

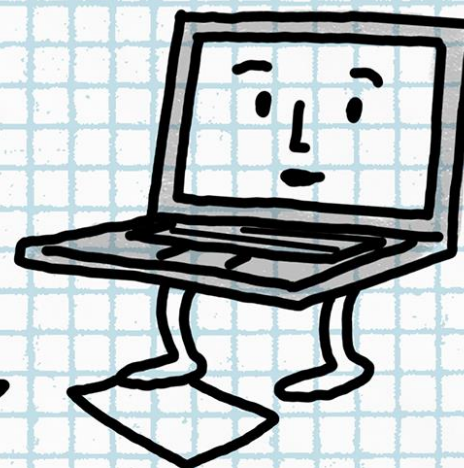
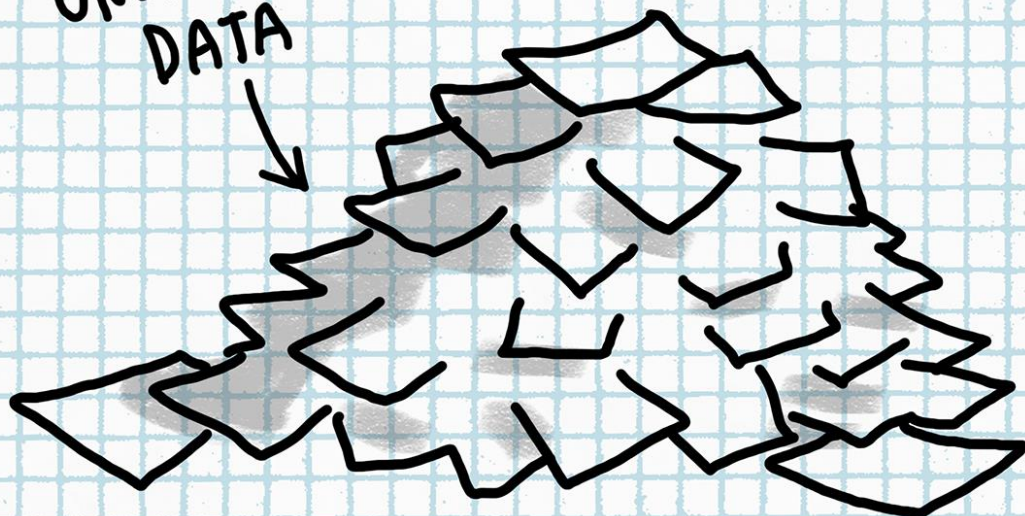
# SELF-SUPERVISED LEARNING

