

Glow: Generative Flow with Invertible 1×1 Convolutions

Нугаманов Эдуард

Москва, 2019

Зачем нужны генеративные модели

- Позволяют моделировать данные
- Извлекают полезные признаки из тонн неразмеченных данных

Виды генеративных моделей

- Likelihood-based:
 - Авторегрессионные (PixelRNN, PixelCNN)
 - VAE
 - Flow-based (NICE, Real NVP, GLOW)
- GANs

Почему Flow-based?

- Точный вывод латентных переменных, расчет и оптимизация точного log-likelihood
- Вывод и сэмплирование эффективно параллелятся
- Получается хорошее скрытое представление
- Обратимые нейросети могут потреблять константное по отношению к глубине количество памяти [[1](#)].

Пример работы GLOW



Flows

- Найти такое преобразование $z = f(x)$ данных, что $P_Z(z)$ – простое
- f - обратимая и $f = f_1 \circ f_2 \circ \dots \circ f_K$ - (normalizing) flow

Обучение Flow-based models

- Оптимизируем:

$$\log p_X(x) = \log p_Z(f(x)) + \sum_{i=1}^K \log \left| \det \left(\frac{dz_i}{dz_{i-1}} \right) \right|$$

(здесь $z_0 = x, z_K = z, z_i = f_i(z_{i-1})$)

Причем $p_Z(z)$ - тоже может меняться на обучении

Требования к потоку

- $\det \left(\frac{dz_i}{dz_{i-1}} \right)$ должен легко вычисляться (треугольная матрица)
- $f^{-1}(z)$ тоже легко вычисляется

Тогда сэмплировать x тоже легко:

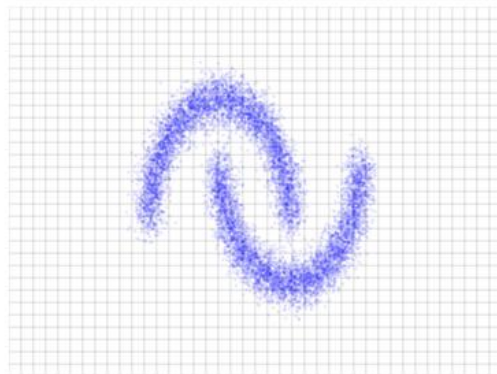
$$\begin{aligned} z &\sim p_Z(z) \\ x &= f^{-1}(z) \end{aligned}$$

Пример

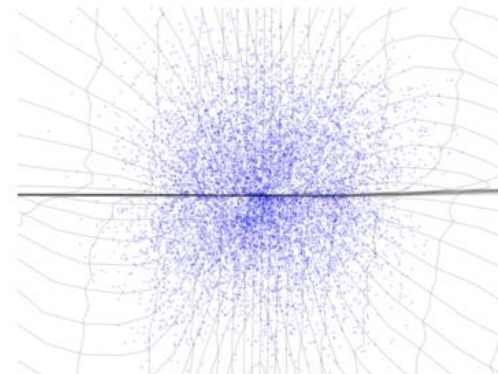
Inference

$$x \sim \hat{p}_X$$
$$z = f(x)$$

Data space \mathcal{X}

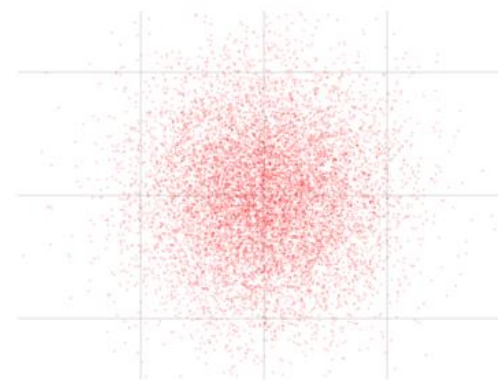
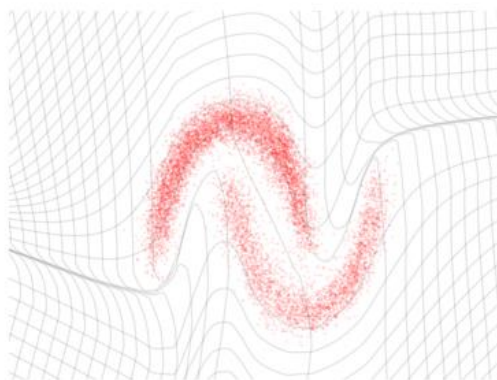


