

# Глубинное обучение

Лекция 8: Оптимизация против нейросетей  
или Adversarial X

Лектор: Максим Рябинин  
Старший исследователь, Yandex Research

Программа ML Residency: [yandex.ru/yaintern/research\\_ml\\_residency](https://yandex.ru/yaintern/research_ml_residency)

ФКН ВШЭ, 2022



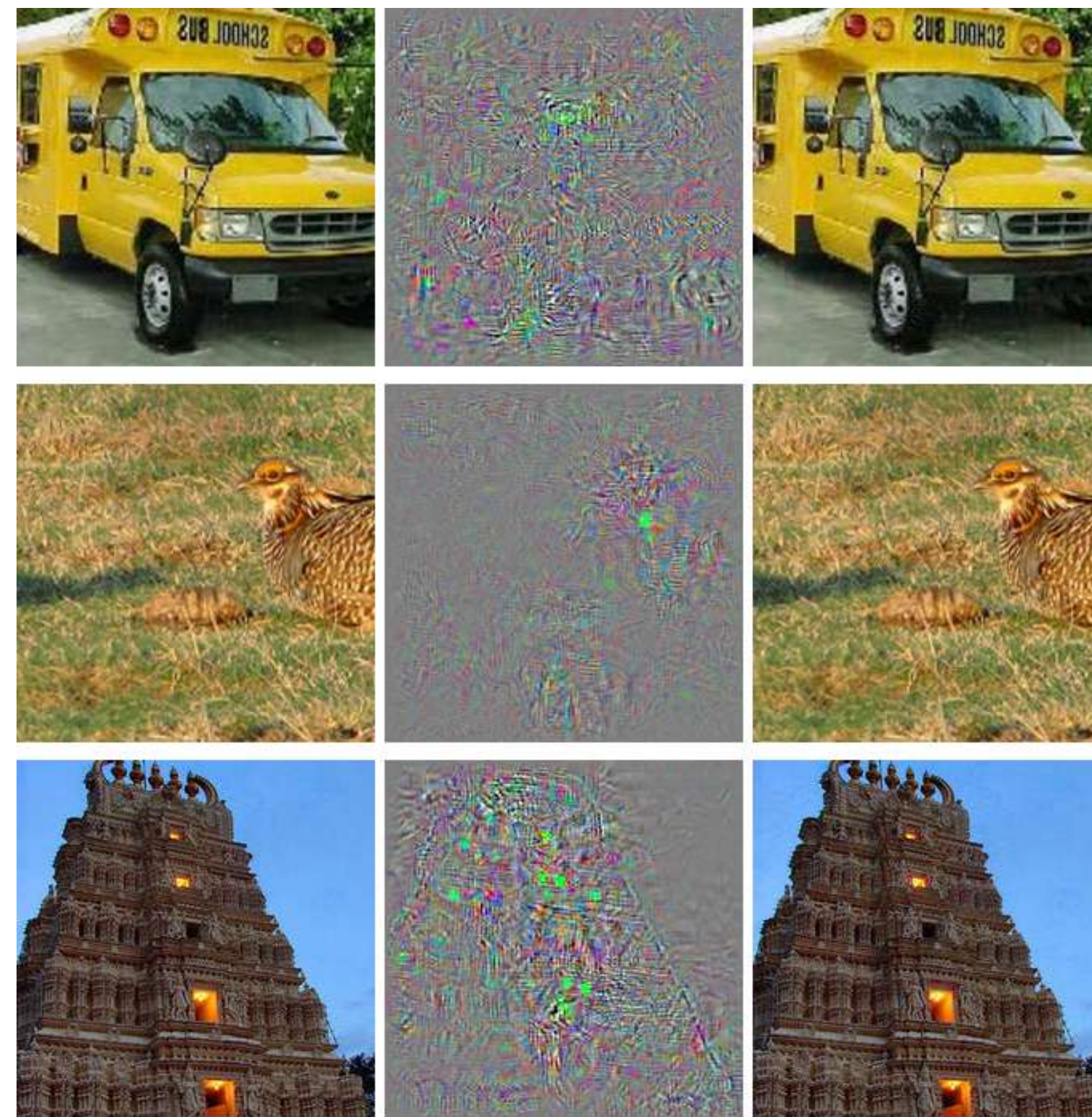
# План лекции

- Противоборствующие/сопоставительные примеры (adversarial examples)
  - Примеры, причины, следствия
  - Способы построения, борьбы, применения
- Адаптация к данным (domain adaptation)
  - Что это? Зачем?
  - Простейшие методы
- Нейросети как функции потерь
  - Generative Adversarial Networks
  - ELECTRA



# Adversarial examples

- Невидимые глазу изменения меняют результат сети!



правильная  
классификация

шум

«ostrich»

[Szegedy et al., 2013]



# Adversarial examples for text

- Изменения текста (дискретные) меняют результат сети!
- Удаления и вставки текста (настоящий текст и опечатки)

The Old Harbor Reservation Parkways are three *historie* roads in the Old Harbor area of Boston. *Some exhibitions of Navy aircrafts were often held here.* They are part of the Boston parkway system designed by Frederick Law Olmsted. They include all of William J. Day Boulevard running from *Castle* Island to Kosciuszko Circle along Pleasure Bay and the Old Harbor shore. The part of Columbia Road from its northeastern end at Farragut Road west to Pacuska Circle (formerly called Preble Circle). Old Harbor Reservation

**83.7% Building    =>    88.7% Means of Transportation**

# Adversarial атаки и защиты от них

- Пусть у нас есть обученная сеть (например, на ImageNet)
  - Сеть и функция потерь  $J_{\theta}(x, l)$
  - Атака оптимизацией ( $\eta = x' - x$ ,  $l'$  – целевая метка для атаки):
$$\min_{x'} c\|\eta\| + J_{\theta}(x', l')$$
$$s.t. \quad x' \in [0, 1].$$
  - Быстрая атака знаком градиента (Fast Gradient Sign Method, FGSM)

$$x' = x - \epsilon \text{sign}(\nabla_x J(\theta, x, l'))$$

- Простые варианты защит:
  - Обучение на атаках, ансамбли, дистилляция сетей

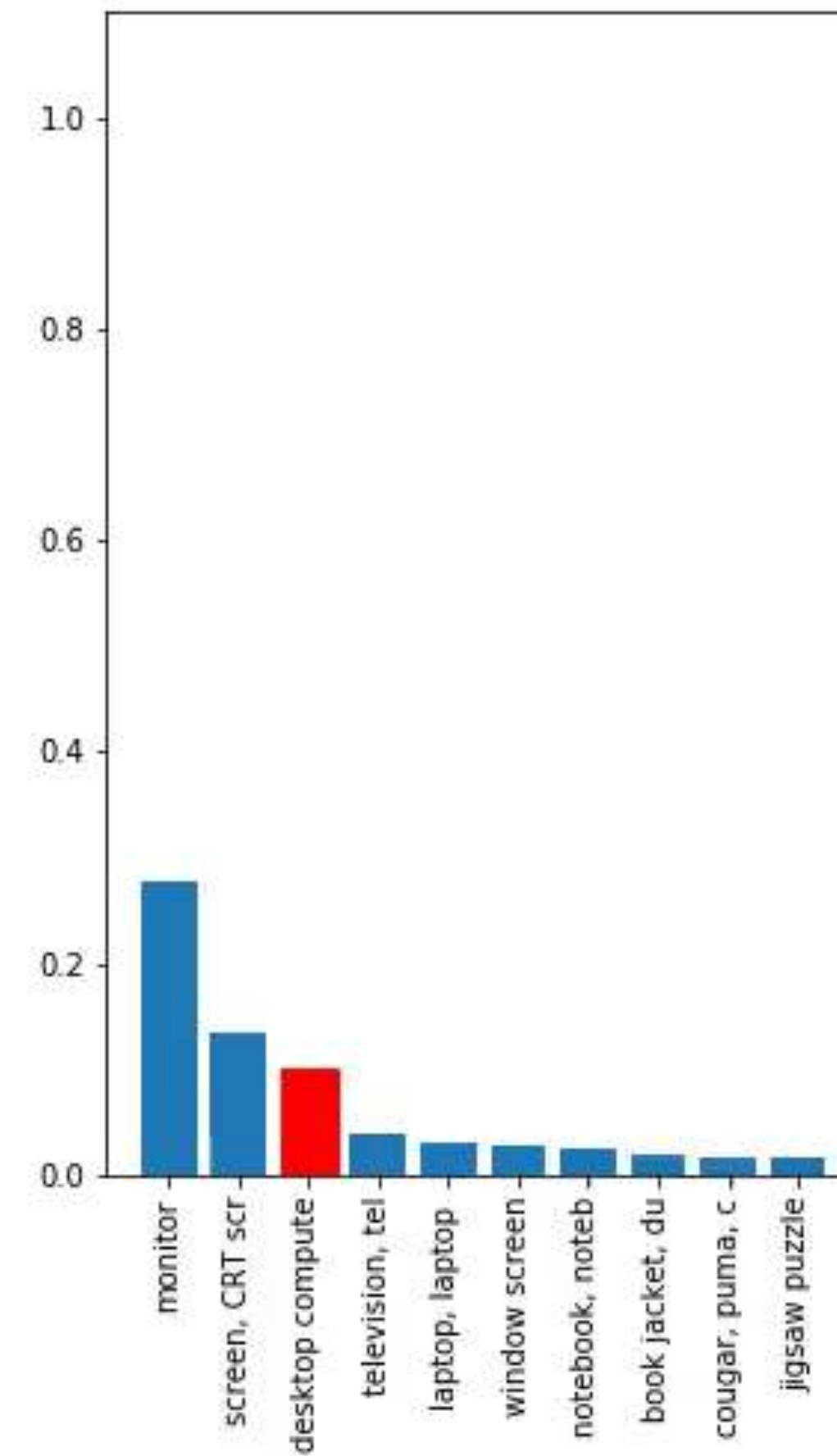
# Что все это значит?

