Variance-reduced Gradient Estimation via Noise-Reuse in Online Evolution Strategies Google DeepMind

(Vicol et al, 2021)

1 sample / W unrolls

sampled so far.

1 sample / K unrolls

Oscar Li, James Harrison, Jascha Sohl-Dickstein, Virginia Smith, Luke Metz

Come here to learn about

a new online gradient estimator

for unrolled computation graphs!

 $x_{t+1} = x_t + a(y_t - x_t)dt$

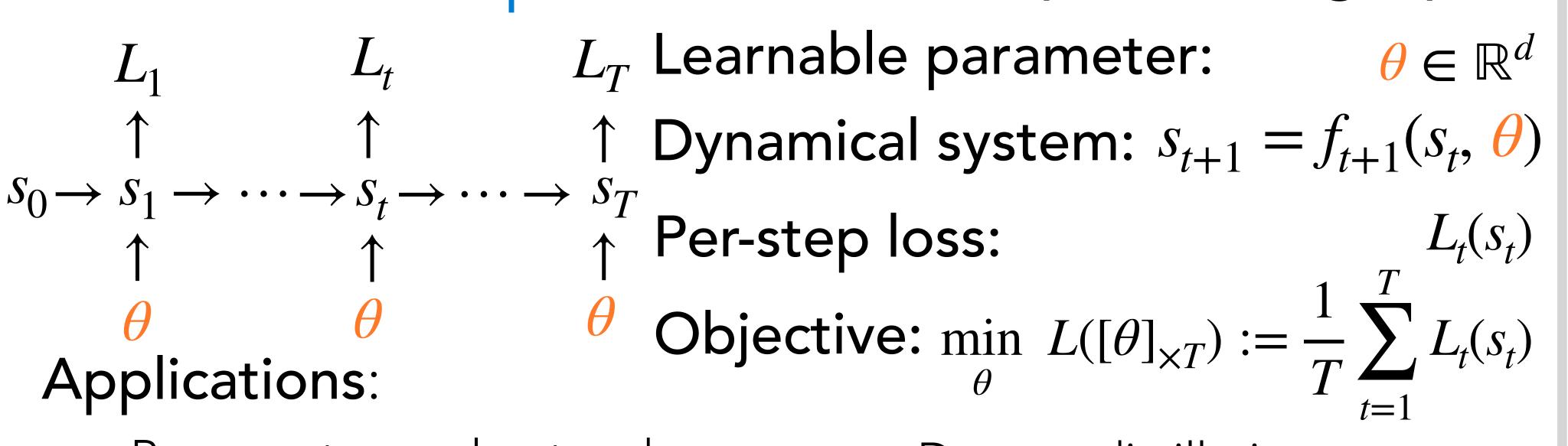
 $y_{t+1} = y_t + [x_t \cdot (r - z_t) - y_t]dt$

 $z_{t+1} = z_t + [x_t \cdot y_t - 8/3 \cdot z_t]dt$

PES W = 100, N = 20000

NRES W = 100, N = 20000

Unrolled computation graphs Online Evolutionary Strategies 1 careable parameter: 0 = md 1 update / W unrolls



- Recurrent neural networks
- Dataset distillation
- Meta-training learned optimizers
 Policy learning in RL
- Q: How best to compute gradients to optimize $L([\theta]_{\times T})$?

(when loss.backward() doesn't work) Loss properties

- Losses with extreme local sensitivity
- non-differentiable losses
 - r black box losses
- e.g., Model-free RL
- e.g., L is the zero-one loss discontinuous losses

Evolution Strategies (FullES)

