## ПОРОЖДАЮЩИЕ СОСТЯЗАТЕЛЬНЫЕ СЕТИ

Сергей Николенко СПбГУ— Санкт-Петербург 30 октября 2025 г.



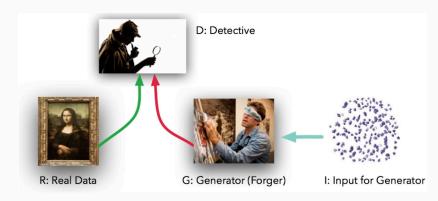


#### Random facts:

- 30 октября в России День памяти жертв политических репрессий; 30 октября 1990 г.
   на Лубянской площади был установлен Соловецкий камень в память о них
- 30 октября 1905 г. Николай II ввёл Высочайший Манифест об усовершенствовании государственного порядка («Манифест 17 октября»)
- 30 октября 1907 г. русский физик Борис Розинг получил патент №18076 на «Способ электрической передачи изображений на расстояние» (то есть телевидение)
- 30 октября 1938 г. Орсон Уэллс поставил на радиостанции CBS «Войну миров», стилизованную под прямой репортаж с места событий; многие американцы поверили, и постановка вызвала настоящую панику; в ходе последовавших разбирательств выяснилось, что около 300 тысяч американцев лично видели марсиан
- 30 октября 1941 г. Франклин Рузвельт одобрил выделение 1 млрд долларов в качестве помощи Советскому Союзу
- 30 октября 1961 г. взорвалась самая мощная бомба в мировой истории: 58-мегатонная водородная бомба («Царь-бомба») на острове Новая Земля

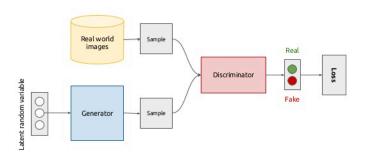
# Порождающие состязательные сети

• Порождающие состязательные сети (Generative adversarial networks, GAN): генератор vs. дискриминатор,  $p_{\mathrm{data}}$  vs.  $p_g$ .

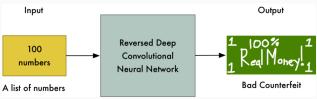


• Порождающие состязательные сети (Generative adversarial networks, GAN): генератор vs. дискриминатор,  $p_{\mathrm{data}}$  vs.  $p_g$ .

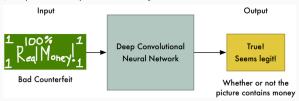
## Generative adversarial networks (conceptual)



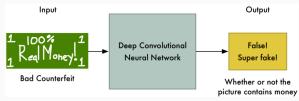
- В идеале должно быть примерно так (Geitgey, 2017):
  - сначала генератор плохой:



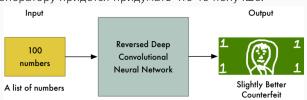
• но и дискриминатор ничего не умеет:



- В идеале должно быть примерно так (Geitgey, 2017):
  - поэтому сначала чуть обучим дискриминатор:



• ...и генератору придётся придумать что-то получше:



• И так они будут продолжать, пока всё не получится.

• Формально, генератор – это функция

$$G = G(\mathbf{z}; \theta_g) : Z \to X,$$

а дискриминатор – это функция

$$D = D(\mathbf{x}; \boldsymbol{\theta}_d) : X \to [0, 1].$$

• Целевая функция для дискриминатора – это бинарная классификация:

$$\mathbf{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \left[ \log D(\mathbf{x}) \right] + \mathbf{E}_{\mathbf{x} \sim p_q(\mathbf{x})} \left[ \log (1 - D(\mathbf{x})) \right],$$

где  $p_g(\mathbf{x}) = G_{\mathbf{z} \sim p_z}(\mathbf{z})$  – распределение, порождённое G; то есть мы берём два мини-батча, с метками 0 из генератора и с метками 1 из данных.

• А генератор обучается обманывать дискриминатор, то есть минимизировать

$$\mathbf{E}_{\mathbf{x} \sim p_g(\mathbf{x})} \left[ \log(1 - D(\mathbf{x})) \right] = \mathbf{E}_{\mathbf{z} \sim p_z(\mathbf{z})} \left[ \log(1 - D(G(\mathbf{z}))) \right].$$

• И теперь, если их объединить, получится типичная минимакс-игра:

$$\min_G \max_D V(D,G)$$
, где  $V(D,G)=\mathrm{E}_{\mathbf{x}\sim n, \ (\mathbf{x})}[\log D(\mathbf{x})]+\mathrm{E}_{z\sim n, \ (z)}[\log (1-D(G(z)))].$ 

/ı

- · Можно доказать, что  $\max_D V(D,G)$  минимизируется именно тогда, когда  $p_q = p_{\text{data}}.$
- Важное свойство: для фиксированного генератора G оптимальное распределение дискриминатора D это распределение

$$D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})}.$$

• (Упражнение: попробуйте это доказать)

• Глобальный минимум критерия достигается тогда и только тогда, когда  $p_g = p_{\mathrm{data}}$ ; сам критерий для оптимального D эквивалентен

$$\mathrm{KL}\left(p_{\mathrm{data}} \middle\| \frac{p_{\mathrm{data}} + p_g}{2}\right) + \mathrm{KL}\left(p_g \middle\| \frac{p_{\mathrm{data}} + p_g}{2}\right),$$

это т.н. дивергенция Йенсена–Шеннона (Jensen–Shannon divergence).

- . Но это всё результаты для ошибки генератора  $\mathrm{E}_{z\sim p_z(z)}[\log(1-D(G(z)))].$
- А это не самая лучшая функция ошибки для генератора...

- Чем плохо минимизировать  $\mathbf{E}_{z \sim p_z(z)}[\log(1 D(G(z)))]$ ?
  - перекрёстная энтропия хороша в классификаторе: если ответ неправильный, она не насыщается, а если правильный, насыщается и не меняется;
  - тут получается, что когда дискриминатор хорошо работает, его градиент обнуляется...
  - ...и градиент генератора тоже обнуляется!
- Поэтому на самом деле минимизируют

$$-\mathbf{E}_{z \sim p_z(z)}[\log D(G(\mathbf{z}))],$$

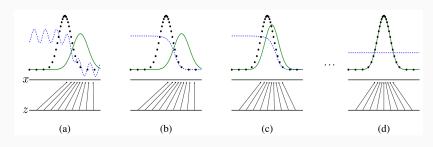
т.е. максимизируем вероятность того, что D ошибся, а не минимизируем вероятность его правильного ответа.

