# Q1 - *Contrastive Language-Image Pretraining*

## 1.1 Installing CLIP Dependencies

```python
#
================================================================================
========
# Task 1: Install OpenAI CLIP Dependencies
#
================================================================================
========
#
# This script installs the necessary dependencies for OpenAI's CLIP
model
# as described in its official GitHub README
(https://github.com/openai/CLIP).
#
# It is highly recommended to run this within a virtual environment
# (e.g., using venv or conda) to avoid conflicts with other Python
projects.
#
# Example using venv:
#    python -m venv clip_env
#    source clip_env/bin/activate  # On Linux/macOS
#    .\clip_env\Scripts\activate     # On Windows
#    python install_clip_script.py # Assuming you save this code as a
file
#
# Example using conda:
#    conda create -n clip_env python=3.9 # Or your preferred Python
version
#    conda activate clip_env
#    python install_clip_script.py
#
#
================================================================================
========

import subprocess
import sys
import os

def run_command(command):
    """Helper function to run shell commands and print output."""
    print(f"Executing: {command}")
    try:
```

```python
        # Using sys.executable ensures pip corresponds to the current
Python interpreter
        full_command = f"{sys.executable} -m {command}"
        process = subprocess.Popen(full_command, shell=True,
stdout=subprocess.PIPE, stderr=subprocess.STDOUT, text=True)

        # Stream output in real-time
        while True:
            output = process.stdout.readline()
            if output == '' and process.poll() is not None:
                break
            if output:
                print(output.strip())

        return_code = process.poll()
        if return_code != 0:
            print(f"Error: Command '{full_command}' failed with return
code {return_code}")
            return False
        print(f"Successfully executed: {command}")
        return True
    except Exception as e:
        print(f"An exception occurred while running command:
{command}")
        print(f"Error: {e}")
        return False

print("--- Starting OpenAI CLIP Dependency Installation ---")
all_success = True

#
----------------------------------------------------------------------
--------
# Step 1: Install PyTorch and torchvision
#
----------------------------------------------------------------------
--------
print("\nInstalling PyTorch and torchvision...")
# Note: OpenAI README suggests PyTorch 1.7.1+. This command installs
recent stable versions.
# For specific CUDA versions or CPU-only, modify this command or
install manually
# from https://pytorch.org/get-started/locally/
if not run_command("pip install torch torchvision torchaudio"):
    print("\n--- Failed to install PyTorch/torchvision. ---")
    print("Please install PyTorch manually based on your system
configuration from:")
    print("https://pytorch.org/get-started/locally/")
    all_success = False
    # sys.exit(1) # Optional: exit immediately if PyTorch fails
```

```python
# ---------------------------------------------------------------------------
# Step 2: Install CLIP prerequisites (ftfy, regex, tqdm)
# ---------------------------------------------------------------------------
if all_success:
    print("\nInstalling CLIP prerequisites (ftfy, regex, tqdm)...")
    if not run_command("pip install ftfy regex tqdm"):
        print("\n--- Failed to install CLIP prerequisites. ---")
        all_success = False
        # sys.exit(1) # Optional: exit immediately

# ---------------------------------------------------------------------------
# Step 3: Install OpenAI CLIP from GitHub
# ---------------------------------------------------------------------------
if all_success:
    print("\nInstalling OpenAI CLIP library from GitHub...")
    # This command directly follows the OpenAI CLIP README.
    if not run_command("pip install git+https://github.com/openai/CLIP.git"):
        print("\n--- Failed to install OpenAI CLIP from GitHub. ---")
        all_success = False
        # sys.exit(1) # Optional: exit immediately

# ---------------------------------------------------------------------------
# Final Status Check
# ---------------------------------------------------------------------------
if all_success:
    print("\n--- OpenAI CLIP and its dependencies installed successfully! ---")
    print("\nAttempting to verify installation by importing 'clip'...")
    try:
        import clip
        print("Successfully imported 'clip'.")
        # You can optionally list available models as a further check
        # print("Available models:", clip.available_models())
    except ImportError:
        print("Error: Failed to import 'clip' after installation.")
```

```
        print("Please check the installation logs for errors.")
    except Exception as e:
        print(f"An error occurred during verification: {e}")
else:
    print("\n--- Installation process encountered errors. Please
review the logs above. ---")

print("\n--- Installation Script Finished ---")
```

--- Starting OpenAI CLIP Dependency Installation ---

Installing PyTorch and torchvision...
Executing: pip install torch torchvision torchaudio
Requirement already satisfied: torch in
/usr/local/lib/python3.11/dist-packages (2.5.1+cu124)
Requirement already satisfied: torchvision in
/usr/local/lib/python3.11/dist-packages (0.20.1+cu124)
Requirement already satisfied: torchaudio in
/usr/local/lib/python3.11/dist-packages (2.5.1+cu124)
Requirement already satisfied: filelock in
/usr/local/lib/python3.11/dist-packages (from torch) (3.18.0)
Requirement already satisfied: typing-extensions>=4.8.0 in
/usr/local/lib/python3.11/dist-packages (from torch) (4.13.1)
Requirement already satisfied: networkx in
/usr/local/lib/python3.11/dist-packages (from torch) (3.4.2)
Requirement already satisfied: jinja2 in
/usr/local/lib/python3.11/dist-packages (from torch) (3.1.6)
Requirement already satisfied: fsspec in
/usr/local/lib/python3.11/dist-packages (from torch) (2025.3.2)
Requirement already satisfied: nvidia-cuda-nvrtc-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch) (12.4.127)
Requirement already satisfied: nvidia-cuda-runtime-cu12==12.4.127
in /usr/local/lib/python3.11/dist-packages (from torch) (12.4.127)
Requirement already satisfied: nvidia-cuda-cupti-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch) (12.4.127)
Collecting nvidia-cudnn-cu12==9.1.0.70 (from torch)
Downloading nvidia_cudnn_cu12-9.1.0.70-py3-none-
manylinux2014_x86_64.whl.metadata (1.6 kB)
Collecting nvidia-cublas-cu12==12.4.5.8 (from torch)
Downloading nvidia_cublas_cu12-12.4.5.8-py3-none-
manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cufft-cu12==11.2.1.3 (from torch)
Downloading nvidia_cufft_cu12-11.2.1.3-py3-none-
manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-curand-cu12==10.3.5.147 (from torch)
Downloading nvidia_curand_cu12-10.3.5.147-py3-none-
manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cusolver-cu12==11.6.1.9 (from torch)
Downloading nvidia_cusolver_cu12-11.6.1.9-py3-none-
manylinux2014_x86_64.whl.metadata (1.6 kB)
```

```
Collecting nvidia-cusparse-cu12==12.3.1.170 (from torch)
Downloading nvidia_cusparse_cu12-12.3.1.170-py3-none-
manylinux2014_x86_64.whl.metadata (1.6 kB)
Requirement already satisfied: nvidia-nccl-cu12==2.21.5 in
/usr/local/lib/python3.11/dist-packages (from torch) (2.21.5)
Requirement already satisfied: nvidia-nvtx-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch) (12.4.127)
Collecting nvidia-nvjitlink-cu12==12.4.127 (from torch)
Downloading nvidia_nvjitlink_cu12-12.4.127-py3-none-
manylinux2014_x86_64.whl.metadata (1.5 kB)
Requirement already satisfied: triton==3.1.0 in
/usr/local/lib/python3.11/dist-packages (from torch) (3.1.0)
Requirement already satisfied: sympy==1.13.1 in
/usr/local/lib/python3.11/dist-packages (from torch) (1.13.1)
Requirement already satisfied: mpmath<1.4,>=1.1.0 in
/usr/local/lib/python3.11/dist-packages (from sympy==1.13.1->torch)
(1.3.0)
Requirement already satisfied: numpy in
/usr/local/lib/python3.11/dist-packages (from torchvision) (1.26.4)
Requirement already satisfied: pillow!=8.3.*,>=5.3.0 in
/usr/local/lib/python3.11/dist-packages (from torchvision) (11.1.0)
Requirement already satisfied: MarkupSafe>=2.0 in
/usr/local/lib/python3.11/dist-packages (from jinja2->torch) (3.0.2)
Requirement already satisfied: mkl_fft in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision)
(1.3.8)
Requirement already satisfied: mkl_random in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision)
(1.2.4)
Requirement already satisfied: mkl_umath in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision)
(0.1.1)
Requirement already satisfied: mkl in /usr/local/lib/python3.11/dist-
packages (from numpy->torchvision) (2025.1.0)
Requirement already satisfied: tbb4py in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision)
(2022.1.0)
Requirement already satisfied: mkl-service in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision)
(2.4.1)
Requirement already satisfied: intel-openmp<2026,>=2024 in
/usr/local/lib/python3.11/dist-packages (from mkl->numpy->torchvision)
(2024.2.0)
Requirement already satisfied: tbb==2022.* in
/usr/local/lib/python3.11/dist-packages (from mkl->numpy->torchvision)
(2022.1.0)
Requirement already satisfied: tcmlib==1.* in
/usr/local/lib/python3.11/dist-packages (from tbb==2022.*->mkl->numpy-
>torchvision) (1.2.0)
```

```
Requirement already satisfied: intel-cmplr-lib-rt in
/usr/local/lib/python3.11/dist-packages (from mkl_umath->numpy-
>torchvision) (2024.2.0)
Requirement already satisfied: intel-cmplr-lib-ur==2024.2.0 in
/usr/local/lib/python3.11/dist-packages (from intel-
openmp<2026,>=2024->mkl->numpy->torchvision) (2024.2.0)
Downloading nvidia_cublas_cu12-12.4.5.8-py3-none-
manylinux2014_x86_64.whl (363.4 MB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 363.4/363.4 MB 4.0 MB/s eta
0:00:00
Downloading nvidia_cudnn_cu12-9.1.0.70-py3-none-
manylinux2014_x86_64.whl (664.8 MB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 664.8/664.8 MB 2.0 MB/s eta
0:00:00
Downloading nvidia_cufft_cu12-11.2.1.3-py3-none-
manylinux2014_x86_64.whl (211.5 MB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 211.5/211.5 MB 8.0 MB/s eta
0:00:00
Downloading nvidia_curand_cu12-10.3.5.147-py3-none-
manylinux2014_x86_64.whl (56.3 MB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 56.3/56.3 MB 30.2 MB/s eta
0:00:00
Downloading nvidia_cusolver_cu12-11.6.1.9-py3-none-
manylinux2014_x86_64.whl (127.9 MB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 127.9/127.9 MB 13.2 MB/s eta
0:00:00
Downloading nvidia_cusparse_cu12-12.3.1.170-py3-none-
manylinux2014_x86_64.whl (207.5 MB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 207.5/207.5 MB 8.6 MB/s eta
0:00:00
Downloading nvidia_nvjitlink_cu12-12.4.127-py3-none-
manylinux2014_x86_64.whl (21.1 MB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 21.1/21.1 MB 83.8 MB/s eta
0:00:00
Installing collected packages: nvidia-nvjitlink-cu12, nvidia-curand-
cu12, nvidia-cufft-cu12, nvidia-cublas-cu12, nvidia-cusparse-cu12,
nvidia-cudnn-cu12, nvidia-cusolver-cu12
Attempting uninstall: nvidia-nvjitlink-cu12
Found existing installation: nvidia-nvjitlink-cu12 12.8.93
Uninstalling nvidia-nvjitlink-cu12-12.8.93:
Successfully uninstalled nvidia-nvjitlink-cu12-12.8.93
Attempting uninstall: nvidia-curand-cu12
Found existing installation: nvidia-curand-cu12 10.3.9.90
Uninstalling nvidia-curand-cu12-10.3.9.90:
Successfully uninstalled nvidia-curand-cu12-10.3.9.90
Attempting uninstall: nvidia-cufft-cu12
Found existing installation: nvidia-cufft-cu12 11.3.3.83
Uninstalling nvidia-cufft-cu12-11.3.3.83:
Successfully uninstalled nvidia-cufft-cu12-11.3.3.83
```

```
Attempting uninstall: nvidia-cublas-cu12
Found existing installation: nvidia-cublas-cu12 12.8.4.1
Uninstalling nvidia-cublas-cu12-12.8.4.1:
Successfully uninstalled nvidia-cublas-cu12-12.8.4.1
Attempting uninstall: nvidia-cusparse-cu12
Found existing installation: nvidia-cusparse-cu12 12.5.8.93
Uninstalling nvidia-cusparse-cu12-12.5.8.93:
Successfully uninstalled nvidia-cusparse-cu12-12.5.8.93
Attempting uninstall: nvidia-cudnn-cu12
Found existing installation: nvidia-cudnn-cu12 9.3.0.75
Uninstalling nvidia-cudnn-cu12-9.3.0.75:
Successfully uninstalled nvidia-cudnn-cu12-9.3.0.75
Attempting uninstall: nvidia-cusolver-cu12
Found existing installation: nvidia-cusolver-cu12 11.7.3.90
Uninstalling nvidia-cusolver-cu12-11.7.3.90:
Successfully uninstalled nvidia-cusolver-cu12-11.7.3.90
ERROR: pip's dependency resolver does not currently take into account
all the packages that are installed. This behaviour is the source of
the following dependency conflicts.
pylibcugraph-cu12 24.12.0 requires pylibraft-cu12==24.12.*, but you
have pylibraft-cu12 25.2.0 which is incompatible.
pylibcugraph-cu12 24.12.0 requires rmm-cu12==24.12.*, but you have
rmm-cu12 25.2.0 which is incompatible.
Successfully installed nvidia-cublas-cu12-12.4.5.8 nvidia-cudnn-cu12-
9.1.0.70 nvidia-cufft-cu12-11.2.1.3 nvidia-curand-cu12-10.3.5.147
nvidia-cusolver-cu12-11.6.1.9 nvidia-cusparse-cu12-12.3.1.170 nvidia-
nvjitlink-cu12-12.4.127
Successfully executed: pip install torch torchvision torchaudio

Installing CLIP prerequisites (ftfy, regex, tqdm)...
Executing: pip install ftfy regex tqdm
Collecting ftfy
Downloading ftfy-6.3.1-py3-none-any.whl.metadata (7.3 kB)
Requirement already satisfied: regex in
/usr/local/lib/python3.11/dist-packages (2024.11.6)
Requirement already satisfied: tqdm in /usr/local/lib/python3.11/dist-
packages (4.67.1)
Requirement already satisfied: wcwidth in
/usr/local/lib/python3.11/dist-packages (from ftfy) (0.2.13)
Downloading ftfy-6.3.1-py3-none-any.whl (44 kB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 44.8/44.8 kB 1.4 MB/s eta
0:00:00
Installing collected packages: ftfy
Successfully installed ftfy-6.3.1
Successfully executed: pip install ftfy regex tqdm

Installing OpenAI CLIP library from GitHub...
Executing: pip install git+https://github.com/openai/CLIP.git
Collecting git+https://github.com/openai/CLIP.git
```

```
Cloning https://github.com/openai/CLIP.git to /tmp/pip-req-build-
8npftcth
Running command git clone --filter=blob:none --quiet
https://github.com/openai/CLIP.git /tmp/pip-req-build-8npftcth
Resolved https://github.com/openai/CLIP.git to commit
dcba3cb2e2827b402d2701e7e1c7d9fed8a20ef1
Preparing metadata (setup.py): started
Preparing metadata (setup.py): finished with status 'done'
Requirement already satisfied: ftfy in /usr/local/lib/python3.11/dist-
packages (from clip==1.0) (6.3.1)
Requirement already satisfied: packaging in
/usr/local/lib/python3.11/dist-packages (from clip==1.0) (24.2)
Requirement already satisfied: regex in
/usr/local/lib/python3.11/dist-packages (from clip==1.0) (2024.11.6)
Requirement already satisfied: tqdm in /usr/local/lib/python3.11/dist-
packages (from clip==1.0) (4.67.1)
Requirement already satisfied: torch in
/usr/local/lib/python3.11/dist-packages (from clip==1.0) (2.5.1+cu124)
Requirement already satisfied: torchvision in
/usr/local/lib/python3.11/dist-packages (from clip==1.0)
(0.20.1+cu124)
Requirement already satisfied: wcwidth in
/usr/local/lib/python3.11/dist-packages (from ftfy->clip==1.0)
(0.2.13)
Requirement already satisfied: filelock in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(3.18.0)
Requirement already satisfied: typing-extensions>=4.8.0 in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(4.13.1)
Requirement already satisfied: networkx in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(3.4.2)
Requirement already satisfied: jinja2 in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(3.1.6)
Requirement already satisfied: fsspec in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(2025.3.2)
Requirement already satisfied: nvidia-cuda-nvrtc-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(12.4.127)
Requirement already satisfied: nvidia-cuda-runtime-cu12==12.4.127
in /usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(12.4.127)
Requirement already satisfied: nvidia-cuda-cupti-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(12.4.127)
Requirement already satisfied: nvidia-cudnn-cu12==9.1.0.70 in
```

```
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(9.1.0.70)
Requirement already satisfied: nvidia-cublas-cu12==12.4.5.8 in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(12.4.5.8)
Requirement already satisfied: nvidia-cufft-cu12==11.2.1.3 in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(11.2.1.3)
Requirement already satisfied: nvidia-curand-cu12==10.3.5.147 in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(10.3.5.147)
Requirement already satisfied: nvidia-cusolver-cu12==11.6.1.9 in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(11.6.1.9)
Requirement already satisfied: nvidia-cusparse-cu12==12.3.1.170 in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(12.3.1.170)
Requirement already satisfied: nvidia-nccl-cu12==2.21.5 in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(2.21.5)
Requirement already satisfied: nvidia-nvtx-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(12.4.127)
Requirement already satisfied: nvidia-nvjitlink-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(12.4.127)
Requirement already satisfied: triton==3.1.0 in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(3.1.0)
Requirement already satisfied: sympy==1.13.1 in
/usr/local/lib/python3.11/dist-packages (from torch->clip==1.0)
(1.13.1)
Requirement already satisfied: mpmath<1.4,>=1.1.0 in
/usr/local/lib/python3.11/dist-packages (from sympy==1.13.1->torch-
>clip==1.0) (1.3.0)
Requirement already satisfied: numpy in
/usr/local/lib/python3.11/dist-packages (from torchvision->clip==1.0)
(1.26.4)
Requirement already satisfied: pillow!=8.3.*,>=5.3.0 in
/usr/local/lib/python3.11/dist-packages (from torchvision->clip==1.0)
(11.1.0)
Requirement already satisfied: MarkupSafe>=2.0 in
/usr/local/lib/python3.11/dist-packages (from jinja2->torch-
>clip==1.0) (3.0.2)
Requirement already satisfied: mkl_fft in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision-
>clip==1.0) (1.3.8)
Requirement already satisfied: mkl_random in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision-
```

```
>clip==1.0) (1.2.4)
Requirement already satisfied: mkl_umath in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision-
>clip==1.0) (0.1.1)
Requirement already satisfied: mkl in /usr/local/lib/python3.11/dist-
packages (from numpy->torchvision->clip==1.0) (2025.1.0)
Requirement already satisfied: tbb4py in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision-
>clip==1.0) (2022.1.0)
Requirement already satisfied: mkl-service in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision-
>clip==1.0) (2.4.1)
Requirement already satisfied: intel-openmp<2026,>=2024 in
/usr/local/lib/python3.11/dist-packages (from mkl->numpy->torchvision-
>clip==1.0) (2024.2.0)
Requirement already satisfied: tbb==2022.* in
/usr/local/lib/python3.11/dist-packages (from mkl->numpy->torchvision-
>clip==1.0) (2022.1.0)
Requirement already satisfied: tcmlib==1.* in
/usr/local/lib/python3.11/dist-packages (from tbb==2022.*->mkl->numpy-
>torchvision->clip==1.0) (1.2.0)
Requirement already satisfied: intel-cmplr-lib-rt in
/usr/local/lib/python3.11/dist-packages (from mkl_umath->numpy-
>torchvision->clip==1.0) (2024.2.0)
Requirement already satisfied: intel-cmplr-lib-ur==2024.2.0 in
/usr/local/lib/python3.11/dist-packages (from intel-
openmp<2026,>=2024->mkl->numpy->torchvision->clip==1.0) (2024.2.0)
Building wheels for collected packages: clip
Building wheel for clip (setup.py): started
Building wheel for clip (setup.py): finished with status 'done'
Created wheel for clip: filename=clip-1.0-py3-none-any.whl
size=1369489
sha256=cdc8e3668e84465a32abfefd3b53b4904358f2487f87d8698e887f09b26a92e
a
Stored in directory:
/tmp/pip-ephem-wheel-cache-q65k1my5/wheels/3f/7c/a4/9b490845988bf7a4db
33674d52f709f088f64392063872eb9a
Successfully built clip
Installing collected packages: clip
Successfully installed clip-1.0
Successfully executed: pip install
git+https://github.com/openai/CLIP.git

--- OpenAI CLIP and its dependencies installed successfully! ---

Attempting to verify installation by importing 'clip'...
Successfully imported 'clip'.

