Amortized Proximal Optimization: Meta-Learning a Parametric Update Rule via a Proximal Meta-Objective

Juhan Bae*, Paul Vicol*, Jeff Z. Haochen, Roger Grosse

Slides by Paul Vicol





Outline

- Background: Proximal Point Method
 - Different types of discrepancy measures (*weight-space* and *function-space*)
 - Connections Between Proximal Optimization and 2nd-Order Optimization
- Amortized Proximal Optimization (APO)
 - Algorithm
 - Computation & Memory Cost
- APO for Structured Preconditioner Adaptation
- APO for Learning Rate Adaptation
- Experimental Results

Many algorithms in machine learning can be interpreted as approximations to an idealized algorithm called the *proximal point method (PPM)*

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \arg\min_{\mathbf{u} \in \mathbb{R}^m} \mathcal{J}_{\mathcal{B}^{(t)}}(\mathbf{u}) + \lambda D(\mathbf{u}, \boldsymbol{\theta}^{(t)})$$

Minimize the loss on the current minibatch

While staying close to the current iterate

• Many algorithms in machine learning can be interpreted as approximations to an idealized algorithm called the *proximal point method (PPM)*

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \arg\min_{\mathbf{u} \in \mathbb{R}^{m}} \mathcal{J}_{\mathcal{B}^{(t)}}(\mathbf{u}) + \lambda D(\mathbf{u}, \boldsymbol{\theta}^{(t)})$$

$$\underbrace{\text{Minimize the loss on}}_{\text{the current minibatch}} \qquad \text{While staying close to}$$

$$\underbrace{\text{While staying close to}}_{\text{the current iterate}}$$

- How do we measure proximity to the previous iterate, $D(\mathbf{u}, \boldsymbol{\theta}^{(t)})$?
- Several classic algorithms can be interpreted as *balancing between two types of proximity*:

$$\begin{split} & \frac{\text{Function-Space Discrepancy (FSD)}}{\mathcal{D}_{\mathrm{F}}(\mathbf{u}, \boldsymbol{\theta}^{(t)}) = \mathbb{E}_{(\tilde{\mathbf{x}}, \cdot) \sim \mathcal{D}} \left[\rho(f(\tilde{\mathbf{x}}; \mathbf{u}), f(\tilde{\mathbf{x}}; \boldsymbol{\theta}^{(t)})) \right]} \end{split}$$

Weight-Space Discrepancy (WSD)

 $D_{\mathrm{W}}(\mathbf{u}, \boldsymbol{\theta}^{(t)}) = 1/2 \|\mathbf{u} - \boldsymbol{\theta}^{(t)}\|_{2}^{2}$

• To gain intuition, let's see what happens when we minimize the proximal objective *exactly*.

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \arg\min_{\mathbf{u} \in \mathbb{R}^m} \mathcal{J}_{\mathcal{B}^{(t)}}(\mathbf{u}) + \lambda_{\text{FSD}} \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathcal{D}}[D_{\text{F}}(\mathbf{u}, \boldsymbol{\theta}^{(t)}, \tilde{\mathbf{x}})] + \lambda_{\text{WSD}} D_{\text{W}}(\mathbf{u}, \boldsymbol{\theta}^{(t)})$$

• To gain intuition, let's see what happens when we minimize the proximal objective *exactly*.

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \arg\min_{\mathbf{u} \in \mathbb{R}^m} \mathcal{J}_{\mathcal{B}^{(t)}}(\mathbf{u}) + \lambda_{\text{FSD}} \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathcal{D}}[D_{\mathbf{F}}(\mathbf{u}, \boldsymbol{\theta}^{(t)}, \tilde{\mathbf{x}})] + \lambda_{\text{WSD}} D_{\text{W}}(\mathbf{u}, \boldsymbol{\theta}^{(t)})$$



Makes a small change to the parameters but a global change to the outputs.

Fits the current example but also changes the predictions on other data.

• To gain intuition, let's see what happens when we minimize the proximal objective *exactly*.

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \arg\min_{\mathbf{u} \in \mathbb{R}^m} \mathcal{J}_{\mathcal{B}^{(t)}}(\mathbf{u}) + \lambda_{\mathrm{FSD}} \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathcal{D}}[D_{\mathrm{F}}(\mathbf{u}, \boldsymbol{\theta}^{(t)}, \tilde{\mathbf{x}})] + \lambda_{\mathrm{WSD}} D_{\mathrm{W}}(\mathbf{u}, \boldsymbol{\theta}^{(t)})$$



Makes a small change to the parameters but a global change to the outputs.