Федорова Анна, БПМИ-191

# ПРОБЛЕМА РАЗМЕТКИ ДАННЫХ

## SELF-SUPERVISED LEARNING

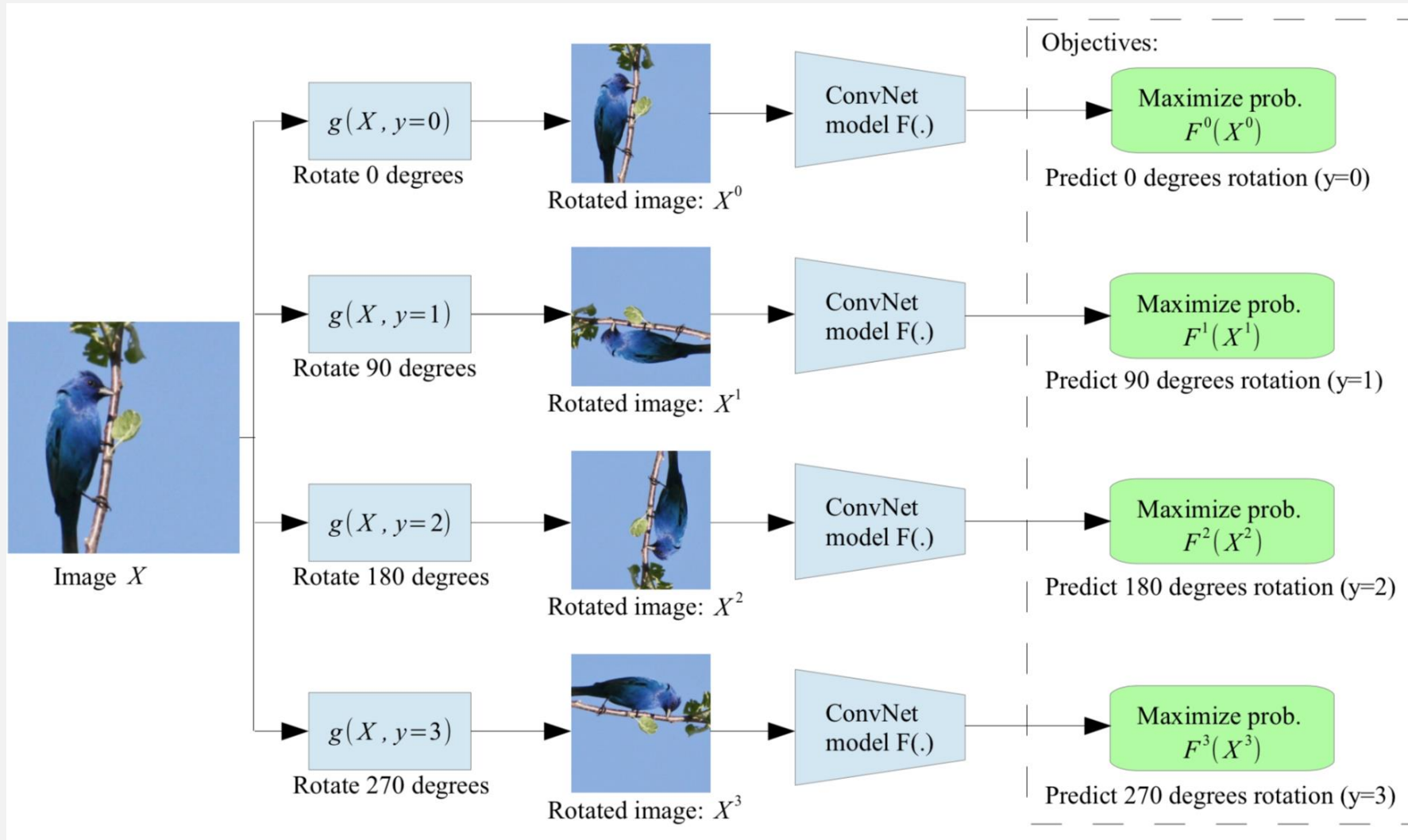
UNLABELED  
DATA



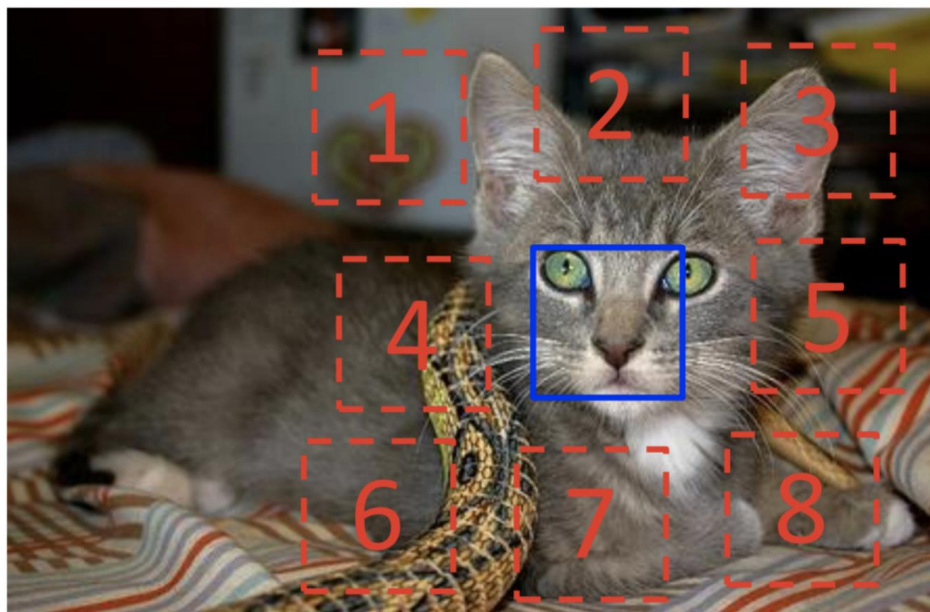




# ПОВОРОТЫ

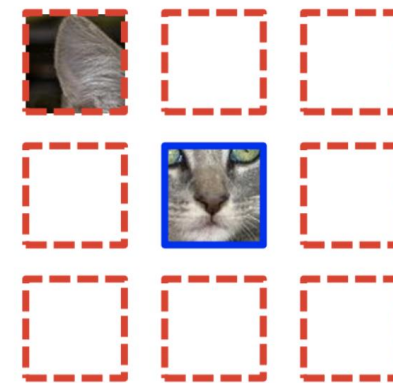


# ВЗАИМОРАСПОЛОЖЕНИЕ ФРАГМЕНТОВ



$$X = (\text{cat face}, \text{cat ear}); Y = 3$$

Example:



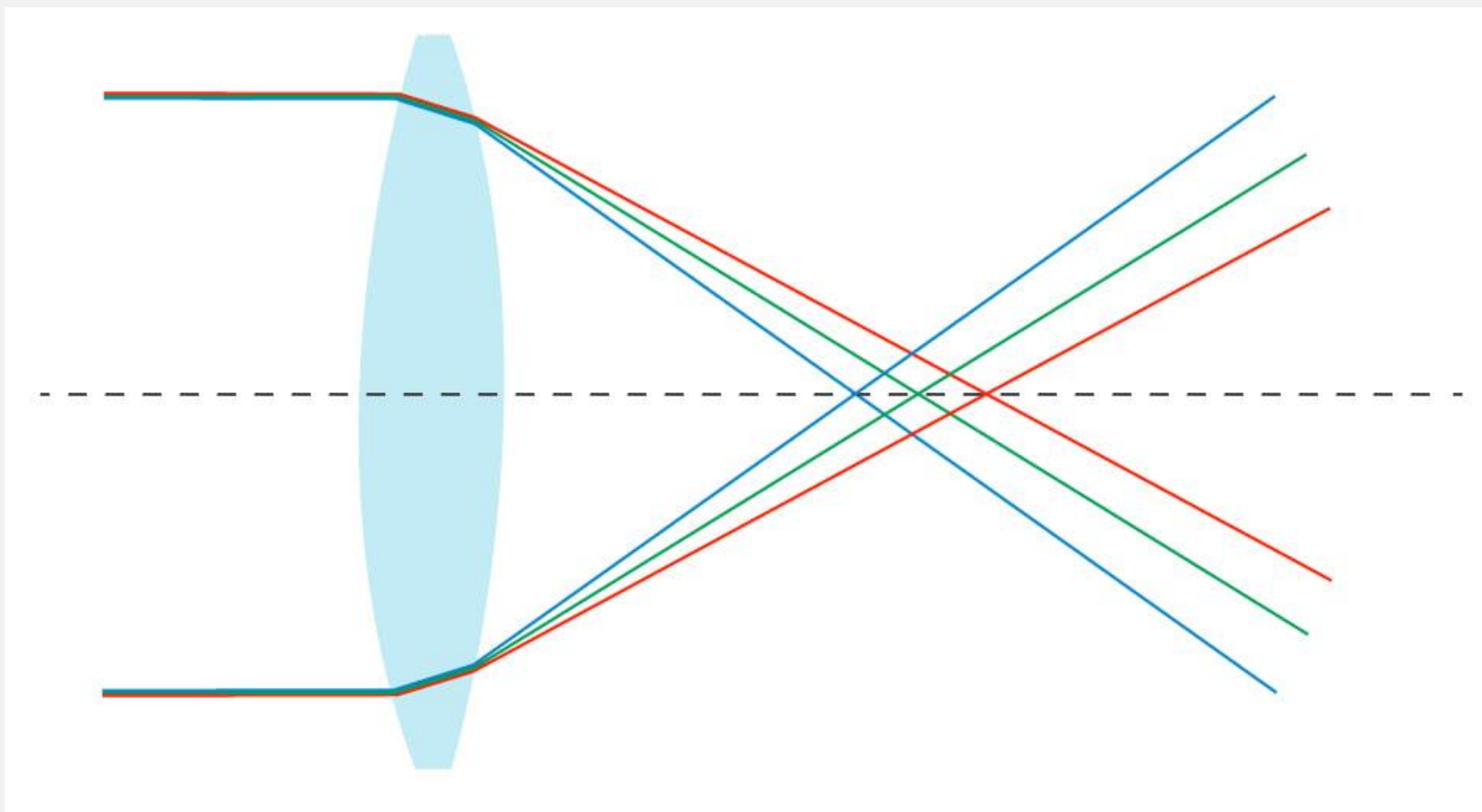
Question 1:



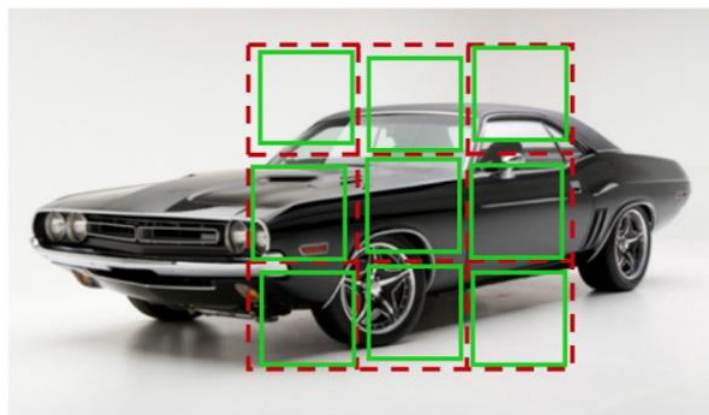
Question 2:



# ПОДВОХ ЦВЕТОВЫХ КАНАЛОВ



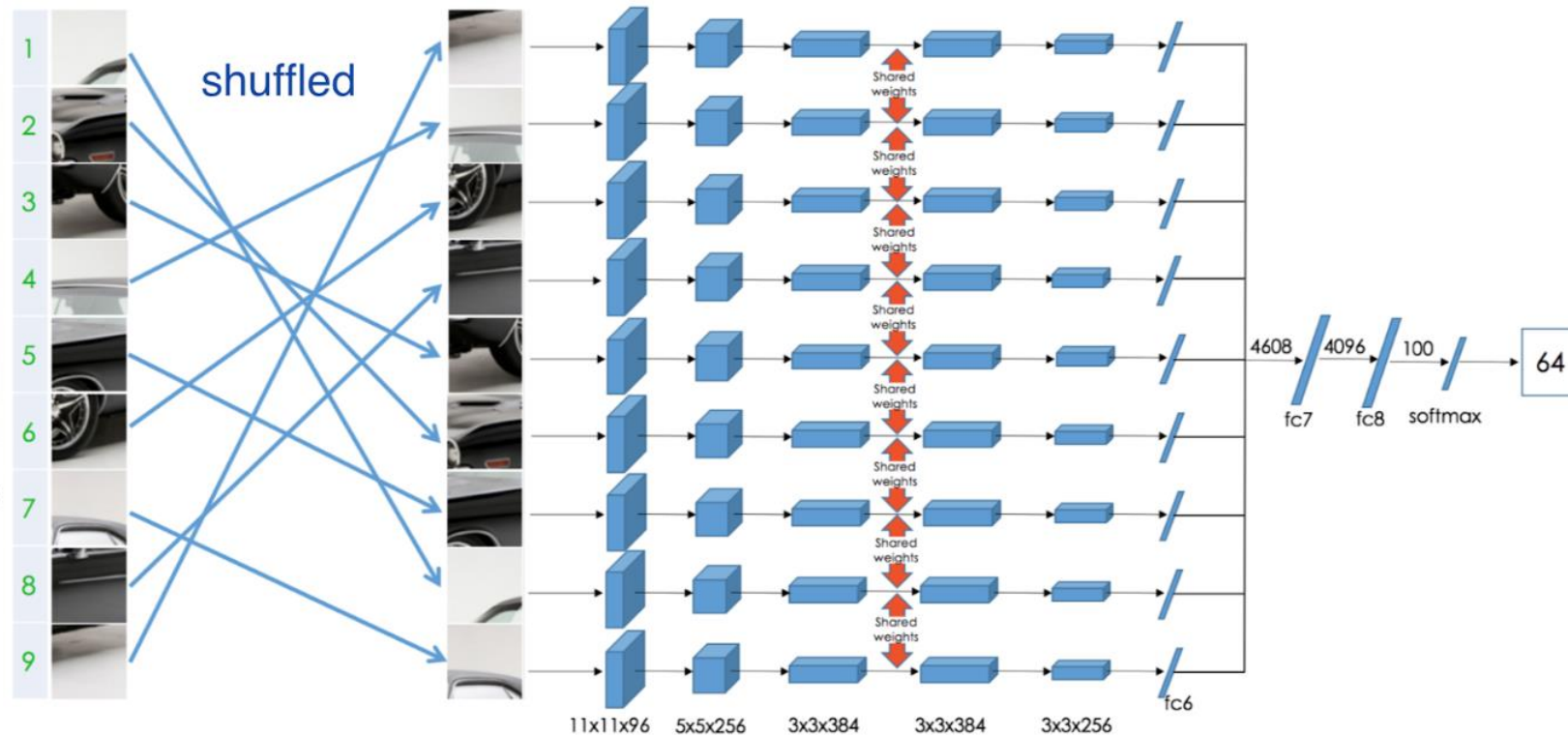
# ПАЗЛЫ



Permutation Set

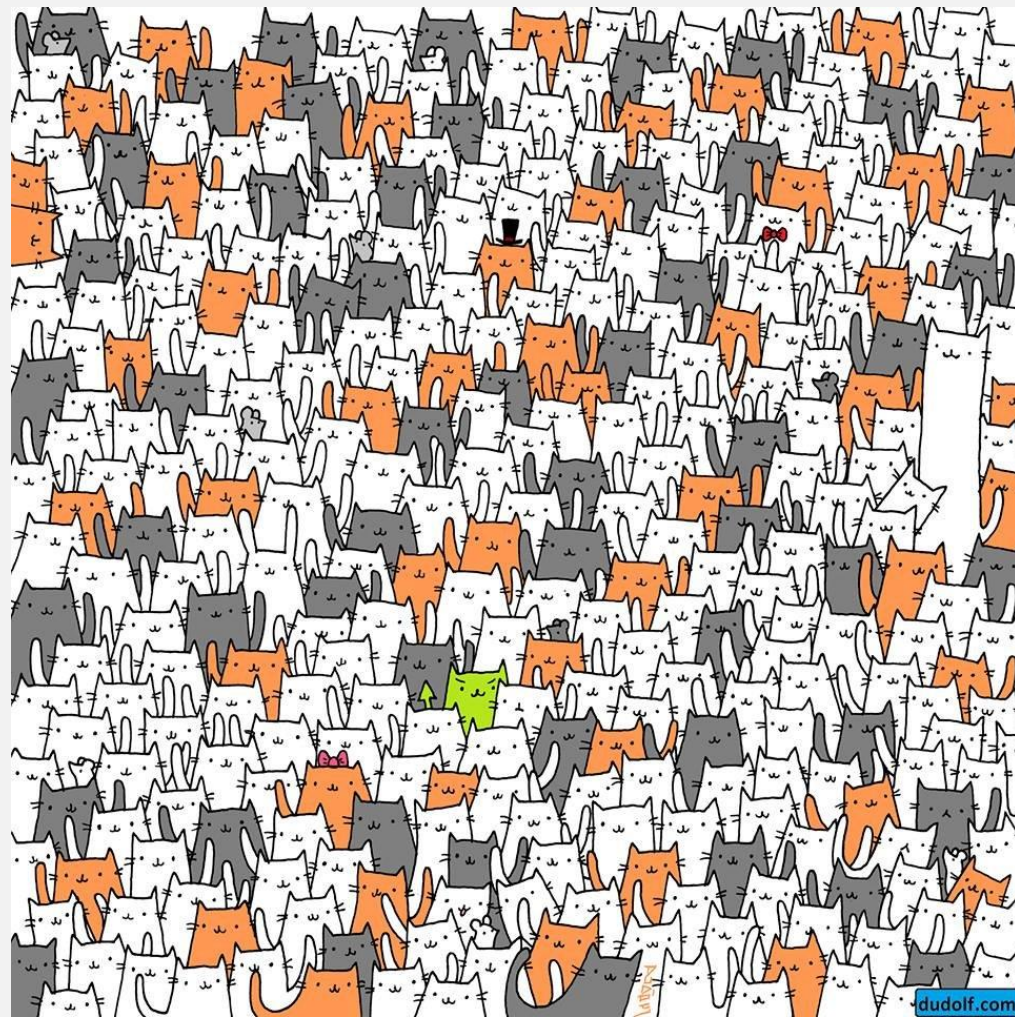
index	permutation
64	9,4,6,8,3,2,5,1,7

Reorder patches according to the selected permutation



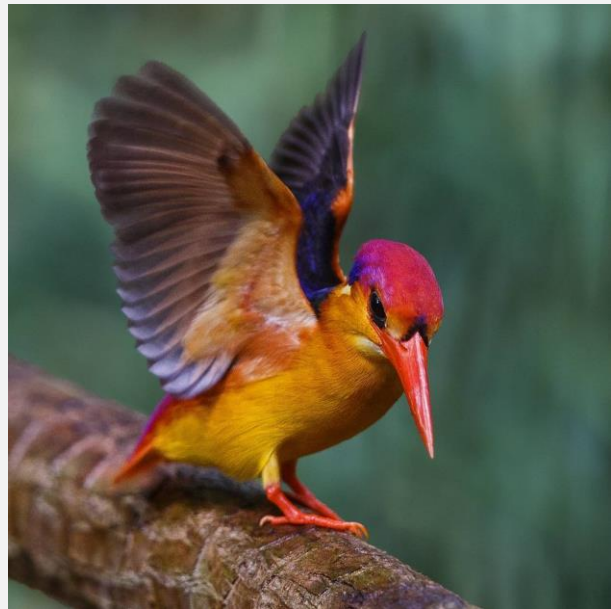
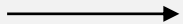
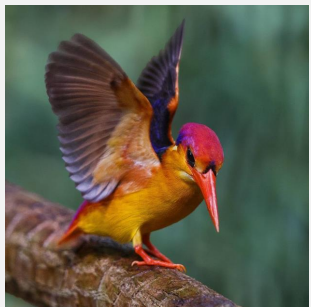


# ПОДСЧЕТ ВИЗУАЛЬНЫХ ПРИМИТИВОВ

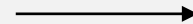




# ПРЕОБРАЗОВАНИЕ ИЗОБРАЖЕНИЙ



Масштабирование



Плитка 2x2

## ФОРМУЛЫ ДЛЯ ОБУЧЕНИЯ

$$\phi(x) = \phi(D \circ x) = \sum_{i=1}^4 \phi(T_i \circ x)$$

Функция, которую ищет модель

$D$  – обратно масштабирует изображение  
 $T_i$  – возвращает  $i$ -ый элемент плитки  $2 \times 2$