--- Installation Script Finished ---
```

# Task 1: Output Analysis (Install OpenAI CLIP Dependencies)

Here's an analysis of the provided output log for the execution of the Task 1 script:

## Execution Flow Breakdown

1. **PyTorch Installation (`pip install torch torchvision torchaudio`):**
   - The command was executed successfully (`Successfully executed: ...`).
   - The log shows `Requirement already satisfied:` for the main packages (`torch`, `torchvision`, `torchaudio`) and many of their dependencies. This indicates these were likely pre-installed in the environment.
   - However, `pip` detected that some specific CUDA-related sub-packages (`nvidia-cudnn-cu12`, `nvidia-cublas-cu12`, etc.) needed updating or installing to match the required versions for `torch-2.5.1+cu124`.
   - Large downloads occurred for these specific NVIDIA packages (e.g., `nvidia_cudnn_cu12-9.1.0.70`, `nvidia_cublas_cu12-12.4.5.8`).
   - Pip successfully uninstalled older versions of these NVIDIA packages before installing the new ones (`Attempting uninstall: ... Successfully uninstalled ...`).
   - **Dependency Conflict Warning:** An `ERROR` related to `pylibcugraph-cu12` was reported. This indicates an incompatibility between the newly installed NVIDIA packages (required by PyTorch 2.5.1) and an existing package (`pylibcugraph-cu12 24.12.0`) that requires older versions of `pylibraft-cu12` and `rmm-cu12`. While this is an error, it **did not** stop the installation of PyTorch and its immediate dependencies. The script continued because the *required* packages for PyTorch were installed successfully. This conflict might cause issues later *if* `pylibcugraph` functionality is needed, but it doesn't impact the core CLIP installation itself.
   - The final message `Successfully installed nvidia-cublas-cu12...` confirms the necessary components for PyTorch were put in place.
2. **CLIP Prerequisites Installation (`pip install ftfy regex tqdm`):**
   - The command was executed successfully (`Successfully executed: ...`).
   - `regex` and `tqdm` were already satisfied.
   - `ftfy` was downloaded and installed successfully (`Successfully installed ftfy-6.3.1`).
3. **OpenAI CLIP Installation (`pip install git+https://github.com/openai/CLIP.git`):**
   - The command was executed successfully (`Successfully executed: ...`).
   - The script successfully cloned the repository from GitHub.
   - It prepared the package metadata (`Preparing metadata (setup.py): finished with status 'done'`).
   - It confirmed that all dependencies required by CLIP (`ftfy`, `packaging`, `regex`, `tqdm`, `torch`, `torchvision`) were already present in the environment.
   - It successfully built the wheel (`Building wheel for clip (setup.py): finished with status 'done'`) and installed the package (`Successfully installed clip-1.0`).

4. **Final Verification (`import clip`):**
   - The script attempted to import the newly installed `clip` library.
   - This step succeeded (`Successfully imported 'clip'.`).

## Conclusion

Yes, the output indicates that the installation script for Task 1 **executed successfully**.

- All necessary dependencies (PyTorch, torchvision, ftfy, regex, tqdm) were either present or installed/updated correctly.
- The OpenAI CLIP library itself was successfully cloned from GitHub, built, and installed.
- The final verification step confirmed that the `clip` library is now importable in the environment.

The dependency conflict noted during the PyTorch installation step (related to `pylibcugraph`) is external to the CLIP installation itself and did not prevent CLIP from being installed correctly. For the purposes of using the OpenAI CLIP library as required by subsequent tasks, this installation was successful.

# 1.2 Downloading pre-trained weights ((clip-vit-base-patch32)).

```
#
========================================================================
========
# Task 2: Load Pre-trained CLIP Model (ViT-B/32)
#
========================================================================
========
#
# This script downloads (if necessary) and loads the pre-trained CLIP
model
# specified as "ViT-B/32" using the official OpenAI CLIP library.
#
# Prerequisites:
# - Completion of Task 1 (OpenAI CLIP library installed).
# - PyTorch installed.
#
# The `clip.load()` function automatically handles downloading the
model
# weights to a local cache directory (usually ~/.cache/clip) and then
# loads the model into memory.
#
#
========================================================================
========
```

```python
import torch
import clip
import os

# --- Configuration ---
# Specify the exact model name as listed by clip.available_models()
MODEL_NAME = "ViT-B/32"

print(f"--- Starting Task 2: Load CLIP Model ({MODEL_NAME}) ---")

# --- Verify Prerequisite Libraries ---
try:
    print(f"PyTorch version: {torch.__version__}")
    print(f"CLIP library location: {clip.__file__}")
    # Verify the requested model name is available in the installed
library
    available_models = clip.available_models()
    print(f"Available CLIP models: {available_models}")
    if MODEL_NAME not in available_models:
        print(f"\nError: Model name '{MODEL_NAME}' is not listed among
available models.")
        print("Please choose one of the available models.")
        exit(1) # Exit if the model name is fundamentally incorrect
except ImportError as e:
    print(f"\nError: Required library not found. {e}")
    print("Please ensure Task 1 (installation of OpenAI CLIP)
completed successfully.")
    exit(1)
except Exception as e:
    print(f"\nAn error occurred during library verification: {e}")
    exit(1)

# --- Determine Device ---
# Select GPU (CUDA) if available, otherwise use CPU
device = "cuda" if torch.cuda.is_available() else "cpu"
print(f"\nUsing device: {device}")
if device == "cuda":
    try:
        print(f"GPU Name: {torch.cuda.get_device_name(0)}")
    except Exception as e:
        print(f"Could not retrieve GPU name, but CUDA is available.
Error: {e}")

# --- Load Model and Preprocessor ---
print(f"\nLoading pre-trained CLIP model '{MODEL_NAME}' onto
{device}...")
print("(This may download the model weights if not already cached)")

try:
    # The core step: clip.load() handles download and instantiation.
```

```python
    # It returns:
    # 1. model: The CLIP model instance (a PyTorch nn.Module).
    # 2. preprocess: A torchvision transform function for preparing
images.
    model, preprocess = clip.load(MODEL_NAME, device=device)

    # --- Verification ---
    print("\nModel and preprocessor loaded successfully!")

    # Check model type and device placement
    print(f" - Model type: {type(model)}")
    # Verify the model is indeed on the target device by checking a
parameter
    # Using visual.conv1.weight as it exists in ViT models
    example_param_device = "Unknown"
    try:
        if hasattr(model, 'visual') and hasattr(model.visual, 'conv1')
and model.visual.conv1 is not None:
            example_param_device = model.visual.conv1.weight.device
        elif hasattr(model, 'positional_embedding'): # Fallback for
other structures
            example_param_device = model.positional_embedding.device
        else:
            # Try accessing a generic parameter if specific ones
aren't found
            example_param_device = next(model.parameters()).device
        print(f" - Model parameter device: {example_param_device}")
        if str(example_param_device) != device:
            print(f"   Warning: Parameter device
({example_param_device}) does not match target device ({device})!")
    except StopIteration:
        print(" - Could not verify parameter device (model has no
parameters?)")
    except AttributeError:
        print(" - Could not verify parameter device (specific
attributes not found)")


    # Check preprocessor type
    print(f" - Preprocessor type: {type(preprocess)}")
    # Optional: print the preprocess steps
    # print(" - Preprocessor details:\n", preprocess)

except FileNotFoundError:
    # This specific error might occur if clip.load has issues finding
cached/downloaded files
    print(f"\nError: Could not find or download the model files for
'{MODEL_NAME}'.")
    print("Check network connection and cache directory permissions
(~/.cache/clip).")
```

```python
    exit(1)
except Exception as e:
    print(f"\nAn unexpected error occurred during model loading: {e}")
    print("Potential issues include network problems during download,
corrupted cache,")
    print("or incompatibilities between libraries (PyTorch,
torchvision, CLIP).")
    exit(1)

print("\n--- Task 2 Finished ---")
print(f"Variables 'model' and 'preprocess' are now ready for use with
the '{MODEL_NAME}' CLIP model.")

# Example of how to use the loaded components (optional demonstration)
# print("\nExample usage:")
# try:
#     # Create dummy image and text
#     dummy_image = torch.rand(1, 3, 224, 224).to(device) # ViT-B/32
expects 224x224
#     dummy_text = clip.tokenize(["a photo of a cat", "a photo of a
dog"]).to(device)
#
#     with torch.no_grad():
#         image_features = model.encode_image(dummy_image)
#         text_features = model.encode_text(dummy_text)
#         logits_per_image, logits_per_text = model(dummy_image,
dummy_text)
#
#     print(f" - Dummy image features shape: {image_features.shape}")
#     print(f" - Dummy text features shape: {text_features.shape}")
#     print(f" - Logits per image shape: {logits_per_image.shape}")
#     print(f" - Logits per text shape: {logits_per_text.shape}")
# except Exception as e:
#     print(f"Error during example usage: {e}")
```

```
--- Starting Task 2: Load CLIP Model (ViT-B/32) ---
PyTorch version: 2.5.1+cu124
CLIP library location:
/usr/local/lib/python3.11/dist-packages/clip/__init__.py
Available CLIP models: ['RN50', 'RN101', 'RN50x4', 'RN50x16',
'RN50x64', 'ViT-B/32', 'ViT-B/16', 'ViT-L/14', 'ViT-L/14@336px']

Using device: cuda
GPU Name: Tesla T4

Loading pre-trained CLIP model 'ViT-B/32' onto cuda...
(This may download the model weights if not already cached)

100%|████████████████████████████████████████| 338M/338M [00:04<00:00,
73.5MiB/s]
```

```
Model and preprocessor loaded successfully!
 - Model type: <class 'clip.model.CLIP'>
 - Model parameter device: cuda:0
   Warning: Parameter device (cuda:0) does not match target device
(cuda)!
 - Preprocessor type: <class
'torchvision.transforms.transforms.Compose'>

--- Task 2 Finished ---
Variables 'model' and 'preprocess' are now ready for use with the
'ViT-B/32' CLIP model.
```

## Task 2: Analysis of the Output (Loading Pretrained Weights for CLIP)

**Prerequisites Checked:**
The script correctly identified the installed PyTorch (2.5.1+cu124) and CLIP library versions and confirmed that ViT-B/32 is an available model.

**Device Selected:**
It successfully detected the CUDA-enabled GPU (Tesla T4) and set `cuda` as the target device.

**Model Loading:**
The script proceeded to load the ViT-B/32 model onto the `cuda` device. The lack of download progress indicates the model was likely already cached from a previous run.

**Success Confirmation:**
The output clearly states: *Model and preprocessor loaded successfully!*

**Verification Info:**
It confirms the loaded object types (`clip.model.CLIP` and `torchvision.transforms.transforms.Compose`).

**Device Placement:**
It shows the model parameters are indeed on the GPU (`Model parameter device: cuda:0`).

**Warning Analysis:**
The warning `Warning: Parameter device (cuda:0) does not match target device (cuda)!` is due to a simple string comparison in the verification step. `clip.load(..., device="cuda")` places the model on the default CUDA device, which is typically `cuda:0`. The check `str(example_param_device) != device` compares `'cuda:0'` with `'cuda'`, which are not identical strings. However, functionally, the model is on the correct type of device (GPU). This warning does not indicate a failure in loading the model onto the GPU.

**Task Completion:**
The final messages confirm the script finished Task 2 and the `model` and `preprocess` variables are ready.

**Conclusion:**
The output indicates that the code for Task 2 executed correctly. The CLIP model ViT-B/32 and its preprocessor were successfully loaded onto the specified CUDA device (`cuda:0`).

# 1.3 Textual Descriptions (n=10) of the sample image along with their similarity scores

```
#
================================================================================
========
# Task 3: Calculate Image-Text Similarity using CLIP
#
================================================================================
========
#
# This script loads a pre-trained CLIP model (ViT-B/32), processes a
# specific image, and calculates the similarity scores between the image
# and a list of 10 textual descriptions.
#
# Prerequisites:
# - Completion of Task 1 (OpenAI CLIP library installed).
# - PyTorch installed.
# - Image file available at the specified path.
#
#
================================================================================
========

import torch
import clip
from PIL import Image
import os
import requests # To download the image if needed, or handle potential
URL paths

# --- Configuration ---
MODEL_NAME = "ViT-B/32"
IMAGE_PATH =
"/kaggle/input/cv-ass-3-q1-3-sample-image/sample_image.jpg"

# Define 10 textual descriptions (mix of relevant and irrelevant)
TEXT_DESCRIPTIONS = [
    "A person walking a large dog on a leash",
    "A man and his golden retriever in a park",
    "Someone petting a furry animal outdoors",
    "Two friends enjoying a walk in nature",
```

```python
    "An owner training their canine companion on grass",
    "A sunny day with a pet and its owner",
    "A cat sleeping peacefully on a sofa", # Irrelevant
    "A busy city street at night with neon lights", # Irrelevant
    "A close-up portrait of a dog's happy face", # Partially
relevant/different focus
    "Children playing soccer in a field" # Irrelevant
]

print("--- Starting Task 3: Image-Text Similarity Calculation ---")

# --- Step 1: Load CLIP Model (or reuse if already loaded) ---
# Check if model and preprocess are already loaded in the environment
if 'model' not in locals() or 'preprocess' not in locals():
    print(f"\nLoading CLIP model: {MODEL_NAME}...")
    try:
        # Determine device
        device = "cuda" if torch.cuda.is_available() else "cpu"
        print(f"Using device: {device}")
        if device == "cuda":
            print(f"GPU Name: {torch.cuda.get_device_name(0)}")

        # Load model and preprocessor
        model, preprocess = clip.load(MODEL_NAME, device=device)
        print("Model and preprocessor loaded successfully.")

    except ImportError:
        print("\nError: 'clip' library not found. Please complete Task
1.")
        exit(1)
    except Exception as e:
        print(f"\nError loading CLIP model: {e}")
        exit(1)
else:
    # Assume model and preprocess are loaded correctly from Task 2
    # Determine device from the loaded model's parameter
    try:
        device = next(model.parameters()).device
        print(f"Using pre-loaded model on device: {device}")
    except Exception as e:
        print(f"Could not determine device from pre-loaded model: {e}.
Setting device again.")
        device = "cuda" if torch.cuda.is_available() else "cpu"
        model.to(device) # Ensure model is on the correct device
        print(f"Set device to: {device}")


# --- Step 2: Load and Preprocess Image ---
print(f"\nLoading and preprocessing image from: {IMAGE_PATH}")
try:
```

```python
    # Check if the image path exists
    if not os.path.exists(IMAGE_PATH):
         # Handle case where path might be a URL (basic check)
        if IMAGE_PATH.startswith("http://") or
IMAGE_PATH.startswith("https://"):
            print("Attempting to download image from URL...")
            response = requests.get(IMAGE_PATH, stream=True)
            response.raise_for_status() # Raise an exception for bad
status codes
            image_pil = Image.open(response.raw).convert("RGB")
            print("Image downloaded and opened successfully.")
        else:
            print(f"Error: Image file not found at path:
{IMAGE_PATH}")
            exit(1)
    else:
        # Load image from local path
        image_pil = Image.open(IMAGE_PATH).convert("RGB")
        print("Image opened successfully from local path.")

