# Supplementary materials

## ECHO: Adaptive Correction for Subgraph-wise Sampling with Lightweight Hyperparameter Search

## A   PROOF OF THEOREM 3.3

By the smoothness of $\mathcal{L}(\theta_t)$ , we have

$$\mathbb{E}\left[\mathcal{L}(\theta_{t+1})\right] \leq \mathbb{E}\left[\mathcal{L}(\theta_t)\right] + \mathbb{E}\left[\langle \nabla \mathcal{L}(\theta_t), \theta_{t+1} - \theta_t \rangle\right] \\ + \frac{L_f}{2}\mathbb{E}\left[\|\theta_{t+1} - \theta_t\|^2\right] \tag{1}$$

according to the update rule of subgraph-wise sampling training

$$\theta_{t+1} = \theta_t - \eta \nabla \mathcal{L}_{\mathcal{P}}(\theta_t)$$

by taking the norm on the both side, we have

$$\mathbb{E}\left[\|\theta_{t+1} - \theta_t\|^2\right] = \eta^2 \mathbb{E}\left[\|\nabla \mathcal{L}_{\mathcal{P}}(\theta_t)\|^2\right] \tag{2}$$

We can give the upper bound for the second term on the right of Eq.1

$$\begin{aligned}
&\mathbb{E}\left[\langle \nabla \mathcal{L}(\theta_t), \theta_{t+1} - \theta_t \rangle\right] \\
&= -\eta \mathbb{E}\left[\langle \nabla \mathcal{L}(\theta_t), \nabla \mathcal{L}_{\mathcal{P}}(\theta_t) \rangle\right] \\
&= -\eta \langle \nabla \mathcal{L}(\theta_t), \mathbb{E}\left[\nabla \mathcal{L}_{\mathcal{P}}(\theta_t)\right]\rangle \\
&\overset{(a)}{=} -\frac{\eta}{2}\left(\left\|\nabla \mathcal{L}(\theta^t)\right\|^2 + \|\mathbb{E}\left[\nabla \mathcal{L}_{\mathcal{P}}(\theta_t)\right]\|^2\right) \\
&\quad + \frac{\eta}{2}\left(\|\nabla \mathcal{L}(\theta_t) - \mathbb{E}\left[\nabla \mathcal{L}_{\mathcal{P}}(\theta_t)\right]\|^2\right) \\
&\overset{(b)}{\leq} -\frac{\eta}{2}\left(\left\|\nabla \mathcal{L}(\theta^t)\right\|^2 + \|\mathbb{E}\left[\nabla \mathcal{L}_{\mathcal{P}}(\theta_t)\right]\|^2 - 2\kappa_1^2 - 4\kappa_2^2\right)
\end{aligned} \tag{3}$$

where (a) is due to $2\langle x, y \rangle = \|x\|^2 + \|y\|^2 - \|x - y\|^2$
where (b)

$$\begin{aligned}
&\|\nabla \mathcal{L}(\theta_t) - \mathbb{E}\left[\nabla \mathcal{L}_{\mathcal{P}}(\theta_t)\right]\|^2 \\
&= \mathbb{E}\left[\left\|\nabla \mathcal{L}(\theta_t) - \nabla \mathcal{L}_{\mathcal{P}}^{\text{full}}(\theta_t) + \nabla \mathcal{L}_{\mathcal{P}}^{\text{full}}(\theta_t) - \nabla \mathcal{L}_{\mathcal{P}}(\theta_t)\right\|^2\right] \\
&\leq 2\mathbb{E}\left[\left\|\nabla \mathcal{L}_{\mathcal{P}}^{\text{full}}(\theta_t) - \nabla \mathcal{L}_{\mathcal{P}}(\theta_t)\right\|^2\right] \\
&\quad + 4\mathbb{E}\left[\left\|\nabla \mathcal{L}_{\mathcal{P}}^{\text{full}}(\theta_t) - \nabla \mathcal{L}_{\mathcal{B}}(\theta_t)\right\|^2\right] \\
&\quad + 4\|\nabla \mathcal{L}(\theta_t) - \mathbb{E}\left[\nabla \mathcal{L}_{\mathcal{B}}(\theta_t)\right]\|^2 \\
&\overset{(c)}{\leq} 2\kappa_1^2 + 4\kappa_2^2
\end{aligned} \tag{4}$$

