

Self-Supervised Cross-View Correspondence with Predictive Cycle Consistency

Anonymous CVPR submission

Paper ID 13577

Abstract

Learning self-supervised visual correspondence is a long-studied task fundamental to visual understanding and human perception. However, existing correspondence methods largely focus on small image transformations, such as object tracking in high-framerate videos or learning pixel-to-pixel mappings between images with high view overlap. This severely limits their application in dynamic multi-view settings such as robot imitation learning. In this work, we introduce Predictive Cycle Consistency for learning object correspondence between extremely disjoint views of a scene without paired segmentation data. Our technique bootstraps object correspondence pseudolabels from raw image segmentations using conditional grayscale colorization and a cycle-consistency refinement prior. We then train deep ViTs on these pseudolabels, which we use to generate higher-quality pseudolabels and iteratively train better correspondence models. We demonstrate the performance of our method under both extreme in-the-wild camera view changes and across large temporal gaps in video. Our approach beats all prior supervised and prior SotA self-supervised correspondence models on the EgoExo4D correspondence benchmark (+6.7 IoU Exo Query) and the prior SotA self-supervised methods SiamMAE and DINO V1&V2 on the DAVIS-2017 and LVOS datasets across large frame gaps.

1. Introduction

At the core of learning is the discovery of recurring patterns. This is the purpose of visual correspondence: given multiple inputs to a scene—whether different camera angles, video frames, or other perspectives—how can we determine which objects are the same and which are different?

In recent years, the introduction of large-scale datasets and powerful model architectures has led to strong results in learning visual correspondence without the need for costly labeled object pairings. However, existing self-supervised object correspondence methods [18, 28, 52, 54] have overwhelmingly focused on domains with continu-

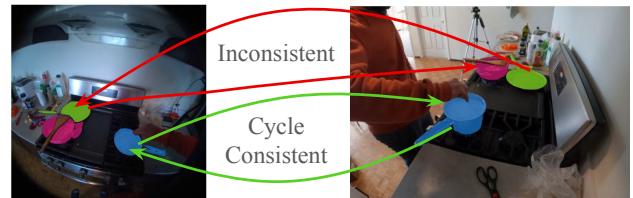


Figure 1. We extract cycle-consistent correspondence at the object level to guide self-supervision across extreme view changes

ous or small transformations, such as Video Object Segmentation (VOS) in continuous videos, or dense pixel-to-pixel methods which necessarily assume that object surfaces overlap between views. This limits the application of self-supervised correspondence to *discontinuous* inputs, such as the first-and-third person EgoExo4D dataset [16], where objects are small, frequently interacted with, and camera pose or depth changes between views are unavailable. Learning correspondence in these challenging, real-world domains is crucial for tasks such as robot imitation learning [42, 44], object reidentification [13, 46], and scene understanding [16]. For example, in robot imitation learning, visual correspondence allows a robot to map observed actions from a third-person perspective to its own first-person actions.

The fundamental challenge of learning visual correspondence across discontinuities is that a system must take into account an object and the surrounding environment holistically to answer the question of where *that* mug is instead of where *a* mug is. A particularly powerful approach for encouraging a holistic understanding in difficult data sources is predictive learning: training a model to de-corrupt a corrupted input. Predictive learning has shown strong results in a wide range of domains, such as natural language processing [9, 39], image classification [2, 21], and generative modeling [22]. Recently, predictive approaches have seen competitive results in object correspondence and representation learning through tasks such as image-conditioned grayscale coloration [52] and conditional masked autoencoders [18, 37].

However, existing predictive methods for correspon-

038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068

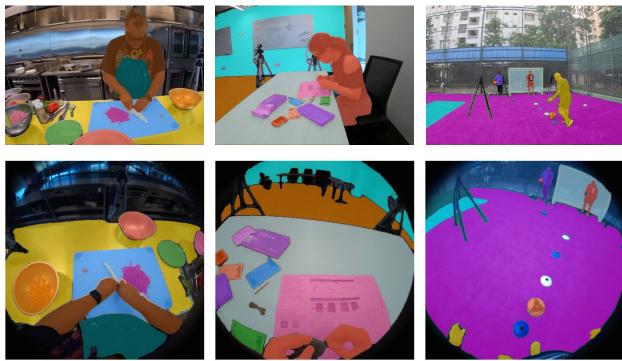


Figure 2. Self-supervised object correspondence on EgoExo4D between Ego views (bottom row) and Exo views (top row) without any paired data. Objects with no match are not colored.

dence struggle to generalize effectively across challenging scenarios (Table 1). State-of-the-art approaches typically rely on K-Nearest-Neighbor matching of emergent learned representations, identifying corresponding regions by comparing the most similar model embeddings across views. Although the pretraining objectives of these methods promote spatial awareness within a scene, the resulting representations often entangle semantic object information with the spatial cues necessary for correspondence. This entanglement reduces robustness against semantically similar distractor objects, limiting generalization and performance in more demanding scenarios.