- Обучение похоже по сути на ЕМ-алгоритм попеременно:
  - $\cdot$  фиксируем G, обновляем дискриминатор с функцией ошибки

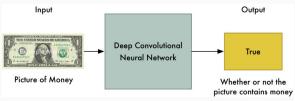
$$\mathbf{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \left[ \log D(\mathbf{x}) \right] + \mathbf{E}_{\mathbf{x} \sim p_g(\mathbf{x})} \left[ \log (1 - D(\mathbf{x})) \right],$$

 $\cdot$  фиксируем D, обновляем генератор с функцией ошибки

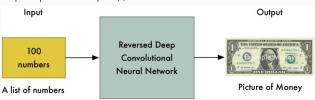
$$-\mathrm{E}_{\mathbf{x} \sim p_{\mathrm{data}}(\mathbf{x})}[\log D(\mathbf{x})],$$



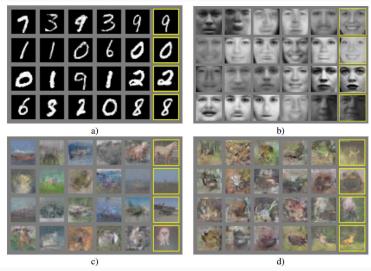
- В идеале должно быть примерно так (Geitgey, 2017):
  - дискриминатор пытается распознавать:



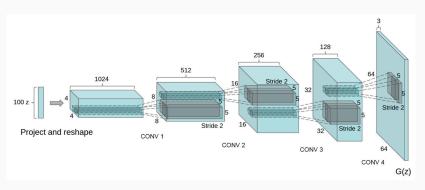
• а генератор хочет порождать:



· Сэмплирование из GAN (Goodfellow et al., 2014):



- DCGAN Deep Convolutional Generative Adversarial Networks (Radford et al., 2016).
- Несколько новых трюков (свёртки с пропусками вместо пулинга, везде batchnorm, убираем полносвязные слои, Adam...)



· Спальни (LSUN) после одной эпохи:



· Спальни (LSUN) после пяти эпох:



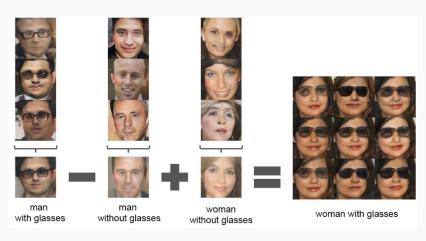
• Прогулка по пространству признаков:



• Убираем фильтры, распознающие окна:



• Векторная арифметика (как в word2vec):

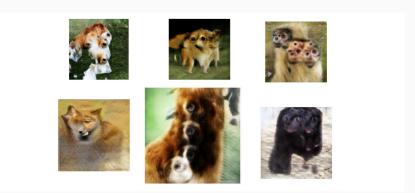


· (Geitgey, 2017): pixel art из игр NES

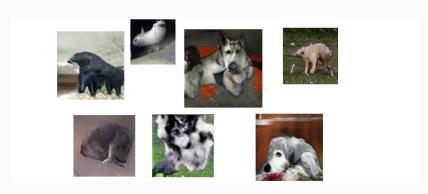


• АНИМаЦИЯ: https://medium.com/@ageitgey/abusing-generative-adversarial-networks-to-make-8-bit-pixel-art-e45d9b96cee7

• Но не всегда, конечно, работает – проблемы с подсчётом:



• Но не всегда, конечно, работает – проблемы с перспективой (right?):



## ПРОГРЕСС

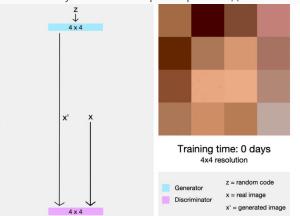
- Так работали самые первые GAN'ы.
- Но мы с тех пор прошли долгий путь:



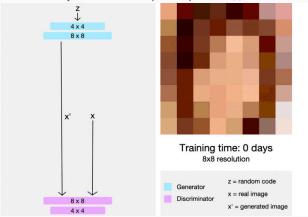
• Как же так получилось?

**ENLARGE YOUR GAN** 

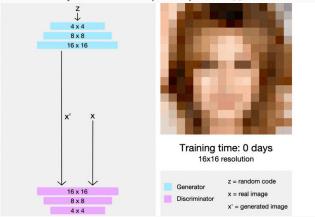
- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$
- И так постепенно увеличивают размерность до  $1024 \times 1024$



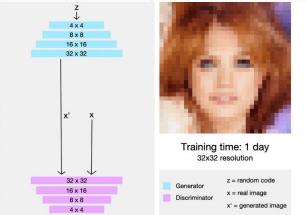
- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$
- И так постепенно увеличивают размерность до  $1024 \times 1024$



- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$
- $\cdot$  И так постепенно увеличивают размерность до  $1024 \times 1024$



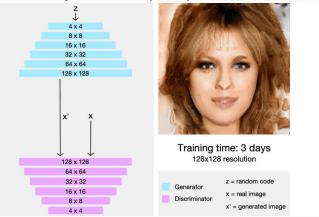
- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$
- И так постепенно увеличивают размерность до  $1024 \times 1024$



- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$
- И так постепенно увеличивают размерность до  $1024 \times 1024$



- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$



- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$



- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$



- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$



## ProGAN

- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$



- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$



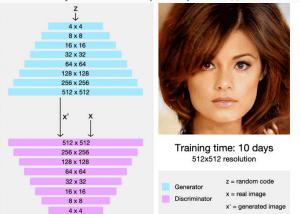
- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$

• И так постепенно увеличивают размерность до  $1024 \times 1024$ 

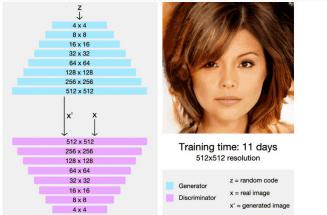


- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$

• И так постепенно увеличивают размерность до  $1024 \times 1024$ 



- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$
- И так постепенно увеличивают размерность до  $1024 \times 1024$



- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$

• И так постепенно увеличивают размерность до  $1024 \times 1024$ 



- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$

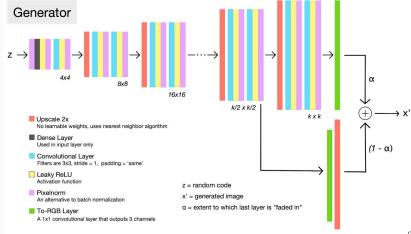
• И так постепенно увеличивают размерность до  $1024 \times 1024$ 



- Важный прорыв пришёл из NVIDIA (Karras et al., 2017)
- Progressive growing (ProGAN): на каждом этапе обучаем GAN создавать картинки вчетверо больше, начиная от  $4\times 4$
- И так постепенно увеличивают размерность до  $1024 \times 1024$

