



THE UNIVERSITY
of EDINBURGH

GRADIENT BASED HPO OVER LONG HORIZONS

Paul Micaelli and Amos Storkey

OVERVIEW

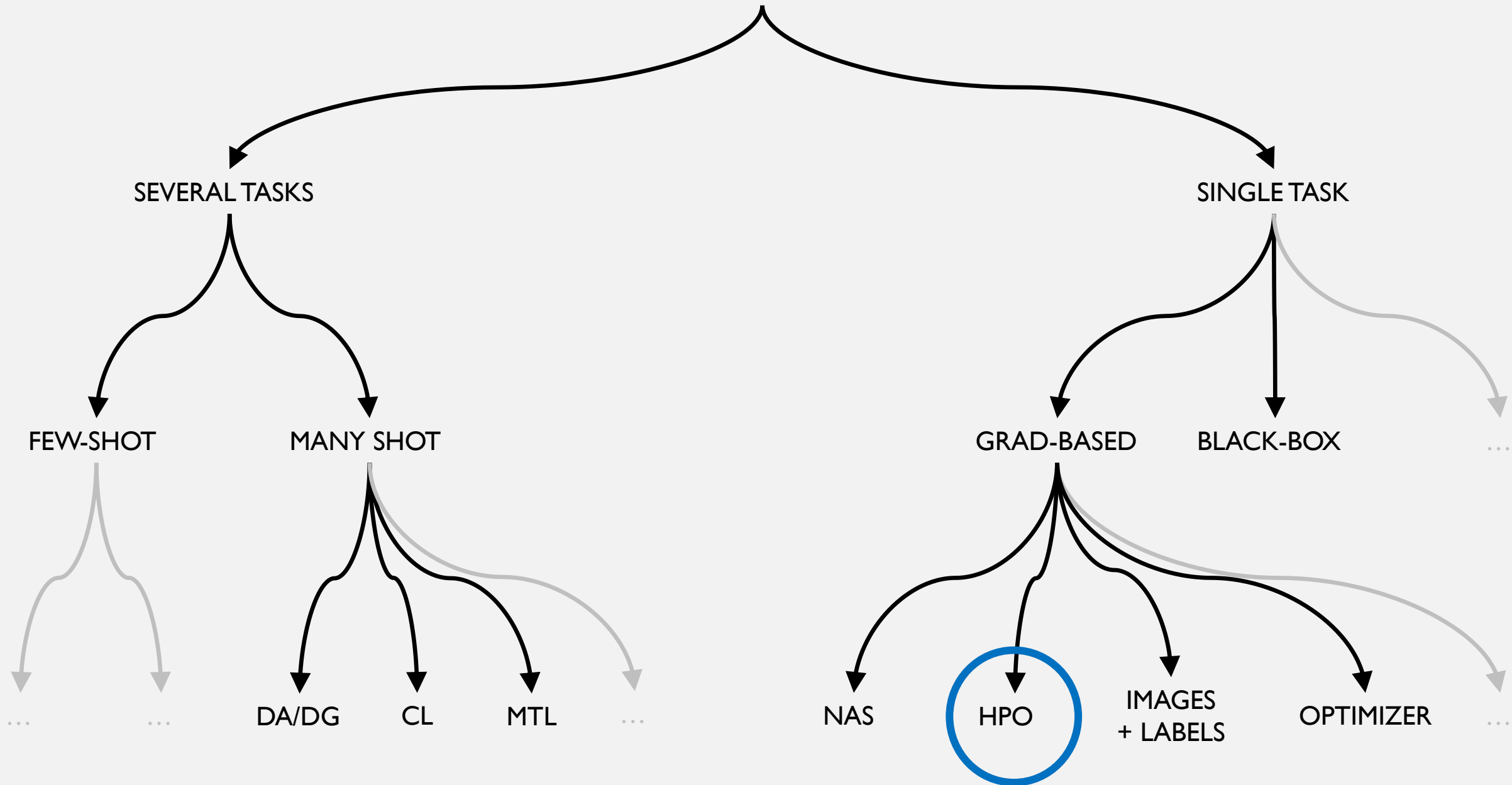
- Meta-learning & HPO
- Critique of BPTT and greediness for long inner loops
- Forward mode differentiation
- Hyperparameter sharing

META-LEARNING

- Learning to learn
- Reduces human “algorithm engineering”
- *Meta-Learning in Neural Networks: A Survey*

Hospedales et al.

