---

**Algorithm 1** Conditional EBM Training Algorithm

---

**Input:** data dist $p_D(\boldsymbol{x})$, relational scene descriptions $R_D(\boldsymbol{r})$, step size $\lambda$, number of steps $K$, data augmentation $D(\cdot)$, stop gradient operator $\Omega(\cdot)$, EBM $E_\theta(\cdot)$, Encoder Enc$(\cdot)$, Parser $P(\cdot)$

$\mathcal{B} \leftarrow \varnothing$

**while** not converged **do**

$\quad \boldsymbol{x}_i^+ \sim p_D$

$\quad R_i \sim R_D$

$\quad \tilde{\boldsymbol{x}}_i^0 \sim \mathcal{B}$ with 99% probability and $\mathcal{U}$ otherwise

$\quad X \sim \mathcal{B}$ for nearest neighbor entropy calculation

$\quad \triangleright$ *Parse a relational scene description:*

$\quad \{\boldsymbol{r}_1, \ldots \boldsymbol{r}_m\} \leftarrow P(R_i)$

$\quad \triangleright$ *Apply data augmentation to sample:*

$\quad \tilde{\boldsymbol{x}}_i^0 = D(\tilde{\boldsymbol{x}}_i^0)$

$\quad \triangleright$ *Generate sample using Langevin dynamics:*

$\quad$ **for** sample step $k = 1$ to $K$ **do**

$\quad\quad \tilde{\boldsymbol{x}}_i^{k-1} = \Omega(\tilde{\boldsymbol{x}}_i^{k-1})$

$\quad\quad \tilde{\boldsymbol{x}}^k \leftarrow \tilde{\boldsymbol{x}}^{k-1} - \nabla_{\boldsymbol{x}} \sum_{j=1}^m E_\theta(\tilde{\boldsymbol{x}}^{k-1} \mid \text{Enc}(\boldsymbol{r}_j)) + \omega, \ \omega \sim \mathcal{N}(0, \sigma)$

$\quad$ **end for**

$\quad \triangleright$ *Generate two variants of $\boldsymbol{x}^-$ with and without gradient propagation:*

$\quad \boldsymbol{x}_i^- = \Omega(\tilde{\boldsymbol{x}}_i^k)$

$\quad \hat{\boldsymbol{x}}_i^- = \tilde{\boldsymbol{x}}_i^k$

$\quad \triangleright$ *Optimize objective $\mathcal{L}_{CD} + \mathcal{L}_{KL}$ wrt $\theta$:*

$\quad \mathcal{L}_{\text{CD}} = \frac{1}{N} \sum_i \sum_{j=1}^m (E_\theta(\boldsymbol{x}_i^+ \mid \text{Enc}(\boldsymbol{r}_j) - E_\theta(\boldsymbol{x}_i^- \mid \text{Enc}(\boldsymbol{r}_j))$

$\quad \mathcal{L}_{\text{KL}} = \sum_{j=1}^m E_{\Omega(\theta)}(\hat{\boldsymbol{x}}_i^- \mid \text{Enc}(\boldsymbol{r}_j)) - \log(NN(\hat{\boldsymbol{x}}_i^-, X)$

$\quad \triangleright$ *Optimize objective $\mathcal{L}_{CD} + \mathcal{L}_{KL}$ wrt $\theta$:*

$\quad \Delta\theta \leftarrow \nabla_\theta(\mathcal{L}_{\text{CD}} + \mathcal{L}_{\text{KL}})$

$\quad$ Update $\theta$ based on $\Delta\theta$ using Adam optimizer

$\quad \triangleright$ *Update replay buffer $\mathcal{B}$*

$\quad \mathcal{B} \leftarrow \mathcal{B} \cup \tilde{\boldsymbol{x}}_i^-$

**end while**

---

---

**Algorithm 2** Image-to-text Retrieval

---

**Input:** input image $\boldsymbol{x}$, relational scene descriptions $\{R_1, \ldots, R_n\}$, EBM $E_\theta(\cdot)$, Parser $P(\cdot)$, Encoder Enc$(\cdot)$, output energy list $\mathcal{O}$, caption prediction $\mathcal{C}$

