

AutoSimulate: (Quickly) Learning Synthetic Data Generation

Harkirat Singh Behl, AG Baydin, Ran Gal, Philip Torr, Vibhav Vineet

Neural Networks are Data Hungry

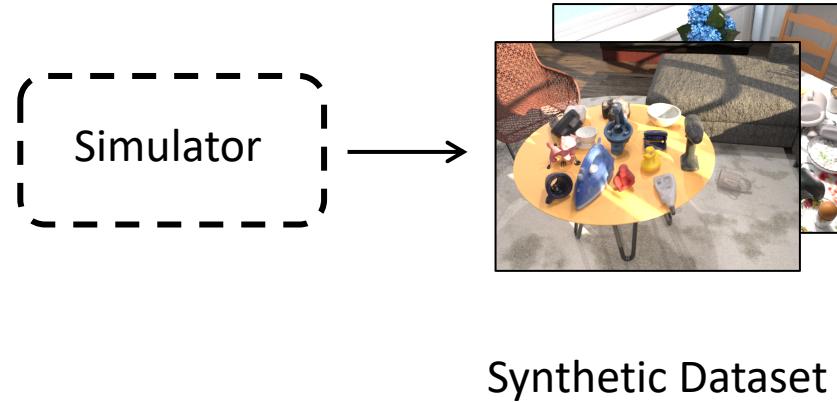


Synthetic Dataset

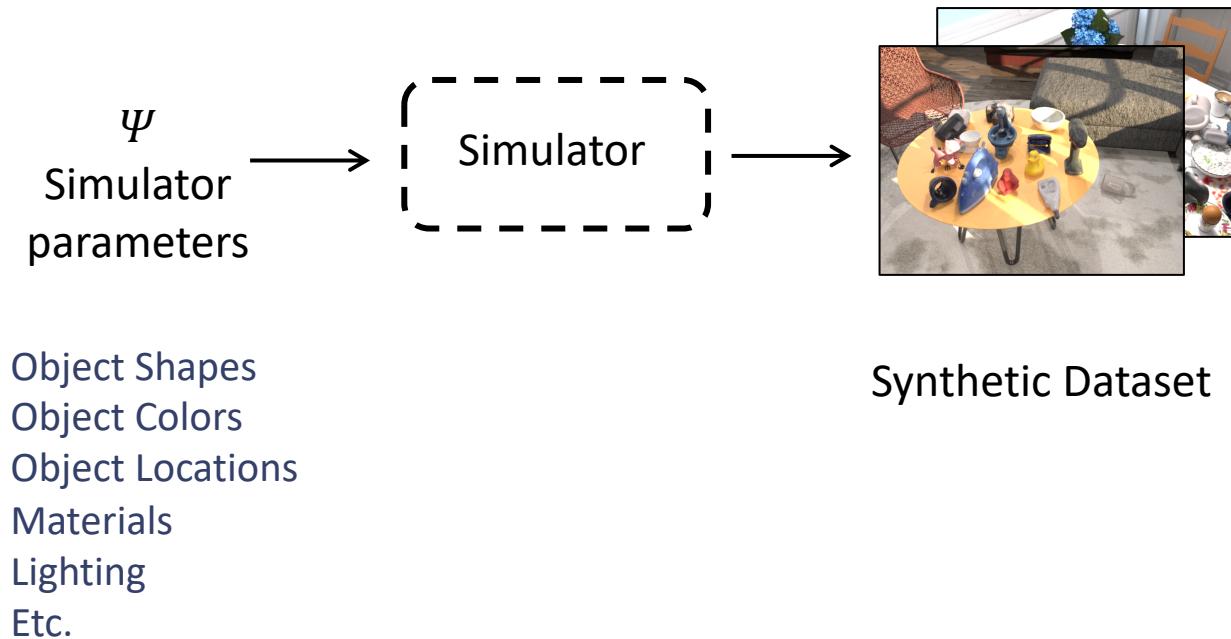


Training on
Synthetic Data

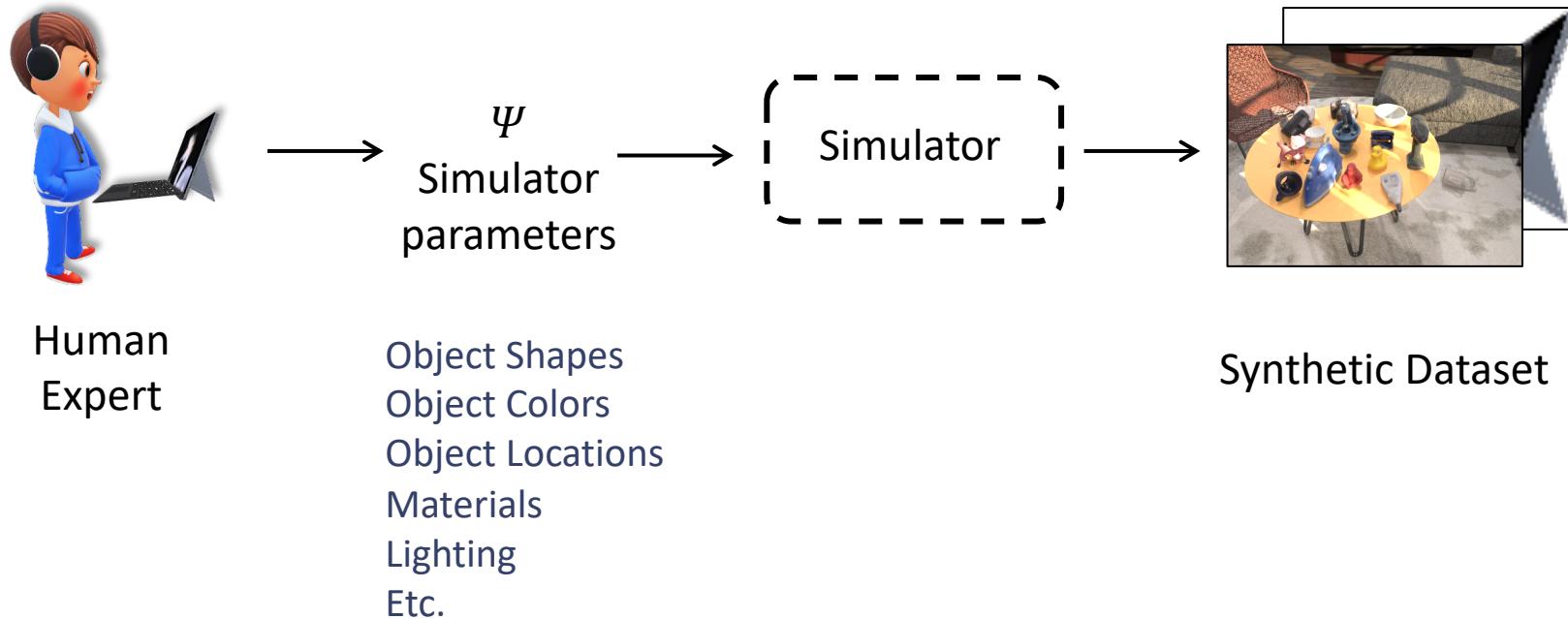
Synthetic Data Generation



Synthetic Data Generation



Manual Synthetic Data Generation



What's Wrong with Manual Synthetic Data Generation?



NEEDS SIGNIFICANT
HUMAN EFFORT



EXPENSIVE



CAN BE SUB-OPTIMAL

What's Wrong with ~~Manual~~ Synthetic Data Generation?



