# Dual Learning: Theoretical Study and an Algorithmic Extension (Supplementary Document)

## Appendix A. Derivation

We leverage the negative logarithmic probability to measure the differences between the original $x^{(2)}$ and the reconstructed one. Let $\mathcal{R}(x^{(2)})$ denote the event that after passing the loop $S_2 \rightarrow S_k \rightarrow S_1 \rightarrow S_2$, $x^{(2)}$ is reconstructed to $x^{(2)}$. We have that

$$\ln \Pr(\mathcal{R}(x^{(2)})) = \sum_{x^{(k)} \in S_k} \sum_{x^{(1)} \in S_1} \ln \Pr(x^{(2)}, x^{(1)}, x^{(k)}|$$

starting from $x^{(2)}$, applied by $\theta_{2k}, \theta_{k1}, \theta_{12}$ sequentially)

$$= \sum_{x^{(k)} \in S_k} \sum_{x^{(1)} \in S_1} \ln \left( \Pr(x^{(1)}, x^{(k)}|x^{(2)}; \theta_{2k}, \theta_{k1}) \cdot \Pr(x^{(2)}|x^{(1)}; \theta_{12}) \right) \tag{1}$$

$$\geq \sum_{x^{(k)} \in S_k} \sum_{x^{(1)} \in S_1} \Pr(x^{(1)}, x^{(k)}|x^{(2)}; \theta_{2k}, \theta_{k1}) \cdot \ln \Pr(x^{(2)}|x^{(1)}; \theta_{12})$$

$$= \sum_{x^{(k)} \in S_k} \sum_{x^{(1)} \in S_1} \Pr(x^{(k)}|x^{(2)}; \theta_{2k}, \theta_{k1}) \Pr(x^{(1)}|x^{(k)}, x^{(2)}; \theta_{2k}, \theta_{k1}) \cdot \ln \Pr(x^{(2)}|x^{(1)}; \theta_{12})$$

$$= \sum_{x^{(k)} \in S_k} \sum_{x^{(1)} \in S_1} \Pr(x^{(k)}|x^{(2)}; \theta_{2k}) \Pr(x^{(1)}|x^{(k)}; \theta_{k1}) \cdot \ln \Pr(x^{(2)}|x^{(1)}; \theta_{12}) \tag{2}$$

$$= \mathbb{E}_{x^{(k)} \sim \Pr(\cdot|x^{(2)}; \theta_{2k})} \mathbb{E}_{x^{(1)} \sim \Pr(\cdot|x^{(k)}; \theta_{k1})} \ln \Pr(x^{(2)}|x^{(1)}; \theta_{12}). \tag{3}$$

In Eqn.(1), the first $\Pr$ represents the jointly probability that $x^{(2)}$ can be translated into $x^{(k)}$ with $\theta_{2k}$, and the the obtained $x^{(k)}$ can be translated into $x^{(1)}$ with $\theta_{k1}$; the second $\Pr$ represents the probability that given $x^{(1)}$, it can be translated back to $x^{(2)}$ with $\theta_{12}$.