

GAN評価方法と アニメーション生成

5/10/2019

酒井一徳

➤ 評価方法

- IS
- m-IS
- FID
- SWD
- MMD
- SSIM

➤ 動画像生成

- vid2vid
- everybody dance now
- cartoonGAN

評価方法

- データセットに忠実かつ多様な生成。
- 人にわかりやすい特徴表現。
- 歪み、攻撃に敏感。
- 計算が少ない。

- ImageNet で学習済みの Inception Net を用いた指標。
 - 大きいほど良い。

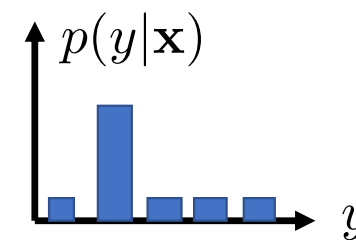
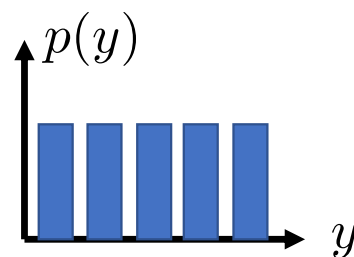
$$\exp(\mathbb{E}_{\mathbf{x}}[\text{KL}(p(y|\mathbf{x})||p(y))]) = \exp(H(y) - \mathbb{E}_{\mathbf{x}}[H(y|\mathbf{x})])$$

Inception Net
の出力

ラベル間の
多様性

質の高さ

理想状態



Inception Net

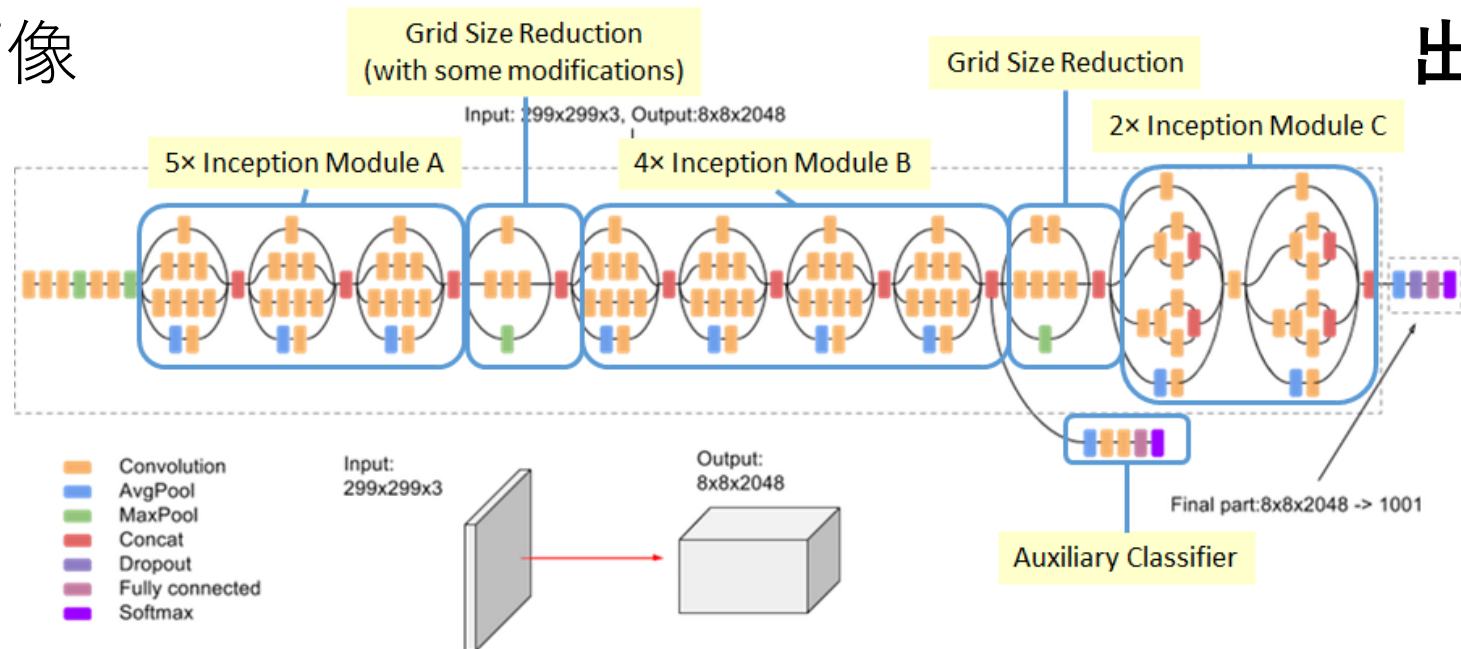
6

入力: 画像

出力: 確率値

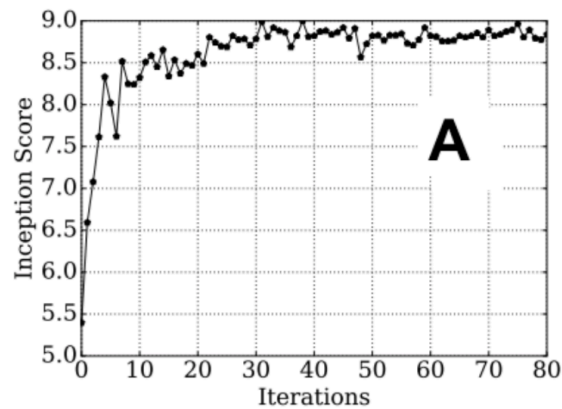
\mathbf{X}

$p(y|\mathbf{x})$

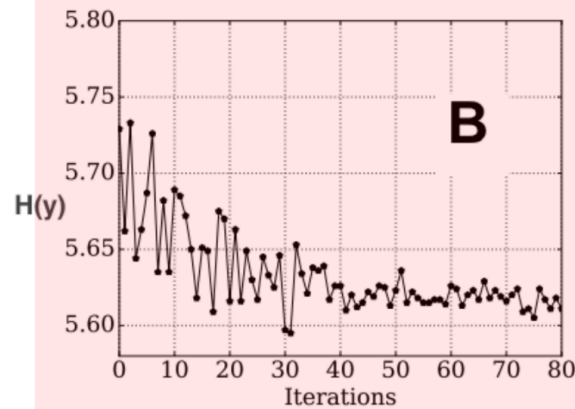


- 生成の多様性を測れない。
- mode collapse 発見できない。
- 画像の鮮明度が重要視される。
- 評価にNN使うのいいの？
 - 別ドメインのデータで学習したNNを評価に用いることの意義が不透明。
 - 画像以外のデータはこれだというNN(識別モデル)がない。

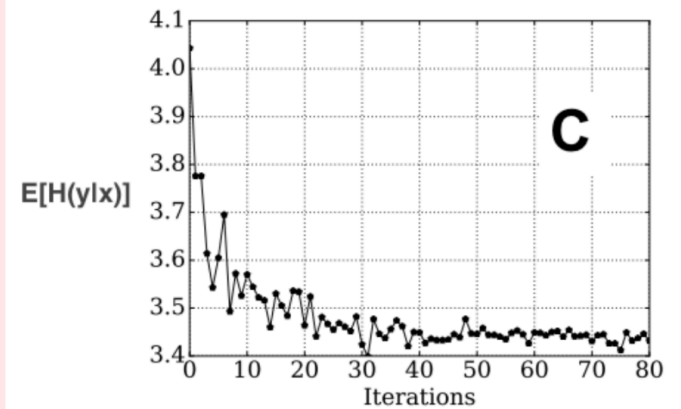
Inception Score



$H(y)$

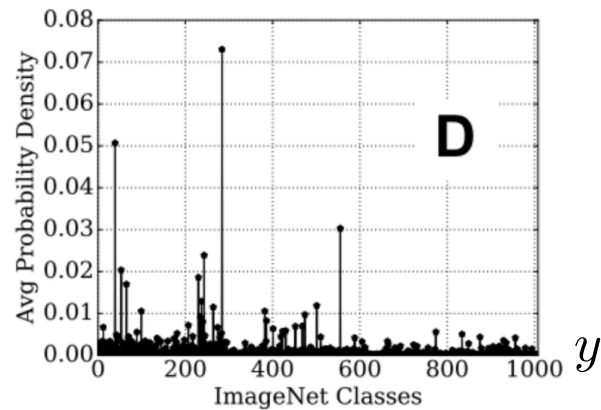


$H(y|\mathbf{x})$



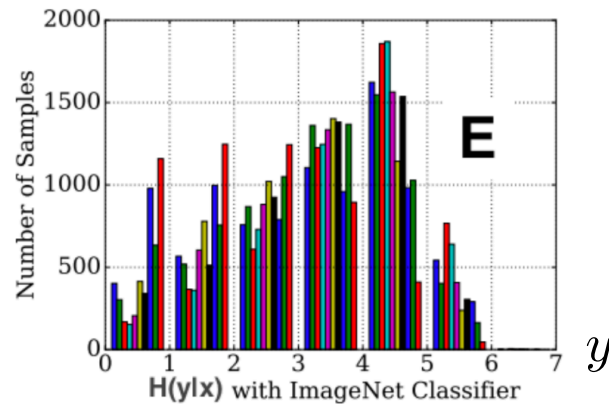
なんで下がるの

- 評価に使うモデルを何で学習するかに超依存。

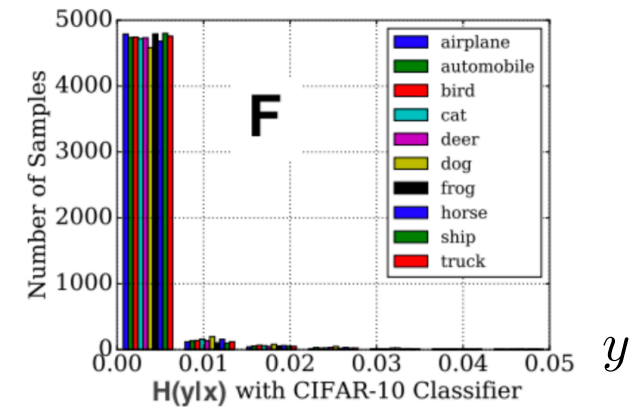


ImageNetで学習した
Inception Netに
CIFAR-10を入れたときの出力。

$$p(y)$$



ImageNetで学習した
Inception Netに
生成データを入れたもの



CIFAR-10で学習した
Inception Netに
生成データを入れたもの

- モード崩壊に鈍感。
 - 各クラスのサンプルを一つ記憶すればISは高くなる。

$$\exp(\mathbb{E}_{\mathbf{x}}[\text{KL}(p(y|\mathbf{x})\|p(y))]) = \exp(H(y) - \mathbb{E}_{\mathbf{x}}[H(y|\mathbf{x})])$$