An effective approach for encouraging a spatially-consistent mapping between views is cycle consistency, which leverages the fact that correspondence is generally invertible across time and view changes [47, 54, 62]. Objects that can be tracked forward in a video can also be tracked backward; camera view changes can be undone. Existing applications of cycle consistency bootstrap long-term correspondence relations via local techniques like optical flow or video palindromes. However, these techniques rely on continuous input data such as high-framerate videos and require hand-crafted biases such as local neighborhoods to restrict a large search space. As a result, cycle consistency approaches have fallen out of favor compared to more modern approaches in challenging scenarios.

In this work, we propose Predictive Cycle Consistency (PCC) for the task of self-supervised object correspondence under extreme viewpoint changes in space (camera angle changes) and time (gaps between video frames). Our approach extends cycle-consistency to operate at the *object level* rather than focusing on features or image patches, allowing for robust correspondence in challenging scenarios. We do this by building on existing predictive approaches [18] and [52] to train a directional correspondence model on the simple, asymmetric task of grayscale col-

orization of a target image conditioned on a colorful source view. Using this, we extract object pairs that correspond to each other when running the correspondence model in *both* directions—source to target and target to source—to generate high-quality paired pseudolabels. We then train deep ViTs on these pseudolabels, which we use to generate higher-quality pseudolabels and iteratively train better correspondence models. All together, our approach combines the power of deep transformers trained with a predictive pretext task for scene understanding with a method to extract refined object boundaries.

Our approach learns object correspondences across a wide range of challenging scenarios, such as matching highly occluded objects (apron and hands, Fig. 2 left) and distinguishing between semantically similar objects (pieces of paper from a COVID test, Fig. 2 middle). We validate our method on several object correspondence benchmarks, achieving superior performance over all previous *labeled* approaches and state-of-the-art unlabeled correspondence methods on the EgoExo4D Object correspondence benchmark [16]. Additionally, our approach sets a new state of the art on widely used video tracking datasets DAVIS-2017 [38] and LVOS [23] under high viewpoint-change conditions. We commit to making our code and checkpoints open source upon acceptance. Our contributions are as follows:

1. We introduce Predictive Cycle Consistency, a technique that combines the powerful representation learning of predictive approaches with the refinement of cycle-consistency for self-supervised correspondence.
2. We propose a pseudo-labeling method that incorporates cycle-consistency on top of existing correspondence models to iteratively refine self-supervised object correspondence outputs.
3. We obtain state-of-the-art results on a suite of correspondence tasks from EgoExo4D, DAVIS, and LVOS.

2. Related Work

Dense Visual Correspondence. The visual correspondence task is fundamental to human visual perception and has a long established history in computer vision. Early techniques for correspondence focused on dense pixel-to-pixel correspondence using classical techniques such as optical flow [3, 24, 27, 30, 35, 48, 57]. In recent years, dense correspondence learning has advanced with deep learning methods that address challenging conditions, including significant viewpoint changes. However, existing dense methods for correspondence either require large amounts of costly labeled data to train, use camera pose [53], depth information [11, 41], or extract annotations through simulated data [50, 51]. Furthermore, implicit to dense correspondence is a restriction on allowed view changes: different views must contain overlapping object surfaces. As a result, our work focuses on the more semantically grounded task of

157 object level correspondence given segmentation masks.

158 **Self-Supervised Object Correspondence** Human vision
159 excels at establishing visual correspondences across
160 space and time without direct supervision, even in the face
161 of occlusions, distractor objects, and object transformations.
162 Inspired by this capability, much work has focused on
163 self-supervised learning of object correspondence through
164 video and scene data.

165 *Contrastive approaches* [20, 61] for object tracking
166 learn correspondence by creating ground-truth annotations
167 through applying strong image transformations. This pow-
168 erful inductive bias has been extraordinarily effective in
169 tracking, but these approaches are criticized for the large
170 amount of engineering and hand-crafting necessary to cre-
171 ate robust augmentations. *Cycle-consistency* [47, 54, 62]
172 emerged as an early technique for learning deep self-
173 supervised object correspondence by leveraging the invert-
174 ability of object correspondence through techniques such as
175 training on video palindromes. However, cycle-consistency
176 as the sole training objective has struggled on its own to
177 maintain stability over long videos or through occlusions,
178 causing it to be replaced by more powerful representa-
179 tion learning techniques [26, 61]. *Tracking-by-Matching*
180 approaches are a recent introduction to learning zero-shot
181 video instance segmentation. These approaches build on
182 top of supervised image segmentation models [56] to cre-
183 ate object-level representations. These representations ei-
184 ther transfer from labeled data across domains [7] or are
185 learned through maintaining invariance under hand-crafted
186 augmentations [33, 34].

