

DINOv2 照片识别测试报告

测试时间: 2025年2月13日

模型: DINOv2 ViT-S/14

设备: NVIDIA GeForce RTX 3090 Ti (CUDA)

测试照片数: 4张

📸 测试照片概览

本次测试使用了4张世界著名地标建筑的照片，检验 DINOv2 对真实世界图像的特征提取能力。

序号	文件名	地标名称	国家/地区	图片尺寸	文件大小
1	OIP-C (1).webp	悉尼歌剧院	澳大利亚 悉尼	474×316	17.7 KB
2	OIP-C.webp	大本钟	英国 伦敦	474×315	13.6 KB
3	Top100AttractionsInTheWorldA9.webp	世界景点合成图	全球	1024×1024	181.4 KB
4	istock000070396403medium.webp	泰姬陵	印度 阿格拉	1240×930	122.0 KB

地标详细介绍

1. 悉尼歌剧院 (Sydney Opera House)

- 建造时间: 1959-1973年
- 建筑师: 约恩·乌松

- **特色:** 贝壳状屋顶设计，20世纪最具特色的建筑之一
- **世界遗产:** 2007年列入联合国教科文组织世界遗产名录

2. 大本钟 (Big Ben)

- **正式名称:** 伊丽莎白塔 (Elizabeth Tower)
- **建造时间:** 1843-1859年
- **高度:** 96米
- **特色:** 哥特复兴式建筑，世界最著名的钟楼之一

3. 世界景点合成图

- **包含地标:** 埃菲尔铁塔、自由女神像、罗马斗兽场、万里长城等
- **用途:** 旅游宣传海报
- **特点:** 多地标组合设计图

4. 泰姬陵 (Taj Mahal)

- **建造时间:** 1632-1653年
- **建造者:** 莫卧儿皇帝沙贾汗
- **特色:** 白色大理石陵墓，穆斯林艺术的瑰宝
- **世界遗产:** 1983年列入世界遗产名录

模型配置

DINOv2 ViT-S/14 模型规格

参数	值
模型类型	Vision Transformer (ViT)
模型大小	Small (S)
Patch 大小	14×14 像素
输入尺寸	224×224 像素
特征维度	384 维
总参数量	22,056,576 (约2200万)

参数	值
Token 数	257 tokens/张 (256 patches + 1 class token)
预训练数据	LVD-142M 数据集
训练方式	自监督学习 (DINO + iBOT)

预处理流程

```

transforms.Compose([
    transforms.Resize(256),      # 短边调整为256
    transforms.CenterCrop(224),   # 中心裁剪224×224
    transforms.ToTensor(),       # 转为张量
    transforms.Normalize(         # ImageNet标准化
        mean=[0.485, 0.456, 0.406],
        std=[0.229, 0.224, 0.225]
    )
])

```



特征提取结果

每张图片的特征统计

照片	特征均值	特征标准差	特征最小值	特征最大值	特征范数
悉尼歌剧院	0.0036	2.4352	-6.5824	7.4447	47.7193
大本钟	0.0203	2.5159	-6.7478	7.6881	49.3033
世界景点合成	0.0209	2.4871	-6.6572	7.4984	48.7388
泰姬陵	-0.0157	2.5282	-6.8312	7.7842	49.5441

分析:

- 特征范数在 47.7-49.5 之间，说明各图片特征激活强度相近

- 泰姬陵的特征范数最高（49.54），可能与其独特的白色大理石建筑特征有关
- 所有特征的均值接近0，符合预训练模型的归一化特性

相似度分析

余弦相似度矩阵

余弦相似度范围：-1（完全不同）到 1（完全相同）

照片	悉尼歌剧院	大本钟	世界景点合成	泰姬陵
悉尼歌剧院	1.000	0.077	0.219	0.117
大本钟	0.077	1.000	0.179	0.079
世界景点合成	0.219	0.179	1.000	0.055
泰姬陵	0.117	0.079	0.055	1.000

关键发现

1. 最高相似度配对

悉尼歌剧院 \leftrightarrow 世界景点合成图 (相似度: 0.219)

- **原因分析:** 世界景点合成图中可能包含悉尼歌剧院元素，或两者共享类似的蓝天+建筑构图
- **特征共享:** 水边建筑、独特屋顶设计、蓝天背景

2. 最低相似度配对

世界景点合成图 \leftrightarrow 泰姬陵 (相似度: 0.055)

- **原因分析:**
 - 世界景点合成图是多元素组合设计图
 - 泰姬陵是单景点摄影，具有独特的对称花园和水池倒影
 - 视觉风格差异极大（现代设计 vs 古典建筑摄影）

3. 整体相似度分析

统计指标	数值
平均相似度	0.1210
最高相似度	0.2194
最低相似度	0.0549
标准差	0.0523

结论:

- 所有照片对的相似度都较低 (< 0.25)，说明 DINOv2 能够区分不同的地标建筑
 - 地标建筑虽然同属"旅游景点"类别，但视觉特征差异明显
 - 这符合预期：不同文化背景、建筑风格、摄影角度产生不同的视觉特征
-

