
Supplementary Material: Improved Multimodal Deep Learning with Variation of Information

Kihyuk Sohn, Wenling Shang and Honglak Lee
University of Michigan Ann Arbor, MI, USA
{kihyuks, shangw, honglak}@umich.edu

S1 Derivation of Equation (4)

The NLL objective function can be written as

$$\begin{aligned}
 2\mathcal{L}^{\text{NLL}}(\theta) &= -2\mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(X, Y)] \\
 &= -\mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(X|Y) + \log P_{\theta}(Y)] - \mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(Y|X) + \log P_{\theta}(X)] \\
 &= -\mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(X|Y) + \log P_{\theta}(Y|X)] - \mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(X) + \log P_{\theta}(Y)] \\
 &= \mathcal{L}^{\text{VI}}(\theta) - \mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(X)] - \mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(Y)] \tag{S1}
 \end{aligned}$$

$$\begin{aligned}
 &= \mathcal{L}^{\text{VI}}(\theta) + \underbrace{\mathbb{E}_{P_{\mathcal{D}}} \left[\log \frac{P_{\mathcal{D}}(X)}{P_{\theta}(X)} \right]}_{KL(P_{\mathcal{D}}(X) \| P_{\theta}(X))} + \underbrace{\mathbb{E}_{P_{\mathcal{D}}} \left[\log \frac{P_{\mathcal{D}}(Y)}{P_{\theta}(Y)} \right]}_{KL(P_{\mathcal{D}}(Y) \| P_{\theta}(Y))} \\
 &\quad \underbrace{-\mathbb{E}_{P_{\mathcal{D}}} [\log P_{\mathcal{D}}(X)] - \mathbb{E}_{P_{\mathcal{D}}} [\log P_{\mathcal{D}}(Y)]}_{C_1} \tag{S2}
 \end{aligned}$$

$$= \mathcal{L}^{\text{VI}}(\theta) + KL(P_{\mathcal{D}}(X) \| P_{\theta}(X)) + KL(P_{\mathcal{D}}(Y) \| P_{\theta}(Y)) + C_1 \tag{S3}$$

where Equation (S1) holds by the definition of $\mathcal{L}^{\text{VI}}(\theta)$. Note that C_1 is independent of θ . Similarly, we can rewrite the MinVI objective as

$$\mathcal{L}^{\text{VI}}(\theta) = -\mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(X|Y) + \log P_{\theta}(Y|X)] \tag{S4}$$

$$\begin{aligned}
 &= \mathbb{E}_{P_{\mathcal{D}}} \left[\log \frac{P_{\mathcal{D}}(X|Y)}{P_{\theta}(X|Y)} \right] + \mathbb{E}_{P_{\mathcal{D}}} \left[\log \frac{P_{\mathcal{D}}(Y|X)}{P_{\theta}(Y|X)} \right] \\
 &\quad \underbrace{-\mathbb{E}_{P_{\mathcal{D}}} [\log P_{\mathcal{D}}(X|Y)] - \mathbb{E}_{P_{\mathcal{D}}} [\log P_{\mathcal{D}}(Y|X)]}_{C_2} \tag{S5}
 \end{aligned}$$

where in Equation (S5), we have

$$\mathbb{E}_{P_{\mathcal{D}}} \left[\log \frac{P_{\mathcal{D}}(X|Y)}{P_{\theta}(X|Y)} \right] = \sum_y P_{\mathcal{D}}(y) \mathbb{E}_{P_{\mathcal{D}}(X|y)} \left[\log \frac{P_{\mathcal{D}}(X|y)}{P_{\theta}(X|y)} \right] \tag{S6}$$

$$= \mathbb{E}_{P_{\mathcal{D}}(Y)} [KL(P_{\mathcal{D}}(X|Y) \| P_{\theta}(X|Y))] \tag{S7}$$

Finally, we have

$$\begin{aligned}
 \mathcal{L}^{\text{VI}}(\theta) &= \mathbb{E}_{P_{\mathcal{D}}(X)} [KL(P_{\mathcal{D}}(Y|X) \| P_{\theta}(Y|X))] + \\
 &\quad \mathbb{E}_{P_{\mathcal{D}}(Y)} [KL(P_{\mathcal{D}}(X|Y) \| P_{\theta}(X|Y))] + C_2. \tag{S8}
 \end{aligned}$$

C_2 is independent of θ and by setting $C = C_1 + C_2$, we derive the Equation (4).

