

Low-Variance Gradient Estimation in Unrolled Computation Graphs with ES-Single

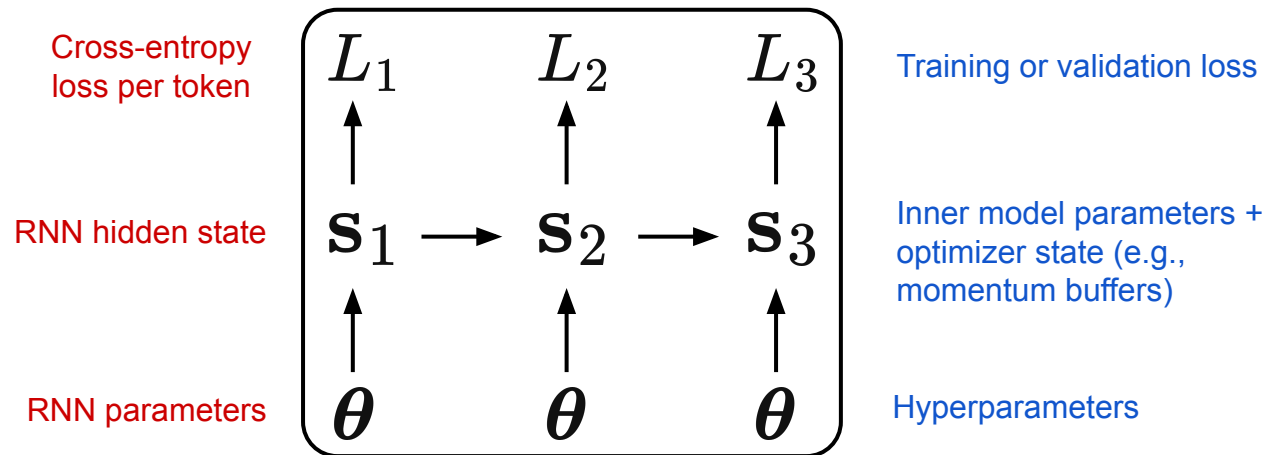
Paul Vicol, Zico Kolter, Kevin Swersky

Motivation

RNN Training

Unrolled Computation Graph

Hyperparameter Optimization



- For all these tasks, the objective is:
$$L(\theta) = \sum_{t=1}^T L_t(\mathbf{s}_t, \theta)$$

➡ We need the gradient $\nabla_{\theta} L(\theta)$

Existing Approaches

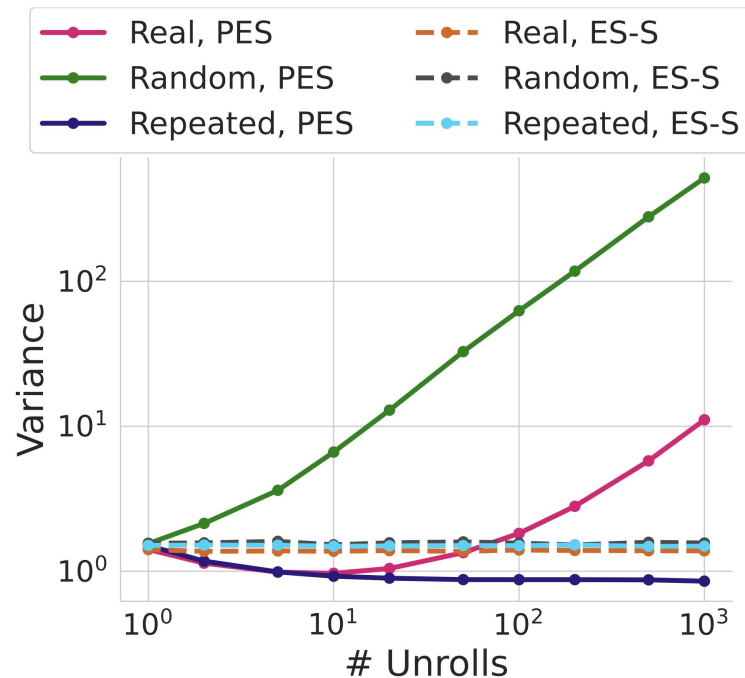
- Challenges with both short and long unrolls of the inner problem
 - **Short unrolls** → *truncation bias*
 - **Long unrolls** → *chaotic outer loss landscapes*
- *Evolution Strategies (ES)* computes an estimate of the gradient of a *smoothed loss*

$$g^{ES} = \frac{1}{\sigma^2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma^2 I)} [\epsilon L(\boldsymbol{\theta} + \epsilon)]$$

- *Smoothing overcomes chaos*
- But applying ES to full unrolls of the inner problem is expensive — slow updates
- Naively applying ES to truncated unrolls leads to bias
- *Persistent Evolution Strategies (PES)* computes unbiased gradient estimates using truncated unrolls

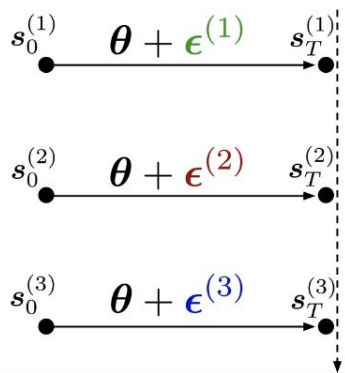
Variance Comparison

- PES is unbiased, but its *variance increases with the number of partial unrolls per inner problem*
- *We introduce ES-Single*: an algorithm for unbiased gradient estimation using partial unrolls
- Simpler and easier to implement than PES
- Has *constant variance* w.r.t. the number of partial unrolls per inner problem



ES-Single Computation Graph

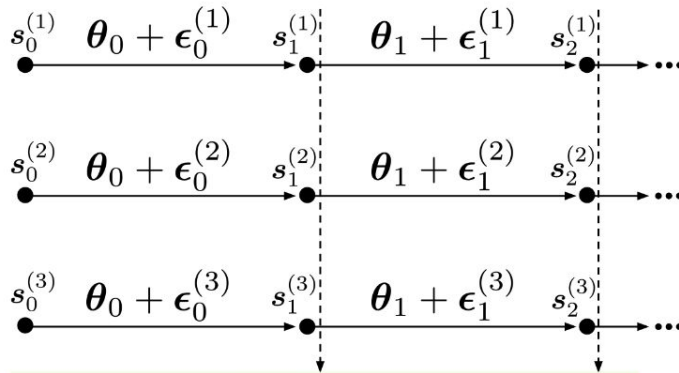
Full-Unroll ES



$$\theta \leftarrow \theta - \alpha \sum_{i=1}^N \epsilon_i \left(\sum_{t=1}^T L_t(\theta + \epsilon^{(i)}) \right)$$