Latent space \mathcal{Z}

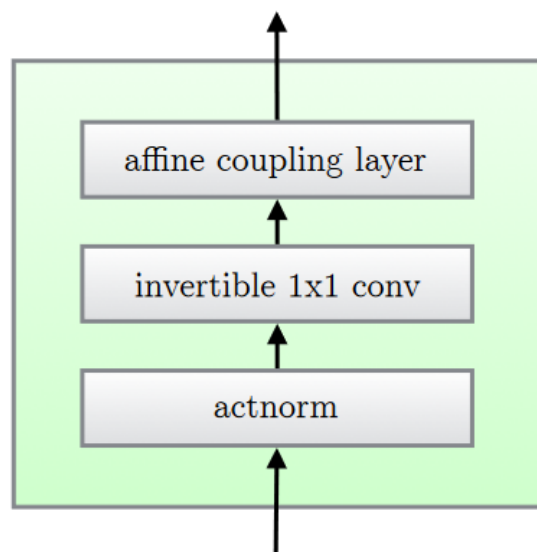


Generation

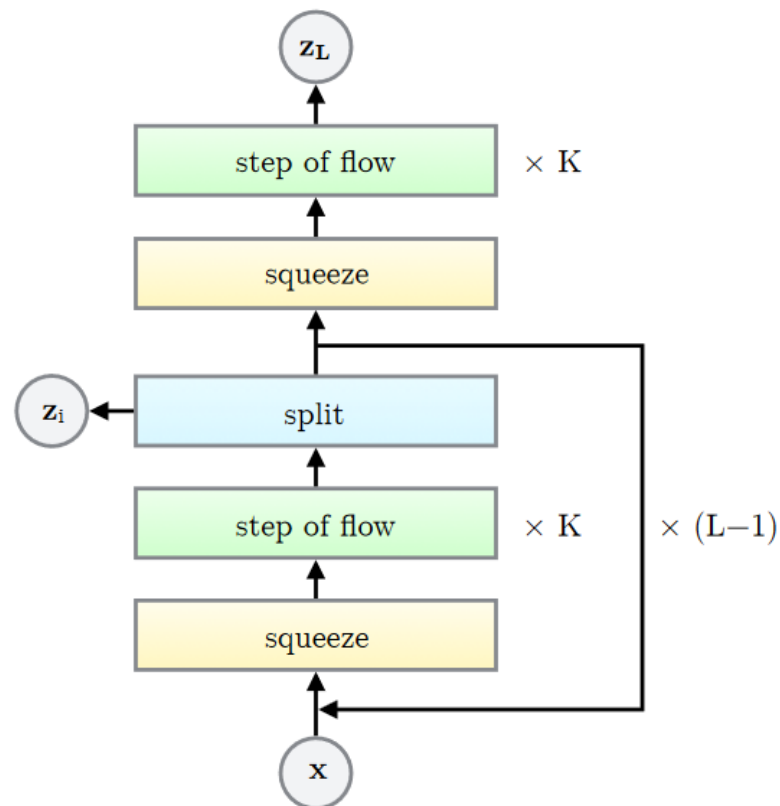
$$z \sim p_Z$$
$$x = f^{-1}(z)$$



Компоненты GLOW



(a) One step of our flow.



(b) Multi-scale architecture (Dinh et al., 2016).

Actnorm

- Обычно для лучшего обучения глубоких моделей используем Batch normalization
- Недостаток BN: на маленьких батчах много шума, а для больших изображений нужны маленькие батчи
- Activation normalization: scaling и bias параметры.
- Инициализация: так, что активации имеют 0 среднее и 1 дисперсию
- Дальше обучаем как обычные параметры

Coupling layers

- **General coupling layer:** $\dim(x) = D$; I_1, I_2 - разделение множества индексов $[1, D]$, $|I_1| = d$ (в GLOW делим каналы пополам).
- Тогда $y = (y_{I_1}, y_{I_2})$, где
$$\begin{aligned} y_{I_1} &= x_{I_1} \\ y_{I_2} &= g\left(x_{I_2}; m(x_{I_1})\right) \end{aligned}$$

Где g – *coupling law*, обратимое отображение относительно первого аргумента.

