

In this exercise you will experiment with semantic segmentation, image captioning and image-text retrieval.

## Submitting

Create a report as PDF (max 2 pages). Upload your report and code as zip to ILIAS, the zip should usually be less than 10MB in size if you adhere to the following rules:

Please do not upload model checkpoints (`.pth` files), tensorboard files (`events.out.tfevents...`) or parts of the dataset. Do not upload the jupyter notebooks with images in them (`.ipynb`), as they will be large. If you want to upload images, add them to the report.

## Setup

### Pool computers

If you want to train on the computers in the PC pool, make sure you applied for an increase in disk space quota. Otherwise, the space in your home folder will not be enough.

How to use the pool:

- Connect to the login server of the technical faculty: `ssh username@login.informatik.uni-freiburg.de`
- Do not run anything on the main server. Connect to one of the pool computer there instead. Preferably use `tfpool25-46` or `tfpool51-63` as they provide the best GPUs. `ssh tfpoolXX` - replace XX with a 2 digit number of the computer you want to connect. Before running make sure no one else is using the selected computer and the GPU is free.
- See GPU load: `nvidia-smi`
- See logged in accounts: `who`
- View running processes: `top`

Start a screen session to keep your code running in the background:

- Start the session  
`screen -S t01`  
`# run your experiment`
- Detach from screen using: `ctrl+a+d`
- Login back into screen:  
`screen -ls`  
`screen -r t01`
- Write down on which computer you started your screen session.
- `ctrl+esc` to enter copy mode if you need to scroll up, navigate with arrow/page keys, copy output etc

### Installing python using miniconda

We recommend installing python inside a conda environment, since this way the installation will be independent of the python installed on the system. Below are the installation instructions for Ubuntu, on other operating systems the workflow is similar.

- Download the latest Miniconda3 Linux 64-bit installer: <https://docs.conda.io/en/latest/miniconda.html> Usually the 64-bit version is correct. Use `wget [link]` to download a file in the terminal.
- Open the terminal and `cd` into the folder containing `Miniconda3-latest-Linux-x86_64.sh`.
- Run `bash Miniconda3-latest-Linux-x86_64.sh` to install it into your home directory. When asked by the installer to initialize Miniconda3, answer yes. This way, your `~/.bashrc` file will be modified so you have conda available at startup.

- Close and reopen your terminal window to make sure the changes take effect.
- Test your installation by running `conda --version` to see the its version.

### Setup conda environment

Next we will set up conda with a new environment named `lab` which uses Python 3.9. Run the following commands in your terminal (macOS, Linux) or miniconda prompt (Windows):

```
conda --version
conda update -n base -c defaults conda -y
conda create --name lab python=3.9 -y
conda activate lab
conda env list
```

The last command should show you all installed environments (including `base` and the activated `lab` environment). This way, you have a separate environment to install everything for this course and avoid possible interference with/from other python projects.

**Note:** Whenever you start the miniconda prompt/terminal for this course, run `conda activate lab` to switch to the correct environment.

### Setup the code

1. Make sure your conda environment is activated and you are using python 3.9.
2. Download the code from ILIAS and cd into the folder.
3. Install PyTorch with GPU support:  
`conda install pytorch torchvision pytorch-cuda=11.8 -c pytorch -c nvidia -y`
4. Install packages: `pip install -U -r requirements.txt`

### Setup the data and model checkpoints

#### On the pool

If you are working on the pool you should already have access to the data and checkpoint. Run `ls /project/dl2023s/lmbcvtst/public` to check if you can see the folder.

#### Locally

If you want to run the exercise on your own hardware you will need to setup the data and model checkpoint locally.