187 **Visual Representation Learning for Correspondence**
188 Current state-of-the-art approaches for both supervised and
189 unsupervised correspondence learning rely on strong base
190 models to construct clean visual representations. In the
191 past few years, these base models have moved away from
192 ResNets [19] towards larger models such as Vision Trans-
193 formers (ViTs) [10] trained using self-supervised rep-
194 resentation learning strategies on large datasets. For ex-
195 ample, the DINO model family [4, 8, 36] emphasizes non-
196 contrastive invariance under augmentation. Other methods
197 use exponentially moving average models [1, 17] or siamese
198 networks [5]. In cases involving multiple views, large-
199 scale contrastive models effectively align embedding spaces
200 across domains, achieving strong results in view-invariant
201 learning and language-image matching [14, 25, 40, 58].

202 Masked Autoencoding [21] trains high-quality visual
203 representations by randomly masking patches of an im-
204 age at a high ratio and reconstructing the missing areas.
205 A recent advancement in this approach for correspondence
206 learning is SiamMAE [18], which adapts the infilling task
207 by asymmetrically masking a future frame and predicting
208 it based on an unmasked past frame. This approach, also
209 used in cross-view masked completion models [59, 60], has

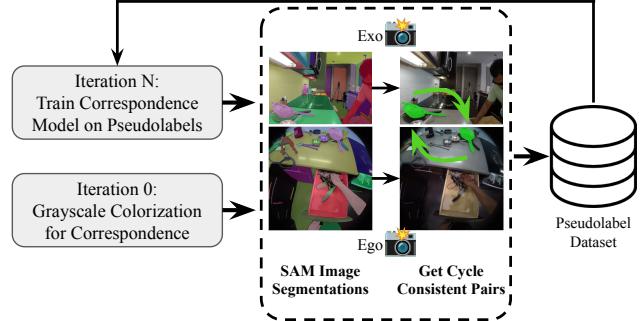


Figure 3. **High-Level Approach.** We 1) Enumerate all objects in each view, 2) Create pseudolabels that are cycle-consistent according to a previous correspondence model iteration, and 3) Use these pseudolabels to train a new correspondence model.

achieved state-of-the-art performance in generating robust representations for object correspondence. In this work, we extend existing masked autoencoder architectures to enhance correspondence across extreme viewpoint changes.

3. Approach

We tackle the task of discovering object correspondence between two images $\mathcal{I}_1 \in \mathbb{R}^{3 \times H_1 \times W_1}$ and $\mathcal{I}_2 \in \mathbb{R}^{3 \times H_2 \times W_2}$ of the same spatiotemporal scene. Specifically, we want to create a model takes in a binary segmentation mask \mathcal{O}_1 of an object in image \mathcal{I}_1 and outputs \mathcal{O}_2 , the corresponding object segmentation mask in \mathcal{I}_2 , or indicates that no correspondence exists.

At a high level, our correspondence pipeline (Fig. 3) operates by 1) Enumerating all objects in each image with existing segmentation models, 2) Using our proposed PCC technique to discover a subset of corresponding objects in each view, 3) Generating correspondence pseduolabels with these mined correspondences, and 4) Iteratively training new models and generating new pseudo labels to learn high-quality self-supervised visual correspondence.

3.1. Image Segmentation

We break down finding corresponding objects between two images into two subtasks: 1) Enumerate all objects in each image (Object detection and Image Segmentation) and 2) Discover which of these objects correspond to one another (Correspondence). This cascaded approach to correspondence has recently obtained state-of-the-art results in video object segmentation with low data or textual object descriptions [7, 34]. There is a great amount of prior work in object detection and image segmentation, with popular approaches being SAM [32] in the supervised setting and UnSAM [56] and CutLER [55] in the self-supervised setting. We follow prior object-centric correspondence work [7, 33, 34] and use SAM to segment an image into distinct object regions,

244 although we suspect self-supervised approaches could be
 245 applied for a completely unsupervised pipeline. The rest of
 246 this paper focuses on the task of (2), learning object corre-
 247 spondence without labels given object segmentations.

248 3.2. Grayscale Colorization

249 Grayscale colorization [52] was an early yet remarkably ef-
 250 fective method for self-supervised object tracking in con-
 251 tinuous videos. However, as computational resources grew
 252 and advanced architectures like Vision Transformers (ViTs)
 253 became prevalent, grayscale colorization fell out of favor,
 254 replaced by more complex and general pretext tasks such
 255 as masked autoencoding [12, 18, 21, 29, 37, 49]. To ad-
 256 dress the challenge of learning correspondence under ex-
 257 treme discontinuities, we reinvestigate the grayscale col-
 258 orization approach and find that it reemerges as strikingly
 259 successful technique.

