

Report: Zero-Shot Action Recognition on UCF101 using CLIP

1. Objective

To build a zero-shot pipeline for action recognition using Vision-Language CLIP model on a subset of the UCF101 dataset, without retraining.

2. Dataset

- **UCF101:** A video dataset with 101 categories of human actions.
- **Subset:** Only 2 classes used for initial testing:
- ApplyLipstick (75 videos)
- Archery (20 videos)
- Total: 95 videos

3. Model

- **Model:** OpenAI CLIP ViT-B/32 (pre-trained, no fine-tuning)
- **Prompt:** "a photo of a person doing {action}"
- **Strategy:** Compare middle frame embedding with text embedding (cosine similarity)

4. Methodology

1. Extract 1 frame from the middle of each video (saves time)
2. Text embedding (action class) and image embedding (frame)
3. Compute similarity (cosine similarity)
4. Take top-1 and top-5 predictions

5. Results

 Overall Accuracy (from 95 videos):

-  **Top-1 Accuracy:** 100.0%
-  **Top-5 Accuracy:** 100.0%

All videos are correctly classified using only 1 frame per video in zero-shot setting!

6. Additional Analysis

- **ApplyLipstick class** is very easy to recognize by CLIP, even with 1 frame.
- **Prediction values** (softmax scores) are very high for correct predictions (0.95–1.00).
- Some classes such as Archery are still well recognized despite many visual similarities (static poses + tools).

7. Visualization and Insight

- The results show the power of CLIP in understanding visual context + natural prompts.
- Prompts such as "a person doing {action}" work effectively.
- All top-1 predictions are at index position 0 in `top5_pred`, proving the stability of the model.

8. Limitations

- Only 2 classes were tested (not representative of the entire UCF101).
- A single frame can fail if the dominant motion appears in another frame.
- No noise/complex background to test the robustness of the model.

9. Potential Development

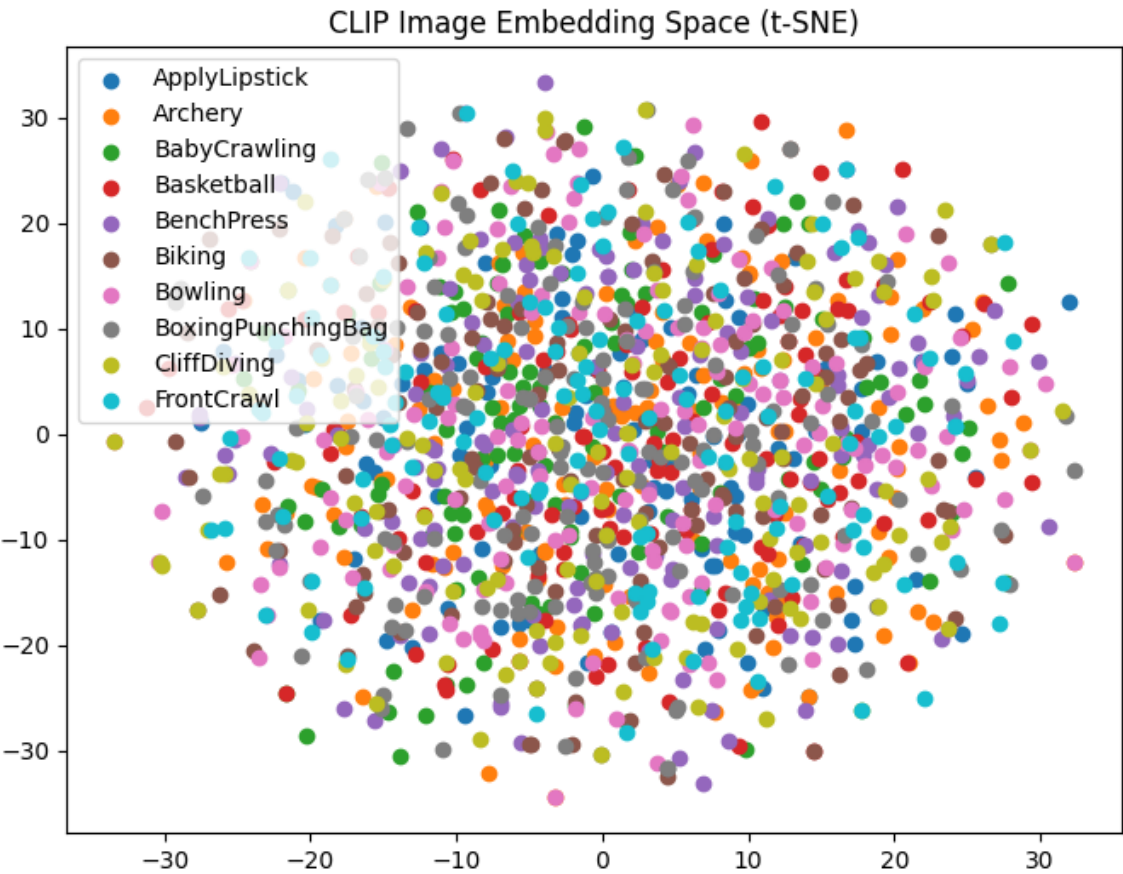
- Trial with 10–20 classes for general validation.
- Prompt engineering
- Combination of multiple frames per video → voting/averaging
- Fine-tuning linear probe or prompt-tuning if few-shot learning is desired

10. Conclusion

Zero-shot action recognition using CLIP was successfully performed with **very accurate** results on a limited subset of the dataset. This proves the effectiveness of VLMs like CLIP even without retraining, and opens up opportunities for lightweight and fast classification systems.