Coupling layers: смысл

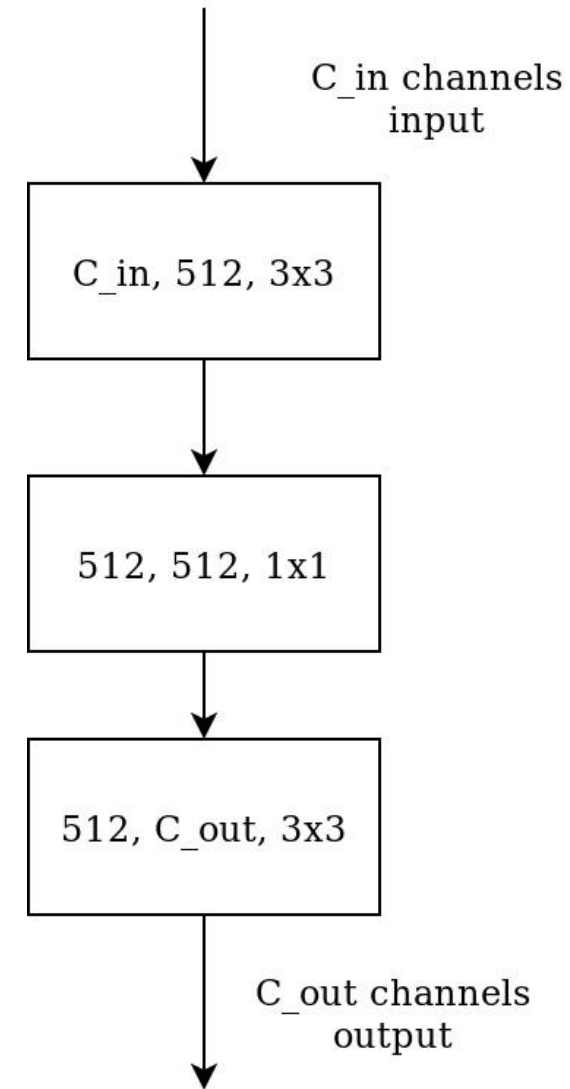
- Якобиан: $\frac{dy}{dx} = \begin{bmatrix} I_d & 0 \\ \frac{dy_{I_2}}{dx_{I_1}} & \frac{dy_{I_2}}{dx_{I_2}} \end{bmatrix} = \frac{dy_{I_2}}{dx_{I_2}}$, зависит только от g ,
- Значит, t может быть достаточно сложной (NN)

Coupling layers: виды

- Additive coupling layer[2]: $g(a; m(b)) = a + m(b);$
- Affine coupling layer[3, 4]: $g(a; m(b)) = a \odot \exp s(b) + t(b)$

Affine coupling layer: NN

- Нелинейное отображение, например, shallow CNN, как в ResNet



Комбинирование coupling layers

Проблема: некоторые компоненты x остаются неизменными

Решения:

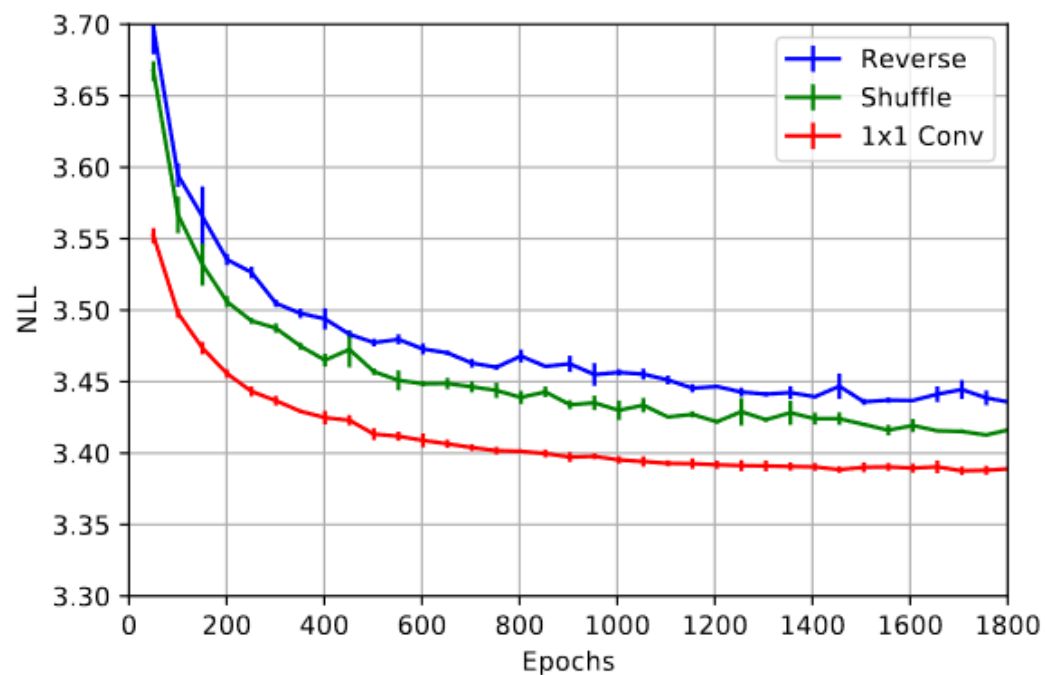
- Менять местами I_1 и I_2 (эквивалентно обратному порядку каналов);
- Случайная (фиксированная) перестановка каналов
- 1x1 invertible convolutions

Invertible 1x1 convs

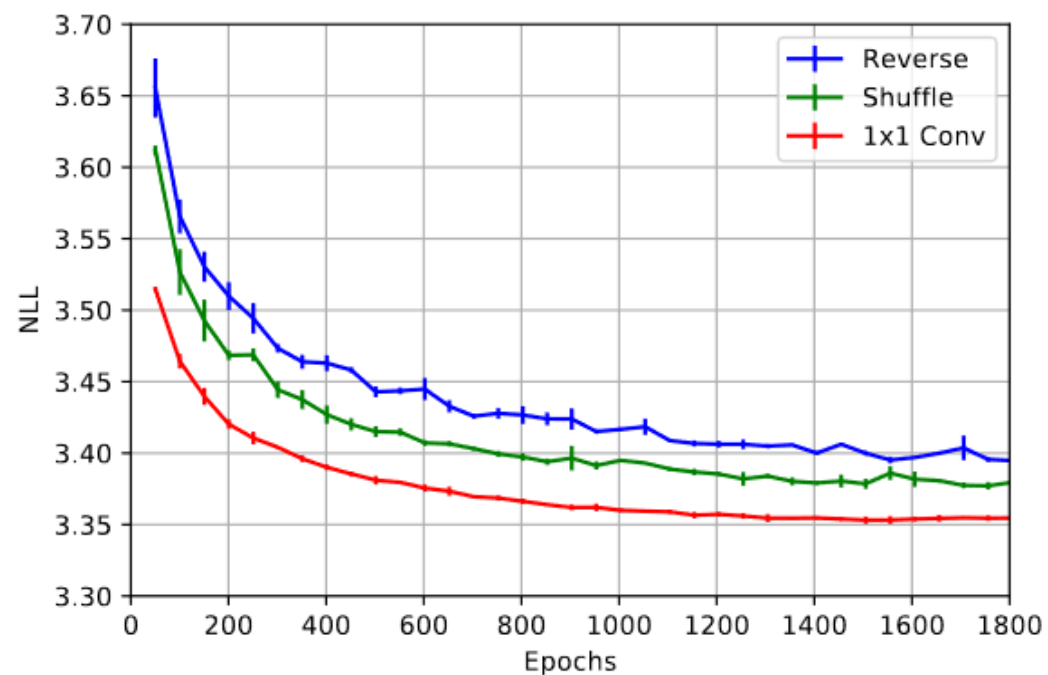
- Обобщение перестановок, когда число $n_{in\ channels} = n_{out\ channels}$
- Пусть x – тензор размера $h \times w \times c$, W – матрица весов свертки размера $c \times c$

$$\log \left| \det \left(\frac{d \text{conv2D}(x; W)}{dx} \right) \right| = h \cdot w \cdot \log |\det(W)|$$

Сравнение подходов



(a) Additive coupling.