- Adversarial examples – область с большими следствиями
  - Обзор: arXiv: 1712.07107 [Yuan et al., 2017]
- Атаки обобщаются! [Papernot et al., 2016]
- Виды атак:
  - White-box vs. black-box
  - Targeted vs. non-targeted
- Другие алгоритмы ML тоже можно атаковать!



# Adversarial example можно напечатать!

- Можно строить устойчивые атаки!



# Adversarial training

- Обучение на adversarial-атаках полезно не только для устойчивости к ним же, но и для общего качества модели
- Free Large-Batch (FreeLB):  
градиентный подъем по шуму, накапливаем градиенты на каждой из итераций
- Стабильно улучшает качество на NLP-бенчмарках

---

**Algorithm 1** “Free” Large-Batch Adversarial Training (FreeLB- $K$ )

---

**Require:** Training samples  $X = \{(\mathbf{Z}, y)\}$ , perturbation bound  $\epsilon$ , learning rate  $\tau$ , ascent steps  $K$ , ascent step size  $\alpha$

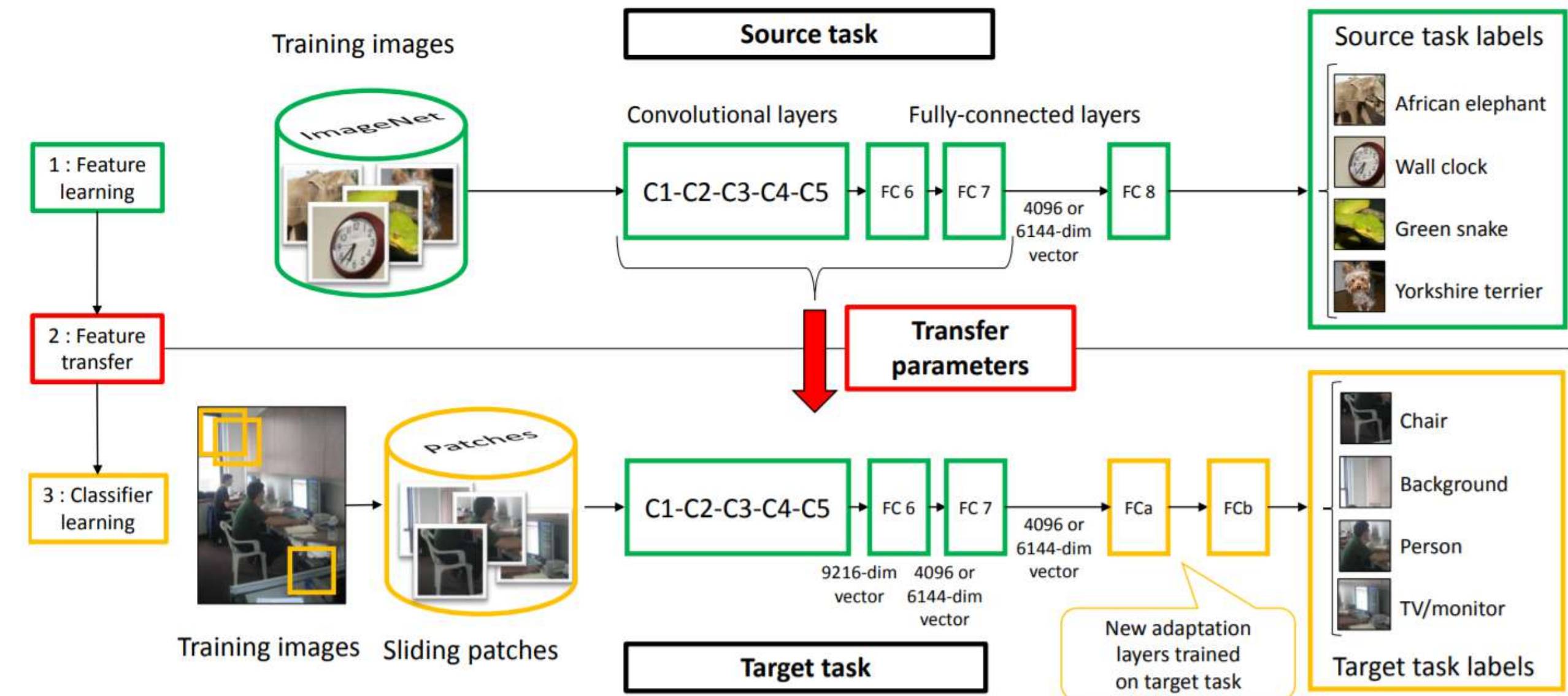
```
1: Initialize  $\theta$ 
2: for epoch = 1 ...  $N_{ep}$  do
3:   for minibatch  $B \subset X$  do
4:      $\delta_0 \leftarrow \frac{1}{\sqrt{N_\delta}} U(-\epsilon, \epsilon)$ 
5:      $\mathbf{g}_0 \leftarrow 0$ 
6:     for  $t = 1 \dots K$  do
7:       Accumulate gradient of parameters  $\theta$ 
8:        $\mathbf{g}_t \leftarrow \mathbf{g}_{t-1} + \frac{1}{K} \mathbb{E}_{(\mathbf{Z}, y) \in B} [\nabla_{\theta} L(f_{\theta}(\mathbf{X} + \delta_{t-1}), y)]$ 
9:       Update the perturbation  $\delta$  via gradient ascend
10:       $\mathbf{g}_{adv} \leftarrow \nabla_{\delta} L(f_{\theta}(\mathbf{X} + \delta_{t-1}), y)$ 
11:       $\delta_t \leftarrow \Pi_{\|\delta\|_F \leq \epsilon} (\delta_{t-1} + \alpha \cdot \mathbf{g}_{adv} / \|\mathbf{g}_{adv}\|_F)$ 
12:    end for
13:     $\theta \leftarrow \theta - \tau \mathbf{g}_K$ 
14:  end for
15: end for
```

---



# Адаптация к данным (domain adaptation)

- При обучении сетей часто сбор данных – узкое место
- Идея: обучиться на других данных и использовать эту сеть
- Формула CV 2013-2017:
  - Взять SOTA-модель, обученную на ImageNet
  - Откинуть голову сети, дообучиться на свою задачу

