

# Giải thích chi tiết ZeroShotImageClassificationPipeline

## Tổng quan

`ZeroShotImageClassificationPipeline` là một class trong thư viện Transformers của Hugging Face, được thiết kế để thực hiện phân loại hình ảnh không cần huấn luyện trước (zero-shot classification). Điều này có nghĩa là bạn có thể phân loại hình ảnh với các nhãn mà mô hình chưa từng được huấn luyện trên đó.

## Cách thức hoạt động

### Nguyên lý cốt lõi

Pipeline này sử dụng mô hình CLIP (Contrastive Language-Image Pre-Training) để:

1. Mã hóa hình ảnh thành vector đặc trưng
2. Mã hóa các nhãn văn bản (được định dạng thành câu) thành vector đặc trưng
3. So sánh độ tương đồng giữa vector hình ảnh và các vector văn bản
4. Trả về xác suất cho mỗi nhãn dựa trên độ tương đồng

### Ví dụ cụ thể

```
python
```

```
from transformers import pipeline
```

```
classifier = pipeline(model="google/siglip-so400m-patch14-384")
```

```
result = classifier(  
    "https://huggingface.co/datasets/Narsil/image_dummy/raw/main/parrots.png",  
    candidate_labels=["animals", "humans", "landscape"],  
)
```

```
# Kết quả: [{'score': 0.965, 'label': 'animals'}, ...]
```

## Cấu trúc Class

### Kế thừa (Method Resolution Order)

ZeroShotImageClassificationPipeline

- |— transformers.pipelines.base.Pipeline (class cha chính)
- |— transformers.pipelines.base.\_ScikitCompat (tương thích với scikit-learn)
- |— abc.ABC (Abstract Base Class)
- |— transformers.utils.hub.PushToHubMixin (chức năng push lên Hub)
- |— builtins.object (class gốc Python)

## Phân biệt các loại tham số

### 1. Tham số khởi tạo (Constructor Parameters)

Đây là các tham số được truyền vào khi **TẠO** pipeline (lúc khởi tạo object):

python

```
# Các tham số này được sử dụng khi khởi tạo pipeline
classifier = pipeline(
    model="google/siglip-so400m-patch14-384", # Tham số khởi tạo
    device=0,                                # Tham số khởi tạo
    torch_dtype="auto",                      # Tham số khởi tạo
    batch_size=4                             # Tham số khởi tạo
)
```

#### Đặc điểm:

- Được thiết lập **MỘT LẦN** khi tạo pipeline
- Ảnh hưởng đến cách pipeline hoạt động trong suốt vòng đời của nó
- Không thể thay đổi mà không tạo lại pipeline

#### Tham số cốt lõi (Core Parameters)

- **model:** Mô hình AI được sử dụng (CLIP, SigLIP, etc.)
- **image\_processor:** Bộ xử lý hình ảnh để tiền xử lý dữ liệu
- **framework:** Chọn "pt" (PyTorch) hoặc "tf" (TensorFlow)

#### Tham số hiệu năng (Performance Parameters)

- **device:** Thiết bị chạy (-1 cho CPU, số dương cho GPU CUDA)
- **torch\_dtype:** Độ chính xác (torch.float16, torch.bfloat16, "auto")
- **batch\_size:** Kích thước batch cho xử lý hàng loạt
- **num\_workers:** Số worker cho DataLoader

## 2. Tham số đầu vào (Input Parameters)

Đây là các tham số được truyền vào **MỖI LẦN GỌI** pipeline (khi sử dụng):

```
python

# Các tham số này được sử dụng mỗi lần gọi classifier
result = classifier(
    image="path/to/image.jpg",           # Tham số đầu vào
    candidate_labels=["dog", "cat", "bird"], # Tham số đầu vào
    hypothesis_template="This shows a {}", # Tham số đầu vào
    timeout=10                           # Tham số đầu vào
)
```

### Đặc điểm:

- Được truyền vào **MỖI LẦN** sử dụng pipeline
- Có thể thay đổi linh hoạt giữa các lần gọi
- Quyết định đầu vào cụ thể và cách xử lý cho lần thực thi đó

### Chi tiết tham số đầu vào:

- **image**: Hình ảnh cần phân loại (URL, đường dẫn, PIL.Image, hoặc list)
- **candidate\_labels**: Danh sách các nhãn ứng viên (bắt buộc)
- **hypothesis\_template**: Template để định dạng nhãn (mặc định: "This is a photo of {")
- **timeout**: Thời gian chờ tối đa khi tải ảnh từ web