META-LEARNING



GRADIENT BASED HPO

- Useful to formalize as a constrained optimization

$$\lambda^* = \operatorname{argmin}_{\lambda} \mathcal{L}_{\text{val}}(\theta_T(\lambda), \mathcal{D}_{\text{val}})$$

$$\text{subject to } \theta_{t+1} = \Phi(\mathcal{L}_{\text{train}}(\theta_t(\lambda), \mathcal{D}_{\text{train}}), \lambda)$$

Outer Loop

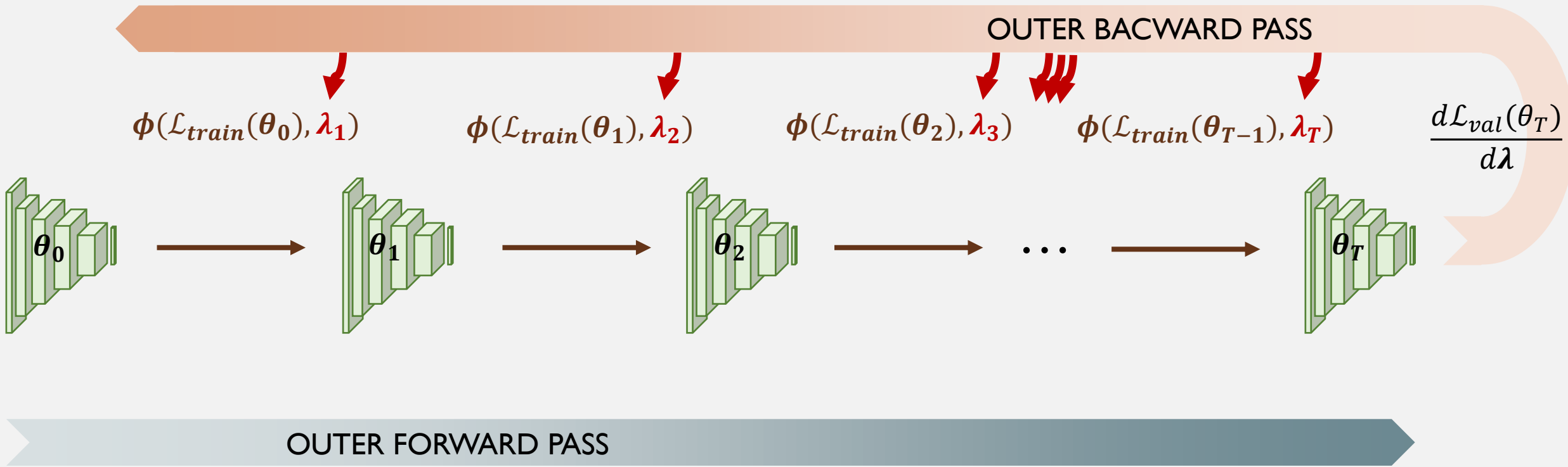
Inner Loop

Learn **hyperparameters** λ

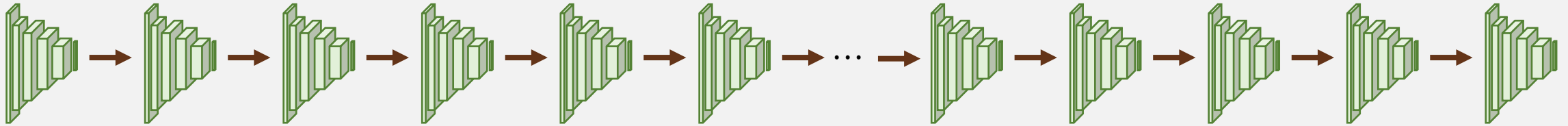
... such that the **network weights** θ_T after T steps of **optimizer** ϕ on the train loss
... also minimize the validation loss

Ultimately what we want is the **hypergradient**: $\frac{d\mathcal{L}_{\text{val}}(\theta_T)}{d\lambda}$

USING BPTT




CHALLENGES WHEN T IS LARGE: MEMORY



- In BPTT, you need to store each inner step in memory
- Ok for few-shot learning (e.g. MAML where $T \sim 5$ inner steps)
- But for problems like CIFAR-10, we need $T \sim 10^4$ inner steps