NEEDS SIGNIFICANT
HUMAN EFFORT

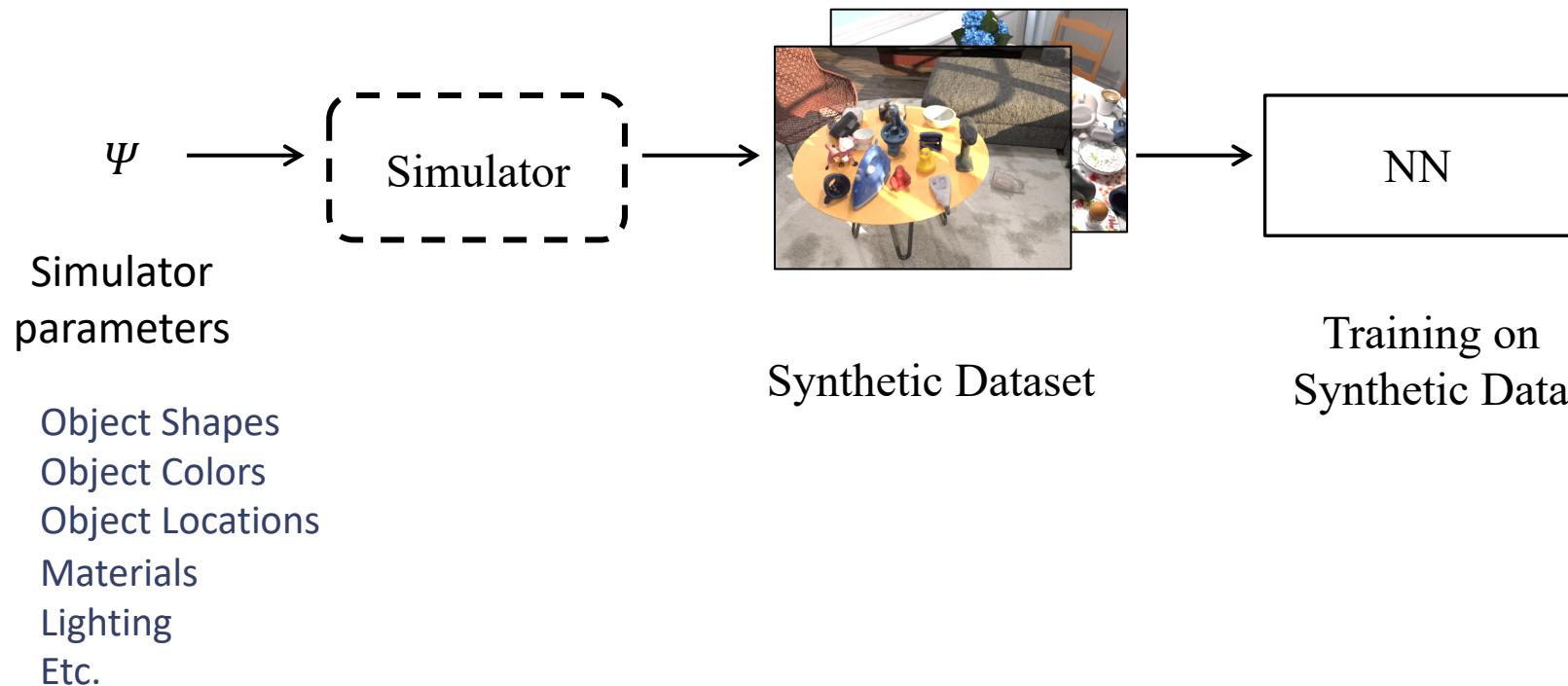


EXPENSIVE

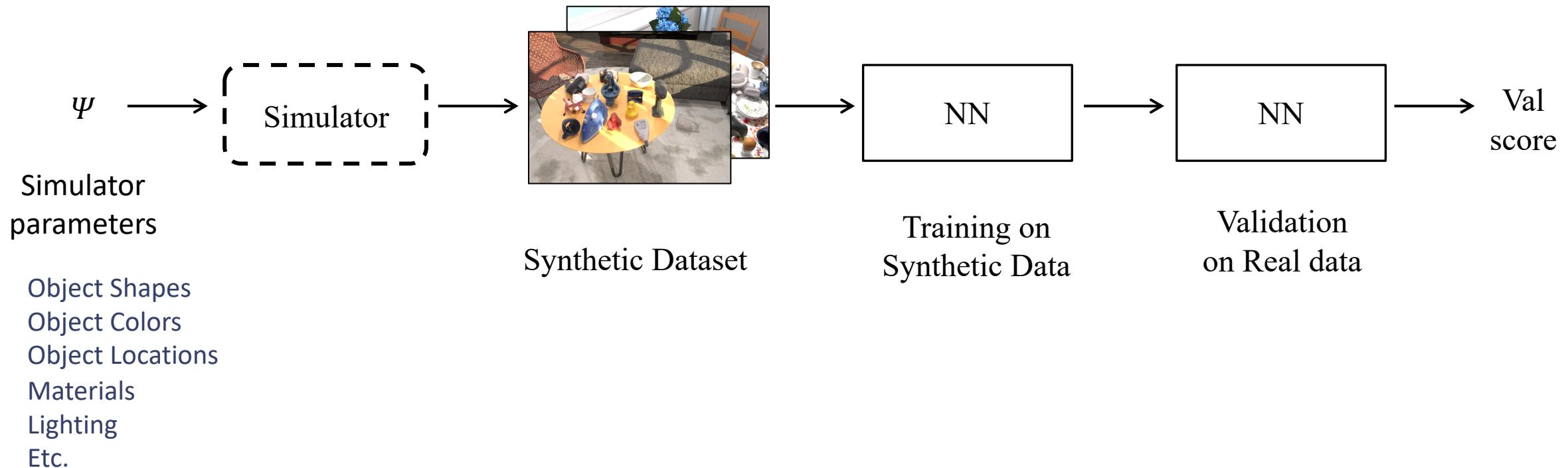


CAN BE SUB-OPTIMAL

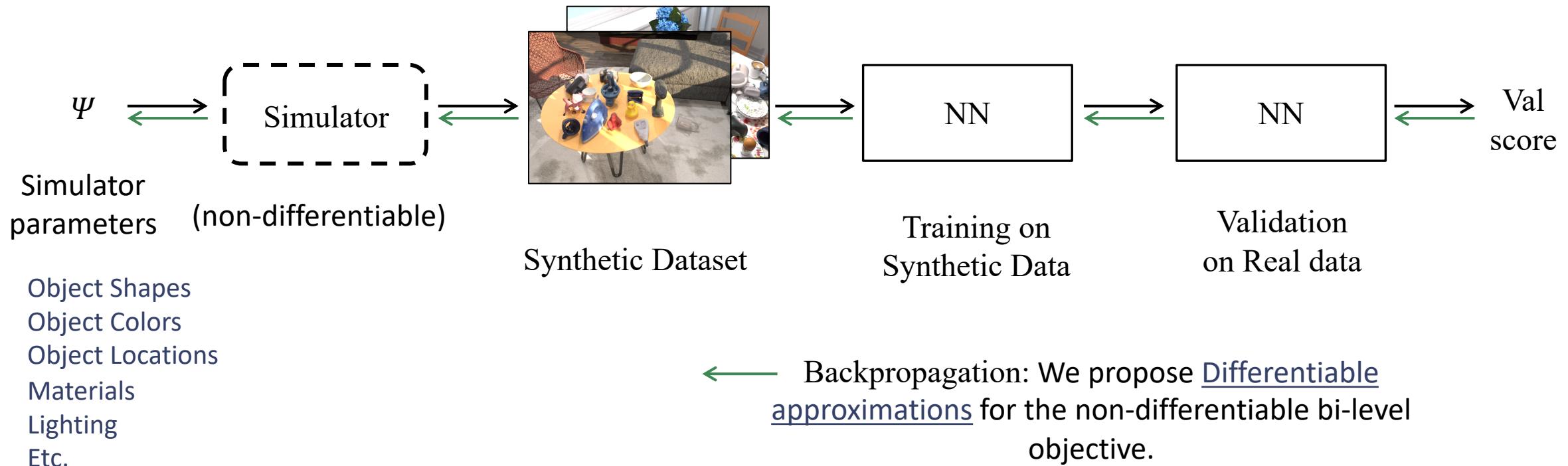
AutoSimulate: Learning Synthetic Data Generation



AutoSimulate: Learning Synthetic Data Generation



AutoSimulate: Learning Synthetic Data Generation



Problem Formulation

- Objective is to find optimal simulator parameters as

$$\min_{\psi} \quad \mathcal{L}_{\text{val}}(\hat{\theta}(\psi)) \tag{1a}$$

$$s.t. \quad \hat{\theta}(\psi) \in \arg \min_{\theta} \mathcal{L}_{\text{train}}(\theta, \psi). \tag{1b}$$

$\mathcal{L}_{\text{val}}(\hat{\theta}(\psi))$: validation loss

$\mathcal{L}_{\text{train}}(\theta, \psi)$: training loss

$\hat{\theta}(\psi)$: neural network parameters after training on data from simulator ψ

- Eqns. (1a) and (1b) represent a bi-level optimization problem.

Problem Formulation



Computing gradient of the bi-level optimization is difficult as:

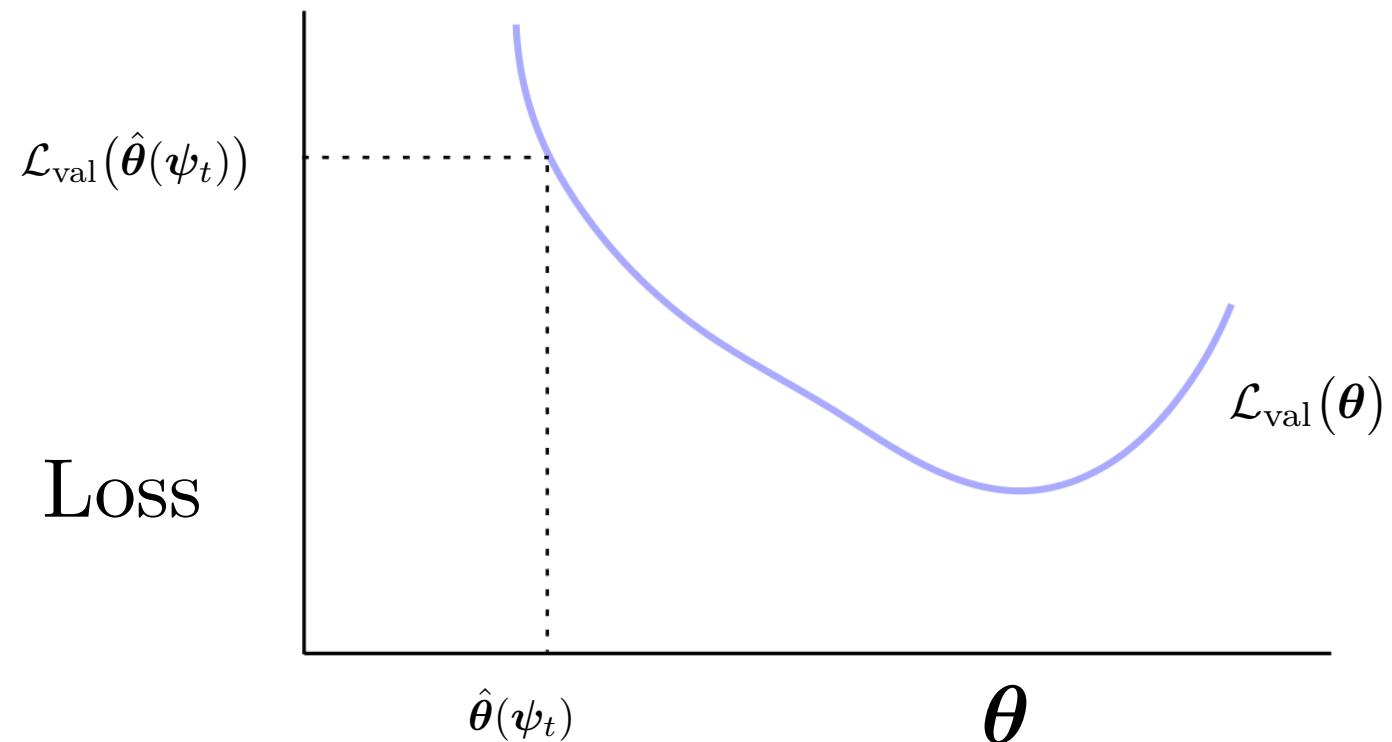
Simulator can be non-differentiable
Backpropagating through NN training is impracticable



We propose differentiable approximations for the non-differentiable bi-level objective