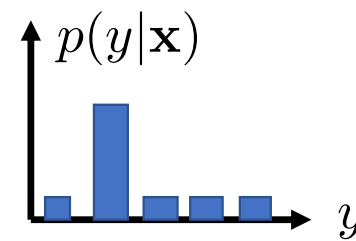
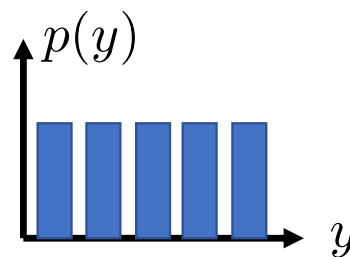
Inception Net
の出力

ラベル間
の多様性

質の高さ

多様性を鑑みる
項はない

理想状態



- クラス内の多様性を考慮したい。
 - 以下を最大化。

$$\exp \left(\mathbb{E}_{\mathbf{x}_i} \left[\mathbb{E}_{\mathbf{x}_j} \left[(\text{KL} (P(y|\mathbf{x}_i) \| P(y|\mathbf{x}_j))) \right] \right] \right)$$

同じクラスの違うサンプル

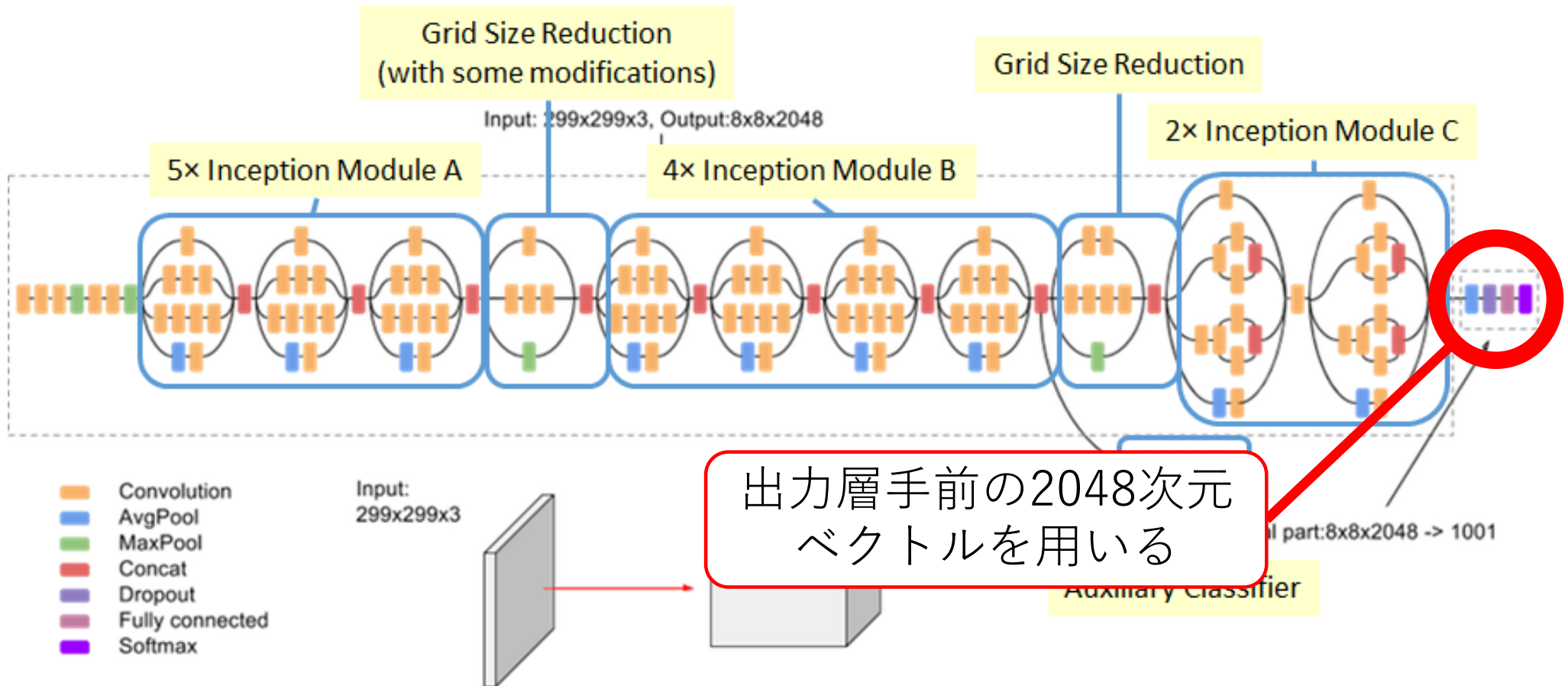
- クロスエントロピーから着想。

- ImageNet で学習済みの Inception Net を用いた指標。
- 実データと生成データの分布間の Fréchet (Wasserstein-2) 距離。
 - 小さいほど良い。

$$FID(r, g) = \|\boldsymbol{\mu}_r - \boldsymbol{\mu}_g\|_2^2 + \text{Tr} \left(\boldsymbol{\Sigma}_r + \boldsymbol{\Sigma}_g - 2(\boldsymbol{\Sigma}_r \boldsymbol{\Sigma}_g)^{\frac{1}{2}} \right)$$

where $\boldsymbol{\mu}_r = \frac{1}{N} \sum_{i=1}^N \mathbf{h}_r^{(i)}$ Inception Net の
中間層の出力(2048次元)

$$\boldsymbol{\Sigma}_r = \frac{1}{N} \sum_{i=1}^N (\mathbf{h}_r^{(i)} - \boldsymbol{\mu})(\mathbf{h}_r^{(i)} - \boldsymbol{\mu})^T$$



- 実データをどれだけ忠実に表現できてるか指標が欲しい。
- IS は実データを評価に使わない。

生成データ

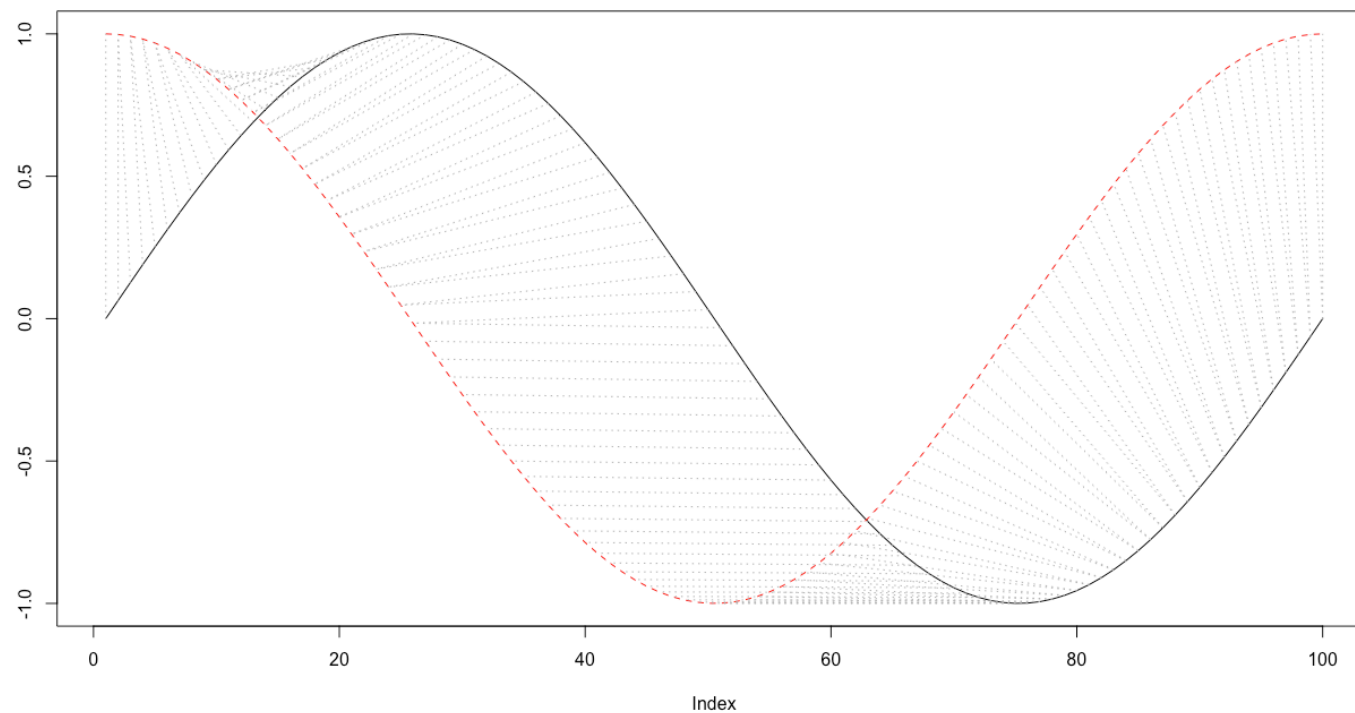
$$\exp(\mathbb{E}_{\mathbf{x}}[\text{KL}(p(y|\mathbf{x})\|p(y))]) = \exp(H(y) - \mathbb{E}_{\mathbf{x}}[H(y|\mathbf{x})])$$

