

462 **Supplementary Materials**

463 **A Proofs of Theorem 1**

464 We basically follow the arguments in the convergence of Euclidean SAM [50], but the details are totally different.

465 **Lemma 1** (Properties of Retraction Smoothness). *Let $w^* = R_w(\rho \text{grad}\mathcal{L}(w))$ and the curve $\gamma(t) = R_w(t\eta)$ with the endpoints $\gamma(0) = w$ and $\gamma(1) = w^*$. Then, we have*

$$\langle \mathcal{T}(\gamma)_{w^*}^w \text{grad}\mathcal{L}(w^*), \text{grad}\mathcal{L}(w) \rangle \geq (1 - \rho L_R) \|\text{grad}\mathcal{L}(w)\|_w^2$$

467 *Proof.* By the condition (C-5), we have

$$\|\mathcal{T}(\gamma)_{w^*}^w \text{grad}\mathcal{L}(w^*) - \text{grad}\mathcal{L}(w)\|_w \leq L_R \|\eta\|_w$$

468 By the Cauchy-Schwarz inequality, we have

$$\begin{aligned} |\langle \mathcal{T}(\gamma)_{w^*}^w \text{grad}\mathcal{L}(w^*), \eta \rangle_w| &\leq \|\mathcal{T}(\gamma)_{w^*}^w \text{grad}\mathcal{L}(w^*) - \text{grad}\mathcal{L}(w)\|_w \|\eta\|_w \\ &\leq L_R \|\eta\|_w^2 \\ &\leq \rho^2 L_R \|\text{grad}\mathcal{L}(w)\|_w^2 \end{aligned}$$

469 Therefore, we obtain

$$\langle \mathcal{T}(\gamma)_{w^*}^w \text{grad}\mathcal{L}(w^*), \rho \text{grad}\mathcal{L}(w) \rangle_w \geq -\rho^2 L_R \|\text{grad}\mathcal{L}(w)\|_w^2$$

470 Removing the constant ρ , the above inequality becomes

$$\langle \mathcal{T}(\gamma)_{w^*}^w \text{grad}\mathcal{L}(w^*), \text{grad}\mathcal{L}(w) \rangle_w \geq -\rho L_R \|\text{grad}\mathcal{L}(w)\|_w^2$$

471 Lastly, we arrive at the final result as

$$\begin{aligned} \langle \mathcal{T}(\gamma)_{w^*}^w \text{grad}\mathcal{L}(w^*), \text{grad}\mathcal{L}(w) \rangle_w &= \langle \mathcal{T}(\gamma)_{w^*}^w \text{grad}\mathcal{L}(w^*) - \text{grad}\mathcal{L}(w), \text{grad}\mathcal{L}(w) \rangle_w + \|\text{grad}\mathcal{L}(w)\|_w^2 \\ &\geq (1 - \rho L_R) \|\text{grad}\mathcal{L}(w)\|_w^2 \end{aligned}$$

472 \square

473 In the next lemma, we will show the alignment of the true Riemannian gradient and the true Riemannian SAM
474 gradient.

475 **Lemma 2** (Alignment of the true Riemannian gradient and the true Riemannian SAM gradient). *Let us
476 denote the stochastic Riemannian gradient at time t by $\text{grad}\mathcal{L}_t(w) = \frac{1}{b} \sum_{i \in J_t} \text{grad}\mathcal{L}(w; x_i) \in T_w \mathcal{M}$ and
477 $w^{adv} = R_w(\rho \text{grad}\mathcal{L}_t(w))$. Further, let $\gamma(t) = R_w(t\eta)$ be a retraction curve with $\gamma(0) = w$ and $\gamma(1) = w^{adv}$.
478 Then, we have the following inequality*

$$\mathbb{E} \left[\langle \mathcal{T}(\gamma)_{w^{adv}}^w \text{grad}\mathcal{L}_t(w^{adv}), \text{grad}\mathcal{L}(w) \rangle_w \right] \geq \left(\frac{1}{2} - \rho L_R - 3\rho^2 L_R^2 \right) \|\text{grad}\mathcal{L}(w)\|_w^2 - \frac{2\rho^2 L_R^2 \sigma^2}{b}$$