Fits the current example but also changes the predictions on other data. Makes a *sharp change to the outputs* to fit the current example while maintaining the previous outputs on other data.

• To gain intuition, let's see what happens when we minimize the proximal objective *exactly*.

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \arg\min_{\mathbf{u} \in \mathbb{R}^m} \mathcal{J}_{\mathcal{B}^{(t)}}(\mathbf{u}) + \lambda_{\text{FSD}} \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathcal{D}}[D_{\text{F}}(\mathbf{u}, \boldsymbol{\theta}^{(t)}, \tilde{\mathbf{x}})] + \lambda_{\text{WSD}} D_{\text{W}}(\mathbf{u}, \boldsymbol{\theta}^{(t)})$$



Makes a small change to the parameters but a global change to the outputs.

Fits the current example but also changes the predictions on other data. Makes a *sharp change to the outputs* to fit the current example while maintaining the previous outputs on other data. Including both FSD and WSD terms yields a *smoother update to the function*.

Connections to Second-Order Optimization

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \arg\min_{\mathbf{u} \in \mathbb{R}^m} \mathcal{J}_{\mathcal{B}^{(t)}}(\mathbf{u}) + \lambda_{\text{FSD}} \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathcal{D}}[D_{\text{F}}(\mathbf{u}, \boldsymbol{\theta}^{(t)}, \tilde{\mathbf{x}})] + \lambda_{\text{WSD}} D_{\text{W}}(\mathbf{u}, \boldsymbol{\theta}^{(t)})$$

- Minimizing the proximal objective exactly in each iteration is *uneconomical*
- Various first- and second-order optimization algorithms can be interpreted as *minimizing approximations of the proximal objective*
 - Taking 1st or 2nd order Taylor expansions of the loss or FSD terms allows for *closed-form solutions*

Method	Loss Approx.	\mathbf{FSD}	\mathbf{WSD}
Gradient Descent	1^{st} -order	-	1
Hessian-Free	2^{nd} -order	_	\checkmark
Natural Gradient	1^{st} -order	2^{nd} -order	×
Proximal Point	Exact	Exact	1

Amortized Proximal Optimization (APO)

• Consider an *update rule* u parameterized by a vector of *optimization parameters* ϕ

$$\boldsymbol{\theta}^{(t+1)} \leftarrow u(\boldsymbol{\theta}^{(t)}, \boldsymbol{\phi}, \mathcal{B}^{(t)})$$

• Examples:

$$u_{\text{SGD}}(\boldsymbol{\theta}, \eta, \mathcal{B}) = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \mathcal{J}_{\mathcal{B}}(\boldsymbol{\theta})$$
$$u_{\text{Precond}}(\boldsymbol{\theta}, \mathbf{P}, \mathcal{B}) = \boldsymbol{\theta} - \mathbf{P} \nabla_{\boldsymbol{\theta}} \mathcal{J}_{\mathcal{B}}(\boldsymbol{\theta})$$

Amortized Proximal Optimization (APO)

• Consider an *update rule* u parameterized by a vector of *optimization parameters* ϕ

$$\boldsymbol{\theta}^{(t+1)} \leftarrow u(\boldsymbol{\theta}^{(t)}, \boldsymbol{\phi}, \mathcal{B}^{(t)})$$

• Examples:

$$u_{\text{SGD}}(\boldsymbol{\theta}, \eta, \mathcal{B}) = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \mathcal{J}_{\mathcal{B}}(\boldsymbol{\theta})$$
$$u_{\text{Precond}}(\boldsymbol{\theta}, \mathbf{P}, \mathcal{B}) = \boldsymbol{\theta} - \mathbf{P} \nabla_{\boldsymbol{\theta}} \mathcal{J}_{\mathcal{B}}(\boldsymbol{\theta})$$

• APO tunes the optimization parameters phi to *minimize a proximal meta-objective* $\mathcal{Q}(m{\phi})$:

$$\mathcal{Q}(\boldsymbol{\phi}) = \mathbb{E}_{\mathcal{B}\sim\mathcal{D}} \Big[\mathcal{J}_{\mathcal{B}}(u(\boldsymbol{\theta},\boldsymbol{\phi},\mathcal{B})) + \lambda_{\mathrm{FSD}} \mathbb{E}_{(\tilde{\mathbf{x}},\cdot)\sim\mathcal{D}} \left[D_{\mathrm{F}}(u(\boldsymbol{\theta},\boldsymbol{\phi},\mathcal{B}),\boldsymbol{\theta},\tilde{\mathbf{x}}) \right] + \frac{\lambda_{\mathrm{WSD}}}{2} \| u(\boldsymbol{\theta},\boldsymbol{\phi},\mathcal{B}) - \boldsymbol{\theta} \|^2 \Big]$$