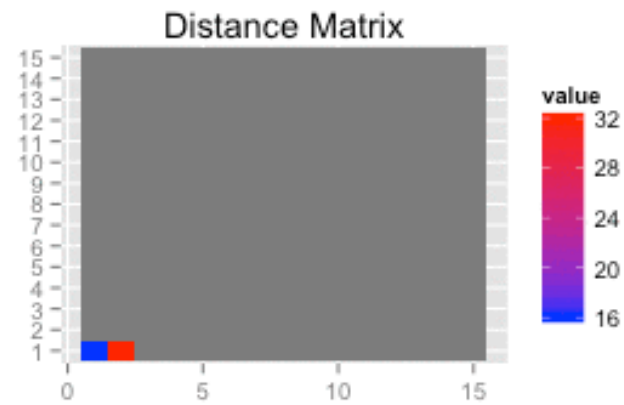
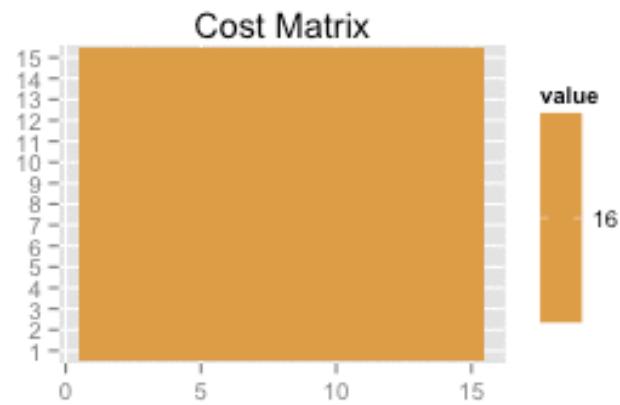
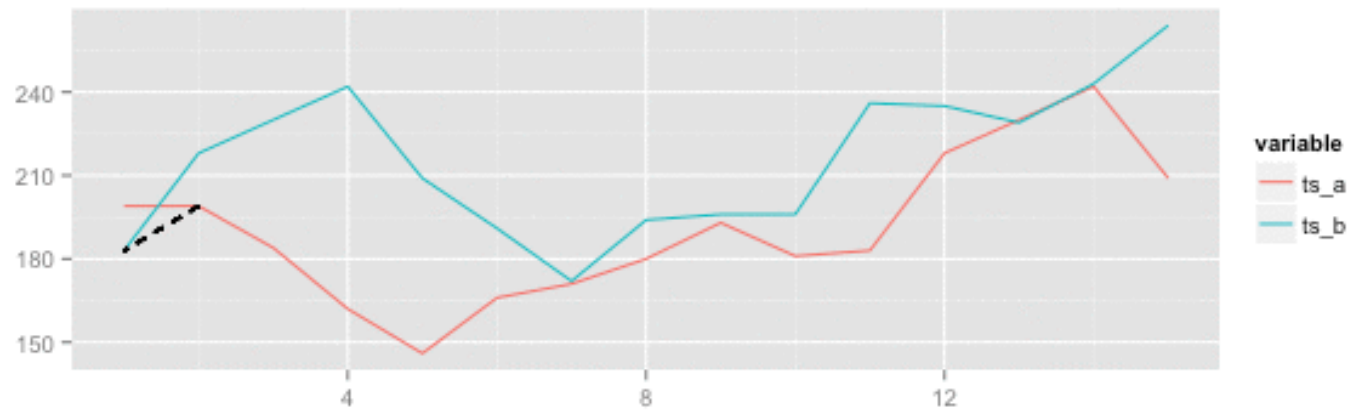
Inception Net
の出力

Fréchet (Wassestein-2) 距離とは？

15

- 曲線間に定義される距離。
- 離散では DTW (Dynamic Time Warping) として知られる。
 - 2系列の各点の距離を総当たりし、その系列間の最短パスがDTW。





<http://rtokei.tech/memo/memodynamic-time-warping/>

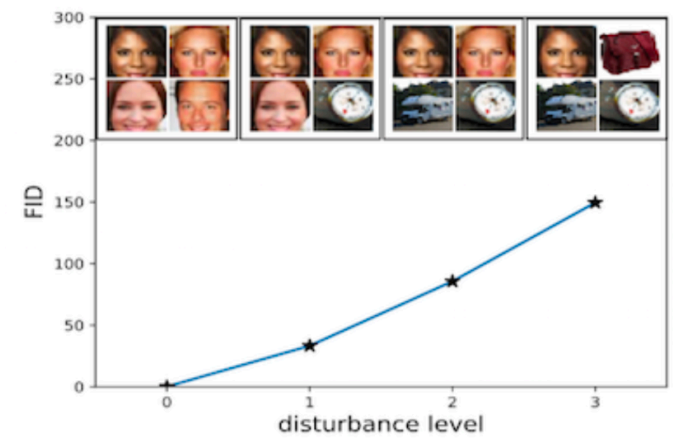
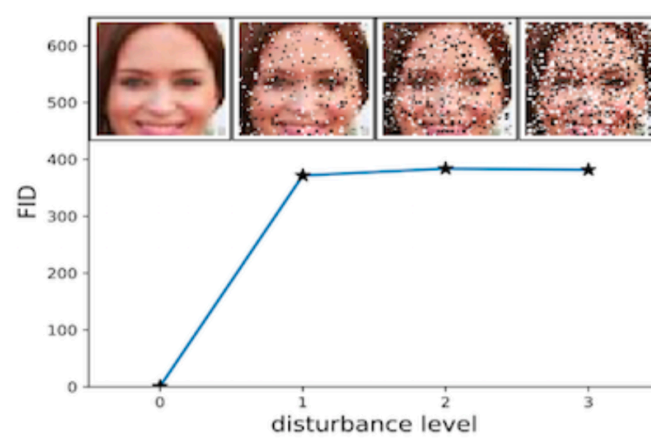
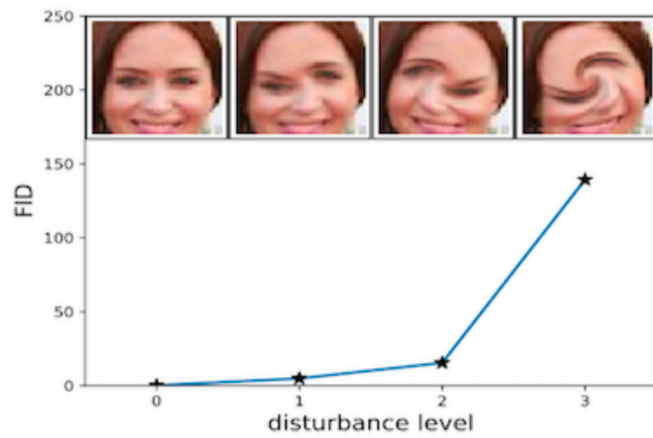
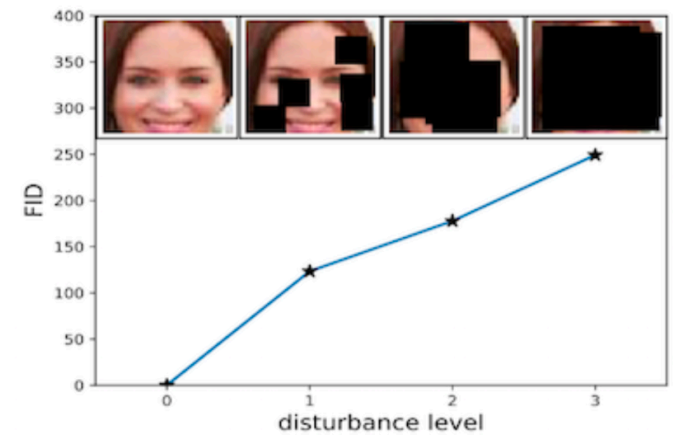
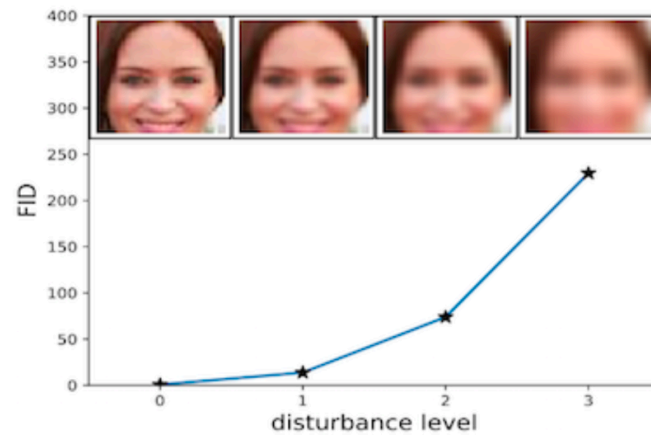
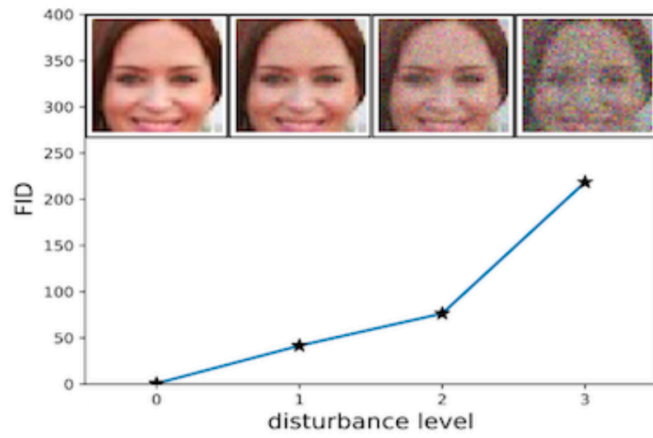
- Inception Net の出力がガウスに従うことを仮定している。
 - ガウス分布間の Fréchet 距離は以下で与えられる。

