

Supplementary Document for “Online Ensemble Multi-kernel Learning Adaptive to Non-stationary and Adversarial Environments”

A Proof of Lemma 1

To prove Lemma 1, we introduce two intermediate lemmata as follows.

Lemma 5 Consider Assumptions 1 and 2 are satisfied, and define f_p^* as the best static solution in (18) with $\mathcal{F}_p := \{f | f(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{z}_p(\mathbf{x}), \forall \boldsymbol{\theta} \in \mathbb{R}^{2D}\}$. If $\{f_{p,t}(\mathbf{x}_t)\}$ denotes the sequence of estimates generated by online kernel learning algorithm with a pre-selected kernel p , the following bound holds true w.p.1, i.e.,

$$\sum_{t=1}^T \ell_t(f_{p,t}(\mathbf{x}_t)) - \sum_{t=1}^T \ell_t(f_p^*(\mathbf{x}_t)) \leq \frac{\|\boldsymbol{\theta}_p^*\|^2}{2\eta} + \frac{\eta L^2 T}{2} \quad (32)$$

where η is the learning rate, L is the Lipschitz constant in Assumption 2, and $\boldsymbol{\theta}_p^*$ is the corresponding weights supporting the best estimator $f_p^*(\mathbf{x}) = (\boldsymbol{\theta}_p^*)^\top \mathbf{z}_p(\mathbf{x})$.

Proof: Similar to the regret analysis of online gradient descent [25], using (11) for any fixed $\boldsymbol{\theta}$, we find

$$\begin{aligned} \|\boldsymbol{\theta}_{p,t+1} - \boldsymbol{\theta}\|^2 &= \|\boldsymbol{\theta}_{p,t} - \eta \nabla \mathcal{L}(\boldsymbol{\theta}_{p,t}^\top \mathbf{z}_p(\mathbf{x}_t), y_t) - \boldsymbol{\theta}\|^2 \\ &= \|\boldsymbol{\theta}_{p,t} - \boldsymbol{\theta}\|^2 + \eta^2 \|\nabla \mathcal{L}(\boldsymbol{\theta}_{p,t}^\top \mathbf{z}_p(\mathbf{x}_t), y_t)\|^2 - 2\eta \nabla^\top \mathcal{L}(\boldsymbol{\theta}_{p,t}^\top \mathbf{z}_p(\mathbf{x}_t), y_t)(\boldsymbol{\theta}_{p,t} - \boldsymbol{\theta}). \end{aligned} \quad (33)$$

Meanwhile, the convexity of the loss under Assumption 1 implies that

$$\mathcal{L}(\boldsymbol{\theta}_{p,t}^\top \mathbf{z}_p(\mathbf{x}_t), y_t) - \mathcal{L}(\boldsymbol{\theta}^\top \mathbf{z}_p(\mathbf{x}_t), y_t) \leq \nabla^\top \mathcal{L}(\boldsymbol{\theta}_{p,t}^\top \mathbf{z}_p(\mathbf{x}_t), y_t)(\boldsymbol{\theta}_{p,t} - \boldsymbol{\theta}). \quad (34)$$

Plugging (34) into (33) and rearranging terms yields

$$\mathcal{L}(\boldsymbol{\theta}_{p,t}^\top \mathbf{z}_p(\mathbf{x}_t), y_t) - \mathcal{L}(\boldsymbol{\theta}^\top \mathbf{z}_p(\mathbf{x}_t), y_t) \leq \frac{\|\boldsymbol{\theta}_{p,t} - \boldsymbol{\theta}\|^2 - \|\boldsymbol{\theta}_{p,t+1} - \boldsymbol{\theta}\|^2}{2\eta} + \frac{\eta}{2} \|\nabla \mathcal{L}(\boldsymbol{\theta}_{p,t}^\top \mathbf{z}_p(\mathbf{x}_t), y_t)\|^2. \quad (35)$$