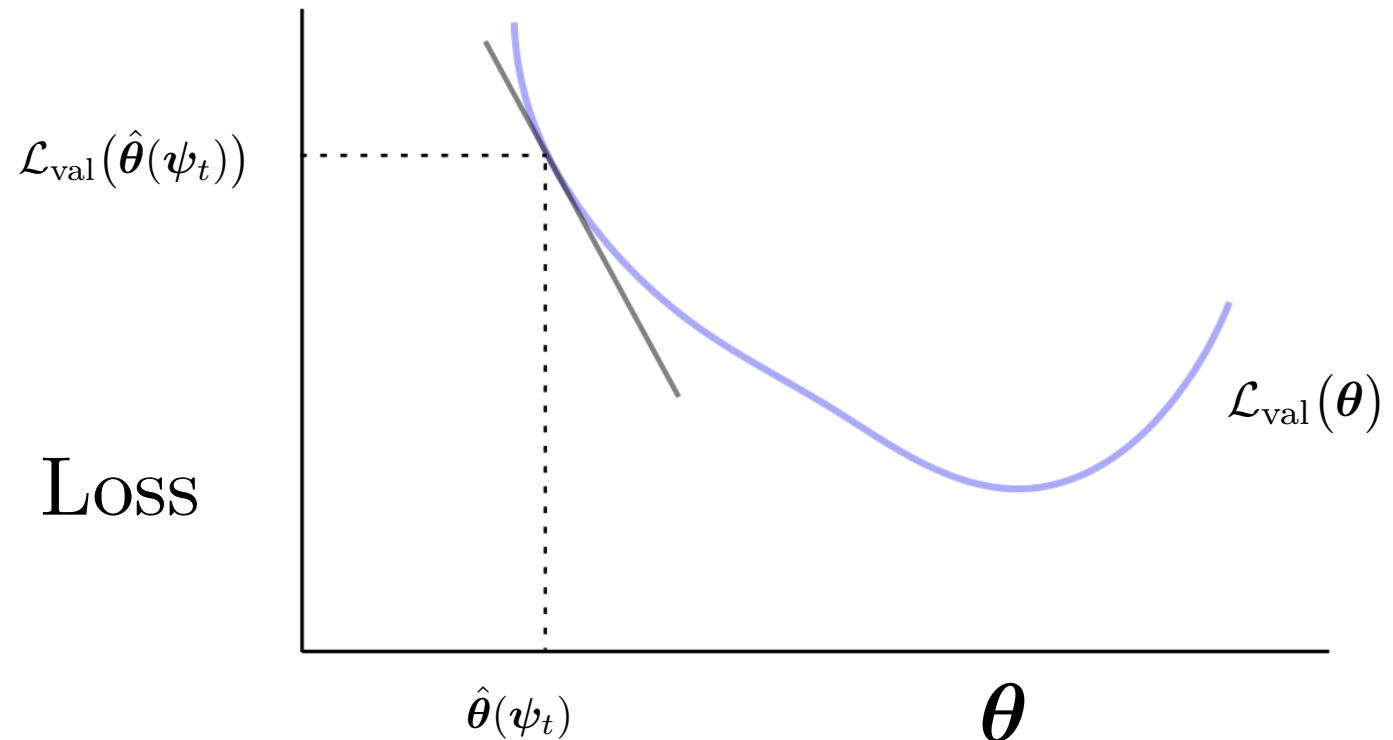
AutoSimulate

Objective: Minimize $\mathcal{L}_{\text{val}}(\hat{\theta}(\psi))$ with respect to ψ



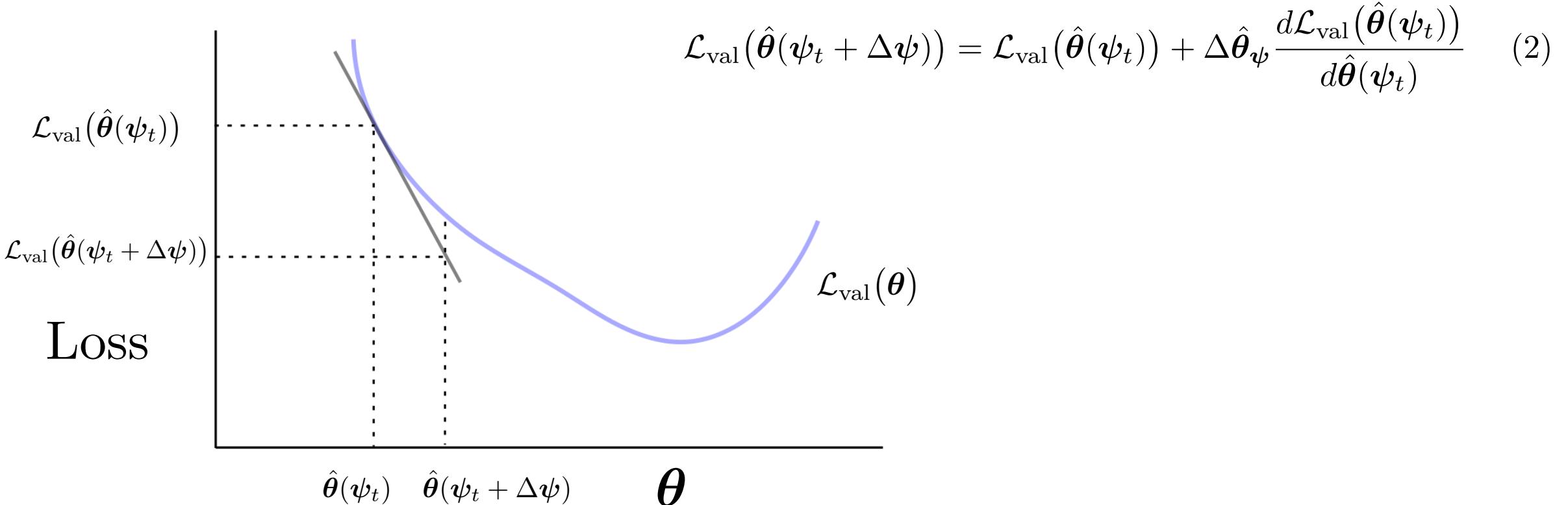
AutoSimulate

Objective: Minimize $\mathcal{L}_{\text{val}}(\hat{\theta}(\psi))$ with respect to ψ



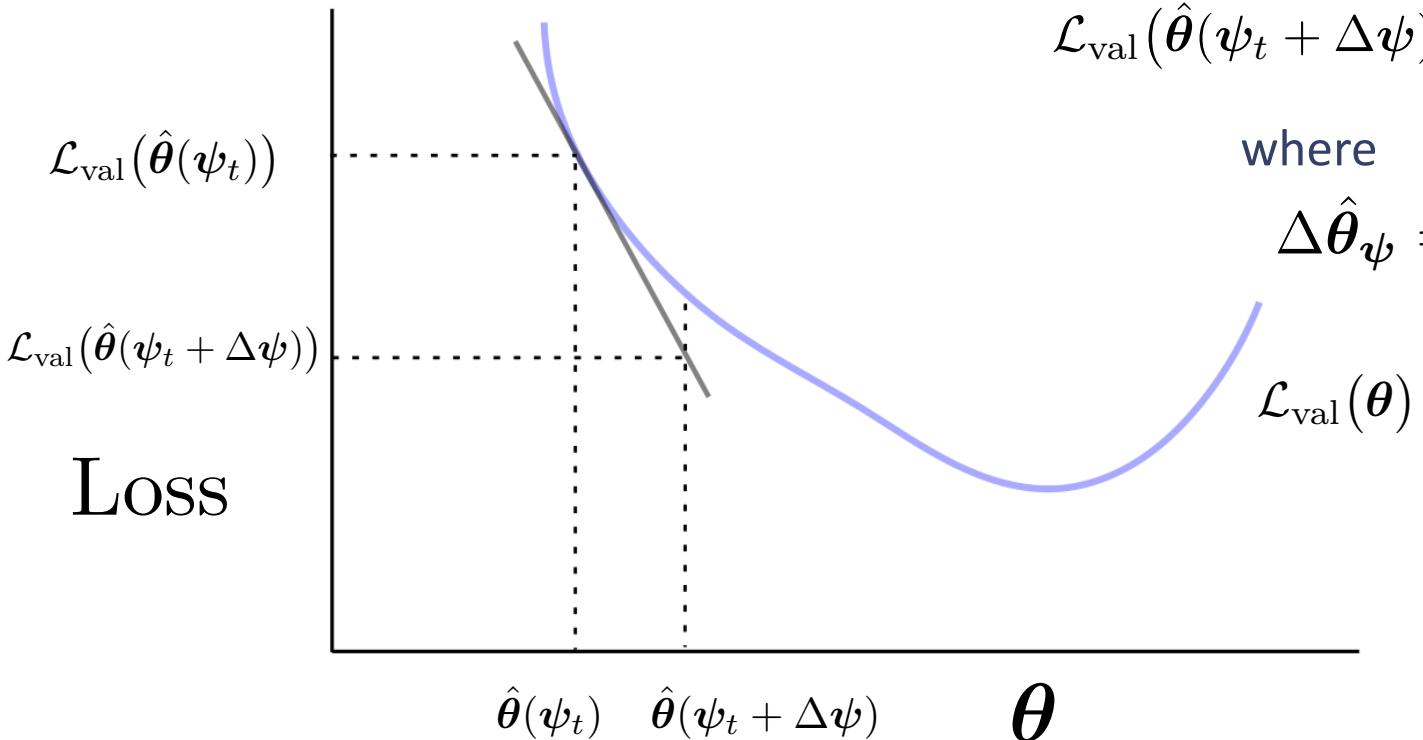
AutoSimulate

Objective: Minimize $\mathcal{L}_{\text{val}}(\hat{\theta}(\psi))$ with respect to ψ



AutoSimulate

Objective: Minimize $\mathcal{L}_{\text{val}}(\hat{\theta}(\psi))$ with respect to ψ