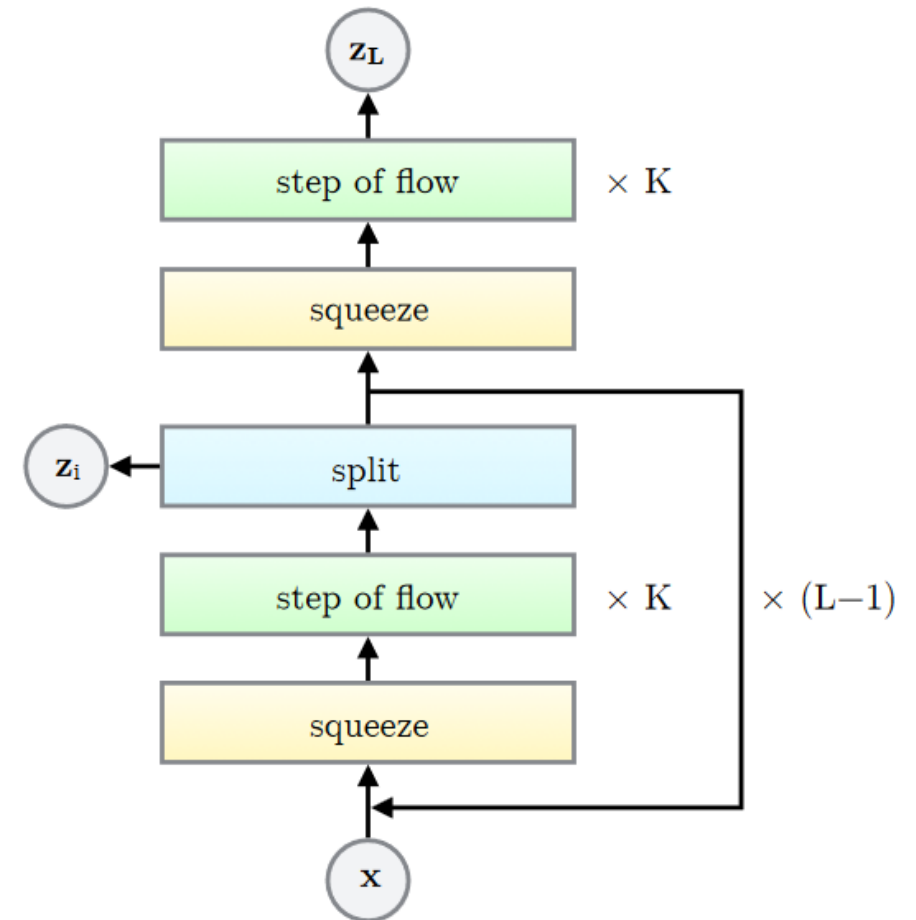
(b) Affine coupling.

