Counterfactual Critic Multi-Agent Training for Scene Graph Generation — Supplementary Document —

This supplementary document is organized as follows:

- Section 1 provides the details of some simplified functions in Agent Communication and Visual Relationship Detection.
- Section 2 provides the detailed proof of the CMAT convergence, which guarantees that the proposed CMAT method can converge to a locally optimal policy.
- Section 3 provides the detailed derivation of Eq. (6), *i.e.*, $\nabla_{\theta} J \approx \sum_{i=1}^{n} \nabla_{\theta} \log p_i^T (v_i^T | h_i^T; \theta) Q(H^T, V^T).$
- Section 4 shows more qualitative results of CMAT compared with the strong baseline MOTIFS [4] in SGDet setting.

1. Details of Some Simplified Functions

We demonstrate the details of some omitted functions in Eq. (1), Eq. (2), Eq. (3) and Eq. (4).

1.1. F_s and F_e in Extract Module

$$\begin{aligned} \boldsymbol{h}_{i}^{t} &= \mathrm{LSTM}(\boldsymbol{h}_{i}^{t-1}, [\boldsymbol{x}_{i}^{t}, \boldsymbol{e}_{i}^{t-1}]), \\ \boldsymbol{s}_{i}^{t} &= \boldsymbol{s}_{i}^{t-1} + \boldsymbol{W}_{h} \boldsymbol{h}_{i}^{t}, \\ \boldsymbol{v}_{i}^{t} &\sim \boldsymbol{p}_{i}^{t} = \mathrm{softmax}(\boldsymbol{s}_{i}^{t}), \\ \boldsymbol{e}_{i}^{t} &= \sum_{\tilde{v}} \boldsymbol{p}_{i}^{t}(\tilde{v}) \mathbf{E}[\tilde{v}], \end{aligned}$$

$$(11)$$

where $h_i^t \in \mathbb{R}^h$ is the hidden state of LSTM, $x_i^t \in \mathbb{R}^d$ is the time-step input feature and $s_i^t \in \mathbb{R}^{|\mathcal{C}|}$ is the object class confidence. $\mathbf{E}[\tilde{v}] \in \mathbb{R}^e$ is the embedding of class label $\tilde{v} \in \mathcal{C}$ and $e_i^t \in \mathbb{R}^e$ is the soft-weighted embedding of class label based on probabilities p_i^t . $W_h \in \mathbb{R}^{h \times |\mathcal{C}|}$ is a learnable matrix. and [,] is concatenate operation.

1.2. F_{m*} in Message Module

$$\boldsymbol{m}_{j}^{t} = \boldsymbol{W}_{u}\boldsymbol{h}_{j}^{t}, \ \boldsymbol{m}_{ij}^{t} = \boldsymbol{W}_{p}\boldsymbol{h}_{ij}^{t}$$
 (12)

where $\boldsymbol{m}_{j}^{t} \in \mathbb{R}^{h}$ and $\boldsymbol{m}_{ij}^{t} \in \mathbb{R}^{h}$ is the unary message and pairwise message, respectively. $\boldsymbol{h}_{ij}^{t} \in \boldsymbol{R}^{d}$ is the pairwsie feature between agent *i* and agent *j*. $\boldsymbol{W}_{u} \in \mathbb{R}^{h \times h}$, $\boldsymbol{W}_{p} \in \mathbb{R}^{d \times h}$ are learnable mapping matrices.