CHALLENGES WHEN T IS LARGE: GRADIENT DEGRADATION

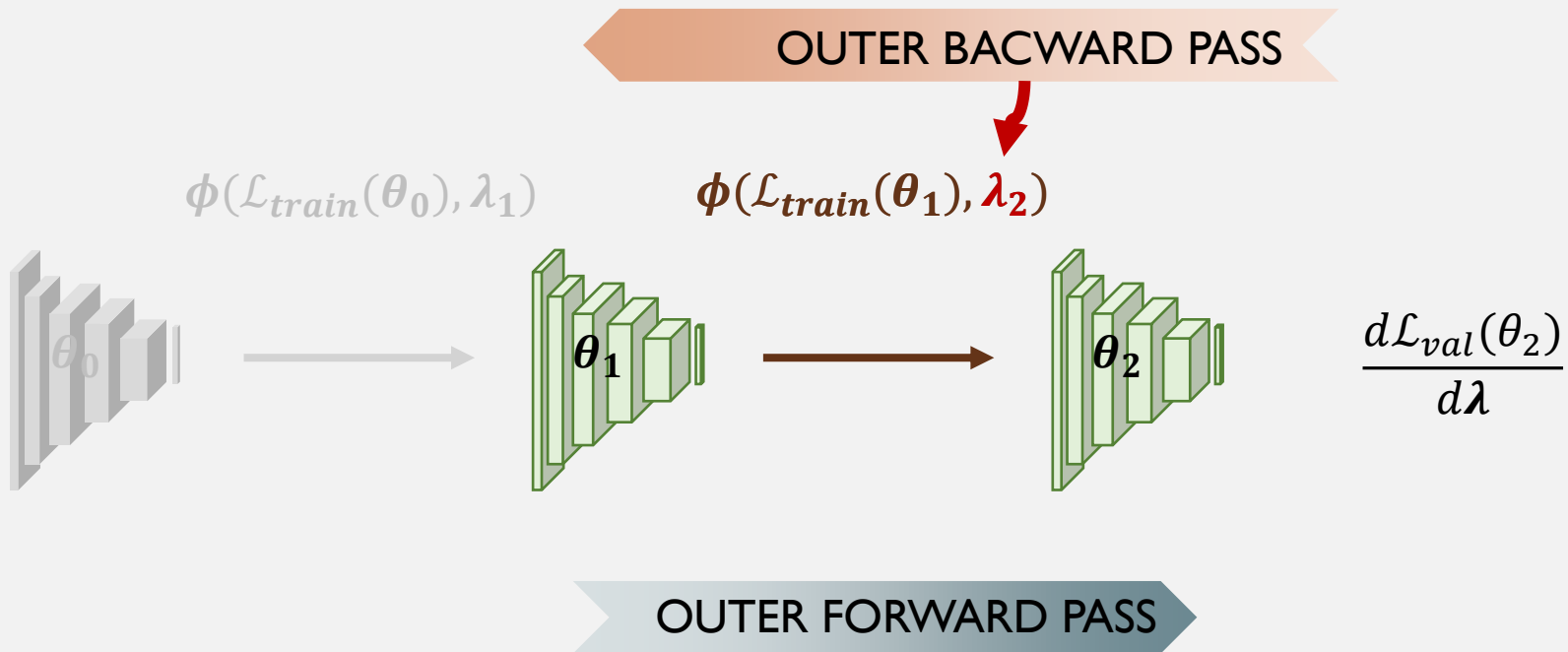
- Broad issue that arises when a parameter influences some scalar in a chaotic fashion (such as a long composition of non-linear ops)
- Vanishing + exploding hypergradients = high variance 

SOLUTION = GREEDINESS..?





- We can take H steps before updating hyperparameters, with $H \ll T$
- H is the *horizon*

SOLUTION = GREEDINESS..?