479 *Proof.* Let $w^* = R_w(\rho \text{grad}\mathcal{L}(w))$ evaluated on the loss function. We first add and subtract
480 $\langle \mathcal{T}(\zeta)_{w^*}^w \text{grad}\mathcal{L}_t(w^*), \text{grad}\mathcal{L}(w) \rangle_w$ where $\zeta(t) = R_w(t\xi)$ is a retraction curve where $\zeta(0) = w$ and
481 $\zeta(1) = w^*$.

$$\begin{aligned} \langle \mathcal{T}(\gamma)_{w^{adv}}^w \text{grad}\mathcal{L}_t(w^{adv}), \text{grad}\mathcal{L}(w) \rangle_w &= \underbrace{\langle \mathcal{T}(\gamma)_{w^{adv}}^w \text{grad}\mathcal{L}_t(w^{adv}) - \mathcal{T}(\zeta)_{w^*}^w \text{grad}\mathcal{L}_t(w^*), \text{grad}\mathcal{L}(w) \rangle_w}_{T_1} \\ &\quad + \underbrace{\langle \mathcal{T}(\zeta)_{w^*}^w \text{grad}\mathcal{L}_t(w^*), \text{grad}\mathcal{L}(w) \rangle_w}_{T_2} \end{aligned}$$

482 We will bound two terms, T_1 and T_2 , separately. Regarding the term T_1 , we derive

$$\begin{aligned} -T_1 &= -\langle \mathcal{T}(\gamma)_{w^{adv}}^w \text{grad}\mathcal{L}_t(w^{adv}) - \mathcal{T}(\zeta)_{w^*}^w \text{grad}\mathcal{L}_t(w^*), \text{grad}\mathcal{L}(w) \rangle_w \\ &\leq \frac{1}{2} \left\| \mathcal{T}(\gamma)_{w^{adv}}^w \text{grad}\mathcal{L}_t(w^{adv}) - \mathcal{T}(\zeta)_{w^*}^w \text{grad}\mathcal{L}_t(w^*) \right\|_w^2 + \frac{1}{2} \|\text{grad}\mathcal{L}(w)\|_w^2 \\ &\leq \left\| \mathcal{T}(\gamma)_{w^{adv}}^w \text{grad}\mathcal{L}_t(w^{adv}) - \text{grad}\mathcal{L}_t(w) \right\|_w^2 + \left\| \mathcal{T}(\zeta)_{w^*}^w \text{grad}\mathcal{L}_t(w^*) - \text{grad}\mathcal{L}_t(w) \right\|_w^2 + \frac{1}{2} \|\text{grad}\mathcal{L}(w)\|_w^2 \\ &\leq L_R^2 \|\rho \text{grad}\mathcal{L}_t(w)\|_w^2 + L_R^2 \|\rho \text{grad}\mathcal{L}(w)\|_w^2 + \frac{1}{2} \|\text{grad}\mathcal{L}(w)\|_w^2 \\ &\leq \rho^2 L_R^2 \left(2 \|\text{grad}\mathcal{L}_t(w) - \text{grad}\mathcal{L}(w)\|_w^2 + 2 \|\text{grad}\mathcal{L}(w)\|_w^2 \right) + \left(\frac{1}{2} + \rho^2 L_R^2 \right) \|\text{grad}\mathcal{L}(w)\|_w^2 \\ &\leq \frac{2\rho^2 L_R^2 \sigma^2}{b} + \left(\frac{1}{2} + 3\rho^2 L_R^2 \right) \|\text{grad}\mathcal{L}(w)\|_w^2 \end{aligned}$$

483 From the above inequality, we could finally bound the term T_1 as

$$T_1 \geq -\frac{2\rho^2 L_R^2 \sigma^2}{b} - \left(\frac{1}{2} + 3\rho^2 L_R^2\right) \|\text{grad}\mathcal{L}(w)\|_w^2$$

484 Regarding the term T_2 , we just use the lemma as

$$T_2 = \langle \mathcal{T}(\zeta)_{w^*}^w \text{grad}\mathcal{L}_t(w^*), \text{grad}\mathcal{L}(w) \rangle_w \geq (1 - \rho L_R) \|\text{grad}\mathcal{L}(w)\|_w^2$$