1.3. F_{att*} and F_{u*} in Update Module

$$u_{j}^{t} = \boldsymbol{w}_{u}[\boldsymbol{h}_{i}^{t}, \boldsymbol{h}_{j}^{t}], \ \alpha_{j}^{t} = \exp(u_{j}^{t}) / \sum_{k} \exp(u_{k}^{t}),$$

$$u_{ij}^{t} = \boldsymbol{w}_{p}[\boldsymbol{h}_{i}^{t}, \boldsymbol{h}_{ij}^{t}], \ \alpha_{ij}^{t} = \exp(u_{ij}^{t}) / \sum_{k} \exp(u_{ik}^{t})$$

$$\boldsymbol{x}_{i}^{t+1} = \boldsymbol{W}_{x}(\text{ReLU}(\boldsymbol{h}_{i}^{t} + \sum_{j} \alpha_{j}^{t} \boldsymbol{m}_{j}^{t} + \sum_{j} \alpha_{ij}^{t} \boldsymbol{m}_{ij}^{t}))$$

$$\boldsymbol{h}_{ij}^{t+1} = \text{ReLU}(\boldsymbol{h}_{ij}^{t} + \boldsymbol{W}_{s} \boldsymbol{h}_{i}^{t+1} + \boldsymbol{W}_{e} \boldsymbol{h}_{j}^{t+1})$$
(13)

where α_j^t and α_{ij}^t are attention weights to fuse different messages, $\boldsymbol{w}_u \in \mathbb{R}^{2h}$, $\boldsymbol{w}_p \in \mathbb{R}^{h+d}$, $\boldsymbol{W}_x \in \mathbb{R}^{h \times d}$, $\boldsymbol{W}_s \in \mathbb{R}^{h \times d}$, and $\boldsymbol{W}_e \in \mathbb{R}^{h \times d}$ are learnable mapping matrices.

1.4. *F_r* in Visual Relationship Detection

$$\boldsymbol{z}_{i} = \boldsymbol{W}_{o}[\boldsymbol{h}_{i}^{T}, \mathbf{E}[\boldsymbol{v}_{i}^{T}]], \ \boldsymbol{z}_{j} = \boldsymbol{W}_{o}[\boldsymbol{h}_{j}^{T}, \mathbf{E}[\boldsymbol{v}_{j}^{T}]],$$
$$\boldsymbol{p}_{ij} = \operatorname{softmax}([\boldsymbol{z}_{i}, \boldsymbol{z}_{j}] \odot \boldsymbol{W}_{r} \boldsymbol{z}_{ij} + \boldsymbol{w}_{\boldsymbol{v}_{i}^{T}, \boldsymbol{v}_{j}^{T}}), \qquad (14)$$
$$\boldsymbol{r}_{ij} = \arg \max_{r \in \mathcal{R}} \boldsymbol{p}_{ij}(r),$$

where $W_o \in \mathbb{R}^{(h+e)\times z}$, $W_r \in \mathbb{R}^{z\times 2z}$ are transformation matrices, $z_{ij} \in \mathbb{R}^z$ is the predicated visual feature between agent *i* and *j*, \odot is a fusing function ¹, and $w_{v_i^T, v_j^T} \in \mathbb{R}^{|\mathcal{C}|}$ is the bias vector specific to head and tail labels as in [4].

Predicate Visual Features z_{ij} . For the predicate visual features, we used RoIAlign to pool the union box of subject and object, and resized the union box feature to $7 \times 7 \times 512$. Following [4, 1], we used a $14 \times 14 \times 2$ binary feature map to model the geometric spatial position of subject and object, with one channel per box. We applied two convolutional layers on this binary feature map and obtained a new $7 \times 7 \times 512$ spatial position feature map. We added this position feature map with the previous resized union box feature, and applied two fully-connected layers to obtain the final predicate visual feature.

¹Different functions get comparable performance. In our experiments, we follow [5]: $\mathbf{x} \odot \mathbf{y} = \text{ReLU}(\mathbf{W}_x \mathbf{x} + \mathbf{W}_y \mathbf{y}) - (\mathbf{W}_x \mathbf{x} - \mathbf{W}_y \mathbf{y})^2$.