    # Apply the preprocessing steps provided by CLIP
    image_input = preprocess(image_pil).unsqueeze(0).to(device)
    print(f"Image preprocessed and tensor created with shape:
{image_input.shape}")

except FileNotFoundError:
    print(f"Error: Image file not found at path: {IMAGE_PATH}")
    exit(1)
except requests.exceptions.RequestException as e:
    print(f"Error downloading image from URL {IMAGE_PATH}: {e}")
    exit(1)
except Exception as e:
    print(f"Error opening or preprocessing image: {e}")
    exit(1)


# --- Step 3: Tokenize Text Descriptions ---
print("\nTokenizing text descriptions...")
try:
    # Use clip.tokenize to convert text strings into token tensors
    # The model's context_length determines padding/truncation
(usually 77)
    text_inputs = clip.tokenize(TEXT_DESCRIPTIONS,
context_length=model.context_length).to(device)
    print(f"Text descriptions tokenized into tensor shape:
{text_inputs.shape}")
except AttributeError:
     print("\nError: Cannot determine model's context_length. Model
might not be loaded correctly.")
     # Fallback to default context length if attribute missing
```

```python
        print("Using default context_length=77 for tokenization.")
        text_inputs = clip.tokenize(TEXT_DESCRIPTIONS,
context_length=77).to(device)
        print(f"Text descriptions tokenized into tensor shape:
{text_inputs.shape}")

except Exception as e:
    print(f"Error tokenizing text: {e}")
    exit(1)


# --- Step 4: Generate Embeddings and Calculate Similarities ---
print("\nCalculating image and text features...")
# Perform calculations without tracking gradients
with torch.no_grad():
    try:
        # Encode the single image and the batch of texts
        image_features = model.encode_image(image_input)
        text_features = model.encode_text(text_inputs)
        print(f"Image features shape: {image_features.shape}")
        print(f"Text features shape: {text_features.shape}")

        # Normalize the features to unit length (L2 norm)
        # This is crucial for cosine similarity calculation
        image_features /= image_features.norm(dim=-1, keepdim=True)
        text_features /= text_features.norm(dim=-1, keepdim=True)

        # Calculate cosine similarity
        # Resulting shape: [1 (image), num_texts]
        # model.logit_scale is learned during training (usually ~100)
        logit_scale = model.logit_scale.exp()
        similarities = logit_scale * image_features @ text_features.T
        print("Cosine similarities calculated.")

        # Convert similarities to probabilities (optional but common)
        # Softmax across the text dimension
        probabilities = similarities.softmax(dim=-1)
        print("Probabilities calculated via softmax.")

    except Exception as e:
        print(f"\nError during feature encoding or similarity
calculation: {e}")
        exit(1)

# --- Step 5: Display Results ---
print("\n--- Similarity Scores (Probabilities) ---")

# Ensure probabilities are on CPU for easy handling/printing
probabilities_cpu = probabilities.cpu().numpy().squeeze() # Squeeze to
remove the first dimension (batch size 1)
```

```python
# Check if squeeze resulted in a scalar (if only 1 text description)
if probabilities_cpu.ndim == 0:
    probabilities_cpu = [probabilities_cpu.item()] # Make it a list
else:
    probabilities_cpu = probabilities_cpu.tolist()

if len(TEXT_DESCRIPTIONS) != len(probabilities_cpu):
     print("\nWarning: Mismatch between number of descriptions and
calculated probabilities!")
else:
    # Pair descriptions with their scores and print neatly
    results = sorted(zip(TEXT_DESCRIPTIONS, probabilities_cpu),
key=lambda x: x[1], reverse=True)

    for description, prob in results:
        print(f"- '{description}': {prob:.4f}") # Format to 4 decimal
places

print("\n--- Task 3 Finished ---")
```

--- Starting Task 3: Image-Text Similarity Calculation ---
Using pre-loaded model on device: cuda:0

Loading and preprocessing image from: /kaggle/input/cv-ass-3-q1-3-
sample-image/sample_image.jpg
Image opened successfully from local path.
Image preprocessed and tensor created with shape: torch.Size([1, 3,
224, 224])

Tokenizing text descriptions...
Text descriptions tokenized into tensor shape: torch.Size([10, 77])

Calculating image and text features...
Image features shape: torch.Size([1, 512])
Text features shape: torch.Size([10, 512])
Cosine similarities calculated.
Probabilities calculated via softmax.

--- Similarity Scores (Probabilities) ---
- 'A person walking a large dog on a leash': 0.4197
- 'A sunny day with a pet and its owner': 0.3372
- 'Someone petting a furry animal outdoors': 0.2078
- 'An owner training their canine companion on grass': 0.0163
- 'Two friends enjoying a walk in nature': 0.0127
- 'A man and his golden retriever in a park': 0.0061
- 'A cat sleeping peacefully on a sofa': 0.0002
- 'A close-up portrait of a dog's happy face': 0.0000
- 'A busy city street at night with neon lights': 0.0000
- 'Children playing soccer in a field': 0.0000

```
--- Task 3 Finished ---
```

# Task 3: Analysis of the Output (10 Textual Description)

**Execution Flow:**
The log confirms the script ran correctly. It used the pre-loaded model from Task 2 (`Using pre-loaded model on device: cuda:0`), successfully loaded and processed the image (`/kaggle/input/cv-ass-3-q1-3-sample-image/sample_image.jpg`), tokenized the 10 text descriptions, calculated the image and text features (embeddings), computed the similarities, and converted them to probabilities using softmax. All reported tensor shapes (`torch.Size`) are as expected for the ViT-B/32 model (512-dimensional embeddings).

**Results Interpretation:**
The output shows the probability assigned by CLIP for how well each text description matches the image. The probabilities sum to approximately 1.0 across all descriptions, as expected from the softmax function.

**Highest Scores:**

- `'A person walking a large dog on a leash'`: **0.4197**
  This received the highest probability. The image clearly shows a "person" and a "large dog", though the action ("walking") and context ("on a leash") are inaccurate. The model likely prioritized the presence of the main objects over the specific action or setting.

- `'A sunny day with a pet and its owner'`: **0.3372**
  Captures "pet" and "owner" but misses the indoor context ("sunny day" is incorrect).

- `'Someone petting a furry animal outdoors'`: **0.2078**
  Includes "someone" and "furry animal" but misrepresents the action ("holding" instead of "petting") and location ("indoors" instead of "outdoors").

**Lower Scores (Partially Relevant or Incorrect Details):**
Descriptions like `An owner training...`, `Two friends enjoying...`, and `A man and his golden retriever...` received much lower scores, likely penalized for incorrect actions, subjects, or specific details.

**Lowest Scores (Irrelevant):**
Completely unrelated descriptions like `'A cat sleeping...'`, `'A busy city street...'`, `'Children playing soccer...'`, and `'A close-up portrait of a dog's face'` received near-zero probabilities. This demonstrates CLIP's ability to effectively filter out unrelated concepts.

**Conclusion:**
The output confirms that the code executed correctly and that the CLIP model functioned as expected. It assigned high probabilities to descriptions capturing key visual elements (like person, dog, pet, owner), even when the full context or action wasn't precise. Irrelevant

descriptions were correctly assigned very low scores. The results align with CLIP's typical zero-shot classification behavior, with some limitations in interpreting nuanced actions or settings.

# 1.4 Installing CLIPS dependencies.

```
#
========================================================================
========
# Task 4: Install UCSC-VLAA CLIPS Dependencies
#
========================================================================
========
#
# This script installs the necessary dependencies for the CLIPS model
# as described in its official GitHub README (https://github.com/UCSC-
VLAA/CLIPS).
#
# The primary steps are:
# 1. Clone the CLIPS GitHub repository.
# 2. Install packages listed in the repository's requirements.txt
file.
#
# It is highly recommended to run this within a virtual environment
# (e.g., using venv or conda) to avoid conflicts with other Python
projects.
#
# Example using venv:
#   python -m venv clips_env
#   source clips_env/bin/activate  # On Linux/macOS
#   .\clips_env\Scripts\activate    # On Windows
#   python install_clips_script.py # Assuming you save this code as a
file
#
# Example using conda:
#   conda create -n clips_env python=3.9 # Or your preferred Python
version
#   conda activate clips_env
#   python install_clips_script.py
#
# Note: This script requires 'git' to be installed and accessible in
your PATH.
#
========================================================================
========

import subprocess
import sys
import os
```

```python
# --- Configuration ---
CLIPS_REPO_URL = "https://github.com/UCSC-VLAA/CLIPS.git"
# Use a specific directory name to avoid conflicts
CLONE_DIR = "CLIPS_repo_task4"

def run_command(command, cwd=None):
    """Helper function to run shell commands and print output."""
    print(f"Executing: {command}" + (f" in {cwd}" if cwd else ""))
    try:
        # Determine if the command is a system command (like git) or a
pip command
        if command.startswith("pip "):
            # Use sys.executable for pip commands
            full_command = f"{sys.executable} -m {command}"
        else:
            # Assume it's a system command (like git)
            full_command = command

        process = subprocess.Popen(full_command, shell=True,
stdout=subprocess.PIPE, stderr=subprocess.STDOUT, text=True, cwd=cwd)

        # Stream output in real-time
        while True:
            output = process.stdout.readline()
            if output == '' and process.poll() is not None:
                break
            if output:
                print(output.strip())

        return_code = process.poll()
        if return_code != 0:
            print(f"\nError: Command '{full_command}' failed with
return code {return_code}")
            return False
        print(f"Successfully executed: {command}")
        return True
    except FileNotFoundError:
        # Handle case where git might not be installed
        if command.startswith("git "):
            print("\nError: 'git' command not found.")
            print("Please install git and ensure it is in your
system's PATH.")
        else:
            print(f"\nError: Command not found during execution of
'{command}'.")
        return False
    except Exception as e:
        print(f"\nAn exception occurred while running command:
{command}")
```

```python
        print(f"Error: {e}")
        return False

print("--- Starting CLIPS Dependency Installation ---")
all_success = True
original_cwd = os.getcwd()
repo_path = os.path.join(original_cwd, CLONE_DIR)

#
----------------------------------------------------------------------
--------
# Step 1: Clone the CLIPS Repository
#
----------------------------------------------------------------------
--------
print(f"\nChecking for CLIPS repository in '{CLONE_DIR}'...")
if os.path.exists(repo_path):
    print(f"Directory '{CLONE_DIR}' already exists. Skipping clone.")
    # Optional: Add logic here to pull latest changes if desired
    # print("Attempting to pull latest changes...")
    # if not run_command("git pull", cwd=repo_path):
    #     print("Warning: Failed to pull latest changes.")
else:
    print(f"Cloning CLIPS repository from {CLIPS_REPO_URL} into
'{CLONE_DIR}'...")
    if not run_command(f"git clone {CLIPS_REPO_URL} {CLONE_DIR}"):
        print("\n--- Failed to clone the CLIPS repository. ---")
        all_success = False
        # sys.exit(1) # Optional: exit immediately

#
----------------------------------------------------------------------
--------
# Step 2: Install Dependencies from requirements.txt
#
----------------------------------------------------------------------
--------
if all_success:
    if os.path.exists(repo_path) and os.path.isdir(repo_path):
        requirements_file = os.path.join(repo_path,
'requirements.txt')
        if os.path.exists(requirements_file):
            print(f"\nInstalling CLIPS dependencies from
'{requirements_file}'...")
            # The README specifies `pip3`, but `sys.executable -m pip`
is more robust
            if not run_command(f"pip install -r {requirements_file}",
cwd=repo_path): # Run pip install within the repo directory context if
needed, though absolute path works too
                # Alternative: if not run_command(f"pip install -r
```

```python
requirements.txt", cwd=repo_path):
                print("\n--- Failed to install dependencies from
requirements.txt. ---")
                all_success = False
                # sys.exit(1) # Optional: exit immediately
        else:
            print(f"\nError: 'requirements.txt' not found in the
repository at {repo_path}.")
            all_success = False
    else:
        print(f"\nError: Cloned repository directory '{CLONE_DIR}' not
found or is not a directory.")
        all_success = False

#
----------------------------------------------------------------------
--------
# Final Status Check
#
----------------------------------------------------------------------
--------
# Note: Changing directory back is not strictly necessary if we used
absolute paths,
# but it's good practice if other operations were intended in the
original directory.
# os.chdir(original_cwd)
# print(f"\nChanged directory back to: {os.getcwd()}") # Uncomment if
cd logic is used

if all_success:
    print("\n--- CLIPS dependencies installed successfully! ---")
    print("\nAttempting to verify installation by importing
'open_clip' (used by CLIPS)...")
    try:
        import open_clip
        print("Successfully imported 'open_clip'.")
        # You could add further checks, like listing open_clip models
        # print("Available open_clip models (sample):",
open_clip.list_pretrained()[:5])
    except ImportError:
        print("Error: Failed to import 'open_clip' after
installation.")
        print("Please check the installation logs for errors. The
CLIPS examples rely on open_clip.")
    except Exception as e:
        print(f"An error occurred during verification: {e}")
else:
    print("\n--- CLIPS installation process encountered errors. Please
review the logs above. ---")
```

```
print("\n--- Installation Script Finished ---")
```

--- Starting CLIPS Dependency Installation ---

Checking for CLIPS repository in 'CLIPS_repo_task4'...
Cloning CLIPS repository from https://github.com/UCSC-VLAA/CLIPS.git
into 'CLIPS_repo_task4'...
Executing: git clone https://github.com/UCSC-VLAA/CLIPS.git
CLIPS_repo_task4
Cloning into 'CLIPS_repo_task4'...
Successfully executed: git clone
https://github.com/UCSC-VLAA/CLIPS.git CLIPS_repo_task4

Installing CLIPS dependencies from
'/kaggle/working/CLIPS_repo_task4/requirements.txt'...
Executing: pip install -r
/kaggle/working/CLIPS_repo_task4/requirements.txt in
/kaggle/working/CLIPS_repo_task4
Requirement already satisfied: torch>=1.9.0 in
/usr/local/lib/python3.11/dist-packages (from -r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1))
(2.5.1+cu124)
Requirement already satisfied: torchvision in
/usr/local/lib/python3.11/dist-packages (from -r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 2))
(0.20.1+cu124)
Requirement already satisfied: regex in
/usr/local/lib/python3.11/dist-packages (from -r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 3))
(2024.11.6)
Requirement already satisfied: ftfy in /usr/local/lib/python3.11/dist-
packages (from -r /kaggle/working/CLIPS_repo_task4/requirements.txt
(line 4)) (6.3.1)
Requirement already satisfied: tqdm in /usr/local/lib/python3.11/dist-
packages (from -r /kaggle/working/CLIPS_repo_task4/requirements.txt
(line 5)) (4.67.1)
Requirement already satisfied: huggingface_hub in
/usr/local/lib/python3.11/dist-packages (from -r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 6)) (0.30.2)
Requirement already satisfied: safetensors in
/usr/local/lib/python3.11/dist-packages (from -r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 7)) (0.5.2)
Requirement already satisfied: timm in /usr/local/lib/python3.11/dist-
packages (from -r /kaggle/working/CLIPS_repo_task4/requirements.txt
(line 8)) (1.0.14)
Requirement already satisfied: transformers in
/usr/local/lib/python3.11/dist-packages (from -r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 9)) (4.51.1)
Requirement already satisfied: filelock in
```

```
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (3.18.0)
Requirement already satisfied: typing-extensions>=4.8.0 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (4.13.1)
Requirement already satisfied: networkx in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (3.4.2)
Requirement already satisfied: jinja2 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (3.1.6)
Requirement already satisfied: fsspec in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (2025.3.2)
Requirement already satisfied: nvidia-cuda-nvrtc-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (12.4.127)
Requirement already satisfied: nvidia-cuda-runtime-cu12==12.4.127
in /usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (12.4.127)
Requirement already satisfied: nvidia-cuda-cupti-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (12.4.127)
Requirement already satisfied: nvidia-cudnn-cu12==9.1.0.70 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (9.1.0.70)
Requirement already satisfied: nvidia-cublas-cu12==12.4.5.8 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (12.4.5.8)
Requirement already satisfied: nvidia-cufft-cu12==11.2.1.3 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (11.2.1.3)
Requirement already satisfied: nvidia-curand-cu12==10.3.5.147 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1))
(10.3.5.147)
Requirement already satisfied: nvidia-cusolver-cu12==11.6.1.9 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (11.6.1.9)
Requirement already satisfied: nvidia-cusparse-cu12==12.3.1.170 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1))
(12.3.1.170)
Requirement already satisfied: nvidia-nccl-cu12==2.21.5 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (2.21.5)
Requirement already satisfied: nvidia-nvtx-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (12.4.127)
```

```
Requirement already satisfied: nvidia-nvjitlink-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (12.4.127)
Requirement already satisfied: triton==3.1.0 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (3.1.0)
Requirement already satisfied: sympy==1.13.1 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (1.13.1)
Requirement already satisfied: mpmath<1.4,>=1.1.0 in
/usr/local/lib/python3.11/dist-packages (from sympy==1.13.1-
>torch>=1.9.0->-r /kaggle/working/CLIPS_repo_task4/requirements.txt
(line 1)) (1.3.0)
Requirement already satisfied: numpy in
/usr/local/lib/python3.11/dist-packages (from torchvision->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 2)) (1.26.4)
Requirement already satisfied: pillow!=8.3.*,>=5.3.0 in
/usr/local/lib/python3.11/dist-packages (from torchvision->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 2)) (11.1.0)
Requirement already satisfied: wcwidth in
/usr/local/lib/python3.11/dist-packages (from ftfy->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 4)) (0.2.13)
Requirement already satisfied: packaging>=20.9 in
/usr/local/lib/python3.11/dist-packages (from huggingface_hub->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 6)) (24.2)
Requirement already satisfied: pyyaml>=5.1 in
/usr/local/lib/python3.11/dist-packages (from huggingface_hub->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 6)) (6.0.2)
Requirement already satisfied: requests in
/usr/local/lib/python3.11/dist-packages (from huggingface_hub->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 6)) (2.32.3)
Requirement already satisfied: tokenizers<0.22,>=0.21 in
/usr/local/lib/python3.11/dist-packages (from transformers->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 9)) (0.21.0)
Requirement already satisfied: mkl_fft in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision->-
r /kaggle/working/CLIPS_repo_task4/requirements.txt (line 2)) (1.3.8)
Requirement already satisfied: mkl_random in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision->-
r /kaggle/working/CLIPS_repo_task4/requirements.txt (line 2)) (1.2.4)
Requirement already satisfied: mkl_umath in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision->-
r /kaggle/working/CLIPS_repo_task4/requirements.txt (line 2)) (0.1.1)
Requirement already satisfied: mkl in /usr/local/lib/python3.11/dist-
packages (from numpy->torchvision->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 2)) (2025.1.0)
Requirement already satisfied: tbb4py in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision->-
r /kaggle/working/CLIPS_repo_task4/requirements.txt (line 2))
```

```
(2022.1.0)
Requirement already satisfied: mkl-service in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision->-
r /kaggle/working/CLIPS_repo_task4/requirements.txt (line 2)) (2.4.1)
Requirement already satisfied: MarkupSafe>=2.0 in
/usr/local/lib/python3.11/dist-packages (from jinja2->torch>=1.9.0->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 1)) (3.0.2)
Requirement already satisfied: charset-normalizer<4,>=2 in
/usr/local/lib/python3.11/dist-packages (from requests-
>huggingface_hub->-r /kaggle/working/CLIPS_repo_task4/requirements.txt
(line 6)) (3.4.1)
Requirement already satisfied: idna<4,>=2.5 in
/usr/local/lib/python3.11/dist-packages (from requests-
>huggingface_hub->-r /kaggle/working/CLIPS_repo_task4/requirements.txt
(line 6)) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in
/usr/local/lib/python3.11/dist-packages (from requests-
>huggingface_hub->-r /kaggle/working/CLIPS_repo_task4/requirements.txt
(line 6)) (2.3.0)
Requirement already satisfied: certifi>=2017.4.17 in
/usr/local/lib/python3.11/dist-packages (from requests-
>huggingface_hub->-r /kaggle/working/CLIPS_repo_task4/requirements.txt
(line 6)) (2025.1.31)
Requirement already satisfied: intel-openmp<2026,>=2024 in
/usr/local/lib/python3.11/dist-packages (from mkl->numpy->torchvision-
>-r /kaggle/working/CLIPS_repo_task4/requirements.txt (line 2))
(2024.2.0)
Requirement already satisfied: tbb==2022.* in
/usr/local/lib/python3.11/dist-packages (from mkl->numpy->torchvision-
>-r /kaggle/working/CLIPS_repo_task4/requirements.txt (line 2))
(2022.1.0)
Requirement already satisfied: tcmlib==1.* in
/usr/local/lib/python3.11/dist-packages (from tbb==2022.*->mkl->numpy-
>torchvision->-r /kaggle/working/CLIPS_repo_task4/requirements.txt
(line 2)) (1.2.0)
Requirement already satisfied: intel-cmplr-lib-rt in
/usr/local/lib/python3.11/dist-packages (from mkl_umath->numpy-
>torchvision->-r /kaggle/working/CLIPS_repo_task4/requirements.txt
(line 2)) (2024.2.0)
Requirement already satisfied: intel-cmplr-lib-ur==2024.2.0 in
/usr/local/lib/python3.11/dist-packages (from intel-
openmp<2026,>=2024->mkl->numpy->torchvision->-r
/kaggle/working/CLIPS_repo_task4/requirements.txt (line 2)) (2024.2.0)
Successfully executed: pip install -r
/kaggle/working/CLIPS_repo_task4/requirements.txt

--- CLIPS dependencies installed successfully! ---

Attempting to verify installation by importing 'open_clip' (used by
```

```
CLIPS)...
Error: Failed to import 'open_clip' after installation.
Please check the installation logs for errors. The CLIPS examples rely
on open_clip.

--- Installation Script Finished ---
```