$$FID(r, g) = \|\boldsymbol{\mu}_r - \boldsymbol{\mu}_g\|_2^2 + \text{Tr} \left(\boldsymbol{\Sigma}_r + \boldsymbol{\Sigma}_g - 2(\boldsymbol{\Sigma}_r \boldsymbol{\Sigma}_g)^{\frac{1}{2}} \right)$$

where $\boldsymbol{\mu}_r = \frac{1}{N} \sum_{i=1}^N \mathbf{h}_r^{(i)}$ Inception Net の
中間層の出力(2048次元)

$$\boldsymbol{\Sigma}_r = \frac{1}{N} \sum_{i=1}^N (\mathbf{h}_r^{(i)} - \boldsymbol{\mu})(\mathbf{h}_r^{(i)} - \boldsymbol{\mu})^T$$

FID の挙動



- Inception Net の中間層がガウスに従うことの仮定。
- IS と同様に NN を使うことの意義は不透明。

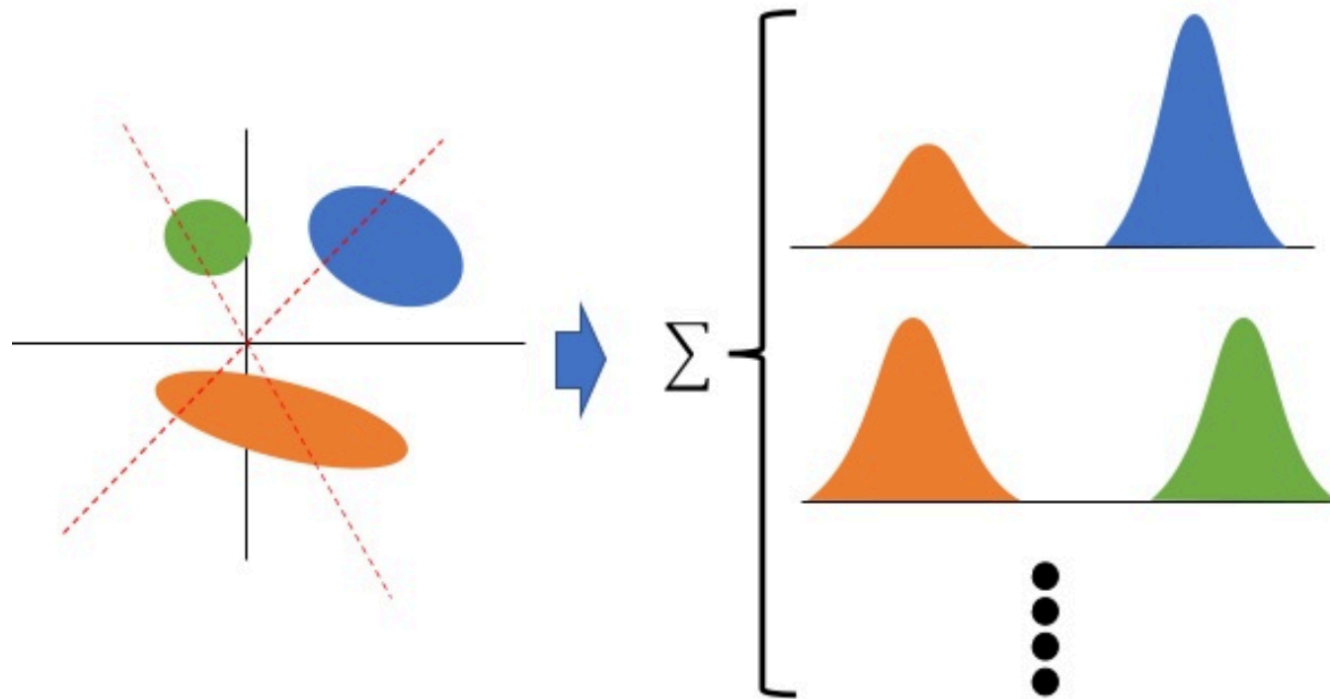
➤ Wasserstein Distance

- 高次元では計算量増大の上、近似になる。

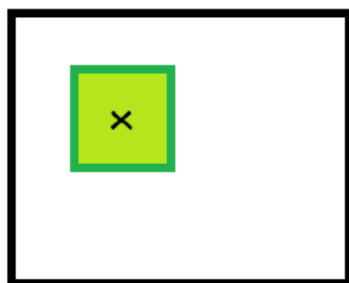
$$W(P_r, P_g) \propto \max_f \mathbb{E}_{\mathbf{x} \sim P_r} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim P_g} [f(\mathbf{x})]$$

$$\hat{W}(\mathbf{x}_{test}, \mathbf{x}_g) = \frac{1}{N} \sum_{i=1}^N \hat{f}(\mathbf{x}_{test}[i]) - \frac{1}{N} \sum_{i=1}^N \hat{f}(\mathbf{x}_g[i])$$

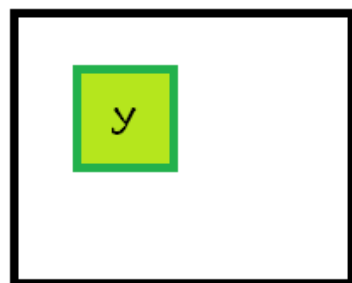
- 1次元ではclosed formで計算が楽。
➤ Sliceしよう。



- 画質の評価基準の一つ。
 - MSEなどは画像全体が少し違う場合と、局所的に大きく違う場合を見分けられない。



原画像



比較画像

画像の局所領域の平均と分散で計算。

輝度

コントラスト

構造

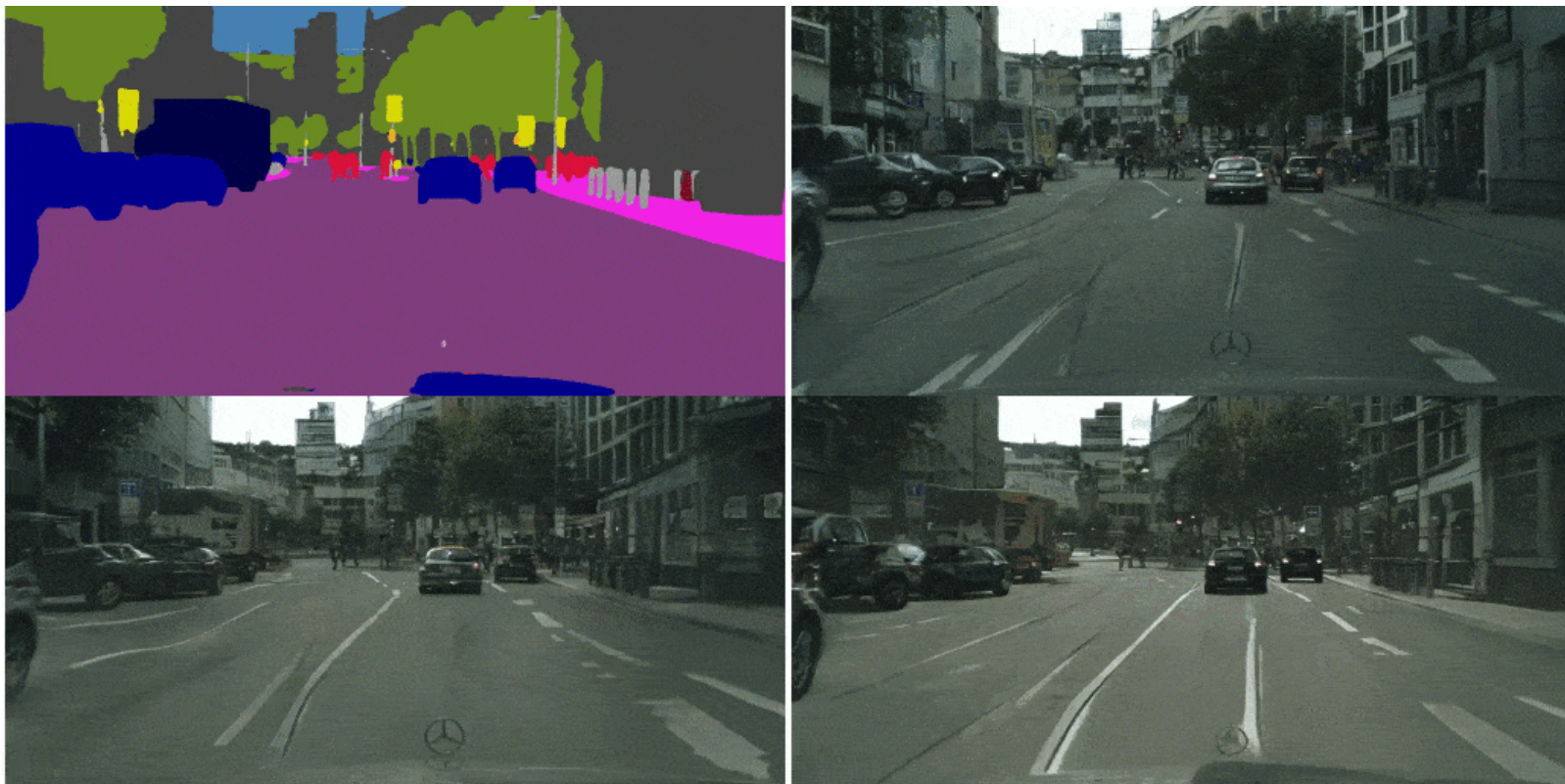
$$I(x, y) = \frac{2\mu_x\mu_y + C_1}{\mu_x^2 + \mu_y^2 + C_1} \quad C(x, y) = \frac{2\sigma_x\sigma_y + C_2}{\sigma_x^2 + \sigma_y^2 + C_2} \quad S(x, y) = \frac{\sigma_{xy} + C_3}{\sigma_x\sigma_y + C_3}$$