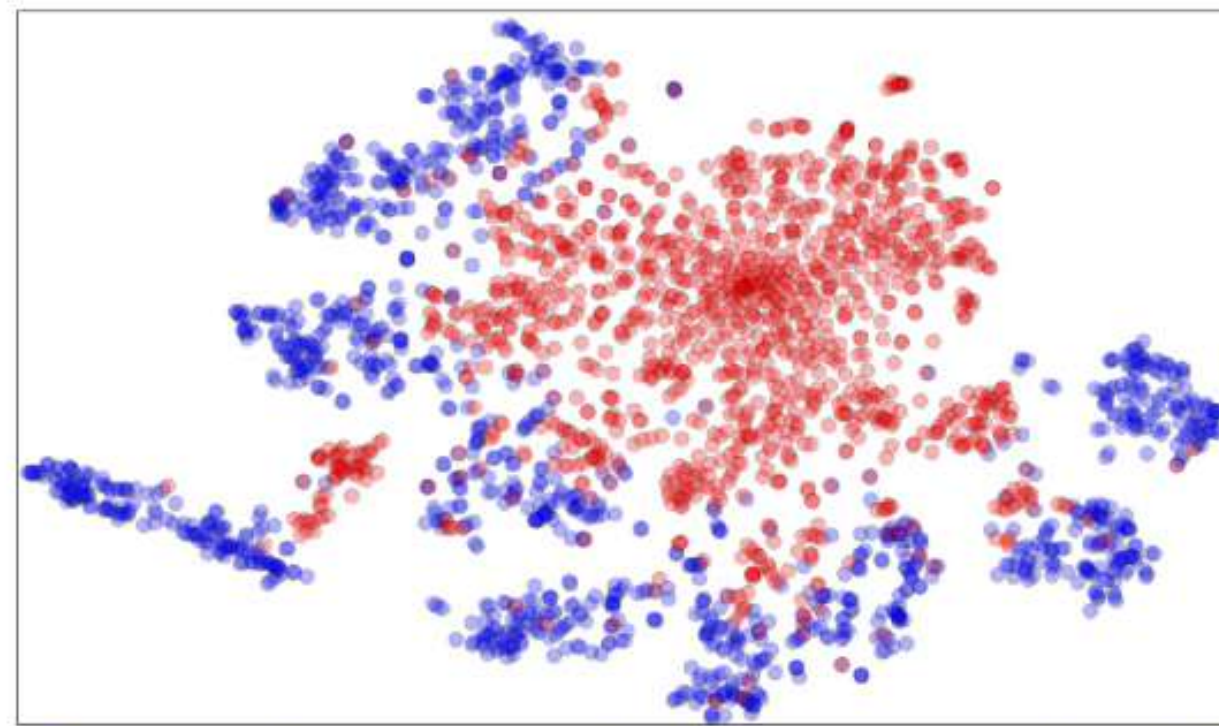


# Адаптация к данным (domain adaptation)

- При обучении сетей часто сбор данных – узкое место
- Идея – обучиться на других данных и использовать эту сеть
- Формула CV 2013-2017:
  - Взять SOTA-модель, обученную на ImageNet
  - Откинуть голову сети, дообучиться на свою задачу
- Формула работает хуже, если данные слишком различны
  - Явление называют domain/distribution shift
- Доменная адаптация
  - Есть много размеченных данных в source domain
  - Есть много неразмеченных данных в target domain
  - Как это использовать?

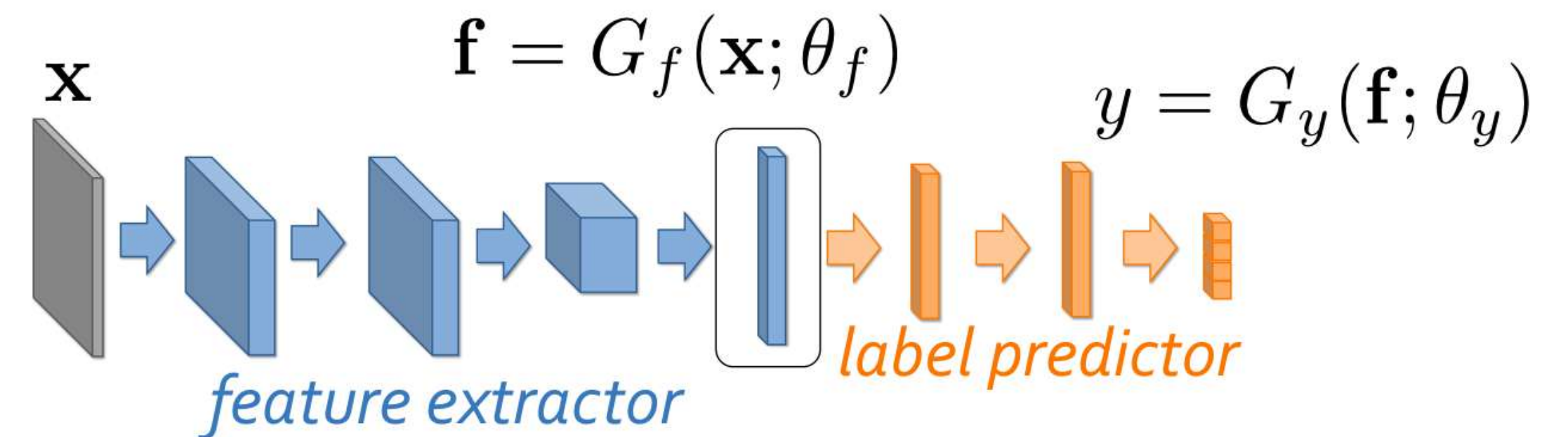
# Как убрать разницу (domain shift)?

- Где проявляется разница?
- На промежуточных слоях разные значения признаков



$$S(\mathbf{f}) = \{G_f(\mathbf{x}; \theta_f) \mid \mathbf{x} \sim S(\mathbf{x})\}$$

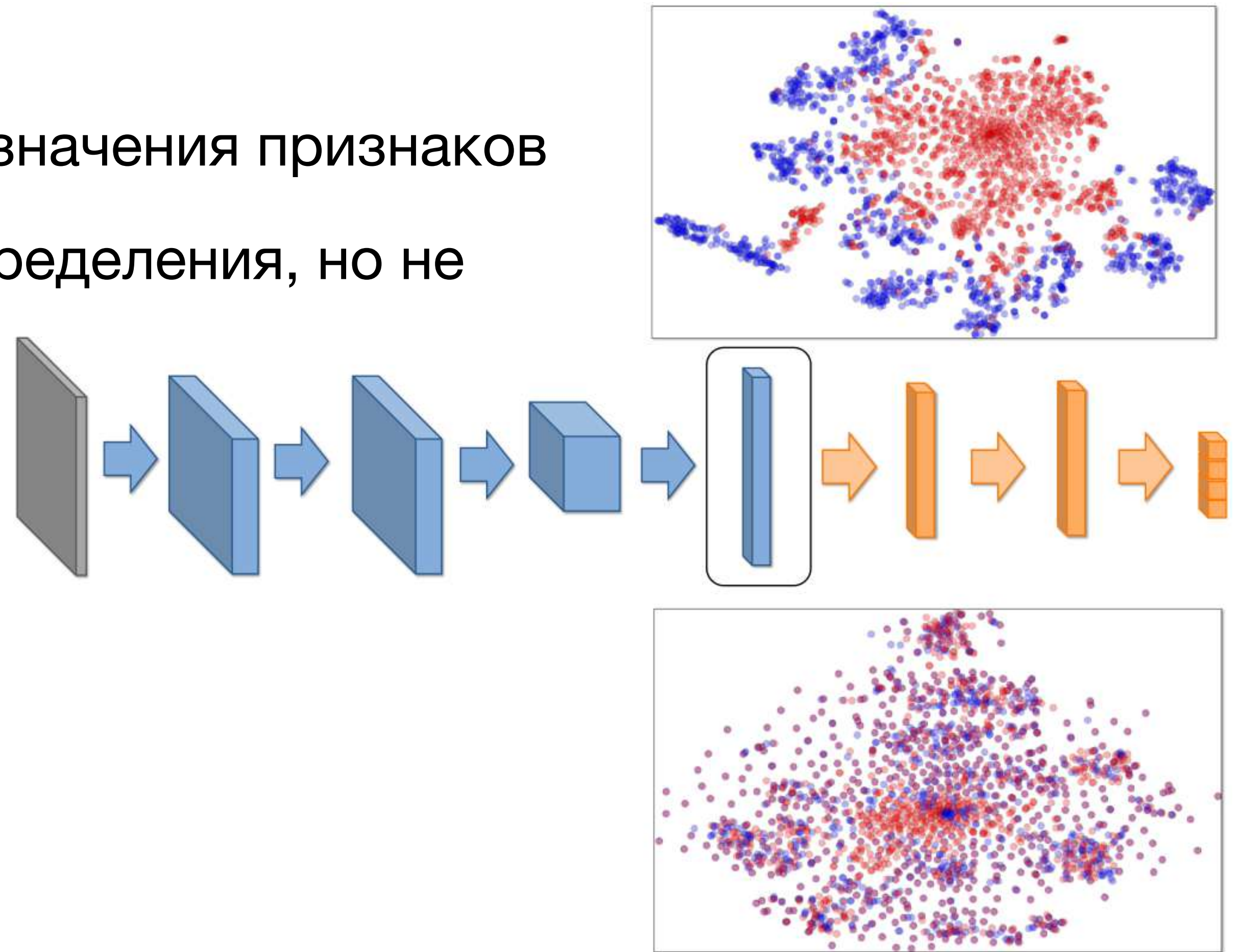
$$T(\mathbf{f}) = \{G_f(\mathbf{x}; \theta_f) \mid \mathbf{x} \sim T(\mathbf{x})\}$$





