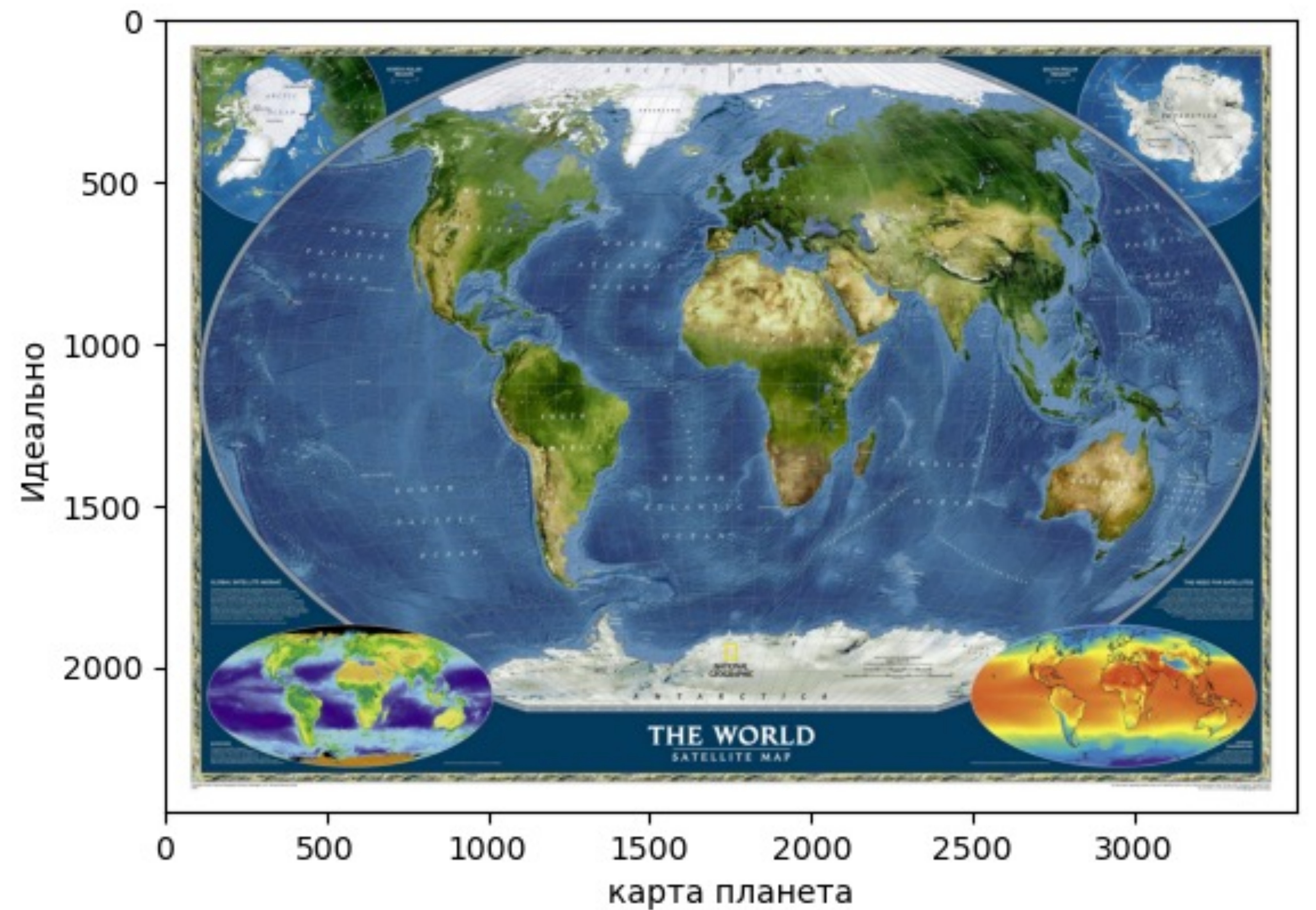