260 Our architecture for grayscale colorization is depicted
 261 in Figure 4 (a). Conditioned on our source image \mathcal{I}_1 we
 262 optimize to reconstruct \mathcal{I}_2 given its grayscale counterpart
 263 \mathcal{I}_2^g . We adapt our architecture from prior correspondence
 264 work [18, 60] and use a two-stage ViT [10] that 1) passes
 265 \mathcal{I}_1 and \mathcal{I}_2 separately through an encoder module \mathbf{E} to ob-
 266 tain $\mathbf{E}(\mathcal{I}_1)$, $\mathbf{E}(\mathcal{I}_2^g)$, and 2) passes the encoded grayscale im-
 267 age $\mathbf{E}(\mathcal{I}_2^g)$ into a decoder \mathbf{D} that merges the colorful image
 268 conditioning using cross attention: $\mathbf{D}(\mathbf{E}(\mathcal{I}_1), \mathbf{E}(\mathcal{I}_2^g))$. For
 269 optimization, we use a simple mean-squared-error loss in
 270 the RGB color space.

271 3.3. Visual Correspondence from Colorization

272 We now propose an approach to extract object correspon-
 273 dence from the inputs and outputs of conditional coloriza-
 274 tion models based on a simple observation: if we augment
 275 an object’s color in the source view, we would expect the
 276 output of the colorization model to change the object’s color
 277 in the target view as well.

278 Our approach is demonstrated in Figure 5. Given a bi-
 279 nary object segmentation map \mathcal{O}_1 in \mathcal{I}_1 we seek to obtain
 280 the corresponding segmentation mask \mathcal{O}_2 in view \mathcal{I}_2 . We
 281 perform two forward passes of our grayscale colorization
 282 model, \mathbf{F} . First, we perform grayscale cross-view infilling
 283 of \mathcal{I}_2^g given a colorized \mathcal{I}_1 from before: $y = \mathbf{F}(\mathcal{I}_2^g, \mathcal{I}_1)$. For
 284 our second model pass, we augment \mathcal{I}_1 at the location of the
 285 segmentation mask \mathcal{O}_1 by adding a constant color offset to
 286 each channel. We denote this augmented image as \mathcal{I}'_1 , and
 287 the RGB augmentation as a vector $c \in \mathbb{R}^3$. We then colorize
 288 \mathcal{I}_2^g conditioned on this augmented image $y' = \mathbf{F}(\mathcal{I}_2^g, \mathcal{I}'_1)$.
 289 To identify where objects correspond, we then take the av-
 290 erage absolute difference between y and y' over each color
 291 channel and normalize to sum to one, outputting an object
 292 segmentation heatmap $\mathcal{H} \in \mathbb{R}_{+}^{H_2 \times W_2}$:

$$\mathcal{H}_{ij} = \frac{\sum_{c=1}^3 |\mathbf{F}(\mathcal{I}_2^g, \mathcal{I}_1)_{cij} - \mathbf{F}(\mathcal{I}_2^g, \mathcal{I}'_1)_{cij}|}{\sum_{c,k,l} |\mathbf{F}(\mathcal{I}_2^g, \mathcal{I}_1)_{ckl} - \mathbf{F}(\mathcal{I}_2^g, \mathcal{I}'_1)_{ckl}|} \quad (1) \quad 293$$

294 Importantly, this approach is entirely blackbox with re-
 295 spect to the grayscale colorization model. As a result,
 296 the spatial awareness of objects captured during grayscale
 297 model pretraining naturally emerges during extraction. This
 298 is not necessarily true for masked autoencoder approaches
 299 like SiamMAE, where internal model representations are
 300 used to extract correspondences that may not sufficiently
 301 encode cross-view information to disentangle semantically
 302 similar distractors.

3.4. Generating Cycle-Consistent Correspondence

303 In this section, our goal is to create a set of pseudolabels
 304 containing corresponding objects between two im-
 305 ages \mathcal{I}_1 and \mathcal{I}_2 given their respective object segmentations
 306 $\{\mathcal{O}_{1,i}\}_{i=1}^{N_1}$ and $\{\mathcal{O}_{2,j}\}_{j=1}^{N_2}$ output by SAM, where $\mathcal{I}_1, \mathcal{I}_2$
 307 have N_1, N_2 segmented objects respectively. Following
 308 Sec. 3.3, we generate an augmented $\mathcal{I}'_{1,i}$ for each segmen-
 309 tation mask $\mathcal{O}_{1,i}$ and calculate the output correspondence
 310 heatmap in \mathcal{I}_2 , which we write as $\mathcal{H}_i^{1 \rightarrow 2}$. Here, “1 → 2”
 311 represents predicting \mathcal{I}_2 conditioned on \mathcal{I}_1 . To reduce vari-
 312 ance, we use the same augmentation vector c for each aug-
 313 mentation. By caching the unaugmented pass, calculating
 314 all $\mathcal{H}_i^{1 \rightarrow 2}$ takes $N_1 + 1$ forward passes.