Full-Unroll ES samples a perturbation for each particle, and *runs a full unroll for T steps using perturbed outer parameters*

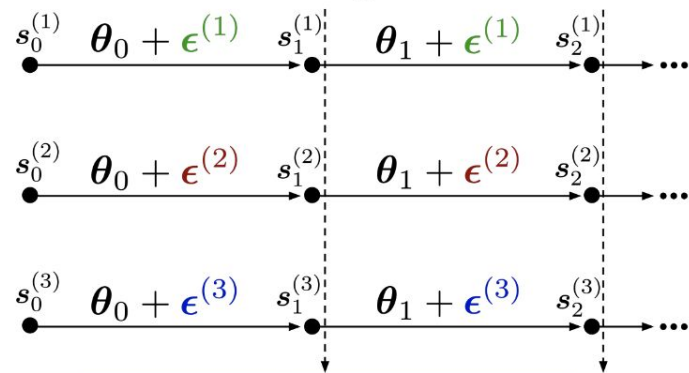
PES



$$\theta_{t+1} \leftarrow \theta_t - \alpha \sum_{i=1}^N \left(\sum_{\tau=1}^t \epsilon_{\tau}^{(i)} \right) L_t(s_t^{(i)}, \theta_t + \epsilon_t^{(i)})$$

PES samples a new perturbation for each particle in each unroll, and sums the perturbations experienced by each particle up to the current point in the inner problem

ES-Single



$$\theta_{t+1} \leftarrow \theta_t - \alpha \sum_{i=1}^N \epsilon^{(i)} L_t(s_t^{(i)}, \theta_t + \epsilon^{(i)})$$

ES-Single samples a single perturbation per particle at the start of each inner problem—keeping it fixed for the duration of the problem—and does not sum perturbations over time.

ES-Single Algorithm

Algorithm 1 Truncated Evolution Strategies (ES) applied to partial unrolls of a computation graph.

Input: s_0 , initial state
 K , truncation length for partial unrolls
 N , number of particles
 σ , standard deviation of perturbations
 α , learning rate for outer optimization

Initialize $s = s_0$

while inner problem not finished **do**

$$\hat{g}^{\text{ES}} \leftarrow \mathbf{0}$$

for $i = 1, \dots, N$ **do**

$$\epsilon^{(i)} = \begin{cases} \text{draw from } \mathcal{N}(0, \sigma^2 I) & i \text{ odd} \\ -\epsilon^{(i-1)} & i \text{ even} \end{cases}$$

$$\hat{L}_K^{(i)} \leftarrow \text{unroll}(s, \theta + \epsilon^{(i)}, K)$$

$$\hat{g}^{\text{ES}} \leftarrow \hat{g}^{\text{ES}} + \epsilon^{(i)} \hat{L}_K^{(i)}$$

end for

$$\hat{g}^{\text{ES}} \leftarrow \frac{1}{N\sigma^2} \hat{g}^{\text{ES}}$$

$$s \leftarrow \text{unroll}(s, \theta, K)$$

$$\theta \leftarrow \theta - \alpha \hat{g}^{\text{ES}}$$

end while

Algorithm 2 ES with a single perturbation per particle re-applied in each truncated unroll (ES-Single).

Input: s_0 , initial state
 K , truncation length for partial unrolls
 N , number of particles
 σ , standard deviation of perturbations
 α , learning rate for outer optimization

Initialize $s^{(i)} = s_0$ for $i \in \{1, \dots, N\}$

for $i = 1, \dots, N$ **do**

$$\epsilon^{(i)} = \begin{cases} \text{draw from } \mathcal{N}(0, \sigma^2 I) & i \text{ odd} \\ -\epsilon^{(i-1)} & i \text{ even} \end{cases}$$

end for

while inner problem not finished **do**

$$\hat{g}^{\text{ES-Single}} \leftarrow \mathbf{0}$$

for $i = 1, \dots, N$ **do**

$$s^{(i)}, \hat{L}_K^{(i)} \leftarrow \text{unroll}(s^{(i)}, \theta + \epsilon^{(i)}, K)$$

$$\hat{g}^{\text{ES-Single}} \leftarrow \hat{g}^{\text{ES-Single}} + \epsilon^{(i)} \hat{L}_K^{(i)}$$

end for

$$\hat{g}^{\text{ES-Single}} \leftarrow \frac{1}{N\sigma^2} \hat{g}^{\text{ES-Single}}$$

$$\theta \leftarrow \theta - \alpha \hat{g}^{\text{ES-Single}}$$

end while

ES-Single Properties

- ES-Single is *mathematically equivalent* to Full-Unroll ES, but *differs algorithmically*
 - ES-Single has the *same bias and variance characteristics as Full-Unroll ES*

Bias

Proposition 3.1 (ES-Single is unbiased). *Assume that $L(\boldsymbol{\theta})$ is quadratic and $\nabla_{\boldsymbol{\theta}}L(\boldsymbol{\theta})$ exists. Then, the ES-Single gradient estimator is unbiased, that is, $\text{bias}(\hat{\mathbf{g}}^{ES\text{-}Single}) = \mathbb{E}_{\epsilon} [\hat{\mathbf{g}}^{ES\text{-}Single}] - \nabla_{\boldsymbol{\theta}}L(\boldsymbol{\theta}) = 0$.*