485 Hence, we arrive at

$$\mathbb{E} \left[\langle \mathcal{T}(\gamma)_{w^{adv}}^w \text{grad}\mathcal{L}_t(w^{adv}), \text{grad}\mathcal{L}(w) \rangle_w \right] \geq \left(\frac{1}{2} - \rho L_R - 3\rho^2 L_R^2 \right) \|\text{grad}\mathcal{L}(w)\|_w^2 - \frac{2\rho^2 L_R^2 \sigma^2}{b}$$

486 \square

487 According to Algorithm 1, we follow the notation as

$$\begin{aligned} \text{grad}\mathcal{L}_t(w) &= \frac{1}{b} \sum_{i \in I_t} \text{grad}\ell_i(w) \\ w_t^{adv} &= \mathbf{R}_{w_t}(\rho \text{grad}\mathcal{L}_t(w_t)) \end{aligned}$$

488 We assume the stochastic m -SAM where the same batch is used for both inner and outer updates.

489 **Lemma 3** (Descent inequality). *Under the assumptions in Theorem 1, we have*

$$\mathbb{E}[\mathcal{L}(w_{t+1})] \leq \mathbb{E}[\mathcal{L}(w_t)] - \frac{3\alpha}{8} \mathbb{E}[\|\text{grad}\mathcal{L}(w_t)\|_{w_t}^2] + \frac{\alpha^2 L_S^2 \sigma^2}{b} + \frac{2\alpha^2 L_S^3 \rho^2 \sigma^2}{b} + \frac{2\alpha \rho^3 L_R^2 \sigma^2}{b}$$

490 *Proof.* Using the condition (C-4), we have

$$\begin{aligned} \mathcal{L}(w_{t+1}) &= \mathcal{L} \left(\mathbf{R}_{w_t} \left(-\alpha \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}) \right) \right) \\ &\leq \mathcal{L}(w_t) - \alpha \left\langle \text{grad}\mathcal{L}(w_t), \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}) \right\rangle_{w_t} + \frac{\alpha^2 L_S}{2} \left\| \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}) \right\|_{w_t}^2 \end{aligned}$$

491 For the last term in RHS, we can bound as

$$\begin{aligned} \left\| \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}) \right\|_{w_t}^2 &= -\|\text{grad}\mathcal{L}(w_t)\|_{w_t}^2 + \left\| \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}) - \text{grad}\mathcal{L}(w_t) \right\|_{w_t}^2 \\ &\quad + 2 \left\langle \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}), \text{grad}\mathcal{L}(w_t) \right\rangle_{w_t} \end{aligned}$$