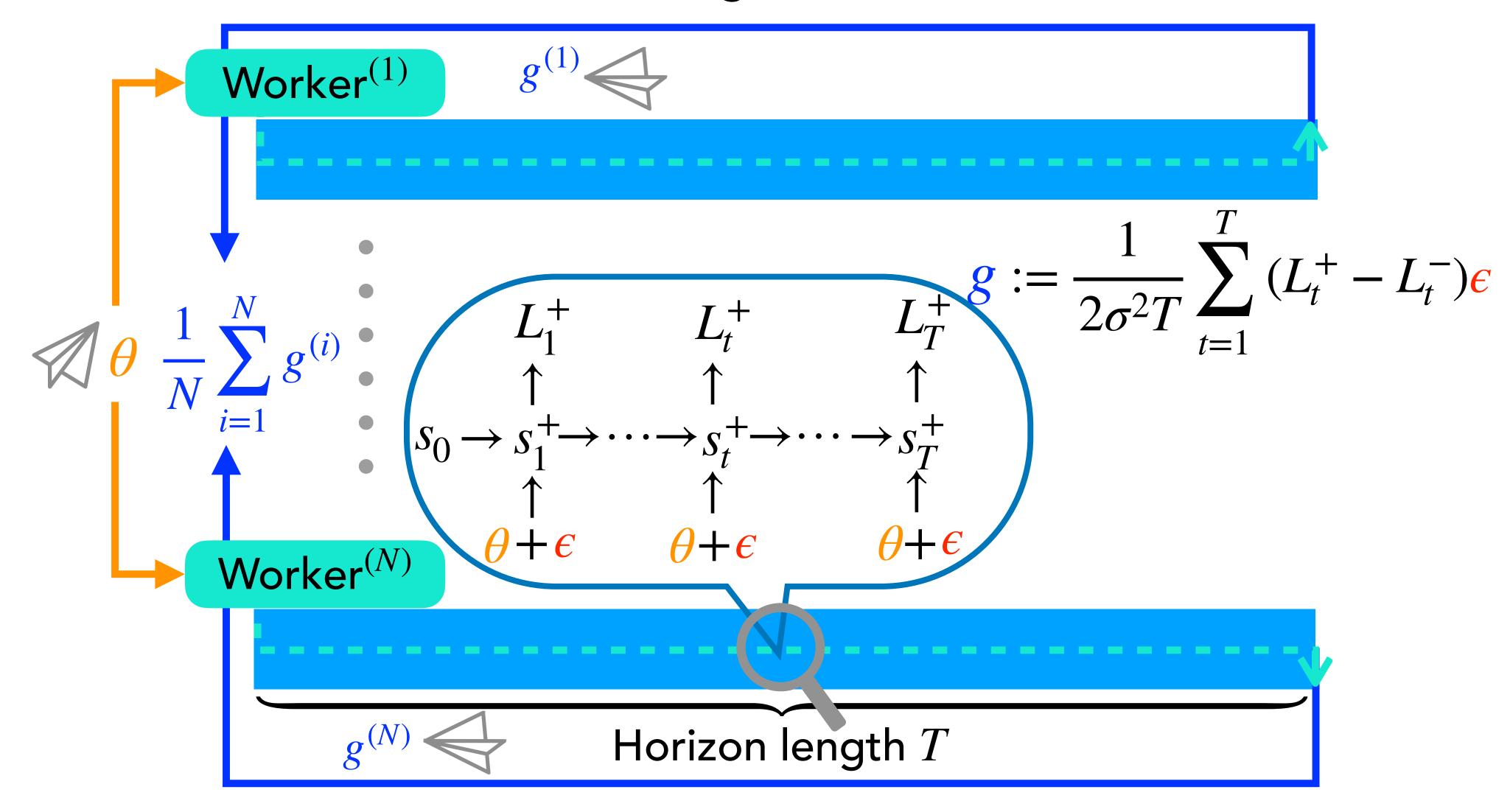
FullES smoothing objective:

$$\min_{\alpha} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I_{d \times d})} L([\theta + \epsilon]_{\times T})$$

FullES gradient estimator: θ $\frac{1}{2\sigma^2} \left[L([\theta + \epsilon]_{\times T}) - L([\theta - \epsilon]_{\times T}) \right] \epsilon \text{, where } \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I_{d \times d})$

1 update / T unrolls

FullES Training Protocol



PES smoothing objective:

Generalized PES (ours)