315 We then define the similarity between $\mathcal{H}_i^{1 \rightarrow 2}$ and $\mathcal{O}_{2,j}$ to
 316 be the amount of weight $\mathcal{H}_i^{1 \rightarrow 2}$ places on the segmentation
 317 mask of $\mathcal{O}_{2,j}$, equal to the Frobenius inner product:

$$\text{Sim}(\mathcal{H}_i^{1 \rightarrow 2}, \mathcal{O}_{2,j}) = \langle \mathcal{H}_i^{1 \rightarrow 2}, \mathcal{O}_{2,j} \rangle_F \quad (2) \quad 319$$

320 Then, we pair every object $\mathcal{O}_{2,j}$ with the most similar
 321 heatmap $\mathcal{H}_i^{1 \rightarrow 2}$, which we write as $\mathbf{P}_j^{1 \rightarrow 2}$:

$$\mathbf{P}_j^{1 \rightarrow 2} = \arg \max_i \text{Sim}(\mathcal{H}_i^{1 \rightarrow 2}, \mathcal{O}_{2,j}) \quad (3) \quad 322$$

323 There is an important difference between taking the
 324 argmax over the source object augmentation heatmaps and
 325 taking the argmax over the objects in the target image. Out-
 326 puts of grayscale colorization can be correlated with many
 327 factors such as object semantics (the sky is probably blue
 328 but a balloon can be any color) or lighting conditions, which
 329 may mean that some objects in the target view are more sus-
 330 ceptible to having their predicted colors change with respect
 331 to augmented inputs than others. As a result, the object with
 332 the greatest color change in \mathcal{I}_2 is often the same for many
 333 $\mathcal{H}_i^{1 \rightarrow 2}$. By taking the most similar $\mathcal{H}_i^{1 \rightarrow 2}$ for each $\mathcal{O}_{2,j}$, we
 334 normalize over the easiness of each object to be changed
 335 as the result of an augmentation in \mathcal{I}_1 during grayscale col-
 336 orization.

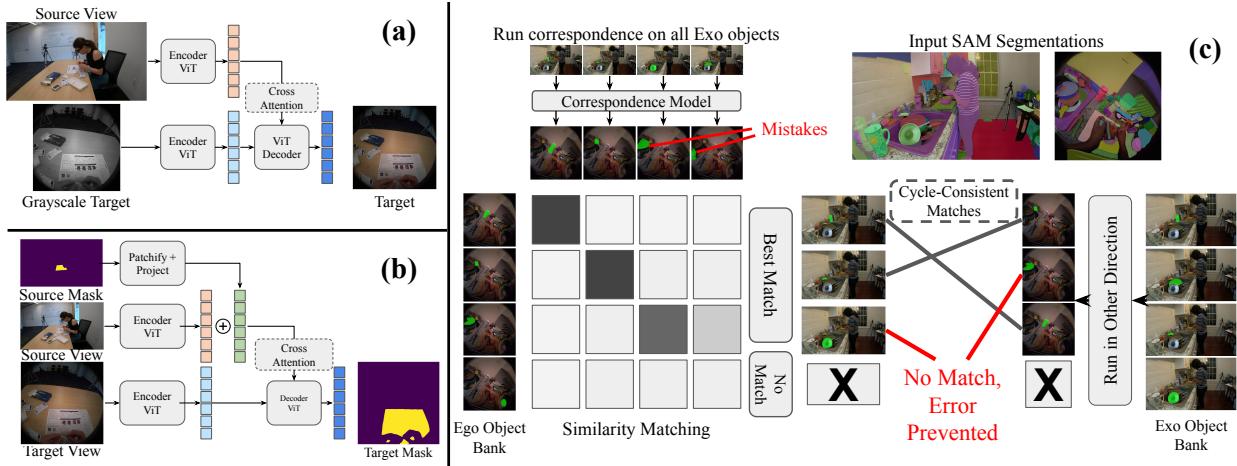


Figure 4. Model Architecture and Cycle Consistency. (a) Our architecture for grayscale color completion, using a shared Encoder for each view and a Decoder with Cross and Self Attention. (b) Our architecture for correspondence, which inserts a projection to represent the mask. (c) Our cycle-consistency pipeline. We run correspondence using the model from (a) or (b) on all objects in each view, find the best matching objects, and extract cycle-consistent matches.

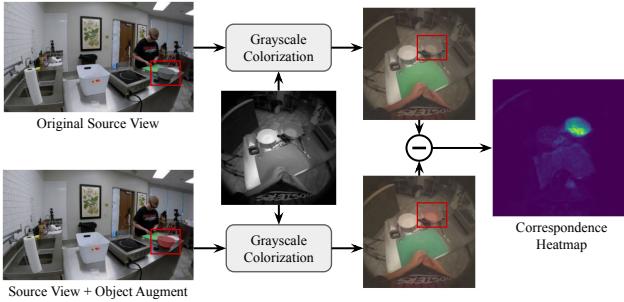


Figure 5. A simple method to extract object correspondence from a deep conditional grayscale colorization model.