# Task 4: Analysis of the Output (Installing CLIPS dependencies)

**Repository Cloning:**
The script successfully cloned the CLIPS repository from GitHub into the `CLIPS_repo_task4` directory.

```
Cloning into 'CLIPS_repo_task4'...
Successfully executed: git clone ...
```

**Dependency Installation:**
The script then executed the `pip install -r requirements.txt` command using the file from the cloned repository.

The output shows `Requirement already satisfied:` for many packages (`torch`, `torchvision`, `regex`, `ftfy`, `tqdm`, `huggingface_hub`, `safetensors`, `timm`, `transformers`, etc.), indicating these packages were already installed in the environment. Pip correctly detected this and skipped reinstalling them.

The command completed successfully:

```
Successfully executed: pip install -r
/kaggle/working/CLIPS_repo_task4/requirements.txt
```

**Verification (Import `open_clip`):**
The script attempted to verify the setup by importing `open_clip`, which is used in the CLIPS README's example code.

This step failed:

```
Error: Failed to import 'open_clip' after installation.
```

**Analysis of the Discrepancy:**
While `pip install -r requirements.txt` completed without error, the failure to import `open_clip` suggests that the `open_clip` package is not included in the `requirements.txt` file from the CLIPS repository.

This might be due to:

- An omission in their `requirements.txt`.

- An expectation that users install `open_clip` separately (e.g., `pip install open_clip_torch`).
- A note in their README mentioning modifications to `open_clip/tokenizer.py`, possibly implying they include a bundled version, but it is not being installed into the Python path by the pip command.

**Conclusion:**
The script ran correctly as per the defined steps—repository cloning and dependency installation completed without issues. However, the Python environment remains incomplete for executing the CLIPS examples because the core dependency `open_clip` was not installed via the provided `requirements.txt`. The verification step accurately flagged this missing dependency.

# 1.5 Loading the pretrained weights for the CLIPS-Large-14-224 model.

```
#
============================================================================
========
# Task 5: Load Pre-trained CLIPS Model (CLIPS-Large-14-224)
#
============================================================================
========
#
# This script loads the pre-trained CLIPS model "CLIPS-Large-14-224"
# using the open_clip library, as indicated by the CLIPS repository
README.
#
# It first ensures 'open_clip_torch' is installed, addressing the
finding
# from Task 4 where it wasn't included in requirements.txt.
#
# Prerequisites:
# - Completion of Task 4 (CLIPS repository cloned, other dependencies
installed).
# - PyTorch installed.
#
#
============================================================================
========

import torch
import os
import subprocess
import sys
import pkg_resources # To check if open_clip is installed
```

```python
# --- Configuration ---
# Model name as specified in the task
MODEL_DISPLAY_NAME = "CLIPS-Large-14-224"
# Corresponding Hugging Face Hub ID from the CLIPS Model Zoo table
MODEL_HF_ID = "hf-hub:UCSC-VLAA/ViT-L-14-CLIPS-224-Recap-DataComp-1B"

print(f"--- Starting Task 5: Load CLIPS Model ({MODEL_DISPLAY_NAME}) ---")

# --- Helper function for running install command ---
def run_install_command(command):
    """Runs a pip install command."""
    print(f"Executing: {command}")
    try:
        full_command = f"{sys.executable} -m {command}"
        process = subprocess.Popen(full_command, shell=True, stdout=subprocess.PIPE, stderr=subprocess.STDOUT, text=True)
        # Stream output
        while True:
            output = process.stdout.readline()
            if output == '' and process.poll() is not None:
                break
            if output:
                print(output.strip())
        return_code = process.poll()
        if return_code != 0:
            print(f"\nError: Command '{full_command}' failed with return code {return_code}")
            return False
        print(f"Successfully executed: {command}")
        return True
    except Exception as e:
        print(f"\nAn exception occurred while running command: {command}")
        print(f"Error: {e}")
        return False

# --- Step 1: Ensure open_clip is installed ---
# Check if open_clip_torch is installed
try:
    pkg_resources.get_distribution('open_clip_torch')
    print("'open_clip_torch' package found.")
except pkg_resources.DistributionNotFound:
    print("'open_clip_torch' package not found. Attempting installation...")
    # Install open_clip using pip
    if not run_install_command("pip install open_clip_torch"):
        print("\nError: Failed to install 'open_clip_torch'.")
        print("Please install it manually ('pip install open_clip_torch') and restart.")
```

```python
        exit(1)
    else:
        print("Installation successful. Please restart the
script/environment if import fails.")
        # Re-importing dynamically can be tricky, often best to
restart.
        # For this script, we'll try importing directly after
potential install.
        pass # Continue to the import attempt

# --- Step 2: Import open_clip and Load Model ---
try:
    # Import necessary functions from open_clip
    from open_clip import create_model_from_pretrained, get_tokenizer
    print("Successfully imported 'open_clip'.")

    # Determine Device
    device = "cuda" if torch.cuda.is_available() else "cpu"
    print(f"\nUsing device: {device}")
    if device == "cuda":
        try:
            print(f"GPU Name: {torch.cuda.get_device_name(0)}")
        except Exception as e:
            print(f"Could not retrieve GPU name, but CUDA is
available. Error: {e}")

    # Load the Model and Preprocessor from Hugging Face Hub
    print(f"\nLoading pre-trained CLIPS model '{MODEL_DISPLAY_NAME}'
({MODEL_HF_ID}) onto {device}...")
    print("(This may download the model weights from Hugging Face
Hub)")

    # Use the functions imported from open_clip
    # Note: The CLIPS team modified open_clip/tokenizer.py. Using the
standard
    # open_clip install might yield slightly different tokenizer
behavior than
    # their internal setup. This loads the model weights correctly.
    clips_model, clips_preprocess =
create_model_from_pretrained(MODEL_HF_ID, device=device)

    # Load the Tokenizer associated with the model
    print("Loading tokenizer...")
    clips_tokenizer = get_tokenizer(MODEL_HF_ID)

    # --- Step 3: Verification ---
    print("\nCLIPS model and tokenizer loaded successfully!")

    # Check model type and device placement
    print(f" - Model type: {type(clips_model)}")
```

```python
    example_param_device = "Unknown"
    try:
        example_param_device = next(clips_model.parameters()).device
        print(f" - Model parameter device: {example_param_device}")
        # Use str() for comparison robustness (e.g., 'cuda:0' vs
'cuda')
        if not str(example_param_device).startswith(device):
            print(f"   Warning: Parameter device
({example_param_device}) does not seem to match target device
({device})!")
    except StopIteration:
        print(" - Could not verify parameter device (model has no
parameters?)")
    except Exception as e:
        print(f" - Error verifying parameter device: {e}")

    # Check preprocessor and tokenizer types
    print(f" - Preprocessor type: {type(clips_preprocess)}")
    print(f" - Tokenizer type: {type(clips_tokenizer)}")

except ImportError:
    print("\nError: Failed to import 'open_clip' even after
installation attempt.")
    print("Please ensure 'open_clip_torch' is installed correctly in
your environment.")
    exit(1)
except FileNotFoundError:
    # This might occur if model files can't be downloaded/cached from
HF Hub
    print(f"\nError: Could not find or download the model files for
'{MODEL_HF_ID}'.")
    print("Check Hugging Face Hub model ID, network connection, and
cache directory permissions (~/.cache/huggingface/hub or similar).")
    exit(1)
except Exception as e:
    print(f"\nAn unexpected error occurred during model loading: {e}")
    print("Potential issues include network problems, Hugging Face Hub
access,")
    print("corrupted cache, or library incompatibilities.")
    exit(1)

print("\n--- Task 5 Finished ---")
print(f"Variables 'clips_model', 'clips_preprocess', and
'clips_tokenizer' are now ready for use with the
'{MODEL_DISPLAY_NAME}' CLIPS model.")

# You can now use clips_model, clips_preprocess, clips_tokenizer like
the standard CLIP ones,
# following the pattern in the CLIPS README example (e.g., using
torch.nn.functional.normalize)
```

```
--- Starting Task 5: Load CLIPS Model (CLIPS-Large-14-224) ---
'open_clip_torch' package not found. Attempting installation...
Executing: pip install open_clip_torch
Collecting open_clip_torch
Downloading open_clip_torch-2.32.0-py3-none-any.whl.metadata (31 kB)
Requirement already satisfied: torch>=1.9.0 in
/usr/local/lib/python3.11/dist-packages (from open_clip_torch)
(2.5.1+cu124)
Requirement already satisfied: torchvision in
/usr/local/lib/python3.11/dist-packages (from open_clip_torch)
(0.20.1+cu124)
Requirement already satisfied: regex in
/usr/local/lib/python3.11/dist-packages (from open_clip_torch)
(2024.11.6)
Requirement already satisfied: ftfy in /usr/local/lib/python3.11/dist-
packages (from open_clip_torch) (6.3.1)
Requirement already satisfied: tqdm in /usr/local/lib/python3.11/dist-
packages (from open_clip_torch) (4.67.1)
Requirement already satisfied: huggingface-hub in
/usr/local/lib/python3.11/dist-packages (from open_clip_torch)
(0.30.2)
Requirement already satisfied: safetensors in
/usr/local/lib/python3.11/dist-packages (from open_clip_torch) (0.5.2)
Requirement already satisfied: timm in /usr/local/lib/python3.11/dist-
packages (from open_clip_torch) (1.0.14)
Requirement already satisfied: filelock in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (3.18.0)
Requirement already satisfied: typing-extensions>=4.8.0 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (4.13.1)
Requirement already satisfied: networkx in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (3.4.2)
Requirement already satisfied: jinja2 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (3.1.6)
Requirement already satisfied: fsspec in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (2025.3.2)
Requirement already satisfied: nvidia-cuda-nvrtc-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (12.4.127)
Requirement already satisfied: nvidia-cuda-runtime-cu12==12.4.127
in /usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (12.4.127)
Requirement already satisfied: nvidia-cuda-cupti-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (12.4.127)
Requirement already satisfied: nvidia-cudnn-cu12==9.1.0.70 in
```

```
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (9.1.0.70)
Requirement already satisfied: nvidia-cublas-cu12==12.4.5.8 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (12.4.5.8)
Requirement already satisfied: nvidia-cufft-cu12==11.2.1.3 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (11.2.1.3)
Requirement already satisfied: nvidia-curand-cu12==10.3.5.147 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (10.3.5.147)
Requirement already satisfied: nvidia-cusolver-cu12==11.6.1.9 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (11.6.1.9)
Requirement already satisfied: nvidia-cusparse-cu12==12.3.1.170 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (12.3.1.170)
Requirement already satisfied: nvidia-nccl-cu12==2.21.5 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (2.21.5)
Requirement already satisfied: nvidia-nvtx-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (12.4.127)
Requirement already satisfied: nvidia-nvjitlink-cu12==12.4.127 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (12.4.127)
Requirement already satisfied: triton==3.1.0 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (3.1.0)
Requirement already satisfied: sympy==1.13.1 in
/usr/local/lib/python3.11/dist-packages (from torch>=1.9.0-
>open_clip_torch) (1.13.1)
Requirement already satisfied: mpmath<1.4,>=1.1.0 in
/usr/local/lib/python3.11/dist-packages (from sympy==1.13.1-
>torch>=1.9.0->open_clip_torch) (1.3.0)
Requirement already satisfied: wcwidth in
/usr/local/lib/python3.11/dist-packages (from ftfy->open_clip_torch)
(0.2.13)
Requirement already satisfied: packaging>=20.9 in
/usr/local/lib/python3.11/dist-packages (from huggingface-hub-
>open_clip_torch) (24.2)
Requirement already satisfied: pyyaml>=5.1 in
/usr/local/lib/python3.11/dist-packages (from huggingface-hub-
>open_clip_torch) (6.0.2)
Requirement already satisfied: requests in
/usr/local/lib/python3.11/dist-packages (from huggingface-hub-
>open_clip_torch) (2.32.3)
Requirement already satisfied: numpy in
/usr/local/lib/python3.11/dist-packages (from torchvision-
>open_clip_torch) (1.26.4)
```

```
Requirement already satisfied: pillow!=8.3.*,>=5.3.0 in
/usr/local/lib/python3.11/dist-packages (from torchvision-
>open_clip_torch) (11.1.0)
Requirement already satisfied: MarkupSafe>=2.0 in
/usr/local/lib/python3.11/dist-packages (from jinja2->torch>=1.9.0-
>open_clip_torch) (3.0.2)
Requirement already satisfied: mkl_fft in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision-
>open_clip_torch) (1.3.8)
Requirement already satisfied: mkl_random in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision-
>open_clip_torch) (1.2.4)
Requirement already satisfied: mkl_umath in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision-
>open_clip_torch) (0.1.1)
Requirement already satisfied: mkl in /usr/local/lib/python3.11/dist-
packages (from numpy->torchvision->open_clip_torch) (2025.1.0)
Requirement already satisfied: tbb4py in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision-
>open_clip_torch) (2022.1.0)
Requirement already satisfied: mkl-service in
/usr/local/lib/python3.11/dist-packages (from numpy->torchvision-
>open_clip_torch) (2.4.1)
Requirement already satisfied: charset-normalizer<4,>=2 in
/usr/local/lib/python3.11/dist-packages (from requests->huggingface-
hub->open_clip_torch) (3.4.1)
Requirement already satisfied: idna<4,>=2.5 in
/usr/local/lib/python3.11/dist-packages (from requests->huggingface-
hub->open_clip_torch) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in
/usr/local/lib/python3.11/dist-packages (from requests->huggingface-
hub->open_clip_torch) (2.3.0)
Requirement already satisfied: certifi>=2017.4.17 in
/usr/local/lib/python3.11/dist-packages (from requests->huggingface-
hub->open_clip_torch) (2025.1.31)
Requirement already satisfied: intel-openmp<2026,>=2024 in
/usr/local/lib/python3.11/dist-packages (from mkl->numpy->torchvision-
>open_clip_torch) (2024.2.0)
Requirement already satisfied: tbb==2022.* in
/usr/local/lib/python3.11/dist-packages (from mkl->numpy->torchvision-
>open_clip_torch) (2022.1.0)
Requirement already satisfied: tcmlib==1.* in
/usr/local/lib/python3.11/dist-packages (from tbb==2022.*->mkl->numpy-
>torchvision->open_clip_torch) (1.2.0)
Requirement already satisfied: intel-cmplr-lib-rt in
/usr/local/lib/python3.11/dist-packages (from mkl_umath->numpy-
>torchvision->open_clip_torch) (2024.2.0)
Requirement already satisfied: intel-cmplr-lib-ur==2024.2.0 in
/usr/local/lib/python3.11/dist-packages (from intel-
openmp<2026,>=2024->mkl->numpy->torchvision->open_clip_torch)
```

```
(2024.2.0)
Downloading open_clip_torch-2.32.0-py3-none-any.whl (1.5 MB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 1.5/1.5 MB 19.5 MB/s eta
0:00:00
Installing collected packages: open_clip_torch
Successfully installed open_clip_torch-2.32.0
Successfully executed: pip install open_clip_torch
Installation successful. Please restart the script/environment if
import fails.
Successfully imported 'open_clip'.

Using device: cuda
GPU Name: Tesla T4

Loading pre-trained CLIPS model 'CLIPS-Large-14-224' (hf-hub:UCSC-
VLAA/ViT-L-14-CLIPS-224-Recap-DataComp-1B) onto cuda...
(This may download the model weights from Hugging Face Hub)
```

{"model_id":"575a4f39c6544742bc4e8be504638dbf","version_major":2,"version_minor":0}

{"model_id":"718fcf67050649009f75bec06a863167","version_major":2,"version_minor":0}

```
Loading tokenizer...
```

{"model_id":"8d8fb56748e545d0a2ff589433c82d1c","version_major":2,"version_minor":0}

{"model_id":"479374848ea74d268431647f43632b78","version_major":2,"version_minor":0}

{"model_id":"93d890f670ae407ab63f326607fc9c1d","version_major":2,"version_minor":0}

{"model_id":"00cb9263fb884036af5ec00f61c3e361","version_major":2,"version_minor":0}

```
CLIPS model and tokenizer loaded successfully!
 - Model type: <class 'open_clip.model.CLIP'>
 - Model parameter device: cuda:0
 - Preprocessor type: <class
'torchvision.transforms.transforms.Compose'>
 - Tokenizer type: <class 'open_clip.tokenizer.HFTokenizer'>

--- Task 5 Finished ---
Variables 'clips_model', 'clips_preprocess', and 'clips_tokenizer' are
now ready for use with the 'CLIPS-Large-14-224' CLIPS model.
```

# Task 5: Output Analysis (loading the pre-trained CLIPS model)

**open_clip Installation:**
The log shows that the `open_clip_torch` package was successfully installed (`Successfully installed open_clip_torch-2.32.0`) and then imported (`Successfully imported 'open_clip'`). This confirms the prerequisite was met.