S2 Proof of Theorem 2.1

Proposition S2.1 ([1, 2]). *Let \mathcal{X} be a finite state space. Let irreducible transition matrices T_n and T converge to $\pi_n(X)$ and $\pi(X)$, respectively, where $\pi(X) = P_{\mathcal{D}}(X)$ is a data-generating distribution of X . If T_n converges to T in the induced matrix norm, which is denoted by $\|\cdot\|$, then $\pi_n(X)$ converges to $P_{\mathcal{D}}(X)$ in l^2 norm.*

Proof. Let $|\mathcal{X}|$ be the number of states. For simplicity, we denote $\pi = \pi(X)$ and $\pi_n = \pi_n(X)$. Since π is a stationary distribution of irreducible transition matrix T , π is uniquely defined and it satisfies the following:

$$T\pi = \pi, \mathbf{1}^\top \pi = 1. \quad (\text{S9})$$

Combining above two equations, we have

$$\underbrace{\begin{bmatrix} T_{1,1} - 1 & T_{1,2} & \cdots & T_{1,|\mathcal{X}|} \\ T_{2,1} & T_{2,2} - 1 & \cdots & T_{2,|\mathcal{X}|} \\ \vdots & \cdots & \cdots & \vdots \\ T_{|\mathcal{X}|-1,1} & \cdots & \cdots & T_{|\mathcal{X}|-1,|\mathcal{X}|-1} - 1 \\ 1 & 1 & \cdots & 1 \end{bmatrix}}_{=\tilde{T}} \pi = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \quad (\text{S10})$$

Since π exists and unique, \tilde{T} is invertible and we have

$$\pi = \tilde{T}^{-1} [0 \ 0 \ \cdots \ 1]^\top \quad (\text{S11})$$

and similarly,

$$\pi_n = \tilde{T}_n^{-1} [0 \ 0 \ \cdots \ 1]^\top \quad (\text{S12})$$

Since T_n (entrywise) converges to T , T_n^{-1} also converges to T^{-1} . Therefore, we conclude π_n converges to $\pi = P_{\mathcal{D}}(X)$. \square

Now, we provide a proof of Theorem 2.1.

Proof of Theorem 2.1. To prove the convergence of marginal distributions, it is sufficient to show the convergence of transition operators. Since $|\mathcal{X}|$ and $|\mathcal{Y}|$ are finite, for any $\epsilon > 0$, there exists N such that $\forall n \geq N$, with probability at least $1 - \epsilon$, $\forall x \in \mathcal{X}, \forall y \in \mathcal{Y}$,

$$|P_{\theta_n}(y|x) - P_{\mathcal{D}}(y|x)| < \epsilon, |P_{\theta_n}(x|y) - P_{\mathcal{D}}(x|y)| < \epsilon$$

The transition operators are defined as follows:

$$\begin{aligned} T_n^{\mathcal{Y}}(y[t]|y[t-1]) &= \sum_{x \in \mathcal{X}} P_{\theta_n}(y[t]|x) P_{\theta_n}(x|y[t-1]), \\ T^{\mathcal{Y}}(y[t]|y[t-1]) &= \sum_{x \in \mathcal{X}} P_{\mathcal{D}}(y[t]|x) P_{\mathcal{D}}(x|y[t-1]) \end{aligned}$$

where $P_{\theta_n}(x|y)$ and $P_{\theta_n}(y|x)$ are derived from the joint distribution $P_{\theta_n}(x, y)$ and similarly for data-generating distribution, $P_{\mathcal{D}}(x|y)$ and $P_{\mathcal{D}}(y|x)$ are derived from $P_{\mathcal{D}}(x, y)$. Then, for $n \geq N$, we have, for any $y_t, y_{t-1} \in \mathcal{Y}$, with probability at least $1 - \epsilon$,

$$\begin{aligned} & \left| T_n^{\mathcal{Y}}(y_t|y_{t-1}) - T^{\mathcal{Y}}(y_t|y_{t-1}) \right| \\ & \leq \left| \sum_{x \in \mathcal{X}} P_{\theta_n}(y_t|x) P_{\theta_n}(x|y_{t-1}) - P_{\mathcal{D}}(y_t|x) P_{\mathcal{D}}(x|y_{t-1}) \right| \\ & \leq |\mathcal{X}| \max_{x \in \mathcal{X}} \left| P_{\theta_n}(y_t|x) P_{\theta_n}(x|y_{t-1}) - P_{\mathcal{D}}(y_t|x) P_{\mathcal{D}}(x|y_{t-1}) \right| \\ & \leq |\mathcal{X}| (2\epsilon) \end{aligned} \quad (\text{S13})$$