492 Again, we have

$$\begin{aligned} \mathcal{L}(w_{t+1}) &\leq \mathcal{L}(w_t) - \alpha \left\langle \text{grad}\mathcal{L}(w_t), \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}) \right\rangle_{w_t} + \frac{\alpha^2 L_S}{2} \left\| \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}) \right\|_{w_t}^2 \\ &= \mathcal{L}(w_t) - \frac{\alpha^2 L_S}{2} \|\text{grad}\mathcal{L}(w_t)\|_{w_t}^2 + \frac{\alpha^2 L_S}{2} \left\| \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}) - \text{grad}\mathcal{L}(w_t) \right\|_{w_t}^2 \\ &\quad - \alpha(1 - \alpha L_S) \left\langle \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}), \text{grad}\mathcal{L}(w_t) \right\rangle_{w_t} \\ &\leq \mathcal{L}(w_t) - \frac{\alpha^2 L_S}{2} \|\text{grad}\mathcal{L}(w_t)\|_{w_t}^2 + \alpha^2 L_S \left\| \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}) - \text{grad}\mathcal{L}(w_t) \right\|_{w_t}^2 \\ &\quad + \alpha^2 L_S \|\text{grad}\mathcal{L}_t(w_t) - \text{grad}\mathcal{L}(w_t)\|_{w_t}^2 - \alpha(1 - \alpha L_S) \left\langle \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}), \text{grad}\mathcal{L}(w_t) \right\rangle_{w_t} \\ &\leq \mathcal{L}(w_t) - \frac{\alpha^2 L_S}{2} \|\text{grad}\mathcal{L}(w_t)\|_{w_t}^2 + \alpha^2 L_S^3 \rho^2 \|\text{grad}\mathcal{L}_t(w_t)\|_{w_t}^2 + \alpha^2 L_S \|\text{grad}\mathcal{L}_t(w_t) - \text{grad}\mathcal{L}(w_t)\|_{w_t}^2 \\ &\quad - \alpha(1 - \alpha L_S) \left\langle \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}), \text{grad}\mathcal{L}(w_t) \right\rangle_{w_t} \\ &\leq \mathcal{L}(w_t) - \frac{\alpha^2 L_S}{2} \|\text{grad}\mathcal{L}(w_t)\|_{w_t}^2 + 2\alpha^2 L_S^3 \rho^2 \|\text{grad}\mathcal{L}(w_t)\|_{w_t}^2 + 2\alpha^2 L_S^3 \rho^2 \|\text{grad}\mathcal{L}_t(w_t) - \text{grad}\mathcal{L}(w_t)\|_{w_t}^2 \\ &\quad + \alpha^2 L_S \|\text{grad}\mathcal{L}_t(w_t) - \text{grad}\mathcal{L}(w_t)\|_{w_t}^2 - \alpha(1 - \alpha L_S) \left\langle \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}), \text{grad}\mathcal{L}(w_t) \right\rangle_{w_t} \\ &= \mathcal{L}(w_t) - \frac{\alpha^2 L_S(1 - 4L_S^2 \rho^2)}{2} \|\text{grad}\mathcal{L}(w_t)\|_{w_t}^2 + \alpha^2 L_S(1 + 2L_S^2 \rho^2) \|\text{grad}\mathcal{L}_t(w_t) - \text{grad}\mathcal{L}(w_t)\|_{w_t}^2 \\ &\quad - \alpha(1 - \alpha L_S) \left\langle \mathcal{T}(\gamma)_{w_t^{adv}}^w \text{grad}\mathcal{L}_t(w_t^{adv}), \text{grad}\mathcal{L}(w_t) \right\rangle_{w_t} \end{aligned}$$

493 Taking the expectation on both sides, we have

$$\begin{aligned}\mathbb{E}[\mathcal{L}(w_{t+1})] &\leq \mathbb{E}[\mathcal{L}(w_t)] - \frac{\alpha^2 L_S(1 - 4L_S^2\rho^2)}{2} \mathbb{E}[\|\text{grad}\mathcal{L}(w_t)\|_{w_t}^2] + \frac{\alpha^2 L_S(1 + 2L_S^2\rho^2)\sigma^2}{b} \\ &\quad - \alpha(1 - \alpha L_S) \mathbb{E} \left[\left\langle \mathcal{T}(\gamma)_{w_t^{adv}}^{\omega_t} \text{grad}\mathcal{L}_t(w_t^{adv}), \text{grad}\mathcal{L}(w_t) \right\rangle_{w_t} \right] \\ &= \mathbb{E}[\mathcal{L}(w_t)] - \frac{\alpha^2 L_S(1 - 4L_S^2\rho^2)}{2} \mathbb{E}[\|\text{grad}\mathcal{L}(w_t)\|_{w_t}^2] + \frac{\alpha^2 L_S(1 + 2L_S^2\rho^2)\sigma^2}{b} \\ &\quad - \alpha(1 - \alpha L_S) \left[\left(\frac{1}{2} - \rho L_R - 3\rho^2 L_R^2 \right) \mathbb{E}[\|\text{grad}\mathcal{L}(w_t)\|_{w_t}^2] - \frac{2\rho^2 L_R^2 \sigma^2}{b} \right]\end{aligned}$$