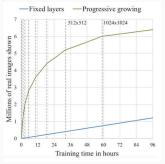


 Именно добавляем новые слои каждый раз; дискриминатор примерно симметричен



9

• Это помогает экономить время обучения, показывать больше данных за то же время



- Давайте посмотрим:
  - https://www.youtube.com/watch?v=G06dEcZ-QTg
  - https://www.youtube.com/watch?v=36lE9tV9vm0
- На других категориях до  $1024 \times 1024$  не доводили, но  $256 \times 256$  выглядит неплохо

• На других категориях до  $1024 \times 1024$  не доводили, но  $256 \times 256$  выглядит неплохо



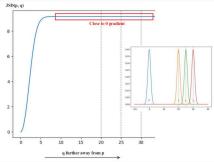
- И всё это обучается довольно стандартным способом: авторы сравнивают функции ошибки от least squares GAN и от Wasserstein GAN...
- · ...wait, what??? Давайте разбираться...

9

# Функции ошибки в GAN

#### **LSGAN**

- · (Mao et al., 2016): Least Squares GAN (LSGAN)
- Проблема с обычными GAN'ами: функция ошибки (дивергенция Йенсена-Шеннона) насыщается, когда распределение генератора далеко от правильного, и ничего не обучается.

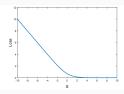


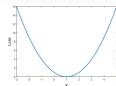
• Это потому, что дискриминатор в обычном GAN работает как классификатор с логистическим сигмоидом.

• Давайте попробуем перейти от сигмоидальной к квадратичной функции ошибки:

$$\begin{split} & \min_{D} V_{\mathrm{LSGAN}}(D) = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\mathrm{data}}} \left[ \left( D(\mathbf{x}) - b \right)^2 \right] + \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} \left[ \left( D(G(\mathbf{z})) - a \right)^2 \right], \\ & \min_{G} V_{\mathrm{LSGAN}}(G) = \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} \left[ \left( D(G(\mathbf{z})) - c \right)^2 \right], \end{split}$$

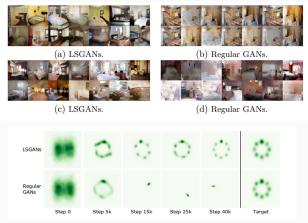
т.е. дискриминатор учится отвечать a для фейков и b для настоящих данных, а генератор хочет убедить его отвечать c на фейках. Обычно берут просто  $a=0,\,b=c=1.$ 





## **LSGAN**

 Mao et al. показали, что LSGAN устойчивее к изменениям архитектуры, меньше страдает от mode collapse; в целом, сейчас квадратичная функция ошибки применяется часто.



- (Chen et al., 2016): InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets
- Важное дополнение: было бы круто добавить disentanglement, чтобы разные признаки имели смысл
- В обычных GAN'ах используют просто вектор шума  ${f z}$ . В InfoGAN мы его разбиваем на части:
  - z это действительно несжимаемый шум, источник энтропии;
  - $\mathbf{c} = c_1 c_2 \dots c_L$  это латентный код.
- Можно в обычный GAN попробовать это добавить, чтобы генератор стал  $G(\mathbf{z},\mathbf{c})$ , но надо что-то добавить, чтобы  $\mathbf{c}$  было важно и G не мог бы просто забить на него.

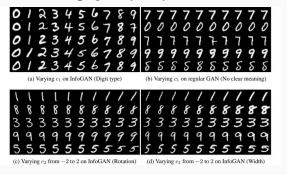
- InfoGAN: будем требовать, чтобы была большая взаимная информация (mutual information)  $I(\mathbf{c}, G(\mathbf{z}, \mathbf{c}))$ , т.е. в распределении  $G(\mathbf{z}, \mathbf{c})$  должно остаться много от  $\mathbf{c}$ .
- Формально вариационная оценка:

$$\begin{split} I(c;G(z,c)) &= H(c) - H(c|G(z,c)) \\ &= \mathbb{E}_{x \sim G(z,c)} [\mathbb{E}_{c' \sim P(c|x)} [\log P(c'|x)]] + H(c) \\ &= \mathbb{E}_{x \sim G(z,c)} [\underbrace{D_{\mathrm{KL}}(P(\cdot|x) \parallel Q(\cdot|x))}_{\geq 0} + \mathbb{E}_{c' \sim P(c|x)} [\log Q(c'|x)]] + H(c) \\ &\geq \mathbb{E}_{x \sim G(z,c)} [\mathbb{E}_{c' \sim P(c|x)} [\log Q(c'|x)]] + H(c) \end{split}$$

• И эту нижнюю оценку мы через reparametrization trick параметризуем нейросетью Q; она обычно почти целиком совпадает с D.

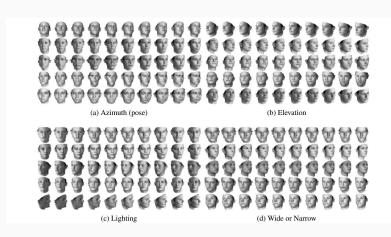
#### INFOGAN

• Пример на MNIST с  ${f c}=c_1c_2c_3$ , где  $c_1$  – дискретная метка с 10 категориями (без supervision! с равномерным априорным распределением), и  $c_2,c_3\sim U[-1,1]$ :



 $\cdot$  На самом деле по коду  $c_1$  классификатор даёт ошибку 5% на MNIST без всякого supervision.

# • 3D-лица:



# INFOGAN

## · 3D-стулья:



#### Wasserstein GAN

- · Wasserstein GAN (Arjovsky et al., 2017):
  - что такое обучить распределение вероятностей?
  - это значит минимизировать  $\mathrm{KL}(p_{\mathrm{data}} \| p_{\mathrm{model}})$ ;
  - но для этого нужно, чтобы  $p_{\mathrm{model}}$  существовала, а это неверно, если распределение сосредоточено на подмногообразиях малой размерности; вряд ли многообразие  $p_{\mathrm{model}}$  будет существенно пересекаться с многообразием  $p_{\mathrm{data}}$ , и KL будет не определено (бесконечно);
  - обычно мы добавляем шум (гауссовский, например), все модели в ML с шумом;
  - но для порождающих моделей шум мешает, делает порождённые примеры размытыми;
  - поэтому обычно шум используют для подсчёта ML, но убирают во время порождения, т.е. шум это неправильно, но надо для ML.