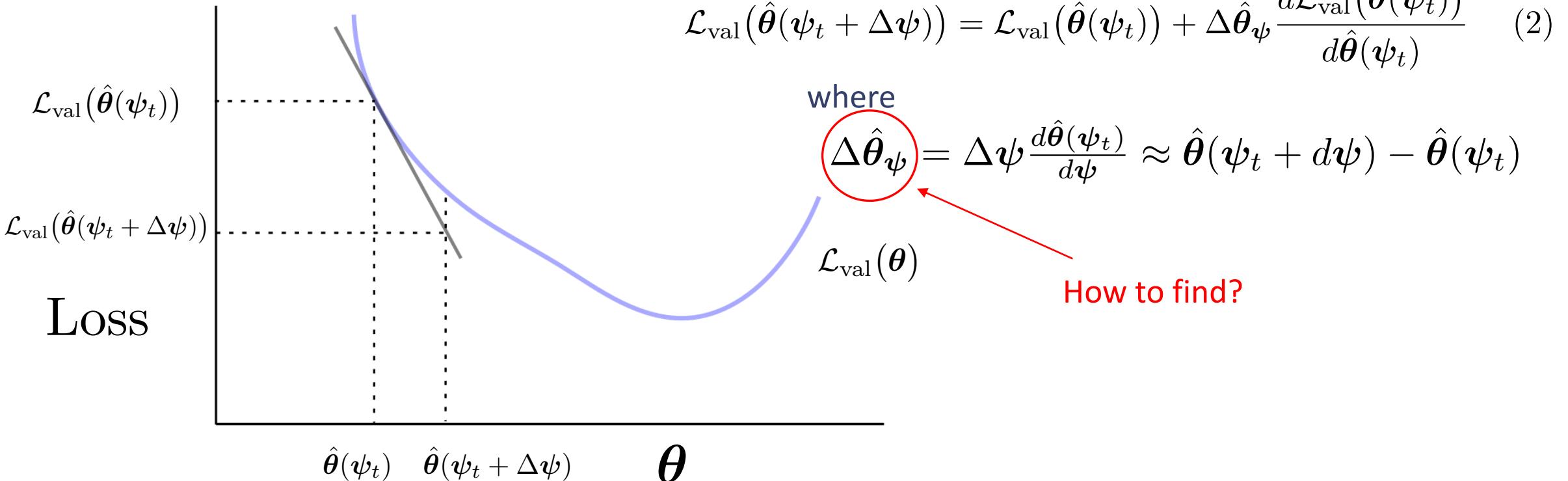
$$\mathcal{L}_{\text{val}}(\hat{\theta}(\psi_t + \Delta\psi)) = \mathcal{L}_{\text{val}}(\hat{\theta}(\psi_t)) + \Delta\hat{\theta}_\psi \frac{d\mathcal{L}_{\text{val}}(\hat{\theta}(\psi_t))}{d\hat{\theta}(\psi_t)} \quad (2)$$

where

$$\Delta\hat{\theta}_\psi = \Delta\psi \frac{d\hat{\theta}(\psi_t)}{d\psi} \approx \hat{\theta}(\psi_t + d\psi) - \hat{\theta}(\psi_t)$$

AutoSimulate

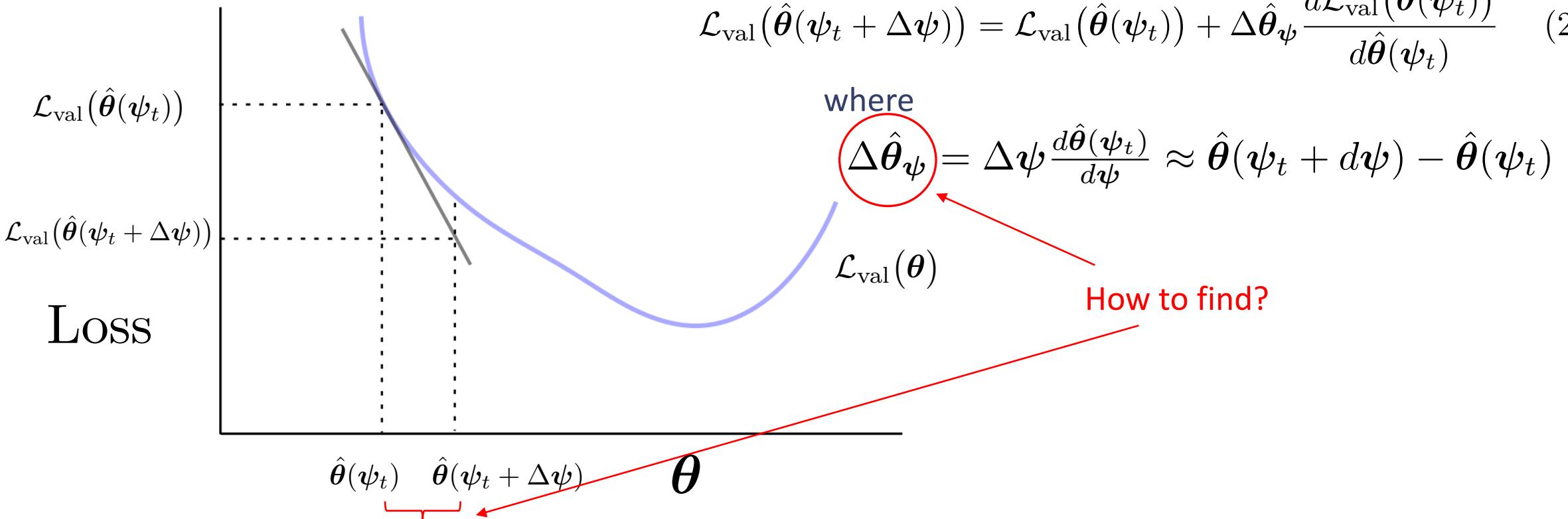
Objective: Minimize $\mathcal{L}_{\text{val}}(\hat{\theta}(\psi))$ with respect to ψ



AutoSimulate

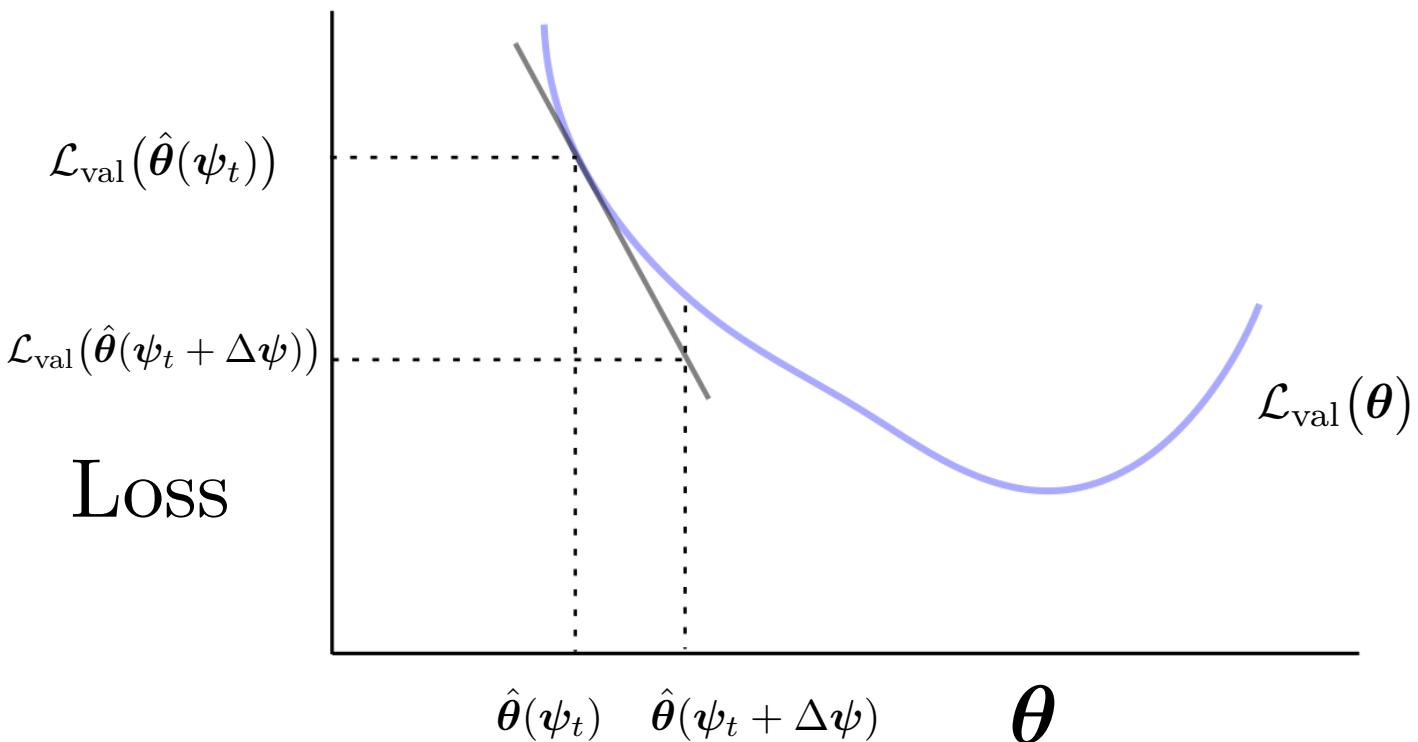
Objective: Minimize $\mathcal{L}_{\text{val}}(\hat{\theta}(\psi))$ with respect to ψ

$$\mathcal{L}_{\text{val}}(\hat{\theta}(\psi_t + \Delta\psi)) = \mathcal{L}_{\text{val}}(\hat{\theta}(\psi_t)) + \Delta\hat{\theta}_\psi \frac{d\mathcal{L}_{\text{val}}(\hat{\theta}(\psi_t))}{d\hat{\theta}(\psi_t)} \quad (2)$$



