	This is a self-contained notebook that shows how to download and run CLIP models, calculate the similarity between arbitrary image and text inputs, and perform zero-shot image classifications. The next convil install the clip package and its dependencies, and check if PyTorch 1.7.1 or later is installed. ! pip install ftfy regex tqdm ! pip install git+https://github.com/openai/CLIP.git
	Collecting ftfy Downloading ftfy-6.1.1-py3-none-any.whl (53 kB)
	Running command git clone -q https://github.com/openai/CLIP.git /tmp/pip-req-build-ri9zw6nw Requirement already satisfied: ftfy in /usr/local/lib/python3.7/dist-packages (from clip==1.0) (6.1.1) Requirement already satisfied: regex in /usr/local/lib/python3.7/dist-packages (from clip==1.0) (2019.12.20) Requirement already satisfied: tqdm in /usr/local/lib/python3.7/dist-packages (from clip==1.0) (4.63.0) Requirement already satisfied: torch in /usr/local/lib/python3.7/dist-packages (from clip==1.0) (1.10.0+cu111) Requirement already satisfied: torchvision in /usr/local/lib/python3.7/dist-packages (from clip==1.0) (0.11.1+cu111) Requirement already satisfied: wcwidth>=0.2.5 in /usr/local/lib/python3.7/dist-packages (from ftfy->clip==1.0) (0.2.5) Requirement already satisfied: typing-extensions in /usr/local/lib/python3.7/dist-packages (from torch->clip==1.0) (3.10.0.2) Requirement already satisfied: numpy in /usr/local/lib/python3.7/dist-packages (from torchvision->clip==1.0) (7.1.2)
2]:	Building wheels for collected packages: clip Building wheel for clip (setup.py) done Created wheel for clip: filename=clip-1.0-py3-none-any.whl size=1369221 sha256=a03df40b14e04b8d01a72bb6ec53d12a4c9fbb1d4a2b39bda95697569c318e9e Stored in directory: /tmp/pip-ephem-wheel-cache-5bqsmcp3/wheels/fd/b9/c3/5b4470e35ed76e174bff77c92f91da82098d5e35fd5bc8cdac Successfully built clip Installing collected packages: clip Successfully installed clip-1.0 import numpy as np import torch
	<pre>print("Torch version:", torchversion) torch_version = torchversionsplit(".") assert (int(torch_version[0]) == 1 and int(torch_version[1]) >= 7) or int(torch_version[0]) > 1, "PyTorch 1.7.1 or later is required" Torch version: 1.10.0+cu111 Loading the model</pre>
3]: 3]:	<pre>clip.available_models() will list the names of available CLIP models. import clip clip.available_models() ['RN50', 'RN101',</pre>
4]:	'RN50x16', 'RN50x64', 'ViT-B/32', 'ViT-B/16', 'ViT-L/14'] model, preprocess = clip.load("ViT-B/32") model.cuda().eval() input_resolution = model.visual.input_resolution
	<pre>context_length = model.context_length vocab_size = model.vocab_size print("Model parameters:", f"{np.sum([int(np.prod(p.shape)) for p in model.parameters()]):,}") print("Input resolution:", input_resolution) print("Context length:", context_length) print("Vocab size:", vocab_size) 100% 338M/338M [00:05<00:00, 59.9MiB/s]</pre>
	Model parameters: 151,277,313 Input resolution: 224 Context length: 77 Vocab size: 49408 Image Preprocessing We resize the input images and center-crop them to conform with the image resolution that the model expects. Before doing so, we will normalize the pixel intensity using the dataset mean and standard deviation.
5]:	The second return value from clip.load() contains a torchvision Transform that performs this preprocessing. preprocess Compose(Resize(size=224, interpolation=bicubic, max_size=None, antialias=None) CenterCrop(size=(224, 224))
	<pre></pre>
6]:	tensor([[49406, 3306, 1002, 256, 49407, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
	Setting up input images and texts We are going to feed 8 example images and their textual descriptions to the model, and compare the similarity between the corresponding features. The tokenizer is case-insensitive, and we can freely give any suitable textual descriptions. import os import skimage
	<pre>import IPython.display import matplotlib.pyplot as plt from PIL import Image import numpy as np from collections import OrderedDict import torch %matplotlib inline %config InlineBackend.figure_format = 'retina'</pre>
	<pre># images in skimage to use and their textual descriptions descriptions = { "page": "a page of text about segmentation", "chelsea": "a facial photo of a tabby cat", "astronaut": "a portrait of an astronaut with the American flag", "rocket": "a rocket standing on a launchpad", "motorcycle_right": "a red motorcycle standing in a garage", "camera": "a person looking at a camera on a tripod",</pre>
8]:	<pre>"horse": "a black-and-white silhouette of a horse", "coffee": "a cup of coffee on a saucer" } original_images = [] images = [] texts = [] plt.figure(figsize=(16, 5))</pre>
	<pre>for filename in [filename for filename in os.listdir(skimage.data_dir) if filename.endswith(".png") or filename.endswith(".jpg")]: name = os.path.splitext(filename)[0] if name not in descriptions: continue image = Image.open(os.path.join(skimage.data_dir, filename)).convert("RGB") #print(imagedictkeys()) #print(imagesize) #image_sequence = image.getdata()</pre>
	<pre>#image_array = np.array(image_sequence) #print(image_array.shape) plt.subplot(2, 4, len(images) + 1) plt.imshow(image) plt.title(f"{filename}\n{descriptions[name]}") plt.xticks([]) plt.yticks([]) original_images.append(image) images.append(preprocess(image))</pre>
	images.append(preprocess(image)) texts.append(descriptions[name]) plt.tight_layout() motorcycle_right.png a red motorcycle standing in a garage a facial photo of a tabby cat a page.png a page.png a page.png a page of text about segmentation Region-based segmentation Let us first determine markers of the coins and the background. These markers are pixels that we can label
	Building features We normalize the images, tokenize each text input, and run the forward pass of the model to get the image and text features.