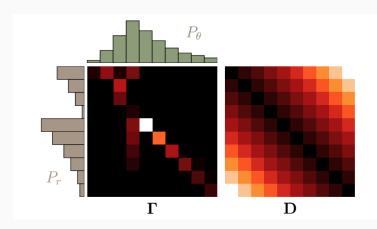
#### Wasserstein GAN

- Wasserstein GAN (Arjovsky et al., 2017): давайте посмотрим на другие метрики близости между  $p_{\rm data}$  и  $p_{
  m model}$ .
- Earth Mover distance, или Wasserstein-1:

$$W(p_{\mathrm{data}}, p_{\mathrm{model}}) = \inf_{\gamma \in \Pi(p_{\mathrm{data}}, p_{\mathrm{model}})} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \gamma} \left[ \| \mathbf{x} - \mathbf{y} \| \right],$$

где  $\Pi(p_{\mathrm{data}},p_{\mathrm{model}})$  – множество совместных распределений  $\gamma(\mathbf{x},\mathbf{y})$ , маргиналы которых –  $p_{\mathrm{data}}$  и  $p_{\mathrm{model}}$ .

• Т.е.  $\gamma(\mathbf{x},\mathbf{y})$  показывает, сколько надо «массы» перераспределить от  $\mathbf{x}$  к  $\mathbf{y}$ , чтобы превратить  $p_{\mathrm{data}}$  в  $p_{\mathrm{model}}$ .



• Пример: рассмотрим  $Z\sim U[0,1]$ ,  $\mathbb{P}_0$  – распределение (0,Z) на  $\mathbb{R}^2$ , и  $g_{\theta}(z)=(\theta,z)$ , где  $\theta$  – параметр. Тогда

$$\begin{split} \bullet \ & W(\mathbb{P}_0,\mathbb{P}_\theta) = |\theta|, \\ \bullet \ & JS(\mathbb{P}_0,\mathbb{P}_\theta) = \begin{cases} \log 2 & \text{if } \theta \neq 0 \ , \\ 0 & \text{if } \theta = 0 \ , \end{cases} \\ \bullet \ & KL(\mathbb{P}_\theta \| \mathbb{P}_0) = KL(\mathbb{P}_0 \| \mathbb{P}_\theta) = \begin{cases} +\infty & \text{if } \theta \neq 0 \ , \\ 0 & \text{if } \theta = 0 \ , \end{cases} \\ \bullet \ & \text{and } \delta(\mathbb{P}_0,\mathbb{P}_\theta) = \begin{cases} 1 & \text{if } \theta \neq 0 \ , \\ 0 & \text{if } \theta = 0 \ . \end{cases} \end{split}$$

• Т.е. только ЕМ-расстояние действительно куда-то сходится при  $\theta \to 0$ , а это ведь очень похоже на обучение распределения, сосредоточенного на подмногообразии.

### Wasserstein GAN

- Можно доказать, что  $W(p_{\mathrm{data}}, p_{\mathrm{model}})$  непрерывная функция от параметров  $\theta$  распределения  $p_{\mathrm{model}}$  при достаточно слабых условиях.
- А двойственность Монжа-Канторовича (Kantorovich-Rubinstein duality) говорит, что инфимум

$$W(p_{\text{data}}, p_{\text{model}}) = \inf_{\gamma \in \Pi(p_{\text{data}}, p_{\text{model}})} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \gamma} \left[ \| \mathbf{x} - \mathbf{y} \| \right]$$

эквивалентен

$$W(p_{\mathrm{data}}, p_{\mathrm{model}}) = \sup_{\|f\|_L \leq 1} \left( \mathbb{E}_{\mathbf{x} \sim p_{\mathrm{data}}} \left[ f(\mathbf{x}) \right] - \mathbb{E}_{\mathbf{x} \sim p_{\mathrm{model}}} \left[ f(\mathbf{x}) \right] \right),$$

где супремум берётся по всем функциям с липшицевой константой  $\leq 1$ .

- · А мы хотим обучить порождающую модель  $p_{\mathrm{model}} = g_{\theta}(\mathbf{z}).$
- Давайте всё параметризуем нейросетями.

- Сеть  $f_{\mathbf{w}}$  для функции f и сеть для  $g_{\theta}$ . Тогда:
  - · для данного  $g_{ heta}$  обучим веса  $f_{\mathbf{w}}$ , максимизируя

$$\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[ f(\mathbf{x}) \right] - \mathbb{E}_{\mathbf{x} \sim p_{\text{model}}} \left[ f(\mathbf{x}) \right];$$

 $\cdot$  для данного  $f_{\mathbf{w}}$  подсчитаем

$$\begin{split} \nabla_{\theta} W(p_{\mathrm{data}}, p_{\mathrm{model}}) &= \nabla_{\theta} \left( \mathbb{E}_{\mathbf{x} \sim p_{\mathrm{data}}} \left[ f(\mathbf{x}) \right] - \mathbb{E}_{\mathbf{x} \sim p_{\mathrm{model}}} \left[ f(\mathbf{x}) \right] \right) = \\ &= - \mathbb{E}_{\mathbf{z} \sim Z} \left[ \nabla_{\theta} f_{\mathbf{w}}(g_{\theta}(\mathbf{z})) \right]. \end{split}$$

- А чтобы  $f_{\mathbf{w}}$  оставалось липшицевым, можно просто weight clipping сделать для градиентов.

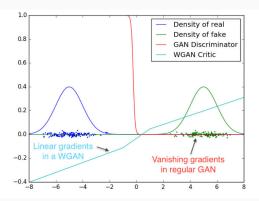
#### Wasserstein GAN

#### • Итого алгоритм обучения:

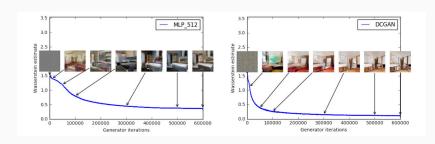
```
Algorithm 1 WGAN, our proposed algorithm. All experiments in the paper used
the default values \alpha = 0.00005, c = 0.01, m = 64, n_{\text{critic}} = 5.
Require: : \alpha, the learning rate. c, the clipping parameter. m, the batch size.
      n_{\text{critic}}, the number of iterations of the critic per generator iteration.
Require: : w_0, initial critic parameters. \theta_0, initial generator's parameters.
  1: while \theta has not converged do
           for t = 0, ..., n_{\text{critic}} do
                Sample \{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r a batch from the real data.
  3:
                Sample \{z^{(i)}\}_{i=1}^m \sim p(z) a batch of prior samples. g_w \leftarrow \nabla_w \left[\frac{1}{m}\sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m}\sum_{i=1}^m f_w(g_\theta(z^{(i)}))\right]
  4:
                w \leftarrow w + \alpha \cdot \text{RMSProp}(w, q_w)
                w \leftarrow \text{clip}(w, -c, c)
           end for
  8-
           Sample \{z^{(i)}\}_{i=1}^m \sim p(z) a batch of prior samples.
           g_{\theta} \leftarrow -\nabla_{\theta} \frac{1}{m} \sum_{i=1}^{m} f_{w}(g_{\theta}(z^{(i)}))
10:
           \theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, q_{\theta})
11:
12: end while
```