$$\text{SSIM}(x, y) = I(x, y)^\alpha C(x, y)^\beta S(x, y)^\gamma$$

- 画像ならISかFID。
 - でもNNを評価に使う意義はやっぱり不透明。
 - 評価するNNの学習まで気を使う必要がある。
- 音楽生成では人へのアンケートが基本っぽい。
 - 画像でもアンケートは基本してる。

動画像生成

- セマンティックセグメンテーション等がなされたソース系列から目的の系列を生成。(pix2pixの動画版)



ソース系列: $\mathbf{s}_1^T \equiv \{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_T\}$

ターゲット系列: $\mathbf{x}_1^T \equiv \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$

生成系列: $\tilde{\mathbf{x}}_1^T \equiv \{\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2, \dots, \tilde{\mathbf{x}}_T\}$

➤ 以下の分布の獲得が目標

$$p(\tilde{\mathbf{x}}_1^T | \mathbf{s}_1^T) = p(\mathbf{x}_1^T | \mathbf{s}_1^T)$$

➤ 目的関数

$$\max_D \min_G E_{(\mathbf{x}_1^T, \mathbf{s}_1^T)} [\log D(\mathbf{x}_1^T, \mathbf{s}_1^T)] + E_{\mathbf{s}_1^T} [\log (1 - D(G(\mathbf{s}_1^T), \mathbf{s}_1^T))]$$

- マルコフ性を仮定(論文ではL=2)。

$$p(\tilde{\mathbf{x}}_1^T | \mathbf{s}_1^T) = \prod_{t=1}^T p(\tilde{\mathbf{x}}_t | \tilde{\mathbf{x}}_{t-L}^{t-1}, \mathbf{s}_{t-L}^t)$$

過去L時刻までの生成系列
+ 現在も含んだL+1時刻のソース系列

- 以下のようなGenerator を定義

$$\tilde{\mathbf{x}}_t = F(\tilde{\mathbf{x}}_{t-L}^{t-1}, \mathbf{s}_{t-L}^t)$$

$$F(\tilde{\mathbf{x}}_{t-L}^{t-1}, \mathbf{s}_{t-L}^t) = (\mathbf{1} - \tilde{\mathbf{m}}_t) \odot \tilde{\mathbf{w}}_{t-1}(\tilde{\mathbf{x}}_{t-1}) + \tilde{\mathbf{m}}_t \odot \tilde{\mathbf{h}}_t$$

$$F(\tilde{\mathbf{x}}_{t-L}^{t-1}, \mathbf{s}_{t-L}^t) = (\mathbf{1} - \tilde{\mathbf{m}}_t) \odot \tilde{\mathbf{w}}_{t-1}(\tilde{\mathbf{x}}_{t-1}) + \tilde{\mathbf{m}}_t \odot \tilde{\mathbf{h}}_t$$

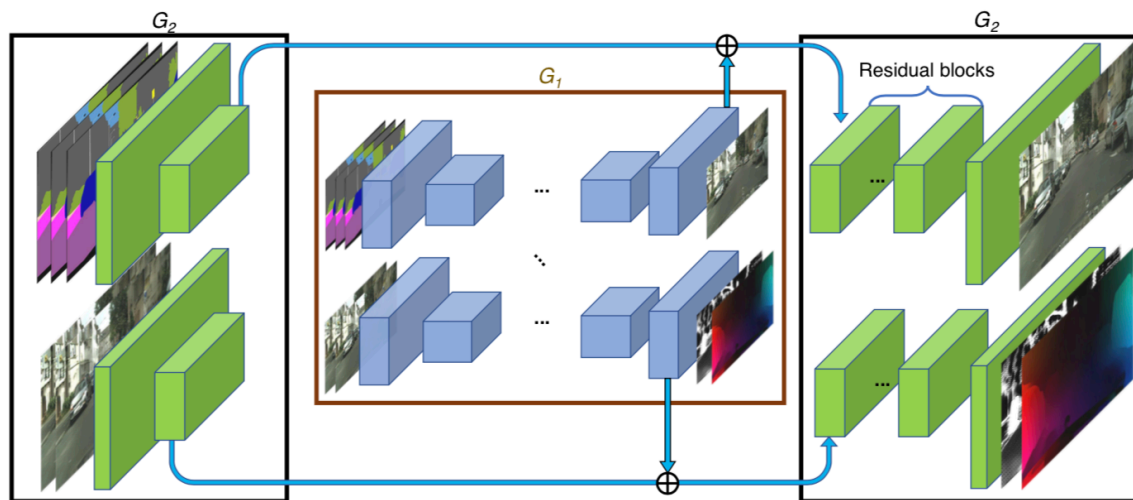
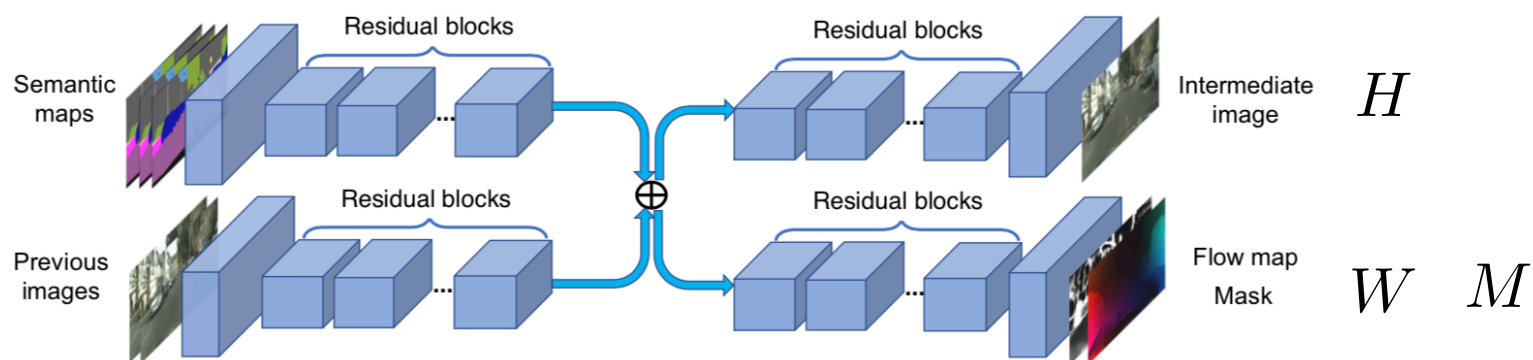
- $\tilde{\mathbf{w}}_{t-1} = W(\tilde{\mathbf{x}}_{t-L}^{t-1}, \mathbf{s}_{t-L}^t)$: estimated optical flow ($\mathbf{x}_{t-1} \rightarrow \mathbf{x}_t$)
- $\tilde{\mathbf{h}}_t = H(\tilde{\mathbf{x}}_{t-L}^{t-1}, \mathbf{s}_{t-L}^t)$: hallucinated image
- $\tilde{\mathbf{m}}_t = M(\tilde{\mathbf{x}}_{t-L}^{t-1}, \mathbf{s}_{t-L}^t)$: occlusion mask

全部NN(ResNet)で構成。
動画の冗長性を考慮する。

vid2vid の工夫 (Generator)

30

- 低解像度のGと高解像度のGを設計し、組み合わせる。



- 2種類の判別器を用いる。
 - 画像の判別器 D_I と、動画の判別器 D_V 。
- 最終的な目的関数:

$$\min_F \left(\max_{D_I} \mathcal{L}_I(F, D_I) + \max_{D_V} \mathcal{L}_V(F, D_V) \right) + \lambda_W \mathcal{L}_W(F)$$

(論文では $\lambda_W = 10$)

$\mathcal{L}_W(F)$: 生成系列と真の系列のVGG net で得られる中間層の
特徴量の L1 loss。(学習速度を上げる目的)