9]: 0]:	<pre>We normalize the images, tokenize each text input, and run the forward pass of the model to get the image and text features. image_input = torch.tensor(np.stack(images)).cuda() text_tokens = clip.tokenize(["This is " + desc for desc in texts]).cuda() with torch.no_grad(): image_features = model.encode_image(image_input).float() text_features = model.encode_text(text_tokens).float()</pre>
	Calculating cosine similarity We normalize the features and calculate the dot product of each pair. image_features /= image_features.norm(dim=-1, keepdim=True) text_features /= text_features.norm(dim=-1, keepdim=True) similarity = text_features.cpu().numpy() @ image_features.cpu().numpy().T
2]:	<pre>count = len(descriptions) plt.figure(figsize=(20, 14)) plt.imshow(similarity, vmin=0.1, vmax=0.3) # plt.colorbar() plt.yticks(range(count), texts, fontsize=18) plt.xticks([])</pre>
	<pre>for i, image in enumerate(original_images): plt.imshow(image, extent=(i - 0.5, i + 0.5, -1.6, -0.6), origin="lower") for x in range(similarity.shape[1]): for y in range(similarity.shape[0]): plt.text(x, y, f"{similarity[y, x]:.2f}", ha="center", va="center", size=12) for side in ["left", "top", "right", "bottom"]: plt.gca().spines[side].set_visible(False) plt.xlim([-0.5, count - 0.5])</pre>
2]:	plt.ylim([count + 0.5, -2]) plt.title("Cosine similarity between text and image features", size=20) Text(0.5, 1.0, 'Cosine similarity between text and image features') Cosine similarity between text and image features Cosine similarity between text and image features Co
	a red motorcycle standing in a garage - 0.32 0.15 0.14 0.13 0.12 0.16 0.15
	a facial photo of a tabby cat - 0.12 0.31 0.12 0.17 0.16 0.12 0.15 0.12 a page of text about segmentation - 0.16 0.20 0.35 0.20 0.20 0.16 0.20 0.15 a cup of coffee on a saucer - 0.12 0.18 0.14 0.29 0.17 0.12 0.15 0.15
	a person looking at a camera on a tripod - 0.16 0.21 0.19 0.14 0.30 0.21 0.20 0.19 a rocket standing on a launchpad - 0.16 0.18 0.17 0.14 0.20 0.30 0.17 0.19
	a black-and-white silhouette of a horse - 0.17 0.15 0.17 0.15 0.21 0.15 0.35 0.11 a portrait of an astronaut with the American flag - 0.15 0.17 0.13 0.15 0.17 0.22 0.16 0.28
	Zero-Shot Image Classification You can classify images using the cosine similarity (times 100) as the logits to the softmax operation.
3]: 4]:	<pre>from torchvision.datasets import CIFAR100 cifar100 = CIFAR100(os.path.expanduser("~/.cache"), transform=preprocess, download=True) Downloading https://www.cs.toronto.edu/~kriz/cifar-100-python.tar.gz to /root/.cache/cifar-100-python.tar.gz Extracting /root/.cache/cifar-100-python.tar.gz to /root/.cache text_descriptions = [f"This is a photo of a {label}" for label in cifar100.classes]</pre>
5]:	<pre>with torch.no_grad(): text_features = model.encode_text(text_tokens).float() text_features /= text_features.norm(dim=-1, keepdim=True) text_probs = (100.0 * image_features @ text_features.T).softmax(dim=-1) top_probs, top_labels = text_probs.cpu().topk(5, dim=-1)</pre>
6]:	<pre>plt.figure(figsize=(16, 16)) for i, image in enumerate(original_images): plt.subplot(4, 4, 2 * i + 1) plt.imshow(image) plt.axis("off")</pre>
	<pre>plt.subplot(4, 4, 2 * i + 2) y = np.arange(top_probs.shape[-1]) plt.grid() plt.barh(y, top_probs[i]) plt.gca().invert_yaxis() plt.gca().set_axisbelow(True) plt.yticks(y, [cifar100.classes[index] for index in top_labels[i].numpy()]) plt.xlabel("probability")</pre> <pre>plt.subplots_adjust(wspace=0.5)</pre>
	motorcycle - bicycle - bicycle - boy - boy -
	sweet_pepper
	Region-based segmentation Let us first determine markers of the coins and the background. These markers are pixels that we can label unambiguously as either object or background. Here, the markers are found at the two extreme parts of the histogram of grey values: The markers are found at the two extreme parts of the histogram of grey values: Table - ta
	probability man boy telephone plain
	voman - 0.0 0.1 0.2 0.3 0.4 probability
	plain dinosaur kangaroo oo
	Section II: Now let's do a Scavenger Hunt! We want you to figure out what caption best describes the image below. We will run your caption against images in ImageNet and display the image with the highest network probability. The goal is that you caption paired with the image below will give the highest network output.