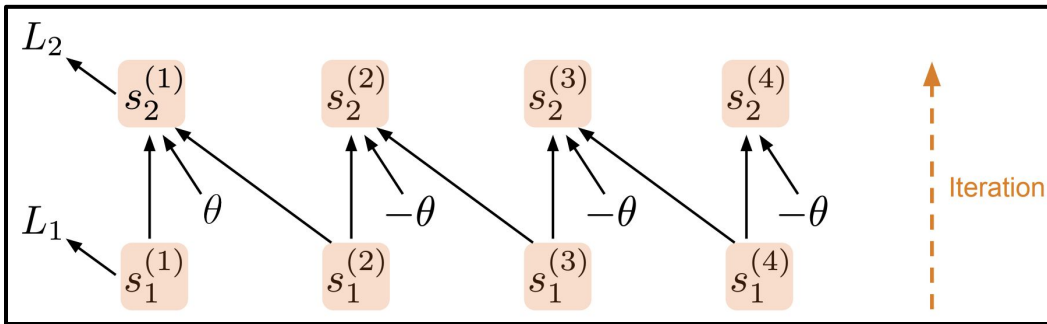
Proof. The proof is provided in Appendix D.1. \square

Variance

Proposition 3.2 (ES-Single Variance). *The total variance of ES-Single is $\text{tr}(\text{Var}(\hat{\mathbf{g}}^{ES\text{-}Single})) = (P + 1)\|\nabla_{\boldsymbol{\theta}}L(\boldsymbol{\theta})\|^2$, where P is the dimensionality of $\boldsymbol{\theta}$.*

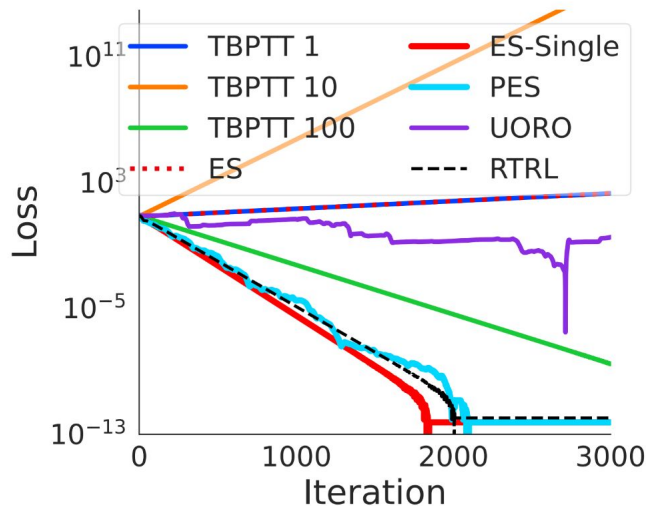
Proof. The proof is provided in Appendix D.2. \square

Influence Balancing: Task Setup

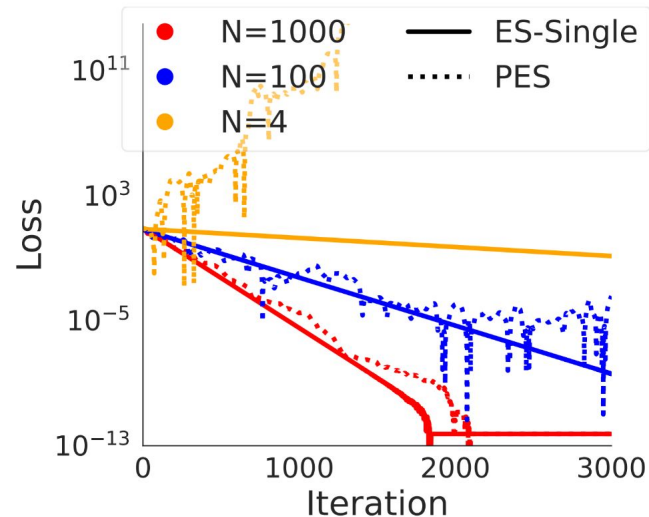


- **Task:** *Influence balancing*, introduced by Tallec & Olivier (2017)
 - Tune a scalar parameter θ that has a *negative influence in the short term*, but a *positive influence in the long term*
 - Designed such that *truncated methods move in the wrong direction*

Influence Balancing: Results



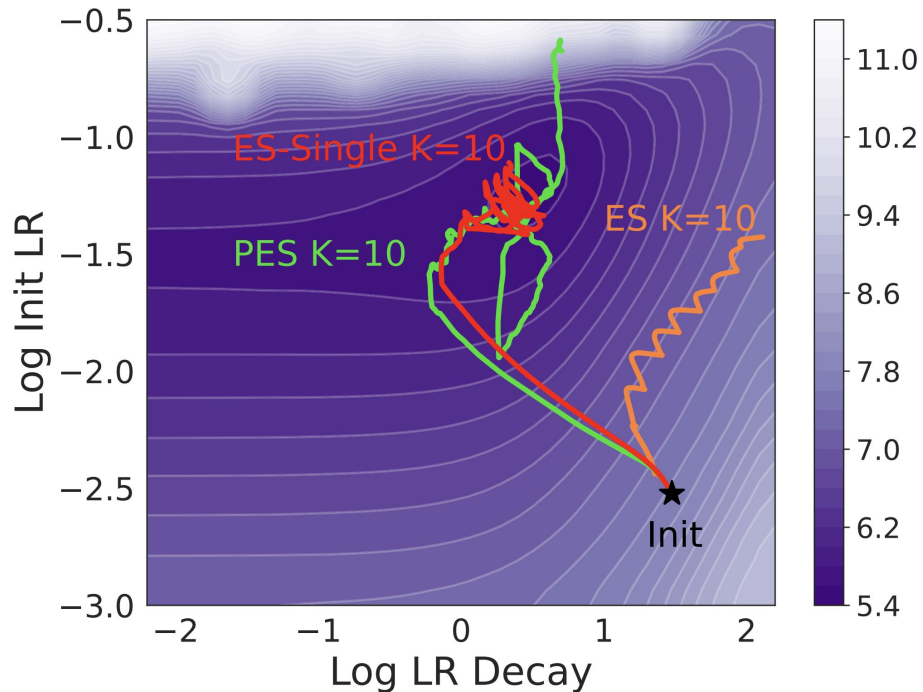
ES-Single is unbiased: it behaves like RTRL when using many particles