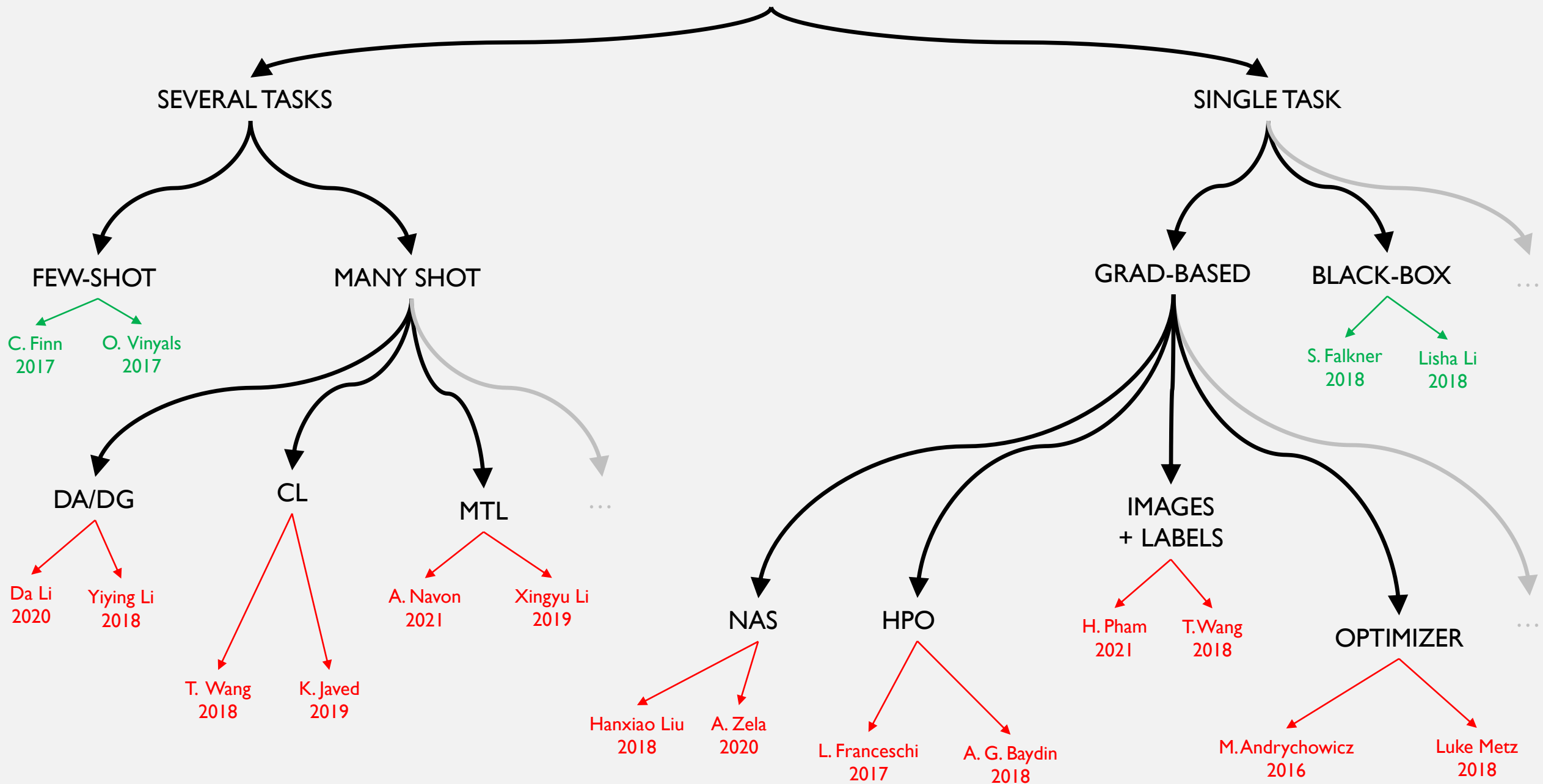
- For instance if $H = 1$:



SOLUTION = GREEDINESS..?

- Solves memory issue of BPTT 
- Solves gradient degradation 
- Improves computational cost 
- Solves for the wrong objective 

META-LEARNING



SOLUTION = FORWARD MODE DIFFERENTIATION...?