• By adapting a parametric update rule, we can *amortize the cost of minimizing the proximal objective* over the course of training

APO Algorithm



APO Algorithm



Algorithm 1 Amortized Proximal Optimization (APO) — Meta-Learning the Optimization Parameters ϕ

Require: θ (initial model parameters), ϕ (initial optimization parameters), K (meta-update interval), α (meta-LR)

Require: λ_{WSD} (weight-space discrepancy term weighting), λ_{FSD} (function-space discrepancy term weighting)

while not converged, iteration t do

 $\begin{array}{lll} \mathcal{B} \sim \mathcal{D}_{\text{train}} & \triangleright \text{ Sample mini-batch to compute the gradient and loss term} \\ \textbf{if} \ t \ \operatorname{mod} \ K = 0 \ \textbf{then} & \triangleright \text{ Perform meta-update every } K \ \textbf{iterations} \\ \mathcal{B}' \sim \mathcal{D}_{\text{train}} & \triangleright \text{ Sample additional mini-batch to compute the FSD term} \\ \boldsymbol{\theta}'(\boldsymbol{\phi}) := u(\boldsymbol{\theta}, \boldsymbol{\phi}, \mathcal{B}) & \triangleright \text{ Compute the 1-step lookahead parameters} \\ \mathcal{Q}(\boldsymbol{\phi}) := \mathcal{J}_{\mathcal{B}} \left(\boldsymbol{\theta}'(\boldsymbol{\phi})\right) + \lambda_{\text{FSD}} / |\mathcal{B}'| \sum_{(\tilde{\mathbf{x}}, \cdot) \in \mathcal{B}'} \mathcal{D}_{\text{F}}(\boldsymbol{\theta}'(\boldsymbol{\phi}), \boldsymbol{\theta}, \tilde{\mathbf{x}}) + \lambda_{\text{WSD}} / 2 \| \boldsymbol{\theta}'(\boldsymbol{\phi}) - \boldsymbol{\theta} \|_{2}^{2} \triangleright \text{ Meta-objective} \\ \boldsymbol{\phi} \leftarrow \boldsymbol{\phi} - \alpha \nabla_{\boldsymbol{\phi}} \mathcal{Q}(\boldsymbol{\phi}) & \triangleright \text{ Update optimizer parameters (e.g. LR or preconditioner)} \\ \textbf{end if} \\ \boldsymbol{\theta} \leftarrow u(\boldsymbol{\theta}, \boldsymbol{\phi}, \mathcal{B}) & \triangleright \text{ Update model parameters} \\ \textbf{end while} \end{array}$

Computation and Memory Cost of APO



- **<u>Computation</u>**: Computing the meta-gradient requires 3 forward passes + a backward pass through the 1-step unrolled computation graph.
 - But we only perform a meta-update once every K iterations
- <u>Memory:</u> APO requires 2x the model memory for the 1-step unroll

- APO can be used to learn the *preconditioning matrix for second-order optimization*
- How should we *parameterize* the preconditioner?

Two desiderata:

We want the preconditioner to be *positive semi-definite*, to maintain descent directions



We want the preconditioner to be tractable to store and compute with

- APO can be used to learn the *preconditioning matrix for second-order optimization*
- How should we *parameterize* the preconditioner?

Two desiderata: 1 We want the preconditioner to be *positive semi-definite*, to maintain descent directions 2 We want the preconditioner to be tractable to store and compute with

Parameterization	Positive Semi-Definite	Tractable to Store in Memory
Р	*	*

- APO can be used to learn the *preconditioning matrix for second-order optimization*
- How should we *parameterize* the preconditioner?

Two desiderata: 1 We want the preconditioner to be *positive semi-definite*, to maintain descent directions 2 We want the preconditioner to be tractable to store and compute with

Parameterization	Positive Semi-Definite	Tractable to Store in Memory
Р	*	*
$\mathbf{M}\mathbf{M}^\top$	\checkmark	*

- APO can be used to learn the *preconditioning matrix for second-order optimization*
- How should we *parameterize* the preconditioner?