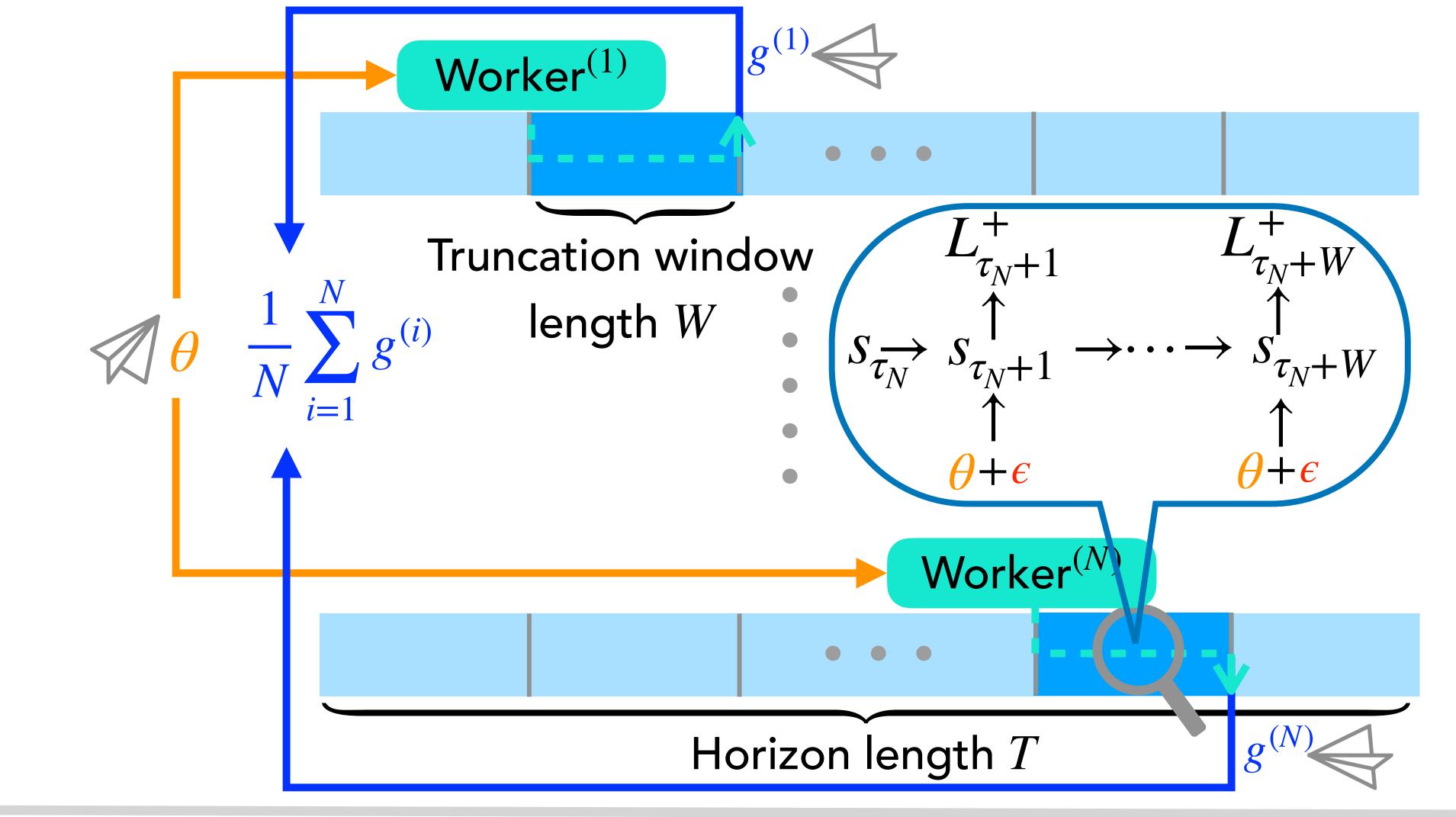
GPES smoothing objective:

Accumulate

Accumulate

d noise

Online Evolution Strategies Training Protocol



 $\min_{\theta} \mathbb{E}_{\epsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \sigma^2 I_{d \times d})} L([\theta + \epsilon_1]_{\times W}, ..., [\theta + \epsilon_{T/W}]_{\times W})$

PESWorker

 $g := \frac{1}{2\sigma^2 W} \left(\sum_{t} (L_t^+ - L_t^-) \right) \sum_{t} \epsilon_i$

 $\min_{\theta} \mathbb{E}_{\epsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \sigma^2 I_{d \times d})} L([\theta + \epsilon_1]_{\times K}, \dots, [\theta + \epsilon_{\lceil T/K \rceil}]_{\times r(T,K)})$