where (c) is due to $\nabla \mathcal{L}(\theta_t) = \mathbb{E}\left[\nabla \mathcal{L}_{\mathcal{B}}(\theta_t)\right]$ and follows the definition of $\kappa_1^2, \kappa_2^2$.

Combining Eq.1,2,3 gives us

$$\begin{aligned}
&\mathbb{E}\left[\mathcal{L}(\theta_{t+1})\right] \\
&\leq \mathbb{E}\left[\mathcal{L}(\theta_t)\right] + \frac{L_f}{2}\eta^2 \mathbb{E}\left[\|\nabla \mathcal{L}_{\mathcal{P}}(\theta_t)\|^2\right] \\
&\quad - \frac{\eta}{2}\left(\|\nabla \mathcal{L}(\theta_t)\|^2 + \|\mathbb{E}\left[\nabla \mathcal{L}_{\mathcal{P}}(\theta_t)\right]\|^2 - 2\kappa_1^2 - 4\kappa_2^2\right)
\end{aligned} \tag{5}$$

Organizing the above formula gives us

$$\begin{aligned}
\|\nabla \mathcal{L}(\theta_t)\|^2 &\leq \frac{2}{\eta}\left(\mathbb{E}\left[\mathcal{L}(\theta_t)\right] - \mathbb{E}\left[\mathcal{L}(\theta_{t+1})\right]\right) \\
&\quad - (1 - L_f \eta)\mathbb{E}\left[\|\nabla \mathcal{L}_{\mathcal{P}}(\theta_t)\|^2\right] \\
&\quad + 2\kappa_1^2 + 4\kappa_2^2
\end{aligned} \tag{6}$$

Summing over $t \in \{0, \cdots, T-1\}$ and dividing both side by $T$, we get,

$$\begin{aligned}
&\frac{1}{T}\sum_{t=0}^{T-1} \mathbb{E}\|[\nabla \mathcal{L}(\theta_t)]\|^2 \\
&\leq \frac{2}{T\eta}\sum_{t=0}^{T-1}(\mathcal{L}(\theta_t) - \mathcal{L}(\theta_{(t+1)})) - \frac{1}{T}\sum_{t=0}^{T-1}(1 - L_f \eta)\mathbb{E}\left[\|\nabla \mathcal{L}_{\mathcal{P}}(\theta_t)\|^2\right] \\
&\quad + 2\kappa_1^2 + 4\kappa_2^2
\end{aligned} \tag{7}$$

Since $\max_t (\nabla \mathcal{L}(\theta_t) - \nabla \mathcal{L}(\theta_{(t+1)})) \leq \nabla \mathcal{L}(\theta_0) - \nabla \mathcal{L}(\theta^*)$

$$\begin{aligned}
&\frac{1}{T}\sum_{t=0}^{T-1} \mathbb{E}\|[\nabla \mathcal{L}(\theta_t)]\|^2 \\
&\leq \frac{2}{\eta}((\nabla \mathcal{L}\theta_0) - \nabla \mathcal{L}(\theta^*)) - \frac{1}{T}\sum_{t=0}^{T-1}(1 - L_f \eta)\mathbb{E}\left[\|\nabla \mathcal{L}_{\mathcal{P}}(\theta_t)\|^2\right] \\
&\quad + 2\kappa_1^2 + 4\kappa_2^2
\end{aligned} \tag{8}$$

if we choose $\eta = \frac{1}{\sqrt{T}}$, where $0 < \eta < \frac{1}{L_f}$, then we have,

$$\frac{1}{T}\sum_{t=0}^{T-1} \mathbb{E}\|[\nabla \mathcal{L}(\theta_t)]\|^2 = O\left(\frac{1}{\sqrt{T}}\right) + O\left(\kappa_1^2 + 2\kappa_2^2\right) \tag{9}$$