AutoSimulate

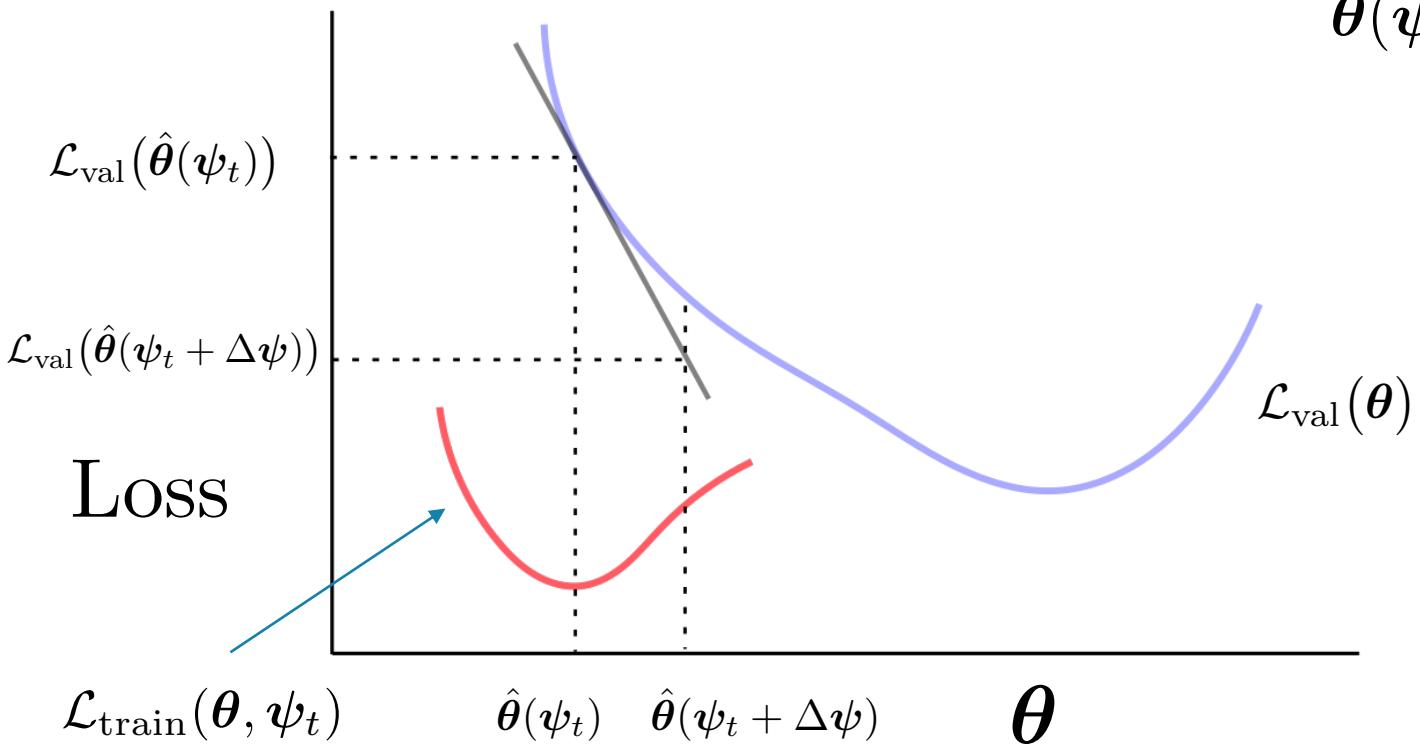
Find: $\Delta\hat{\theta}_\psi \approx \hat{\theta}(\psi_t + d\psi) - \hat{\theta}(\psi_t)$



AutoSimulate

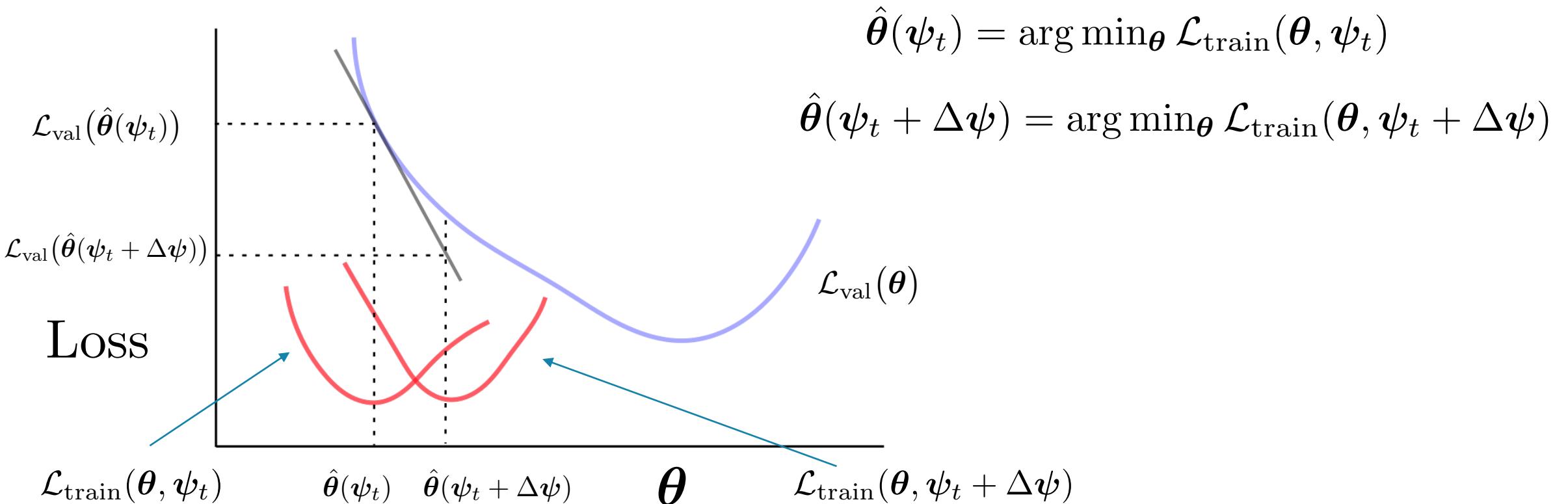
Find: $\Delta\hat{\theta}_\psi \approx \hat{\theta}(\psi_t + d\psi) - \hat{\theta}(\psi_t)$

$$\hat{\theta}(\psi_t) = \arg \min_{\theta} \mathcal{L}_{\text{train}}(\theta, \psi_t)$$



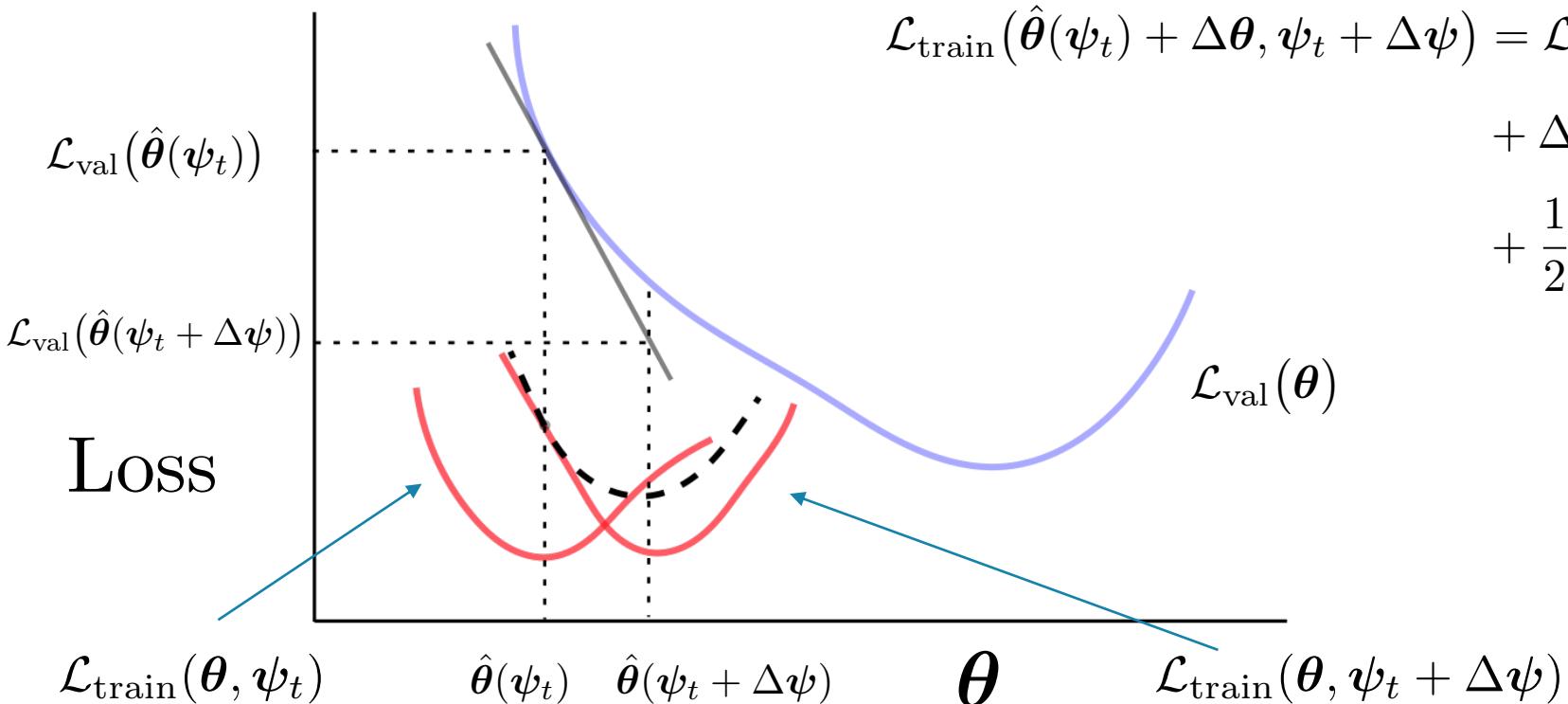
AutoSimulate

Find: $\Delta\hat{\theta}_\psi \approx \hat{\theta}(\psi_t + d\psi) - \hat{\theta}(\psi_t)$