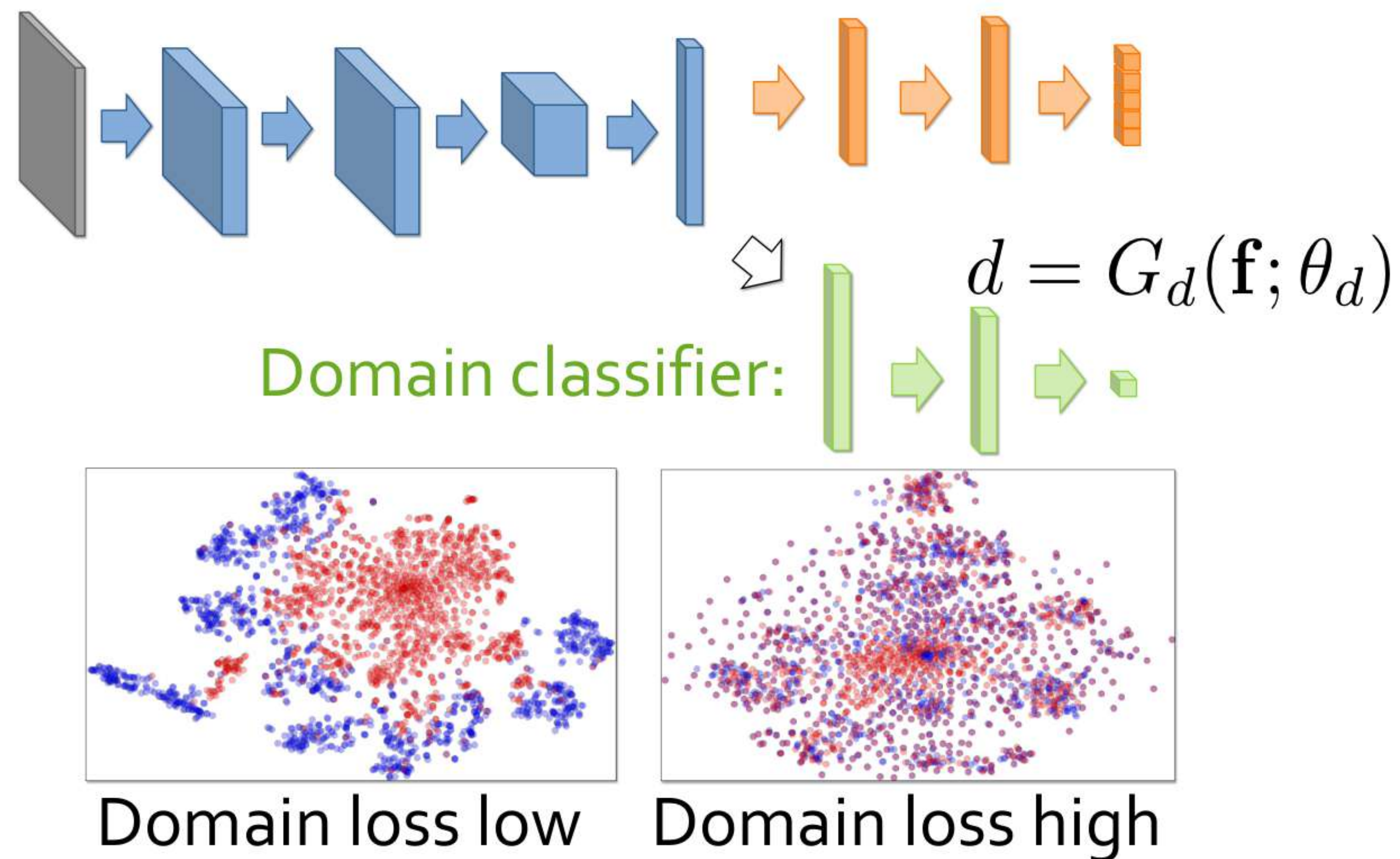
# Как убрать разницу (domain shift)?

- Где проявляется разница?
- На промежуточных слоях разные значения признаков
- Цель – получить одинаковые распределения, но не сломать классификацию



# Как убрать разницу (domain shift)?

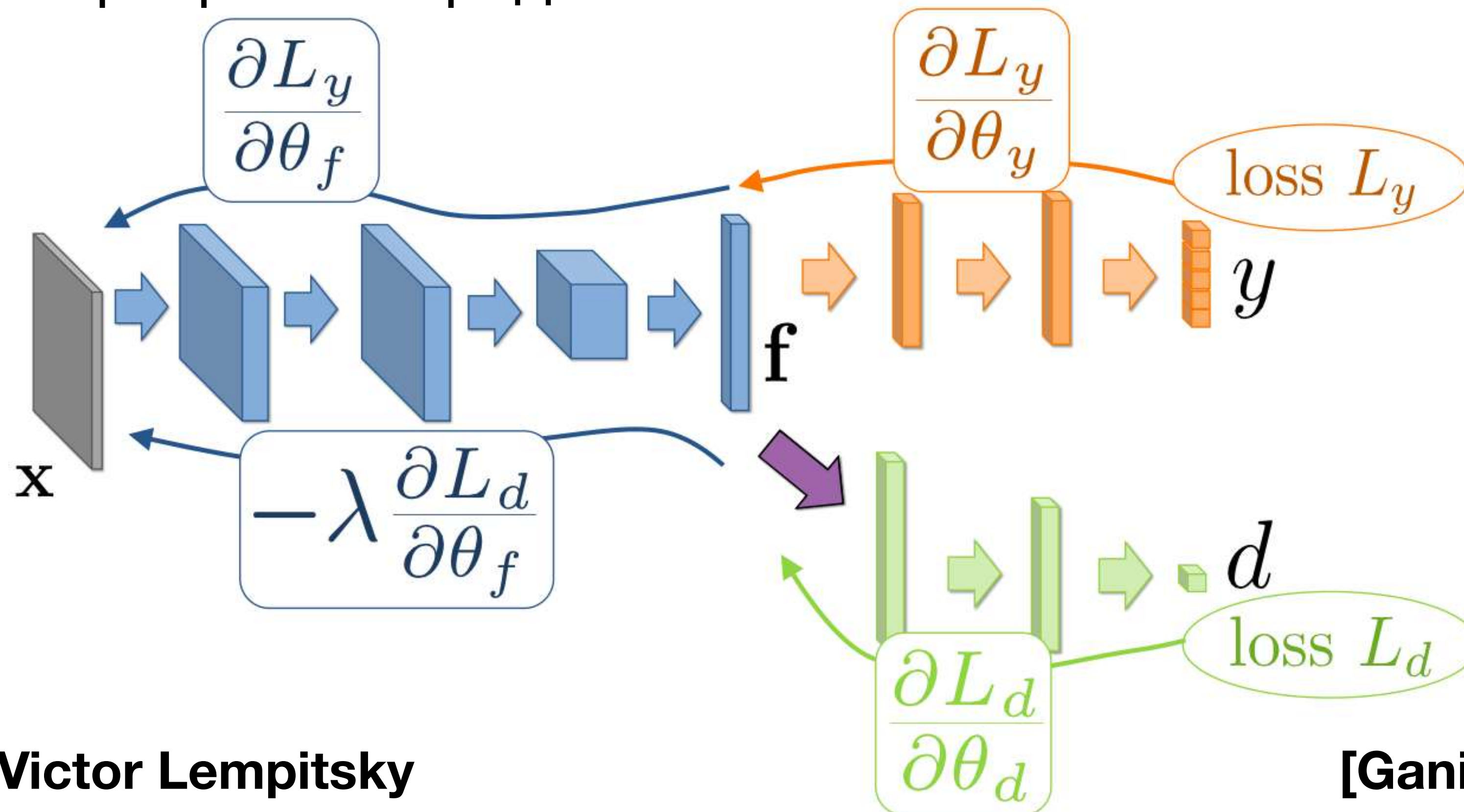
- Идея – использовать классификатор доменов





# Как убрать разницу (domain shift)?

- Идея – использовать классификатор доменов
- ~~Совместное обучение классификаторов?~~
- Нужно инвертировать градиент!