494 For sufficiently large number of total iteration T , the condition $\rho \leq \frac{1}{4L}$ is easily satisfied where $\tilde{L} =$
495 $\max\{L_R, L_S\}$ (defined in Theorem 1). Hence, we obtain

$$\begin{aligned}\frac{3\alpha}{8} \mathbb{E}[\|\text{grad}\mathcal{L}(w_t)\|_{w_t}^2] &\leq \mathbb{E}[\mathcal{L}(w_t)] - \mathbb{E}[\mathcal{L}(w_{t+1})] + \frac{\alpha^2 L_S(1 + 2L_S^2\rho^2)\sigma^2}{b} + \frac{2\alpha(1 - \alpha L_S)\rho^3 L_R^2 \sigma^2}{b} \\ &\leq \mathbb{E}[\mathcal{L}(w_t)] - \mathbb{E}[\mathcal{L}(w_{t+1})] + \frac{\alpha^2 L_S \sigma^2}{b} + \frac{2\alpha^2 L_S^3 \rho^2 \sigma^2}{b} + \frac{2\alpha \rho^3 L_R^2 \sigma^2}{b}\end{aligned}$$

496 By telescoping the above inequality from $t = 0 \sim T - 1$, we arrive at

$$\mathbb{E}[\|\text{grad}\mathcal{L}(\tilde{w})\|_w^2] = \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\text{grad}\mathcal{L}(w_t)\|_{w_t}^2] \leq \frac{8\Delta}{3\alpha T} + \frac{8\alpha^2 L_S \sigma^2}{3b} + \frac{16\alpha^2 L_S^3 \rho^2 \sigma^2}{3b} + \frac{16\alpha \rho^3 L_R^2 \sigma^2}{3b}$$

497 Under the step size condition $\alpha_t = \frac{1}{\sqrt{T}\tilde{L}}$ and $\rho_t = \frac{1}{T^{1/6}\tilde{L}}$, we finally get

$$\mathbb{E}[\|\text{grad}\mathcal{L}(\tilde{w})\|_{\tilde{w}}^2] \leq \frac{Q_1 \tilde{L} \Delta}{\sqrt{T}} + \frac{Q_2 \sigma^2}{b\sqrt{T}} + \frac{Q_3 \sigma^2}{bT^{5/6}}$$

498 for appropriate constants $\{Q_i\}_{i=1}^3$. □

499 B Hyperparameter Details

500 We use the almost same hyperparameters in the study [51] and implement our experiments in Section 5 upon
 501 its official implementation. For completeness, we summarize the hyperparameter configurations in Table 4 and
 502 Table 5.

Table 4: Hyperparameter configurations for knowledge graph completion.

| Dimension | WN18RR | | FB15k-237 | |
|------------------|--------|---------|-----------|---------|
| | 32 | β | 32 | β |
| Batch Size | 1000 | 1000 | 500 | 500 |
| Negative Samples | 50 | 50 | 50 | 50 |
| Margin | 8.0 | 8.0 | 8.0 | 8.0 |
| Epochs | 1000 | 1000 | 500 | 500 |
| Max Norm | 1.5 | 2.5 | 1.5 | 1.5 |
| Max Scaler | 3.5 | 2.5 | 2.5 | 2.5 |
| Learning Rate | 0.005 | 0.003 | 0.003 | 0.003 |
| Gradient Norm | 0.5 | 0.5 | 0.5 | 0.5 |

Table 5: Hyperparameter configurations for machine translation.

| Hyperparameter | IWSLT'14 | WMT'14 |
|-----------------------------|----------|--------|
| GPU Numbers | 4 | 4 |
| Embedding Dimension d | 64 | 64 |
| Feed-forward Dimension | 256 | 256 |
| Batch Type | Token | Token |
| Batch Size | 10240 | 10240 |
| Gradient Accumulation Steps | 1 | 1 |
| Training Steps | 40000 | 200000 |
| Dropout | 0.0 | 0.1 |
| Attention Dropout | 0.1 | 0.0 |
| Max Gradient Norm | 0.5 | 0.5 |
| Warmup Steps | 8000 | 6000 |
| Decay Method | noam | noam |
| Label Smoothing | 0.1 | 0.1 |
| Layer Number | 6 | 6 |
| Head Number | 4 | 8 |
| Learning Rate | 5.0 | 5.0 |
| Adam β_2 | 0.998 | 0.998 |