 $\mathsf{GPES}_{K=2W}\mathsf{Worker}$

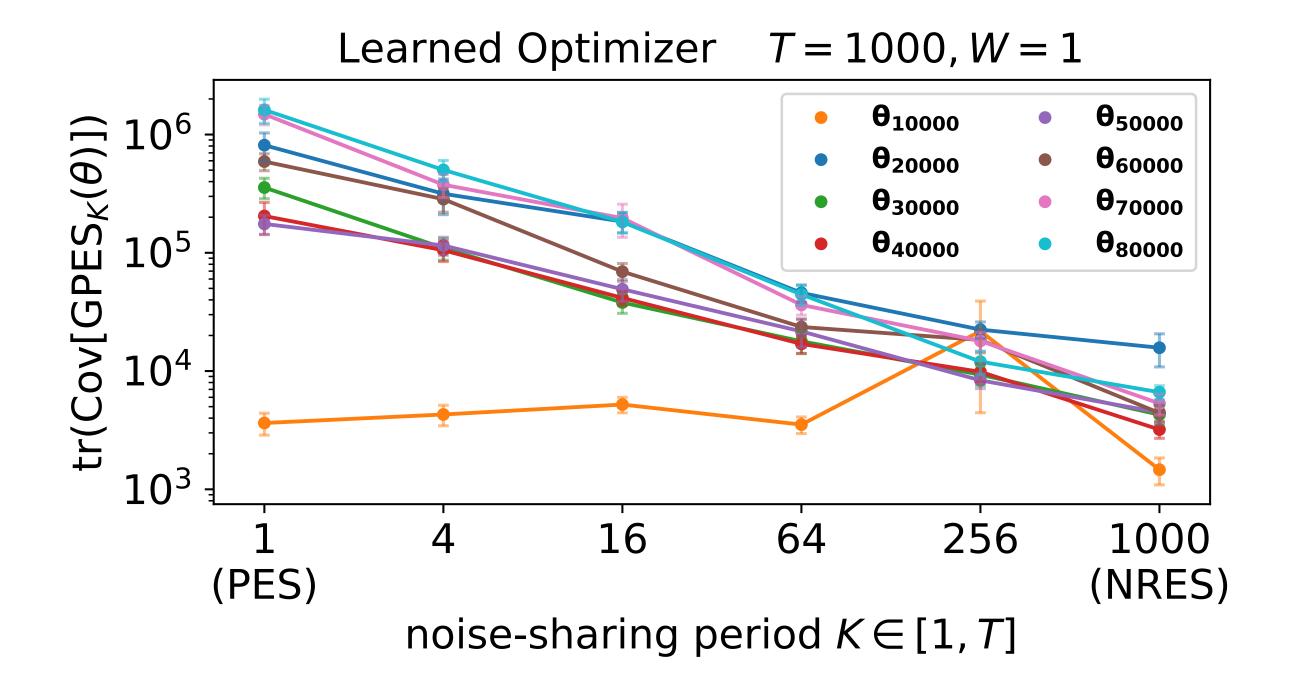
K divides au

Q: Which K to choose for $GPES_K$?