# Как убрать разницу (domain shift)?

- Идея – использовать классификатор доменов
- ~~Совместное обучение классификаторов?~~
- Нужно инвертировать градиент!

- Pytorch code:

```
class GradReverse(Function):  
    def forward(self, x):  
        return x.view_as(x)  
    def backward(self, grad_output):  
        return -lambda * grad_output  
  
def grad_reverse(x):  
    return GradReverse()(x)
```

- В обычном слое сети:

```
def forward(self, x):  
    x = grad_reverse(x)
```

# Как убрать разницу (domain shift)?

- Идея – использовать классификатор доменов
- ~~Совместное обучение классификаторов?~~
- Нужно инвертировать градиент!
- Интерпретация через седловую точку:

$$\begin{aligned} E(\theta_f, \theta_y, \theta_d) &= \sum_{\substack{i=1..N \\ d_i=0}} L_y(G_y(G_f(\mathbf{x}_i; \theta_f); \theta_y), y_i) - \lambda \sum_{i=1..N} L_d(G_d(G_f(\mathbf{x}_i; \theta_f); \theta_d), y_i) \\ &= \sum_{\substack{i=1..N \\ d_i=0}} L_y^i(\theta_f, \theta_y) - \lambda \sum_{i=1..N} L_d^i(\theta_f, \theta_d) \end{aligned}$$

- Оптимизация:
$$(\hat{\theta}_f, \hat{\theta}_y) = \arg \min_{\theta_f, \theta_y} E(\theta_f, \theta_y, \hat{\theta}_d)$$
$$\hat{\theta}_d = \arg \max_{\theta_d} E(\hat{\theta}_f, \hat{\theta}_y, \theta_d)$$

# Generative Adversarial Networks

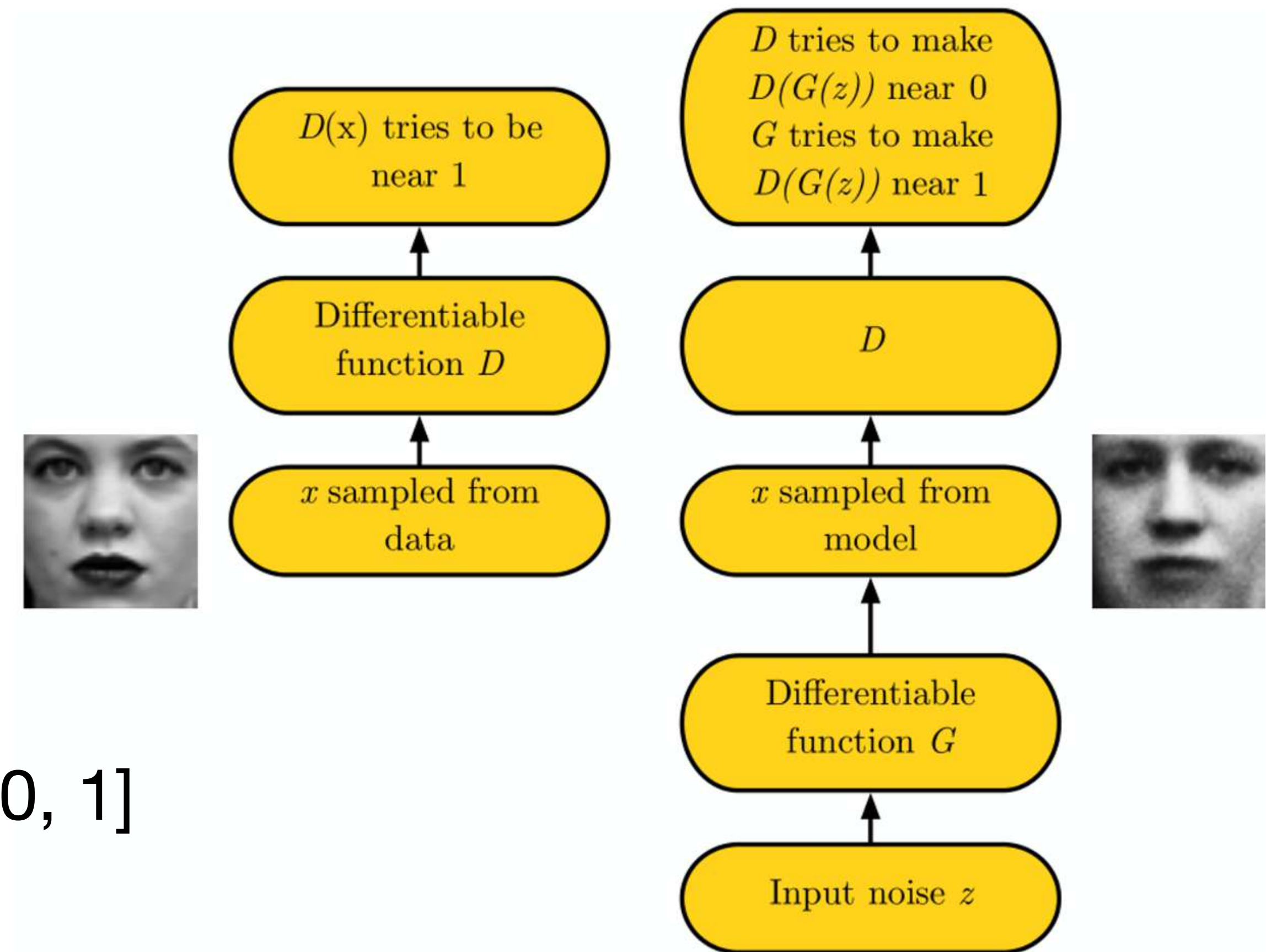
- Генеративные модели (обычно для изображений)
- Основная идея: вместо определения целевой функции (правдоподобие, ошибка реконструкции) целевая функция **обучается вместе с данными**
- Генератор – сеть, синтезирующая картинки из шума
- Дискриминатор – сеть, отличающая настоящие от синтезированных



# Generative Adversarial Networks

Две сети:

- $G$  – генератор, выдает изображение
- $D$  – дискриминатор, выдает число из  $[0, 1]$



# Generative Adversarial Networks

**Седловая точка:**

$$\min_G \max_D \mathbb{E}_{x \sim \text{data}} [\log D(x)] + \mathbb{E}_{z \sim \text{noise}} [\log(1 - D(G(z)))]$$

**Обучение**

**дискриминатора:**

$$\max_D \mathbb{E}_{x \sim \text{data}} [\log D(x)] + \mathbb{E}_{z \sim \text{noise}} [\log(1 - D(G(z)))]$$