- Calculate gradient components during forward pass
- Applying it to HPO:

reverse mode

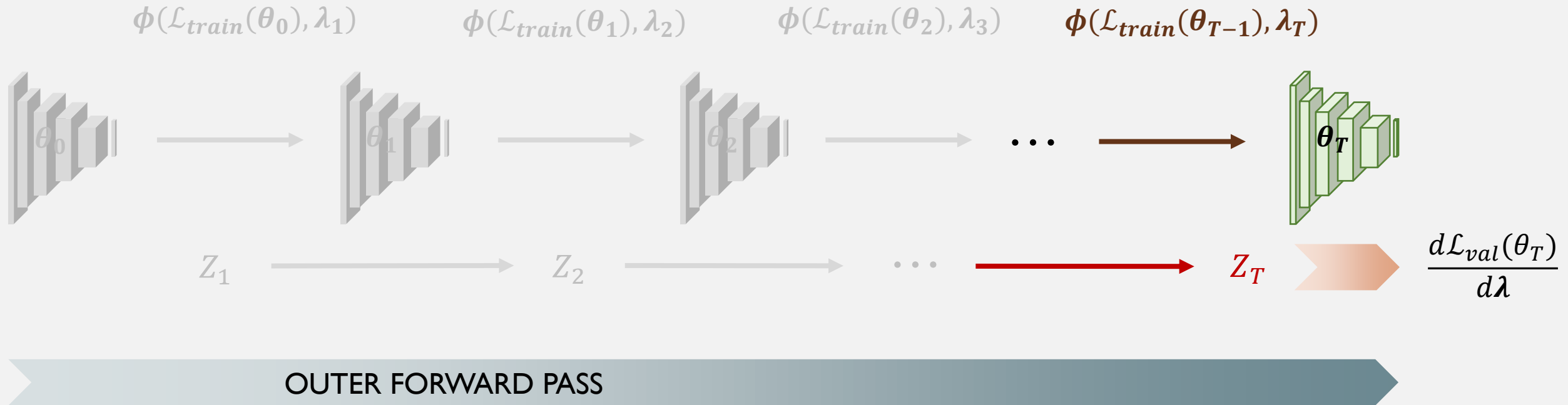
$$\frac{d\mathcal{L}_{\text{val}}}{d\lambda} = \frac{\partial \mathcal{L}_{\text{val}}}{\partial \theta_T} \frac{d\theta_T}{d\lambda}$$

forward mode

$$\frac{d\theta_t}{d\lambda} = \frac{\partial \theta_t}{\partial \theta_{t-1}} \Big|_{\lambda} \frac{d\theta_{t-1}}{d\lambda} + \frac{\partial \theta_t}{\partial \lambda} \Big|_{\theta_{t-1}}$$