## B   PROOF OF THEOREM 4.2

For convenience, we use $\mathcal{T} = \{1, \cdots, T\}$ to denote the total training epochs. Let $\mathcal{T}_s \subseteq \mathcal{T}$ to denote subgraph-wise sampling epochs and $\mathcal{T}_c \subseteq \mathcal{T}$ to denote correction epochs.

By the smoothness of $\mathcal{L}(\theta_t)$ , we have

$$\mathbb{E}\left[\mathcal{L}(\theta_{t+1})\right] \leq \mathbb{E}\left[\mathcal{L}(\theta_t)\right] + \mathbb{E}\left[\langle \nabla \mathcal{L}(\theta_t), \theta_{t+1} - \theta_t \rangle\right] \\ + \frac{L_f}{2}\mathbb{E}\left[\|\theta_{t+1} - \theta_t\|^2\right] \tag{10}$$

For $t \in \mathcal{T}_c$, the update rule is:

$$\theta_{t+1} = \theta_t - \gamma \tilde{\nabla} \mathcal{L}_{\mathcal{B}}(\theta_t) \tag{11}$$

By taking the norm on the both side, we have

$$\mathbb{E}\left[\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|^2\right] = \gamma^2 \mathbb{E}\left[\left\|\tilde{\nabla}\mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_t)\right\|^2\right] \quad (12)$$

We can give the upper bound for the second term on the right of Eq.10

$$\mathbb{E}\left[\langle \nabla\mathcal{L}(\boldsymbol{\theta}_t), \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t \rangle\right]$$

$$= -\gamma\mathbb{E}\left[\langle \nabla\mathcal{L}(\boldsymbol{\theta}_t), \tilde{\nabla}\mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_t) \rangle\right]$$

$$= -\frac{\gamma}{2}\mathbb{E}\left[\|\nabla\mathcal{L}(\boldsymbol{\theta}_t)\|^2 + \left\|\tilde{\nabla}\mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_t)\right\|^2 - \left\|\nabla\mathcal{L}(\boldsymbol{\theta}_t) - \tilde{\nabla}\mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_t)\right\|^2\right]$$

$$\overset{\leq}{_{(d)}} -\frac{\gamma}{2}\left(\mathbb{E}\left[\|\nabla\mathcal{L}(\boldsymbol{\theta}_t)\|^2\right] + \mathbb{E}\left[\left\|\tilde{\nabla}\mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_t)\right\|^2\right] - \sigma_{\text{bias}}^2\right)$$

$$(13)$$

Where (d) is due to Assumption ??.

Combining Eq.10,12,13 gives us

$$\mathbb{E}\left[\mathcal{L}(\boldsymbol{\theta}_{t+1})\right]$$

$$\leq \mathbb{E}\left[\mathcal{L}(\boldsymbol{\theta}_t)\right] - \frac{\gamma}{2}\left(\mathbb{E}\|\nabla\mathcal{L}(\boldsymbol{\theta}_t)\|^2 + \mathbb{E}\left\|\tilde{\nabla}\mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_t)\right\|^2 - \sigma_{\text{bias}}^2\right) \quad (14)$$

$$+ \frac{L_f}{2}\gamma^2\mathbb{E}\left[\left\|\tilde{\nabla}\mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_t)\right\|^2\right]$$

Reorganizing the above equation gives us

$$\mathbb{E}\left[\|\nabla\mathcal{L}(\boldsymbol{\theta}_t)\|^2\right] \leq \frac{2}{\gamma}\left(\mathbb{E}\left[\mathcal{L}(\boldsymbol{\theta}_t)\right] - \mathbb{E}\left[\mathcal{L}(\boldsymbol{\theta}_{t+1})\right]\right)$$

$$+ (L_f\gamma - 1)\mathbb{E}\left[\left\|\tilde{\nabla}\mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_t)\right\|^2\right] + \sigma_{\text{bias}}^2 \quad (15)$$