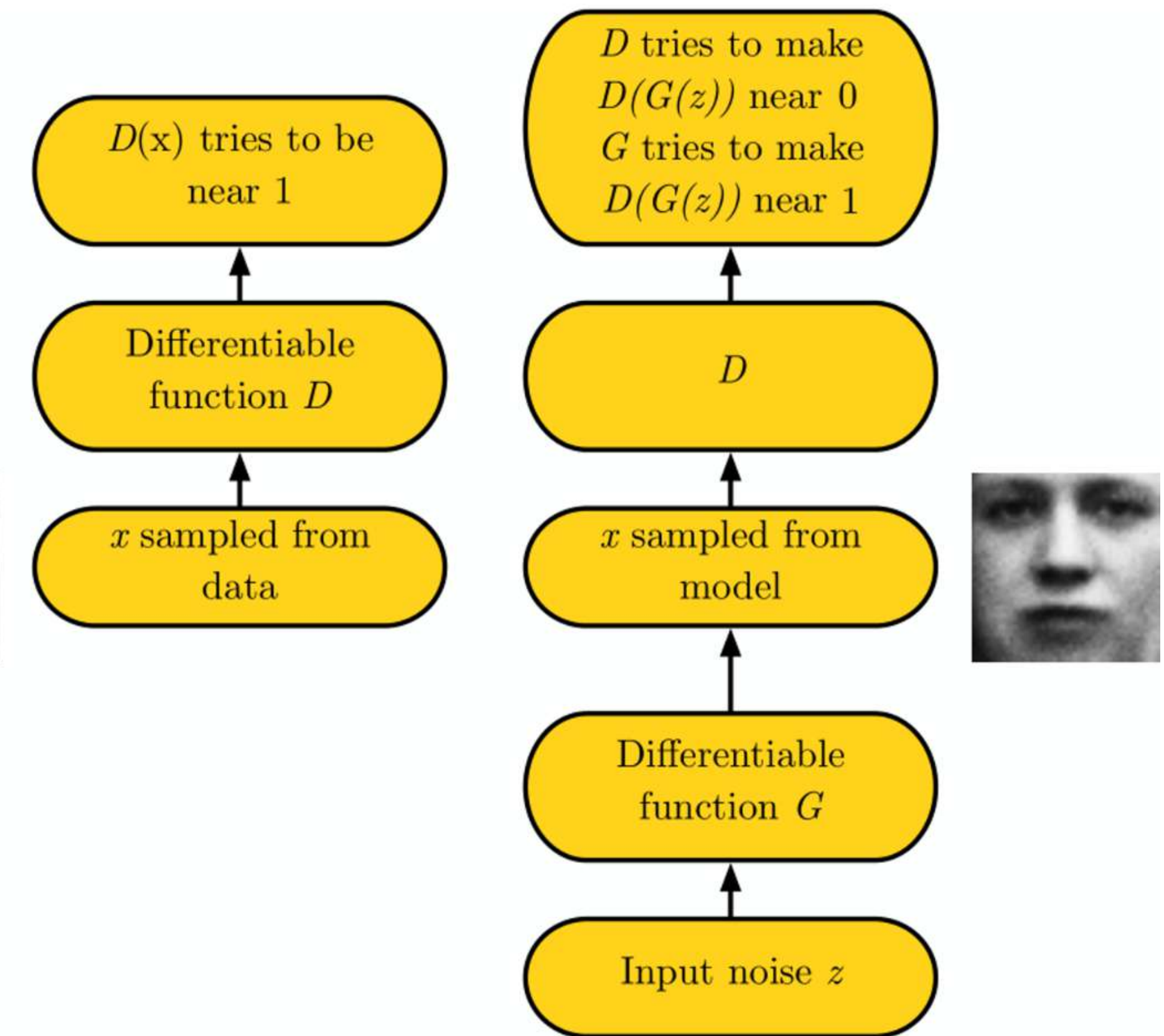
соответствует минимизации log loss

**Обучение генератора:**

~~$$\min_G \mathbb{E}_{z \sim \text{noise}} [\log(1 - D(G(z)))]$$~~

$$\max_G \mathbb{E}_{z \sim \text{noise}} [\log D(G(z))]$$

Затухает градиент, когда  
плохой G (начало обучения)

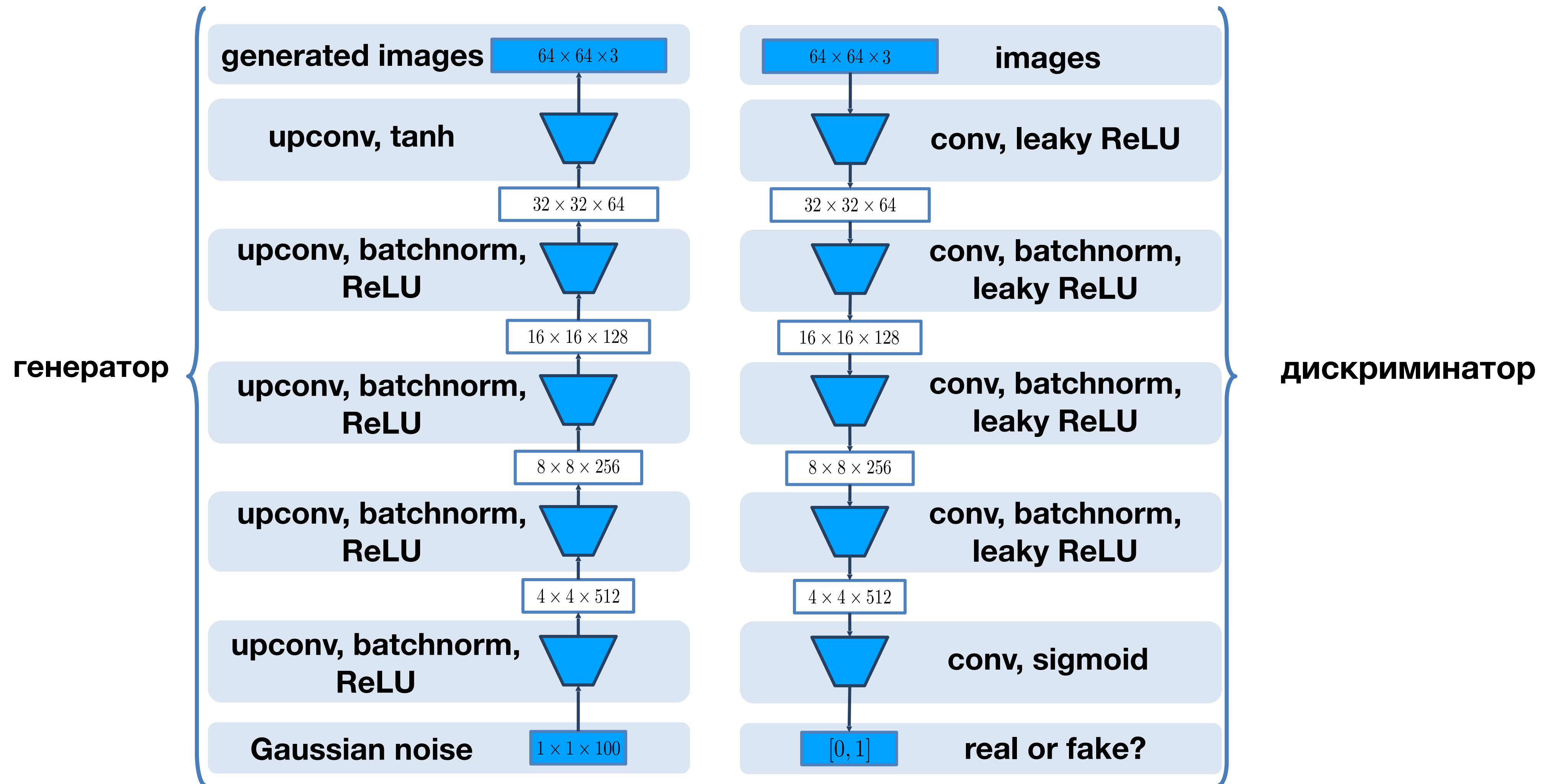


# GANs работают?

- Часто модель сложно заставить работать
- Много тонкостей реализации – очень важен исходный код
- Требуется очень много вычислительных ресурсов
- BigGAN [Brock et al., 2018], StyleGANv2 [Karras et al., 2021], StyleGANv3 (aka Alias-Free GAN) [Karras et al., 2021]