Функции, обращающие наши преобразования

## ФУНКЦИЯ ОШИБКИ

$$L_{feat} = \|\phi(D \circ x) - \sum_{i=1}^4 \phi(T_i \circ x)\|_2^2$$

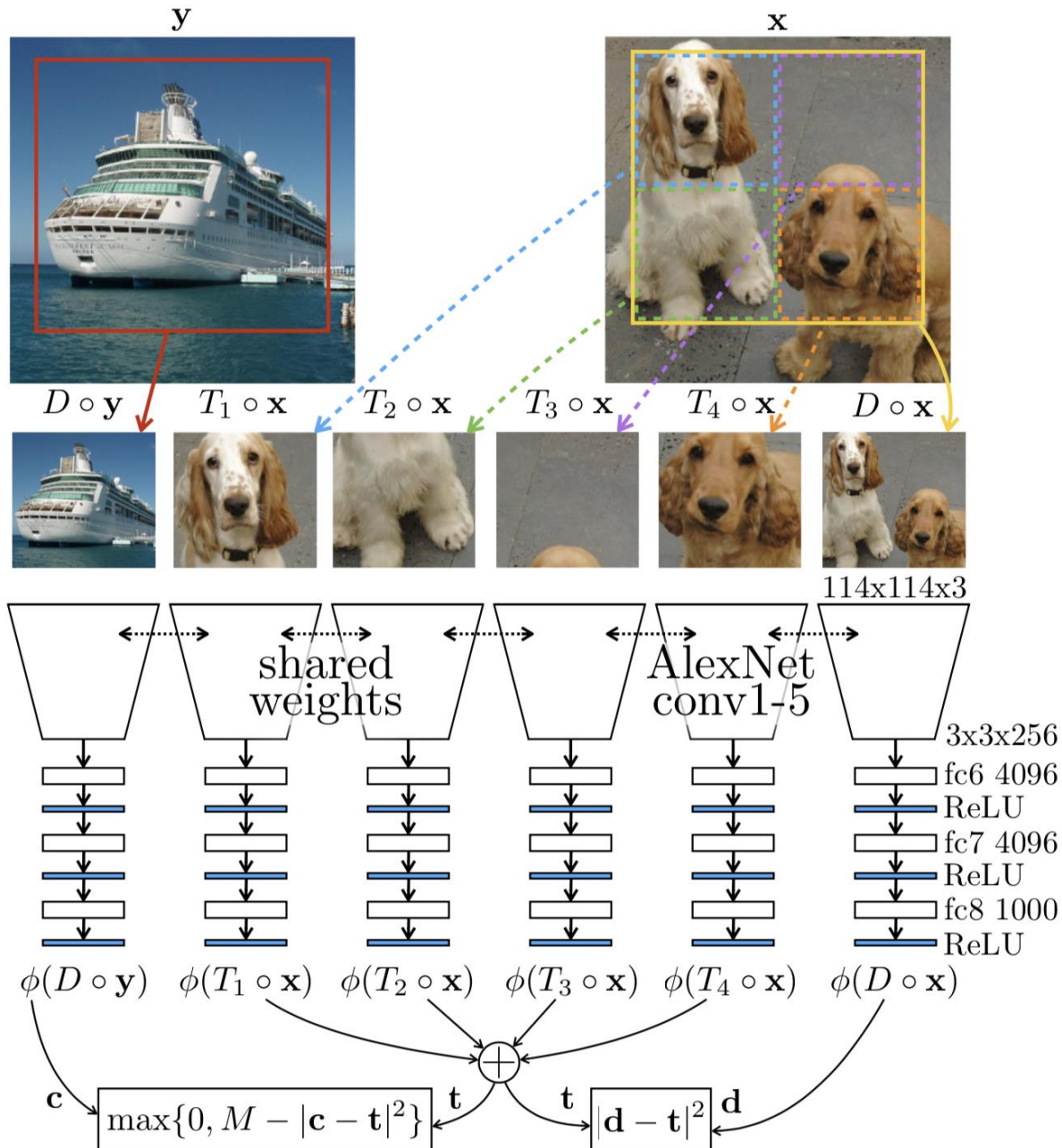
Функция ошибки

$$L_{diff} = \max(0, c - \|\phi(D \circ y) - \sum_{i=1}^4 \phi(T_i \circ x)\|_2^2)$$

Исключаем тождественный ноль

$$L = L_{feat} + L_{diff} = \|\phi(D \circ x) - \sum_{i=1}^4 \phi(T_i \circ x)\|_2^2 + \max(0, c - \|\phi(D \circ y) - \sum_{i=1}^4 \phi(T_i \circ x)\|_2^2)$$