For subgraph-wise sampling epochs $t \in \mathcal{T}_m$, we have Eq.6. Summing over $t \in \{1, \cdots, T\}$ and combing Eq.6 and Eq.15, we have,

$$\sum_{t=1}^{T}\mathbb{E}\left[\|\nabla\mathcal{L}(\boldsymbol{\theta}_t)\|^2\right] = \sum_{t\in\mathcal{T}_c}\mathbb{E}\left[\|\nabla\mathcal{L}(\boldsymbol{\theta}_t)\|^2\right] + \sum_{t\in\mathcal{T}_m}\mathbb{E}\left[\|\nabla\mathcal{L}(\boldsymbol{\theta}_t)\|^2\right] \quad (16)$$

Reorganizing the above equation gives us

$$\sum_{t=1}^{T}\mathbb{E}\left[\|\nabla\mathcal{L}(\boldsymbol{\theta}_t)\|^2\right]$$

$$\leq \sum_{t\in\mathcal{T}_c}\left[\frac{2}{\gamma}\left(\mathbb{E}\left[\mathcal{L}(\boldsymbol{\theta}_t)\right] - \mathbb{E}\left[\mathcal{L}(\boldsymbol{\theta}_{t+1})\right]\right)\right]$$

$$+ \sum_{t\in\mathcal{T}_s}\left[\frac{2}{\eta}\left(\mathbb{E}\left[\mathcal{L}(\boldsymbol{\theta}_t)\right] - \mathbb{E}\left[\mathcal{L}(\boldsymbol{\theta}_{t+1})\right]\right)\right]$$

$$+ \sum_{t\in\mathcal{T}_c}(L_f\gamma - 1)\mathbb{E}\left[\left\|\tilde{\nabla}\mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_t)\right\|^2\right] \quad (17)$$

$$+ \sum_{t\in\mathcal{T}_s}(L_f\eta - 1)\mathbb{E}\left[\left\|\nabla\mathcal{L}_{\mathcal{P}}^{\text{local}}(\boldsymbol{\theta}_t)\right\|^2\right]$$

$$+ \sum_{t\in\mathcal{T}_c}\sigma_{\text{bias}}^2 + \sum_{t\in\mathcal{T}_s}\left(2\kappa_1^2 + 4\kappa_2^2\right)$$

We aim to select appropriate $\mathcal{T}_c, \mathcal{T}_s$ to satisfy the following inequality

$$\sum_{t\in\mathcal{T}_c}\sigma_{\text{bias}}^2 + \sum_{t\in\mathcal{T}_s}\left(2\kappa_1^2 + 4\kappa_2^2\right) \leq \sum_{t\in\mathcal{T}_s}(1 - L_f\eta)\mathbb{E}\left[\|\nabla\mathcal{L}_{\mathcal{P}}(\boldsymbol{\theta}_t)\|^2\right]$$

$$+ \sum_{t\in\mathcal{T}_c}(1 - L_f\gamma)\mathbb{E}\left[\left\|\tilde{\nabla}\mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_t)\right\|^2\right] \quad (18)$$

Let $G_s = \min_{t\in\mathcal{T}_s}\mathbb{E}\left[\|\nabla\mathcal{L}_{\mathcal{P}}(\boldsymbol{\theta}_t)\|^2\right], G_c = \min_{t\in\mathcal{T}_c}\mathbb{E}\left[\left\|\tilde{\nabla}\mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}_t)\right\|^2\right]$. After rearranging Eq.18, we have

$$|\mathcal{T}_c| \geq \frac{T((L_f\eta - 1)G_s + 2\kappa_1^2 + 4\kappa_2^2)}{(1 - L_f\gamma)G_c - (1 - L_f\eta)G_s + 2\kappa_1^2 + 4\kappa_2^2 - \sigma_{\text{bias}}^2} \quad (19)$$

