

Meeting 6.1 - Extra tasks - HDF5

Pirms šī sagatavot numpy mmap dataset!

Dataset:

<https://www.kaggle.com/hsankesara/flickr-image-dataset>

1. Implementēt data pre-processor script, izmantojot HDF5 formātu un pēc tam uztaisīt, ka torch.data.utils.Dataset to izmanto, lai

getitem ielādētu no šī faila pointera

<https://docs.h5py.org/en/stable/>

Šim formātam nav nepieciešams blakus vēl viens fails, jo labels var glabāt kā dictionary key

2. Implementēt data pre-processor + dataset script, izmantojot cupy memory mapping

<https://github.com/cupy/cupy/issues/3431>

<https://stackoverflow.com/questions/57752516/how-to-use-cuda-pinned-zero-copy-memory-for-a-memory-mapped-file>

3. Pa tiešo file based-dataset (https://pytorch.org/tutorials/beginner/basics/data_tutorial.html)

```
1 import os
2 import pandas as pd
3 from torchvision.io import read_image
4
5 class CustomImageDataset(Dataset):
6     def __init__(self, annotations_file, img_dir, transform=None, target_transform=None):
7         self.img_labels = pd.read_csv(annotations_file)
8         self.img_dir = img_dir
9         self.transform = transform
10        self.target_transform = target_transform
11
12    def __len__(self):
13        return len(self.img_labels)
14
15    def __getitem__(self, idx):
16        img_path = os.path.join(self.img_dir, self.img_labels.iloc[idx, 0])
17        image = read_image(img_path)
18        label = self.img_labels.iloc[idx, 1]
19        if self.transform:
20            image = self.transform(image)
21        if self.target_transform:
22            label = self.target_transform(label)
23        return image, label
```

4. Notestēt apmācību uz 10 epochs, salīdzinot izpildes ātrumu ar `time.time()` mmap, hdf5, cupy mmap un file-based