# Архитектура DCGAN

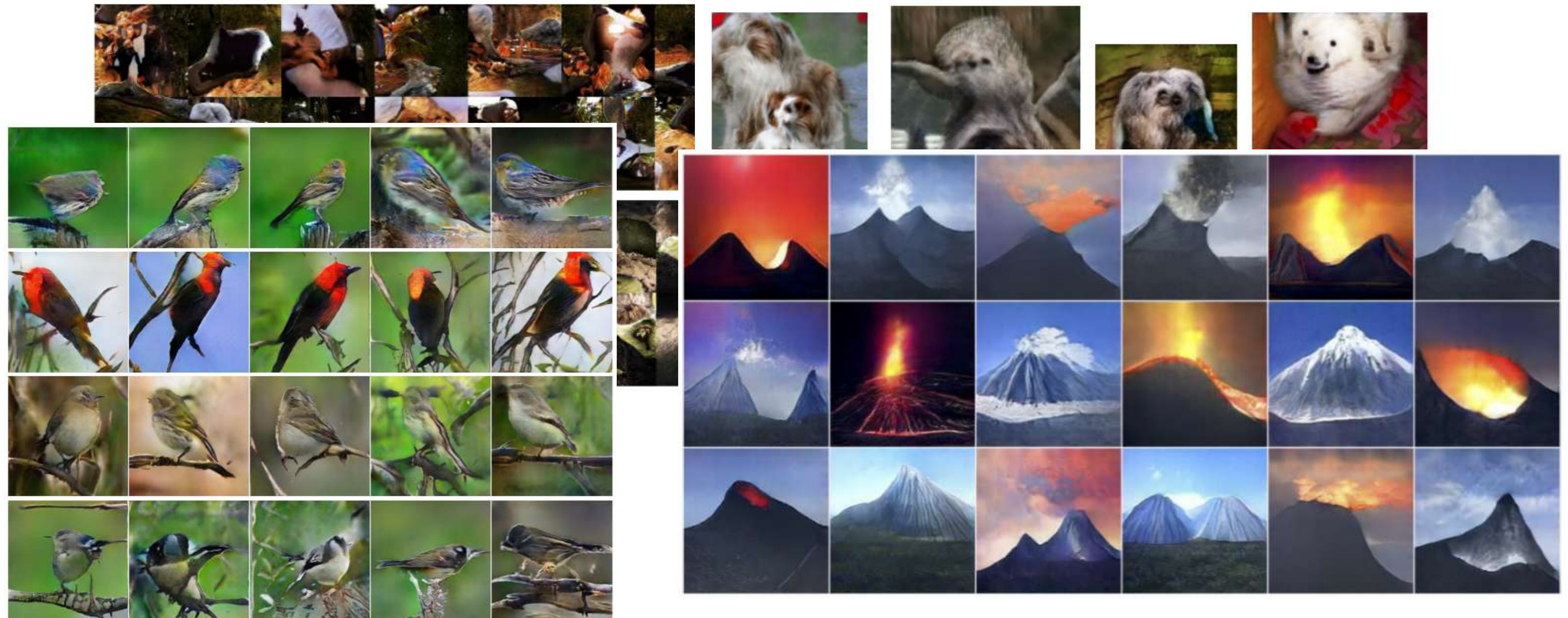


[Radford et al. 2014]



# Что генерируют GANs?

- Много хайпа
- Примеры из разных GAN-методов:



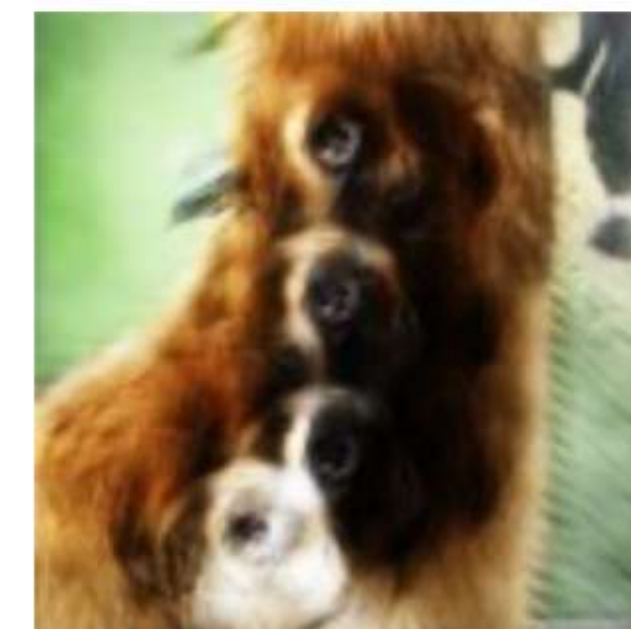
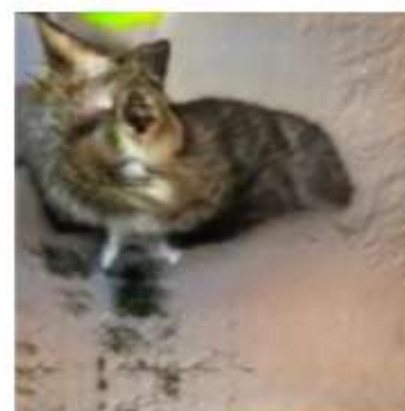
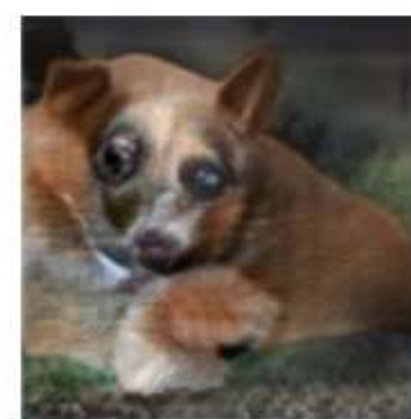
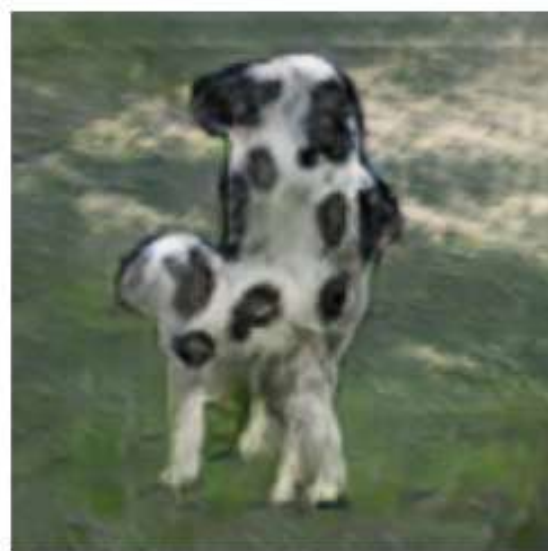
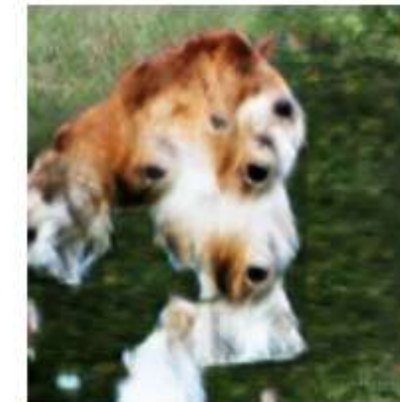
Images from Goodfellow (2016)



# Что генерируют GANs?

- Сложно генерировать реалистичные изображения
- Но естественные изображения – сложный объект