Suppose Eq.19 holds, we have

$$\sum_{t=1}^{T}\mathbb{E}\left[\|\nabla\mathcal{L}(\boldsymbol{\theta}_t)\|^2\right] \leq \sum_{t\in\mathcal{T}_c}\left[\frac{2}{\gamma}\left(\mathbb{E}\left[\mathcal{L}(\boldsymbol{\theta}_t)\right] - \mathbb{E}\left[\mathcal{L}(\boldsymbol{\theta}_{t+1})\right]\right)\right]$$

$$+ \sum_{t\in\mathcal{T}_s}\left[\frac{2}{\eta}\left(\mathbb{E}\left[\mathcal{L}(\boldsymbol{\theta}_t)\right] - \mathbb{E}\left[\mathcal{L}(\boldsymbol{\theta}_{t+1})\right]\right)\right] \quad (20)$$

Let $\eta = \gamma = \frac{1}{\sqrt{T}}$, we will have

$$\sum_{t=1}^{T}\mathbb{E}\left[\|\nabla\mathcal{L}(\boldsymbol{\theta}_t)\|^2\right] \leq \sum_{t\in\mathcal{T}}\left[\frac{2}{\sqrt{T}}\left(\mathbb{E}\left[\mathcal{L}(\boldsymbol{\theta}_t)\right] - \mathbb{E}\left[\mathcal{L}(\boldsymbol{\theta}_{t+1})\right]\right)\right] \quad (21)$$

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\left[\|\nabla\mathcal{L}(\boldsymbol{\theta}_t)\|^2\right] \leq \frac{2}{\sqrt{T}}\left(\mathcal{L}(\boldsymbol{\theta}_0) - \mathcal{L}(\boldsymbol{\theta}^*)\right) \quad (22)$$

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\left[\|\nabla\mathcal{L}(\boldsymbol{\theta}_t)\|^2\right] = O\left(\frac{1}{\sqrt{T}}\right) \quad (23)$$

## C ADDITIONAL EXPERIMENTAL RESULTS

### C.1 Effectiveness of AES

Figure 1 and Figure 2 respectively plot the test accuracy and validation loss on all tested datasets for both models as complements to Figure 9 and Figure 10 in the paper. Consistent with the statement in the paper, AES achieves nearly the optimal trade-off between accuracy and time efficiency in all tested cases. Also, AES shows the fastest convergence among tested correction strategies.
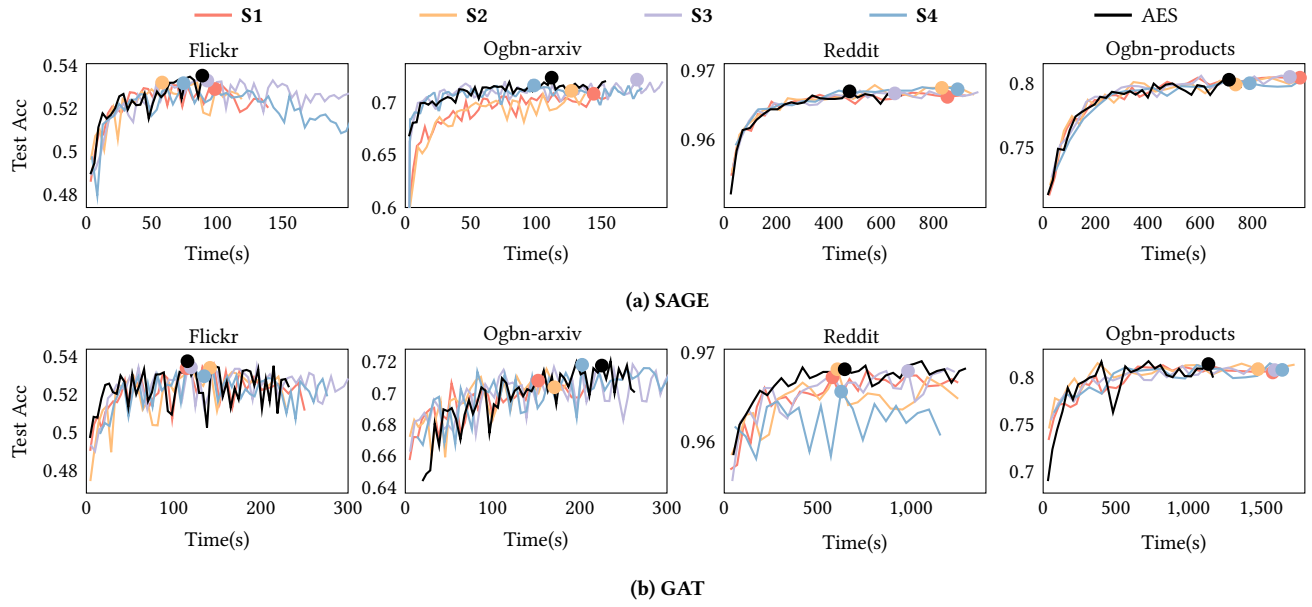
(a) SAGE