Суммарно о каждой операции

Description	Function	Reverse Function	Log-determinant
Actnorm. See Section 3.1.	$\forall i, j : \mathbf{y}_{i,j} = \mathbf{s} \odot \mathbf{x}_{i,j} + \mathbf{b}$	$\forall i, j : \mathbf{x}_{i,j} = (\mathbf{y}_{i,j} - \mathbf{b})/\mathbf{s}$	$h \cdot w \cdot \text{sum}(\log \mathbf{s})$
Invertible 1×1 convolution. $\mathbf{W} : [c \times c]$. See Section 3.2.	$\forall i, j : \mathbf{y}_{i,j} = \mathbf{W}\mathbf{x}_{i,j}$	$\forall i, j : \mathbf{x}_{i,j} = \mathbf{W}^{-1}\mathbf{y}_{i,j}$	$h \cdot w \cdot \log \det(\mathbf{W}) $ or $h \cdot w \cdot \text{sum}(\log \mathbf{s})$ (see eq. 10))
Affine coupling layer. See Section 3.3 and (Dinh et al., 2014)	$\mathbf{x}_a, \mathbf{x}_b = \text{split}(\mathbf{x})$ $(\log \mathbf{s}, \mathbf{t}) = \text{NN}(\mathbf{x}_b)$ $\mathbf{s} = \exp(\log \mathbf{s})$ $\mathbf{y}_a = \mathbf{s} \odot \mathbf{x}_a + \mathbf{t}$ $\mathbf{y}_b = \mathbf{x}_b$ $\mathbf{y} = \text{concat}(\mathbf{y}_a, \mathbf{y}_b)$	$\mathbf{y}_a, \mathbf{y}_b = \text{split}(\mathbf{y})$ $(\log \mathbf{s}, \mathbf{t}) = \text{NN}(\mathbf{y}_b)$ $\mathbf{s} = \exp(\log \mathbf{s})$ $\mathbf{x}_a = (\mathbf{y}_a - \mathbf{t})/\mathbf{s}$ $\mathbf{x}_b = \mathbf{y}_b$ $\mathbf{x} = \text{concat}(\mathbf{x}_a, \mathbf{x}_b)$	$\text{sum}(\log(\mathbf{s}))$

Multi-scale architecture[[3](#)]

- **Squeeze:** разделить изображение на блоки $2 \times 2 \times c$, reshape to $1 \times 1 \times 4c$, обмен пространственной информацией на глубину каналов
- **Factor out:** Для экономии памяти и времени после каждого шага потока отделяем часть каналов, как «готовые»
- $z = (z_1, z_2, \dots, z_L)$



Результаты (bits per dim)

Table 2: Best results in bits per dimension of our model compared to RealNVP.

Model	CIFAR-10	ImageNet 32x32	ImageNet 64x64	LSUN (bedroom)	LSUN (tower)	LSUN (church outdoor)
RealNVP	3.49	4.28	3.98	2.72	2.81	3.08
Glow	3.35	4.09	3.81	2.38	2.46	2.67