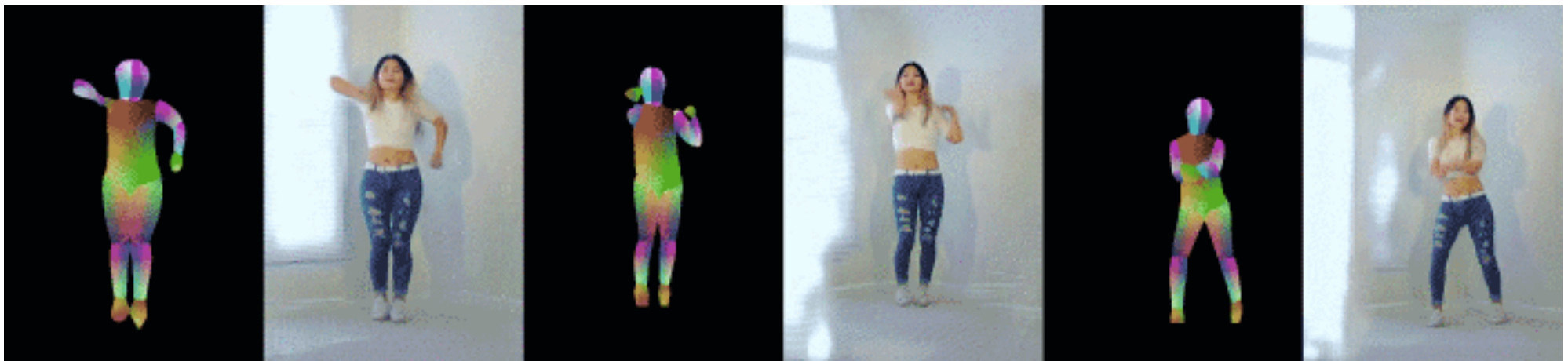
➤ labelの操作。



➤ スケッチをソースに。



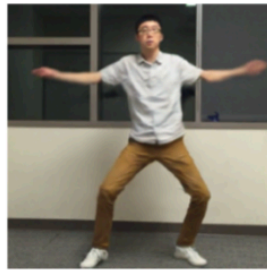
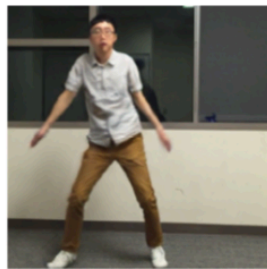
➤ poseをソースに。



➤ 評価にはFID、生成物の評価はアンケートで。

- 深さ方向の情報はないため、方向転換する車などの生成は難しい(セマンティックセグメンテーションの場合)。
- 色や質感の一貫性は保証できない(形状も上記の場合などでだめ)。