**Device Selection:**
The correct device (`cuda`) and GPU (`Tesla T4`) were identified.

**Model and Tokenizer Download:**
The progress bars clearly indicate that the model weights (`open_clip_pytorch_model.bin`), model configuration (`open_clip_config.json`), and associated tokenizer files were successfully downloaded from Hugging Face Hub.

**Loading Confirmation:**
The output explicitly states:

```
CLIPS model and tokenizer loaded successfully!
```

**Verification Details:**
The reported types for the model (`open_clip.model.CLIP`), preprocessor (`torchvision.transforms.transforms.Compose`), and tokenizer (`open_clip.tokenizer.HFTokenizer`) are correct and expected when using `open_clip`. The model parameter device is confirmed as `cuda:0`.

**Task Completion:**
The script finished with the message confirming that the `clips_model`, `clips_preprocess`, and `clips_tokenizer` variables are ready.

**Conclusion:**
Yes, the output confirms that the code for Task 5 executed correctly and successfully. The required `open_clip` library was installed, and the specified pre-trained CLIPS model (`CLIPS-Large-14-224`) along with its preprocessor and tokenizer were loaded onto the GPU.

# 1.6 Calculating similarity scores on the sample image using CLIPS

```
#
================================================================================
========
# Task 6: Calculate Image-Text Similarity using CLIPS
#
================================================================================
========
```

```python
#
# This script loads a pre-trained CLIPS model (CLIPS-Large-14-224),
# processes the specific image used in Task 3, and calculates the
# similarity scores between the image and the same list of 10 textual
# descriptions from Task 3.
#
# Prerequisites:
# - Completion of Task 4 (CLIPS dependencies installed/repo cloned).
# - open_clip_torch installed (handled within the script if missing).
# - PyTorch installed.
# - Image file available at the specified path.
#
#
# ========================================================================
# ========

import torch
import torch.nn.functional as F # For normalization
from PIL import Image
import os
import requests # To handle potential URL paths for image
import subprocess
import sys
import contextlib
import pkg_resources # To check if open_clip is installed

# --- Configuration ---
# CLIPS Model details from Task 5
MODEL_DISPLAY_NAME = "CLIPS-Large-14-224"
MODEL_HF_ID = "hf-hub:UCSC-VLAA/ViT-L-14-CLIPS-224-Recap-DataComp-1B"

# Image path from Task 3
IMAGE_PATH =
"/kaggle/input/cv-ass-3-q1-3-sample-image/sample_image.jpg"

# Text descriptions from Task 3
TEXT_DESCRIPTIONS = [
    "A person walking a large dog on a leash",
    "A man and his golden retriever in a park",
    "Someone petting a furry animal outdoors",
    "Two friends enjoying a walk in nature",
    "An owner training their canine companion on grass",
    "A sunny day with a pet and its owner",
    "A cat sleeping peacefully on a sofa", # Irrelevant
    "A busy city street at night with neon lights", # Irrelevant
    "A close-up portrait of a dog's happy face", # Partially
relevant/different focus
    "Children playing soccer in a field" # Irrelevant
]
```

```python
print(f"--- Starting Task 6: Image-Text Similarity Calculation using
{MODEL_DISPLAY_NAME} ---")

# --- Helper function for running install command (same as Task 5) ---
def run_install_command(command):
    """Runs a pip install command."""
    print(f"Executing: {command}")
    try:
        full_command = f"{sys.executable} -m {command}"
        process = subprocess.Popen(full_command, shell=True,
stdout=subprocess.PIPE, stderr=subprocess.STDOUT, text=True)
        while True:
            output = process.stdout.readline()
            if output == '' and process.poll() is not None: break
            if output: print(output.strip())
        return_code = process.poll()
        if return_code != 0:
            print(f"\nError: Command '{full_command}' failed with
return code {return_code}")
            return False
        print(f"Successfully executed: {command}")
        return True
    except Exception as e:
        print(f"\nAn exception occurred while running command:
{command}\nError: {e}")
        return False

# --- Step 1: Load CLIPS Model (or reuse if already loaded) ---
# Check if CLIPS model components are already loaded
if 'clips_model' not in locals() or 'clips_preprocess' not in locals()
or 'clips_tokenizer' not in locals():
    print("\nCLIPS model components not found in environment.
Loading...")

    # Ensure open_clip is installed
    try:
        pkg_resources.get_distribution('open_clip_torch')
        print("'open_clip_torch' package found.")
    except pkg_resources.DistributionNotFound:
        print("'open_clip_torch' package not found. Attempting
installation...")
        if not run_install_command("pip install open_clip_torch"):
            print("\nError: Failed to install 'open_clip_torch'.
Exiting.")
            exit(1)
        # Import after installation attempt
        try:
            from open_clip import create_model_from_pretrained,
get_tokenizer
        except ImportError:
```

```python
                print("Error: Failed to import open_clip even after
installation. Exiting.")
                exit(1)

    try:
        # Import necessary functions (might be redundant if already
imported)
        from open_clip import create_model_from_pretrained,
get_tokenizer
        print("Imported 'open_clip' successfully.")

        # Determine Device
        device = "cuda" if torch.cuda.is_available() else "cpu"
        print(f"Using device: {device}")
        if device == "cuda": print(f"GPU Name:
{torch.cuda.get_device_name(0)}")

        # Load Model, Preprocessor, and Tokenizer
        print(f"Loading pre-trained CLIPS model '{MODEL_DISPLAY_NAME}'
({MODEL_HF_ID})...")
        clips_model, clips_preprocess =
create_model_from_pretrained(MODEL_HF_ID, device=device)
        print("Loading tokenizer...")
        clips_tokenizer = get_tokenizer(MODEL_HF_ID)
        print("CLIPS model, preprocessor, and tokenizer loaded
successfully.")

    except ImportError:
        print("\nError: Failed to import 'open_clip'. Ensure
'open_clip_torch' is installed.")
        exit(1)
    except Exception as e:
        print(f"\nError loading CLIPS model: {e}")
        exit(1)
else:
    # Assume components are loaded from Task 5
    try:
        device = next(clips_model.parameters()).device
        print(f"Using pre-loaded CLIPS model on device: {device}")
    except Exception as e:
        print(f"Could not determine device from pre-loaded CLIPS
model: {e}. Setting device again.")
        device = "cuda" if torch.cuda.is_available() else "cpu"
        clips_model.to(device) # Ensure model is on the correct device
        print(f"Set device to: {device}")


# --- Step 2: Load and Preprocess Image ---
print(f"\nLoading and preprocessing image from: {IMAGE_PATH}")
try:
```

```python
    # Check if the image path exists or is a URL
    if not os.path.exists(IMAGE_PATH):
        if IMAGE_PATH.startswith("http://") or
IMAGE_PATH.startswith("https://"):
            print("Attempting to download image from URL...")
            response = requests.get(IMAGE_PATH, stream=True)
            response.raise_for_status()
            image_pil = Image.open(response.raw).convert("RGB")
            print("Image downloaded and opened successfully.")
        else:
            print(f"Error: Image file not found at path:
{IMAGE_PATH}")
            exit(1)
    else:
        image_pil = Image.open(IMAGE_PATH).convert("RGB")
        print("Image opened successfully from local path.")

    # Apply the CLIPS (open_clip) preprocessing steps
    # Note: CLIPS preprocess might differ slightly from OpenAI CLIP
preprocess
    image_input = clips_preprocess(image_pil).unsqueeze(0).to(device)
    print(f"Image preprocessed using CLIPS preprocessor. Tensor shape:
{image_input.shape}")

except FileNotFoundError:
    print(f"Error: Image file not found at path: {IMAGE_PATH}")
    exit(1)
except requests.exceptions.RequestException as e:
    print(f"Error downloading image from URL {IMAGE_PATH}: {e}")
    exit(1)
except Exception as e:
    print(f"Error opening or preprocessing image: {e}")
    exit(1)


# --- Step 3: Tokenize Text Descriptions using CLIPS Tokenizer ---
print("\nTokenizing text descriptions using CLIPS tokenizer...")
try:
    # Use the CLIPS (open_clip) tokenizer
    # Determining context_length might require checking model config
if not a direct attribute
    context_length = clips_model.context_length if
hasattr(clips_model, 'context_length') else 77 # Default if
unavailable
    text_inputs = clips_tokenizer(TEXT_DESCRIPTIONS,
context_length=context_length).to(device)
    print(f"Text descriptions tokenized into tensor shape:
{text_inputs.shape}")

except Exception as e:
```

```python
        print(f"Error tokenizing text: {e}")
        exit(1)


# --- Step 4: Generate Embeddings and Calculate Similarities using
CLIPS ---
print("\nCalculating image and text features using CLIPS model...")
# Use autocast for potential speedup with mixed precision (common with
open_clip models)
# Use torch.no_grad() as we are only doing inference
with torch.no_grad(), torch.cuda.amp.autocast() if device == 'cuda'
else contextlib.nullcontext():
    try:
        # Encode image and text using CLIPS model
        image_features = clips_model.encode_image(image_input)
        text_features = clips_model.encode_text(text_inputs)
        print(f"CLIPS Image features shape: {image_features.shape}")
        print(f"CLIPS Text features shape: {text_features.shape}")

        # Normalize the features using torch.nn.functional.normalize
        # (as shown in the CLIPS README example)
        image_features = F.normalize(image_features, dim=-1)
        text_features = F.normalize(text_features, dim=-1)
        print("Features normalized.")

        # Calculate cosine similarity and scale (using fixed 100.0 as
per CLIPS example)
        # Unlike original CLIP, open_clip models might not have a
learnable logit_scale
        # The CLIPS example explicitly uses 100.0 * features @
features.T
        temperature = 100.0
        similarities = temperature * image_features @ text_features.T
        print(f"Cosine similarities calculated (scaled by
{temperature}).")

        # Convert similarities to probabilities using softmax
        probabilities = similarities.softmax(dim=-1)
        print("Probabilities calculated via softmax.")

    except Exception as e:
        print(f"\nError during feature encoding or similarity
calculation with CLIPS: {e}")
        exit(1)

# --- Step 5: Display Results ---
print(f"\n--- CLIPS ({MODEL_DISPLAY_NAME}) Similarity Scores
(Probabilities) ---")

# Move probabilities to CPU for printing/analysis
```

```python
probabilities_cpu = probabilities.cpu().float().numpy().squeeze() #
Use float() for compatibility

# Handle scalar case if only one description was used
if probabilities_cpu.ndim == 0:
    probabilities_cpu = [probabilities_cpu.item()]
else:
    probabilities_cpu = probabilities_cpu.tolist()

if len(TEXT_DESCRIPTIONS) != len(probabilities_cpu):
     print("\nWarning: Mismatch between number of descriptions and
calculated probabilities!")
else:
    # Pair descriptions with scores and sort
    results = sorted(zip(TEXT_DESCRIPTIONS, probabilities_cpu),
key=lambda x: x[1], reverse=True)

    for description, prob in results:
        print(f"- '{description}': {prob:.4f}")

print("\n--- Task 6 Finished ---")
```

--- Starting Task 6: Image-Text Similarity Calculation using CLIPS-
Large-14-224 ---
Using pre-loaded CLIPS model on device: cuda:0

Loading and preprocessing image from: /kaggle/input/cv-ass-3-q1-3-
sample-image/sample_image.jpg
Image opened successfully from local path.
Image preprocessed using CLIPS preprocessor. Tensor shape:
torch.Size([1, 3, 224, 224])

Tokenizing text descriptions using CLIPS tokenizer...
Text descriptions tokenized into tensor shape: torch.Size([10, 80])

Calculating image and text features using CLIPS model...
CLIPS Image features shape: torch.Size([1, 768])
CLIPS Text features shape: torch.Size([10, 768])
Features normalized.
Cosine similarities calculated (scaled by 100.0).
Probabilities calculated via softmax.

--- CLIPS (CLIPS-Large-14-224) Similarity Scores (Probabilities) ---
- 'A person walking a large dog on a leash': 0.8172
- 'A sunny day with a pet and its owner': 0.1716
- 'A man and his golden retriever in a park': 0.0047
- 'An owner training their canine companion on grass': 0.0046
- 'Someone petting a furry animal outdoors': 0.0017
- 'A cat sleeping peacefully on a sofa': 0.0001
- 'Two friends enjoying a walk in nature': 0.0001

```
- 'Children playing soccer in a field': 0.0000
- 'A busy city street at night with neon lights': 0.0000
- 'A close-up portrait of a dog's happy face': 0.0000

--- Task 6 Finished ---
```

# Task 6: Output Analysis (CLIPS Model)

## Execution Flow

The log indicates the script ran as expected: it used the pre-loaded CLIPS model, loaded/processed the image, tokenized the text descriptions, encoded features, normalized them, calculated scaled similarities, and generated probabilities via softmax.

## Tensor Shapes

- Image tensor shape `[1, 3, 224, 224]` is correct for the 224x224 input model.
- Text tensor shape `[10, 80]` indicates the tokenizer used a context length of 80, which might be a default or specific setting for this `open_clip` tokenizer configuration (compared to the typical 77 for OpenAI CLIP). This is acceptable.
- Feature shapes `[1, 768]` and `[10, 768]` are correct. The `ViT-L` (Large) architecture used by this CLIPS model outputs 768-dimensional embeddings, distinct from the `ViT-B` (Base) model's 512 dimensions used in Task 3.

## Results Analysis

Let's compare the CLIPS probabilities with the original CLIP (`ViT-B/32`) probabilities from Task 3:

| Description | CLIP (`ViT-B/32`) Prob. | CLIPS (`ViT-L/14`) Prob. | Rank Change | Confidence Change |
|---|---|---|---|---|
| A person walking a large dog on a leash | 0.4197 | **0.8172** | Same (1st) | Much Higher |
| A sunny day with a pet and its owner | 0.3372 | **0.1716** | Same (2nd) | Lower |
| A man and his golden retriever in a park | 0.0061 | 0.0047 | ↑ (6->3) | Lower |
| An owner training their canine companion on grass | 0.0163 | 0.0046 | Same (4th) | Lower |
| Someone petting a furry animal outdoors | 0.2078 | 0.0017 | ↓ (3->5) | Much Lower |
| A cat sleeping peacefully on a sofa | 0.0002 | 0.0001 | ↑ (7->6) | Lower |
| Two friends enjoying a walk in nature | 0.0127 | 0.0001 | ↓ (5->7) | Much Lower |

| Description | CLIP (`ViT-B/32`) Prob. | CLIPS (`ViT-L/14`) Prob. | Rank Change | Confidence Change |
|---|---|---|---|---|
| Children playing soccer in a field | 0.0000 | 0.0000 | Same (Low) | Same |
| A busy city street at night with neon lights | 0.0000 | 0.0000 | Same (Low) | Same |
| A close-up portrait of a dog's happy face | 0.0000 | 0.0000 | Same (Low) | Same |

- **Increased Confidence:** CLIPS is much more confident in the top prediction (`'A person walking a large dog on a leash'`), assigning it over 81% probability compared to CLIP's 42%. This aligns with the idea that CLIPS, potentially benefiting from synthetic captions, might be better at latching onto core concepts (person, large dog) even if details (walking, leash, indoors) are mismatched.
- **Lower Confidence Elsewhere:** The increased confidence in the top choice significantly reduces the probability assigned to other plausible (but still inaccurate) descriptions like `'A sunny day...'` and `'Someone petting...'`.
- **Rank Changes:** There are minor shifts in the lower rankings, but the top two remain the same. `'A man and his golden retriever...'` moved up slightly relative to others, despite being incorrect.
- **Irrelevant Captions:** Both models effectively assign near-zero probability to completely irrelevant captions.

## Conclusion

Yes, the output is **correct** for the execution of the Task 6 code. The CLIPS model (`CLIPS-Large-14-224`) was loaded and used successfully. The results show plausible similarity scores, demonstrating different confidence levels and slightly different rankings compared to the original OpenAI CLIP model, which is expected behavior when using a different model architecture and training methodology (`ViT-L` vs `ViT-B`, CLIPS training vs CLIP training). The higher confidence in the top, partially correct caption is particularly noteworthy.

# Task 7: Commentary on CLIP vs. CLIPS Results

Here's a comparison of the results obtained from the standard OpenAI CLIP model (`ViT-B/32`) in Task 3 and the CLIPS model (`ViT-L/14-224`) in Task 6 for the given image and text descriptions:

## Similarities

1. **Agreement on Best Matches:** Both models identified the same two descriptions as the most relevant (albeit with different confidence levels):
   - `'A person walking a large dog on a leash'` (Rank 1 for both)
   - `'A sunny day with a pet and its owner'` (Rank 2 for both) This indicates a shared basic understanding of the core subject matter (a person with a pet dog).

2. **Rejection of Irrelevant Captions:** Both models effectively assigned near-zero probability scores to clearly irrelevant descriptions like `'A cat sleeping peacefully...'`, `'A busy city street...'`, and `'Children playing soccer...'`. This demonstrates strong performance in distinguishing relevant from completely unrelated concepts.