**Глобальная  
структура:**



**Подсчёт  
частей:**

**Images from Goodfellow (2016)**



# GANs быстро развиваются



2014



2015



2016



2017



2018

<https://youtu.be/XOxxPcy5Gr4>



2020

<https://youtu.be/c-NJtV9Jvp0>



# GANs быстро развиваются: BigGAN

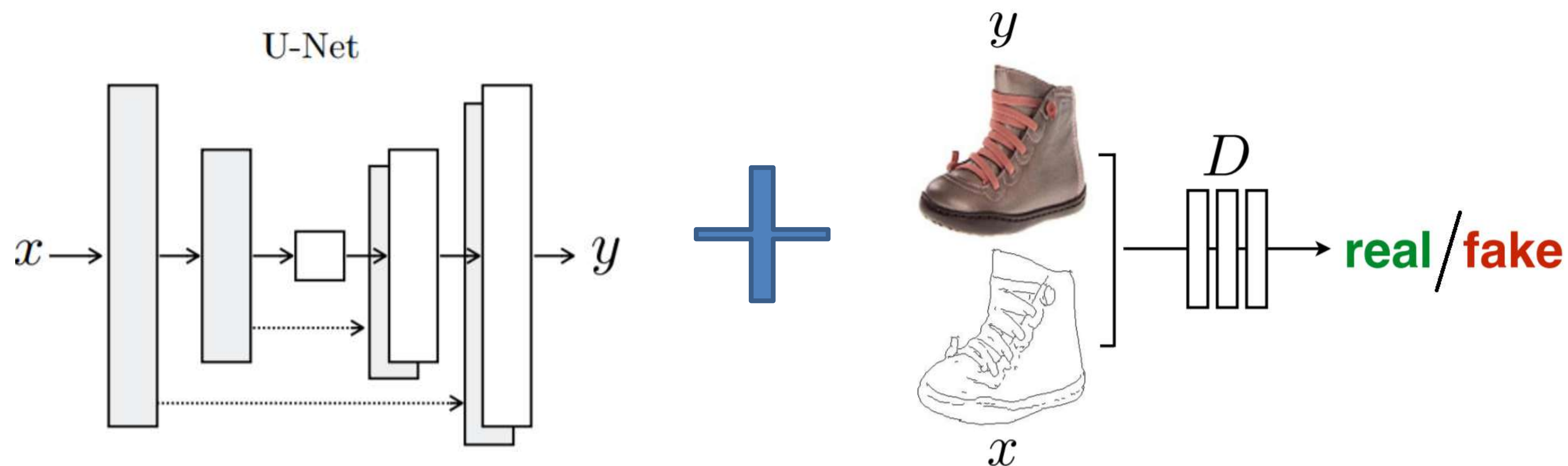


[Brock et al., 2018]



# Adversarial-функции потерь

- Дискриминатор является дифференцируемой функцией потерь!
  - Обучаемой вместе с генератором
- Пример: pix2pix [Isola et al., 2017]

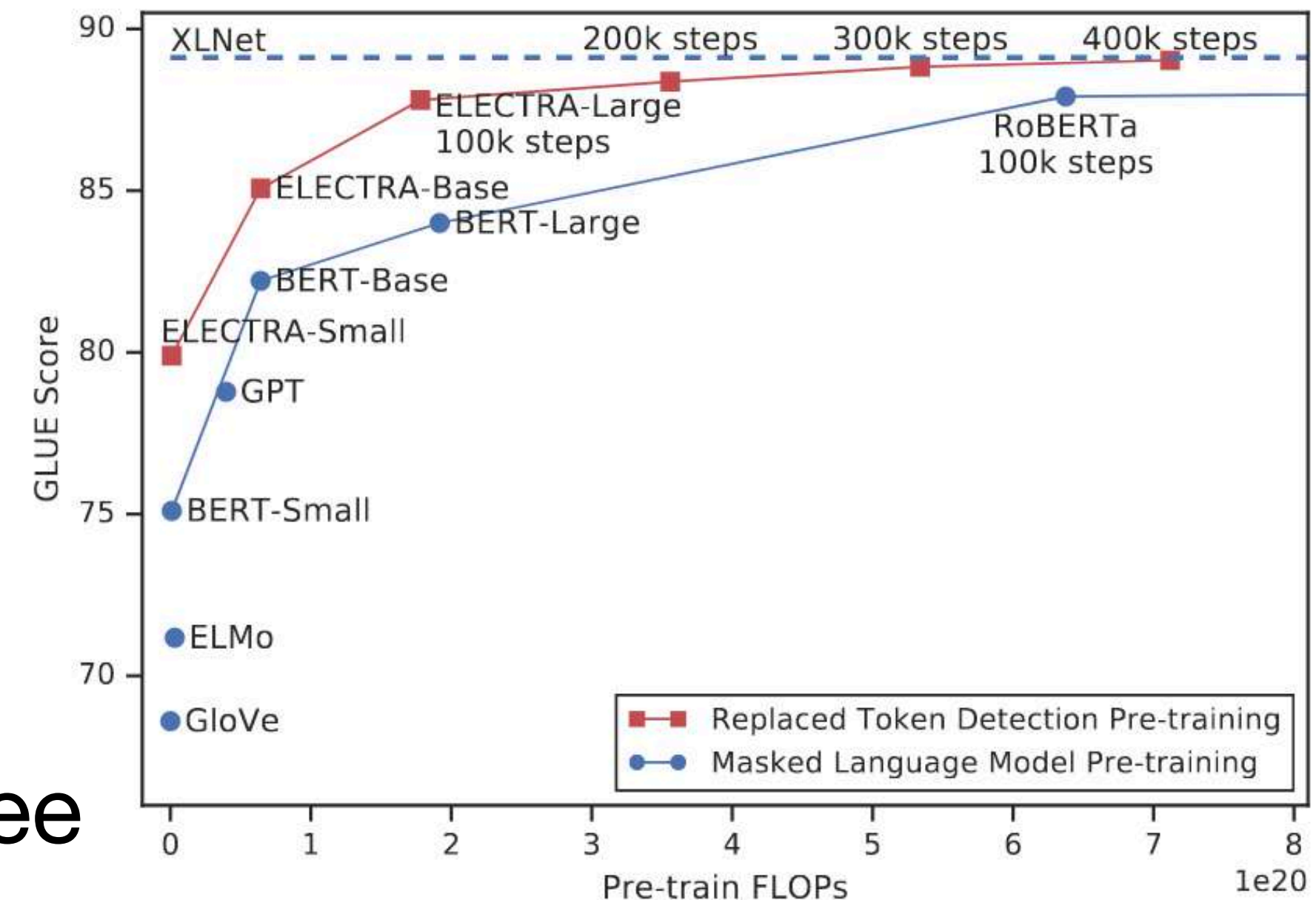
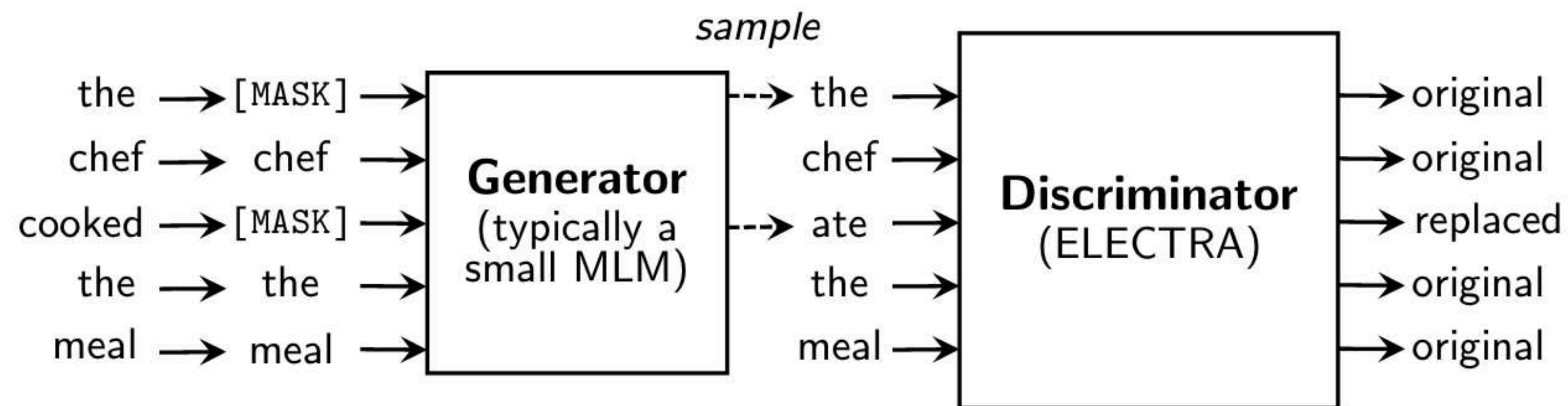


- Используется там, где нет хороших функций потерь
  - Изображения, звук, RL
  - Для текстов – делают, но все сложно



# Bonus: ELECTRA

- Сигнал из нейросети бывает заметно богаче сигнала из данных
- Можно использовать для более эффективного обучения текстовых моделей



- В итоге модель сходится существенно быстрее
- Важно: это не **GAN** (генератор не обманывает)

[Clark et al., 2020]



# Заключение

- Adversarial examples – фундаментальный феномен ML
  - Можно ли бороться?
  - Является ли ключом к пониманию сетей?
  - Adversarial для людей
- Domain adaptation – способ использовать больше данных с метками из другого домена
  - Очень важно на практике
- Adversarial-функции потерь – способ задавать функции потерь для сложных данных

Chihuahua or Muffin?