Two desiderata: We want the preconditioner to be *positive semi-definite*, to maintain descent directions We want the preconditioner to be tractable to store and compute with

Parameterization	Positive Semi-Definite	Tractable to Store in Memory		
Р	*	*		
$\mathbf{M}\mathbf{M}^\top$	\checkmark	*		
$\mathbf{P}_{\mathrm{S}} = (\mathbf{A} \otimes \mathbf{B}) \mathrm{diag}(\mathrm{vec}(\mathbf{S}))^2 (\mathbf{A} \otimes \mathbf{B})$	Т			

• In practice, we use the EKFAC-structured preconditioner:

$$\mathbf{P}_{\mathrm{S}} = (\mathbf{A} \otimes \mathbf{B}) \mathrm{diag}(\mathrm{vec}(\mathbf{S}))^2 (\mathbf{A} \otimes \mathbf{B})^{\top}$$

- EKFAC uses *complicated covariance estimation and matrix inversion* to construct the block matrices
- In contrast, APO *directly meta-learns the matrices*
- APO does not require inverting (or performing eigendecompositions of) the block matrices
- \rightarrow Our structured representation incurs *less computation per iteration than EKFAC.*

APO Can Recover Classic 2nd-Order Methods

• If we make the same assumptions used to derive natural gradient, then the optimal preconditioner P that *minimizes the one-step lookahead meta-objective corresponds to the damped natural gradient update*:

 $\mathbf{P}^{\star} = (\lambda_{\rm FSD}\mathbf{G} + \lambda_{\rm WSD}\mathbf{I})^{-1}$

Theorem 1. Consider an approximation $\hat{Q}(\mathbf{P})$ to the meta-objective (Eq. 8) where the loss term is linearized around the current weights $\boldsymbol{\theta}$ and the FSD term is replaced by its second-order approximation around $\boldsymbol{\theta}$. Denote the gradient on a mini-batch as $\mathbf{g} = \nabla_{\boldsymbol{\theta}} \mathcal{J}_{\mathcal{B}}(\boldsymbol{\theta})$, and assume that the second moment matrix $\mathbb{E}_{\mathcal{B}\sim\mathcal{D}} \left[\mathbf{gg}^{\top}\right]$ is non-singular. Then, the preconditioning matrix which minimizes $\hat{\mathcal{Q}}$ is given by $\mathbf{P}^{\star} = (\lambda_{FSD}\mathbf{G} + \lambda_{WSD}\mathbf{I})^{-1}$, where \mathbf{G} denotes the Hessian of the FSD evaluated at $\boldsymbol{\theta}$.

APO Can Recover the KFAC Update

Corollary 2. Suppose that (1) the assumptions for Theorem 1 are satisfied, (2) the FSD term measures the KL divergence, and (3) $\lambda_{WSD} = 0$ and $\lambda_{FSD} = 1$. Moreover, suppose that the parameters $\boldsymbol{\theta}$ satisfy the KFAC assumptions listed in Appendix H. Then, the optimal solution to the approximate meta-objective recovers the KFAC update, which can be represented using the structured preconditioner in Eq. 9.

Experiments: Poorly-Conditioned Problems

• Using APO to tune the preconditioner works well on *poorly-conditioned tasks where* second-order optimizers are typically required



Experiments: UCI Regression Tasks

- APO-Precond outperforms 1st and 2nd-order baselines (including Shampoo and KFAC) on UCI regression tasks
 - Typically, APO-Precond *converges more quickly and reaches lower training loss*



Experiments: Architectures on CIFAR-10/100

- We evaluated the generalization of a variety of APO-trained models on CIFAR-10/100
 - Tuned LR, weight decay, etc. separately for each task.



Low-Precision (16-bit) Training

- Low-precision training presents a challenge for second-order optimizers such as KFAC
 - These rely on *matrix inverses that may be sensitive to quantization noise*.
- APO does not require inversion, and remains stable

Task	Model	\mathbf{SGDm}	KFAC	APO-P
CIFAR-10 CIFAR-10 CIFAR-100	LeNet ResNet-18 ResNet-18	$75.65 \\ 94.15 \\ 73.53$	$74.95 \\92.72 \\73.12$	$77.25 \\94.79 \\75.47$

LR Adaptation: MNIST

- APO can be used to tune the global learning rate for a variety of *base optimizers*
 - Including SGD, SGDm, RMSprop, Adam
- The lambda hyperparameters for APO need to be searched over, but *each lambda setting yields a learning rate schedule*
- Here, we tuned the LR for training an MLP on MNIST, with SGDm and RMSprop



