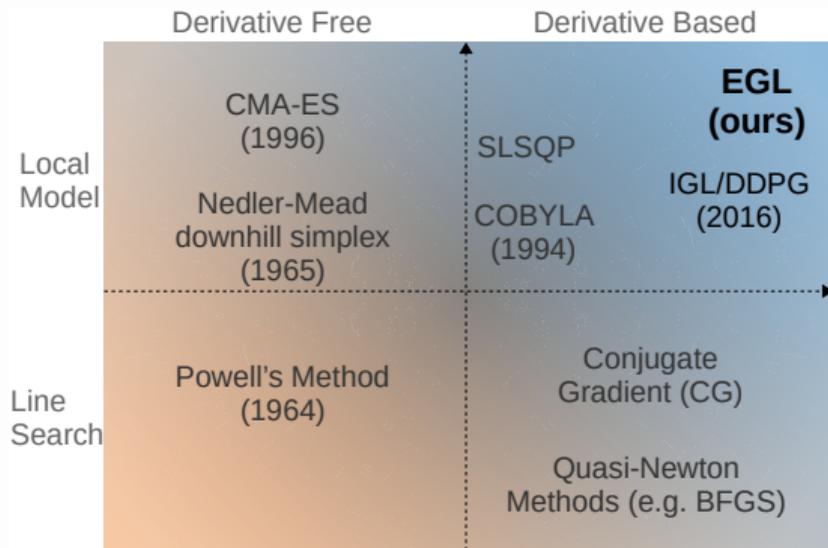
- **Hyperparameter tuning:** all the non-differentiable parameters of a learning algorithm.
- **Reinforcement Learning:** Find the optimal parameters of a policy  $\pi_\theta : s \rightarrow a$  s.t. the expected utility function is maximized

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)]$$

In the real-world, everything is a Black-Box.