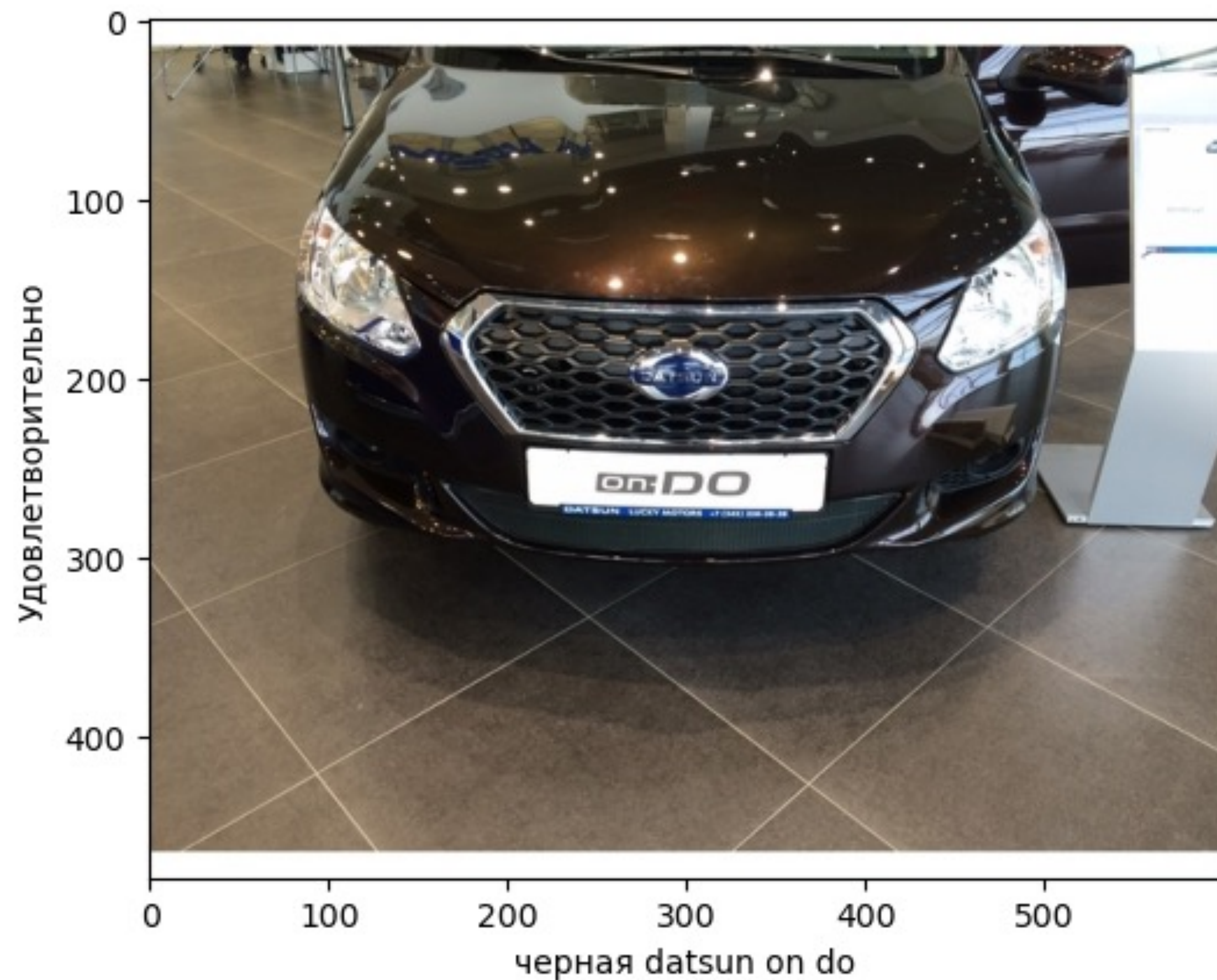