As we assume finite sets \mathcal{X} and \mathcal{Y} , this proves the convergence (in probability) of transition operator $T_n^{\mathcal{Y}}$ to $T^{\mathcal{Y}}$. The same argument holds for the convergence of transition operator $T_n^{\mathcal{X}}$ to $T^{\mathcal{X}}$. With

Proposition S2.1, we proved the convergence of asymptotic marginal distribution $\pi_n(X)$ and $\pi_n(Y)$ to data-generating marginal distributions $P_{\mathcal{D}}(X)$ and $P_{\mathcal{D}}(Y)$, respectively.

Now, let's look at the joint probability distributions $P_{\theta_n}(x, y) = P_{\theta_n}(x|y)P_{\theta_n}(y)$ and similarly, $P_{\mathcal{D}}(x, y) = P_{\mathcal{D}}(x|y)P_{\mathcal{D}}(y)$. As we proved above, the following inequalities hold $\forall n \geq N'$:

$$\left| P_{\theta_n}(y) - P_{\mathcal{D}}(y) \right| < \epsilon, \quad \left| P_{\theta_n}(x|y) - P_{\mathcal{D}}(x|y) \right| < \epsilon \quad (\text{S14})$$

Therefore, using the similar argument in Equation (S13), we have

$$\left| P_{\theta_n}(x, y) - P_{\mathcal{D}}(x, y) \right| < 2\epsilon \quad (\text{S15})$$

and this completes the proof. \square

S3 Retrieval Task

We provide more results of retrieval task with multimodal queries on MIR-Flickr database.



sky, night, clouds, space



2007, beauty, hair, friend



studio, craft, room



puppy



skyline, indiana, 1855mm



portrait, me



portrait, explore, blackwhite, portraits



white, me



blue, night, city, explore, nyc, newyork, lights, newyorkcity, manhattan, ny, skyline, cityscape, twilight, skyscrapers



sunset, explore, sun



sky, night, stars



bw, portrait, blackandwhite, girl, nikon80



home, toys, interior, bed, books, decor



cute, puppy



night, city, river, dark, buildings, skyline



portrait, man, colours



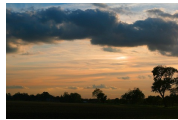
portrait



bw, selfportrait, me, layers



nyc, newyorkcity



sunset, trees, platinumphoto, silhouette



sky, night, mountains, stars, fab



blackandwhite, selfportrait, happy, mac, makeup



chair



puppy



night, reflection, longexposure, buildings, massachusetts, boston



selfportrait, me, 365days, 365, self



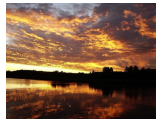
london, uk



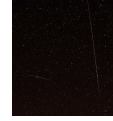
bw, blackandwhite, selfportrait, me, 365days, photoshop, self, face, head, myself, me



city



atardecer, sunset, abigfave, searchthebest, nubes, sol



2007



bw, portrait, photo



design, studio



dog



city, lights, buildings, fireworks, skyscrapers



selfportrait, 365days, 365



bw, halloween



de



night, lights, new, york, exposure, noche, long, pier



sunset



nikon, nature, sky, night, landscape, impressedbeauty, d300, dark, longexposure, colorado, stars