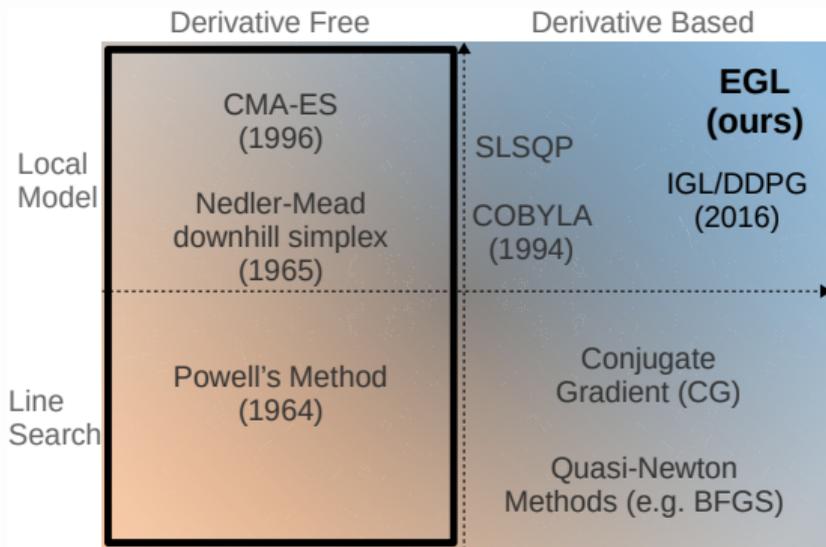


# BBO METHODS: TAXONOMY

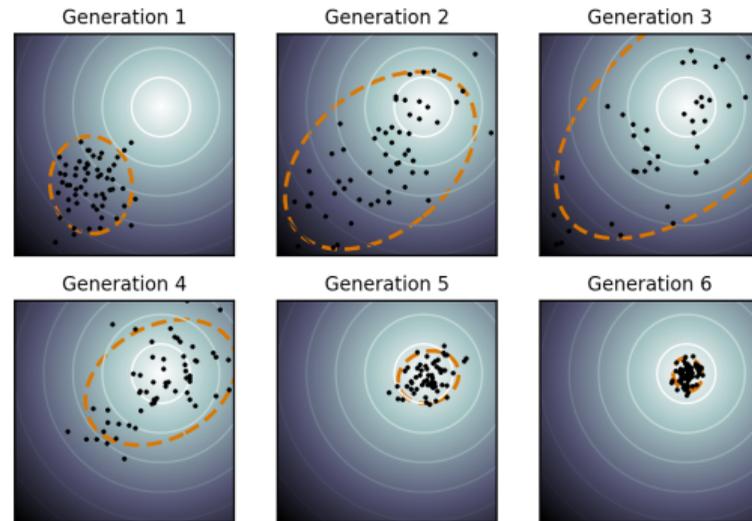


**The Black-Box Optimization Taxonomy**

# BBO METHODS: TAXONOMY

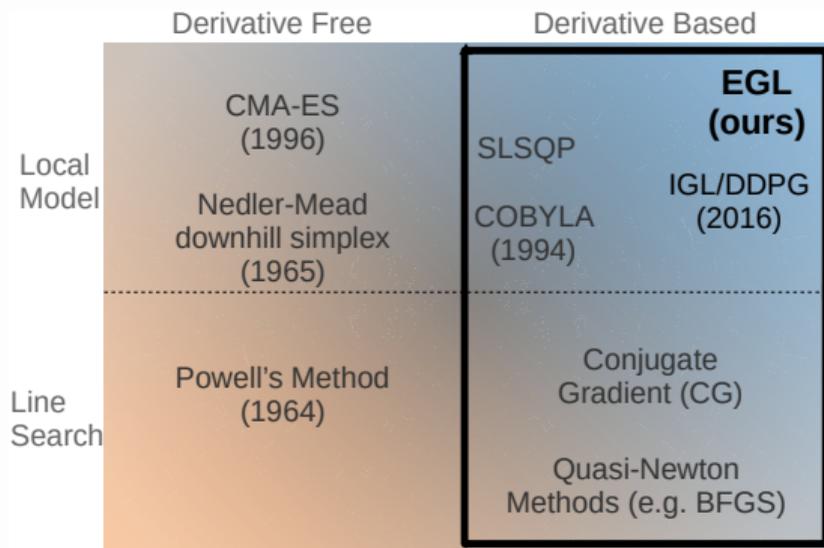


**The Black-Box Optimization Taxonomy**



**Figure 1: CMA-ES algorithm at work (from Wikipedia)**

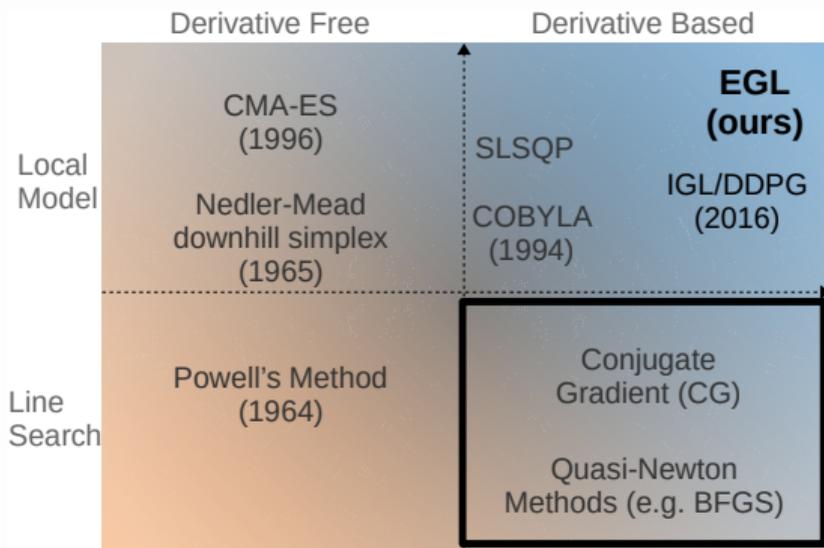
# BBO METHODS: TAXONOMY



**The Black-Box Optimization Taxonomy**

Assume a differentiable function  $f \in \mathcal{C}^1$  and use a gradient estimation to direct the search process.

# BBO METHODS: TAXONOMY

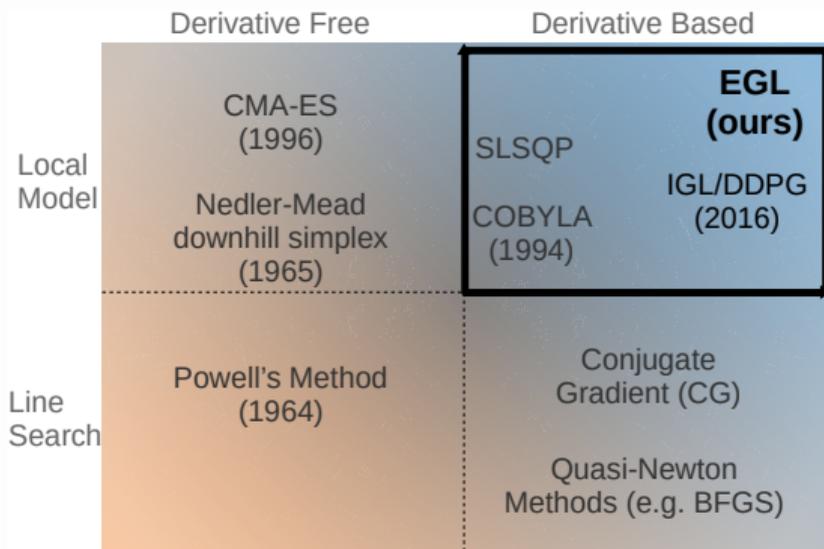


The Black-Box Optimization Taxonomy

In Line-Search methods, the directional derivative may be estimated by numerical methods, e.g.

$$n \cdot \nabla f(x) \approx \frac{f(x + \Delta n) - f(x - \Delta n)}{2\Delta}$$

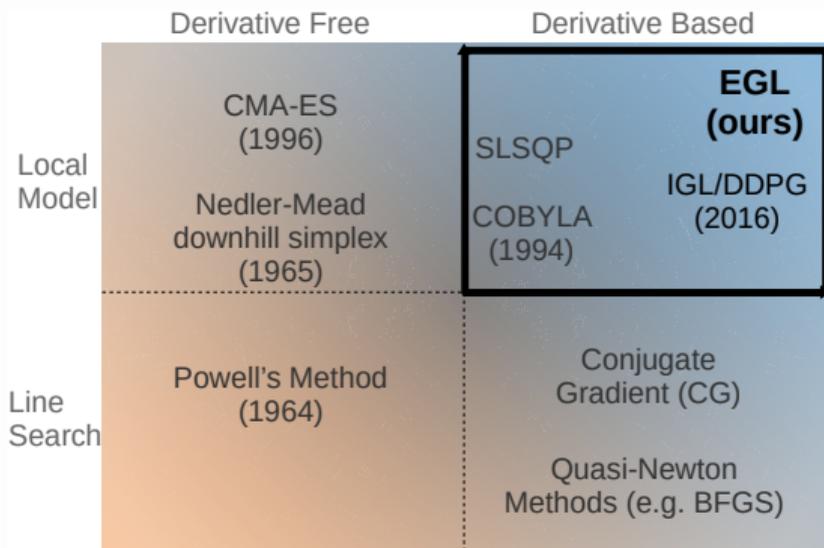
# BBO METHODS: TAXONOMY



**The Black-Box Optimization Taxonomy**

In Model-Based methods, the gradient can be estimated by fitting a parametric model  $f_\theta \approx f$  and following the parametric gradient  $\nabla f_\theta$

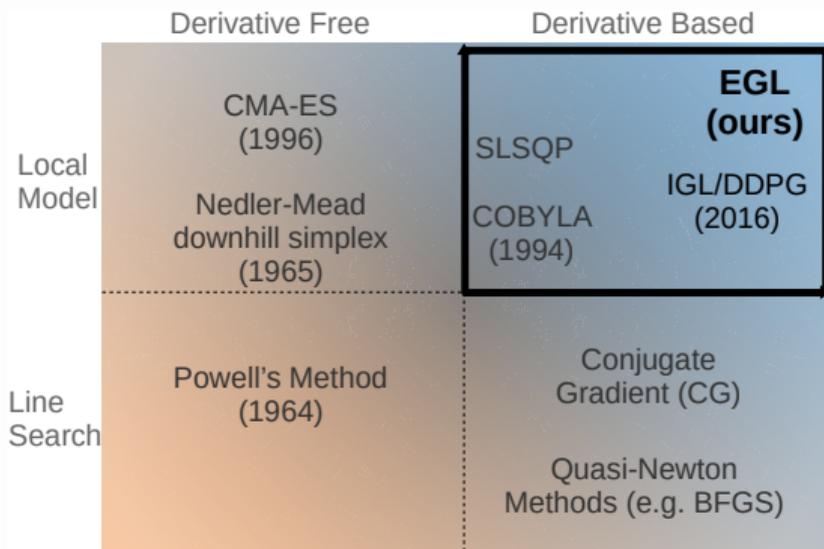
# BBO METHODS: TAXONOMY



**The Black-Box Optimization Taxonomy**

EGL takes Model-Derivative-Based methods forward:

# BBO METHODS: TAXONOMY

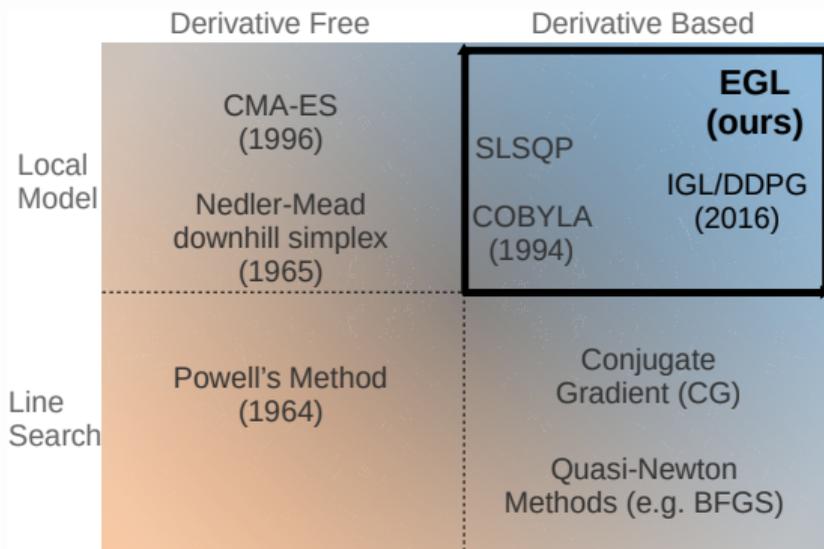


**The Black-Box Optimization Taxonomy**

EGL takes Model-Derivative-Based methods forward:

1. Instead of learning the function and obtain the parametric gradient. It **directly** fits a global model of the gradient from the data.

# BBO METHODS: TAXONOMY



**The Black-Box Optimization Taxonomy**

EGL takes Model-Derivative-Based methods forward:

1. Instead of learning the function and obtain the parametric gradient. It **directly** fits a global model of the gradient from the data.
2. It works with merely locally-integrable functions.

To develop EGL, let's take a closer look at Model-Based methods with Neural-Network parameterization.

# INDIRECT GRADIENT LEARNING (IGL)

Neural-Networks excel in fitting models to data  $f_{\theta} \approx f$

# INDIRECT GRADIENT LEARNING (IGL)

Neural-Networks excel in fitting models to data  $f_\theta \approx f$

1. Given a set of samples around a candidate  $x_k$ , minimize the MSE objective

$$\theta^* = \arg \min_{\theta} \sum_i \|y_i - f_\theta(x_i)\|^2$$

# INDIRECT GRADIENT LEARNING (IGL)

Neural-Networks excel in fitting models to data  $f_\theta \approx f$

1. Given a set of samples around a candidate  $x_k$ , minimize the MSE objective

$$\theta^* = \arg \min_{\theta} \sum_i \|y_i - f_\theta(x_i)\|^2$$

2. Approximate the gradient with the parametric gradient  $\nabla f \approx \nabla f_\theta$

# INDIRECT GRADIENT LEARNING (IGL)

Neural-Networks excel in fitting models to data  $f_\theta \approx f$

1. Given a set of samples around a candidate  $x_k$ , minimize the MSE objective

$$\theta^* = \arg \min_{\theta} \sum_i \|y_i - f_\theta(x_i)\|^2$$

2. Approximate the gradient with the parametric gradient  $\nabla f \approx \nabla f_\theta$
3. Take a gradient-descent step

$$x_{k+1} = x_k - \alpha \nabla f_\theta(x_k)$$

# INDIRECT GRADIENT LEARNING (IGL)

Neural-Networks excel in fitting models to data  $f_\theta \approx f$

1. Given a set of samples around a candidate  $x_k$ , minimize the MSE objective

$$\theta^* = \arg \min_{\theta} \sum_i \|y_i - f_\theta(x_i)\|^2$$

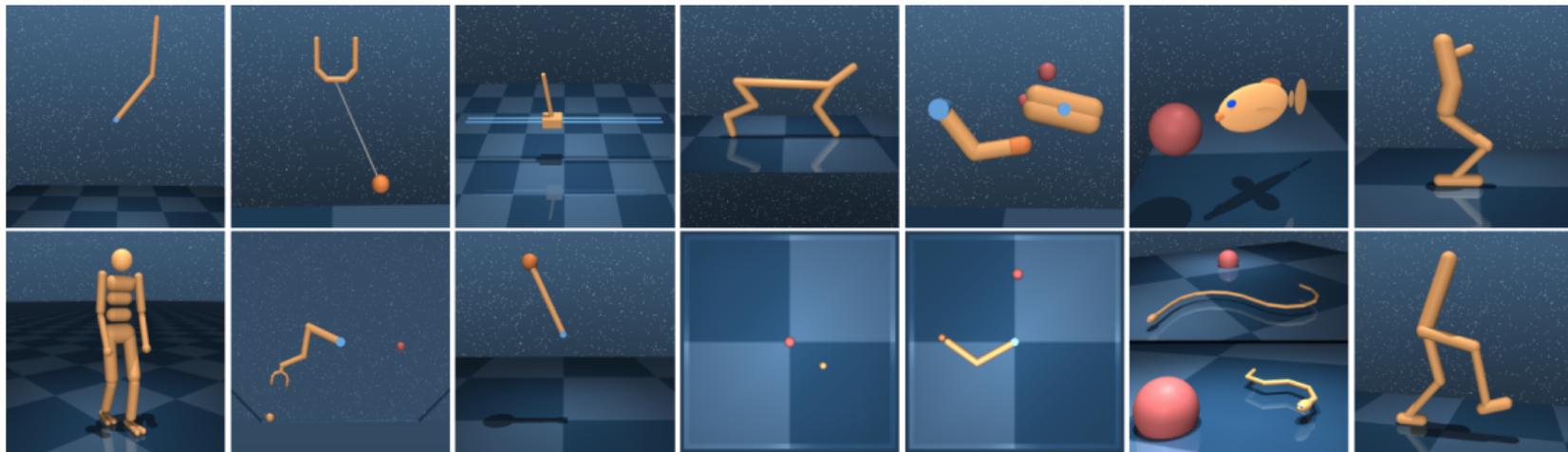
2. Approximate the gradient with the parametric gradient  $\nabla f \approx \nabla f_\theta$
3. Take a gradient-descent step