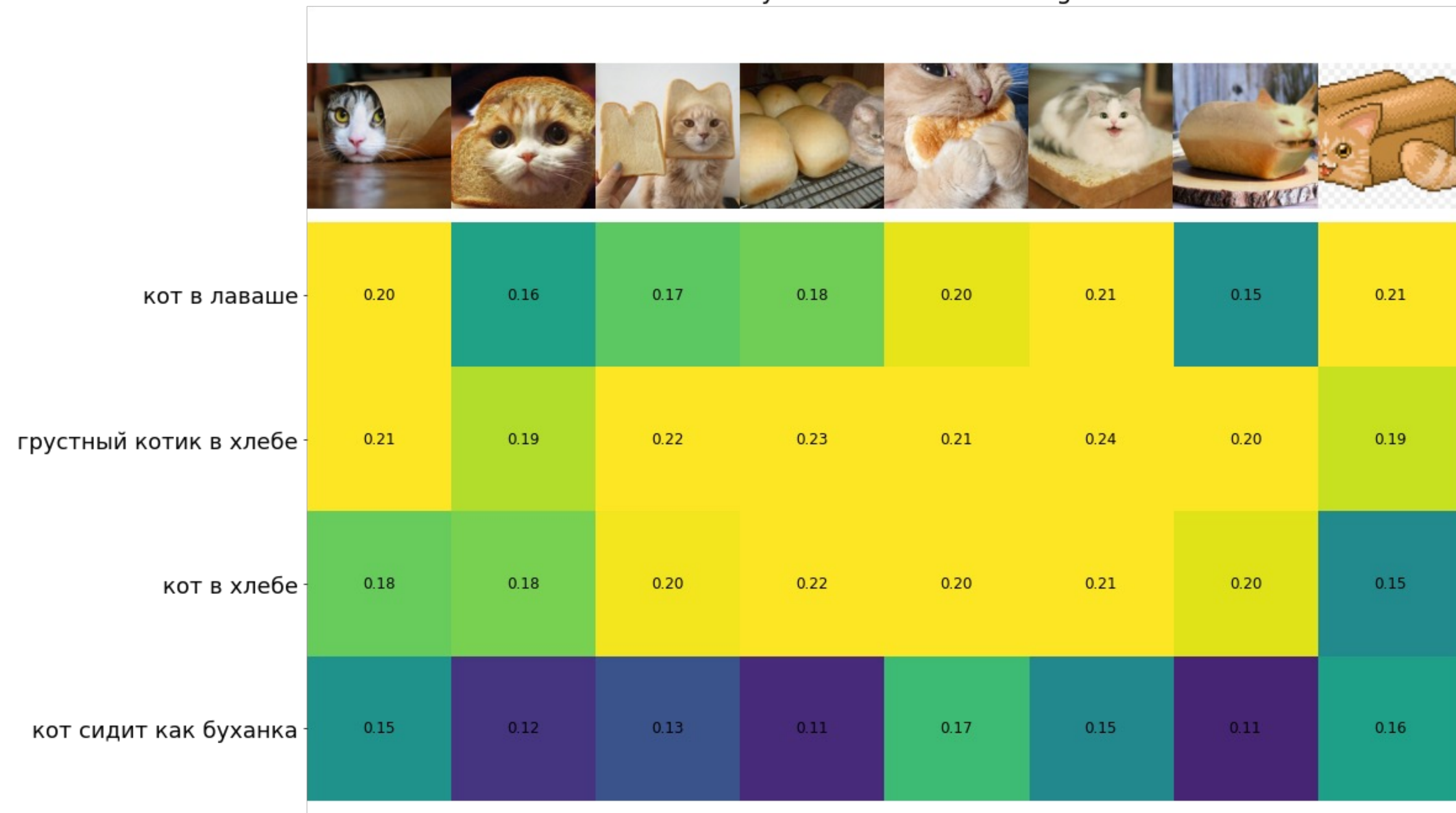
CLIP

Image-text matching



CLIP

Cosine similarity between text and image features



<https://habr.com/ru/companies/sberdevices/articles/564440/>

CLIP

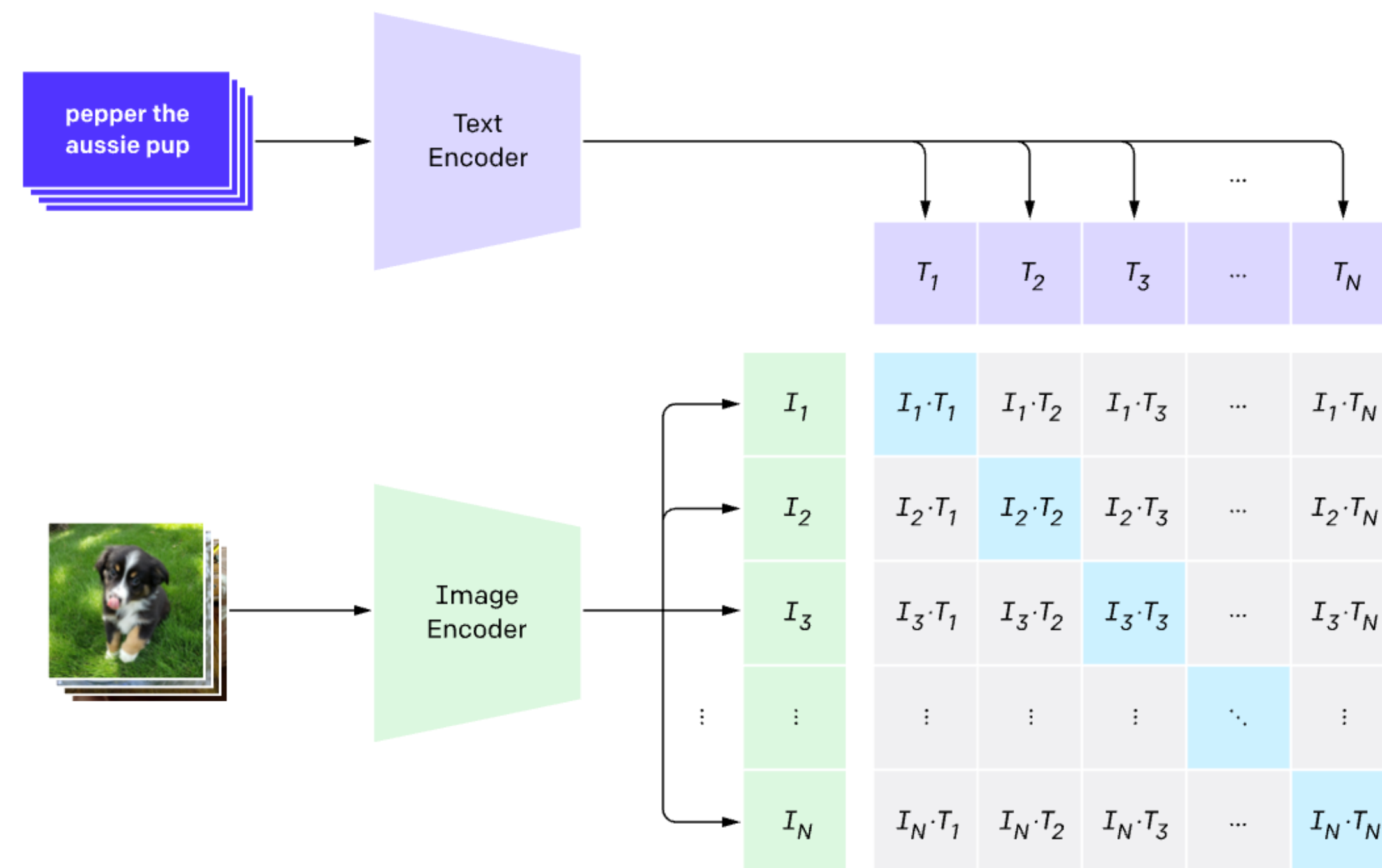
CLIP — это модель, состоящая из двух частей (или нейронных сетей):

1. **Image Encoder** — часть для кодирования изображений и перевода их в общее векторное пространство. В качестве архитектуры в оригинальной работе берутся ResNet разных размеров и Visual Transformer — тоже разных размеров.
2. **Text Encoder** — часть для кодирования текстов и перевода их в общее векторное пространство. В качестве архитектуры в оригинальной работе используется небольшой текстовый Transformer.

