

294 **A Algorithm**

295 In Algorithm 2, we also apply Ouroboros to Adam, an adaptive variant of SGD. We update first
 296 moment vectors and second moment vectors using gradients computed by the Ouroboros algorithm
 297 so that the updates in modules can be parallelized.

Algorithm 2 Ouroboros + Adam

Require:

Initial weights $w^0 = [w_{\mathcal{G}(1)}^0, \dots, w_{\mathcal{G}(K)}^0]$;
 Initial word embedding $V_i^0 = V_o^0$;
 Stepsize: γ ; Small constant $\epsilon = 10^{-8}$;
 Exponential decay: $\beta_1 = 0.9, \beta_2 = 0.999$;
 1_{st} moment vector: $m_{\mathcal{G}(k)}^0 = 0, \forall k, m_V^0 = 0$;
 2_{nd} moment vector: $v_{\mathcal{G}(k)}^0 = 0, \forall k, v_V^0 = 0$;

- 1: **for** $t = 0, 1, 2, \dots, T - 1$ **do**
- 2: **for** $k = 1, \dots, K$ **in parallel do**
- 3: Compute delayed gradient g_k^t for module k following (8):
- 4: Compute mixed gradient g_V^t for embedding layer following (9):
- 5: Update biased first moment estimate:

$$m_{\mathcal{G}(k)}^{t+1} = \beta_1 \cdot m_{\mathcal{G}(k)}^t + (1 - \beta_1) \cdot g_k^t;$$

$$m_V^{t+1} = \beta_1 \cdot m_V^t + (1 - \beta_1) \cdot g_V^t;$$
- 6: Update biased second moment estimate:

$$v_{\mathcal{G}(k)}^{t+1} = \beta_2 \cdot v_{\mathcal{G}(k)}^t + (1 - \beta_2) \cdot (g_k^t)^2;$$

$$v_V^{t+1} = \beta_2 \cdot v_V^t + (1 - \beta_2) \cdot (g_V^t)^2;$$
- 7: Compute bias-correct first moment estimate:

$$\hat{m}_{\mathcal{G}(k)}^{t+1} = m_{\mathcal{G}(k)}^{t+1} / (1 - \beta_1^{t+1});$$

$$\hat{m}_V^{t+1} = m_V^{t+1} / (1 - \beta_1^{t+1});$$
- 8: Compute bias-correct second moment estimate:

$$\hat{v}_{\mathcal{G}(k)}^{t+1} = v_{\mathcal{G}(k)}^{t+1} / (1 - \beta_2^{t+1});$$

$$\hat{v}_V^{t+1} = v_V^{t+1} / (1 - \beta_2^{t+1});$$
- 9: Update weights and embedding layer following Adam:

$$w_{\mathcal{G}(k)}^{t+1} = w_{\mathcal{G}(k)}^t - \gamma \cdot \frac{\hat{m}_{\mathcal{G}(k)}^{t+1}}{\left(\sqrt{\hat{v}_{\mathcal{G}(k)}^{t+1}} + \epsilon\right)};$$

$$V_i^{t+1} = V_o^{t+1} = V_i^t - \gamma \cdot \frac{\hat{m}_V^{t+1}}{\left(\sqrt{\hat{v}_V^{t+1}} + \epsilon\right)};$$
- 10: **end for**
- 11: **end for**
- 12: Output w^s, V_i^s and V_o^s randomly from $\{w^t\}_{t=0}^{T-1}, \{V_i^t\}_{t=0}^{T-1}$ and $\{V_o^t\}_{t=0}^{T-1}$.

298 **B Proof**

299 **Proof to Lemma 1**

300 *Proof:* Let $\tilde{w} = [V, w]$, it is satisfied that:

$$\nabla f(\tilde{w}^t) = \sum_{k=1}^K \nabla f_{\mathcal{G}(k)}(\tilde{w}^t) + \nabla f_{V_i}(\tilde{w}^t) + \nabla f_{V_o}(\tilde{w}^t) \quad (18)$$

301 According to Assumption 1, the following inequality holds that:

$$f(\tilde{w}^{t+1}) \leq f(\tilde{w}^t) + \nabla f(\tilde{w}^t)^T (\tilde{w}^{t+1} - \tilde{w}^t) + \frac{L}{2} \|\tilde{w}^{t+1} - \tilde{w}^t\|_2^2. \quad (19)$$