$$\mathbf{Z}_t = \mathbf{A}_t \mathbf{Z}_{t-1} + \mathbf{B}_t$$

FORWARD MODE DIFFERENTIATION



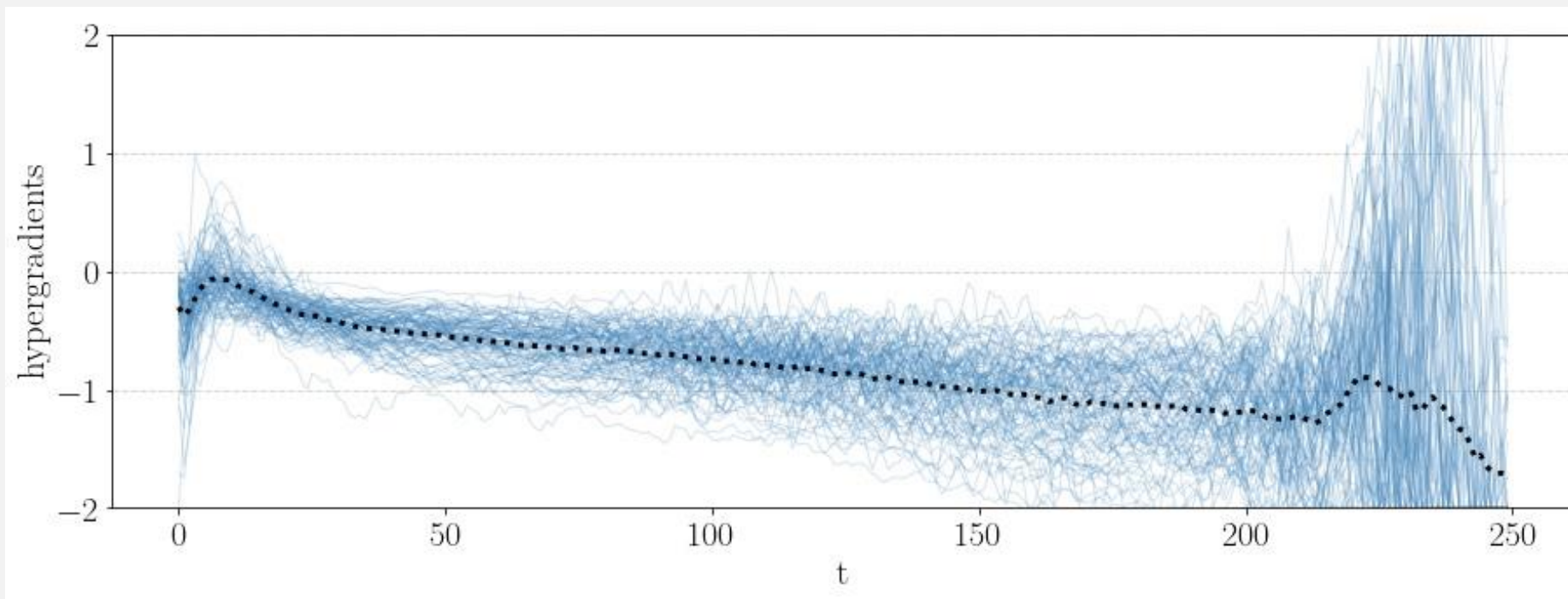
FORWARD MODE DIFFERENTIATION

- Memory cost constant with horizon size
- Solves for correct objective
- Memory cost scales poorly with size of λ
- Doesn't solve gradient degradation



SOLVING GRADIENT DEGRADATION WITH ENSEMBLE AVERAGING

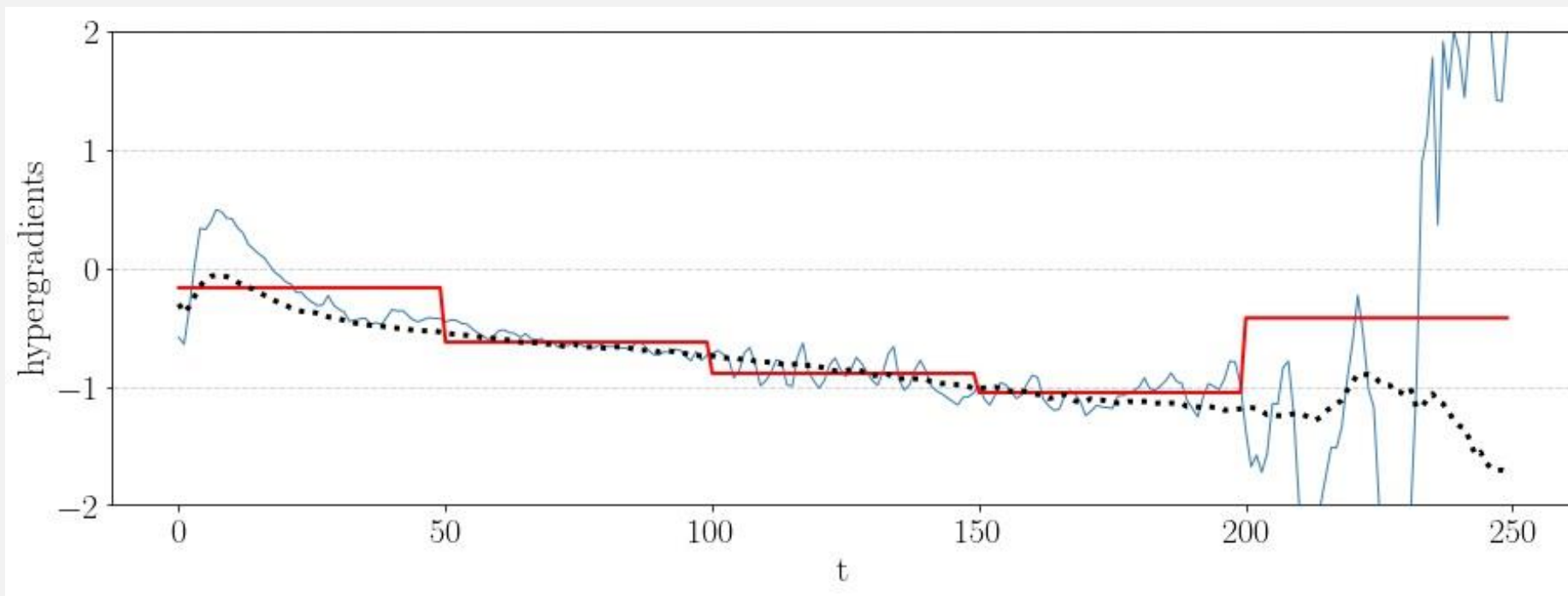
- Each seed gives different hypergradients