LR Adaptation: MNIST

• The APO-adapted LR schedules usually *outperform optimal fixed LRs, and are competitive with manual step decay schedules*



LR Adaptation: MNIST

• The APO-adapted LR schedules usually *outperform optimal fixed LRs, and are competitive with manual step decay schedules*



	CIFAR-10 ResNet-32		${f CIFAR-10}\ {f ResNet-34}$		CIFAR-10 WRN 28-10		CIFAR-100 WRN 28-10					
	Fixed	Decay	APO	Fixed	Decay	APO	Fixed	Decay	APO	Fixed	Decay	APO
\mathbf{SGD}	90.07	93.30	92.71	93.00	93.54	94.27	93.38	94.86	94.85	76.29	77.92	76.87
\mathbf{SGDm}	89.40	93.34	92.75	92.99	95.08	94.47	93.46	95.98	95.50	74.81	81.01	79.33
RMSprop	89.84	91.94	91.28	92.87	93.87	93.97	92.91	93.60	94.22	72.06	76.06	74.17
$\mathbf{A}\mathbf{d}\mathbf{a}\mathbf{m}$	90.45	92.26	91.81	93.23	94.12	93.80	92.81	94.04	93.83	72.01	75.53	76.33

LR Adaptation: Robustness to the Initial LR



LR Adaptation: Robustness to the Initial LR



On Short-Horizon Bias

- APO is a greedy method—it only considers the effect of the optimization parameters through a one-step lookahead → Does APO suffer from short-horizon bias?
- Original short-horizon bias meta-objective: $\mathbb{E}_{\mathcal{B}\sim\mathcal{D}_{train}}[\mathcal{J}_{\mathcal{B}}(u(\theta,\phi,\mathcal{B}))]$

On Short-Horizon Bias

- APO is a greedy method—it only considers the effect of the optimization parameters through a one-step lookahead → Does APO suffer from short-horizon bias?
- Original short-horizon bias meta-objective: $\mathbb{E}_{\mathcal{B}\sim\mathcal{D}_{train}}[\mathcal{J}_{\mathcal{B}}(u(\theta,\phi,\mathcal{B}))]$

$$\begin{array}{|c|c|} \hline \mathbf{Q}(\boldsymbol{\phi}) = \mathcal{J}_{\mathcal{B}}(u(\boldsymbol{\theta},\boldsymbol{\phi},\mathcal{B})) + \lambda_{\mathrm{FSD}}\mathbb{E}_{(\tilde{\mathbf{x}},\cdot)\sim\mathcal{D}}[D_{\mathrm{F}}(u(\boldsymbol{\theta},\boldsymbol{\phi},\mathcal{B}),\boldsymbol{\theta},\tilde{\mathbf{x}})] + \frac{\lambda_{\mathrm{WSD}}}{2}||u(\boldsymbol{\theta},\boldsymbol{\phi},\mathcal{B}) - \boldsymbol{\theta}||_{2}^{2} \\ \hline \mathbf{Q}(\boldsymbol{\phi}) = \mathcal{J}_{\mathcal{B}'}(u(\boldsymbol{\theta},\boldsymbol{\phi},\mathcal{B})) + \lambda_{\mathrm{FSD}}\mathbb{E}_{(\tilde{\mathbf{x}},\cdot)\sim\mathcal{D}}[D_{\mathrm{F}}(u(\boldsymbol{\theta},\boldsymbol{\phi},\mathcal{B}),\boldsymbol{\theta},\tilde{\mathbf{x}})] + \frac{\lambda_{\mathrm{WSD}}}{2}||u(\boldsymbol{\theta},\boldsymbol{\phi},\mathcal{B}) - \boldsymbol{\theta}||_{2}^{2} \\ \end{array}$$