🎯 相似度分组分析

阈值：相似度 > 0.85

结果: 没有超过0.85的相似照片组

这表明：

1. 4张照片都是独特的地标，没有重复或高度相似的图像
2. DINOv2 能够有效区分不同的建筑结构和视觉风格
3. 即使是同一类别（地标建筑），也能捕捉细微的差异

潜在分组（降低阈值至 0.15）

分组	照片	相似度	共同点
组1	悉尼歌剧院、世界景点合成	0.219	都包含悉尼歌剧院元素
组2	大本钟、世界景点合成	0.179	西方地标建筑

⚡ 性能分析

处理速度

指标	数值
总处理时间	0.24 秒
平均每张	59.73 毫秒
吞吐量	16.74 张/秒
GPU 利用率	高

时间分解

模型加载: ~0.5秒 (首次运行, 已缓存)

特征提取: ~0.06秒/张

相似度计算: <0.01秒

总计: ~0.24秒 (4张)

与几何形状测试对比

测试类型	平均处理时间	说明
几何形状 (10张)	23 ms/张	纯色背景, 特征简单
真实照片 (4张)	60 ms/张	复杂场景, 特征丰富

分析: 真实照片处理时间稍长, 可能是因为:

- 图像内容更复杂, 特征计算更密集
- 首次加载模型后缓存效应



DINOv2 特征理解

特征向量特性

1. 高维稀疏性 (384维)
2. 每个维度捕捉不同的视觉概念
3. 稀疏激活：大部分值接近0，少数显著激活
4. 语义编码
5. 不直接对应人类可理解的标签
6. 编码形状、纹理、颜色、结构等视觉信息
7. 对比性
8. 相似视觉内容 → 高余弦相似度
9. 不同视觉内容 → 低余弦相似度

特征可视化概念

虽然无法直接可视化384维向量，但可以理解：

悉尼歌剧院特征 ≈ [贝壳形状, 现代建筑, 蓝色, 水边, 白色, 独特屋顶, ...]

泰姬陵特征 ≈ [圆顶, 古典建筑, 白色大理石, 对称, 花园, 倒影, ...]



应用场景建议

基于本次测试结果，DINOv2 适合以下应用：

1. 图像检索系统

- 能够区分不同地标
- 可用于相似景点推荐

2. 重复图像检测

- 相似度阈值<0.85视为不同图像
- 可用于数据集去重

3. 图像聚类

- 根据视觉相似度自动分组
- 适合相册整理、内容分类

4. 地标识别

- DINOv2 是特征提取器，不是分类器
- 需配合分类头或检索数据库使用



技术细节补充

Token 计算

输入分辨率: 224×224

Patch 大小: 14×14

Patches 数量: $(224/14) \times (224/14) = 16 \times 16 = 256$

Class Token: 1

总 Tokens: $256 + 1 = 257$

4张照片总 Token 数: 1,028 tokens

模型架构

输入图像 (224×224×3)

↓

Patch Embedding (14×14 patches)

↓

Position Embedding + Class Token

↓

Transformer Encoder × 12 layers

↓

Class Token Output (384-dim)

↓

特征向量 [batch, 384]

缓存信息

- **模型文件:** `~/.cache/torch/hub/checkpoints/dinov2_vits14_pretrain.pth`
(85MB)
- **源码目录:** `~/.cache/torch/hub/facebookresearch_dinov2_main/`
- **总缓存大小:** ~90MB

结论

主要发现

1. DINOv2 成功提取了4张地标照片的视觉特征

2. 每张图片生成384维特征向量

3. 特征范数在47-50之间，分布合理

4. 地标建筑间视觉差异明显

5. 平均相似度仅0.12

6. 最高相似度0.22 (悉尼歌剧院 vs 合成图)

7. 能够有效区分不同建筑风格

8. 处理速度快

9. 60ms/张 (GPU)

10. 适合实时应用

11. 自监督特征的有效性

12. 无需标注数据，自动学习视觉表示

13. 捕捉形状、纹理、结构等多维度信息

限制与注意事项

1. **不是分类器**: DINOv2 输出特征，不输出类别标签
2. **相似度相对性**: 0.2的相似度在不同数据集上含义不同
3. **预处理重要**: 必须使用与预训练相同的归一化参数

下一步建议

1. 扩大测试集（更多地标、不同光线条件）
2. 测试其他 DINOv2 变体（ViT-B/14, ViT-L/14）
3. 结合分类器进行地标识别任务
4. 探索特征可视化方法（t-SNE, PCA）



参考信息

- **DINOv2 论文**: "DINOv2: Learning Robust Visual Features without Supervision" (ICCV 2023)
 - **官方仓库**: <https://github.com/facebookresearch/dinov2>
 - **模型卡片**: https://github.com/facebookresearch/dinov2/blob/main/MODEL_CARD.md
-

报告生成时间: 2025年2月13日

测试环境: Ubuntu 22.04, PyTorch 2.7.0, CUDA 12.8