$$x_{k+1} = x_k - \alpha \nabla f_\theta(x_k)$$

4. Sample points around the new candidate and repeat.

# INDIRECT GRADIENT LEARNING (IGL)

Roots in the Deep Deterministic Policy Gradient (DDPG) seminal paper (2016). Have been applied successfully in robotics domains.



**Figure 2:** From the DeepMind Control suite Github

## Pickle I

The gradient is never explicitly learned, nor we obtain any guarantees for the accuracy of its estimation.

# DRAWBACKS OF IGL

## Pickle I

The gradient is never explicitly learned, nor we obtain any guarantees for the accuracy of its estimation.

## Pickle II

For Neural-Networks, the parametric gradient  $\nabla f_{\theta}$  may be discontinuous even if the objective is continuous.

# DRAWBACKS OF IGL

## Pickle I

The gradient is never explicitly learned, nor we obtain any guarantees for the accuracy of its estimation.

## Pickle II

For Neural-Networks, the parametric gradient  $\nabla f_{\theta}$  may be discontinuous even if the objective is continuous.

To overcome these drawbacks of IGL, we wish to learn the gradient explicitly.

# DRAWBACKS OF IGL

## Pickle I

The gradient is never explicitly learned, nor we obtain any guarantees for the accuracy of its estimation.

## Pickle II

For Neural-Networks, the parametric gradient  $\nabla f_{\theta}$  may be discontinuous even if the objective is continuous.

To overcome these drawbacks of IGL, we wish to learn the gradient explicitly.

## Problem

We cannot sample from  $\nabla f(x)$  directly.

# DRAWBACKS OF IGL

## Pickle I

The gradient is never explicitly learned, nor we obtain any guarantees for the accuracy of its estimation.

## Pickle II

For Neural-Networks, the parametric gradient  $\nabla f_\theta$  may be discontinuous even if the objective is continuous.

To overcome these drawbacks of IGL, we wish to learn the gradient explicitly.

## Problem

We cannot sample from  $\nabla f(x)$  directly.

## Solution: Learning a surrogate

Instead of learning  $\nabla f(x)$  we learn the *mean-gradient*  $g_\varepsilon(x)$ : averages over the gradient in a volume  $V_\varepsilon(x)$  s.t.  $\|x' - x\| \leq \varepsilon$  for all  $x' \in V_\varepsilon(x)$ .

# EXPLICIT GRADIENT LEARNING (EGL)

Recall the first order Taylor expression for differentiable functions

$$f(x + \tau) = f(x) + \nabla f(x) \cdot \tau + O(\|\tau\|^2).$$

# EXPLICIT GRADIENT LEARNING (EGL)

Recall the first order Taylor expression for differentiable functions

$$f(x + \tau) = f(x) + \nabla f(x) \cdot \tau + O(\|\tau\|^2).$$

## Definition: The Mean-Gradient

The mean-gradient at  $x$  with  $\varepsilon > 0$  averaging radius is

$$g_\varepsilon(x) = \arg \min_{g \in \mathbb{R}^n} \int_{V_\varepsilon(x)} |g \cdot \tau - f(x + \tau) + f(x)|^2 d\tau$$

where  $V_\varepsilon(x) \subset \mathbb{R}^n$  is a convex subset s.t.  $\|x' - x\| \leq \varepsilon$  for all  $x' \in V_\varepsilon(x)$  and the integral domain is over  $\tau$  s.t.  $x + \tau \in V_\varepsilon(x)$ .

$$g_\varepsilon(x) = \arg \min_{g \in \mathbb{R}^n} \int_{V_\varepsilon(x)} |g \cdot \tau - f(x + \tau) + f(x)|^2 d\tau$$

## Benefit I: Continuity

If  $f(x)$  is continuous in  $V$  s.t.  $V_\varepsilon(x) \subset V$  then the mean-gradient is a continuous function at  $x$ .

# CHARACTERISTICS OF THE MEAN-GRADIENT

$$g_\varepsilon(x) = \arg \min_{g \in \mathbb{R}^n} \int_{V_\varepsilon(x)} |g \cdot \tau - f(x + \tau) + f(x)|^2 d\tau$$