$\mathcal{O} \leftarrow []$

$\triangleright$ *Generate image-caption matching energies iteratively*

**for** number of scene relations descriptions $i = 1$ to $n$ **do**

$\quad \triangleright$ *Parse a relational scene description:*

$\quad \{\boldsymbol{r}_1, \ldots \boldsymbol{r}_m\} \leftarrow P(R_i)$

$\quad \boldsymbol{e}_i = \sum_{j=1}^m E_\theta(\boldsymbol{x} \mid \text{Enc}(\boldsymbol{r}_j))$

$\quad \triangleright$ *output energy list $\mathcal{O}$*

$\quad \mathcal{O}.append(\boldsymbol{e}_i)$

**end for**

$\triangleright$ *Final output:*

$\mathcal{C} = \arg\min \mathcal{O}$

---

---

**Algorithm 3** Image Generation Algorithm

---

**Input:** Relational scene description $R$, number of data augmentation applications $N$, step size $\lambda$, number of steps $K$, data augmentation $D(\cdot)$, EBM $E_\theta(\cdot)$, Parser $P(\cdot)$, Encoder Enc$(\cdot)$

$\tilde{\boldsymbol{x}}^0 \sim \mathcal{U}$

$\triangleright$ *Parse a relational scene description:*

$\{\boldsymbol{r}_1, \ldots \boldsymbol{r}_m\} \leftarrow P(R)$

$\triangleright$ *Generate samples through $N$ iterative steps of data augmentation/Langevin dynamics:*

**for** sample step $n = 1$ to $N$ **do**

$\quad \triangleright$ *Apply data augmentation to samples:*

$\quad \tilde{\boldsymbol{x}}^0 = D(\tilde{\boldsymbol{x}}_i^0)$

$\quad \triangleright$ *Run $K$ steps of Langevin dynamics:*

$\quad$ **for** sample step $k = 1$ to $K$ **do**

$\quad\quad \tilde{\boldsymbol{x}}^k \leftarrow \tilde{\boldsymbol{x}}^{k-1} - \sum_{i=1}^n \nabla_{\boldsymbol{x}} E_\theta(\tilde{\boldsymbol{x}}^{k-1} \mid \text{Enc}(\boldsymbol{r}_i)) + \omega, \ \omega \sim \mathcal{N}(0, \sigma)$

$\quad$ **end for**

$\quad \triangleright$ *Iteratively refine samples:*

$\quad \tilde{\boldsymbol{x}}^0 = \tilde{\boldsymbol{x}}^k$

**end for**

$\triangleright$ *Final output:*

$\boldsymbol{x} = \tilde{\boldsymbol{x}}^0$

---

---

**Algorithm 4** Image Editing Algorithm

---

**Input:** input image $\tilde{\boldsymbol{x}}^0$, relational scene description $R$, number of data augmentation applications $N$, step size $\lambda$, number of steps $K$, data augmentation $D(\cdot)$, EBM $E_\theta(\cdot)$, Parser $P(\cdot)$, Encoder Enc$(\cdot)$

$\triangleright$ *Parse a relational scene description:*

$\{\boldsymbol{r}_1, \ldots \boldsymbol{r}_m\} \leftarrow P(R)$

$\triangleright$ *Generate samples through $N$ iterative steps of data augmentation/Langevin dynamics:*

**for** sample step $n = 1$ to $N$ **do**

$\quad \triangleright$ *Apply data augmentation to samples:*

$\quad \tilde{\boldsymbol{x}}^0 = D(\tilde{\boldsymbol{x}}_i^0)$

$\quad \triangleright$ *Run $K$ steps of Langevin dynamics:*

$\quad$ **for** sample step $k = 1$ to $K$ **do**

$\quad\quad \tilde{\boldsymbol{x}}^k \leftarrow \tilde{\boldsymbol{x}}^{k-1} - \sum_{i=1}^n \nabla_{\boldsymbol{x}} E_\theta(\tilde{\boldsymbol{x}}^{k-1} \mid \text{Enc}(\boldsymbol{r}_i)) + \omega, \ \omega \sim \mathcal{N}(0, \sigma)$

$\quad$ **end for**

$\quad \triangleright$ *Iteratively refine samples:*

$\quad \tilde{\boldsymbol{x}}^0 = \tilde{\boldsymbol{x}}^k$

**end for**

$\triangleright$ *Final output:*

$\boldsymbol{x} = \tilde{\boldsymbol{x}}^0$

---