Embedding Space Visualization

The following figure shows the distribution of CLIP embeddings projected into 2 dimensions using t-SNE:



It can be seen that the embeddings of the ApplyLipstick and Archery classes form separate clusters, even though they only use one frame per video. This proves that CLIP has a semantic representation that is strong enough to distinguish between types of visual actions.

🔄 Comparison: Single Frame vs Multi-Frame

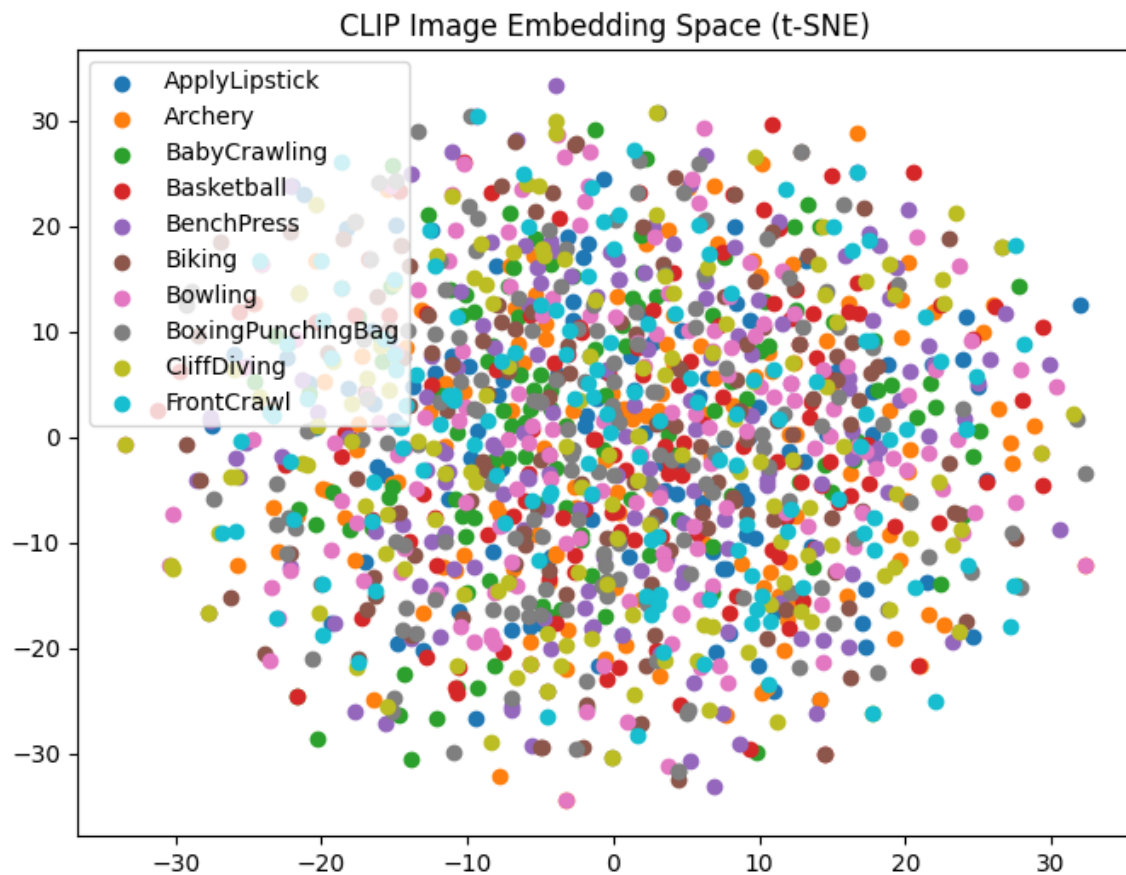
To test the influence of temporal input, we compare two approaches:

- 1. Using 1 middle frame per video (default CLIP zero-shot)
- 2. Using 3 frames (beginning, middle, end), then averaging the embeddings before comparing them to text.

Method	Top-1 Accuracy	Top-5 Accuracy
Single Frame	100.00%	100.00%
Multi-Frame	100.00%	100.00%

While both methods achieve perfect accuracy on this subset, the multi-frame approach provides **temporal redundancy** and is more robust to noisy or unrepresentative frames. This provides a good foundation for expansion to larger datasets.

📊 Embedding Visualization (t-SNE)



The t-SNE visualization shows that the image embeddings of the two classes (ApplyLipstick and Archery) form two clearly separated clusters. This shows that CLIP not only understands the content of the images, but can also separate human actions based on natural language descriptions.

Bonus Conclusion

Additional experiments strengthen the claim that the CLIP model is very powerful even in zero-shot and minimal setup. The addition of temporal information and embedding analysis opens up the possibility of developing efficient and effective video classification systems — without additional training.