#### Wasserstein GAN

• Пример, на котором видно, что в WGAN можно спокойно обучать  $f_{\mathbf{w}}$  до сходимости, а в обычном GAN дискриминатор, может, и не стоит так обучать:



• Качество тесно коррелирует с функцией ошибки; слева тут просто полносвязная сеть в качестве генератора, справа – DCGAN:



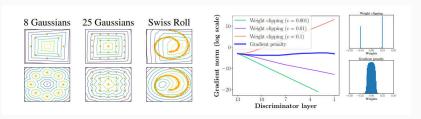
• WGAN устойчивее к изменению архитектуры; вот обычные WGAN и DCGAN (снизу):



• А вот что будет, если убрать из генератора batchnorm и сделать фильтров поменьше (одинаковое число на каждом уровне):



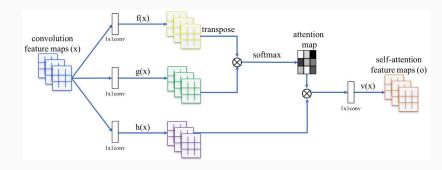
- (Gulrajani et al., 2017): Improved Training of Wasserstein GANs
- Weight clipping в критике на самом деле ведёт к проблемам, не обучаются более высокие моменты распределений и градиенты взрываются/затухают
- Лучше делать мягкий регуляризатор на градиент, типа  $\lambda \mathbb{E}_{\hat{\mathbf{x}} \sim P_{\hat{\mathbf{x}}}} \left[ \left( \| \nabla_{\hat{\mathbf{x}}} D(\hat{\mathbf{x}}) \|_2 1 \right)^2 \right]$



- (Zhang et al., 2019): Self-Attention Generative Adversarial Networks (SAGAN)
- Пытаются решить проблему с локальностью порождения свёрточными сетями



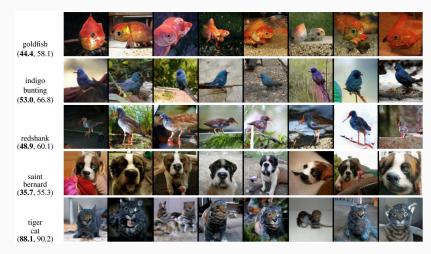
• Self-attention здесь примерно как в Transformer, только переложенный на картинки:



• И получается, что действительно генератор (здесь примеры с последнего слоя) может смотреть достаточно далеко, когда порождает картинку:



• В исходной статье (Zhang et al., 2019) уже получаются хорошие картинки:



• Следующий этап – BigGAN (Brock et al., 2019), улучшил Inception score в три раза...



• Используют SA-GAN (GAN c self-attention) со всеми новыми трюками и специальными архитектурами

• Интерполяции:



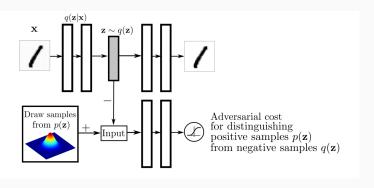
· Ой, а что такое «лучше»?..

- Трудно оценивать качество картинок.
- Inception score (Salimans et al., 2016; вообще важная статья о GAN'ax):
  - берём порождённые картинки х;
  - · применяем стандартный классификатор Inception;
  - · хотим, чтобы у  $p(y \mid \mathbf{x})$  была маленькая энтропия, но разнообразие чтобы было, т.е. чтобы у  $p(y) = \int p(y \mid \mathbf{x} = G(\mathbf{z})) \mathrm{d}\mathbf{z}$  была большая энтропия;
  - т.е. получаем метрику  $\exp(\mathbb{E}_{\mathbf{x}}\mathrm{KL}(p(y\mid \mathbf{x})\|p(y))));$
  - для обучения это трудно приспособить, а вот для оценки качества хорошо.

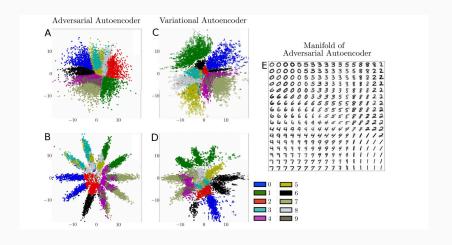
# Архитектуры, основанные на GAN

#### **AAE**

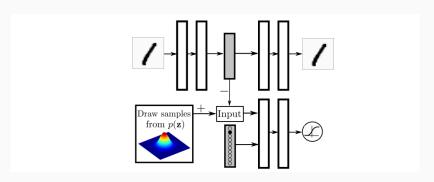
• Соперничающие автокодировщики (adversarial autoencoders; Makhzani et al., 2015): дискриминатор приводит распределение признаков (на скрытом слое) к заданному  $p(\mathbf{z})$ .

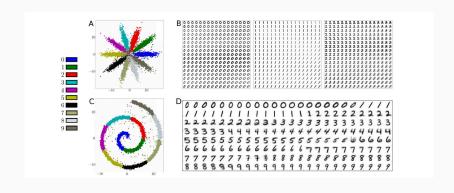


#### **AAE**



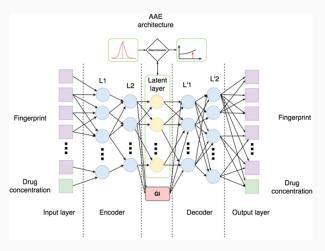
• Можно сделать распределения условными и получить semi-supervised learning:





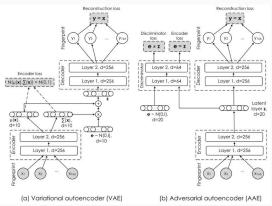
#### **AAE**

• (Kadurin et al., 2017) – ААЕ для порождения молекул в онкологии:



#### AAE

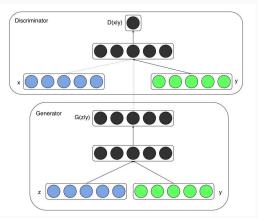
• (Kadurin et al., 2018) – druGAN для порождения молекул; сравнили с VAE:



• (Polikovsky et al., 2018): MOSES – платформа для сравнения порождающих моделей для молекул

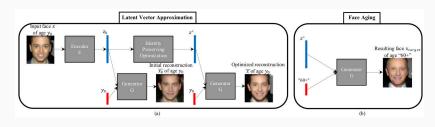
#### Условные GAN'ы

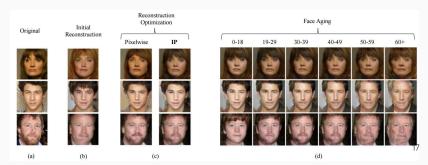
• Условные GAN'ы (conditional GANs; Mirza and Osindero, 2014): G получает случайный вектор и вектор входа, вход – это условие.