Uncomment the last two lines in `paths.py` to set the code to load from local paths **or** set the environment variables to the local paths.

```
# assuming you are in the root of the repository
# create data folder
mkdir data

# download and extract dataset into the data folder
wget http://host.robots.ox.ac.uk/pascal/VOC/voc2012/VOCtrainval_11-May-2012.tar
tar -xf VOCtrainval_11-May-2012.tar -C data
```

```
# copy the captions into the data folder
cp -r datasets/Captions data/VOCdevkit/VOC2012

# download the DINO checkpoint for the segmentation task
mkdir ckpt
wget https://dl.fbaipublicfiles.com/dino/dino_deitsmall16_pretrain/dino_deitsmall16_pretrain.pth -P ckpt

# download the original BLIP checkpoint for the image-text tasks
wget https://storage.googleapis.com/sfr-vision-language-research/BLIP/models/model_base_capfilt_large.pth

# run the postprocessing script, this should create ckpt/blip_model_base.pth
python postprocess_ckpt.py

# optionally remove the downloaded files
rm VOCtrainval_11-May-2012.tar
rm model_base_capfilt_large.pth
```

## Additional resources

If you have difficulties understanding PyTorch, we suggest to check out some of the [tutorials](#) for PyTorch, [learn the basics](#) is a good tutorial to get started.

A very important concept is indexing of tensors, see [numpy indexing](#). The indexing in PyTorch works very similar to numpy.

We recommend using an IDE like VSCode or PyCharm for writing and debugging python code.

## Using browser-based tools

In this exercise, you will **optionally** use browser-based tools like jupyter notebook and tensorboard. The easy way to do this is to run the server and browser on the same machine (either directly sit in the pool or setup everything on your home computer).

However, if your code is running on a *different computer as the browser*, i.e. you are running the code on a pool machine but running the browser on your home computer) you will need to: **1)** start the respective server on the pool machine. **2)** establish a SSH tunnel from your home computer to the tf login node. **3)** establish a SSH tunnel from the tf login node to the pool machine where the server is running. The tunneling command looks like this:

```
ssh -v -N -L ${port}:localhost:${port} "${user}@${targethost}"
```

You might need to change the port of your server application if the port is already in use by another person. Note that everyone who can access the login node will be potentially able to see your server. Jupyter notebook protects this with an access token by default, but tensorboard does not.

## 1 Segmentation

Your task is to implement and train some transformer-based semantic segmentation models by completing the given code. Start from the `train_segmentation.py` script, which builds the data loaders, model, and implements a simple supervised training loop. The basic usage of the script is `python train_segmentation.py <decoder type>`.

The dataset you will use is PASCAL VOC [1], which contains 1464 training samples and 1449 validation samples with fully annotated segmentation masks for 21 semantic categories/classes.

The model definitions for this task are in the `models/segmentation` and `models/vit` folders:

- `models/segmentation/vit_encoder.py` contains the implementation of the vision transformer ViT, which you will use as the encoder of your networks. It inherits from the main definition in `models/vit/vit.py` and adapts it for segmentation. The parameters of the encoder are initialized with pre-trained DINO [2] weights.
- `models/segmentation/transformer_segmentation_head.py` contains the incomplete implementation of the transformer-based segmentation head presented in the lecture. You will have to complete this script.
- `models/segmentation/encoder_decoder.py` contains a wrapper to concatenate both parts of the network.

### 1.1 Linear head

A simple linear is already implemented using a convolutional layer (`torch.nn.Conv2d`) with kernel size 1. Your first task is to train a segmentation model with this head. Report a plot of the validation performance (mIoU) of your model over the epochs, together with a few qualitative examples from the validation set. For displaying the prediction of your networks, overlaid on the images, we provide the utility class `SegmentationDisplay` in `utils/misc.py`. (1 point)

### 1.2 Convolutional head

Based on the implementation of the linear head, you should be able to easily build and train a model with a convolutional head. The exact configuration of the head is up to you, we recommend keeping it single-layer. Report the architectural details of your convolutional head, plot its validation performance, and produce some qualitative results. How does this head perform compared to the linear one? How do you explain the performance comparison? (2 points)

### 1.3 Transformer head

Complete the implementation of this head in `models/segmentation/transformer_segmentation_head.py`. The transformer head should consist of a single transformer block with added class embeddings, as explained in the lecture. For reference, you can look at the implementation of the standard transformer block in `models/vit/vit.py`. As for the other tasks, plot the validation performance and report some qualitative examples. (3 points)