ES-Single is *more stable than PES* when using fewer particles

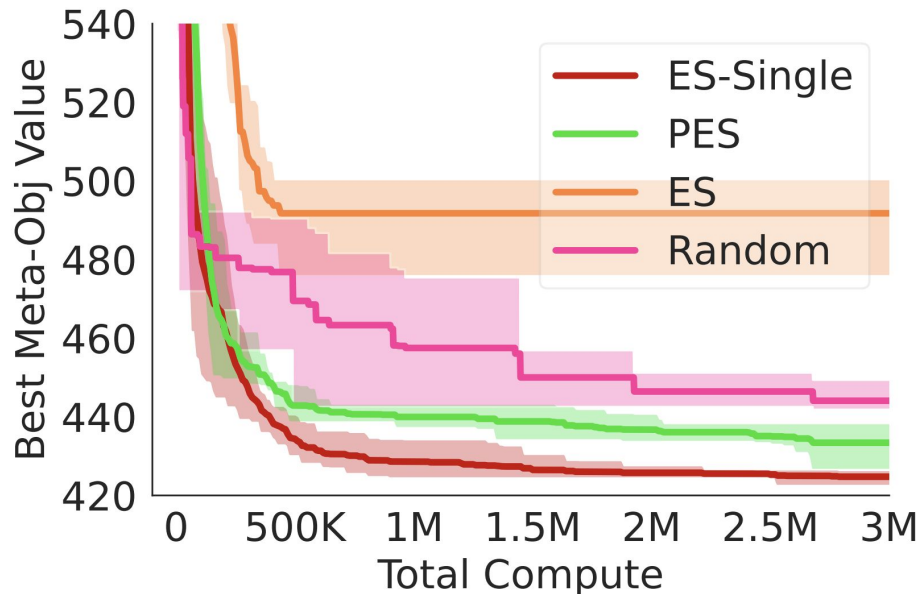
Meta-Optimizing an MNIST LR Schedule

- Meta-optimizing a learning rate schedule for training an MLP on MNIST
 - Tune the *initial learning rate and decay factor*
- ES-Single behaves similarly to PES, but has *more stable convergence at the optimum*



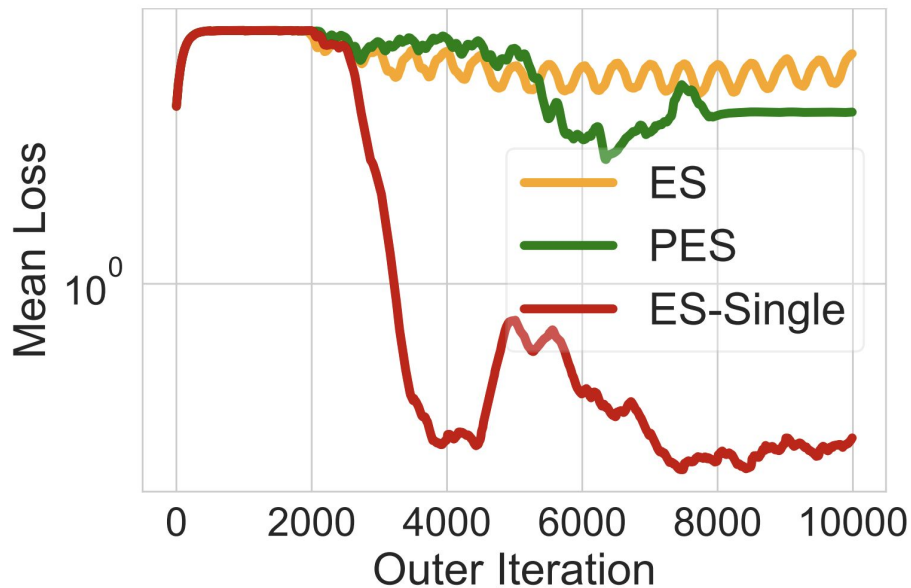
Optimizing Several Continuous & Discrete Hyperparams

- Training a 5-layer MLP on FashionMNIST
- *Optimizing 29 continuous and discrete hyperparameters*
 - Per-parameter block learning rates and momentum coefficients, and the number of hidden units per layer
- *ES-Single reaches lower meta-objective values using less total compute than truncated ES, random search, or PES*



Meta-Training a Learned Optimizer

- Meta-training a learned optimizer targeting the optimization of an MLP on FashionMNIST
- Here, $T=5000$ and $K=10$
- *ES-Single works in some scenarios where PES does not*