#### Условные GAN'ы

· (Antipov et al., 2017) – состарим лицо условным GAN'ом:





#### Условные GAN'ы

 (Ledig et al., 2017) – ESRGAN для superresolution (но об этом можно долго разговаривать)



#### STACKED GANS

- · Часто полезно сделать несколько GAN'ов подряд.
- · (Huang et al., 2017): Stacked Generative Adversarial Networks
- $\cdot$   $G_i$  последовательно обучает более низкоуровневые представления, обращая глубокую сеть из  $E_i$ .
- Новые loss functions:
  - $\cdot$  conditional loss мы обучаем  $\hat{\mathbf{h}}_i$  при условии  $\mathbf{h}_{i+1}$ , и чтобы генератор не игнорировал условие, надо, чтобы восстанавливались  $\mathbf{h}_{i+1}$  через соответствующий encoder:

$$\mathcal{L}^{\mathrm{cond}} = \mathbb{E}_{\mathbf{h}_{i+1} \sim p_{\mathrm{data}}, \mathbf{z}_i \sim p_{\mathbf{z}_i}} \left[ d(E_i(G_i(\mathbf{h}_{i+1}, \mathbf{z}_i)), \mathbf{h}_{i+1}) \right];$$

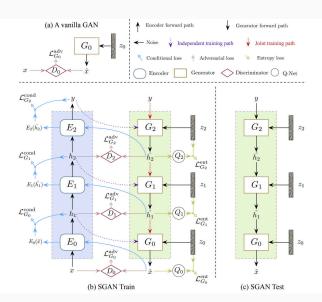
 $\cdot$  entropy loss – но и надо, чтобы  $G_i$  не игнорировал  ${f z}$ , надо добавить diversity; для этого максимизируем энтропию:

$$\mathcal{L}^{\text{ent}} = \mathbb{E}_{\mathbf{z}_i \sim p_{\mathbf{z}_i}} \left[ \mathbb{E}_{\hat{\mathbf{h}}_i \sim G_i(\hat{\mathbf{h}}_i | \mathbf{z}_i)} \left[ -\log Q_i(\mathbf{z}_i \mid \hat{\mathbf{h}}_i] \right],$$

где  $Q_i$  – параметризованное нейросетью приближение для апостериорного распределения  $p_i(\mathbf{z}_i\mid\hat{\mathbf{h}}_i)$  (вариационная оценка).

#### STACKED GANS

#### · SGAN:



#### STACKED GANS

• Получается лучше, чем у предшественников (хотя ещё лучше progressive growing):



(a) SGAN samples (conditioned on (b) Real images (nearest neighbor) labels)



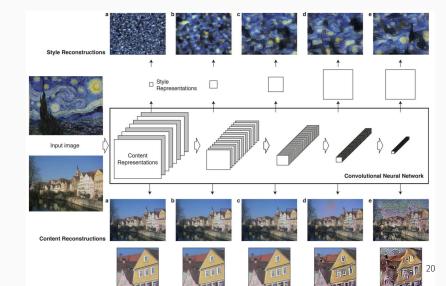
(c) SGAN samples (conditioned on (d) SGAN samples (conditioned generated fc3 features) on generated fc3 features, trained without entropy loss)

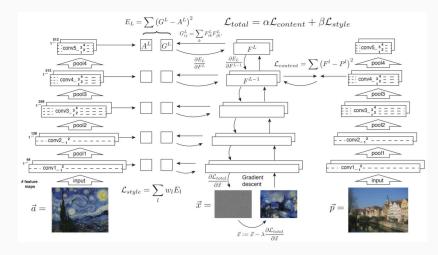
Method	Score
Infusion training [1]	$4.62 \pm 0.06$
ALI [10] (as reported in [63])	$5.34 \pm 0.05$
GMAN [11] (best variant)	$6.00 \pm 0.19$
EGAN-Ent-VI [4]	$7.07 \pm 0.10$
LR-GAN [65]	$7.17 \pm 0.07$
Denoising feature matching [63]	$7.72 \pm 0.13$
DCGAN <sup>†</sup> (with labels, as reported in [61])	6.58
SteinGAN <sup>†</sup> [61]	6.35
Improved GAN <sup>†</sup> [53] (best variant)	$8.09 \pm 0.07$
AC-GAN <sup>†</sup> [43]	$8.25 \pm 0.07$
DCGAN $(\mathcal{L}^{adv})$	$6.16 \pm 0.07$
DCGAN $(\mathcal{L}^{adv} + \mathcal{L}^{ent})$	$5.40 \pm 0.16$
DCGAN $(\mathcal{L}^{adv} + \mathcal{L}^{cond})^{\dagger}$	$5.40 \pm 0.08$
DCGAN $(\mathcal{L}^{adv} + L^{cond} + \mathcal{L}^{ent})^{\dagger}$	$7.16 \pm 0.10$
SGAN-no-joint <sup>†</sup>	$\pmb{8.37} \pm 0.08$
$SGAN^{\dagger}$	$\pmb{8.59} \pm 0.12$
Real data	$11.24 \pm 0.1$

<sup>†</sup> Trained with labels.

Table 1: **Inception Score on CIFAR-10.** SGAN and SGANno-joint outperform previous state-of-the-art approaches.