On Short-Horizon Bias

- APO is a greedy method—it only considers the effect of the optimization parameters through a one-step lookahead → Does APO suffer from short-horizon bias?
- Original short-horizon bias meta-objective: $\mathbb{E}_{\mathcal{B}\sim\mathcal{D}_{train}}[\mathcal{J}_{\mathcal{B}}(u(\boldsymbol{\theta},\boldsymbol{\phi},\mathcal{B}))]$

$$\begin{array}{|c|c|} \hline & \mathcal{Q}(\phi) = \mathcal{J}_{\mathcal{B}}(u(\theta,\phi,\mathcal{B})) + \lambda_{\mathrm{FSD}} \mathbb{E}_{(\tilde{\mathbf{x}},\cdot)\sim\mathcal{D}}[D_{\mathrm{F}}(u(\theta,\phi,\mathcal{B}),\theta,\tilde{\mathbf{x}})] + \frac{\lambda_{\mathrm{WSD}}}{2} ||u(\theta,\phi,\mathcal{B}) - \theta||_{2}^{2} \\ \hline & \mathbf{Q}(\phi) = \mathcal{J}_{\mathcal{B}'}(u(\theta,\phi,\mathcal{B})) + \lambda_{\mathrm{FSD}} \mathbb{E}_{(\tilde{\mathbf{x}},\cdot)\sim\mathcal{D}}[D_{\mathrm{F}}(u(\theta,\phi,\mathcal{B}),\theta,\tilde{\mathbf{x}})] + \frac{\lambda_{\mathrm{WSD}}}{2} ||u(\theta,\phi,\mathcal{B}) - \theta||_{2}^{2} \\ \hline & \mathbf{Q}(\phi) = \mathcal{J}_{\mathcal{B}'}(u(\theta,\phi,\mathcal{B})) + \lambda_{\mathrm{FSD}} \mathbb{E}_{(\tilde{\mathbf{x}},\cdot)\sim\mathcal{D}}[D_{\mathrm{F}}(u(\theta,\phi,\mathcal{B}),\theta,\tilde{\mathbf{x}})] + \frac{\lambda_{\mathrm{WSD}}}{2} ||u(\theta,\phi,\mathcal{B}) - \theta||_{2}^{2} \\ \hline & \mathbf{Q}(\phi) = \mathcal{J}_{\mathcal{B}'}(u(\theta,\phi,\mathcal{B})) + \lambda_{\mathrm{FSD}} \mathbb{E}_{(\tilde{\mathbf{x}},\cdot)\sim\mathcal{D}}[D_{\mathrm{F}}(u(\theta,\phi,\mathcal{B}),\theta,\tilde{\mathbf{x}})] + \frac{\lambda_{\mathrm{WSD}}}{2} ||u(\theta,\phi,\mathcal{B}) - \theta||_{2}^{2} \\ \hline & \mathbf{Q}(\phi) = \mathcal{J}_{\mathcal{B}'}(u(\theta,\phi,\mathcal{B})) + \lambda_{\mathrm{FSD}} \mathbb{E}_{(\tilde{\mathbf{x}},\cdot)\sim\mathcal{D}}[D_{\mathrm{F}}(u(\theta,\phi,\mathcal{B}),\theta,\tilde{\mathbf{x}})] + \frac{\lambda_{\mathrm{WSD}}}{2} ||u(\theta,\phi,\mathcal{B}) - \theta||_{2}^{2} \\ \hline & \mathbf{Q}(\phi) = \mathcal{J}_{\mathcal{B}'}(u(\theta,\phi,\mathcal{B})) + \lambda_{\mathrm{FSD}} \mathbb{E}_{(\tilde{\mathbf{x}},\cdot)\sim\mathcal{D}}[D_{\mathrm{F}}(u(\theta,\phi,\mathcal{B}),\theta,\tilde{\mathbf{x}})] + \frac{\lambda_{\mathrm{WSD}}}{2} ||u(\theta,\phi,\mathcal{B}) - \theta||_{2}^{2} \\ \hline & \mathbf{Q}(\phi) = \mathcal{J}_{\mathcal{B}'}(u(\theta,\phi,\mathcal{B})) + \lambda_{\mathrm{FSD}} \mathbb{E}_{(\tilde{\mathbf{x}},\cdot)\sim\mathcal{D}}[D_{\mathrm{F}}(u(\theta,\phi,\mathcal{B}),\theta,\tilde{\mathbf{x}})] + \frac{\lambda_{\mathrm{WSD}}}{2} ||u(\theta,\phi,\mathcal{B}) - \theta||_{2}^{2} \\ \hline & \mathbf{Q}(\phi) = \mathcal{J}_{\mathcal{B}'}(u(\theta,\phi,\mathcal{B})) + \frac{\lambda_{\mathrm{FSD}}}{2} ||u(\theta,\phi,\mathcal{B})| + \frac{\lambda_{\mathrm{FSD}}}{2} ||u(\theta,\phi,\mathcal{B})| + \frac{\lambda_{\mathrm{WSD}}}{2} ||u(\theta,\phi,\mathcal{B})$$



Thank you!