Линейная интерполяция

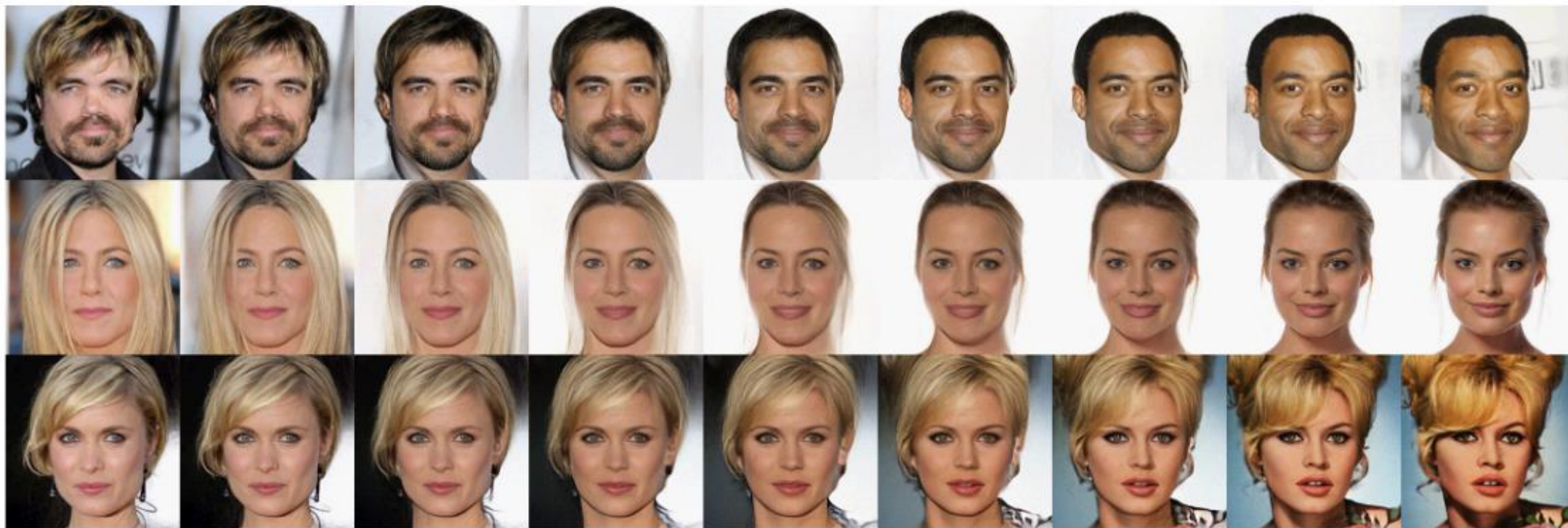


Figure 5: Linear interpolation in latent space between real images

Изменение атрибутов

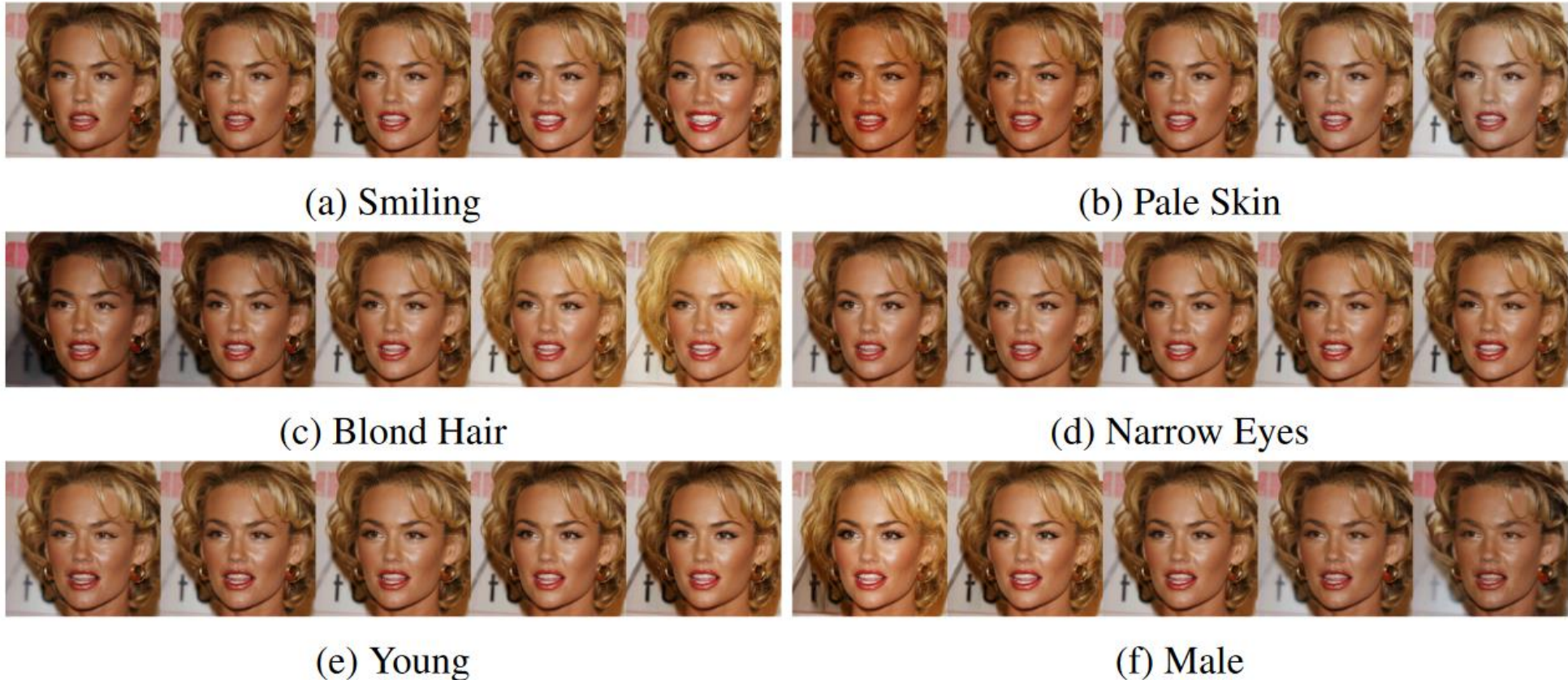


Figure 6: Manipulation of attributes of a face. Each row is made by interpolating the latent code of an image along a vector corresponding to the attribute, with the middle image being the original image