## Differences

1. **Confidence Levels & Distribution:**
   – **CLIPS showed much higher confidence** in its top prediction (81.7%) compared to CLIP (42.0%).
   – This resulted in a **sharper probability distribution** for CLIPS, heavily favoring the top match. CLIP's probabilities were more distributed among the top three somewhat plausible (though inaccurate) descriptions.
2. **Handling of Details vs. Core Concepts:**
   – CLIPS's strong preference for `'A person walking a large dog on a leash'` suggests it might prioritize matching key nouns ("person", "large dog") very strongly, potentially overlooking inaccuracies in the action ("walking" vs. "holding") or context ("leash", "indoors" vs. assumed "outdoors").
   – CLIP, while still ranking this description first, assigned it lower confidence, potentially indicating slightly more sensitivity to these conflicting details, distributing the remaining probability more evenly among other options containing related concepts ("pet", "owner", "furry animal").
3. **Ranking Variations:** While the top 2 ranks matched, there were minor differences in the mid-to-lower ranks for partially relevant descriptions (e.g., `'Someone petting...'` ranked 3rd for CLIP but 5th for CLIPS). This reflects subtle differences in how each model weighs the various elements within the descriptions against the image features.
4. **Underlying Model Differences:** These variations stem from fundamental differences:
   – **Architecture:** CLIPS used a larger `ViT-L/14` model, while CLIP used `ViT-B/32`. Larger models generally have greater capacity.
   – **Training:** CLIP used standard web-scraped image-text pairs. CLIPS employs an enhanced framework, possibly leveraging synthetic captions, aiming for improved zero-shot performance. The increased confidence and focus on core objects observed in CLIPS *might* be a result of this different training strategy.

## Conclusion

Both CLIP and CLIPS successfully identified the main subjects in the image and discarded irrelevant text. However, CLIPS (`ViT-L/14-224`) demonstrated significantly higher confidence in its top prediction compared to the standard CLIP (`ViT-B/32`), concentrating most of the probability mass there. This suggests CLIPS might be more decisive in identifying core concepts, potentially due to its larger architecture and enhanced training methodology, even if it sometimes overlooks finer contextual or action-related details that don't perfectly align with the top-matching concept. CLIP showed a more "cautious" distribution across plausible options.

----------------------------------------------------------------------------------

# Q2 - BLIP (Bootstrapping Language-Image Pre-training)

## 2.1 Installing dependencies and pre-trained weights.

```python
# Cell 1: Installation
import os
print("Installing/Updating transformers and dependencies...")
# Using --quiet to make the output cleaner
!pip install --quiet transformers torch torchvision Pillow accelerate
print("Installation complete.")

# Verify installation and versions
print("\nVerifying package versions...")
!pip show transformers torch Pillow accelerate | grep -E '^Name:|^Version:'

# Set TOKENIZERS_PARALLELISM to false to potentially avoid
warnings/issues in some environments
os.environ["TOKENIZERS_PARALLELISM"] = "false"
print("\nTOKENIZERS_PARALLELISM set to false.")
```

```
Installing/Updating transformers and dependencies...
Installation complete.

Verifying package versions...
Name: transformers
Version: 4.51.1
Name: torch
Version: 2.5.1+cu124
Name: pillow
Version: 11.1.0
Name: accelerate
Version: 1.3.0

TOKENIZERS_PARALLELISM set to false.
```

# Task 1: Output Analysis (installing dependencies and pre-trained weights)

## Environment Setup Summary

### Execution Status
- The `pip install` command completed successfully.

### Package Installation
- Core required libraries were installed or confirmed to be present:
  - `transformers` **v4.51.1**
  - `torch` **v2.5.1+cu124**
  - `Pillow` **v11.1.0**
  - `accelerate` **v1.3.0**

### Dependency Conflicts
- `pip` reported **dependency conflicts** related to:
  - `pylibcugraph-cu12`
  - `pylibraft-cu12`
  - `rmm-cu12`

These are components of the **RAPIDS library suite**, which is often pre-installed in Kaggle GPU environments.

> **Note:** These conflicts do **not** involve the packages used for this task (`transformers`, `torch`, etc.). Hence, they are **unlikely to affect BLIP model functionality**.

### Environment Variable
- `TOKENIZERS_PARALLELISM` was successfully set to `false` to suppress tokenizer-related parallelism warnings.

### Conclusion
- Despite unrelated dependency conflict warnings, all **necessary libraries for the Visual Question Answering (VQA)** task using **BLIP** have been successfully set up.
- The environment is **ready to proceed** to model inference.

# 2.2 For the previous sample image of human and dog, generate an answer to the question "Where is the dog present in the image?".

```
# Cell 2: Visual Question Answering
import torch
from PIL import Image
```

```python
import requests # Import requests, might be useful if loading from URL
later
import os
from transformers import BlipProcessor, BlipForQuestionAnswering

# --- Configuration ---
# Model identifier from Hugging Face Hub for VQA
model_id = "Salesforce/blip-vqa-base"

# Input image path and the specific question for this part
image_path =
"/kaggle/input/cv-ass-3-q1-3-sample-image/sample_image.jpg"
question = "Where is the dog present in the image?" # Question for
this part

# --- Device Setup ---
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print(f"Using device: {device}")

# --- Load Processor and Model ---
print(f"Loading processor and model: {model_id}...")
try:
    # Load the processor (handles image preprocessing and text
tokenization)
    processor = BlipProcessor.from_pretrained(model_id)
    # Load the VQA model
    model =
BlipForQuestionAnswering.from_pretrained(model_id).to(device)
    model.eval() # Set model to evaluation mode
    print("Processor and model loaded successfully.")
except Exception as e:
    print(f"Error loading model or processor: {e}")
    print("Ensure the model ID is correct and internet connectivity is
enabled in Kaggle settings.")
    # Depending on the error, you might need to stop execution
    # import sys
    # sys.exit(1)

# --- Image Loading ---
print(f"Loading image from: {image_path}...")
if not os.path.exists(image_path):
    print(f"ERROR: Image file not found at: {image_path}")
    # Handle error appropriately, e.g., skip the rest of the cell
else:
    try:
        raw_image = Image.open(image_path).convert('RGB')
        print("Image loaded successfully.")

        # --- Preprocessing and Inference ---
        print("Preprocessing image and question...")
```

```python
        # Prepare inputs using the processor
        inputs = processor(raw_image, question,
return_tensors="pt").to(device)
        print("Preprocessing complete.")

        print("Generating answer...")
        with torch.no_grad(): # Disable gradient calculation for
inference
            # Generate output token IDs
            output_ids = model.generate(**inputs, max_new_tokens=20) #
Limit max generated tokens

        print("Decoding answer...")
        # Decode the generated token IDs back to text
        answer = processor.batch_decode(output_ids,
skip_special_tokens=True)[0].strip()
        print("Decoding complete.")

        # --- Output ---
        print("-" * 30)
        print(f"Image Path: {image_path}")
        print(f"Question: {question}")
        print(f"Predicted Answer: {answer}")
        print("-" * 30)

    except FileNotFoundError:
         print(f"Operation halted: Image file not found at
{image_path}")
    except Exception as e:
        print(f"An error occurred during image loading, processing, or
inference: {e}")

# --- End of Cell ---

2025-04-18 12:39:28.300204: E
external/local_xla/xla/stream_executor/cuda/cuda_fft.cc:477] Unable to
register cuFFT factory: Attempting to register factory for plugin
cuFFT when one has already been registered
WARNING: All log messages before absl::InitializeLog() is called are
written to STDERR
E0000 00:00:1744979968.757790      19 cuda_dnn.cc:8310] Unable to
register cuDNN factory: Attempting to register factory for plugin
cuDNN when one has already been registered
E0000 00:00:1744979968.880296      19 cuda_blas.cc:1418] Unable to
register cuBLAS factory: Attempting to register factory for plugin
cuBLAS when one has already been registered
Using a slow image processor as `use_fast` is unset and a slow
processor was saved with this model. `use_fast=True` will be the
default behavior in v4.52, even if the model was saved with a slow
```

```
processor. This will result in minor differences in outputs. You'll
still be able to use a slow processor with `use_fast=False`.

Using device: cuda
Loading processor and model: Salesforce/blip-vqa-base...
```

```
{"model_id":"807dd15c1a9b40f09cf6ada5c057e5b1","version_major":2,"vers
ion_minor":0}

{"model_id":"c21deff51e3748f1afc74b7468174e55","version_major":2,"vers
ion_minor":0}

{"model_id":"db542e17d69742c793c9e6eb1e2b6675","version_major":2,"vers
ion_minor":0}

{"model_id":"ae453c4671944d4ebb53d050bb35879a","version_major":2,"vers
ion_minor":0}

{"model_id":"1746cd6e3a734df6be9ba6b97baa1f3c","version_major":2,"vers
ion_minor":0}

{"model_id":"b41f1ffe5eb343c0b64df98e662a7ba0","version_major":2,"vers
ion_minor":0}

{"model_id":"bae5c272bdf74649b6d3708cfcd63f34","version_major":2,"vers
ion_minor":0}
```

```
Processor and model loaded successfully.
Loading image from:
/kaggle/input/cv-ass-3-q1-3-sample-image/sample_image.jpg...
Image loaded successfully.
Preprocessing image and question...
Preprocessing complete.
Generating answer...
Decoding answer...
Decoding complete.
------------------------------
Image Path: /kaggle/input/cv-ass-3-q1-3-sample-image/sample_image.jpg
Question: Where is the dog present in the image?
Predicted Answer: in man ' s arms
------------------------------
```

# Task 2: Output Analysis (answer to the sample question)

## Device Selection:

The code correctly identified and selected the **cuda device** for **GPU acceleration**.

## Model Loading:

The **BlipProcessor** and **BlipForQuestionAnswering model** (**Salesforce/blip-vqa-base**) were successfully downloaded from the **Hugging Face Hub** (indicated by download progress bars for configs, tokenizer files, and the **1.54G model weights**) and loaded onto the GPU.

## Image Handling:

The image from the specified path (**/kaggle/input/cv-ass-3-q1-3-sample-image/sample_image.jpg**) was loaded successfully.

## Preprocessing & Inference:

The image and question were successfully **preprocessed**, and the model generated an **answer** without errors.

## Decoding & Output:

The generated tokens were **decoded** into the text answer: **"man's arms"**.

## Result:

The model provided a **relevant** and **plausible answer** to the question "**Where is the dog present in the image?**". The answer directly addresses the location of the dog relative to the other main subject often depicted in this common sample image.

## Conclusion:

The script executed successfully from start to finish, performed the **VQA task correctly** using the loaded **BLIP model**, and produced a **meaningful** and **accurate** answer based on the visual content.

# 2.3 "Where is the man present in the image"

```python
# Cell 3: Visual Question Answering (Second Question)
import torch
from PIL import Image
import requests # Import requests, might be useful if loading from URL
later
import os
from transformers import BlipProcessor, BlipForQuestionAnswering
import logging # Import logging to potentially reduce verbosity of
lower-level libraries

# --- Configuration ---
# Set transformers logging level to ERROR to hide informational
messages if desired
# logging.getLogger("transformers").setLevel(logging.ERROR)

# Model identifier from Hugging Face Hub for VQA
```

```python
model_id = "Salesforce/blip-vqa-base"

# Input image path (same as before) and the new question
image_path =
"/kaggle/input/cv-ass-3-q1-3-sample-image/sample_image.jpg"
question = "Where is the man present in the image?" # <<< New question
for this part

# --- Device Setup ---
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print(f"Using device: {device}")

# --- Load Processor and Model (reuse if already loaded in the same
session) ---
# Check if model and processor already exist in the environment to
avoid reloading
if 'model' not in locals() or 'processor' not in locals():
    print(f"Loading processor and model: {model_id}...")
    try:
        processor = BlipProcessor.from_pretrained(model_id)
        model =
BlipForQuestionAnswering.from_pretrained(model_id).to(device)
        model.eval() # Set model to evaluation mode
        print("Processor and model loaded successfully.")
    except Exception as e:
        print(f"Error loading model or processor: {e}")
        print("Ensure the model ID is correct and internet
connectivity is enabled in Kaggle settings.")
        # Stop execution if loading fails
        raise SystemExit(f"Failed to load model/processor: {e}")
else:
    print("Processor and model already loaded.")
    model.to(device) # Ensure model is on the correct device if re-
running cell
    model.eval()      # Ensure model is in eval mode

# --- Image Loading ---
print(f"Loading image from: {image_path}...")
if not os.path.exists(image_path):
    print(f"ERROR: Image file not found at: {image_path}")
    # Handle error appropriately
    raise FileNotFoundError(f"Image file not found: {image_path}")
else:
    try:
        raw_image = Image.open(image_path).convert('RGB')
        print("Image loaded successfully.")

        # --- Preprocessing and Inference ---
        print("Preprocessing image and question...")
        # Prepare inputs using the processor
```

```python
        inputs = processor(raw_image, question,
return_tensors="pt").to(device)
        print("Preprocessing complete.")

        print("Generating answer...")
        with torch.no_grad(): # Disable gradient calculation for
inference
            # Generate output token IDs
            # Added a max_length or max_new_tokens constraint to
prevent overly long/runaway generation
            output_ids = model.generate(**inputs, max_new_tokens=20)

        print("Decoding answer...")
        # Decode the generated token IDs back to text
        answer = processor.batch_decode(output_ids,
skip_special_tokens=True)[0].strip()
        print("Decoding complete.")

        # --- Output ---
        print("-" * 30)
        print(f"Image Path: {image_path}")
        print(f"Question: {question}")
        print(f"Predicted Answer: {answer}")
        print("-" * 30)

    except FileNotFoundError:
        # This specific exception was already handled above, but kept
for structure
        print(f"Operation halted: Image file not found at
{image_path}")
    except Exception as e:
        print(f"An error occurred during image loading, processing, or
inference: {e}")
        # Raise the exception to make it clear execution failed
        raise e

# --- End of Cell ---
```

```
Using device: cuda
Processor and model already loaded.
Loading image from:
/kaggle/input/cv-ass-3-q1-3-sample-image/sample_image.jpg...
Image loaded successfully.
Preprocessing image and question...
Preprocessing complete.
Generating answer...
Decoding answer...
Decoding complete.
------------------------------
Image Path: /kaggle/input/cv-ass-3-q1-3-sample-image/sample_image.jpg
```

```
Question: Where is the man present in the image?
Predicted Answer: living room
----------------------------
```

# Task 3: Output Analysis (VQA – Second Question)
- **Execution Status:** The script completed successfully.
- **Device Selection:** Correctly used the `cuda` device.
- **Model/Processor Loading:** Correctly identified that the model and processor were already loaded from the previous step, saving time.
- **Image Loading:** The same image (`/kaggle/input/cv-ass-3-q1-3-sample-image/sample_image.jpg`) was loaded successfully.
- **Question:** The question processed was "Where is the man present in the image?".
- **Preprocessing & Inference:** These steps completed without errors.
- **Decoding & Output:** The model generated the answer "living room".
- **Result:** The predicted answer "living room" is a highly plausible and contextually appropriate description of the man's location in the image. The visual cues (bookshelf, wooden floor, indoor setting, doorway) strongly suggest a residential room like a living room, study, or den.
- **Conclusion:** The model correctly interpreted the visual scene and provided a relevant and accurate answer to the question about the man's location. The VQA process was successful.

# 2.4 Output and accuracy of task 2 and task 3.

## Comment on Output and Accuracy (Previous Two Questions)

The BLIP VQA model (`Salesforce/blip-vqa-base`) demonstrated **high accuracy and relevance** in answering both questions based on the provided image.

1. **Question 1: "Where is the dog present in the image?"**
   - **Output:** "in man ' s arms"
   - **Accuracy:** This answer is **highly accurate and specific**. The image clearly shows the man holding the large dog in his arms. The model correctly identified the relationship and the location of the dog relative to the man.
2. **Question 2: "Where is the man present in the image?"**
   - **Output:** "living room"
   - **Accuracy:** This answer is **accurate and contextually appropriate**. While the image doesn't explicitly label the room, the presence of a large bookshelf, wooden flooring, and an interior setting strongly implies a residential room like a living room, study, or den. "Living room" is a very plausible inference based on these visual cues. The model successfully interpreted the overall scene context to determine the man's general location.

**Overall:** The model successfully processed the image and understood both questions, providing concise, relevant, and accurate answers. It demonstrated the ability to identify not only the

relative position of objects (dog in arms) but also to infer the broader environmental context (living room).

---------------------------------------------------------

----------------------

# Q3 BLIP vs CLIP

## 3.1 Loading BLIP weights for image captioning.

```python
# Cell for Q3.1: Load BLIP Image Captioning Model

import torch
from PIL import Image
import requests # Good practice import
import os
from transformers import BlipProcessor, BlipForConditionalGeneration

# --- Configuration ---
# Model identifier from Hugging Face Hub for Image Captioning (Base
model)
caption_model_id = "Salesforce/blip-image-captioning-base"

# --- Device Setup ---
# Assuming 'device' was defined in previous cells (Q2)
if 'device' not in locals():
    device = torch.device("cuda" if torch.cuda.is_available() else
"cpu")
    print(f"Device not found in locals, setting to: {device}")
else:
    print(f"Using existing device: {device}")

# --- Load Processor and Model for Captioning ---
print(f"\nLoading processor and model for Image Captioning:
{caption_model_id}...")
try:
    # Use distinct variable names to avoid overwriting VQA components
if needed later
    caption_processor =
BlipProcessor.from_pretrained(caption_model_id)
    caption_model =
BlipForConditionalGeneration.from_pretrained(caption_model_id).to(devi
ce)
    caption_model.eval() # Set model to evaluation mode
    print("BLIP Image Captioning processor and model loaded
```

```
successfully.")

    # --- Verification (Optional) ---
    print(f" - Caption Processor Type: {type(caption_processor)}")
    print(f" - Caption Model Type: {type(caption_model)}")
    param_device = next(caption_model.parameters()).device
    print(f" - Caption Model Parameter device: {param_device}")
    if not str(param_device).startswith(str(device)):
        print(f"   Warning: Parameter device ({param_device}) does
not seem to match target device ({device})!")

except Exception as e:
    print(f"\nError loading BLIP captioning model or processor: {e}")
    print("Ensure the model ID is correct and internet connectivity is
enabled.")
    # Optionally raise the error if this step is critical for
subsequent ones
    # raise e

# --- End of Cell ---
```

Using existing device: cuda

Loading processor and model for Image Captioning: Salesforce/blip-
image-captioning-base...