2. Proof of the Convergence of CMAT

Proof. We denote π_i as the policy of agent *i*, *i.e.*, $\pi_i = p_i^T$ and π as the joint policy of all agents, *i.e.*, $\pi = \{p_1^T, ..., p_n^T\}$. Then, the expected gradient of CMAT is given by (cf. Eq. (11)):

$$\nabla_{\theta} J = \mathbb{E}_{\boldsymbol{\pi}} \left[\sum_{i=1}^{n} \nabla_{\theta} \log \pi_{i}(v_{i}^{T}) A^{i}(H^{T}, V^{T}) \right],$$

$$= \mathbb{E}_{\boldsymbol{\pi}} \left[\sum_{i=1}^{n} \nabla_{\theta} \log \pi_{i}(v_{i}^{T}) (R(H^{T}, V^{T}) - b(H^{T}, V_{-i}^{T})) \right].$$
(15)

where the expection \mathbb{E}_{π} is with respect to the state-action distribution induced by the joint policy π , $b(H^T, V_{-i}^T)$ is the counterfactual baseline in CMAT model, *i.e.*, $b(H^T, V_{-i}^T) = \sum p_i^T (\tilde{v}_i^T) R(H^T, (V_{-i}^T, \tilde{v}_i^T))$.

First, consider the expected contribution of this counterfactual baseline $b(H^T, V_{-i}^T)$,

$$\nabla_{\theta} J_b = \mathbb{E}_{\pi} \left[\sum_{i=1}^n \nabla_{\theta} \log \pi_i(v_i^T) b(H^T, V_{-i}^T) \right].$$
(16)

Let $d^{\pi}(s)$ be the discounted ergodic state distribution as defined by [3]:

$$\nabla_{\theta} J_{b} = \sum_{s} d^{\pi}(s) \sum_{i=1}^{n} \sum_{V_{-i}^{T}} \pi(V_{-i}^{T}) \sum_{v_{i}^{T}} \pi_{i}(v_{i}^{T}) \nabla_{\theta} \log \pi_{i}(v_{i}^{T}) b(H^{T}, V_{-i}^{T})
= \sum_{s} d^{\pi}(s) \sum_{i=1}^{n} \sum_{V_{-i}^{T}} \pi(V_{-i}^{T}) \sum_{v_{i}^{T}} \nabla_{\theta} \pi_{i}(v_{i}^{T}) b(H^{T}, V_{-i}^{T})
= \sum_{s} d^{\pi}(s) \sum_{i=1}^{n} \sum_{V_{-i}^{T}} \pi(V_{-i}^{T}) b(H^{T}, V_{-i}^{T}) \nabla_{\theta} 1
= 0$$
(17)

Thus, this counterfactual baseline does not change the expected gradient. The reminder of the expected policy gradient is given by:

$$\nabla_{\theta} J = \mathbb{E}_{\pi} \left[\sum_{i=1}^{n} \nabla_{\theta} \log \pi_{i}(v_{i}^{T}) R(H^{T}, V^{T}) \right],$$
(18)

$$= \mathbb{E}_{\boldsymbol{\pi}} \left[\nabla_{\boldsymbol{\theta}} \log \prod_{i=1}^{n} \pi_i(v_i^T) R(H^T, V^T) \right].$$
(19)

Writing the joint policy into a product of the independent policies:

$$\pi(V^T) = \prod_{i=1}^n \pi_i(v_i^T),$$
(20)

we have the standard single-agent policy gradient:

$$\nabla_{\theta} J = \mathbb{E}_{\boldsymbol{\pi}} \left[\nabla_{\theta} \log \boldsymbol{\pi}(V^T) R(H^T, V^T) \right].$$
(21)

Konda *et al.* [2] proved that this gradient converges to a local maximum of the expected return J, given that: 1) the policy π is differentiable, 2) the update timescales for π are sufficiently slow. Meanwhile, the parameterization of the policy (*i.e.*, the single-agent joint-action learner is decomposed into independent policies) is immaterial to convergence, as long as it remains differentiable.