Persistent Evolutionary Strategies

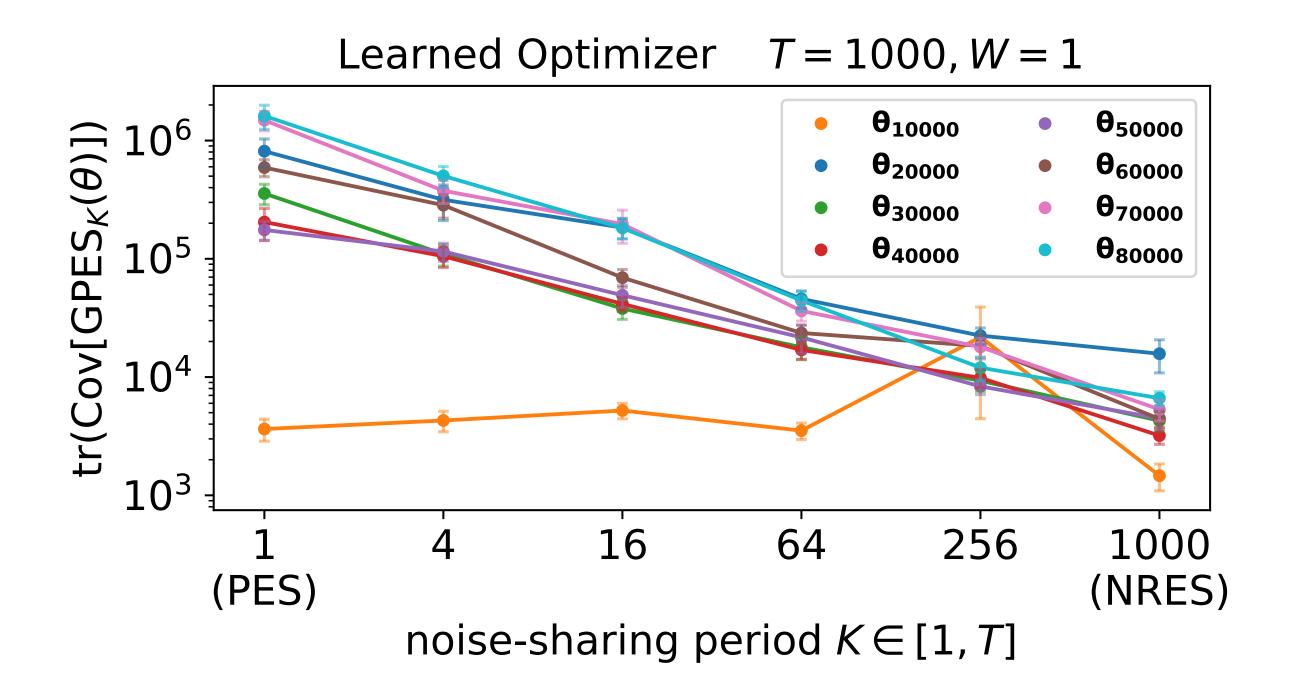
Theorem 5: Analytical characterization of the total variance of $GPES_K$ for all K.

total variance among all **GPES**_K estimators.



Variance characterization

Corollary 8: $GPES_{K=T}$ provably has the lowest



Noise-Reuse Evolution Strategies (our

1 sample / T unrolls

NRESWorker

No need to Sample Reuse

$$g := \frac{1}{2\sigma^2 W} \left(\sum_{t=0}^{\infty} (L_t^+ - L_t^-) \right) \epsilon$$
over time t in the ϵ sampled at the beginning of this episode.

NRES vs. FullES

- The same smoothing objective
- ullet NRES is T/W times more parallelizable than
- Variance comparison:

Theorem 9: Under reasonable assumptions and the same compute budget, NRES has lower variance than FullES.

	total variance tr(Cov[·])	
er. j f θ_j 10^4)	$\frac{W}{T} \sum_{i=1}^{T/W} \text{NRES}_i(\theta_j)$	$\mathrm{FullES}(heta_j)$
1	${f 1.47 \pm 0.38}$	864.34 ± 385.01
2	${f 15.74} \pm 4.90$	513.50 ± 74.71
3	${f 4.28 \pm 0.78}$	201.38 ± 30.18
4	3.20 ± 0.50	684.85 ± 86.90
5	4.47 ± 0.88	181.06 ± 24.97
6	4.42 ± 0.68	288.61 ± 46.77
7	${f 5.33} \pm 1.08$	448.30 ± 62.44
8	6.62 ± 0.89	154.86 ± 28.21
'	•	