AutoSimulate

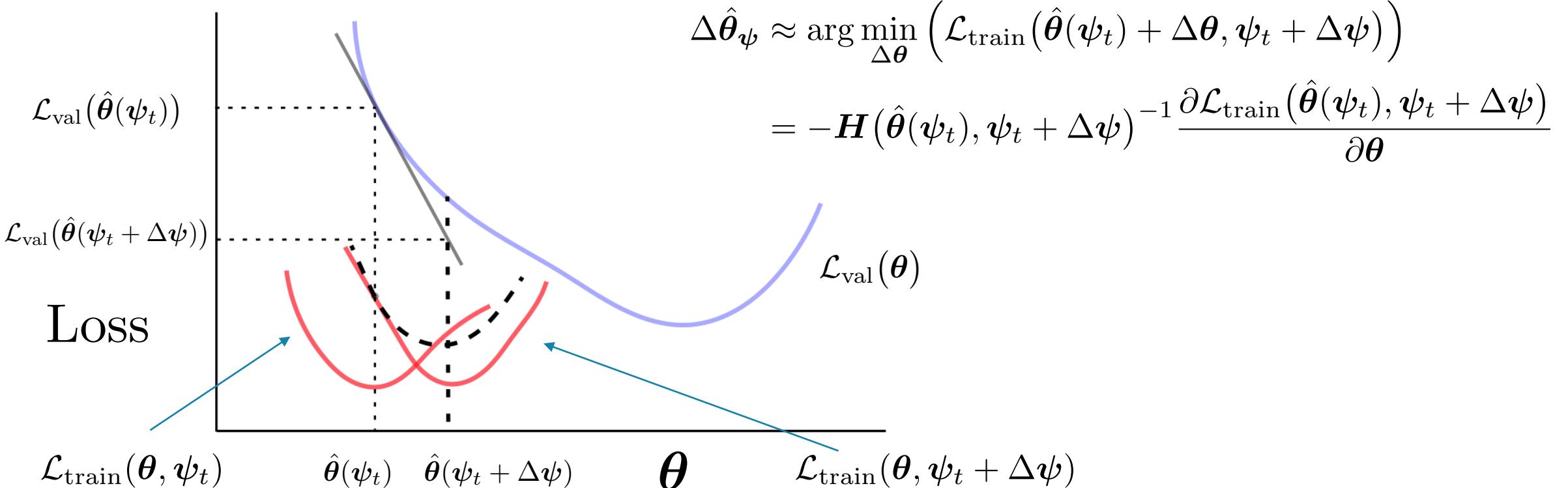
Find: $\Delta\hat{\theta}_\psi \approx \hat{\theta}(\psi_t + d\psi) - \hat{\theta}(\psi_t)$



$$\begin{aligned}\mathcal{L}_{\text{train}}(\hat{\theta}(\psi_t) + \Delta\theta, \psi_t + \Delta\psi) &= \mathcal{L}_{\text{train}}(\hat{\theta}(\psi_t), \psi_t + \Delta\psi) \\ &\quad + \Delta\theta^\top \frac{\partial}{\partial \theta} \mathcal{L}_{\text{train}}(\hat{\theta}(\psi_t), \psi_t + \Delta\psi) \\ &\quad + \frac{1}{2} \Delta\theta^\top H(\hat{\theta}(\psi_t), \psi_t + \Delta\psi) \Delta\theta + \dots\end{aligned}$$

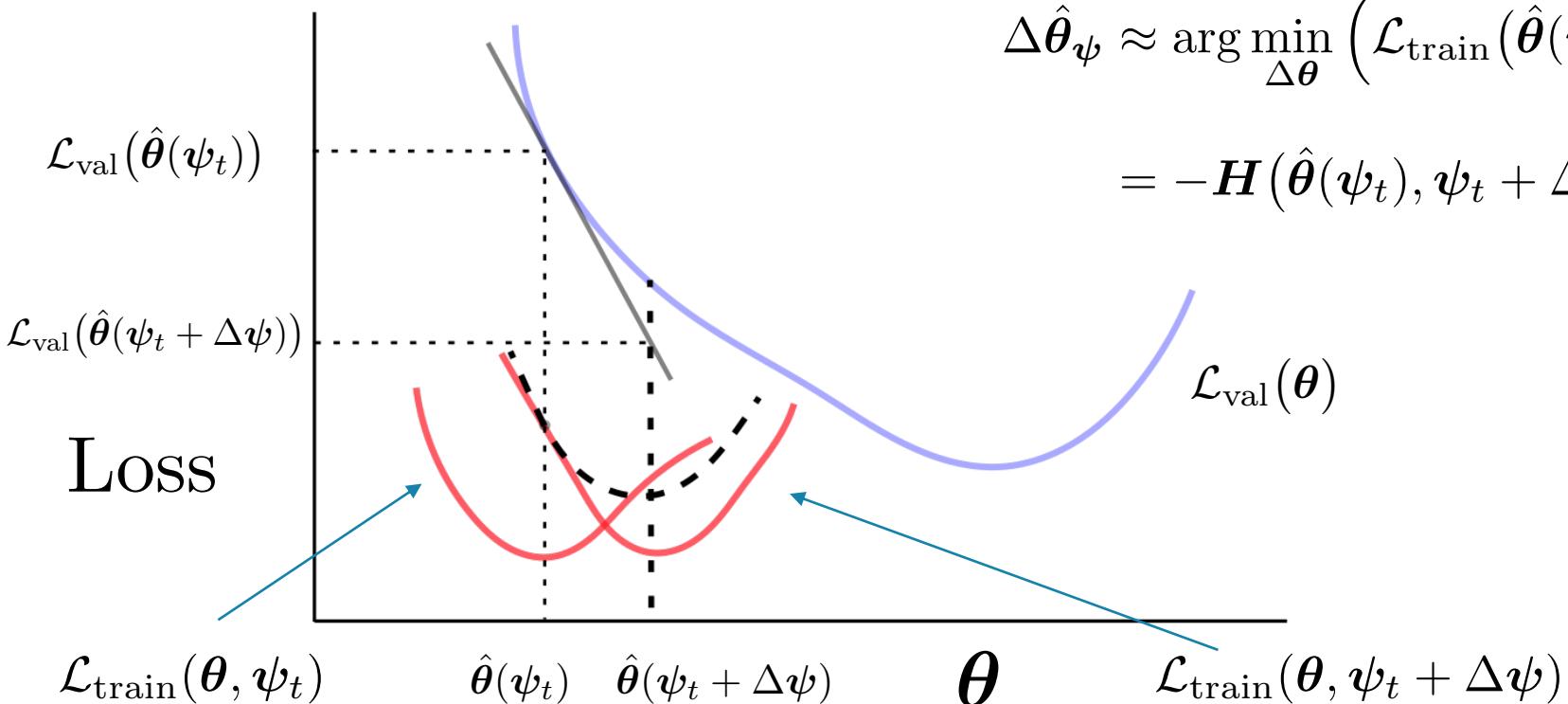
AutoSimulate

Find: $\Delta\hat{\theta}_\psi \approx \hat{\theta}(\psi_t + d\psi) - \hat{\theta}(\psi_t)$



AutoSimulate

Find: $\Delta\hat{\theta}_\psi \approx \hat{\theta}(\psi_t + d\psi) - \hat{\theta}(\psi_t)$



$$\begin{aligned}\Delta\hat{\theta}_\psi &\approx \arg \min_{\Delta\theta} \left(\mathcal{L}_{\text{train}}(\hat{\theta}(\psi_t) + \Delta\theta, \psi_t + \Delta\psi) \right) \\ &= -\mathbf{H}(\hat{\theta}(\psi_t), \psi_t + \Delta\psi)^{-1} \frac{\partial \mathcal{L}_{\text{train}}(\hat{\theta}(\psi_t), \psi_t + \Delta\psi)}{\partial \theta}\end{aligned}$$

Putting back into Eq.2

AutoSimulate

Differentiable Approximation:

$$\tilde{\mathcal{L}}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi_t + \Delta\psi)) = \mathcal{L}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi_t)) - \frac{\partial \mathcal{L}_{\text{train}}(\hat{\boldsymbol{\theta}}(\psi_t), \psi_t + \Delta\psi)}{\partial \boldsymbol{\theta}}^\top \mathbf{H}(\hat{\boldsymbol{\theta}}(\psi_t), \psi_t + \Delta\psi)^{-1} \frac{d\mathcal{L}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi_t))}{d\boldsymbol{\theta}}$$

AutoSimulate

Differentiable Approximation:

$$\tilde{\mathcal{L}}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi_t + \Delta\psi)) = \mathcal{L}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi_t)) - \frac{\partial \mathcal{L}_{\text{train}}(\hat{\boldsymbol{\theta}}(\psi_t), \psi_t + \Delta\psi)}{\partial \boldsymbol{\theta}}^\top \mathbf{H}(\hat{\boldsymbol{\theta}}(\psi_t), \psi_t + \Delta\psi)^{-1} \frac{d\mathcal{L}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi_t))}{d\boldsymbol{\theta}}$$

$$\Rightarrow \tilde{\mathcal{L}}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi)) = \mathcal{L}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi_t)) - \frac{\partial \mathcal{L}_{\text{train}}(\hat{\boldsymbol{\theta}}(\psi_t), \psi)}{\partial \boldsymbol{\theta}}^\top \mathbf{H}(\hat{\boldsymbol{\theta}}(\psi_t), \psi)^{-1} \frac{d\mathcal{L}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi_t))}{d\boldsymbol{\theta}}$$

AutoSimulate

Differentiable Approximation:

$$\tilde{\mathcal{L}}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi)) = \mathcal{L}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi_t)) - \frac{\partial \mathcal{L}_{\text{train}}(\hat{\boldsymbol{\theta}}(\psi_t), \psi)}{\partial \boldsymbol{\theta}}^\top \mathbf{H}(\hat{\boldsymbol{\theta}}(\psi_t), \psi)^{-1} \frac{d\mathcal{L}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi_t))}{d\boldsymbol{\theta}}$$

Derivative:

$$\frac{\partial \tilde{\mathcal{L}}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi))}{\partial \psi} \Big|_{\psi=\psi_t} = - \frac{\partial}{\partial \psi} \left[\frac{\partial \mathcal{L}_{\text{train}}(\hat{\boldsymbol{\theta}}(\psi_t), \psi)}{\partial \boldsymbol{\theta}}^\top \right] \Big|_{\psi=\psi_t} \mathbf{H}(\hat{\boldsymbol{\theta}}(\psi_t), \psi)^{-1} \frac{d\mathcal{L}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi_t))}{d\boldsymbol{\theta}}$$

AutoSimulate

Derivative:

$$\frac{\partial \tilde{\mathcal{L}}_{\text{val}}(\hat{\theta}(\psi))}{\partial \psi} \Big|_{\psi=\psi_t} = -\frac{\partial}{\partial \psi} \left[\frac{\partial \mathcal{L}_{\text{train}}(\hat{\theta}(\psi_t), \psi)}{\partial \theta} \right]^\top \Big|_{\psi=\psi_t} H(\hat{\theta}(\psi_t), \psi)^{-1} \frac{d\mathcal{L}_{\text{val}}(\hat{\theta}(\psi_t))}{d\theta}$$

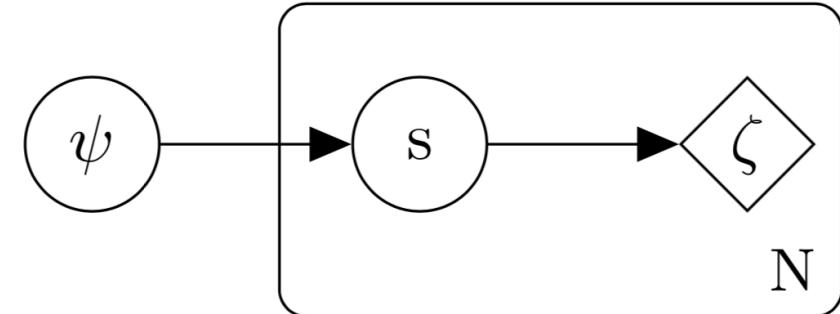
How to find?