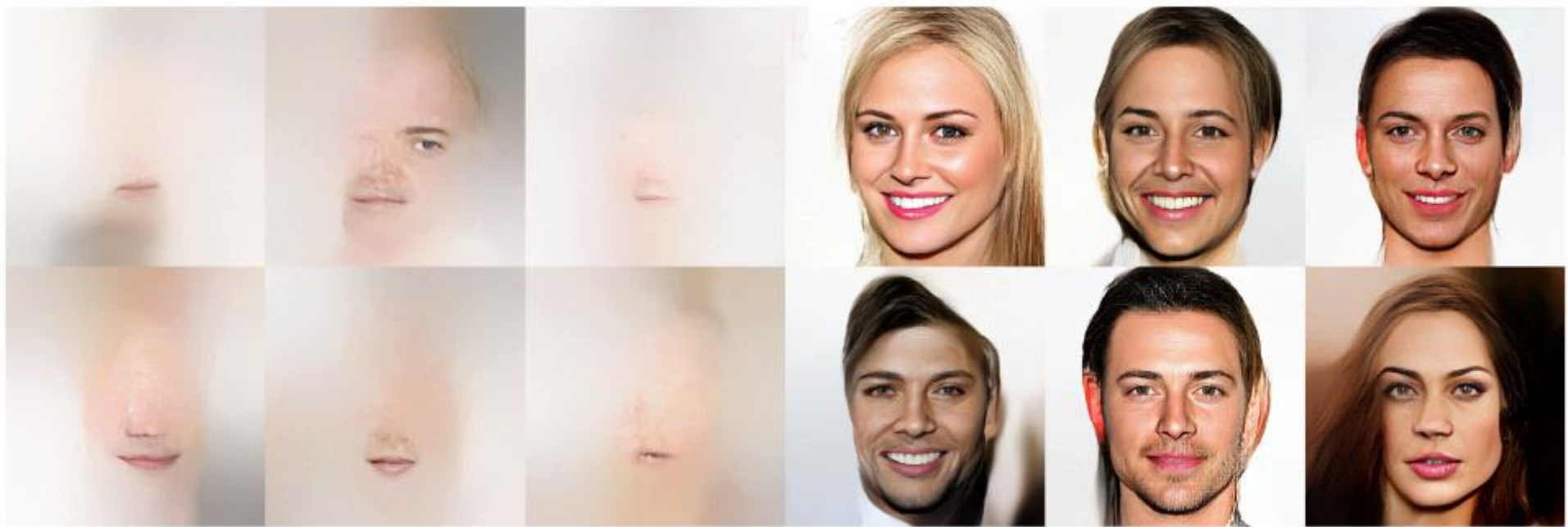
LSUN bedrooms сэмплы



Влияние температуры ($p_{x,T}(x) \propto p_X^{T^2}(x)$)



Влияние глубины модели



References

- [1] Gomez, A. N., Ren, M., Urtasun, R., and Grosse, R. B. (2017). The reversible residual network: Backpropagation without storing activations. <https://arxiv.org/abs/1707.04585>
- [2] Dinh, L., Krueger, D., and Bengio, Y. (2014). Nice: non-linear independent components estimation. <https://arxiv.org/abs/1410.8516>
- [3] Dinh, L., Sohl-Dickstein, J., and Bengio, S. (2016). Density estimation using Real NVP. <https://arxiv.org/abs/1605.08803>
- [4] Kingma, D.P., Dhariwal, P. (2018). Glow: Generative Flow with Invertible 1x1 Convolutions. <https://arxiv.org/abs/1807.03039v2>