CLIP

CLIP выучивает мультимодальное пространство путём совместного обучения Image Encoder и Text Encoder, чтобы максимизировать косинусную близость эмбедингов изображения и текста реальных пар и минимизировать косинусную близость эмбедингов неправильных пар. Авторы оптимизируют симметричную кросс-энтропийную функцию потерь над полученными отношениями близости.

1. Contrastive pre-training



RUCLiP

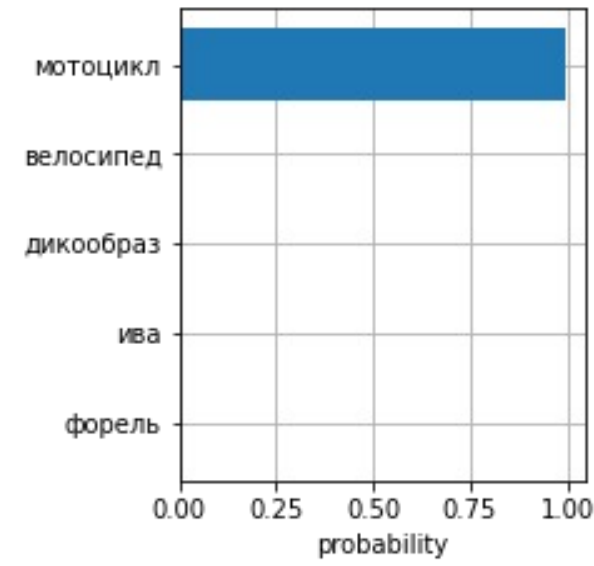
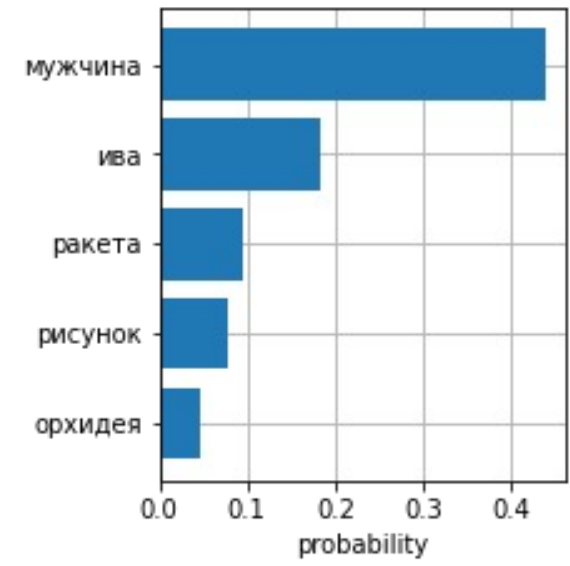
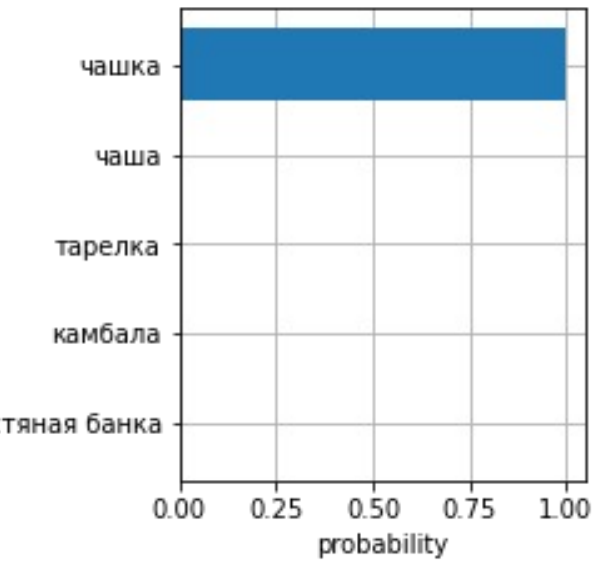
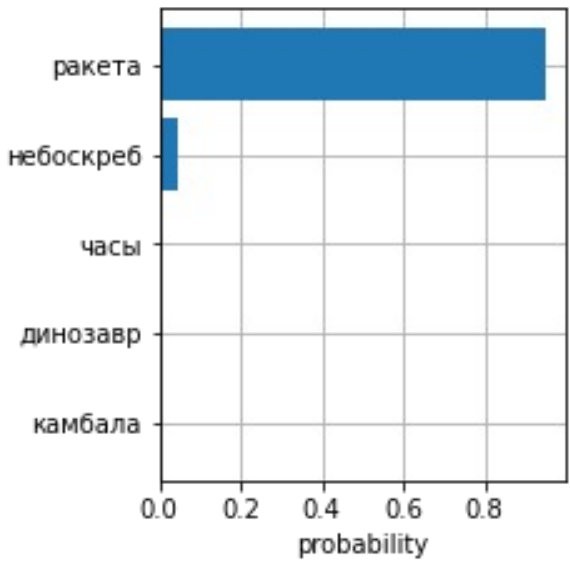
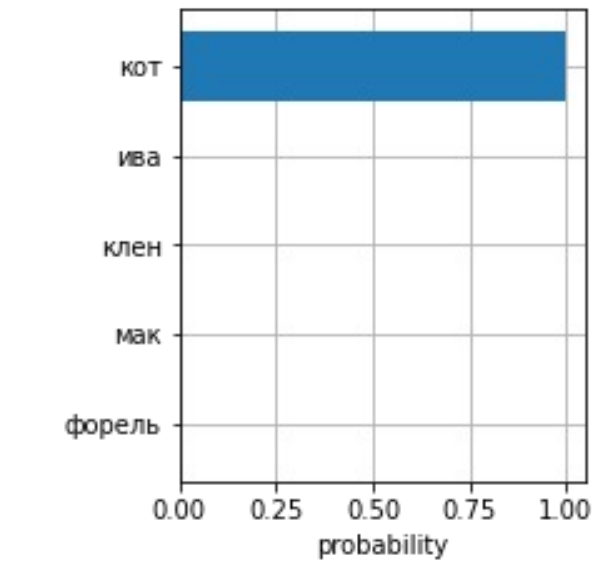
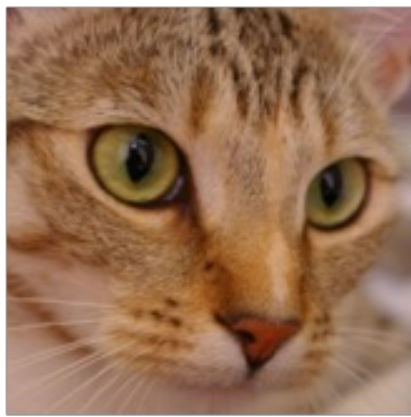
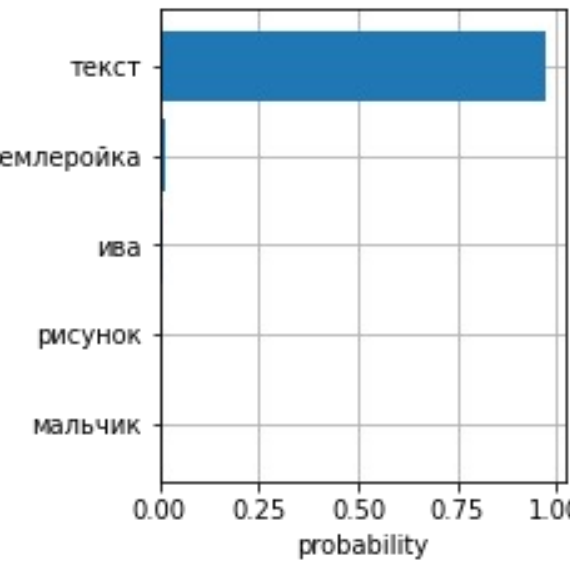
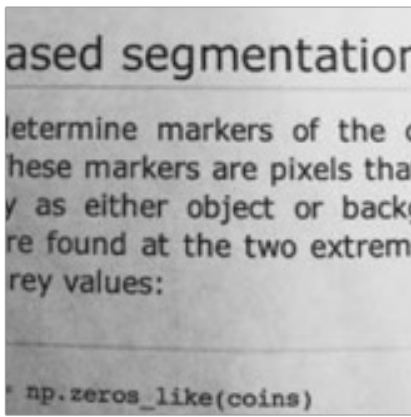
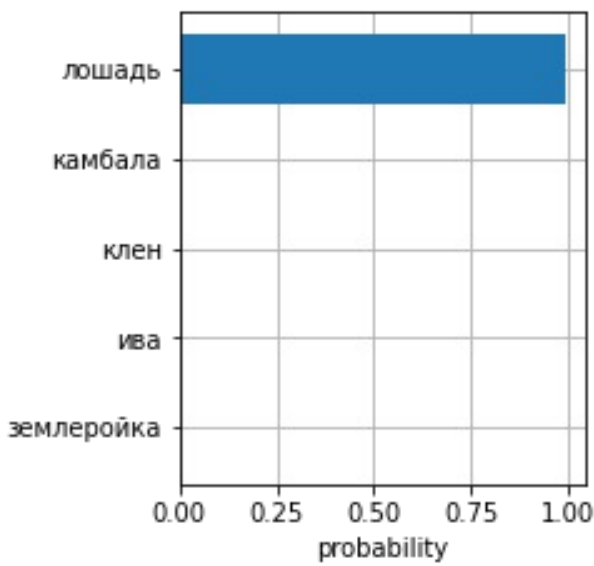
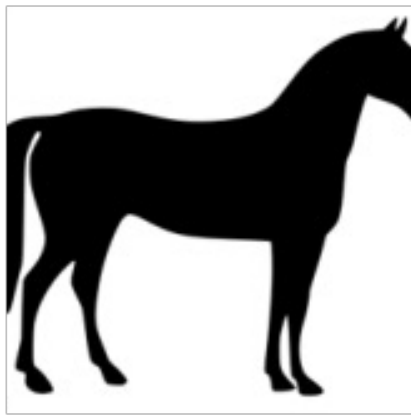