 $T = 1000, W = 1 \quad (\pm 95\% \text{ ci})$

Meta-training learned optimizer

Learning chaotic dynamical system

Lorenz system T = 2000

 $s_t = (x_t, y_t, z_t) \in \mathbb{R}^3$

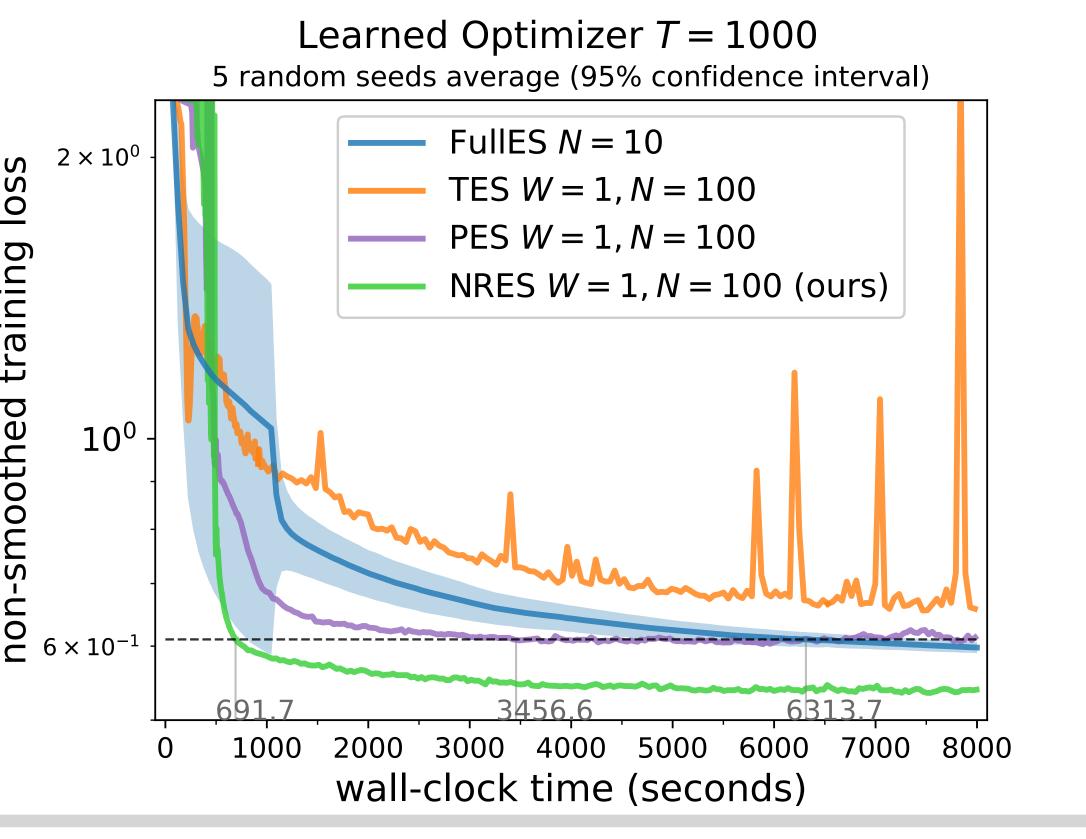
Email: oscarli@cmu.edu

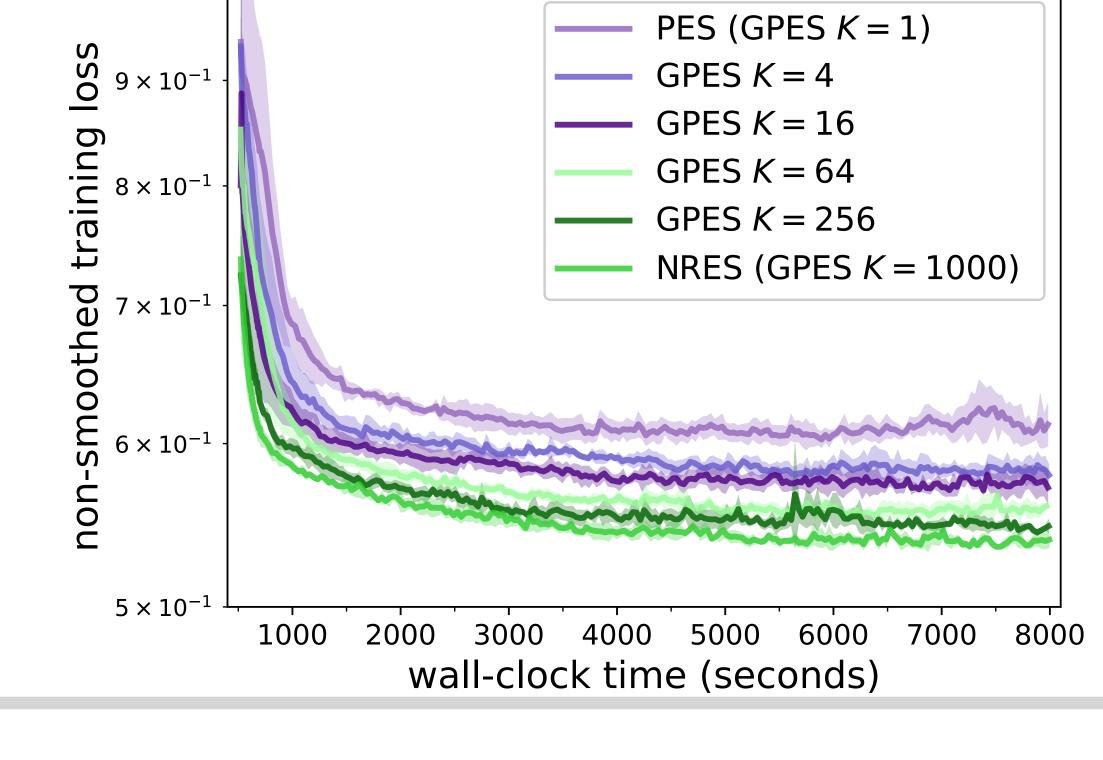
Experiments

 $\theta = (\ln(r), \ln(a))$

 $L_t^s(s_t) := (z_t - z_t^{\mathfrak{GT}})^2$