Summing (35) over $t = 1, \dots, T$, with $f_{p,t}(\mathbf{x}_t) = \boldsymbol{\theta}_{p,t}^\top \mathbf{z}_p(\mathbf{x}_t)$, we arrive at

$$\begin{aligned} & \sum_{t=1}^T \left(\mathcal{L}(f_{p,t}(\mathbf{x}_t), y_t) - \mathcal{L}(\boldsymbol{\theta}^\top \mathbf{z}_p(\mathbf{x}_t), y_t) \right) \\ & \leq \frac{\|\boldsymbol{\theta}_{p,1} - \boldsymbol{\theta}\|^2 - \|\boldsymbol{\theta}_{p,T+1} - \boldsymbol{\theta}\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\nabla \mathcal{L}(\boldsymbol{\theta}_{p,t}^\top \mathbf{z}_p(\mathbf{x}_t), y_t)\|^2 \\ & \stackrel{(a)}{\leq} \frac{\|\boldsymbol{\theta}\|^2}{2\eta} + \frac{\eta L^2 T}{2} \end{aligned} \quad (36)$$

where (a) uses the Lipschitz constant in Assumption 2, the non-negativity of $\|\boldsymbol{\theta}_{p,T+1} - \boldsymbol{\theta}\|^2$, and the initial value $\boldsymbol{\theta}_{p,1} = \mathbf{0}$. The proof is complete by choosing $\boldsymbol{\theta} = \boldsymbol{\theta}_p^* = \sum_{t=1}^T \alpha_{p,t}^* \mathbf{z}_p(\mathbf{x}_t)$ such that $f_p^*(\mathbf{x}_t) = \boldsymbol{\theta}^\top \mathbf{z}_p(\mathbf{x}_t)$ in (36). ■

Lemma 5 demonstrates that the static regret of the Raker approach is upper bounded by some constants, which mainly depend on the stepsize in (16) and the time horizon T .

In addition, we will bound the difference between the loss of the solution obtained from Algorithm 1 and the loss of the best single kernel-based online learning algorithm. Specifically the following lemma holds:

Lemma 6 Without loss of generality, assume the cost $\ell_{p,t}(\cdot) \in [-1, 1]$ in Assumption 2. Let $\{f_{p,t}\}$ be the RF-based function estimators obtained from Algorithm 1, for any kernel p we have

$$\sum_{t=1}^T \sum_{p=1}^P \bar{w}_{p,t} \ell_t(f_{p,t}(\mathbf{x}_t)) - \sum_{t=1}^T \ell_t(f_{p,t}(\mathbf{x}_t)) \leq \eta T + \frac{\ln P}{\eta} \quad (37)$$

where η is the learning rate in (14), and P is the number of kernels in the dictionary.

Proof: Defining $W_t = \sum_{p=1}^P w_{p,t}$, the weight recursion in (14) implies that

$$\begin{aligned}
 W_{t+1} &= \sum_{p=1}^P w_{p,t+1} \\
 &= \sum_{p=1}^P w_{p,t} \exp(-\eta \ell_t(f_{p,t}(\mathbf{x}_t))) \\
 &\leq \sum_{p=1}^P w_{p,t} \left(1 - \eta \ell_t(f_{p,t}(\mathbf{x}_t)) + \eta^2 \ell_t(f_{p,t}(\mathbf{x}_t))^2\right)
 \end{aligned} \tag{38}$$

where the last inequality follows $\exp(-\eta x) \leq 1 - \eta x + \eta^2 x^2$, for $|\eta| \leq 1$. Furthermore, substituting the definition $\bar{w}_{p,t} := w_{p,t} / \sum_{p=1}^P w_{p,t} = w_{p,t} / W_t$ into (38), it follows that

$$\begin{aligned}
 W_{t+1} &\leq \sum_{p=1}^P W_t \bar{w}_{p,t} \left(1 - \eta \ell_t(f_{p,t}(\mathbf{x}_t)) + \eta^2 \ell_t(f_{p,t}(\mathbf{x}_t))^2\right) \\
 &= W_t \left(1 - \eta \sum_{p=1}^P \bar{w}_{p,t} \ell_t(f_{p,t}(\mathbf{x}_t)) + \eta^2 \sum_{p=1}^P \bar{w}_{p,t} \ell_t(f_{p,t}(\mathbf{x}_t))^2\right) \\
 &\stackrel{(a)}{\leq} W_t \exp\left(-\eta \sum_{p=1}^P \bar{w}_{p,t} \ell_t(f_{p,t}(\mathbf{x}_t)) + \eta^2 \sum_{p=1}^P \bar{w}_{p,t} \ell_t(f_{p,t}(\mathbf{x}_t))^2\right)
 \end{aligned} \tag{39}$$

where (a) follows from $1 + x \leq e^x$, $\forall x$. Telescoping (39) from $t = 1$ to T , we have ($W_1 = 1$)

$$W_{T+1} \leq \exp\left(-\eta \sum_{t=1}^T \sum_{p=1}^P \bar{w}_{p,t} \ell_t(f_{p,t}(\mathbf{x}_t)) + \eta^2 \sum_{t=1}^T \sum_{p=1}^P \bar{w}_{p,t} \ell_t(f_{p,t}(\mathbf{x}_t))^2\right). \tag{40}$$