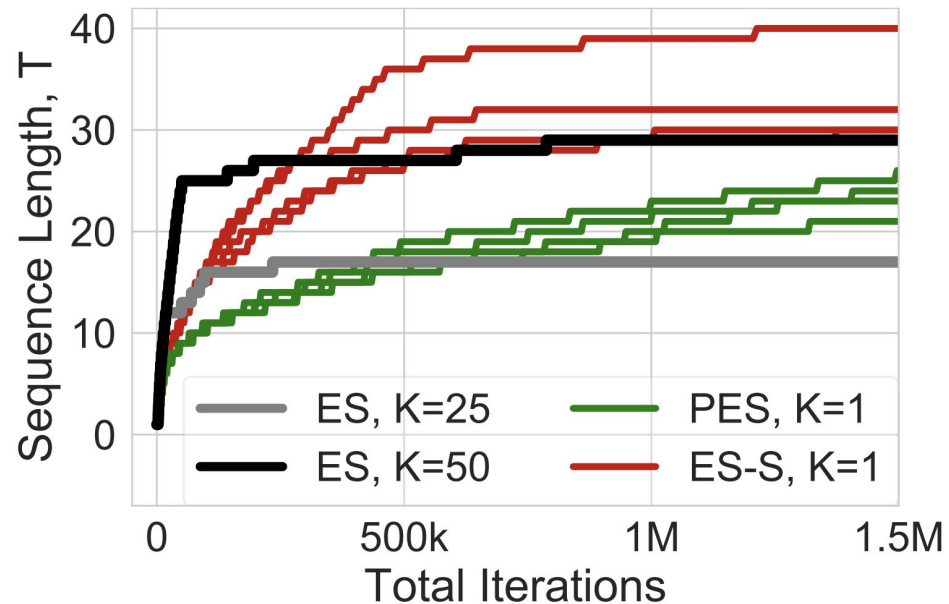
LSTM Copy Task

Copy Task

Input: 001101-----

Output: -----001101

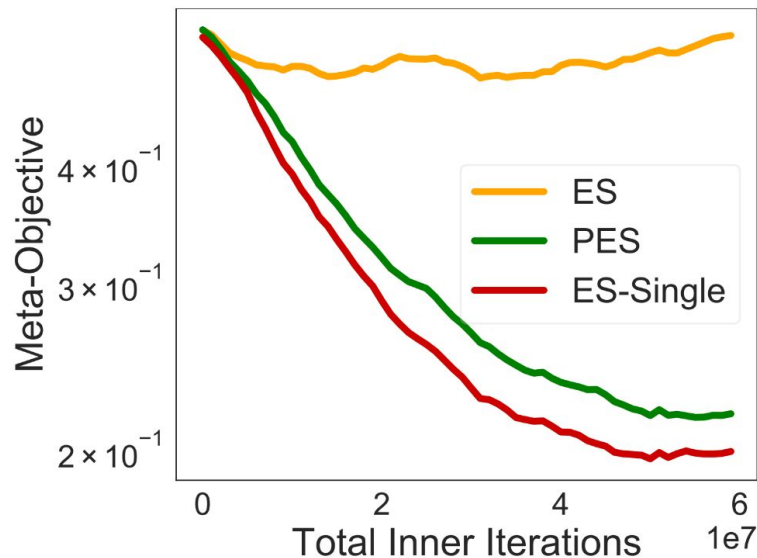
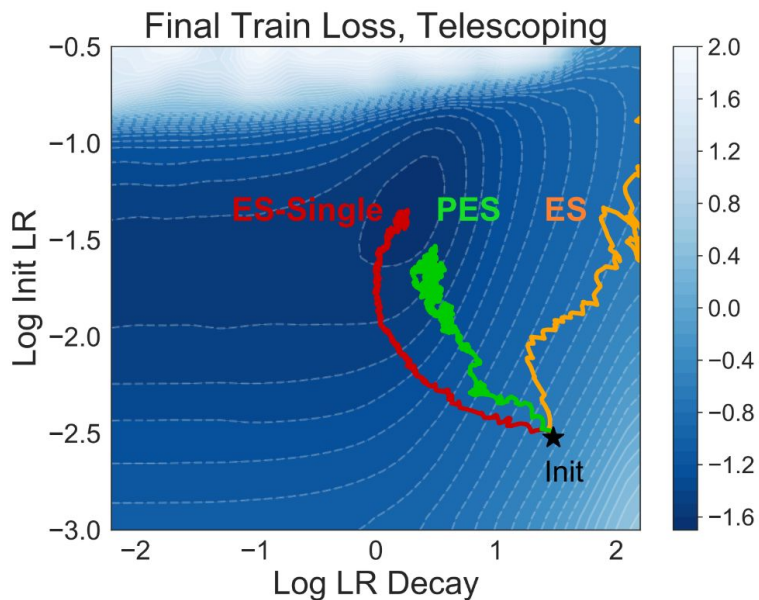
- **Challenge:** *Scaling to longer sequences*
- *Truncated methods* like TBPTT and truncated ES *fail to model long-term dependencies*
- PES with $K=1$ works, but is slow
- *ES-Single with $K=1$ is significantly faster than PES, and reaches larger T*



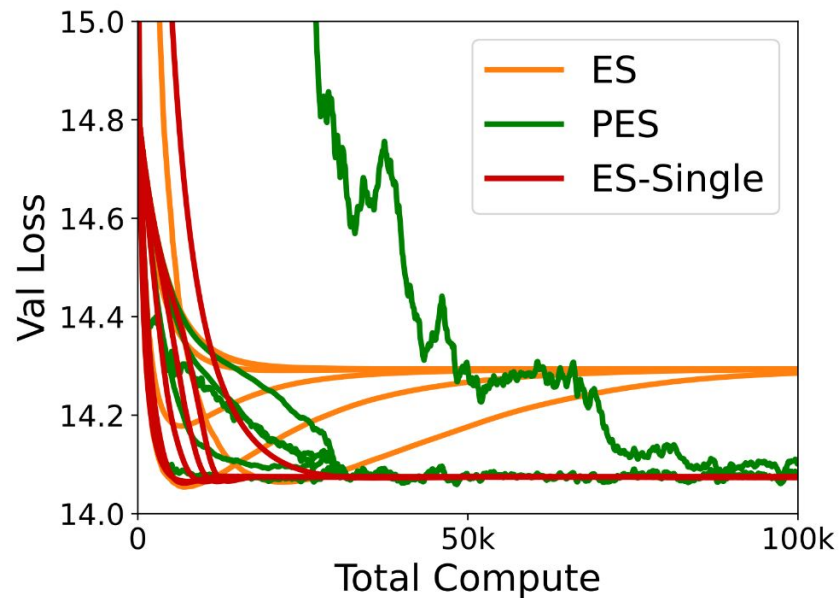
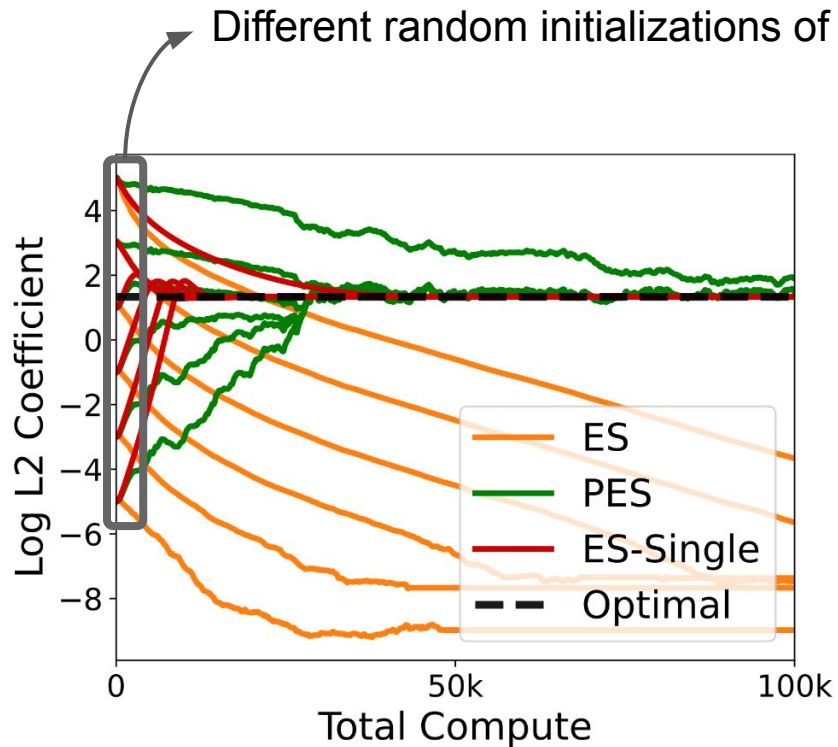
Telescoping Sums