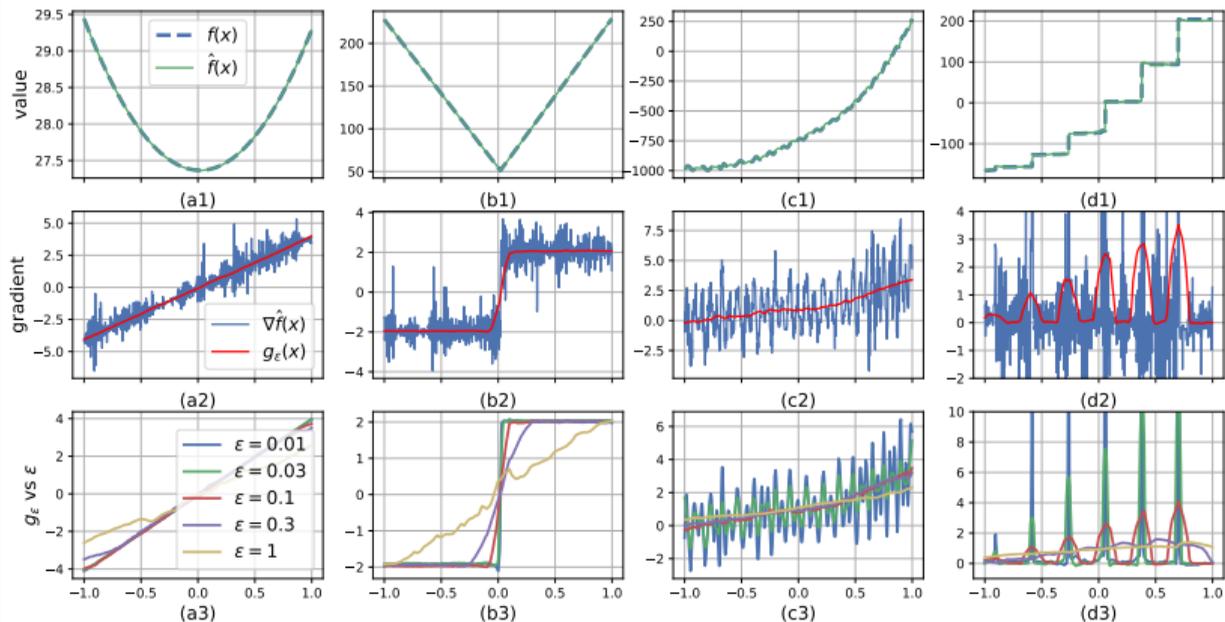
## Benefit I: Continuity

If  $f(x)$  is continuous in  $V$  s.t.  $V_\varepsilon(x) \subset V$  then the mean-gradient is a continuous function at  $x$ .

## Benefit II: Controllable Accuracy

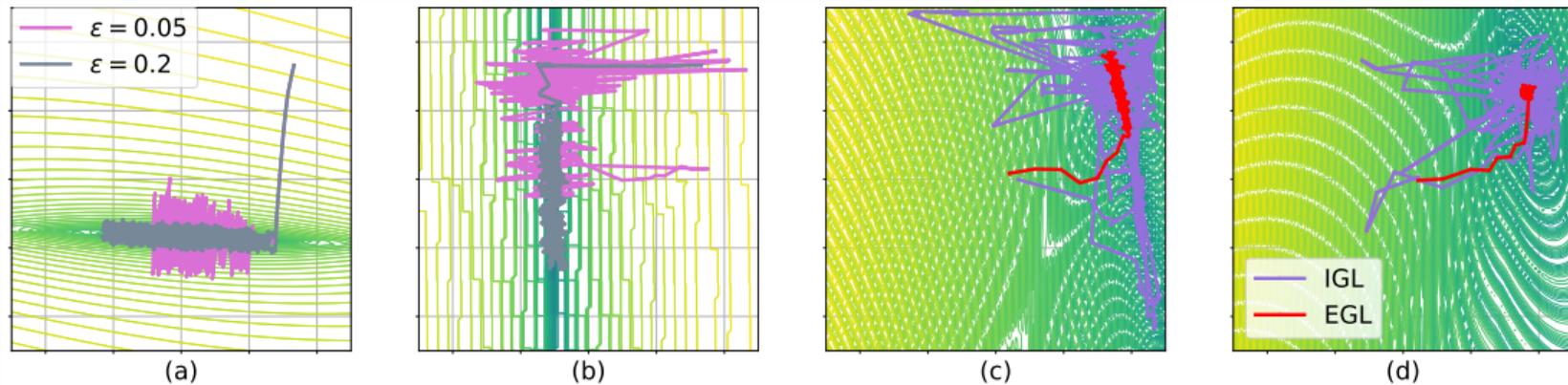
For any differentiable function  $f$  with a continuous gradient, there is  $\kappa_g > 0$ , so that for any  $\varepsilon > 0$  the mean-gradient satisfies  $\|g_\varepsilon(x) - \nabla f(x)\| \leq \kappa_g \varepsilon$  for all  $x \in \Omega$ .

# EXPLICIT GRADIENT LEARNING: EGL vs IGL



**Figure 3:** Comparing indirect gradient learning and explicit gradient learning for 4 typical functions: (a) parabolic; (b) piecewise linear; (c) multiple local minima; (d) step function.

# EXPLICIT GRADIENT LEARNING: EGL vs IGL



**Figure 4:** Visualizing EGL and IGL with different  $\epsilon$  for various 2D problems from COCO test suite.

$$g_\epsilon(x) = \arg \min_{g \in \mathbb{R}^n} \int_{V_\epsilon(x)} |g \cdot \tau - f(x + \tau) + f(x)|^2 d\tau$$

Learning the mean-gradient is done by minimizing the Monte-Carlo approximation of the Mean-Gradient:

Learning the mean-gradient is done by minimizing the Monte-Carlo approximation of the Mean-Gradient:

1. Sample set of pairs of observations  $\mathcal{D}_k = \{(x_i, y_i)\}_{i=1}^m$  s.t.  $x_i \in V_\varepsilon(x_k)$  where  $x_k$  is a candidate solution.

Learning the mean-gradient is done by minimizing the Monte-Carlo approximation of the Mean-Gradient:

1. Sample set of pairs of observations  $\mathcal{D}_k = \{(x_i, y_i)\}_{i=1}^m$  s.t.  $x_i \in V_\varepsilon(x_k)$  where  $x_k$  is a candidate solution.
2. Minimize the loss function

$$\mathcal{L}_{k,\varepsilon}(\theta) = \sum_{i=1}^m \sum_{x_j \in V_\varepsilon(x_i)} |(x_j - x_i) \cdot g_\theta(x_i) - y_j + y_i|^2 \quad (3)$$

$$g_\varepsilon(x) = \arg \min_{g \in \mathbb{R}^n} \int_{V_\varepsilon(x)} |g \cdot \tau - f(x + \tau) + f(x)|^2 d\tau$$

Learning the mean-gradient is done by minimizing the Monte-Carlo approximation of the Mean-Gradient:

1. Sample set of pairs of observations  $\mathcal{D}_k = \{(x_i, y_i)\}_{i=1}^m$  s.t.  $x_i \in V_\varepsilon(x_k)$  where  $x_k$  is a candidate solution.
2. Minimize the loss function

$$\mathcal{L}_{k,\varepsilon}(\theta) = \sum_{i=1}^m \sum_{x_j \in V_\varepsilon(x_i)} |(x_j - x_i) \cdot g_\theta(x_i) - y_j + y_i|^2 \quad (4)$$

## Theorem + Corollary

Given a proper set of samples (denoted as a poised set), any Lipschitz continuous Neural Network that optimizes Eq. (4) is a controllably accurate model.

$$\|\nabla f(x) - g_\theta(x)\| \leq \kappa_g \varepsilon$$

Proposition (Bertsekas (1999), Proposition 1.2.3)

Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $\kappa_f$ -smooth and bounded below. Let  $x_{k+1} = x_k - \alpha \nabla f(x_k)$  and  $\alpha \leq \frac{1}{\kappa_f}$ . Then  $\|\nabla f(x_k)\| \xrightarrow{k \rightarrow \infty} 0$ .