337 We repeat the above process in the reverse direction
 338 to extract $\mathcal{H}_j^{2 \rightarrow 1}$ and $\mathcal{P}_i^{2 \rightarrow 1}$. Finally, we extract cycle-
 339 consistent pseudolabels at the object level by taking all pairs
 340 (i, j) where $P_j^{1 \rightarrow 2} = i$ and $P_i^{2 \rightarrow 1} = j$. This cycle consis-
 341 tency enforces that augmenting O_i in I_1 leads to changes at
 342 O_j in I_2 and augmenting O_j in I_2 leads to changes at O_i
 343 in I_1 during grayscale infilling, which strongly implies cor-
 344 respondence. We then use these extracted pseudolabels on
 345 each image to generate a dataset of labeled correspondence
 346 data with paired images and segmentation masks.

3.5. Correspondence Model

348 Our end goal is to have a correspondence model G than
 349 can take in a source image I_1 , source object segmentation
 350 mask \mathcal{O}_1 , and target image I_2 and output the correspond-
 351 ing object segmentation mask for the target image. To train
 352 such a model on our extracted pseudolabels or ground truth
 353 labeled data, we slightly modify our grayscale colorization

architecture, shown in Fig 4 (b). We patchify the segmentation mask \mathcal{O}_1 , apply a linear projection, and add the output to $E(I_1)$. The decoder then takes in the encoded target view $E(I_2)$ and uses cross attention to query the source view and source segmentation mask. As done in [16] we use a Dice Loss [45] and a BCE Loss to predict the target mask.

3.6. Iterative Pseudolabeling and Training

We repeat our pseudolabeling pipeline by substituting grayscale colorization with a trained correspondence model. Assuming that a correspondence model G outputs a binary mask, we replace the similarity score (Equation 2) with intersection over union (IoU) and make no other changes. We can then use these new pseudolabels to train new correspondence models, which in turn can generate new pseudolabels. In our experiments, we find that running three rounds of pseudolabeling (one grayscale, two iterative) saturates quality.

4. Experiments

For comprehensive evaluation, we train and evaluate our PCC pipeline from start to finish in two settings 1) high camera pose changes and 2) large temporal gaps in videos. We describe our training and evaluation approaches for each setting in Sec. 4.1 and present results in Sec. 4.2 and 4.3.

4.1. Experimental Setup

Datasets To train and evaluate the quality of PCC across extreme camera viewpoint changes, we use the EgoExo4D correspondence dataset [16]. EgoExo4D [16] builds on Ego4D [15] and captures over 1000 hours of complex human actions such as cooking and basketball from several

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383 time-synchronized camera views. The camera views con-
384 sist of one egocentric (ego) view captured by person per-
385 forming the action and one or more exocentric (exo) camera
386 views that capture the holistic action environment from the
387 third person. We train and evaluate on the object corre-
388 spondence benchmark in EgoExo4D, which is unique across cor-
389 respondence datasets in that it focuses on extremely chal-
390 lenging settings with many small objects, high occlusion,
391 and extreme camera angle differences in a dynamic environ-
392 ment. Correspondence is evaluated in two directions, either
393 with taking a query mask in the ego view and segmenting
394 the corresponding object in a paired exo view (ego query),
395 or the inverse task from the exo to the ego view (exo query).

396 To run PCC on video object segmentation across large
397 temporal gaps, we train on the commonly used Kinetics-400
398 dataset [31], consisting of large-scale action videos taken
399 from the web. We evaluate temporal correspondence on
400 DAVIS-2017 [38], which captures densely segmented short
401 videos and LVOS [23], which measures video object seg-
402 mentation across long dynamic videos with an average du-
403 ration of over a minute.

404 **Backbone** We implement our correspondence models
405 using a ViT-B with a patch size of 16 to match the com-
406 pute and parameter count of current state-of-the-art self-
407 supervised correspondence work [18, 37]. We initialize our
408 model from CroCo v2 [59, 60], which pretrains with the
409 task of image conditioned cross-view autoencoding.

410 **Grayscale Model** We train our grayscale colorization
411 models separately for EgoExo4D and Kinetics-400. For
412 EgoExo4D, for each ego view in the dataset, we 1) select
413 a random corresponding exo camera, 2) choose a random
414 synchronized frame from each view, and 3) choose with
415 50% probability whether to grayscale the ego or exo view.
416 For Kinetics-400, we follow SiamMAE and select two ran-
417 dom frames between 4 and 48 frames apart and colorize
418 the grayscale future frame conditioned on the past frame.
419 We train each grayscale colorization model for 60k updates
420 with a batch size of 256, which corresponds to 200 epochs
421 for EgoExo4D and 50 epochs for Kinetics-400. Additional
422 hyperparameters are in Appendix 6.1.

