

### Abstract

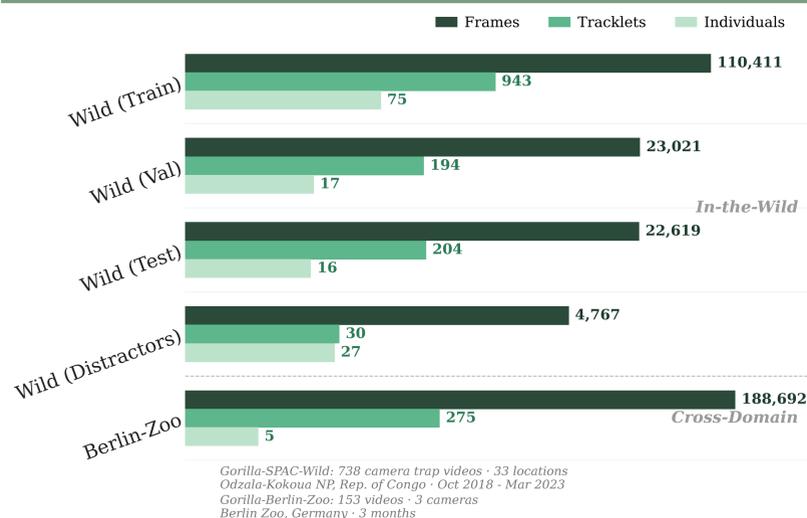
**Monitoring** critically endangered western lowland gorillas **requires immense manual effort** to analyze vast camera trap archives. We present GorillaWatch, an **end-to-end pipeline** for **automated gorilla detection, tracking, and re-identification**, alongside three novel benchmark datasets. Using multi-frame self-supervised pretraining and spatiotemporal constrained clustering, our system enables scalable, non-invasive population monitoring without manual labels.

### Contributions

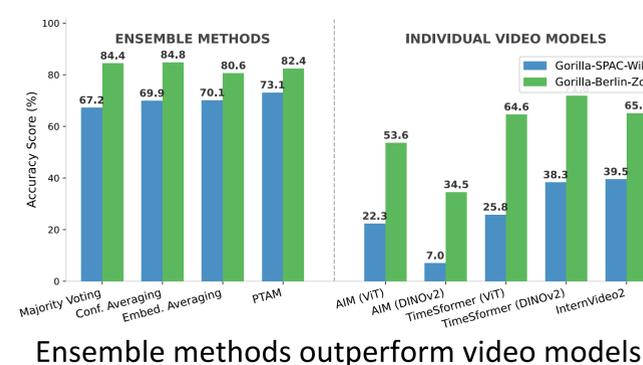
- ① We introduce **three** novel datasets for open-set primate re-identification and tracking.
- ② We propose a multi-frame self-supervised pretraining strategy which leverages temporal consistency in tracklets without manual labels.
- ③ We demonstrate that aggregating features from large-scale image backbones outperforms specialized video architectures.
- ④ We address unsupervised population counting via spatiotemporal constrained clustering, reducing over-segmentation.

# Because every gorilla deserves to be re-KONG-nized.

### Dataset



### Ensemble vs. Video



### AttnLRP - Adaption



### Constrained Clustering

Method	ARI ↑	AMI ↑	# Clusters
<b>Unconstrained</b>			
HAC	0.606	0.728	465
DBSCAN	0.379	0.625	341
HDBSCAN	0.114	0.510	113
<b>Constrained</b>			
HAC	<b>0.837</b>	<b>0.891</b>	<b>17 ✓</b>
DBSCAN	0.586	0.873	13
HDBSCAN	0.184	0.677	204

### Further in the paper

- ① **Tracker Comparison:** Detailed comparison of trackers on our Gorilla-SPAC-Multi-Object-Tracking dataset
- ② **Backbone Zero-Shot Benchmark:** Systematic comparison of pre-trained embedding models.
- ③ **DINOv2 Scaling Analysis:** Accuracy vs. model size (Small → Giant) after fine-tuning
- ④ **Video Architecture Deep-Dive:** Full results, including the effect of replacing standard ViT with DINOv2.

### Pipeline