3. Derivation of Eq. (6)

Based on the policy gradient theorem we provide the detailed derivation of Eq. (6) as follows. We denote the action sequence for agent *i* as $\hat{A}_i = \{\hat{a}_i^1, \hat{a}_i^2, ..., \hat{a}_i^T\}$, and value function $V_{\theta}(\hat{A}_i)$ as the expected future reward of sequence \hat{A}_i . Then the gradient of agent *i* is:

$$\begin{aligned} \nabla_{\theta} J_{i} &= \frac{dV(A_{i})}{d\theta} = \frac{d}{d\theta} \mathbb{E}_{\hat{A}_{i} \sim \pi_{i}^{t}(\hat{A}_{i})} R(\hat{A}_{i}) \\ &= \sum_{\hat{A}_{i}} \frac{d}{d\theta} \left[\pi_{i}^{t}(\hat{a}_{i}^{1}) \pi_{i}^{t}(\hat{a}_{i}^{2}|\hat{a}_{i}^{1}) \dots \pi_{i}^{t}(\hat{a}_{i}^{T}|\hat{a}_{i}^{1}\dots\hat{a}_{i}^{T-1}) \right] R(\hat{A}_{i}) \\ &= \sum_{t=1}^{T} \sum_{\hat{A}_{i}} \pi_{i}^{t}(\hat{A}_{i}^{1\dots t-1}) \frac{d\pi_{i}^{t}(\hat{a}_{i}^{t}|\hat{A}_{i}^{1\dots t-1})}{d\theta} \pi_{i}^{t}(\hat{A}_{i}^{t+1\dots T}|\hat{A}_{i}^{1\dots t}) R(\hat{A}_{i}) \\ &= \sum_{t=1}^{T} \sum_{\hat{A}_{i}} \pi_{i}^{t}(\hat{A}_{i}^{1\dots t-1}) \frac{\pi_{i}^{t}(\hat{a}_{i}^{t}|\hat{A}_{i}^{1\dots t-1})}{d\theta} \sum_{\hat{A}_{i}^{t+1\dots T}} \pi_{i}^{t}(\hat{A}_{i}^{t+1\dots T}|\hat{A}_{i}^{1\dots t}) \sum_{\tau=1}^{T} r_{i}^{\tau}(\hat{a}_{i}^{\tau}; \hat{A}_{i}^{1\dots \tau-1}) \\ &= \sum_{t=1}^{T} \sum_{\hat{A}_{i}^{1\dots t}} \pi_{i}^{t}(\hat{A}_{i}^{1\dots t-1}) \frac{\pi_{i}^{t}(\hat{a}_{i}^{t}|\hat{A}_{i}^{1\dots t-1})}{d\theta} \left[r^{t}(\hat{a}_{i}^{t}; \hat{A}_{i}^{1\dots t-1}) + \sum_{\hat{A}_{i}^{t+1\dots T}} \pi_{i}^{t}(\hat{Y}_{i}^{t+1\dots T}|\hat{Y}_{i}^{1\dots t}) \sum_{\tau=t+1}^{T} r_{i}^{\tau}(\hat{a}_{i}^{\tau}; \hat{A}_{i}^{1\dots \tau-1}) \right] \\ &= \sum_{t=1}^{T} \sum_{\hat{A}_{i}^{1\dots t}} \pi_{i}^{t}(\hat{A}_{i}^{1\dots t-1}) \sum_{a^{t} \in \mathcal{A}} \frac{d\pi_{i}^{t}(a^{t}|\hat{A}_{i}^{1\dots t-1})}{d\theta} Q(a^{t}; \hat{A}_{i}^{1\dots t-1}) \\ &= \mathbb{E}_{\hat{A}_{i} \sim \pi_{i}^{t}(\hat{A}_{i})} \sum_{t=1}^{T} \sum_{a^{t} \in \mathcal{A}} \frac{\pi_{i}^{t}(a^{t}|\hat{A}_{i}^{1\dots t-1})}{d\theta} Q(a^{t}; \hat{A}_{i}^{1\dots t-1}) \end{aligned}$$