# EXPLICIT GRADIENT LEARNING: CONVERGENCE

Proposition (Bertsekas (1999), Proposition 1.2.3)

Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $\kappa_f$ -smooth and bounded below. Let  $x_{k+1} = x_k - \alpha \nabla f(x_k)$  and  $\alpha \leq \frac{1}{\kappa_f}$ . Then  $\|\nabla f(x_k)\| \xrightarrow{k \rightarrow \infty} 0$ .

We prove that EGL converges to a stationary point and to the global optimum if the problem is convex.

# EXPLICIT GRADIENT LEARNING: CONVERGENCE

## Proposition (Bertsekas (1999), Proposition 1.2.3)

Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $\kappa_f$ -smooth and bounded below. Let  $x_{k+1} = x_k - \alpha \nabla f(x_k)$  and  $\alpha \leq \frac{1}{\kappa_f}$ . Then  $\|\nabla f(x_k)\| \xrightarrow{k \rightarrow \infty} 0$ .

We prove that EGL converges to a stationary point and to the global optimum if the problem is convex.

## Convergence of EGL

Suppose a controllable mean-gradient model  $g_\varepsilon$  with error constant  $\kappa_g$ , the gradient descent iteration  $x_{k+1} = x_k - \alpha_k g_{\varepsilon_k}(x_k)$  with  $\alpha_k$  s.t.  $\frac{5\varepsilon_k}{\|\nabla f(x_k)\|} \leq \alpha_k \leq \min\left(\frac{1}{\kappa_g}, \frac{1}{\kappa_f}\right)$  guarantees:

1. Monotonically decreasing steps s.t.  $f(x_{k+1}) \leq f(x_k) - 2.25 \frac{\varepsilon_k^2}{\alpha}$ .
2.  $\|\nabla f(x_k)\| \xrightarrow{k \rightarrow \infty} 0$  for a proper choice of  $\varepsilon_k$ .

# EXPLICIT GRADIENT LEARNING: CONVERGENT ALGORITHM

---

## Algorithm 1: Convergent EGL

---

**Input:**  $x_0, \alpha, \varepsilon, \gamma_\alpha < 1, \gamma_\varepsilon < 1, \bar{\varepsilon}$

$k = 0$

**while**  $\varepsilon \leq \bar{\varepsilon}$  **do**

**Build Model:**

        Collect data  $\{(x_i, y_i)\}_1^m, x_i \in V_\varepsilon(x_k)$

        Learn a local model  $g_\varepsilon(x_k)$

**Gradient Descent:**

$x_{k+1} \leftarrow x_k - \alpha g_\varepsilon(x_k)$

**if**  $f(x_{k+1}) > f(x_k) - 2.25 \frac{\varepsilon^2}{\alpha}$  **then**

$\alpha \leftarrow \gamma_\alpha \alpha$

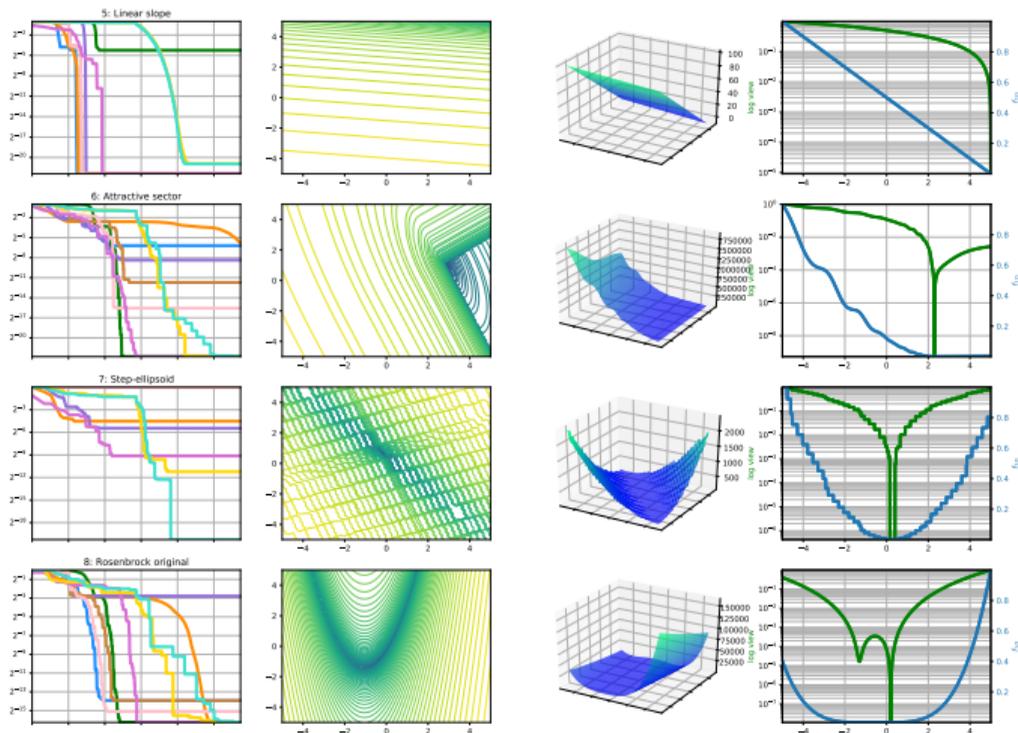
$\varepsilon \leftarrow \gamma_\alpha \gamma_\varepsilon \varepsilon$