## So sánh trực quan

Khía cạnh	Tham số khởi tạo	Tham số đầu vào
Khi nào sử dụng	Lúc tạo pipeline	Mỗi lần gọi pipeline
Tần suất thay đổi	Một lần duy nhất	Có thể thay đổi mỗi lần
Ví dụ	<code>model</code> , <code>device</code> , <code>batch_size</code>	<code>image</code> , <code>candidate_labels</code>
Mục đích	Cấu hình cách pipeline hoạt động	Cung cấp dữ liệu để xử lý
Ảnh hưởng	Toàn bộ vòng đời pipeline	Chỉ cho lần thực thi hiện tại

## Ví dụ minh họa

python

*# BƯỚC 1: Khởi tạo pipeline với tham số khởi tạo*

```
classifier = pipeline(  
    task="zero-shot-image-classification",  
    model="google/siglip-so400m-patch14-384",    # Tham số khởi tạo  
    device=0,                                    # Tham số khởi tạo  
    torch_dtype="float16"                        # Tham số khởi tạo  
)
```

*# BƯỚC 2: Sử dụng pipeline với tham số đầu vào (có thể gọi nhiều lần)*

```
result1 = classifier(  
    image="dog.jpg",                                # Tham số đầu vào  
    candidate_labels=["dog", "cat", "bird"]          # Tham số đầu vào  
)  
  
result2 = classifier(  
    image="landscape.jpg",                          # Tham số đầu vào khác  
    candidate_labels=["mountain", "beach", "city"], # Tham số đầu vào khác  
    hypothesis_template="This landscape shows a {}" # Tham số đầu vào khác  
)
```

## Phương thức chính

**Kết quả trả về:** Danh sách dictionary, mỗi dictionary chứa:

- **label:** Một trong các candidate\_labels
- **score:** Điểm số từ 0-1 (được tính bằng softmax)

## Quy trình xử lý nội bộ

### 1. preprocess()

- Tải và tiền xử lý hình ảnh
- Định dạng các nhãn thành câu hoàn chỉnh bằng hypothesis\_template
- Tokenize văn bản
- Chuẩn bị tensor đầu vào

### 2. \_forward() (kế thừa)

- Chạy mô hình CLIP
- Tính toán độ tương đồng giữa hình ảnh và văn bản

- Trả về logits

### 3. `postprocess()`

- Áp dụng softmax lên logits để có xác suất
- Sắp xếp kết quả theo thứ tự giảm dần của score
- Format thành dictionary cuối cùng

## Ví dụ sử dụng nâng cao

### 1. Phân loại phong cách ảnh

python

```
classifier(  
    "path/to/image.jpg",  
    candidate_labels=["black and white", "photorealist", "painting"],  
)
```

### 2. Tùy chỉnh template

python

```
classifier(  
    image,  
    candidate_labels=["dog", "cat", "bird"],  
    hypothesis_template="This image shows a {}",  
)
```

### 3. Xử lý nhiều ảnh

python

```
classifier(  
    ["image1.jpg", "image2.jpg"],  
    candidate_labels=["indoor", "outdoor"],  
)
```

## Các phương thức khác từ class cha

### Quản lý mô hình

- `save_pretrained()`: Lưu mô hình và tokenizer

- **push\_to\_hub()**: Upload lên Hugging Face Hub
- **device\_placement()**: Context manager cho việc đặt tensor trên device

## Tương thích Scikit-learn

- **predict()**: Interface tương thích với sklearn
- **transform()**: Interface tương thích với sklearn

## Lưu ý quan trọng

### 1. Zero-shot concept

- Không cần huấn luyện thêm cho nhãn mới
- Chất lượng phụ thuộc vào mô hình CLIP được sử dụng
- Hiệu quả với các khái niệm mà mô hình đã thấy trong quá trình pre-training

### 2. Hiệu năng

- Sử dụng GPU sẽ nhanh hơn đáng kể
- Batch processing cho nhiều ảnh
- Cân nhắc torch\_dtype để cân bằng tốc độ và độ chính xác

### 3. Template quan trọng

- Template mặc định "This is a photo of {}" phù hợp cho hầu hết trường hợp
- Có thể tùy chỉnh template để phù hợp với ngữ cảnh cụ thể
- Template tốt sẽ cải thiện độ chính xác

## Mô hình được hỗ trợ

Pipeline này hoạt động với các mô hình CLIP-based như:

- google/siglip-so400m-patch14-384
- openai/clip-vit-base-patch32
- Các mô hình CLIP khác trên Hugging Face Hub

Đây là một công cụ mạnh mẽ cho việc phân loại hình ảnh linh hoạt mà không cần huấn luyện lại mô hình cho từng tập nhãn mới.