heta: learned optimizer s_t : the inner model's parameters and momentum statistics $L_t(s_t)$: the learned model's validation loss.



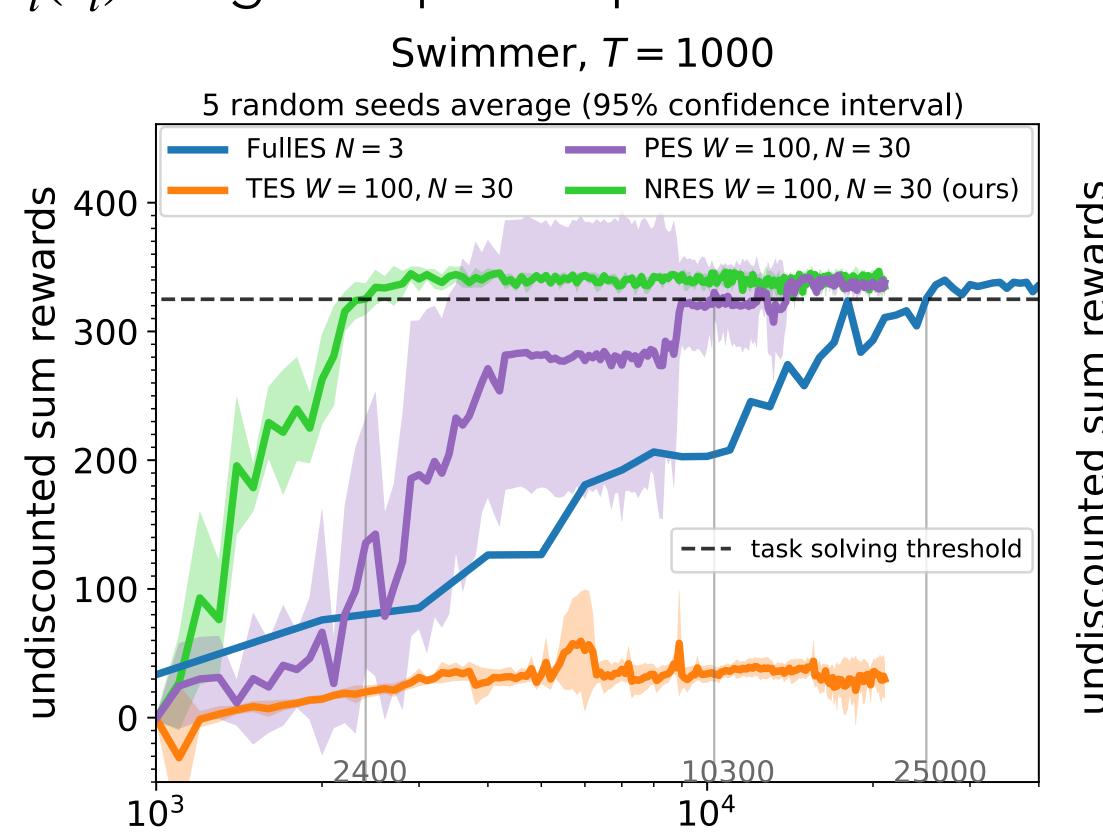


Learned Optimizer T = 1000, W = 1, N = 100

Mujoco policy learning (RL)

heta: parameter of a policy network

 $L_t(s_t)$: negative per-step reward



of sequential environment steps

 s_t : the mujoco environment state.