423 **Predictive Cycle Consistency** We then run Predictive
424 Cycle Consistency to extract paired pseudolabels. We use
425 the ViT-H version of SAM to run image segmentation on
426 video frames to extract $\{\mathcal{O}_{1,i}\}_{i=1}^{N_1}$ and $\{\mathcal{O}_{2,j}\}_{j=1}^{N_2}$. For
427 EgoExo4D, we segment every labeled (ego,exo) frame pair
428 in the correspondence benchmark, which generally amounts
429 to one frame per second. For Kinetics-400, we segment two
430 random frames per video selected 2 seconds apart, and we
431 ablate the length between frames in our experiments.

432 **Correspondence Model** For a fair comparison, we
433 match the compute of the EgoExo4D supervised baselines
434 and train our correspondence models on 5M supervised ex-
435 amples, divided into 10k updates with a batch size of 512.

We train for the same duration on Kinetics-400, and additional hyperparameters are provided in Appendix 6.1. For a given input example, we select a random cycle-consistent pseudolabel as a training target, and with a 50% chance select whether to train with an ego or exo query. On 25% of inputs, we replace the correspondence objective with learning correspondence existence accuracy, where we select a uniformly random SAM image segment, label it positive if it is also a cycle-consistent pseudolabel and negative otherwise.

Baselines To the best of our knowledge, this is the first work to tackle self-supervised correspondence across the uniquely extreme viewpoint changes on EgoExo4D. As a result, we are rigorous about reimplementing the prior SoTA open source correspondence models. As demonstrated in the SiamMAE paper [18], the current best approaches for correspondence without labeled data are deep ViT models. In particular, SiamMAE makes two key observations 1) a small patch size (ViT-X/8) results in significantly stronger correspondence results (+9.5 on DAVIS-17 for ViT-S/8 vs SiamMAE ViT-S/16) and 2) despite deep work in prior correspondence approaches, DINO v1 [4], which focuses on learning robust and general visual representations at scale, outperforms all approaches except SiamMAE. As a result, for prior work we compare against 1) SiamMAE ViT-S/8, 2) SoTA DINO v1 and v2 models, and 3) CroCo v2 (to compare versus our baseline model architecture). For fairness, we continually pretrain SiamMAE and CroCO v2 on EgoExo4D. Additional implementation details are in Appendix 6. Furthermore, although our model does not use SAM at inference, we do use image segmentation for pseudolabel generation. To account for this, we implement a setting where we use SAM ViT-H to select the object with the highest IoU. For the DAVIS-17 and LVOS baselines, we similarly compare against the aforementioned models, however find that SAM did not help because it cut off objects (e.g. selecting a tire instead of a bike), which we discuss more in Appendix 6.

4.2. Results Across Space

We report quantitative results for PCC across extreme viewpoint changes on the EgoExo4D correspondence benchmark [16] in Table 1. Following the EgoExo4D baseline, we measure 1) IoU (denoted as \mathcal{J}_m in VOS datasets), 2) Location Score, representing the mean-squared-error distance between the predicted and ground truth centroid, 3) Contour Accuracy (CA, denoted as \mathcal{F}_m in VOS datasets), and 4) object presence balanced accuracy (Bal. Acc.). We use the official code for EgoExo4D for evaluation, and run all baselines at 480p, as done in both [16] and [18].

Our results demonstrate that PCC outperforms all prior labeled and all prior SoTA self-supervised approaches across all metrics. We observe that the DINO family

Table 1. Results on of the EgoExo4D [16] correspondence benchmark (v1, test set). \diamond : Continual Pretraining on EgoExo4D \clubsuit : Models with access to multiple frames per view. Bold is best, underlined is second best. Top labeled, middle and bottom self-supervised.