302 From the update rule in Algorithm I, we take expectation on both sides and obtain: small

$$\begin{aligned}
& \mathbb{E}[f(\tilde{w}^{t+1})] \\
\leq & f(\tilde{w}^t) - \gamma_t \nabla f(\tilde{w}^t)^T \left(\sum_{k=1}^K \nabla f_{\mathcal{G}(k)}(\tilde{w}^{t-K+k}) + \frac{1}{2} \nabla f_V(\tilde{w}^{t-K+1}) + \frac{1}{2} \nabla f_V(\tilde{w}^t) \right. \\
& \left. + \nabla f(\tilde{w}^t) - \nabla f(\tilde{w}^t) \right) + \frac{L\gamma_t^2}{2} \mathbb{E} \left\| \sum_{k=1}^K \nabla f_{\mathcal{G}(k), x_{i(t-K+k)}}(\tilde{w}^{t-K+k}) \right. \\
& \left. + \frac{1}{2} \nabla f_{V, x_{i(t-K+1)}}(\tilde{w}^{t-K+1}) + \frac{1}{2} \nabla f_{V, x_{i(t)}}(\tilde{w}^t) - \nabla f(\tilde{w}^t) + \nabla f(\tilde{w}^t) \right\|_2^2 \\
= & f(\tilde{w}^t) - \left(\gamma_t - \frac{L\gamma_t^2}{2} \right) \|\nabla f(\tilde{w}^t)\|_2^2 + \frac{L\gamma_t^2}{2} \mathbb{E} \left\| \sum_{k=1}^K \nabla f_{\mathcal{G}(k), x_{i(t-K+k)}}(\tilde{w}^{t-K+k}) \right. \\
& \left. + \frac{1}{2} \nabla f_{V, x_{i(t-K+1)}}(\tilde{w}^{t-K+1}) + \frac{1}{2} \nabla f_{V, x_{i(t)}}(\tilde{w}^t) - \nabla f(\tilde{w}^t) \right\|_2^2 \\
& - (\gamma_t - L\gamma_t^2) \nabla f(\tilde{w}^t)^T \left(\sum_{k=1}^K \nabla f_{\mathcal{G}(k)}(\tilde{w}^{t-K+k}) + \frac{1}{2} \nabla f_V(\tilde{w}^{t-K+1}) + \frac{1}{2} \nabla f_V(\tilde{w}^t) - \nabla f(\tilde{w}^t) \right) \\
= & f(\tilde{w}^t) - \left(\gamma_t - \frac{L\gamma_t^2}{2} \right) \|\nabla f(\tilde{w}^t)\|_2^2 + \underbrace{\frac{L\gamma_t^2}{2} \mathbb{E} \left\| \sum_{k=1}^K \nabla f_{\mathcal{G}(k), x_{i(t-K+k)}}(\tilde{w}^{t-K+k}) - \sum_{k=1}^K \nabla f_{\mathcal{G}(k)}(\tilde{w}^t) \right\|_2^2}_{Q_1} \\
& + \underbrace{\frac{L\gamma_t^2}{2} \mathbb{E} \left\| \frac{1}{2} \nabla f_{V, x_{i(t-K+1)}}(\tilde{w}^{t-K+1}) + \frac{1}{2} \nabla f_{V, x_{i(t)}}(\tilde{w}^t) - \nabla f_V(\tilde{w}^t) \right\|_2^2}_{Q_2} \\
& - (\gamma_t - L\gamma_t^2) \nabla f(\tilde{w}^t)^T \underbrace{\left(\sum_{k=1}^K \nabla f_{\mathcal{G}(k)}(\tilde{w}^{t-K+k}) + \frac{1}{2} \nabla f_V(\tilde{w}^{t-K+1}) + \frac{1}{2} \nabla f_V(\tilde{w}^t) - \nabla f(\tilde{w}^t) \right)}_{Q_3} \tag{20}
\end{aligned}$$

303 where the equalities follow from the unbiased gradient $\mathbb{E}[\nabla f_{x_i}(w)] = \nabla f(w)$ and $[\nabla f_{\mathcal{G}(k)}(w)]_j =$

304 0, $\forall j \notin \mathcal{G}(k)$. Because of $\|x + y\|_2^2 \leq 2\|x\|_2^2 + 2\|y\|_2^2$, we have the upper bound of Q_1 as follows:

$$\begin{aligned}
Q_1 &= \frac{L\gamma_t^2}{2} \mathbb{E} \left\| \sum_{k=1}^K \nabla f_{\mathcal{G}(k), x_{i(t-K+k)}}(\tilde{w}^{t-K+k}) - \sum_{k=1}^K \nabla f_{\mathcal{G}(k)}(\tilde{w}^t) - \sum_{k=1}^K \nabla f_{\mathcal{G}(k)}(\tilde{w}^{t-K+k}) + \sum_{k=1}^K \nabla f_{\mathcal{G}(k)}(\tilde{w}^{t-K+k}) \right\|_2^2 \\
&\leq L\gamma_t^2 \mathbb{E} \underbrace{\left\| \sum_{k=1}^K \nabla f_{\mathcal{G}(k), x_{i(t-K+k)}}(\tilde{w}^{t-K+k}) - \sum_{k=1}^K \nabla f_{\mathcal{G}(k)}(\tilde{w}^{t-K+k}) \right\|_2^2}_{Q_4} \\
&\quad + L\gamma_t^2 \underbrace{\left\| \sum_{k=1}^K \nabla f_{\mathcal{G}(k)}(\tilde{w}^{t-K+k}) - \sum_{k=1}^K \nabla f_{\mathcal{G}(k)}(\tilde{w}^t) \right\|_2^2}_{Q_5}. \tag{21}
\end{aligned}$$

305 Similarly, we get the upper bound of Q_2 as follows:

$$\begin{aligned}
Q_2 &\leq \frac{L\gamma_t^2}{4} \mathbb{E} \|\nabla f_{V, x_{i(t-K+1)}}(\tilde{w}^{t-K+1}) - \nabla f_V(\tilde{w}^t)\|_2^2 + \frac{L\gamma_t^2}{4} \|\nabla f_{V, x_{i(t)}}(\tilde{w}^t) - \nabla f_V(\tilde{w}^t)\|_2^2 \\
&\leq \frac{L\gamma_t^2}{2} \underbrace{\mathbb{E} \|\nabla f_{V, x_{i(t-K+1)}}(\tilde{w}^{t-K+1}) - \nabla f_V(\tilde{w}^{t-K+1})\|_2^2}_{Q_6} + \frac{L\gamma_t^2}{2} \underbrace{\mathbb{E} \|\nabla f_V(\tilde{w}^{t-K+1}) - \nabla f_V(\tilde{w}^t)\|_2^2}_{Q_7} \\
&\quad + \frac{L\gamma_t^2}{4} \|\nabla f_{V, x_{i(t)}}(\tilde{w}^t) - \nabla f_V(\tilde{w}^t)\|_2^2. \tag{22}
\end{aligned}$$

306 Because of $xy \leq \frac{1}{2}\|x\|_2^2 + \frac{1}{2}\|y\|_2^2$, we have:

$$\begin{aligned} Q_3 &\leq \frac{\gamma_t - L\gamma_t^2}{2} \|\nabla f(\tilde{w}^t)\|_2^2 + \frac{\gamma_t - L\gamma_t^2}{2} \left\| \sum_{k=1}^K \nabla f_{\mathcal{G}(k)}(\tilde{w}^{t-K+k}) - \sum_{k=1}^K \nabla f_{\mathcal{G}(k)}(\tilde{w}^t) \right\|_2^2 \\ &\quad + \frac{\gamma_t - L\gamma_t^2}{8} \|\nabla f_V(\tilde{w}^{t-K+1}) - \nabla f_V(\tilde{w}^t)\|_2^2. \end{aligned} \quad (23)$$

307 According to Assumption 2, we can bound Q_4 as follows:

$$\begin{aligned} Q_4 &= \sum_{k=1}^K \mathbb{E} \left\| \nabla f_{\mathcal{G}(k), x_{i(t-K+k)}}(\tilde{w}^{t-K+k}) - \nabla f_{\mathcal{G}(k)}(\tilde{w}^{t-K+k}) \right\|_2^2 \\ &\leq \sum_{k=1}^K \mathbb{E} \left\| \nabla f_{\mathcal{G}(k), x_{i(t-K+k)}}(\tilde{w}^{t-K+k}) \right\|_2^2 \\ &\leq KM, \end{aligned} \quad (24)$$

308 where the first equality follows from the definition of $\nabla f_{\mathcal{G}(k)}(w)$ so that $[\nabla f_{\mathcal{G}(k)}(w)]_j = 0$, $\forall j \notin \mathcal{G}(k)$ and the last inequality is from Assumption 2. Similarly, we can also bound Q_6 as follows:

$$Q_6 \leq M. \quad (25)$$

310 We can get the upper bound of Q_5 :

$$\begin{aligned} Q_5 &= \sum_{k=1}^K \left\| \nabla f_{\mathcal{G}(k)}(\tilde{w}^{t-K+k}) - \nabla f_{\mathcal{G}(k)}(\tilde{w}^t) \right\|_2^2 \\ &\leq \sum_{k=1}^K \left\| \nabla f(\tilde{w}^{t-K+k}) - \nabla f(\tilde{w}^t) \right\|_2^2 \\ &\leq L^2 \sum_{k=1}^K \left\| \sum_{j=\max\{0, t-K+k\}}^{t-1} (\tilde{w}^{j+1} - \tilde{w}^j) \right\|^2 \\ &\leq L^2 \gamma_{\max\{0, t-K+1\}}^2 K \sum_{k=1}^K \sum_{j=\max\{0, t-K+k\}}^{t-1} \left\| \sum_{k=1}^K \nabla f_{\mathcal{G}(k), x_{(j-K+k)}}(\tilde{w}^{j-K+k}) \right. \\ &\quad \left. + \nabla f_{V, x_{(j-K+k)}}(\tilde{w}^{j-K+k}) + \nabla f_{V, x_{(j)}}(\tilde{w}^j) \right\|_2^2 \\ &\leq KL\gamma_t \frac{\gamma_{\max\{0, t-K+1\}}}{\gamma_t} \sum_{k=1}^K \sum_{j=\max\{0, t-K+k\}}^{t-1} \left\| \sum_{k=1}^K \nabla f_{\mathcal{G}(k), x_{(j-K+k)}}(\tilde{w}^{j-K+k}) \right. \\ &\quad \left. + \nabla f_{V, x_{(j-K+k)}}(\tilde{w}^{j-K+k}) + \nabla f_{V, x_{(j)}}(\tilde{w}^j) \right\|_2^2 \\ &\leq L\gamma_t \sigma K^3(K+4)M, \end{aligned} \quad (26)$$

311 where the second inequality is from Assumption I, the fourth inequality follows from that $L\gamma_t \leq 1$
312 and the last inequality follows from $\|z_1 + \dots + z_r\|_2^2 \leq r(\|z_1\|_2^2 + \dots + \|z_r\|_2^2)$, Assumption 2 and
313 $\sigma := \max_t \frac{\gamma_{\max\{0, t-K+1\}}}{\gamma_t}$. Similarly, we can bound Q_7 :

$$Q_7 \leq L\gamma_t \sigma K^2(K+4)M. \quad (27)$$

314 Integrating the upper bound of Q_1 to Q_7 in (21), we have:

$$\mathbb{E}[f(\tilde{w}^{t+1})] - f(\tilde{w}^t) \leq -\frac{\gamma_t}{2} \|\nabla f(\tilde{w}^t)\|^2 + \gamma_t^2 LM_K, \quad (28)$$

315 where we let $M_K = (K + \frac{3}{4})M + \sigma(\frac{K^2}{2} + K^3)(K+4)M$.

316 \square

317 **Proof to Theorem 1**

318 *Proof:* When γ_t is constant and $\gamma_t = \gamma$, we have $\sigma = 1$. Because of the definition of M_K and taking
 319 total expectation of (15) in Lemma 1, we obtain:

$$\mathbb{E}[f(w^{t+1})] - f(w^t) \leq -\frac{\gamma}{2} \mathbb{E} \|\nabla f(w^t)\|_2^2 + \gamma^2 L M_K. \quad (29)$$

320 Summing (29) from $t = 0$ to $T - 1$, we have:

$$\begin{aligned} \mathbb{E}[f(w^T)] - f(w^0) &\leq -\frac{\gamma}{2} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(w^t)\|_2^2 \\ &\quad + T \gamma^2 L M_K. \end{aligned} \quad (30)$$

321 Supposing w^* is the optimal solution for $f(w)$, it holds that $f(w^*) - f(w^0) \leq \mathbb{E}[f(w^T)] - f(w^0)$.
 322 Rearranging inequality (30) and dividing both sides by $\frac{2T}{\gamma}$, we complete the proof.

323 \square

324 **Proof to Theorem 2**

325 *Proof:* $\{\gamma_t\}$ is a diminishing sequence and $\gamma_t = \frac{\gamma_0}{1+t}$, such that $\sigma \leq K$ and $M_K = (K + \frac{3}{4})M +$
 326 $(\frac{K^3}{2} + K^4)(K + 4)M$. Taking total expectation of (15) in Lemma 1 and summing it from $t = 0$ to
 327 $T - 1$, we obtain:

$$\begin{aligned} \mathbb{E}[f(w^T)] - f(w^0) &\leq -\frac{1}{2} \sum_{t=0}^{T-1} \gamma_t \mathbb{E} \|\nabla f(w^t)\|_2^2 \\ &\quad + \sum_{t=0}^{T-1} \gamma_t^2 L M_K. \end{aligned} \quad (31)$$

328 Suppose w^* is the optimal solution for $f(w)$; therefore $f(w^*) - f(w^0) \leq \mathbb{E}[f(w^T)] - f(w^0)$.
 329 Rearranging inequality (31) and dividing both sides by T , we complete the proof.

330 \square