On the other hand, for any p , the following holds true

$$\begin{aligned}
 W_{T+1} \geq w_{p,T+1} &= w_{p,1} \prod_{t=1}^T \exp(-\eta \ell_t(f_{p,t}(\mathbf{x}_t))) \\
 &= w_{p,1} \exp\left(-\eta \sum_{t=1}^T \ell_t(f_{p,t}(\mathbf{x}_t))\right).
 \end{aligned} \tag{41}$$

Combining (40) with (41), we arrive at

$$\exp\left(-\eta \sum_{t=1}^T \sum_{p=1}^P \bar{w}_{p,t} \ell_t(f_{p,t}(\mathbf{x}_t)) + \eta^2 \sum_{t=1}^T \sum_{p=1}^P \bar{w}_{p,t} \ell_t(f_{p,t}(\mathbf{x}_t))^2\right) \geq w_{p,1} \exp\left(-\eta \sum_{t=1}^T \ell_t(f_{p,t}(\mathbf{x}_t))\right). \tag{42}$$

Taking the logarithm on both sides of (42), it follows that (cf. $w_{p,1} = 1/P$)

$$-\eta \sum_{t=1}^T \sum_{p=1}^P \bar{w}_{p,t} \ell_t(f_{p,t}(\mathbf{x}_t)) + \eta^2 \sum_{t=1}^T \sum_{p=1}^P \bar{w}_{p,t} \ell_t(f_{p,t}(\mathbf{x}_t))^2 \geq -\eta \sum_{t=1}^T \ell_t(f_{p,t}(\mathbf{x}_t)) - \ln P \tag{43}$$

which leads to

$$\sum_{t=1}^T \sum_{p=1}^P \bar{w}_{p,t} \ell_t(f_{p,t}(\mathbf{x}_t)) \leq \sum_{t=1}^T \ell_t(f_{p,t}(\mathbf{x}_t)) + \eta \sum_{t=1}^T \sum_{p=1}^P \bar{w}_{p,t} \ell_t(f_{p,t}(\mathbf{x}_t))^2 + \frac{\ln P}{\eta} \tag{44}$$

from which the proof is complete since $\ell_t(f_{p,t}(\mathbf{x}_t))^2 \leq 1$ and $\sum_{p=1}^P \bar{w}_{p,t} = 1$. ■

Moreover, since the loss function $\ell_t(\cdot)$ is convex, the following inequality holds

$$\ell_t \left(\sum_{p=1}^P \bar{w}_{p,t} f_{p,t}(\mathbf{x}_t) \right) \leq \sum_{p=1}^P \bar{w}_{p,t} \ell_t(f_{p,t}(\mathbf{x}_t)). \quad (45)$$

Plugging (45) into (37) in Lemma 6, we arrive at

$$\begin{aligned} \sum_{t=1}^T \ell_t \left(\sum_{p=1}^P \bar{w}_{p,t} f_{p,t}(\mathbf{x}_t) \right) &\leq \sum_{t=1}^T \ell_t(f_{p,t}(\mathbf{x}_t)) + \eta T + \frac{\ln P}{\eta} \\ &\stackrel{(b)}{\leq} \sum_{t=1}^T \ell_t(f_p^*(\mathbf{x}_t)) + \frac{\ln P}{\eta} + \frac{\|\boldsymbol{\theta}_p^*\|^2}{2\eta} + \frac{\eta L^2 T}{2} + \eta T \end{aligned} \quad (46)$$

where (b) follows from $\boldsymbol{\theta}_p^*$ is the optimal solution for any given kernel p , which leads to Lemma 1.