Stochastic Simulator

- The term requires backpropagation through the dataset generation.

$$\frac{\partial}{\partial \psi} \left[\frac{\partial \mathcal{L}_{\text{train}}(\hat{\theta}(\psi_t), \psi)}{\partial \theta} \right]^\top = \frac{\partial}{\partial \psi} \mathbb{E}_{\zeta \sim p_\psi} \left[\frac{\partial}{\partial \theta} l(\zeta, \hat{\theta}(\psi_t)) \right]$$

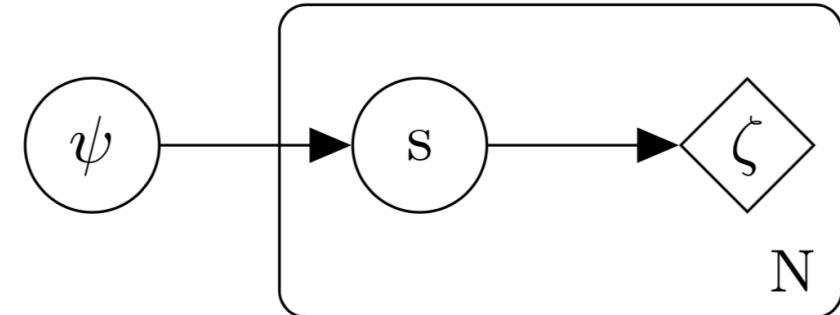
- We assume a stochastic simulator that involves a deterministic renderer.
- We make the stochasticity in the process explicit by separating the stochastic part of the simulator from the deterministic rendering.



Stochastic Simulator

Given the deterministic renderer component $\zeta = r(s)$, we would like to find the optimal values of simulator parameters ψ that parameterize $s \sim q_\psi(s)$ representing the stochastic component, expressing the overall simulator as $\zeta \sim p_\psi(\zeta)$.

$$p_\psi(\zeta) = \int_{s \in \{s | r(s) = \zeta\}} q_\psi(s) ds$$



AutoSimulate

Derivative:

$$\frac{\partial \tilde{\mathcal{L}}_{\text{val}}(\hat{\theta}(\psi))}{\partial \psi} \Big|_{\psi=\psi_t} = -\frac{\partial}{\partial \psi} \left[\frac{\partial \mathcal{L}_{\text{train}}(\hat{\theta}(\psi_t), \psi)}{\partial \theta} \right]^\top \Big|_{\psi=\psi_t} \mathbf{H}(\hat{\theta}(\psi_t), \psi)^{-1} \frac{d\mathcal{L}_{\text{val}}(\hat{\theta}(\psi_t))}{d\theta}$$

How to find?

How to find?

AutoSimulate

Derivative:

$$\frac{\partial \tilde{\mathcal{L}}_{\text{val}}(\hat{\theta}(\psi))}{\partial \psi} \Big|_{\psi=\psi_t} = -\frac{\partial}{\partial \psi} \left[\frac{\partial \mathcal{L}_{\text{train}}(\hat{\theta}(\psi_t), \psi)}{\partial \theta} \right]^\top \Big|_{\psi=\psi_t} \mathbf{H}(\hat{\theta}(\psi_t), \psi)^{-1} \frac{d\mathcal{L}_{\text{val}}(\hat{\theta}(\psi_t))}{d\theta}$$

How to find?

See paper for more details

How to find?

Approximations for $\Delta\hat{\theta}_\psi$

- Earlier we proposed a quadratic approximation.
- To further reduce the compute overhead, we propose more approximations.

Approximation ($\Delta\hat{\theta}_\psi$)	Derivative Term ($\frac{\partial}{\partial \boldsymbol{\psi}} \tilde{\mathcal{L}}_{\text{val}}(\hat{\boldsymbol{\theta}}(\boldsymbol{\psi}))$)
Quadratic $-H(\hat{\boldsymbol{\theta}}(\boldsymbol{\psi}_t), \boldsymbol{\psi})^{-1} \frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{L}_{\text{tr.}}(\hat{\boldsymbol{\theta}}(\boldsymbol{\psi}_t), \boldsymbol{\psi})$	$-\frac{\partial}{\partial \boldsymbol{\psi}} \mathbb{E}_{\zeta \sim p_\psi} [\frac{\partial}{\partial \boldsymbol{\theta}} l(\zeta, \hat{\boldsymbol{\theta}}(\boldsymbol{\psi}_t))]^\top \Big _{\boldsymbol{\psi}=\boldsymbol{\psi}_t} \mathbf{H}(\hat{\boldsymbol{\theta}}(\boldsymbol{\psi}_t), \boldsymbol{\psi}_t)^{-1} \frac{d}{d \boldsymbol{\theta}} \mathcal{L}_{\text{val}}(\hat{\boldsymbol{\theta}}(\boldsymbol{\psi}_t))$
Approx. Quadratic $H(\hat{\boldsymbol{\theta}}(\boldsymbol{\psi}_t), \boldsymbol{\psi}) \frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{L}_{\text{tr.}}(\hat{\boldsymbol{\theta}}(\boldsymbol{\psi}_t), \boldsymbol{\psi})$	$\frac{\partial}{\partial \boldsymbol{\psi}} \mathbb{E}_{\zeta \sim p_\psi} [\frac{\partial}{\partial \boldsymbol{\theta}} l(\zeta, \hat{\boldsymbol{\theta}}(\boldsymbol{\psi}_t))]^\top \Big _{\boldsymbol{\psi}=\boldsymbol{\psi}_t} \mathbf{H}(\hat{\boldsymbol{\theta}}(\boldsymbol{\psi}_t), \boldsymbol{\psi}_t) \frac{d}{d \boldsymbol{\theta}} \mathcal{L}_{\text{val}}(\hat{\boldsymbol{\theta}}(\boldsymbol{\psi}_t))$
Linear $-\frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{L}_{\text{tr.}}(\hat{\boldsymbol{\theta}}(\boldsymbol{\psi}_t), \boldsymbol{\psi})$	$-\frac{\partial}{\partial \boldsymbol{\psi}} \mathbb{E}_{\zeta \sim p_\psi} [\frac{\partial}{\partial \boldsymbol{\theta}} l(\zeta, \hat{\boldsymbol{\theta}}(\boldsymbol{\psi}_t))]^\top \Big _{\boldsymbol{\psi}=\boldsymbol{\psi}_t} \frac{d}{d \boldsymbol{\theta}} \mathcal{L}_{\text{val}}(\hat{\boldsymbol{\theta}}(\boldsymbol{\psi}_t))$
No Val 1	$-\frac{\partial}{\partial \boldsymbol{\psi}} \mathbb{E}_{\zeta \sim p_\psi} [\frac{\partial}{\partial \boldsymbol{\theta}} l(\zeta, \hat{\boldsymbol{\theta}}(\boldsymbol{\psi}_t))]^\top \Big _{\boldsymbol{\psi}=\boldsymbol{\psi}_t}$

AutoSimulate

Derivative:

$$\frac{\partial \tilde{\mathcal{L}}_{\text{val}}(\hat{\theta}(\psi))}{\partial \psi} \Big|_{\psi=\psi_t} = - \frac{\partial}{\partial \psi} \left[\frac{\partial \mathcal{L}_{\text{train}}(\hat{\theta}(\psi_t), \psi)}{\partial \theta} \right]^\top \Big|_{\psi=\psi_t} \mathbf{H}(\hat{\theta}(\psi_t), \psi)^{-1} \frac{d \mathcal{L}_{\text{val}}(\hat{\theta}(\psi_t))}{d \theta}$$

Optimizing Simulator:

$$\psi_{t+1} \leftarrow \psi_t - \alpha \frac{\partial \tilde{\mathcal{L}}_{\text{val}}(\hat{\theta}(\psi))}{\partial \psi} \Big|_{\psi=\psi_t}$$

Algorithm: AutoSimulate

Algorithm 1: AutoSimulate

for number of iterations **do**

 Sample dataset of size K: $D_{\text{train}} \sim p_{\psi_t}(\zeta)$

 Fine-tune model for ϵ epochs on D_{train}

 Compute $\mathbf{H}(\hat{\boldsymbol{\theta}}(\psi_t), \psi_t)^{-1} \frac{d}{d\boldsymbol{\theta}} \mathcal{L}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi_t))$ using CG

 Compute gradient of expectation as $\sum_{k=1}^K \frac{d}{d\psi} \log q_\psi(s_k) \cdot \left[\frac{\partial}{\partial \boldsymbol{\theta}} l(r(s_k), \hat{\boldsymbol{\theta}}(\psi_t)) \right]^\top$