- Can use telescoping sums to target the *final training loss*

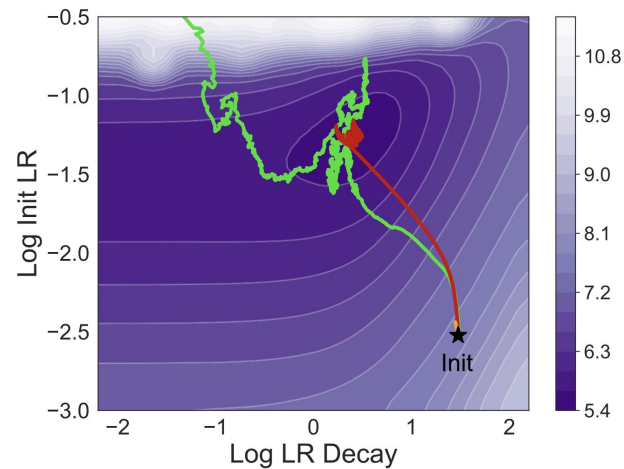
$$\sum_{t=0}^T p_t = (\cancel{L_0} - L_{-1}) + \dots + (L_T - \cancel{L_{T-1}}) = L_T$$



Tuning L_2 Regularization for UCI Regression



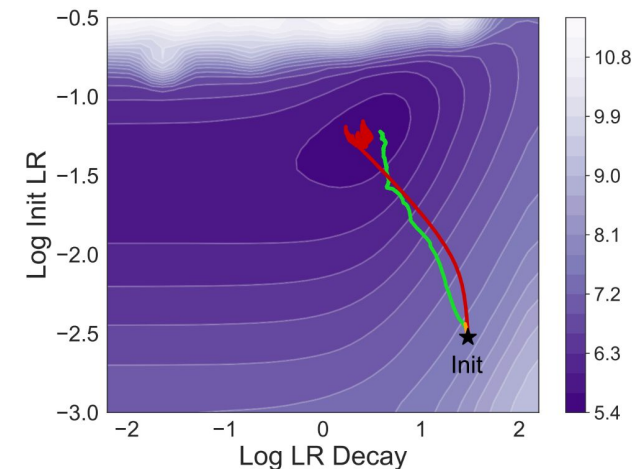
PES vs ES-Single: Stability and Performance



Large LR for PES



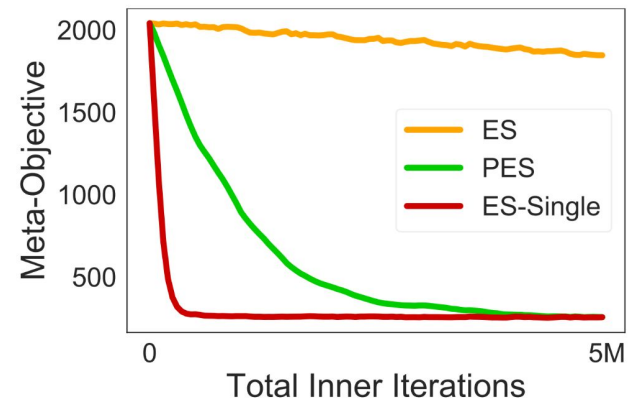
Diverges after initially getting to the optimal region



Small LR for PES

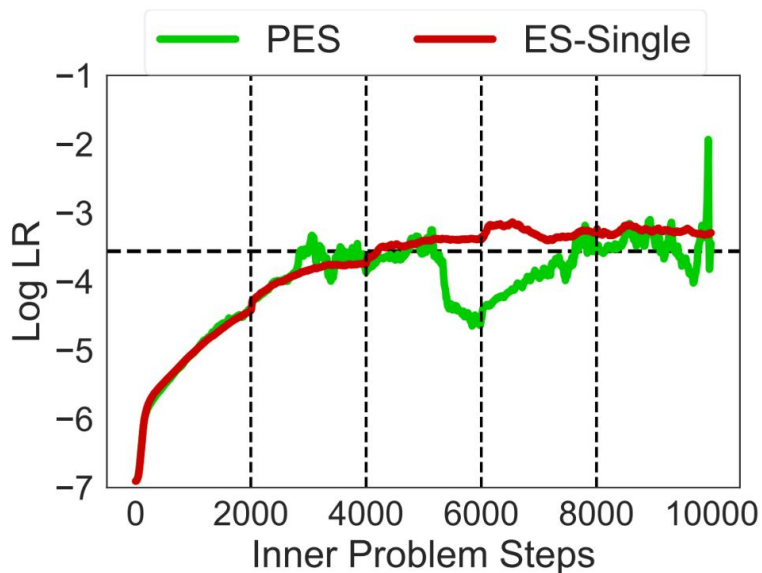


Stable convergence, but very slow compared to ES-Single

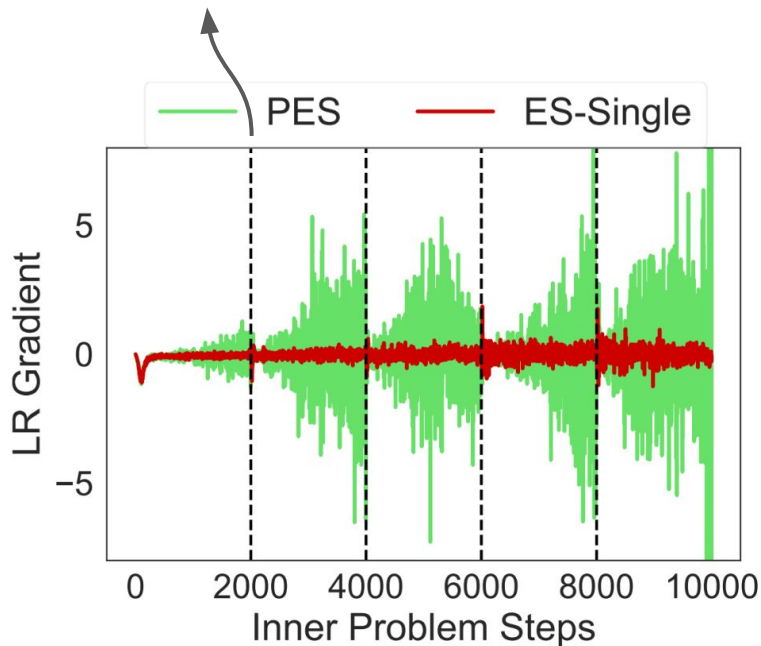


Meta-Gradient Comparison

Vertical lines denote the *start of a new inner problem* (T=2000)