2007, may



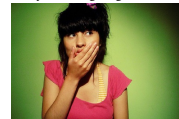
studio, craft



explore



nikon, night, d80, asia, skyline, hongkong, harbour



portrait, girl, woman, birthday



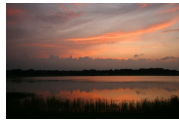
blackandwhite, milan



bw, self



sunset, chicago, tower, skyline, dusk, skyscraper



canon, naturesfinest, 30d

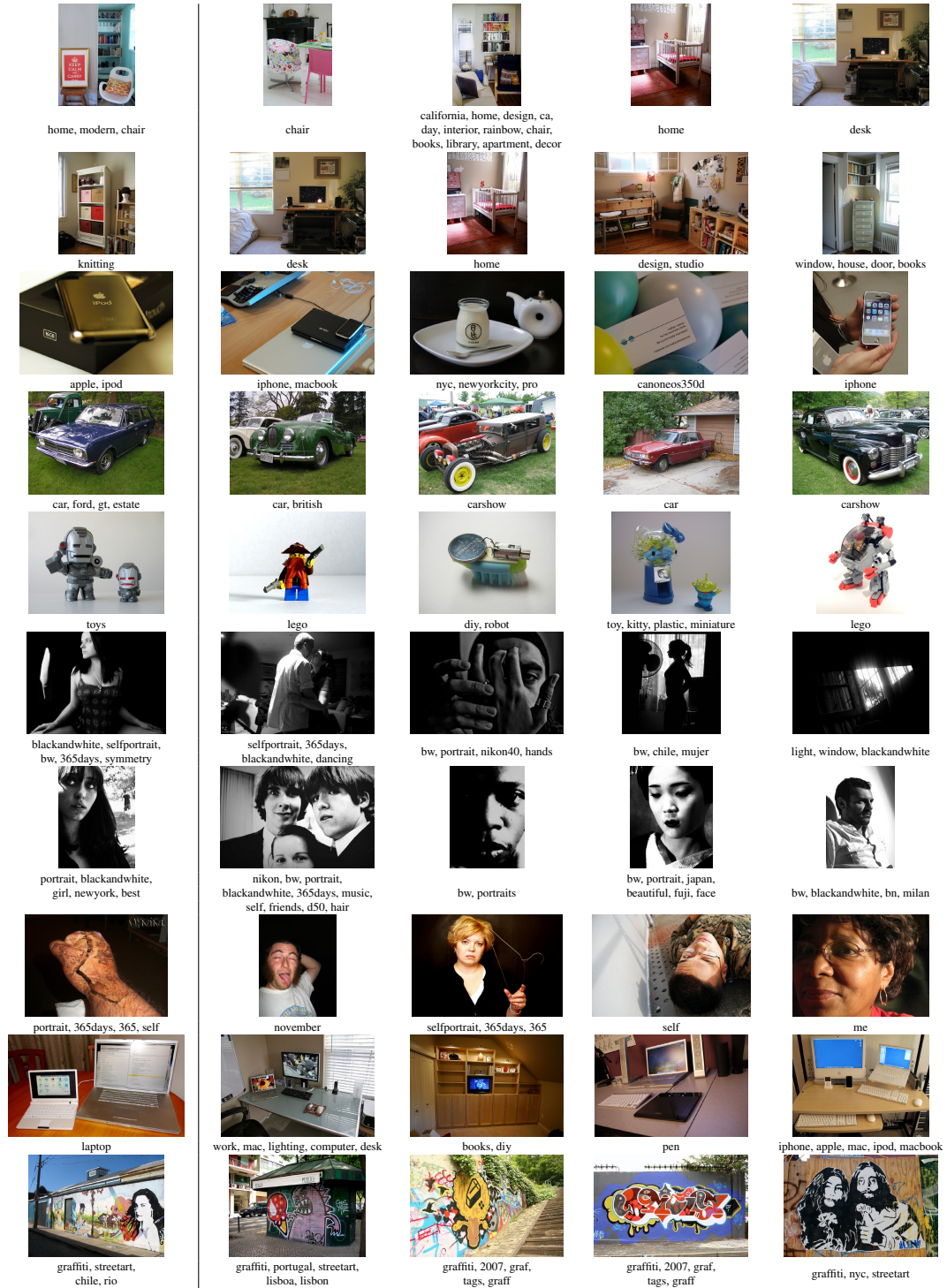


Figure S1: Retrieval results with multimodal queries on MIR-Flickr database. The leftmost image-text pairs are multimodal queries and those in the right side of the bar are retrieved samples with the highest similarities to the query.

References

- [1] Y. Bengio, L. Yao, G. Alain, and P. Vincent. Generalized denoising auto-encoders as generative models. In *NIPS*, 2013.
- [2] Y. Bengio, E. Thibodeau-Laufer, G. Alain, and J. Yosinski. Deep generative stochastic networks trainable by backprop. In *ICML*, 2014.