 Update simulator by descending the gradient

$- \frac{\partial}{\partial \psi} \mathbb{E}_{\zeta \sim p_\psi} \left[\frac{\partial}{\partial \boldsymbol{\theta}} l(\zeta, \hat{\boldsymbol{\theta}}(\psi_t)) \right]^\top \Big|_{\psi=\psi_t} \mathbf{H}(\hat{\boldsymbol{\theta}}(\psi_t), \psi_t)^{-1} \frac{d}{d\boldsymbol{\theta}} \mathcal{L}_{\text{val}}(\hat{\boldsymbol{\theta}}(\psi_t))$

end for

Experiments

- Optimizing photorealistic Arnold renderer (non-differentiable)
- For object detection on real-world dataset LM-O (with 12 classes)



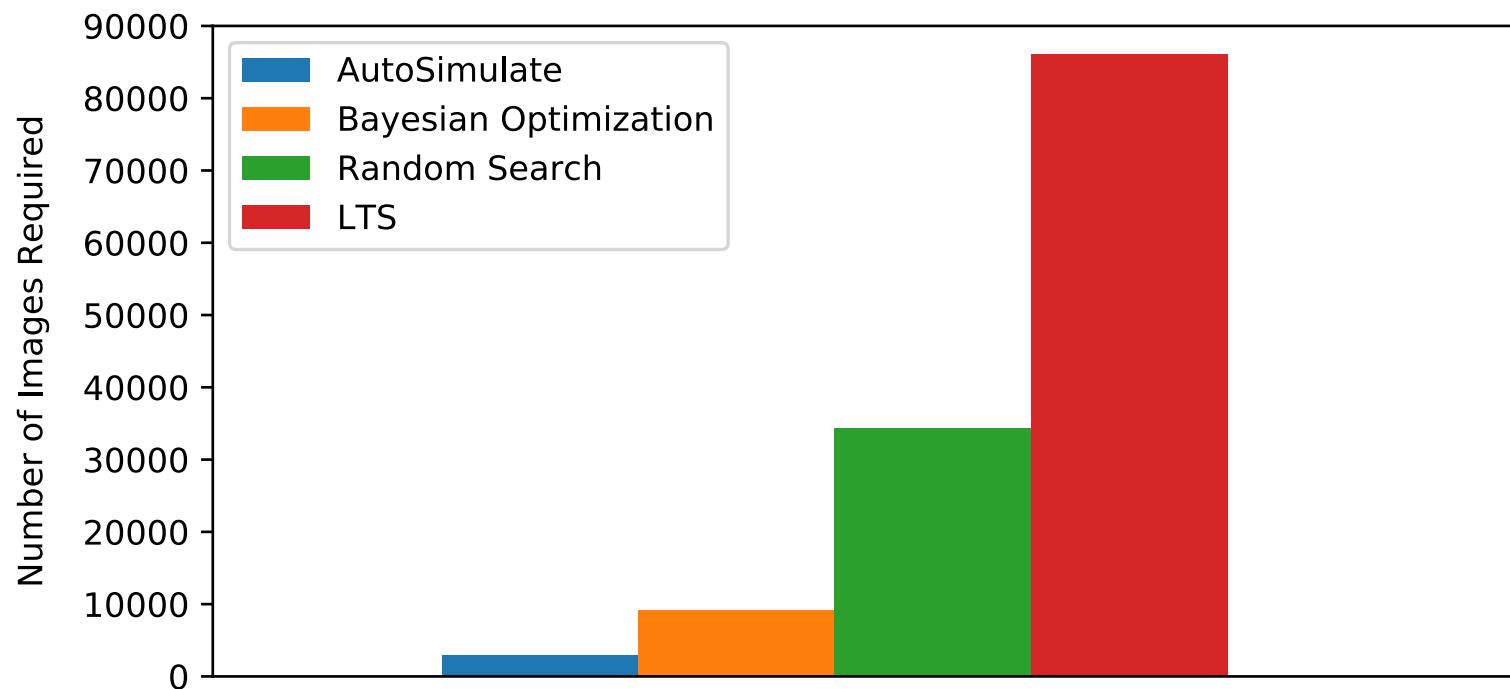
Results: Faster

50x

Results: More Accurate (mAP)

+8.5%

Results: Requires Lesser Data

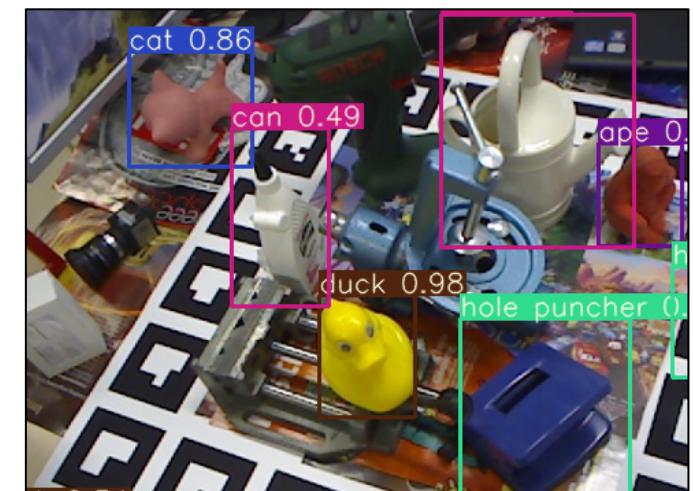
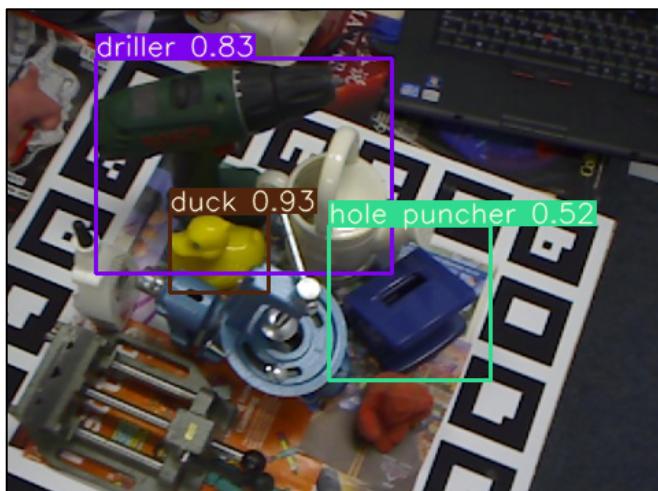
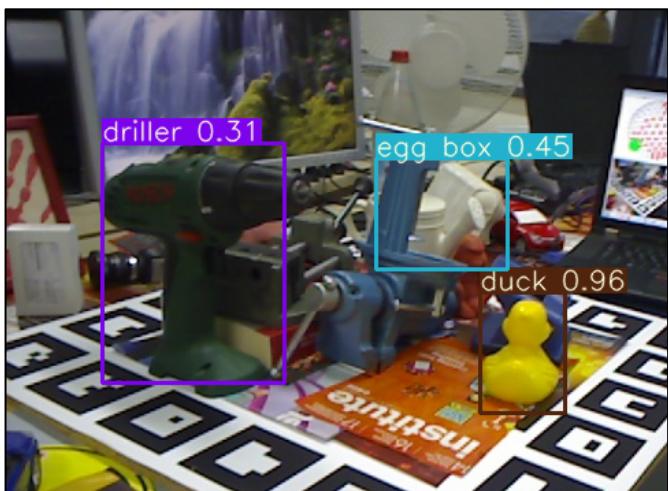


Real World Object Detection

- We run each method for a 1,000 epochs
- Val. mAP: maximum validation mAP, Test mAP: test mAP.
- Images and Time: number of images generated and time spent to reach maximum validation mAP.

Method	Val. mAP	Images	Time(s)	Test mAP
REINFORCE (LTS)	40.2	86,150	114,360	37.2
Bayesian Optimization	39.3	9,200	83,225	37.5
Random Search	40.3	34,300	134,318	37.0
AutoSimulate	37.1	8,950	23,193	36.1
AutoSimulate(Approx Quad)	40.1	2,950	2,321	37.4
AutoSimulate (Linear)	41.4	17,850	30,477	45.9

Results: Object Detections



Ablation: Generalization

- In this ablation, we study whether a simulator trained on a shallow network generalizes to a deeper network.
- Simulator was trained on Yolo and tested on Faster-rcnn.

Method	mAP (Faster-rcnn)	mAP (Yolo)
REINFORCE (LTS)	51.8	37.2
Bayesian Optimization	46.0	37.5
Random Search	50.3	36.8
Ours	55.1	45.9

Ablation: Effect of Freezing Layers

- It is a common practise to train on synthetic data with the initial layers of the network frozen and trained on real data.
- We see the relative performance of different simulator optimization methods under this scenario

Method	0 frozen layers			98 frozen layers			104 frozen layers		
	mAP	Time(s)	Images	mAP	Time(s)	Images	mAP	Time(s)	Images
REINFORCE (LTS)	37.2	114,360	86,150	33.0	114,360	86,150	31.9	145,193	104,600
Bayesian Optimization	37.5	83,225	9,225	31.7	13,940	3,550	31.7	30,538	6,050
Random Search	36.8	134,137	34,300	30.2	8,913	3,500	28.9	73,411	21,650
Ours	45.9	30,477	17,850	37.1	2,321	2,950	35.8	958	1,000

Ablation: Effect of Network Size

- We first examine the effect of network depth on simulator training
- Our method is even faster in optimizing smaller networks.

Method	Yolo-spp			Yolo-Tiny		
	mAP	Time(s)	Images	mAP	Time(s)	Images
REINFORCE (LTS)	37.2	114,360	86,150	24.7	3,475	11,550
Bayesian Optimization	37.5	83,225	9,225	19.5	65,760	35,700
Random Search	36.8	134,137	34,300	20.6	7,319	11,620
Ours	45.9	30,477	17,850	21.2	484	280

Acknowledgements

- Vibhav Vineet (MSR Redmond)
- Tencent
- Ondrej Miksik (Microsoft AI, Zurich), Tomas Hodan (Czech Technical University Prague)
- Emanuel Shalev, Pedro Urbina (Microsoft Hololens Synthetics Team)
- Prof. M. Pawan Kumar (Oxford)

Thank You