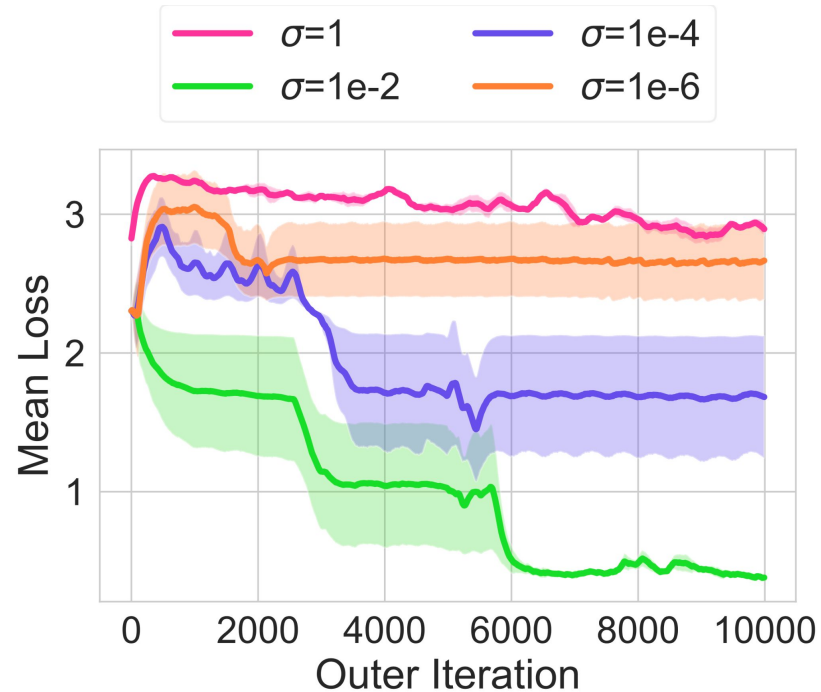
(a) Adaptation of the log-learning rate.



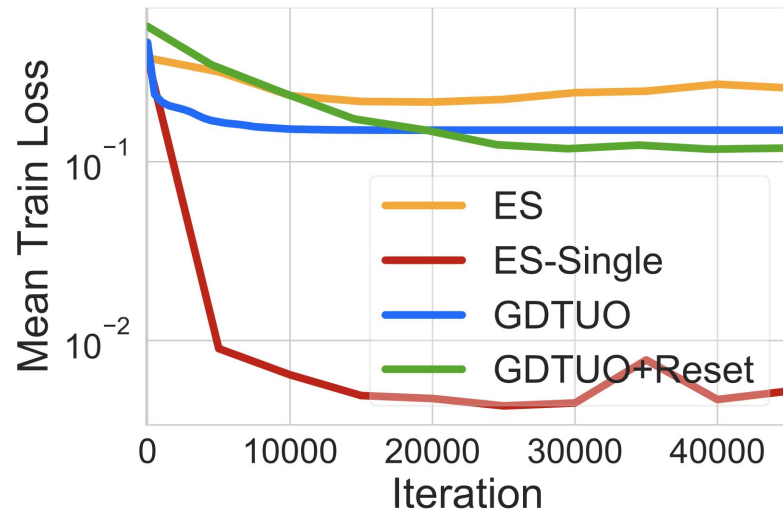
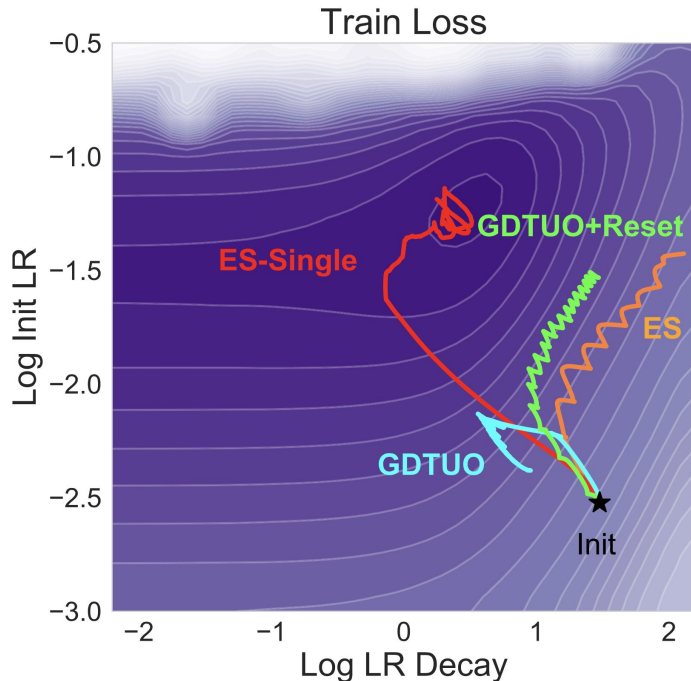
(b) PES and ES-Single meta-gradients over the course of multiple inner problems.