B Proof of Theorem 2

To derive the performance bound relative to the best function estimator $f^*(\mathbf{x}_t)$ in the RKHS, the key step is to bound the error of random feature approximation. For a given shift-invariant kernel κ_p , with probability at least $1 - 2^8 \left(\frac{\sigma_p}{\epsilon}\right)^2 \exp\left(\frac{-D\epsilon^2}{4d+8}\right)$, the point-wise error of $2D$ -dimension random feature approximation is uniformly bounded by [15]

$$\sup_{\mathbf{x}_i, \mathbf{x}_j \in \mathcal{X}} \left| \mathbf{z}_p(\mathbf{x}_i)^\top \mathbf{z}_p(\mathbf{x}_j) - \kappa_p(\mathbf{x}_i, \mathbf{x}_j) \right| < \epsilon \quad (47)$$

where $\epsilon > 0$ is a given constant, D denotes the number of random Fourier features, while d represents the dimension of original datum \mathbf{x} , and $\sigma_p^2 := \mathbb{E}_p[\mathbf{v}^\top \mathbf{v}]$ is the second order moments of the random Fourier features. Henceforth, for the optimal function estimator (18) in \mathcal{H}_p denoted by $f_{\mathcal{H}_p}^*(\mathbf{x}) = \sum_{t=1}^T \alpha_t^* \kappa_p(\mathbf{x}, \mathbf{x}_t)$, defining the corresponding function $\tilde{f}_p^* := \sum_{t=1}^T \alpha_t^* \mathbf{z}_p^\top(\mathbf{x}) \mathbf{z}_p(\mathbf{x}_t) \in \mathcal{F}_p$, we have

$$\begin{aligned} \left| \sum_{t=1}^T \ell_t(\tilde{f}_p^*(\mathbf{x}_t)) - \sum_{t=1}^T \ell_t(f_{\mathcal{H}_p}^*(\mathbf{x}_t)) \right| &\stackrel{(a)}{\leq} \sum_{t=1}^T \left| \ell_t(\tilde{f}_p^*(\mathbf{x}_t)) - \ell_t(f_{\mathcal{H}_p}^*(\mathbf{x}_t)) \right| \\ &\stackrel{(b)}{\leq} \sum_{t=1}^T L \left| \sum_{i=1}^T \alpha_i^* \mathbf{z}_p^\top(\mathbf{x}_i) \mathbf{z}_p(\mathbf{x}_t) - \sum_{i=1}^T \alpha_i^* \kappa_p(\mathbf{x}_i, \mathbf{x}_t) \right| \\ &\stackrel{(c)}{\leq} \sum_{t=1}^T L \sum_{i=1}^T |\alpha_i^*| \left| \mathbf{z}_p^\top(\mathbf{x}_i) \mathbf{z}_p(\mathbf{x}_t) - \kappa_p(\mathbf{x}_i, \mathbf{x}_t) \right| \end{aligned} \quad (48)$$

where (a) follows from the triangle inequality, (b) uses the Lipschitz continuity of the loss function, and (c) is due to the Cauchy-Schwarz inequality. Combining with (47) obtains

$$\left| \sum_{t=1}^T \ell_t(\tilde{f}_p^*(\mathbf{x}_t)) - \sum_{t=1}^T \ell_t(f_{\mathcal{H}_p}^*(\mathbf{x}_t)) \right| \leq \sum_{t=1}^T L \epsilon \sum_{i=1}^T |\alpha_i^*| = \epsilon L T \|f_{\mathcal{H}_p}^*\|_1. \quad (49)$$

Furthermore, since we assume $\kappa_p(\mathbf{x}_i, \mathbf{x}_j) \leq 1, \forall i, j$, the uniform convergence in (47) also implies $\sup_{\mathbf{x}_i, \mathbf{x}_j \in \mathcal{X}} \mathbf{z}_p^\top(\mathbf{x}_i) \mathbf{z}_p(\mathbf{x}_j) \leq 1 + \epsilon$, w.h.p., and it in turn leads to

$$\|\boldsymbol{\theta}_p^*\|^2 = \left\| \sum_{t=1}^T \alpha_t^* \mathbf{z}_p(\mathbf{x}_t) \right\|^2 \leq (1 + \epsilon) \|f_{\mathcal{H}_p}^*\|_1^2 \quad (50)$$

where $\|f_{\mathcal{H}_p}^*\|_1 := \sum_{t=1}^T |\alpha_t^*|$.

Therefore, Lemma 1 together with (49) and (50) leads to the regret of the proposed Raker algorithm compared with the best static *function* in RKHS, i.e.,

$$\begin{aligned}
 & \sum_{t=1}^T \ell_t \left(\sum_{p=1}^P w_{p,t} f_{p,t}(\mathbf{x}_t) \right) - \sum_{t=1}^T \ell_t(f_{\mathcal{H}_p}^*(\mathbf{x}_t)) \\
 &= \sum_{t=1}^T \ell_t \left(\sum_{p=1}^P w_{p,t} f_{p,t}(\mathbf{x}_t) \right) - \sum_{t=1}^T \ell_t(\tilde{f}_p^*(\mathbf{x}_t)) + \sum_{t=1}^T \ell_t(\tilde{f}_p^*(\mathbf{x}_t)) - \sum_{t=1}^T \ell_t(f_{\mathcal{H}_p}^*(\mathbf{x}_t)) \\
 &\leq \frac{\ln P}{\eta} + \frac{\eta L^2 T}{2} + \eta T + \frac{(1+\epsilon)\|f_{\mathcal{H}_p}^*\|_1^2}{2\eta} + \epsilon L T \|f_{\mathcal{H}_p}^*\|_1
 \end{aligned} \tag{51}$$

which completes the proof of Theorem 2.