{"model_id":"743fcd6a42274bd18a603d95d5bae78c","version_major":2,"version_minor":0}

{"model_id":"adca529e760f43cb92c0e57064c3af97","version_major":2,"version_minor":0}

{"model_id":"353f09e853de497fa50b45f4ee604d57","version_major":2,"version_minor":0}

{"model_id":"39f4dbe269464039be79cfe8ecc93a1f","version_major":2,"version_minor":0}

{"model_id":"bb7c92c83a0540a890508b3b7c7517b2","version_major":2,"version_minor":0}

{"model_id":"5c1e8db04e4049b292532d111347f8b6","version_major":2,"version_minor":0}

{"model_id":"56ea6eb1740c477b919b3e2284238ef6","version_major":2,"version_minor":0}

{"model_id":"627f53f63dd6482ab26c98a84b542a5e","version_major":2,"version_minor":0}

BLIP Image Captioning processor and model loaded successfully.
 - Caption Processor Type: <class
```

```
'transformers.models.blip.processing_blip.BlipProcessor'>
 - Caption Model Type: <class
'transformers.models.blip.modeling_blip.BlipForConditionalGeneration'>
 - Caption Model Parameter device: cuda:0
```

## Analysis of Q3.1 Output (Loading BLIP Captioning Model)

- **Device Selection:** The output confirms the code correctly identified and used the existing `cuda:0` device.
- **Model/Processor Loading:** Download progress bars and logs show that the components for the `Salesforce/blip-image-captioning-base` model (preprocessor config, tokenizer files, model config, and model weights - `pytorch_model.bin` or `model.safetensors` at ~990MB) were successfully downloaded from the Hugging Face Hub.
- **Success Confirmation:** The message `BLIP Image Captioning processor and model loaded successfully.` is present, indicating the loading process completed without raising exceptions.
- **Verification:** The printed types (`BlipProcessor`, `BlipForConditionalGeneration`) are correct for BLIP captioning via the `transformers` library. The model's parameters are confirmed to be on the target device (`cuda:0`).
- **Conclusion:** The output clearly indicates that the code **executed successfully**. The pre-trained BLIP image captioning model and its associated processor were correctly loaded onto the GPU and are ready for use.

# 3.2 Generating captions for sample images

```python
# Cell for Q3.2: Generate Captions for Sample Images

import torch
from PIL import Image
import requests
import os
from transformers import BlipProcessor, BlipForConditionalGeneration
import logging

# --- Configuration ---
image_dir = "/kaggle/input/cv-ass-3-q3-sample-images/samples"
valid_image_extensions = ('.jpg', '.jpeg', '.png', '.bmp', '.gif',
'.webp')

# --- Ensure Model and Processor are loaded ---
# Check if variables from Q3.1 exist. If not, attempt to load them.
if 'caption_model' not in locals() or 'caption_processor' not in
locals():
    print("Captioning model/processor not found in environment.
Attempting to load...")
    caption_model_id = "Salesforce/blip-image-captioning-base"
    try:
```

```python
        if 'device' not in locals():
            device = torch.device("cuda" if torch.cuda.is_available()
else "cpu")
            print(f"Device not found, setting to: {device}")
        else:
            print(f"Using existing device: {device}")

        caption_processor =
BlipProcessor.from_pretrained(caption_model_id)
        caption_model =
BlipForConditionalGeneration.from_pretrained(caption_model_id).to(devi
ce)
        caption_model.eval()
        print("Captioning model and processor loaded successfully.")
    except Exception as e:
        print(f"\nError loading BLIP captioning model or processor:
{e}")
        raise SystemExit("Cannot proceed without captioning
model/processor.")
else:
    # Assume model and processor are loaded correctly from Q3.1
    # Ensure they are on the correct device and in eval mode
    try:
        if 'device' not in locals():
            device = torch.device("cuda" if torch.cuda.is_available()
else "cpu")
            print(f"Device not found, setting to: {device}")
        else:
            print(f"Using existing device: {device}")
        caption_model.to(device)
        caption_model.eval()
        print("Using pre-loaded captioning model and processor.")
    except Exception as e:
         print(f"Error ensuring pre-loaded model is ready: {e}")
         raise SystemExit("Cannot proceed.")


# --- Image Processing and Caption Generation ---
generated_captions = {}

print(f"\nScanning directory: {image_dir}")
if not os.path.isdir(image_dir):
    print(f"Error: Directory not found - {image_dir}")
else:
    image_files = [f for f in os.listdir(image_dir) if
f.lower().endswith(valid_image_extensions)]
    print(f"Found {len(image_files)} image files.")

    if not image_files:
        print("No images found in the directory to process.")
```

```python
    else:
        for filename in image_files:
            image_path = os.path.join(image_dir, filename)
            print("-" * 30)
            print(f"Processing: {filename}")
            try:
                # Load Image
                raw_image = Image.open(image_path).convert('RGB')

                # Prepare image for model
                # Option 1: Unconditional captioning (no text prompt)
                inputs = caption_processor(images=raw_image,
return_tensors="pt").to(device)
                pixel_values = inputs.pixel_values

                # Option 2: Conditional captioning (if you wanted to
provide a prompt)
                # text = "a photography of"
                # inputs = caption_processor(raw_image, text,
return_tensors="pt").to(device)

                # Generate caption
                with torch.no_grad():
                    output_ids =
caption_model.generate(pixel_values=pixel_values, max_length=50,
num_beams=3) # Using beam search

                # Decode caption
                # Use decode for single sequence, batch_decode for
list
                caption = caption_processor.decode(output_ids[0],
skip_special_tokens=True)
                caption = caption.strip() # Clean up whitespace

                generated_captions[filename] = caption
                print(f"  Generated Caption: {caption}")

            except FileNotFoundError:
                print(f"  Error: File not found at {image_path}")
            except Exception as e:
                print(f"  Error processing {filename}: {e}")

print("\n--- Caption Generation Complete ---")
# You can access the captions later using the `generated_captions`
dictionary
# print(generated_captions)

    # --- End of Cell ---
```

```
Using existing device: cuda
Using pre-loaded captioning model and processor.

Scanning directory: /kaggle/input/cv-ass-3-q3-sample-images/samples
Found 10 image files.
------------------------------
Processing: ILSVRC2012_test_00000004.jpg
  Generated Caption: a small dog running across a green field
------------------------------
Processing: ILSVRC2012_test_00000022.jpg
  Generated Caption: a small white and brown dog standing next to a
pool
------------------------------
Processing: ILSVRC2012_test_00000023.jpg
  Generated Caption: a man riding a bike in the rain
------------------------------
Processing: ILSVRC2012_test_00000026.jpg
  Generated Caption: a man in a suit and tie sitting on a couch
------------------------------
Processing: ILSVRC2012_test_00000018.jpg
  Generated Caption: a group of kids sitting on a towel by a swimming
pool
------------------------------
Processing: ILSVRC2012_test_00000003.jpg
  Generated Caption: a small brown and white dog walking on a green
carpet
------------------------------
Processing: ILSVRC2012_test_00000019.jpg
  Generated Caption: a small bird sitting on top of a green plant
------------------------------
Processing: ILSVRC2012_test_00000030.jpg
  Generated Caption: a duck drinking water from a pond
------------------------------
Processing: ILSVRC2012_test_00000034.jpg
  Generated Caption: coffee being poured into a cup
------------------------------
Processing: ILSVRC2012_test_00000025.jpg
  Generated Caption: a brown butterfly sitting on top of green leaves

--- Caption Generation Complete ---
```

## Analysis of Q3.2 Output (BLIP Caption Generation)

- **Execution Status:** The script completed successfully, processing all 10 image files found in the specified directory (`/kaggle/input/cv-ass-3-q3-sample-images/samples`).
- **Model Usage:** It correctly utilized the pre-loaded BLIP image captioning model (`Salesforce/blip-image-captioning-base`) and processor from the previous step.

- **Process:** For each image, it successfully loaded, preprocessed, generated output tokens using beam search (`num_beams=3`), and decoded these tokens into natural language captions.
- **Generated Captions:** The following captions were generated for the respective images:
  - `ILSVRC2012_test_00000004.jpg`: "a small dog running across a green field"
  - `ILSVRC2012_test_00000022.jpg`: "a small white and brown dog standing next to a pool"
  - `ILSVRC2012_test_00000023.jpg`: "a man riding a bike in the rain"
  - `ILSVRC2012_test_00000026.jpg`: "a man in a suit and tie sitting on a couch"
  - `ILSVRC2012_test_00000018.jpg`: "a group of kids sitting on a towel by a swimming pool"
  - `ILSVRC2012_test_00000003.jpg`: "a small brown and white dog walking on a green carpet"
  - `ILSVRC2012_test_00000019.jpg`: "a small bird sitting on top of a green plant"
  - `ILSVRC2012_test_00000030.jpg`: "a duck drinking water from a pond"
  - `ILSVRC2012_test_00000034.jpg`: "coffee being poured into a cup"
  - `ILSVRC2012_test_00000025.jpg`: "a brown butterfly sitting on top of green leaves"
- **Quality Assessment:** Based on typical content associated with ILSVRC test images and the user's positive verification, the generated captions appear highly relevant and descriptive of the likely image content. They identify main subjects, actions, and sometimes context (e.g., "green field", "pool", "rain", "couch", "pond"). The level of detail seems appropriate for automatic image captioning.
- **Conclusion:** The BLIP image captioning model performed successfully, generating relevant and descriptive captions for all sample images provided.

## 3.3 Evaluating semantic accuracy of BLIP generated captions using CLIP

```python
# Cell for Q3.3: Evaluate BLIP Captions using OpenAI CLIP

import torch
import torch.nn.functional as F
from PIL import Image
import clip # OpenAI CLIP library
import os
import requests
import logging

# --- Configuration ---
OPENAI_CLIP_MODEL_NAME = "ViT-B/32" # The model used in Q1
# Assuming image_dir was defined in Q3.2
if 'image_dir' not in locals():
    image_dir = "/kaggle/input/cv-ass-3-q3-sample-images/samples"
    print(f"image_dir not found, setting to: {image_dir}")
```

```python
# Check if BLIP generated captions exist from Q3.2
if 'generated_captions' not in locals() or not generated_captions:
    print("Error: 'generated_captions' dictionary not found or is
empty.")
    print("Please ensure the previous cell (Q3.2) generating BLIP
captions ran successfully.")
    raise NameError("Missing generated_captions. Cannot proceed.")
else:
    print(f"Found {len(generated_captions)} generated captions from
BLIP.")

# --- Device Setup ---
# Assuming 'device' was defined in previous cells
if 'device' not in locals():
    device = torch.device("cuda" if torch.cuda.is_available() else
"cpu")
    print(f"Device not found, setting to: {device}")
else:
    print(f"Using existing device: {device}")

# --- Step 1: Load OpenAI CLIP Model ---
# Reload OpenAI CLIP model and preprocessor to ensure we use the
correct one
# Use distinct variable names to avoid conflicts
print(f"\nLoading OpenAI CLIP model: {OPENAI_CLIP_MODEL_NAME}...")
try:
    openai_clip_model, openai_clip_preprocess =
clip.load(OPENAI_CLIP_MODEL_NAME, device=device)
    openai_clip_model.eval() # Set to evaluation mode
    print("OpenAI CLIP model and preprocessor loaded successfully.")
except Exception as e:
    print(f"Error loading OpenAI CLIP model: {e}")
    raise SystemExit("Cannot proceed without OpenAI CLIP model.")

# --- Step 2: Calculate Similarity Scores ---
clip_similarity_scores = {}
print("\nCalculating similarity using OpenAI CLIP...")

# Iterate through the images for which BLIP captions were generated
for filename, blip_caption in generated_captions.items():
    image_path = os.path.join(image_dir, filename)
    print("-" * 30)
    print(f"Processing: {filename}")
    print(f"  BLIP Caption: {blip_caption}")

    if not os.path.exists(image_path):
        print(f"  Error: Image file not found at {image_path}.
Skipping.")
        continue
```

```python
    try:
        # Load and preprocess image using OpenAI CLIP's preprocessor
        raw_image = Image.open(image_path).convert("RGB")
        image_input =
openai_clip_preprocess(raw_image).unsqueeze(0).to(device)

        # Tokenize the BLIP-generated caption using OpenAI CLIP's
tokenizer
        text_input = clip.tokenize([blip_caption]).to(device)

        # Calculate features
        with torch.no_grad():
            image_features =
openai_clip_model.encode_image(image_input)
            text_features = openai_clip_model.encode_text(text_input)

            # Normalize features
            image_features_norm = F.normalize(image_features, p=2,
dim=-1)
            text_features_norm = F.normalize(text_features, p=2, dim=-
1)

            # Calculate cosine similarity (dot product of normalized
features)
            # Result is a tensor, get the scalar value
            similarity = torch.matmul(image_features_norm,
text_features_norm.T).item()

        clip_similarity_scores[filename] = similarity
        print(f"  CLIP Cosine Similarity: {similarity:.4f}")

    except Exception as e:
        print(f"  Error calculating similarity for {filename}: {e}")
        clip_similarity_scores[filename] = None # Indicate failure

print("\n--- OpenAI CLIP Similarity Calculation Complete ---")
# The 'clip_similarity_scores' dictionary holds the results.
# print(clip_similarity_scores)

# --- Interpretation ---
print("\n--- Interpretation ---")
print("Cosine similarity scores range from -1 to 1.")
print("A score closer to 1 indicates high semantic similarity between
the image and the BLIP-generated caption, according to the OpenAI CLIP
model.")
print("A score closer to 0 (or negative) indicates low semantic
similarity.")
print("These scores reflect how well CLIP 'thinks' the generated
caption describes the image.")
# --- End of Cell ---
```

```
Found 10 generated captions from BLIP.
Using existing device: cuda

Loading OpenAI CLIP model: ViT-B/32...
OpenAI CLIP model and preprocessor loaded successfully.

Calculating similarity using OpenAI CLIP...
------------------------------
Processing: ILSVRC2012_test_00000004.jpg
  BLIP Caption: a small dog running across a green field
  CLIP Cosine Similarity: 0.3274
------------------------------
Processing: ILSVRC2012_test_00000022.jpg
  BLIP Caption: a small white and brown dog standing next to a pool
  CLIP Cosine Similarity: 0.3445
------------------------------
Processing: ILSVRC2012_test_00000023.jpg
  BLIP Caption: a man riding a bike in the rain
  CLIP Cosine Similarity: 0.3154
------------------------------
Processing: ILSVRC2012_test_00000026.jpg
  BLIP Caption: a man in a suit and tie sitting on a couch
  CLIP Cosine Similarity: 0.2888
------------------------------
Processing: ILSVRC2012_test_00000018.jpg
  BLIP Caption: a group of kids sitting on a towel by a swimming pool
  CLIP Cosine Similarity: 0.3438
------------------------------
Processing: ILSVRC2012_test_00000003.jpg
  BLIP Caption: a small brown and white dog walking on a green carpet
  CLIP Cosine Similarity: 0.3140
------------------------------
Processing: ILSVRC2012_test_00000019.jpg
  BLIP Caption: a small bird sitting on top of a green plant
  CLIP Cosine Similarity: 0.2722
------------------------------
Processing: ILSVRC2012_test_00000030.jpg
  BLIP Caption: a duck drinking water from a pond
  CLIP Cosine Similarity: 0.3054
------------------------------
Processing: ILSVRC2012_test_00000034.jpg
  BLIP Caption: coffee being poured into a cup
  CLIP Cosine Similarity: 0.2859
------------------------------
Processing: ILSVRC2012_test_00000025.jpg
  BLIP Caption: a brown butterfly sitting on top of green leaves
  CLIP Cosine Similarity: 0.3025

--- OpenAI CLIP Similarity Calculation Complete ---
```

```
--- Interpretation ---
Cosine similarity scores range from -1 to 1.
A score closer to 1 indicates high semantic similarity between the
image and the BLIP-generated caption, according to the OpenAI CLIP
model.
A score closer to 0 (or negative) indicates low semantic similarity.
These scores reflect how well CLIP 'thinks' the generated caption
describes the image.
```

## Analysis of Q3.3 Output (CLIP Evaluation of BLIP Captions)

- **Execution Status:** The script executed successfully.

- **Model Loading:** The OpenAI CLIP model (`ViT-B/32`) and its preprocessor were loaded correctly, ensuring the evaluation was performed using the intended CLIP model.

- **Process:** The script iterated through the 10 images and their corresponding captions previously generated by the BLIP captioning model. For each pair, it successfully:

  – Loaded and preprocessed the image using the OpenAI CLIP preprocessor.
  – Tokenized the BLIP-generated caption using the OpenAI CLIP tokenizer.
  – Encoded both image and text into normalized feature vectors using the OpenAI CLIP model.
  – Calculated the cosine similarity between the image and text features.

- **Similarity Scores:** The calculated cosine similarity scores for each image-caption pair were printed, ranging from approximately `0.2722` (for the bird image) to `0.3445` (for the dog by pool image).

- **Interpretation:**

  – The obtained scores (mostly in the 0.27 to 0.34 range) are moderately positive. On the scale of -1 to 1, these values indicate that the OpenAI CLIP model perceives a decent level of semantic correspondence between the images and the captions generated by BLIP.
  – They are significantly higher than scores expected for random or unrelated captions (which would be near 0 or negative) but are not extremely close to 1. This suggests that while the captions capture relevant elements according to CLIP, they might not be considered perfect or exhaustive matches in CLIP's embedding space. This could be due to differences in focus, detail level, or subtle nuances between how BLIP generates captions and how CLIP interprets image-text correspondence.
  – The variation in scores across images (e.g., slightly higher for the dogs/kids by pool, lower for the bird/man on couch) reflects CLIP's assessment of how well each specific BLIP caption matched its corresponding image.

- **Classroom Comment Relevance:** The analysis correctly computed and reported the cosine similarity, adhering to the TA's clarification (which is applicable here for interpreting CLIP-based similarity).

- **Conclusion:** The script successfully quantified the semantic similarity between the images and their BLIP-generated captions using the OpenAI CLIP model. The resulting moderate cosine similarity scores suggest a reasonable degree of accuracy in the BLIP captions as evaluated by CLIP.

**Classroom comments referred**

# 3.4 Using CLIPS to evaluate scores

```python
# Cell for Q3.4: Evaluate BLIP Captions using CLIPS

import torch
import torch.nn.functional as F
from PIL import Image
# Assuming open_clip was installed and imported in previous cells
(Q1.5/Q1.6)
# If not, the check below will handle it.
import os
import requests
import contextlib # For autocast with cpu option
import pkg_resources
import subprocess
import sys

# --- Configuration ---
# Assuming image_dir was defined in Q3.2
if 'image_dir' not in locals():
    image_dir = "/kaggle/input/cv-ass-3-q3-sample-images/samples"
    print(f"image_dir not found, setting to: {image_dir}")

# Check if BLIP generated captions exist from Q3.2
if 'generated_captions' not in locals() or not generated_captions:
    print("Error: 'generated_captions' dictionary not found or is
empty.")
    print("Please ensure the cell generating BLIP captions (Q3.2) ran
successfully.")
    raise NameError("Missing generated_captions. Cannot proceed.")
else:
    print(f"Found {len(generated_captions)} generated captions from
BLIP.")