- Ensemble average = too expensive in compute / memory 

SOLVING GRADIENT DEGRADATION WITH TIME AVERAGING

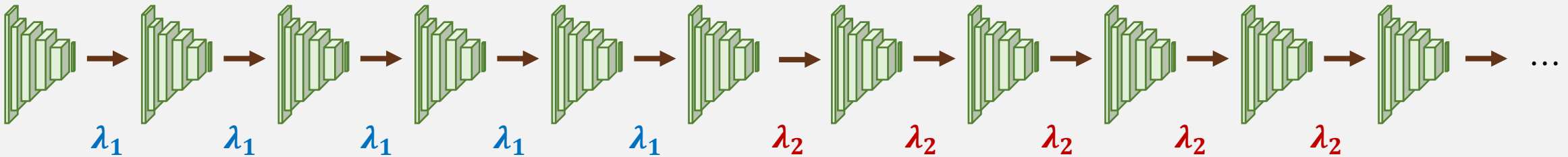
- Time averaging averages hypergradients in window of size W



- Cheap, and easily achieved by sharing hyperparameters



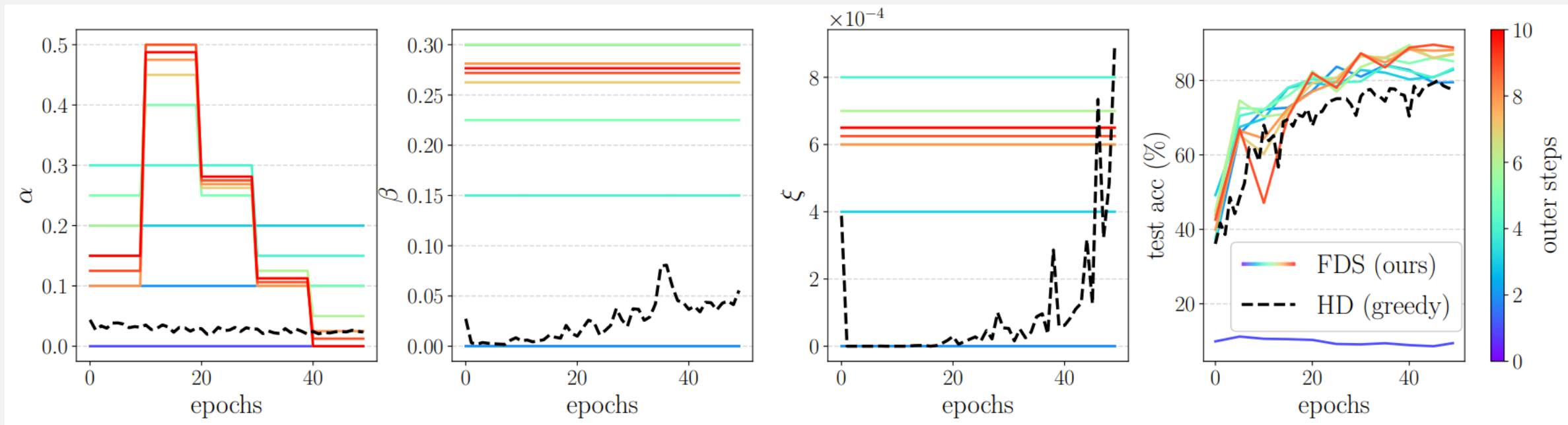
SOLVING GRADIENT DEGRADATION WITH TIME AVERAGING