#### 1.3.1 The role of queries and keys

Modify the attention mechanism in the transformer head to use the same projection layer for both *queries* and *keys*. Do so by completing the class `AttentionSameQK`, which extends the default `Attention`. Train the model again using the command-line argument `--transformer-head-shared-qk`: how does it perform after this modification? What can you conclude from it? (2 points)

### 1.4 Putting it all together

To conclude, report the final validation results of all the model versions you trained, in form of a table. Write a few sentences to summarize the findings of your experiments. (2 points)

## 2 Image-Text

In this exercise you will experiment with the BLIP [3] model. The model is defined in the following folders:

- `models/bert/`: Text encoder and decoder.

- `models/vit/`: Vision encoder.
- `models/blip/`: The model which connects the vision and text parts. `BlipCaption` uses the vision encoder and text decoder for image captioning, `BlipRetrieval` uses the vision encoder and text encoder for image-text retrieval.

## 2.1 Image captioning

Your first task will be to complete the inference code to generate captions for the given VOC dataset.

**Todo:** Run `python eval_captioning.py --help` and familiarize yourself with this evaluation script.

**Todo:** In file `models/blip/blip_caption.py` complete the methods `generate` and `greedy_search`. Use the evaluation script to generate and evaluate captions for the VOC dataset. You should get about 41% BLEU score. (2 points)

**Optional:** Start a jupyter notebook server with the command `jupyter notebook --port 8888` and view the images, reference captions and generated captions in the notebook `captions_visualizer.ipynb`

**Todo:** In file `models/blip/blip_caption.py` complete the method `sampling`. There, Top-K sampling instead of greedy search is used to select the next token when decoding. Evaluate again with Top-K sampling, you should get a lower BLEU score of about 16% for temperature  $\tau = 1.0$  and about 24% for  $\tau = 0.7$ . Why do the results improve with lower temperature? (1 point)

**Todo:** Experiment with different prompts (the default prompt is “a picture of”) and experiment with different decoding parameters (Top-K with different K and temperature or greedy decoding). Can you improve the BLEU score? Try at least 3 new settings and note the hyperparameters, the prompt and the resulting BLEU score in a table for each setting. Add the table to your report. (2 points)

## 2.2 Image-text retrieval

Your second task will be to train and evaluate the retrieval head of the BLIP model.

**Todo:** Complete the forward pass in file `models/blip/blip_retrieval.py`. Run `python eval_retrieval.py` to verify your implementation with the provided checkpoint. You should get about 54% image-to-text R@1. (2 points)

**Todo:** Complete the loss computation in file `train_retrieval.py` function `train_epoch`. Run `python train_retrieval.py` to train the retrieval projection layers from scratch (i.e. from random initialization). You should get about 43% image-to-text R@1. (1 point)

**Optional:** start a tensorboard server `tensorboard --logdir outputs --port 6006` and watch the experiment in the browser.

**Todo:** Now, try finetuning the head instead with `--finetune`. Set learning rate to  $1e-5$ , weight decay to 0 and train for 3 epochs. What score do you get and how can you explain the difference to the score when training from scratch? (1 point)

**Optional:** Run the jupyter notebook `retrieval_search.ipynb` to show qualitative results for the model. Try different search queries. What do you observe?

**Todo:** Experiment with different hyperparameters in the random initialization setting (i.e. without finetuning). Try at least 3 new hyperparameter settings and note the hyperparameters and the resulting image-to-text R@1 score in a table for each setting. Can you improve over the baselines? Add the table to your

report. (1 point)

## References

- [1] *The PASCAL Visual Object Classes Project*. <http://host.robots.ox.ac.uk/pascal/VOC/>.
- [2] Mathilde Caron et al. “Emerging Properties in Self-Supervised Vision Transformers”. In: *Proceedings of the International Conference on Computer Vision (ICCV)*. 2021.
- [3] Junnan Li et al. *BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation*. 2022. arXiv: [2201.12086](https://arxiv.org/abs/2201.12086) [[cs.CV](#)].