7]:	We will download a subset of ImageNet called Tiny ImageNet. Tiny ImageNet has only 200 classes, with each class having 500 trainining images, 50 validation images and 50 test images. ! git clone https://github.com/seshuad/IMagenet
	Cloning into 'IMagenet' remote: Enumerating objects: 120594, done. remote: Total 120594 (delta 0), reused 0 (delta 0), pack-reused 120594 Receiving objects: 100% (120594/120594), 212.68 MiB 25.58 MiB/s, done. Resolving deltas: 100% (1115/1115), done. Checking out files: 100% (120206/120206), done. In order to reduce time and memory consumption, we will only consider the first 1000 images in the test set as the possible search space.
	<pre>import os img_paths = [] for rootdir, subdir, filenames in os.walk("IMagenet/tiny-imagenet-200/test/images"): for file_ in sorted(filenames)[:1000]: img_paths.append(os.path.join(rootdir, file_))</pre> TO DO: change caption below to produce target image
8]:	
8]:	for ima nath in ima naths:
8]:	<pre>for img_path in img_paths: image = Image.open(img_path).convert("RGB") original_images.append(image) images.append(preprocess(image)) image_input = torch.tensor(np.stack(images)).cuda() with torch.no_grad(): image_features = model.encode_image(image_input).float() image_features /= image_features.norm(dim=-1, keepdim=True)</pre>
8]:	<pre>image = Image.open(img_path).convert("RGB") original_images.append(image) images.append(preprocess(image)) image_input = torch.tensor(np.stack(images)).cuda() with torch.no_grad(): image_features = model.encode_image(image_input).float() image_features /= image_features.norm(dim=-1, keepdim=True) text_tokens = clip.tokenize(caption).cuda() with torch.no_grad(): text_features = model.encode_text(text_tokens).float() text_features /= text_features.norm(dim=-1, keepdim=True) text_probs = (100.0 * image_features @ text_features.T).softmax(dim=0).cpu().detach().numpy() highest_prob = np.argmax(text_probs) plt.axis('off')</pre>
8]: 1]:	<pre>image = Image.open(img_path).convert("RGB") original_images.append(image) images.append(preprocess(image)) image_input = torch.tensor(np.stack(images)).cuda() with torch.no_grad(): image_features = model.encode_image(image_input).float() image_features /= image_features.norm(dim=-1, keepdim=True) text_tokens = clip.tokenize(caption).cuda() with torch.no_grad(): text_features = model.encode_text(text_tokens).float() text_features /= text_features.norm(dim=-1, keepdim=True) text_probs = (100.0 * image_features @ text_features.T).softmax(dim=0).cpu().detach().numpy() highest_prob = np.argmax(text_probs)</pre>
8]: 2]:	<pre>image = Image.open(img_path).convert("RGB") original_images.append(image) images.append(preprocess(image)) image_input = torch.tensor(np.stack(images)).cuda() with torch.no_grad(): image_features = model.encode_image(image_input).float() image_features /= image_features.norm(dim=-1, keepdim=True) text_tokens = clip.tokenize(caption).cuda() with torch.no_grad(): text_features = model.encode_text(text_tokens).float() text_features /= text_features.norm(dim=-1, keepdim=True) text_probs = (100.0 * image_features @ text_features.T).softmax(dim=0).cpu().detach().numpy() highest_prob = np.argmax(text_probs) plt.axis('off') plt.imshow(original_images[highest_prob])</pre>
.8]: .2]:	<pre>image = Image.open(img_path).convert("RGB") original_images.append(image) images.append(preprocess(image)) image_input = torch.tensor(np.stack(images)).cuda() with torch.no_grad(): image_features = model.encode_image(image_input).float() image_features /= image_features.norm(dim=-1, keepdim=True) text_tokens = clip.tokenize(caption).cuda() with torch.no_grad(): text_features = model.encode_text(text_tokens).float() text_features /= text_features.norm(dim=-1, keepdim=True) text_probs = (100.0 * image_features @ text_features.T).softmax(dim=0).cpu().detach().numpy() highest_prob = np.argmax(text_probs) plt.axis('off') plt.imshow(original_images[highest_prob])</pre>
.8]: .2]:	image = Image.open(img_path).convert("R68") original_images.append(image) image.sapend(image) image.sapend(image) image.input = torch.temsor(np.stack(images)).cuda() with torch.np.gral():