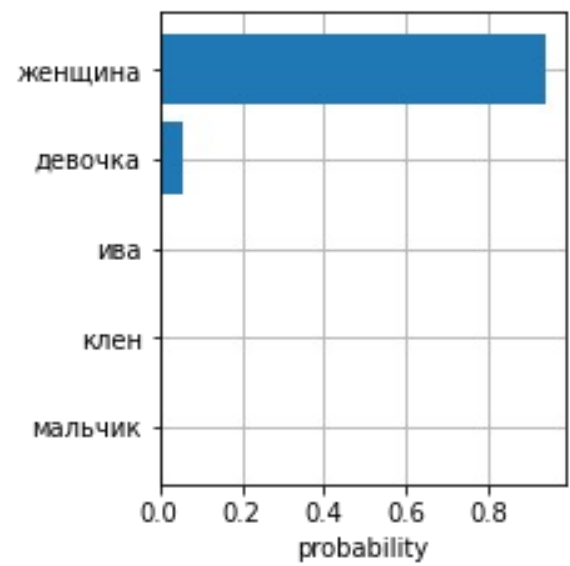
Модель CLiP дообучена на русскоязычных данных - RuCLiP.

Дообучение модели для русского языка происходило на собранных нами датасетах.

Вот некоторые из них:

- ImageNet — переведённый на русский язык;
- Flickr — картинки с русскими описаниями с фотостока;
- Ru-wiki — часть картинок из русской Википедии с описаниями.

RUCLIP



CLIP Demo

<https://sachinruk.github.io/blog/2021-03-07-clip.html>