(b) GAT

Figure 1: Test accuracy over time with different correction strategies.
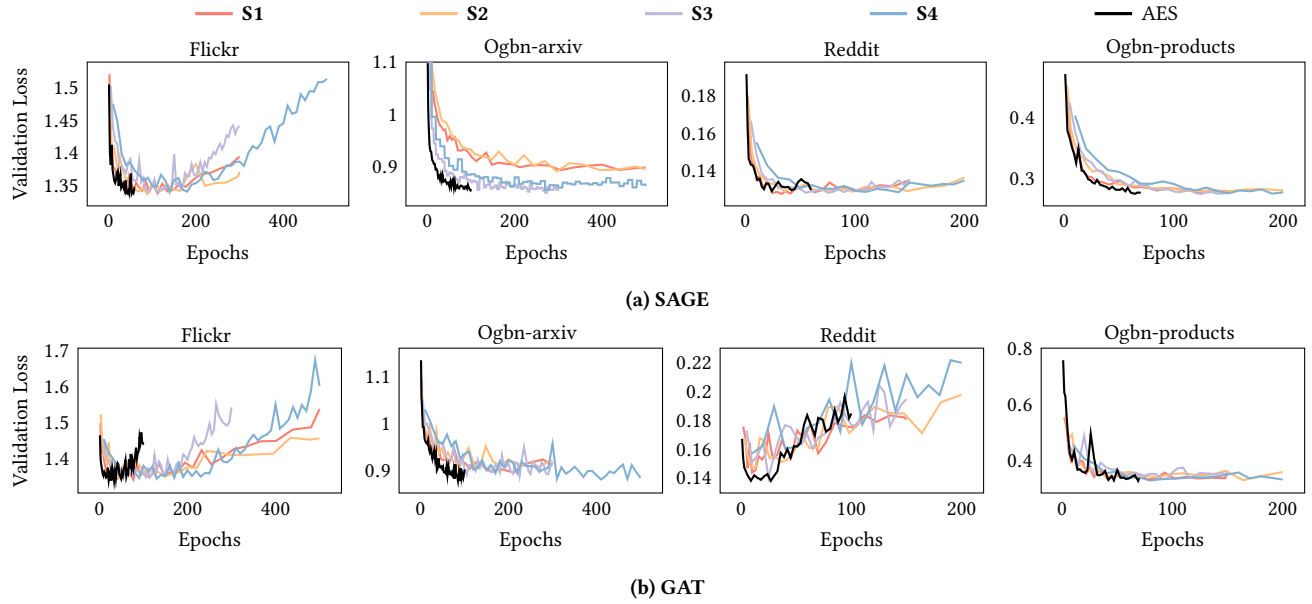


(a) SAGE



(b) GAT

Figure 2: Validation loss over training epochs with different correction strategies.