Further, the gradient for agent i can be simplified as:

$$\nabla_{\theta} J_{i} = \mathbb{E} \left[\sum_{t=1}^{T} \sum_{a_{i}^{t} \in \mathcal{A}} \nabla_{\theta} \pi_{i}^{t}(a_{i}^{t}) Q(s_{i}^{t}, a_{i}^{t}) \right]$$
$$= \mathbb{E} \left[\sum_{t=1}^{T} \sum_{a_{i}^{t} \in \mathcal{A}} \pi_{i}^{t}(a_{i}^{t}) \nabla_{\theta} \log \pi_{i}^{t}(a_{i}^{t}) Q(s_{i}^{t}, a_{i}^{t}) \right]$$
$$\approx \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{i}^{t}(a_{i}^{t}) Q(s_{i}^{t}, a_{i}^{t})$$
(23)

Therefore, for the time step t, the gradient for agent i is $\nabla_{\theta} \log \pi_i^t(a_i^t)Q(s_i^t, a_i^t)$. For multi-agent in a cooperative environment, the state-action function Q should estimate the reward based on the set of all agent state and actions, *i.e.*, $Q(S^t, A^t)$. Then, the gradient for all agents is:

$$\nabla_{\theta} J \approx \sum_{i=1}^{n} \nabla_{\theta} J_{i} = \sum_{i=1}^{n} \nabla_{\theta} \log \pi_{i}^{t}(a_{i}^{t} | s_{i}^{t}) Q(S^{t}, A^{t}).$$

$$(24)$$

In CMAT, we samples actions after T-round agent communication, and the action for agent *i* is v_i^T , the policy function is p_i^T , and the state of agent is h^t , *i.e.*, $S^t = H^t$, $A^t = V^t$. Therefore, the gradient for the cooperative multi-agent in CMAT is:

$$\nabla_{\theta} J \approx \sum_{i=1}^{n} \nabla_{\theta} \log \boldsymbol{p}_{i}^{T}(\boldsymbol{v}_{i}^{T} | \boldsymbol{h}_{i}^{T}; \theta) Q(\boldsymbol{H}^{T}, \boldsymbol{V}^{T})$$
(25)

4. More Qualitative Results

Figure 7 and 8 show more qualitative results of CMAT and MOTIFS in SGDet setting. From the rows where CMAT is better than MOTIFS, we can see that CMAT rarely mistakes at the important hub nodes such as the "surfboard" or "laptop".

This is because CMAT directly optimizes the graph-coherent objective. However, the rows show that the mistakes made by CMAT always come from the imcomplete anntation of CMAT can detect more false positive (the blue color) objects and relationship than MOTIFS. Since the evaluation metric (i.e., Recall@K) is based on the ranking of labeled triplet confidence, thus, detecting more reasonable false positive results with high confidence can worsen the performance.

References

- [1] Bo Dai, Yuqi Zhang, and Dahua Lin. Detecting visual relationships with deep relational networks. In CVPR, 2017. 1
- [2] Vijay R Konda and John N Tsitsiklis. Actor-critic algorithms. In NeurIPS, 2000. 2
- [3] Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *NeurIPS*, 2000. 2
- [4] Rowan Zellers, Mark Yatskar, Sam Thomson, and Yejin Choi. Neural motifs: Scene graph parsing with global context. In *CVPR*, 2018. 1
- [5] Yan Zhang, Jonathon Hare, and Adam Prügel-Bennett. Learning to count objects in natural images for visual question answering. In *ICLR*, 2018. 1



Figure 7: More qualitative results showing comparisons between CMAT and MOTIFS in the SGDet setting. Green boxes are detected boxes with IoU large 0.5 with the ground truth, blue boxes are detected but not labeled, red boxes are ground-truth with no match. Green edges are true positive predicted by each model at the R@20 setting, red edges are false negatives, and blue edges are false positives. Only detected boxes overlapped with ground-truth are shown.



Figure 8: More qualitative results showing comparisons between CMAT and MOTIFS in the SGDet setting continued from Figure 7.