C Proof of Lemma 3

Following from Lemma 1, with $\eta = \frac{1}{\sqrt{T}}$, we have

$$\sum_{t=1}^T \ell_t \left(\sum_{p=1}^P w_{p,t} f_{p,t}(\mathbf{x}_t) \right) - \sum_{t=1}^T \ell_t(f_{p^*}^*(\mathbf{x}_t)) \leq \left(\ln P + \frac{\|\boldsymbol{\theta}_{p^*}^*\|^2}{2} + \frac{L^2}{2} + 1 \right) \sqrt{T} := c_0 \sqrt{T} \tag{52}$$

where the index is defined as $p^* = \arg \min_{p \in \mathcal{P}} \sum_{t=1}^T \ell_t(f_p^*(\mathbf{x}_t))$. Therefore, for at the end of each instance I , we can conclude that the static regret of the Raker learner \mathcal{A}_I is (cf. (28))

$$\text{Reg}_{\mathcal{A}_I}^s(|I|) = \sum_{t \in I} \ell_t(f_t^{(I)}(\mathbf{x}_t)) - \sum_{t \in I} \ell_t(f_{p^*}^*(\mathbf{x}_t)) \leq c_0 \sqrt{|I|} \tag{53}$$

where $f_t^{(I)}(\mathbf{x}_t)$ is defined in (25). To this end, we can sketch the main steps as follows.

For every interval I , the static regret of the ada-Raker learner \mathcal{A} can be decomposed as

$$\text{Reg}_{\mathcal{A}}^s(|I|) = \underbrace{\sum_{t \in I} \ell_t(f_t(\mathbf{x}_t)) - \sum_{t \in I} \ell_t(f_t^{(I)}(\mathbf{x}_t))}_{\mathcal{R}_1} + \underbrace{\sum_{t \in I} \ell_t(f_t^{(I)}(\mathbf{x}_t)) - \sum_{t \in I} \ell_t(f_{p^*}^*(\mathbf{x}_t))}_{\mathcal{R}_2} \tag{54}$$

where \mathcal{R}_1 is the regret of the Ada-Raker learner \mathcal{A} relative to the Raker learner \mathcal{A}_I , and \mathcal{R}_2 is the static regret of \mathcal{A}_I on this interval. Notice that \mathcal{R}_2 can directly follow from (53), while \mathcal{R}_1 can be bounded following the same steps as multiple kernel combinations in Lemma 6. Different from the kernel selections however, the crux here is that the number of Raker learners (experts) here is time-varying, i.e., $|\mathcal{I}(t)|$.

A tight bound can be resolved via the *Sleeping Experts* reformulation [31], meaning that the expert that has never appeared before should be thought of as being *asleep* for all previous rounds. Nevertheless, to get a looser estimate, we assume all the experts (instances $\{\mathcal{A}_I\}$) ever appeared until t are all active; that is, the total number of experts is upper bounded by $t \log t$, since at most $\log t$ experts are run during time t . Following (38)-(44), we have that

$$\mathcal{R}_1 \leq \eta^{(I)} |I| + \frac{\ln(t \log t)}{\eta^{(I)}} = \sqrt{|I|} (1 + \ln t + \ln(\log t)) \leq \sqrt{|I|} (1 + 2 \ln t) \tag{55}$$

where $\eta^{(I)}$ is chosen as $\eta^{(I)} = 1/\sqrt{|I|}$, and $\ln(\log t) \leq \ln(t)$. Together with (53), it follows that for any interval I , we have

$$\text{Reg}_{\mathcal{A}}^s(|I|) = \sqrt{|I|} (1 + c_0 + 2 \ln t) \leq \sqrt{|I|} (1 + c_0 + 2 \ln T). \tag{56}$$

Note that this bound only holds for those intervals (collected in \mathcal{I}) (re)initializing Raker instance \mathcal{A}_I , since the static regret bound (54) holds only at the end of such interval.