## СХЕМА РАБОТЫ НЕЙРОННОЙ СЕТИ

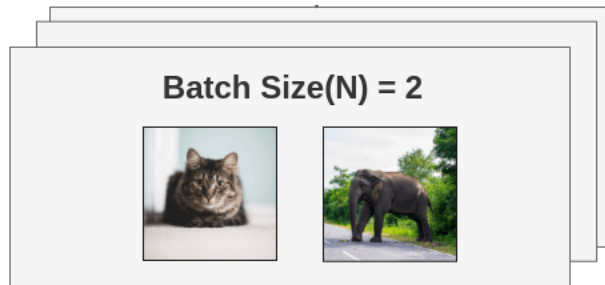
СПАСИБО ЗА ВНИМАНИЕ



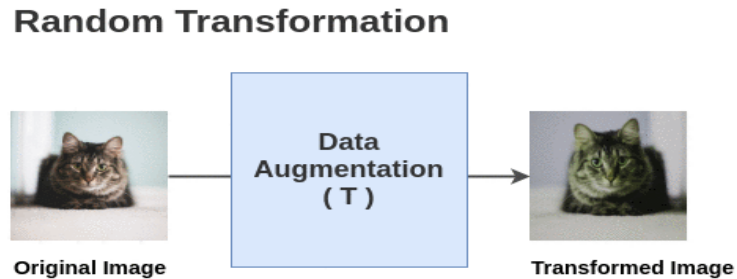
Шаг 0.1) Нам нужны примеры пар изображений, которые похожи, и которые отличаются



Шаг 0.2) Делим весь массив фотографий на батчи

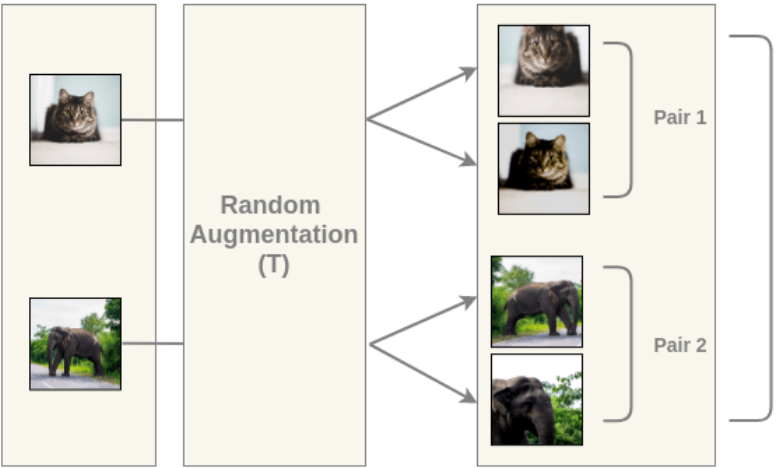


Шаг 1.1) Трансформируем изображения random (crop + flip + color jitter + grayscale) – обрезаем (берём фрагмент), поворачиваем, приводим к серому, смещения зеленого и пурпурного в сторону серого, случайное удаление 2 из 3 цветовых каналов

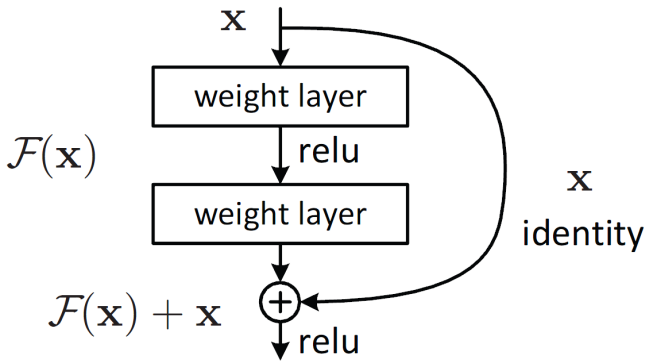




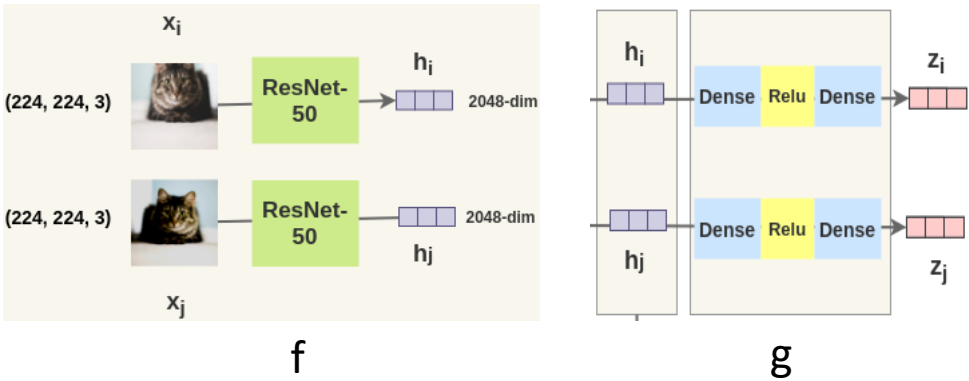
Шаг 1.2) Заменяем исходные изображения в батче их трансформированными копиями