$k \leftarrow k + 1$

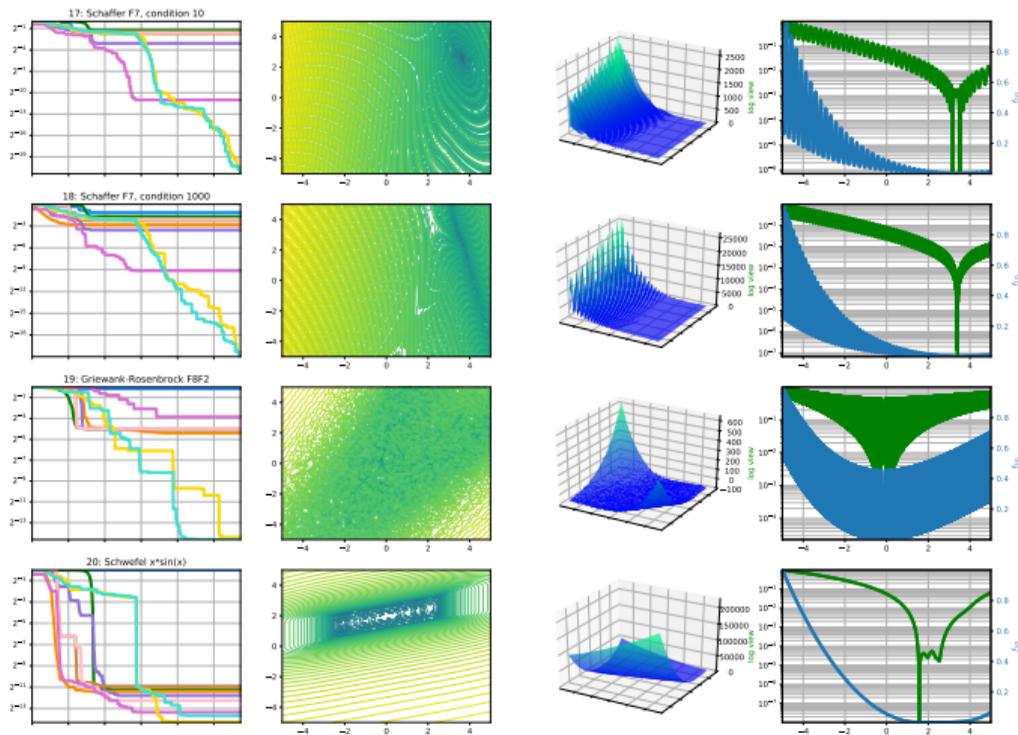
**return**  $x_k$

---

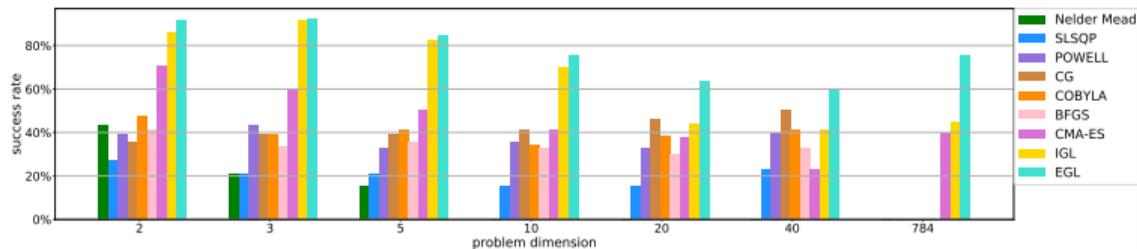
# EVALUATION IN THE COCO TEST SUITE



# EVALUATION IN THE COCO TEST SUITE

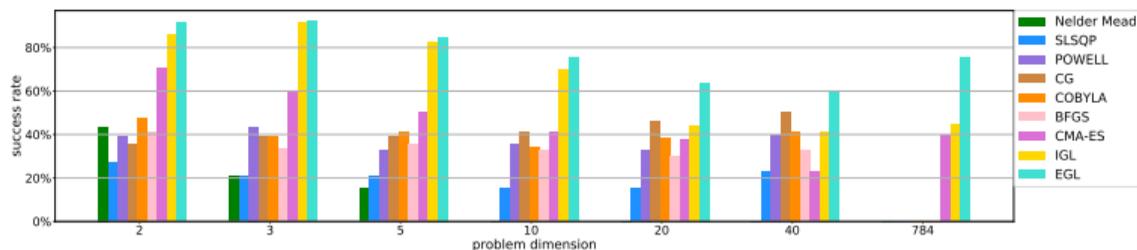


# EVALUATION IN THE COCO TEST SUITE

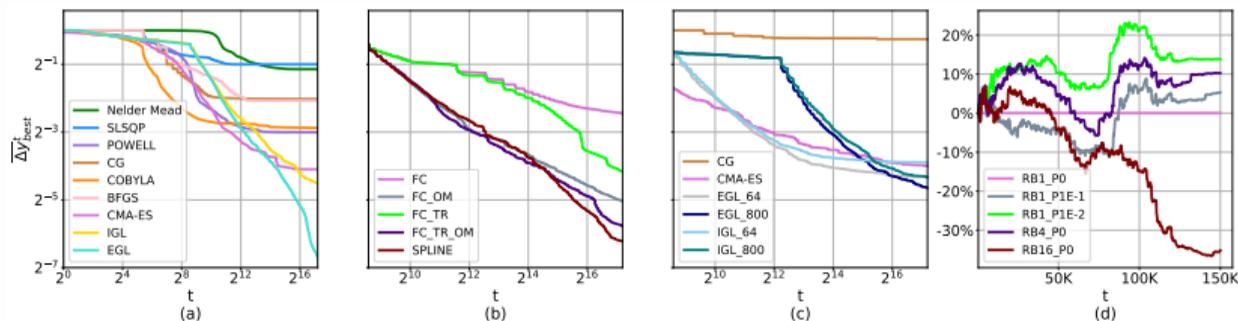


**Figure 5:** Comparing the success rate for a budget  $C = 150 \cdot 10^3$ .

# EVALUATION IN THE COCO TEST SUITE

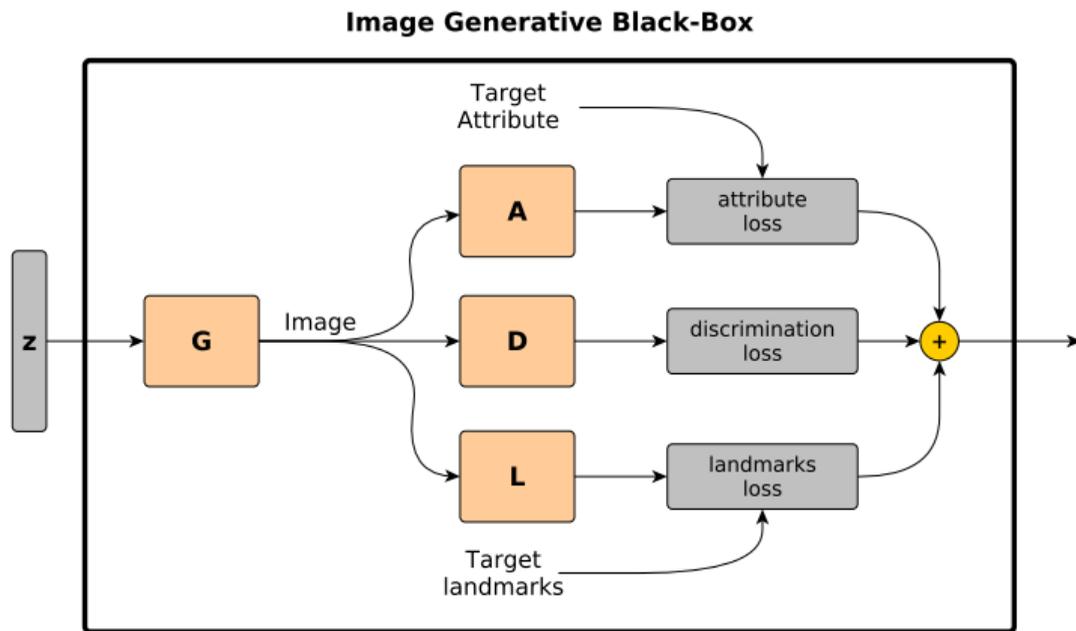


**Figure 5:** Comparing the success rate for a budget  $C = 150 \cdot 10^3$ .



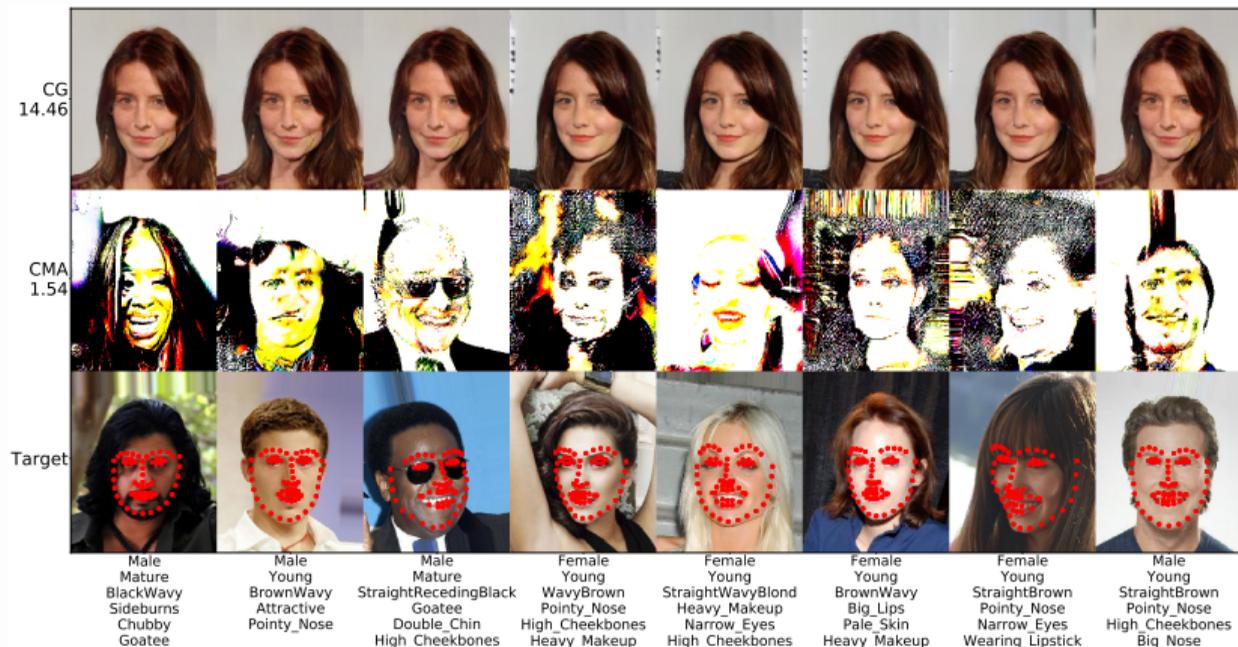
**Figure 6:** The scaled distance  $\Delta y^t_{best}$  as a function of  $t \in [1, \dots, C]$  for: (a) EGL and baselines on 40D, (b) dynamic mapping ablation test, (c) Different  $m$  samples on the 784D set.

# SEARCHING THE LATENT SPACE OF GENERATIVE MODELS



$$f_{al}(z) = \lambda_a \mathcal{L}_a(G(z)) + \lambda_l \mathcal{L}_l(G(z)) + \lambda_g \tanh(D(G(z)))$$

# SEARCHING THE LATENT SPACE OF GENERATIVE MODELS



# SEARCHING THE LATENT SPACE OF GENERATIVE MODELS



# SEARCHING THE LATENT SPACE OF GENERATIVE MODELS



- EGL is a model-based and derivative based BBO method.

- EGL is a model-based and derivative based BBO method.
- EGL can optimize non-convex and noisy functions.

- EGL is a model-based and derivative based BBO method.
- EGL can optimize non-convex and noisy functions.
- EGL converges to a local minimum.

- EGL is a model-based and derivative based BBO method.
- EGL can optimize non-convex and noisy functions.
- EGL converges to a local minimum.
- EGL outperforms existing methods both in a synthetic test-suite and real-world optimization application.

Thank you for your attention