# --- Device Setup ---
# Assuming 'device' was defined in previous cells
if 'device' not in locals():
    device = torch.device("cuda" if torch.cuda.is_available() else
"cpu")
```

```python
        print(f"Device not found, setting to: {device}")
else:
    print(f"Using existing device: {device}")

# --- Step 1: Ensure CLIPS Model Components are Loaded ---
# Check if variables from Q1.6 exist. If not, attempt to load them.
if 'clips_model' not in locals() or 'clips_preprocess' not in locals()
or 'clips_tokenizer' not in locals():
    print("\nCLIPS model components not found in environment.
Attempting to load...")
    # --- Configuration from Q1.5 ---
    MODEL_DISPLAY_NAME = "CLIPS-Large-14-224"
    MODEL_HF_ID = "hf-hub:UCSC-VLAA/ViT-L-14-CLIPS-224-Recap-DataComp-
1B"

    # --- Helper function for running install command (needed if
reloading) ---
    def run_install_command(command):
        print(f"Executing: {command}")
        try:
            full_command = f"{sys.executable} -m {command}"
            process = subprocess.Popen(full_command, shell=True,
stdout=subprocess.PIPE, stderr=subprocess.STDOUT, text=True)
            while True:
                output = process.stdout.readline()
                if output == '' and process.poll() is not None: break
                if output: print(output.strip())
            return_code = process.poll()
            if return_code != 0:
                print(f"\\nError: Command '{full_command}' failed with
return code {return_code}")
                return False
            print(f"Successfully executed: {command}")
            return True
        except Exception as e:
            print(f"\\nAn exception occurred while running command:
{command}\\nError: {e}")
            return False

    # Ensure open_clip is installed
    try:
        pkg_resources.get_distribution('open_clip_torch')
        print("'open_clip_torch' package found.")
    except pkg_resources.DistributionNotFound:
        print("'open_clip_torch' package not found. Attempting
installation...")
        if not run_install_command("pip install open_clip_torch"):
            print("\\nError: Failed to install 'open_clip_torch'.
Exiting.")
            exit(1)
```

```python
    # Import and Load
    try:
        from open_clip import create_model_from_pretrained,
get_tokenizer
        print("Imported 'open_clip' successfully.")
        print(f"Loading pre-trained CLIPS model '{MODEL_DISPLAY_NAME}'
({MODEL_HF_ID})...")
        clips_model, clips_preprocess =
create_model_from_pretrained(MODEL_HF_ID, device=device)
        print("Loading tokenizer...")
        clips_tokenizer = get_tokenizer(MODEL_HF_ID)
        clips_model.eval() # Set to eval mode
        print("CLIPS model, preprocessor, and tokenizer loaded
successfully.")
    except ImportError:
        print("\\nError: Failed to import 'open_clip'. Ensure
'open_clip_torch' is installed.")
        raise SystemExit("Cannot proceed without open_clip.")
    except Exception as e:
        print(f"\\nError loading CLIPS model: {e}")
        raise SystemExit("Cannot proceed without CLIPS model.")
else:
    # Assume components are loaded from Q1.6
    try:
        # Ensure model is on the correct device and in eval mode
        clips_model.to(device)
        clips_model.eval()
        print("Using pre-loaded CLIPS model, preprocessor, and
tokenizer.")
        # Verify device from parameters
        param_device = next(clips_model.parameters()).device
        print(f"CLIPS Model is on device: {param_device}")
        if not str(param_device).startswith(str(device)):
            print(f"   Warning: Parameter device ({param_device})
does not seem to match target device ({device})!")

    except Exception as e:
        print(f"Error ensuring pre-loaded CLIPS model is ready: {e}")
        raise SystemExit("Cannot proceed.")


# --- Step 2: Calculate Similarity Scores using CLIPS ---
clips_similarity_scores = {} # Use a different name to avoid confusion
print("\nCalculating similarity using CLIPS model...")

# Iterate through the images and their BLIP captions
for filename, blip_caption in generated_captions.items():
    image_path = os.path.join(image_dir, filename)
    print("-" * 30)
```

```python
    print(f"Processing: {filename}")
    print(f"  BLIP Caption: {blip_caption}")

    if not os.path.exists(image_path):
        print(f"  Error: Image file not found at {image_path}.
Skipping.")
        continue

    try:
        # Load and preprocess image using CLIPS preprocessor
        raw_image = Image.open(image_path).convert("RGB")
        image_input =
clips_preprocess(raw_image).unsqueeze(0).to(device)

        # Tokenize the BLIP-generated caption using CLIPS tokenizer
        # Determine context length (might differ from OpenAI CLIP)
        context_length = clips_model.context_length if
hasattr(clips_model, 'context_length') else 77
        text_input = clips_tokenizer([blip_caption],
context_length=context_length).to(device)

        # Calculate features using CLIPS model
        # Use autocast for potential speedup with mixed precision
        with torch.no_grad(), torch.cuda.amp.autocast() if
str(device).startswith('cuda') else contextlib.nullcontext():
            image_features = clips_model.encode_image(image_input)
            text_features = clips_model.encode_text(text_input)

            # Normalize features using torch.nn.functional.normalize
            image_features_norm = F.normalize(image_features, dim=-1)
            text_features_norm = F.normalize(text_features, dim=-1)

            # Calculate cosine similarity (dot product of normalized
features)
            similarity = torch.matmul(image_features_norm,
text_features_norm.T).item()

        clips_similarity_scores[filename] = similarity
        print(f"  CLIPS Cosine Similarity: {similarity:.4f}")

    except Exception as e:
        print(f"  Error calculating similarity for {filename} using
CLIPS: {e}")
        clips_similarity_scores[filename] = None # Indicate failure

print("\n--- CLIPS Similarity Calculation Complete ---")
# The 'clips_similarity_scores' dictionary holds the results.
# print(clips_similarity_scores)

# --- Interpretation ---
```

```python
print("\n--- Interpretation ---")
print("Cosine similarity scores range from -1 to 1.")
print("A score closer to 1 indicates high semantic similarity between
the image and the BLIP-generated caption, according to the CLIPS
model.")
print("A score closer to 0 (or negative) indicates low semantic
similarity.")
print("These scores reflect how well CLIPS 'thinks' the generated
caption describes the image.")

# --- End of Cell ---
```

Found 10 generated captions from BLIP.
Using existing device: cuda
Using pre-loaded CLIPS model, preprocessor, and tokenizer.
CLIPS Model is on device: cuda:0

Calculating similarity using CLIPS model...
-------------------------------
Processing: ILSVRC2012_test_00000004.jpg
  BLIP Caption: a small dog running across a green field
  CLIPS Cosine Similarity: 0.1914
-------------------------------
Processing: ILSVRC2012_test_00000022.jpg
  BLIP Caption: a small white and brown dog standing next to a pool
  CLIPS Cosine Similarity: 0.2050
-------------------------------
Processing: ILSVRC2012_test_00000023.jpg
  BLIP Caption: a man riding a bike in the rain
  CLIPS Cosine Similarity: 0.1729
-------------------------------
Processing: ILSVRC2012_test_00000026.jpg
  BLIP Caption: a man in a suit and tie sitting on a couch
  CLIPS Cosine Similarity: 0.1276
-------------------------------
Processing: ILSVRC2012_test_00000018.jpg
  BLIP Caption: a group of kids sitting on a towel by a swimming pool

/tmp/ipykernel_19/1129003262.py:136: FutureWarning:
`torch.cuda.amp.autocast(args...)` is deprecated. Please use
`torch.amp.autocast('cuda', args...)` instead.
  with torch.no_grad(), torch.cuda.amp.autocast() if
str(device).startswith('cuda') else contextlib.nullcontext():

  CLIPS Cosine Similarity: 0.1807
-------------------------------
Processing: ILSVRC2012_test_00000003.jpg
  BLIP Caption: a small brown and white dog walking on a green carpet
  CLIPS Cosine Similarity: 0.1886
-------------------------------

```
Processing: ILSVRC2012_test_00000019.jpg
  BLIP Caption: a small bird sitting on top of a green plant
  CLIPS Cosine Similarity: 0.1792
------------------------------
Processing: ILSVRC2012_test_00000030.jpg
  BLIP Caption: a duck drinking water from a pond
  CLIPS Cosine Similarity: 0.1652
------------------------------
Processing: ILSVRC2012_test_00000034.jpg
  BLIP Caption: coffee being poured into a cup
  CLIPS Cosine Similarity: 0.1182
------------------------------
Processing: ILSVRC2012_test_00000025.jpg
  BLIP Caption: a brown butterfly sitting on top of green leaves
  CLIPS Cosine Similarity: 0.1750

--- CLIPS Similarity Calculation Complete ---

--- Interpretation ---
Cosine similarity scores range from -1 to 1.
A score closer to 1 indicates high semantic similarity between the
image and the BLIP-generated caption, according to the CLIPS model.
A score closer to 0 (or negative) indicates low semantic similarity.
These scores reflect how well CLIPS 'thinks' the generated caption
describes the image.
```

## Analysis of Q3.4 Output (CLIPS Evaluation of BLIP Captions)

- **Execution Status:** The script completed successfully. A minor `FutureWarning` regarding `torch.cuda.amp.autocast` syntax was displayed but did not prevent execution.

- **Model Usage:** The script correctly identified and utilized the pre-loaded CLIPS model (`CLIPS-Large-14-224`), preprocessor, and tokenizer from previous steps. The model was confirmed to be on the `cuda:0` device.

- **Process:** Similar to the previous evaluation (Q3.3), the script iterated through the 10 images and their BLIP-generated captions. For each pair, it performed:

  - Image loading and preprocessing using the `clips_preprocess` function.
  - Text tokenization using the `clips_tokenizer`.
  - Image and text feature encoding using the `clips_model`.
  - L2 normalization of the features.
  - Calculation of the raw cosine similarity between the image and text features.
- **Similarity Scores:** The calculated CLIPS cosine similarity scores were printed for each pair:

  - `ILSVRC2012_test_00000004.jpg`: 0.1914
  - `ILSVRC2012_test_00000022.jpg`: 0.2050

- – `ILSVRC2012_test_00000023.jpg`: 0.1729
- – `ILSVRC2012_test_00000026.jpg`: 0.1276
- – `ILSVRC2012_test_00000018.jpg`: 0.1807
- – `ILSVRC2012_test_00000003.jpg`: 0.1886
- – `ILSVRC2012_test_00000019.jpg`: 0.1792
- – `ILSVRC2012_test_00000030.jpg`: 0.1652
- – `ILSVRC2012_test_00000034.jpg`: 0.1182
- – `ILSVRC2012_test_00000025.jpg`: 0.1750

- **Interpretation:**

  - – The CLIPS similarity scores are all positive but generally lower (ranging from ~0.12 to ~0.20) than those obtained using the OpenAI CLIP model in Q3.3 (which were ~0.27 to ~0.34) for the *exact same image-caption pairs*.
  - – This indicates that the CLIPS model perceives a *lower* degree of semantic similarity between the BLIP-generated captions and the images compared to the standard OpenAI CLIP model.
  - – Despite being lower, the scores are still positive, suggesting CLIPS doesn't find the captions completely unrelated, just less similar than OpenAI CLIP did.
  - – Reasons for this difference could include the different model architectures (`ViT-L` vs. `ViT-B`), different training datasets and objectives (CLIPS incorporating synthetic data), leading to variations in the learned embedding space and how semantic similarity is represented.

- **Classroom Comment Relevance:** The calculation adheres to the requirement of using cosine similarity, as clarified by the TA.

- **Conclusion:** The script successfully used the CLIPS model to evaluate the BLIP-generated captions via cosine similarity. The results show a consistent pattern of positive but lower similarity scores compared to the evaluation performed with the standard OpenAI CLIP model, highlighting potential differences in how these models represent and compare image-text semantics.

**classroom comments referred**

# 3.5 Metrics for Quantifying Alignment Between CLIP/CLIPS and BLIP Outputs

Quantifying the alignment between CLIP/CLIPS (which evaluate image-text similarity) and BLIP (which generates captions) involves measuring how well CLIP/CLIPS "agree" with the captions produced by BLIP for a given image. Here are several metrics suitable for this purpose:

## 1. Cosine Similarity Score (Absolute)
- **What it measures:** Calculates the cosine of the angle between the image embedding and the text embedding (of the BLIP-generated caption) in the shared CLIP/CLIPS embedding space. A score closer to 1 means the embeddings are directionally similar.
- **Calculation:**
  - – Encode the image using the CLIP/CLIPS image encoder.

- Encode the BLIP-generated caption using the corresponding CLIP/CLIPS text encoder.
- Normalize both feature vectors (L2 norm).
- Compute the dot product of the normalized vectors.
- *(This was performed in Q3.3 using OpenAI CLIP and Q3.4 using CLIPS).*

- **Interpretation:** Provides a direct measure of semantic similarity *as perceived by the specific CLIP/CLIPS model*. Higher scores indicate better alignment according to that evaluation model. Scores range from -1 (opposite) to 1 (identical direction).
- **When Most Useful:**
  - For a quick, quantitative check of whether a generated caption is relevant to the image according to CLIP/CLIPS.
  - For comparing the *absolute* level of perceived similarity across different images or captions evaluated by the *same* model (e.g., CLIP score for image A vs. image B).
  - Comparing how different evaluation models (e.g., OpenAI CLIP vs. CLIPS) score the *same* caption for an image.
- **Limitations:**
  - The absolute value can be hard to interpret without context or baseline comparisons.
  - It doesn't inherently compare the BLIP caption against alternative or potentially better captions.

## 2. Rank-Based Metrics (e.g., Recall@k, Mean Reciprocal Rank - MRR)

- **What they measure:** Assess how highly CLIP/CLIPS rank the BLIP-generated caption when compared against a set of *candidate* captions for the same image.
- **Calculation:**
  a. Define a set of candidate captions for an image (e.g., the BLIP caption + several human-written captions + distractors).
  b. Calculate the CLIP/CLIPS cosine similarity score between the image and *each* candidate caption.
  c. Rank the candidate captions based on these similarity scores (highest score = rank 1).
  d. Identify the rank of the BLIP-generated caption within this list.
  e. Compute metrics across a dataset of images:
     - **Recall@k:** Percentage of images where the BLIP caption's rank is within the top-k (e.g., rank <= k). Recall@1 checks if it's ranked as the best match.
     - **MRR:** The average of the reciprocal ranks (1/rank) of the BLIP captions across all images. Rewards ranking higher captions more strongly.
- **Interpretation:** Measures if CLIP/CLIPS consider the BLIP caption to be *relatively* better than other potential descriptions for the image.
- **When Most Useful:**
  - When explicitly comparing the BLIP caption against known ground truth or alternative descriptions.

- To evaluate if the generated caption is considered the *best* or among the *top* descriptions by CLIP/CLIPS.
- Comparing different captioning models based on how well their outputs rank according to a fixed evaluator like CLIP/CLIPS.

- **Limitations:** Requires curating a relevant set of candidate/distractor captions for each image, which can be labor-intensive.

## 3. Correlation Metrics (e.g., Spearman's Rho, Kendall's Tau)

- **What they measure:** The statistical correlation between two sets of rankings or scores over a dataset.
- **Calculation:**
  - **Scenario A (Evaluator Agreement):**
    - i. Calculate CLIP scores and CLIPS scores for the *same set* of image-(BLIP)caption pairs (as derived from Q3.3 & Q3.4).
    - ii. Compute the correlation (e.g., Spearman's rank correlation) between the list of CLIP scores and the list of CLIPS scores.
  - **Scenario B (Alignment with Human Judgment):**
    - i. Collect human ratings (e.g., on a 1-5 scale) for the quality/relevance of the BLIP captions for multiple images.
    - ii. Calculate CLIP/CLIPS similarity scores for the same image-caption pairs.
    - iii. Compute the correlation between the CLIP/CLIPS scores and the human ratings.
- **Interpretation:**
  - *Scenario A:* Measures the consistency between different automatic evaluators (CLIP vs. CLIPS). High positive correlation indicates they tend to agree on the relative quality of captions.
  - *Scenario B:* Measures how well the automatic CLIP/CLIPS scores align with human perception of caption quality.
- **When Most Useful:**
  - *Scenario A:* Assessing the robustness of automatic evaluation methods.
  - *Scenario B:* Validating whether an automatic metric (like CLIP score) is a good proxy for human judgment.
- **Limitations:**
  - Correlation doesn't imply causation or guarantee absolute quality.
  - Requires multiple data points (image-caption pairs) to be meaningful.
  - Scenario B requires collecting human annotations, which is resource-intensive.

## 4. Qualitative Analysis / Error Analysis

- **What it measures:** Subjective assessment of alignment and identification of specific agreement/disagreement patterns by inspecting examples.
- **Calculation:** Manually review:
  - The image.
  - The BLIP-generated caption.
  - The CLIP similarity score.
  - The CLIPS similarity score.

- Look for patterns: cases where scores are high but captions seem poor, scores are low but captions seem good, or where CLIP and CLIPS scores diverge significantly.
- **Interpretation:** Identifies specific strengths, weaknesses, and failure modes in how BLIP generates captions and how CLIP/CLIPS evaluate them. Provides crucial context that numerical scores alone lack.
- **When Most Useful:** Should *always* be used alongside quantitative metrics to understand *why* the scores are the way they are, to gain deeper insights, and to guide model improvements.
- **Limitations:** Subjective, time-consuming, not easily scalable across large datasets.

## Summary
- Use **Cosine Similarity** for direct, per-instance semantic relevance scores.
- Use **Rank-Based Metrics** for comparing against alternatives or ground truth captions.
- Use **Correlation Metrics** for assessing evaluator agreement or alignment with human judgments.
- Always use **Qualitative Analysis** to interpret quantitative results and understand nuances.

**classroom comments referred**: The TA's clarification that cosine similarity is the desired score (for Q1.3) reinforces its use as a primary, direct metric (Metric 1 here) for quantifying alignment as perceived by CLIP/CLIPS.

-------------------------------------------------

----------------------------