everybody dance now



Source Subject

Target Subject 1

Target Subject 2

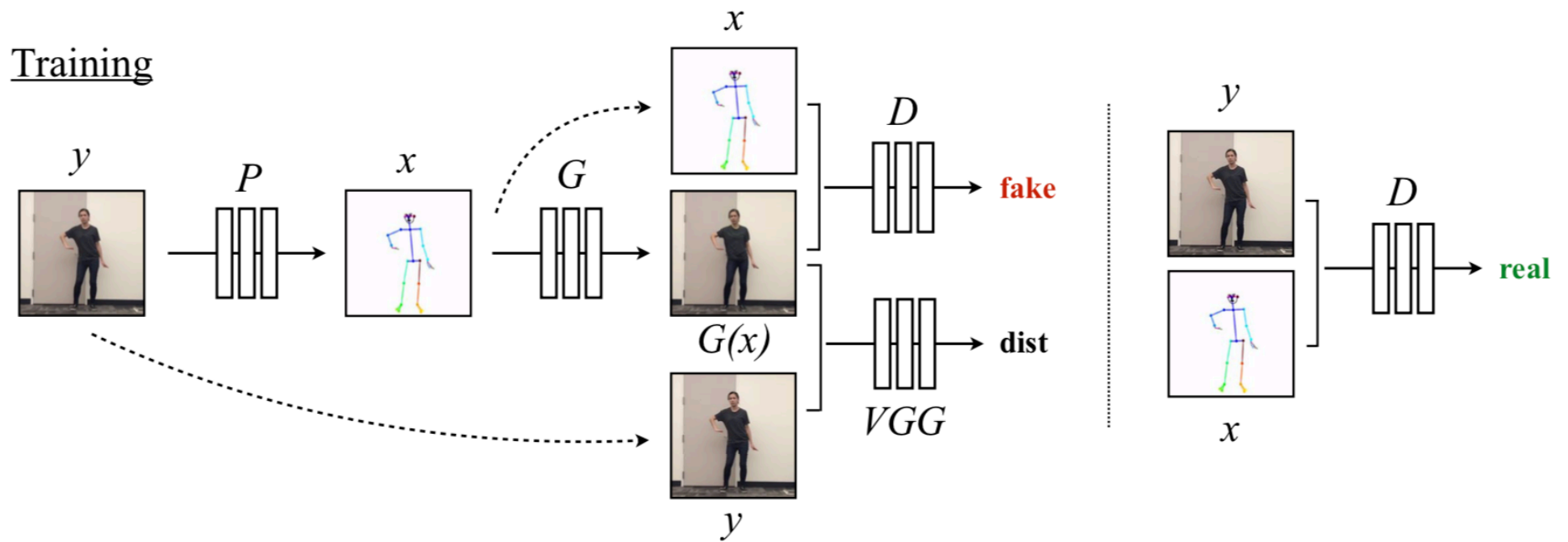


Source Subject

Target Subject 1

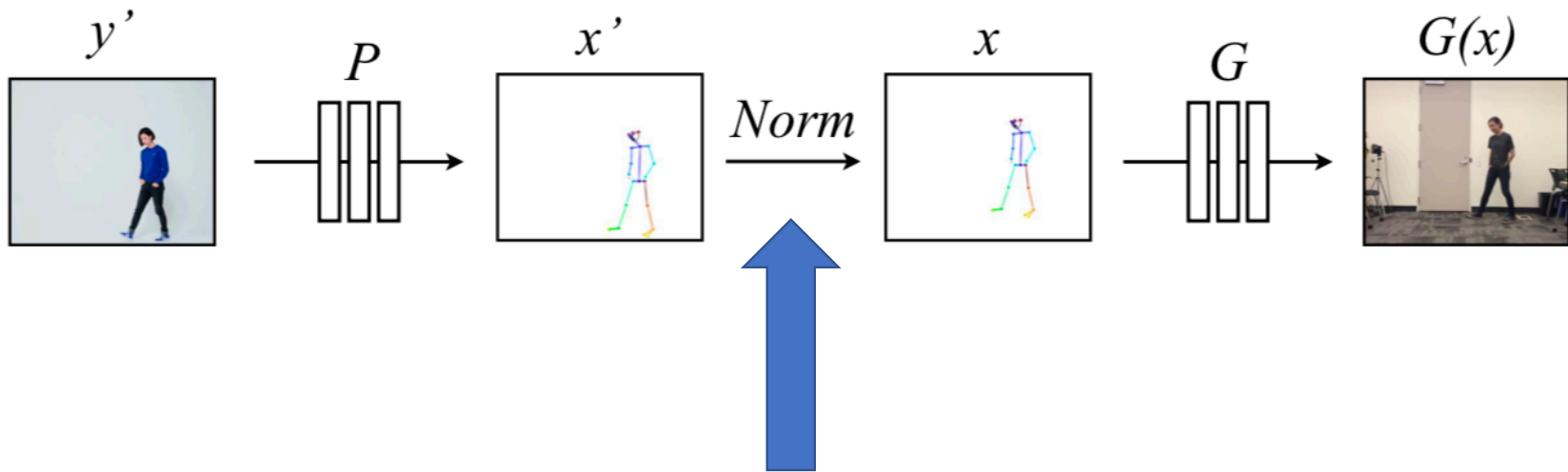
Target Subject 2

openpose と pix2pix



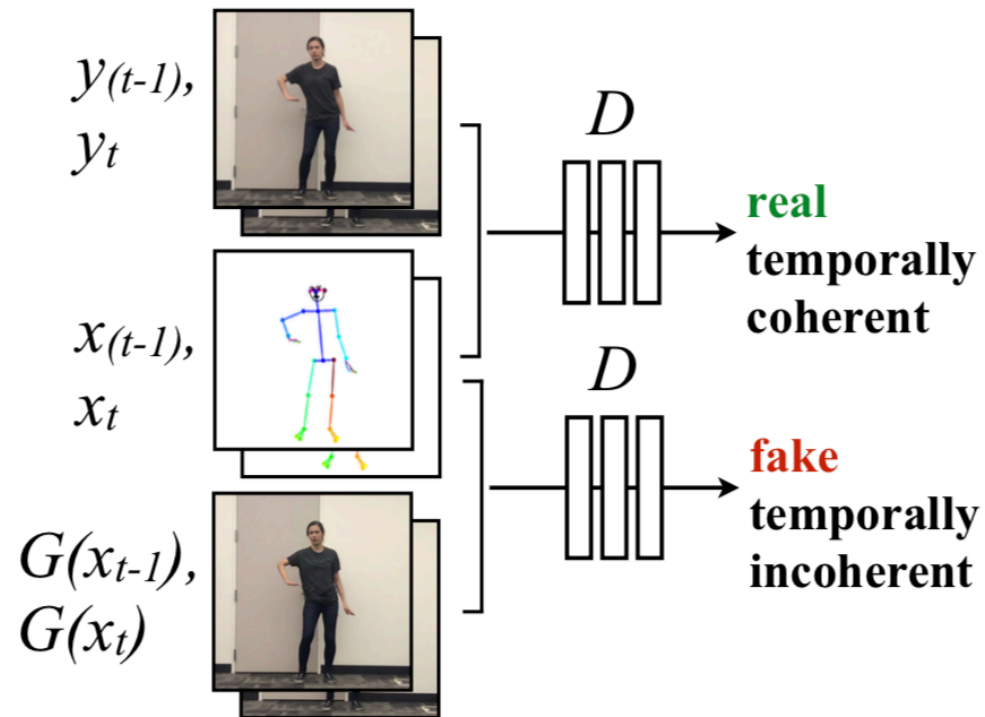
sourceと同じ姿勢のtargetを生成。

Transfer

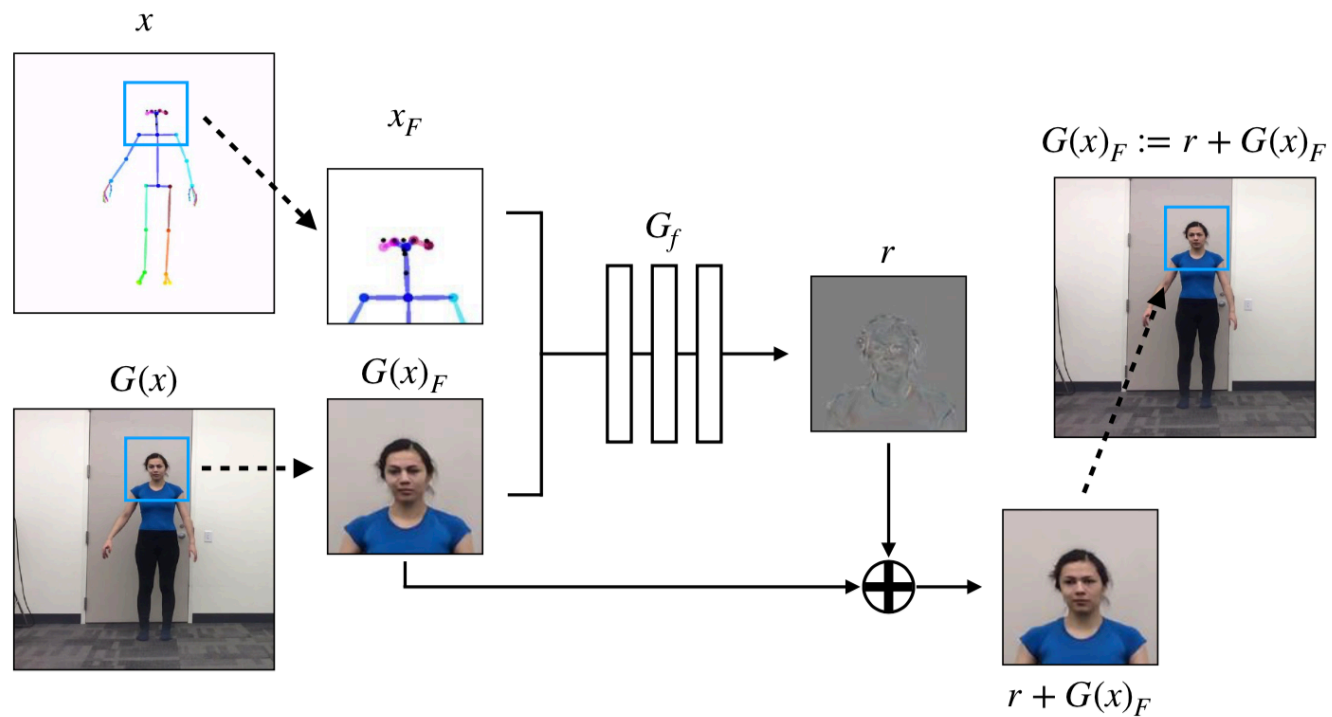


骨格や、床の位置、スケールを補正する。

- 連続するフレームが自然になるように。
 - pix2pixを連続する2枚組で行う。



- 顔部分の生成性能の向上のために。
 - 別のGANで差分を生成。







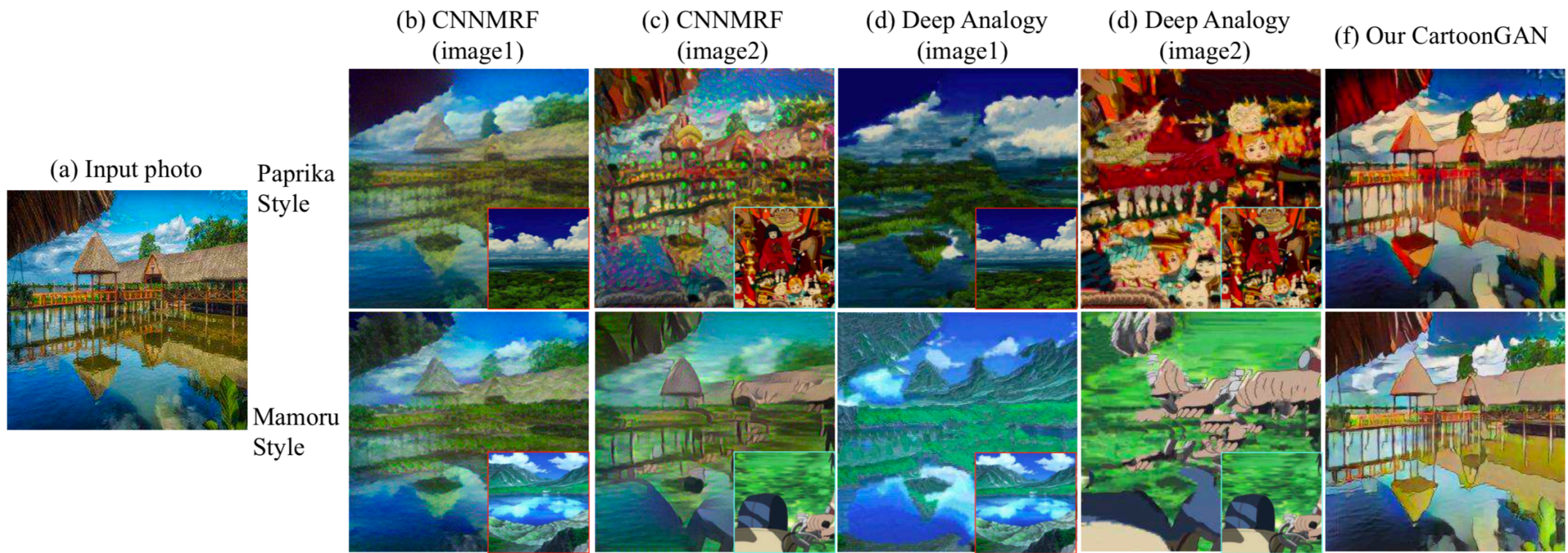
(a) Original scene



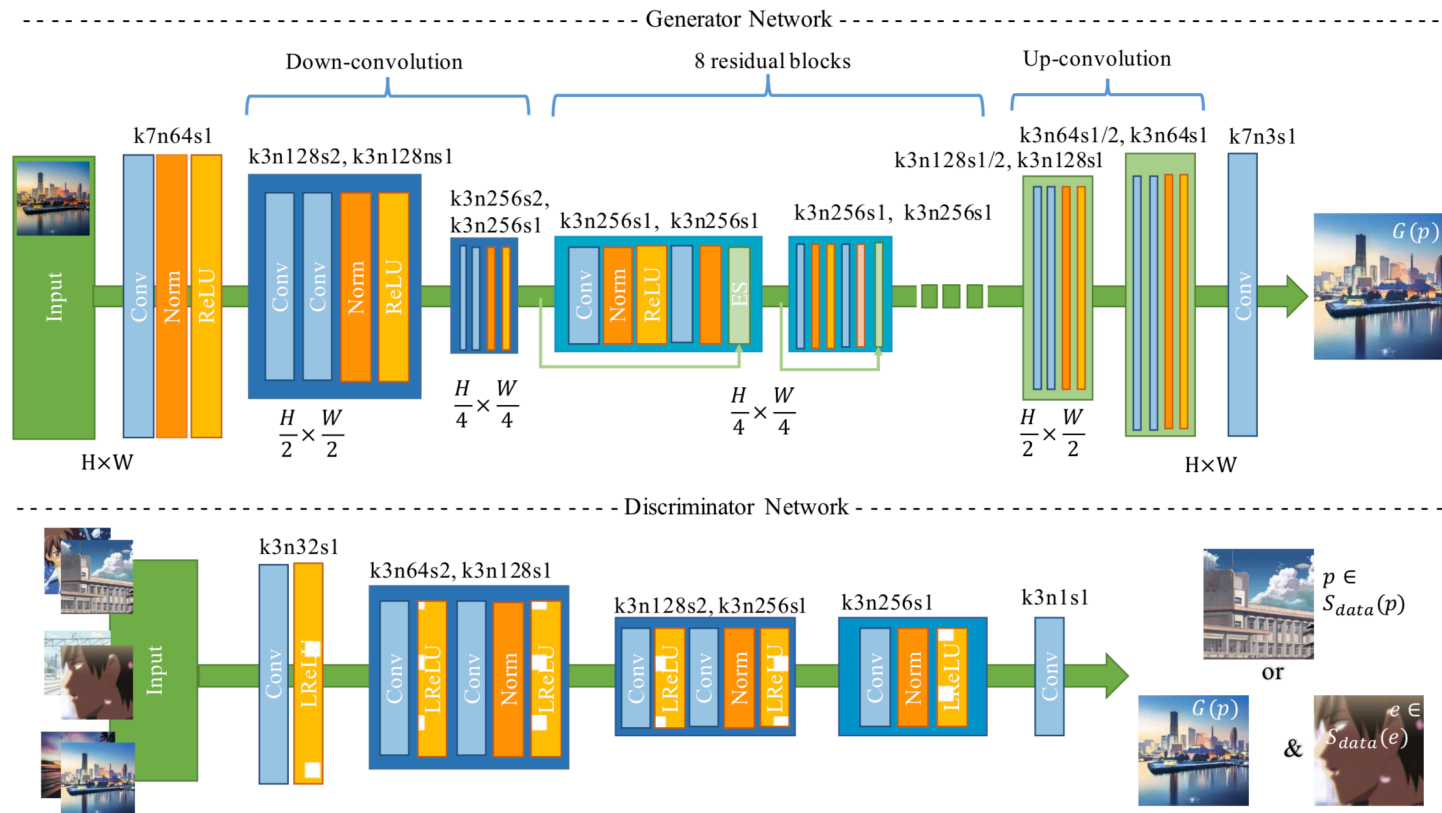
(b) Our result

実写のアニメ調化とは

- 鮮明な輪郭、なめらかな色調、少ないテクスチャ
- 高度な抽象化、単純化



- アニメに合わせた2つの損失関数
- 実写とアニメのデータセットを関連付けずに学習可能
- 収束を早めるための初期化



➤ Discriminatorへの入力は部分画像(パラメータ削減)

$$\mathcal{L}(G, D) = \mathcal{L}_{adv}(G, D) + \omega \mathcal{L}_{con}(G, D)$$

鮮明さ

変換後の内容の保持

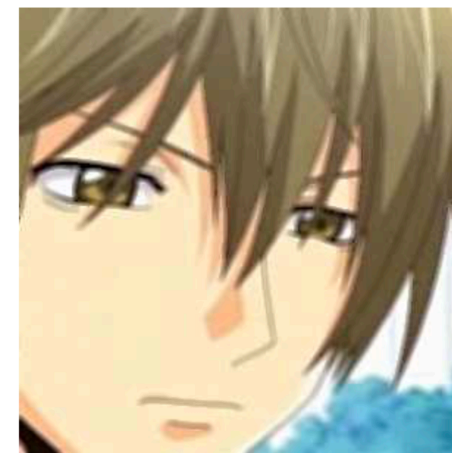
(論文では $w = 10$)

$$\begin{aligned}\mathcal{L}_{adv}(G, D) &= \mathbb{E}_{c_i \sim S_{data}(c)} [\log D(c_i)] && \text{(アニメ画像)} \\ &+ \mathbb{E}_{e_j \sim S_{data}(c)} [\log (1 - D(e_j))] && \text{(ぼけ画像)} \\ &+ \mathbb{E}_{p_k \sim S_{data}(p)} [\log (1 - D(G(p_k)))] && \text{(変換画像)}\end{aligned}$$

正例



負例