The next step is to show that for any interval $I \subseteq \mathcal{T}$, the above sub-linear bound on $\text{Reg}_{\mathcal{A}}^s(|I|)$ holds. This extension can be done whenever the interval set \mathcal{I} is properly designed. Specifically, the example of interval partition given in Section 4.1 does have such desired properties, and the detailed arguments can follow [31, A.2].

D Proof of Theorem 4

To start with, the dynamic regret in (26) can be decomposed by

$$\text{Reg}_{\mathcal{A}}^{\text{d}}(T) := \sum_{t=1}^T \ell_t(f_t(\mathbf{x}_t)) - \sum_{t=1}^T \ell_t(f^*(\mathbf{x}_t)) + \sum_{t=1}^T \ell_t(f^*(\mathbf{x}_t)) - \sum_{t=1}^T \ell_t(f_t^*(\mathbf{x}_t)) \quad (57)$$

where $f^*(\cdot)$ is the best fixed function estimate in (18), and $f_t^*(\cdot)$ is the best dynamic function estimate in (27), both of which belong to the space $\mathcal{F} := \bigcup_{p \in \mathcal{P}} \mathcal{H}_p$. In (57), the first difference term is the static regret of the AdaRaker algorithm, and the second difference term is the relative loss between the best fixed function estimate and the best dynamic solution in the common function space.

Intuitively, if the time horizon T is quite large, then the average static regret will become small, but the gap between two benchmarks is large. With the insights gained from [12, 36], the length of time horizon essentially trades off the values of two terms. Thus, splitting \mathcal{T} into sub-horizons $\{\mathcal{T}_s\}$, $s = 1, \dots, \lfloor T/\Delta T \rfloor$ with each length ΔT , the dynamic regret of AdaRaker can be bounded by

$$\text{Reg}_{\mathcal{A}}^{\text{d}}(T) = \underbrace{\sum_{s=1}^{\lfloor T/\Delta T \rfloor} \sum_{t \in \mathcal{T}_s} (\ell_t(f_t(\mathbf{x}_t)) - \ell_t(f^*(\mathbf{x}_t)))}_{\mathcal{R}_1} + \underbrace{\sum_{s=1}^{\lfloor T/\Delta T \rfloor} \sum_{t \in \mathcal{T}_s} (\ell_t(f^*(\mathbf{x}_t)) - \ell_t(f_t^*(\mathbf{x}_t)))}_{\mathcal{R}_2} \quad (58)$$

where \mathcal{R}_1 can be bounded under AdaRaker from Lemma 3, and \mathcal{R}_2 that depends on the variability of the environments $\mathbb{V}(\{\ell_t\})$, can be bounded by [12, Prop. 2]

$$\mathcal{R}_2 \leq 2\Delta T \mathbb{V}(\{\ell_t\}_{t \in \mathcal{T}_s}). \quad (59)$$

Together with Lemma 3, it follows that

$$\begin{aligned} \text{Reg}_{\mathcal{A}}^{\text{d}}(T) &\leq \sum_{s=1}^{\lfloor T/\Delta T \rfloor} \left((C_0 + C_1 \ln T) \sqrt{\Delta T} + 2\Delta T \mathbb{V}(\{\ell_t\}_{t \in \mathcal{T}_s}) \right) \\ &= (C_0 + C_1 \ln T) \frac{T}{\sqrt{|\Delta T|}} + 2|\Delta T| \mathbb{V}(\{\ell_t\}_{t=1}^T). \end{aligned} \quad (60)$$

Since (29) in Lemma 3 holds for any interval $\Delta T \subseteq \mathcal{T}$, then selecting ΔT so that $|\Delta T| = (T/\mathbb{V}(\{\ell_t\}_{t=1}^T))^{\frac{2}{3}}$, it follows that

$$\text{Reg}_{\mathcal{A}}^{\text{d}}(T) \leq (C_0 + C_1 \ln T) T^{\frac{2}{3}} \mathbb{V}^{\frac{1}{3}}(\{\ell_t\}_{t=1}^T) + 2T^{\frac{2}{3}} \mathbb{V}^{\frac{1}{3}}(\{\ell_t\}_{t=1}^T). \quad (61)$$

The additional approximation error relative to the function in $\mathcal{F} := \bigcup_{p \in \mathcal{P}} \mathcal{H}_p$ can be derived following (48)-(49), from which the proof is complete.