# Case study: перенос стиля





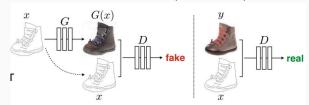


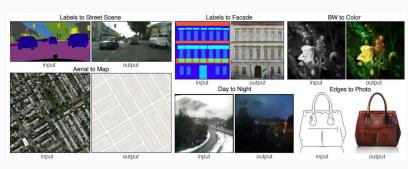




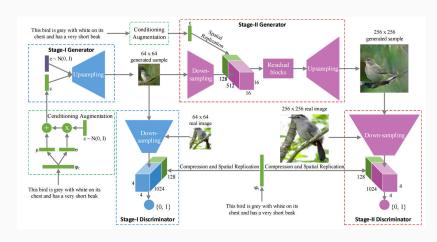


• pix2pix (Isola et al., 2017): отлично работает с paired dataset





- Стиль не обязан быть изображением!
- · StackGAN (Zhang et al., 2016) по тексту породим картинку:



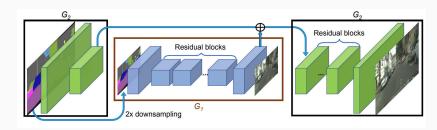
• Вот результаты реального порождения (best of 16):

Text description	This bird is red and brown in color, with a stubby beak	The bird is short and stubby with yellow on its body	A bird with a medium orange bill white body gray wings and webbed feet	This small black bird has a short, slightly curved bill and long legs	A small bird with varying shades of brown with white under the eyes	A small yellow bird with a black crown and a short black pointed beak	This small bird has a white breast, light grey head, and black wings and tail
64x64 GAN-INT-CLS [22]				A	A.	A	0
128x128 GAWWN [20]						P	
256x256 StackGAN				1	A.		

• pix2pixHD (Wang et al., 2017) – то же самое в высоком разрешении



• Генератор теперь из двух частей, вторая улучшает результат первой



• Ещё кое-какие трюки, в общем, лучше получается:



- А параллельно с этим Huang и Belongie придумали AdaIN (adaptive instance normalization):
  - · batch normalization использует статистики по мини-батчу:

$$BN(x) = \gamma \left(\frac{x - \mu(x)}{\sigma(x)}\right) + \beta$$

- instance normalization  $\mathrm{IN}(x)$  это то же самое, но статистики по каждому каналу и каждой картинке по отдельности
- conditional instance normalization это когда  $\gamma$  и  $\beta$  обучаются для каждого стиля:

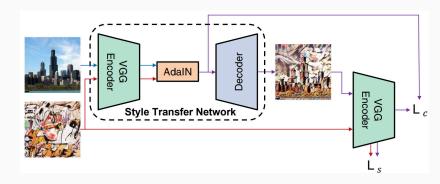
$$CIN(x;s) = \gamma_s \left(\frac{x - \mu(x)}{\sigma(x)}\right) + \beta_s$$

• AdaIN – это то же самое, но мы просто берём параметры от другой картинки, которая задаёт стиль:

AdaIN
$$(x, y) = \sigma(y) \left( \frac{x - \mu(x)}{\sigma(x)} \right) + \mu(y)$$

21

• Теперь сама сеть супер-простая – AdalN действует в feature space какого-нибудь автокодировщика:



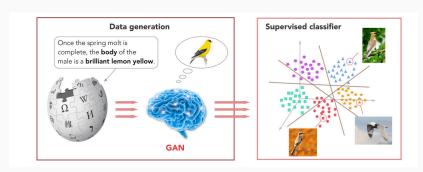
• И получается круто:



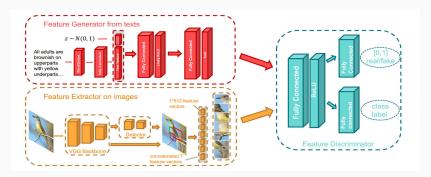
· А ещё можно интерполировать и задавать «силу» стиля:



- (Zhu et al., CVPR 2018): A Generative Adversarial Approach for Zero-Shot Learning from Noisy Texts
- Для 2018 это были вообще безумные вещи можно обучать классификатор картинок на текстах из википедии:



• Логичная архитектура:



• Тут важно, что для классификатора не надо порождать картинки, можно просто признаки порождать.

• И ведь работает:

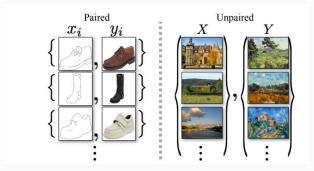


• Можно делать zero-shot retrieval: искать картинки птичек, которых модель никогда не видела



#### **CYCLEGAN**

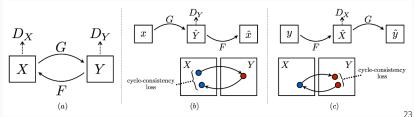
- · CycleGAN (Zhu et al., 2017)
  - · что делать, если данные unpaired?
  - например, для style transfer:



#### CYCLEGAN

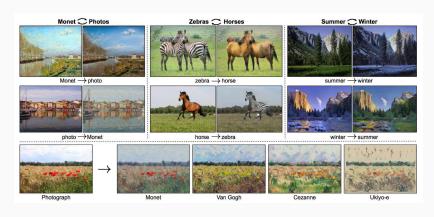
- CvcleGAN (Zhu et al., 2017)
- · Основная идея в том, что после двойного цикла style transfer должно опять получиться то же самое,  $G(F(\mathbf{x})) = \mathbf{x}$
- $\cdot$  Общая функция потерь складывается из двух GAN'ов для G и F и cycle loss:

$$\begin{split} \mathcal{L}(G, F, D_X, D_Y) = & \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) + \mathcal{L}_{\text{GAN}}(F, D_X, Y, X) \\ & + \lambda \mathcal{L}_{\text{cyc}}(G, F) \end{split}$$



#### **CYCLEGAN**

- · CycleGAN (Zhu et al., 2017)
- И получается очень хороший style transfer без всяких выровненных данных



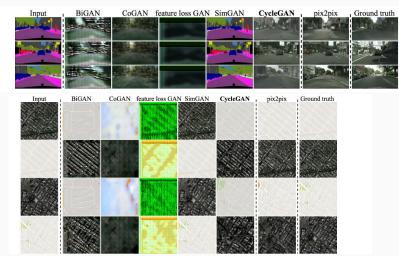
## **CYCLEGAN**

- · CycleGAN (Zhu et al., 2017)
- Можно посмотреть на реконструкции тоже



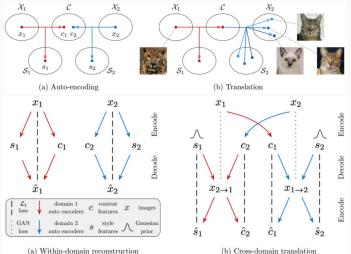
### **CYCLEGAN**

- · CycleGAN (Zhu et al., 2017)
- А можно сравнить с предшественниками



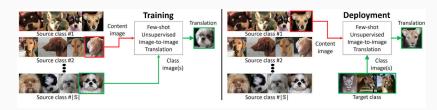
### MUNIT IN FUNIT

- Но всем этим моделям (кроме AdaIN) нужны большие датасеты в каждом стиле
- · Следующий шаг MUNIT (Huang et al., 2018)