Importance of Smoothing

- Ablation over the perturbation scale σ , to optimize an MLP learned optimizer.
- *Small perturbation scales* lead to behavior similar to gradient-based methods, which may *get stuck in sub-optimal local minima in chaotic loss landscapes*.
 - For $\sigma = 1e-6$, meta-optimization fails to make progress
- In contrast, *using an appropriate scale, $\sigma=1e-2$, leads to stable convergence*

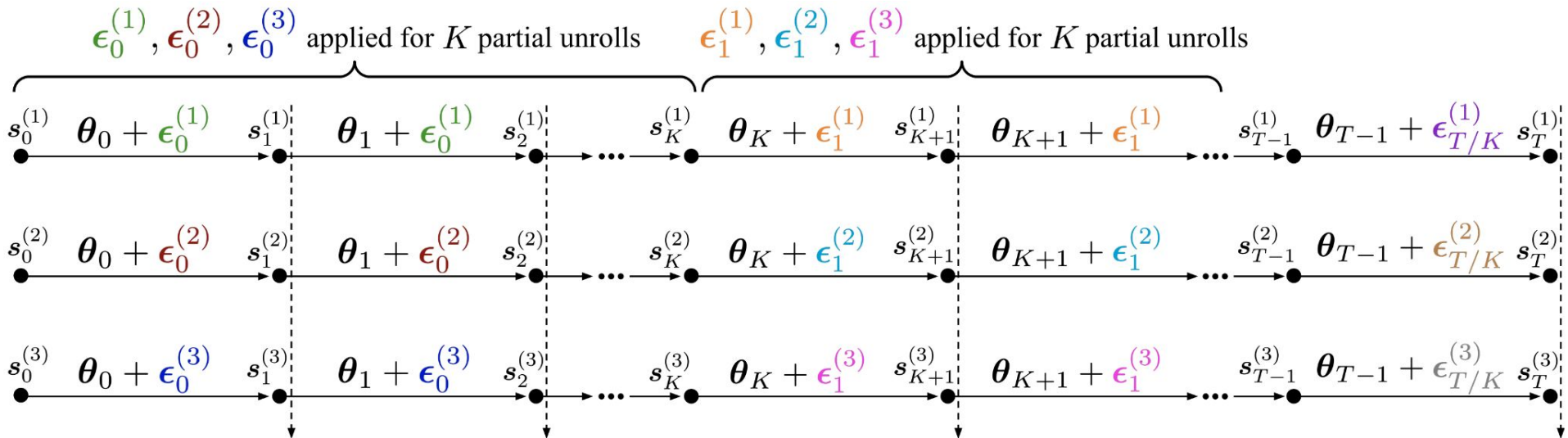


Comparison to a Gradient-Based Heuristic



- Here, we compared to *“Gradient Descent: The Ultimate Optimizer” (GDTUO)*
 - This is based on *Hypergradient Descent (HD)*, which adapts optimizer hyperparameters via a *1-step lookahead meta-objective*
- GDTUO is a gradient-based analogue to vanilla truncated ES, and behaves like ES

Generalization of PES and ES-Single



- This computation graph *generalizes ES-Single and PES*
- Uses the *same perturbation for K sequential truncated unrolls*
- After K unrolls, it adds the current perturbation to the perturbation accumulator, and samples a new perturbation

Generalization of PES and ES-Single: Algorithm

Algorithm 3 Truncated Evolution Strategies (ES) applied to partial unrolls of a computation graph.

Input: s_0 , initial state
 K , truncation length for partial unrolls
 N , number of particles
 σ , standard deviation of perturbations
 α , learning rate for outer optimization

Initialize $s = s_0$

while inner problem not finished **do**

$\hat{g}^{\text{ES}} \leftarrow \mathbf{0}$

for $i = 1, \dots, N$ **do**

$\epsilon^{(i)} = \begin{cases} \text{draw from } \mathcal{N}(0, \sigma^2 I) & i \text{ odd} \\ -\epsilon^{(i-1)} & i \text{ even} \end{cases}$

$\hat{L}_K^{(i)} \leftarrow \text{unroll}(s, \theta + \epsilon^{(i)}, K)$

$\hat{g}^{\text{ES}} \leftarrow \hat{g}^{\text{ES}} + \epsilon^{(i)} \hat{L}_K^{(i)}$

end for

$\hat{g}^{\text{ES}} \leftarrow \frac{1}{N\sigma^2} \hat{g}^{\text{ES}}$

$s \leftarrow \text{unroll}(s, \theta, K)$

$\theta \leftarrow \theta - \alpha \hat{g}^{\text{ES}}$

end while

Algorithm 4 Generalization of ES-Single and PES, with an arbitrary re-sampling interval M .

Input: s_0 , initial state
 K , truncation length for partial unrolls
 M , re-sampling interval
 N , number of particles
 σ , standard deviation of perturbations
 α , learning rate for outer optimization

Initialize $s^{(i)} = s_0$ for $i \in \{1, \dots, N\}$

Initialize $\xi^{(i)} \leftarrow \mathbf{0}$ for $i \in \{1, \dots, N\}$

while inner problem not finished, iteration j **do**

if $j \bmod M = 0$ **then**

for $i = 1, \dots, N$ **do**

$\epsilon^{(i)} = \begin{cases} \text{draw from } \mathcal{N}(0, \sigma^2 I) & i \text{ odd} \\ -\epsilon^{(i-1)} & i \text{ even} \end{cases}$

$\xi^{(i)} \leftarrow \xi^{(i)} + \epsilon^{(i)}$

end for

end if

$\hat{g}^{\text{ES-Gen}} \leftarrow \mathbf{0}$

for $i = 1, \dots, N$ **do**

$s^{(i)}, \hat{L}_K^{(i)} \leftarrow \text{unroll}(s^{(i)}, \theta + \epsilon^{(i)}, K)$

$\hat{g}^{\text{ES-Gen}} \leftarrow \hat{g}^{\text{ES-Gen}} + \xi^{(i)} \hat{L}_K^{(i)}$

end for

$\hat{g}^{\text{ES-Gen}} \leftarrow \frac{1}{N\sigma^2} \hat{g}^{\text{ES-Gen}}$

$\theta \leftarrow \theta - \alpha \hat{g}^{\text{ES-Gen}}$

end while

Conclusion

- ES-Single is a simple method for gradient estimation in unrolled computation graphs
- **Key idea:** sample a *single perturbation per particle at the start of each inner problem, and keep it fixed over all partial unrolls of the problem*
- ES-Single has *constant variance with respect to the number of partial unrolls* per inner problem
 - Addresses a key challenge faced by PES, and makes it scalable to long inner problems with short unrolls
- Empirically, *ES-Single outperforms PES on a range of tasks*, including hyperparameter optimization and RNN training



<https://colab.research.google.com/drive/1fqSzwaIXfJKbYTntEFfNbc2UuTXwEw0A?usp=sharing>

Thank you!