Что такое ResNet или «остаточная сеть», объяснять не буду, вот красивая картинка:



Шаг 2) «Сжимаем» изображение, например с помощью сети ResNet-50 (можно использовать и другие)

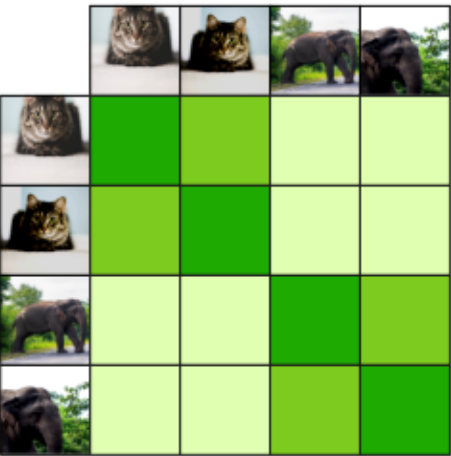


Запомните

$$\text{similarity}(x_i, x_j) = \text{cosine similarity}(z_i, z_j)$$

Шаг 3.0) Считаем Cosine Similarity

$$s_{i,j} = \frac{z_i^T z_j}{(\tau ||z_i|| ||z_j||)}$$



Напоминание:

$$\text{similarity}(x_i, x_j) = \text{cosine similarity} \left( \begin{bmatrix} z_i \end{bmatrix}, \begin{bmatrix} z_j \end{bmatrix} \right)$$

Шаг 3.1) Для каждой (!) пары изображений считаем Softmax



Pair 1

Softmax =

$$\frac{e^{\text{similarity}(\text{cat}, \text{cat})}}{e^{\text{similarity}(\text{cat}, \text{cat})} + e^{\text{similarity}(\text{cat}, \text{elephant})} + e^{\text{similarity}(\text{cat}, \text{elephant})}}$$

Шаг 4.1) Посчитаем  $I(\text{image1}, \text{image2})$  по формуле ниже

$$l(i, j) = -\log \frac{\exp(s_{i,j})}{\sum_{k=1}^{2N} l_{[k \neq i]} \exp(s_{i,k})}$$

$$I(\text{cat1}, \text{cat2}) = -\log \left( \frac{e^{\text{similarity}(\text{cat1}, \text{cat2})}}{e^{\text{similarity}(\text{cat1}, \text{cat2})} + e^{\text{similarity}(\text{cat1}, \text{elephant1})} + e^{\text{similarity}(\text{cat1}, \text{elephant2})}} \right)$$

Шаг 4.2) Также посчитаем  $I(\text{image2}, \text{image1})$

Interchanged

$$I(\text{cat2}, \text{cat1}) = -\log \left( \frac{e^{\text{similarity}(\text{cat2}, \text{cat1})}}{e^{\text{similarity}(\text{cat2}, \text{cat1})} + e^{\text{similarity}(\text{cat2}, \text{elephant1})} + e^{\text{similarity}(\text{cat2}, \text{elephant2})}} \right)$$



Шаг 4.3) Теперь можем посчитать L для нашего батча

$$L = \frac{1}{2N} \sum_{k=1}^N [l(2k-1, 2k) + l(2k, 2k-1)]$$

Pair 1 Loss (k=1)                      Pair 2 Loss (k=2)

$$L = \frac{[l(\text{cat}_1, \text{cat}_2) + l(\text{cat}_2, \text{cat}_1)] + [l(\text{ele}_1, \text{ele}_2) + l(\text{ele}_2, \text{ele}_1)]}{2 * 2}$$

Шаг 4.4) Обновляем параметры сети и возвращаемся к шагу 1.1 с новым батчем

---

**Algorithm 1** SimCLR's main learning algorithm.

---

**input:** batch size  $N$ , constant  $\tau$ , structure of  $f, g, \mathcal{T}$ .  
**for** sampled minibatch  $\{\mathbf{x}_k\}_{k=1}^N$  **do**  
    **for all**  $k \in \{1, \dots, N\}$  **do**  
        draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$   
        # the first augmentation  
         $\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$   
         $\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$  # representation  
         $\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$  # projection  
        # the second augmentation  
         $\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$   
         $\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$  # representation  
         $\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$  # projection  
    **end for**  
    **for all**  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  **do**  
         $s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$  # pairwise similarity  
    **end for**  
    **define**  $\ell(i, j)$  **as**  $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$   
     $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$   
    update networks  $f$  and  $g$  to minimize  $\mathcal{L}$   
**end for**  
**return** encoder network  $f(\cdot)$ , and throw away  $g(\cdot)$

---

P.S.S. Сравнение с другими методами (также из статьи)