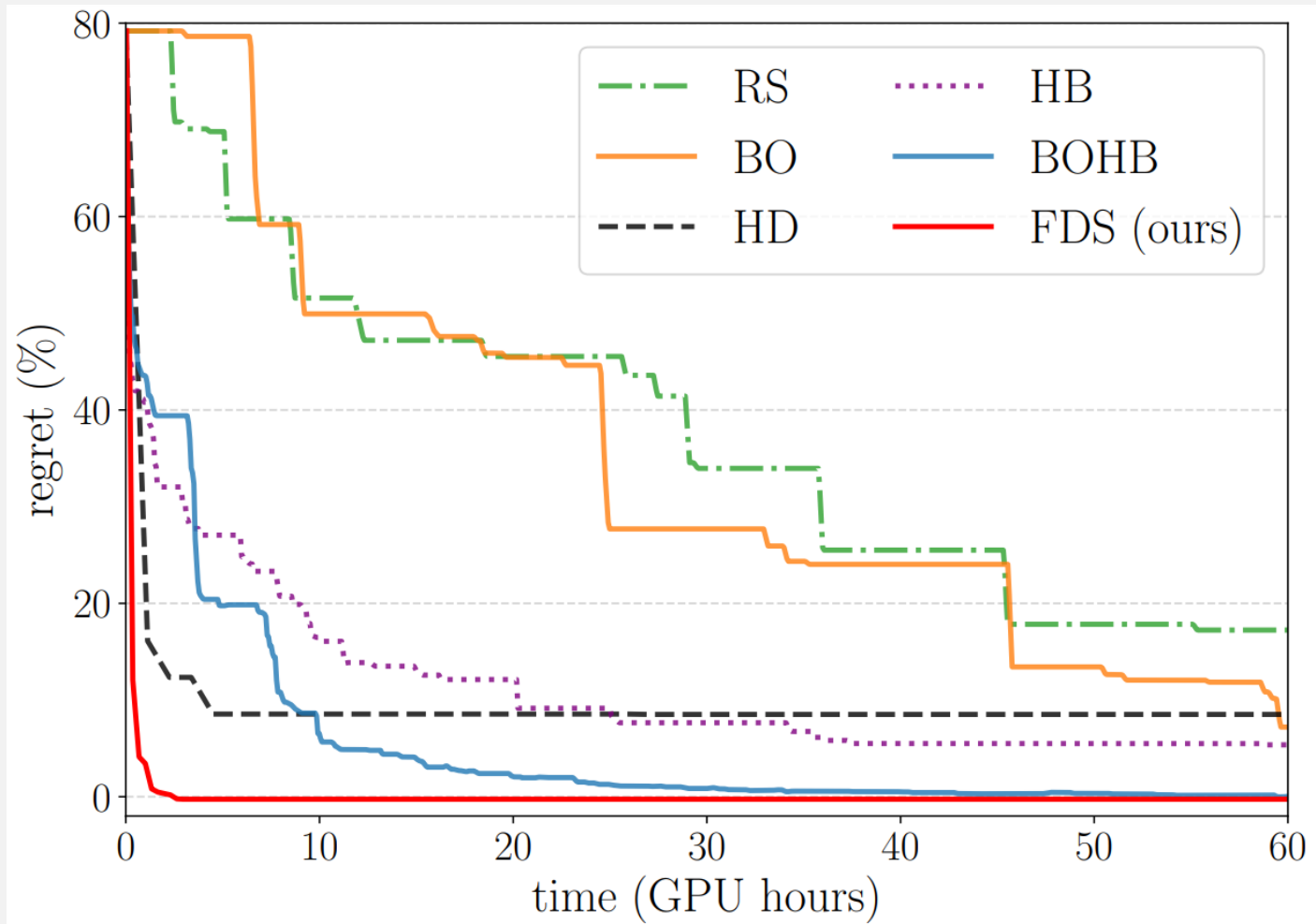
- What window size should we use?
- Optimal W = tradeoff between **variance reduction** and **bias increase**

$$\text{MSE}_W \leq \frac{\text{MSE}_1}{W} + \frac{\epsilon^2(W^2 - 1)}{12}$$

RESULTS: LEARNING SCHEDULES ON CIFAR10



RESULTS: REGRET



CONCLUSION

- You can differentiate through $\sim 10^4$ gradient steps by sharing contiguous hyperparameters in forward mode differentiation
- Future work:
 - Use automatic forward mode differentiation
 - Tackle harder problems/datasets