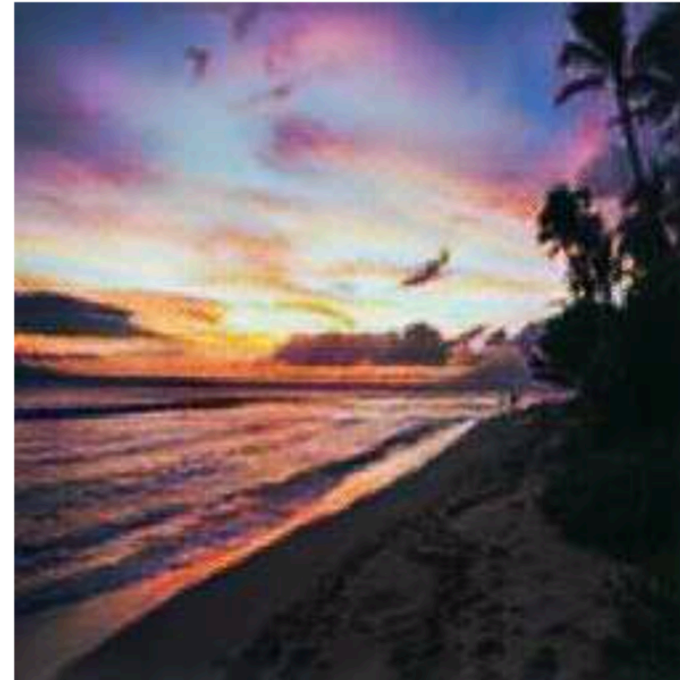
$$\mathcal{L}_{con}(G, D) = \mathbb{E}_{p_i \sim S_{data}(p)} [\|VGG_l(G(p_i)) - VGG_l(p_i)\|_1]$$

- VGGの中間層の特徴量のL1。
- L1使うのが著者曰く工夫ポイント。理由不明。

- 学習初期は $\mathcal{L}_{con}(G, D)$ のみで学習 (ただの復元誤差)



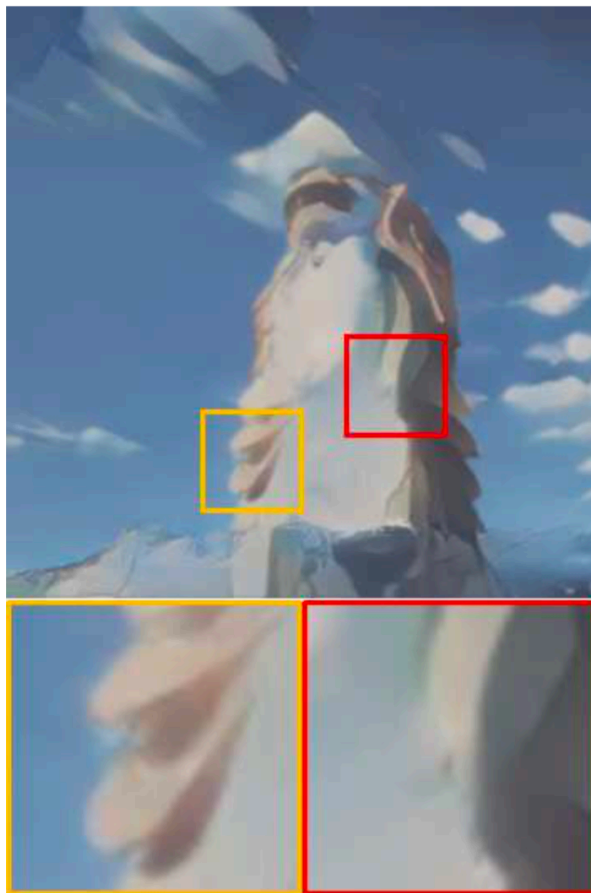
(a) Original photo



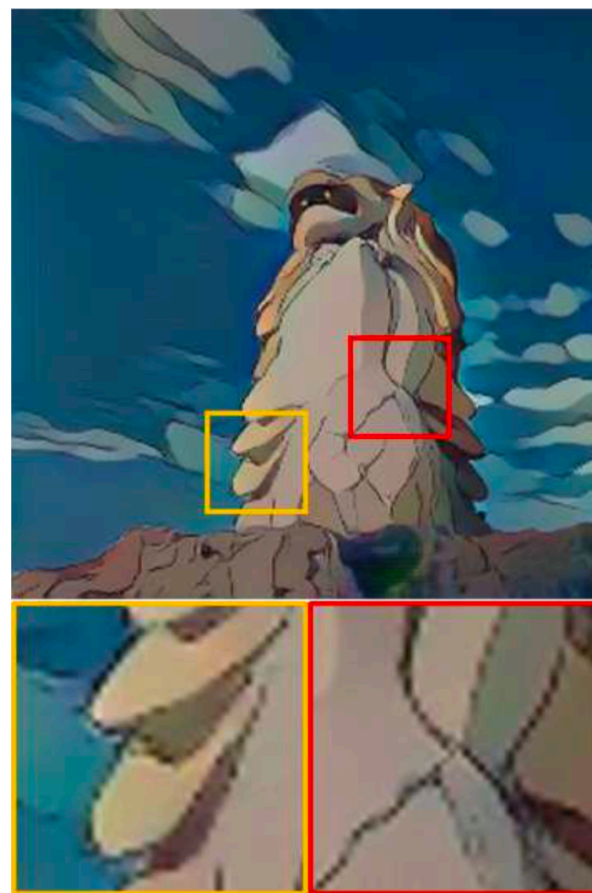
(b) Image after initialization



(a) NST



(b) CycleGAN



(c) CartoonGAN



※これは画像ごとに適用しただけ

<https://qiita.com/t-Asai/items/9e199db5fd0574f2ff8b>

本家 slideshare

<https://www.slideshare.net/hamadakoichi/anime-generation>

- Shmelkov, K., Schmid, C., & Alahari, K. (2018). How good is my GAN?. In *Proceedings of the European Conference on Computer Vision (ECCV)* (pp. 213-229).
- Wang, T. C., Liu, M. Y., Zhu, J. Y., Liu, G., Tao, A., Kautz, J., & Catanzaro, B. (2018). Video-to-video synthesis. *arXiv preprint arXiv:1808.06601*.
- Chan, C., Ginosar, S., Zhou, T., & Efros, A. A. (2019). Everybody dance now. In *Proceedings of the IEEE International Conference on Computer Vision* (pp. 5933-5942).
- Chen, Y., Lai, Y. K., & Liu, Y. J. (2018). Cartoongan: Generative adversarial networks for photo cartoonization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 9465-9474).