#### MUNIT u FUNIT

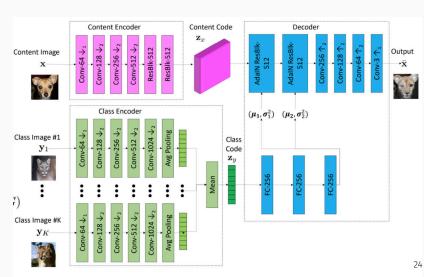
• FUNIT (Liu et al., 2019) может делать трансфер между многими стилями, определяемыми несколькими картинками





#### MUNIT u FUNIT

• Идея: условный генератор и multitask adversarial дискриминатор

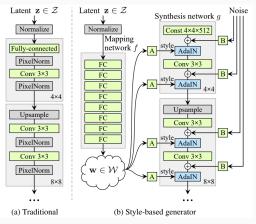


• StyleGAN от NVIDIA (Karras et al., 2019) порождает лица с условиями, взятыми из других лиц:

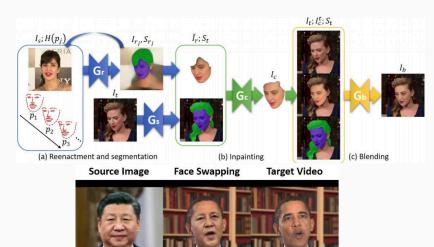
https://www.youtube.com/watch?v=kSLJriaOumA



• Смысл в том, чтобы подавать стиль на вход на разных уровнях генератора, и использовать его в AdalN:



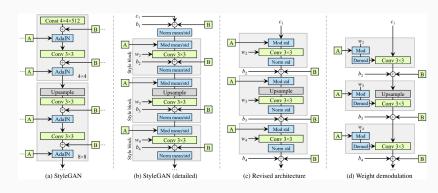
· Deepfakes тоже так работают: FSGAN (Nirkin et al., 2019)



- StyleGAN2 (Karras et al., 2019) NVIDIA смотрела на то, как дальне улучшать качество картинок, получающихся из StyleGAN
- Неочевидные проблемы например, вот этот артефакт происходит от instance normalization:



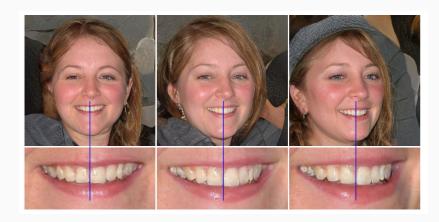
· Поэтому они переделали архитектуры сетей в StyleGAN2:



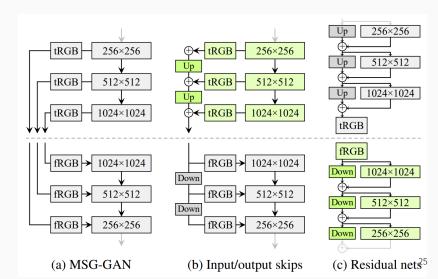
• Заменяют нормализацию на демодуляцию, т.е. чуть другое преобразование:



• Есть артефакты и от progressive growing:



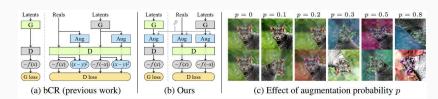
• Поэтому вместо этого сравнивают архитектуры, которые добиваются того же эффекта внутри сети:



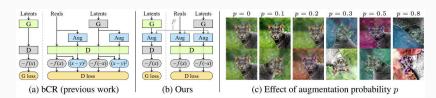
• Ну и, конечно, ещё лучше получается:



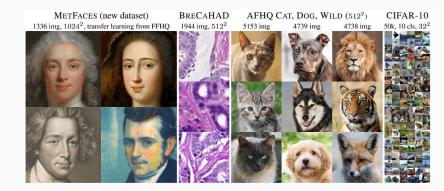
- Следующая работа, (Karras et al., Oct. 2020), обращается к тому, что делать, когда данных мало
- Тут очень важны аугментации
- Раньше аугментации применялись к тому, что видит D, плюс регуляризатор на то, чтобы они не «протекали», т.е. чтобы дискриминатор не обращал внимания на аугментации



- $\cdot$  Но тогда аугментации могут «протекать» всё равно, ведь теперь для D не важно, была аугментация или нет
- · Karras et al. применяют аугментации просто стохастически
- Вторая новизна адаптивная аугментация, т.е. p настраивается динамически в зависимости от того, насколько сильный оверфиттинг мы видим



• Получается отлично, особенно при transfer learning, т.е. если начать с конфигурации, обученной на другом датасете:



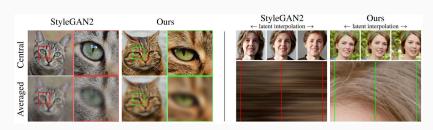
• Вот так делал StyleGAN2:



· А вот так новый StyleGAN Ada:



- · (Karras et al., 2021): StyleGAN 3
- Улучшают latent interpolation (прогулку по пространству признаков) и борются с texture sticking (чтобы изменения в пространстве признаков были более непрерывными)



 Для этого переходят к непрерывным картам признаков! Не будем подробно об этом, но начинается DSP: нужны свёртки с дельта-функциями, чтобы перейти от непрерывного к дискретному, а потом низкочастотные фильтры, чтобы перейти обратно

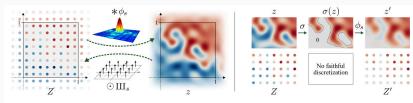
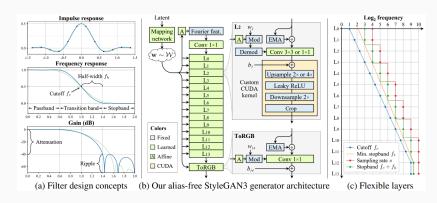
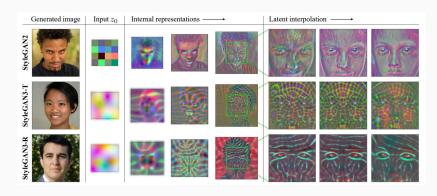


Figure 2: Left: Discrete representation Z and continuous representation z are related to each other via convolution with ideal interpolation filter  $\phi_s$  and pointwise multiplication with Dirac comb  $\text{III}_s$ . Right: Nonlinearity  $\sigma$ , ReLU in this example, may produce arbitrarily high frequencies in the continuous-domain  $\sigma(z)$ . Low-pass filtering via  $\phi_s$  is necessary to ensure that Z' captures the result.

 В самой архитектуре тоже дискретные вещи заменяются на непрерывные



• В результате получается, что некоторые признаки кодируют фазу, а не только амплитуду сигналов



# Спасибо!

# Спасибо за внимание!