Method	Backbone	Ego Query			CA. \uparrow	Exo Query			
		Bal. Acc. \uparrow	IoU \uparrow	Loc. Score \downarrow		Bal. Acc. \uparrow	IoU \uparrow	Loc. Score \downarrow	
XSegTx [16]	SegSwap [43] + ViT-B	62.63	13.88	0.154	0.239	74.6	21.8	0.133	0.265
XMem \clubsuit [6]	ResNet-50 [19] + Memory	42.33	13.07	0.312	0.182	56.96	10.2	0.249	0.125
XView-XMem \clubsuit [16]	XMem + ViT-B	53.28	22.14	0.176	0.325	59.36	23.56	0.186	0.308
XView-XMem (+ XSegTx) \clubsuit [16]	XMem + SegSwa + ViT-B	54.61	22.5	0.139	0.347	52.28	19.39	0.208	0.255
Ours Supervised	ViT-B/16	<u>74.7</u>	<u>38.41</u>	<u>0.037</u>	0.603	88.45	<u>43.70</u>	0.049	<u>0.555</u>
Ours Supervised + PCC	ViT-B/16	76.9	39.01	0.033	<u>0.600</u>	87.23	47.06	<u>0.054</u>	0.590
SiamMAE \diamond [18]	ViT-S/8	\emptyset	12.24	0.170	0.185	\emptyset	14.18	0.159	0.198
CrocoV2 \diamond [59]	ViT-B/16	\emptyset	7.14	0.200	0.138	\emptyset	9.56	0.164	0.136
DINO [4]	ViT-B/8	\emptyset	12.55	0.153	0.178	\emptyset	15.94	0.137	0.246
DINOv2+Registers [36] [8]	ViT-B/14	\emptyset	20.26	0.125	0.299	\emptyset	24.6	0.169	0.307
SiamMAE \diamond +SAM	ViT-S/8+ViT-H	\emptyset	17.97	0.180	0.254	\emptyset	24.05	0.143	0.315
DINOv2+Registers+SAM	ViT-B/14+ViT-H	\emptyset	28.92	0.153	0.365	\emptyset	34.78	0.123	0.433
Grayscale Coloration + SAM \diamond	ViT-B/16+ViT-H	\emptyset	20.82	0.110	0.311	\emptyset	19.50	0.109	0.276
PCC Iter 1	ViT-B/16	\emptyset	26.41	<u>0.085</u>	0.396	\emptyset	34.35	0.090	0.436
PCC Iter 2	ViT-B/16	65.22	29.98	0.083	0.446	66.40	<u>40.41</u>	<u>0.079</u>	<u>0.502</u>
PCC Iter 3	ViT-B/16	60.66	<u>29.89</u>	0.094	<u>0.432</u>	67.90	41.45	0.071	0.508

of models, which learns robustness against strong image augmentations, outperforms the predictive SiamMAE and CroCo V2 approaches. Although applying SAM at inference time strongly improves self-supervised approaches in IoU and Contour Accuracy, it harms spatial accuracy as measured by Location Score, where PCC dramatically outperforms all prior work. We find that our self-supervised method performs relatively better in the Exo Query setting, where input objects are small, but output objects are generally larger. We find that iteratively training PCC strongly improves results, although performance largely saturates after the second iteration.

We further apply our correspondence model architecture (Sec. 3.5) on the labeled training data split of EgoExo4D, achieving state-of-the-art results compared to prior supervised approaches. When combining the supervised data from EgoExo4D with PCC Iteration 3 pseudolabels, we observe improvements over the supervised baseline, particularly in the Exo Query setting, suggesting that PCC captures complementary information to the EgoExo4D labels.

We show qualitative results of PCC on EgoExo4D in 6. As shown in rows 3 and 4, our model can locate corresponding objects even under significant occlusions. However, it sometimes struggles to produce fine-grained masks. Examples of pseudolabels on EgoExo4D are visible in 2.

4.3. Results Across Time

We present quantitative results for PCC across large temporal gaps between video frames on the DAVIS-17 dataset [38] in Table 2. We follow the evaluation methodology in [11] and evaluate correspondences on all video frame pairs in DAVIS with a temporal gap of 20 frames. Evaluation is restricted to only cover objects present in both views, and report \mathcal{J} & \mathcal{F}_m , \mathcal{J}_m , and \mathcal{F}_m . We report quantitative results

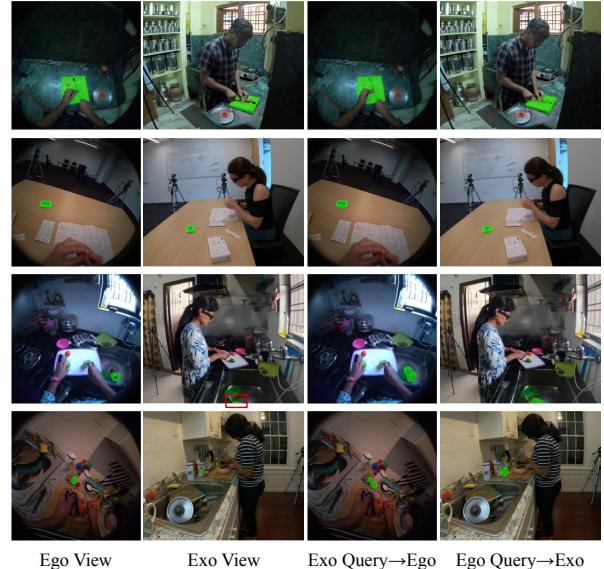


Figure 6. Qualitative results of PCC on the EgoExo4D correspondence benchmark.

on LVOS in Figure 7. Enabled by the large lengths of LVOS videos, we report results across a range of temporal gaps between 5 frames (the LVOS annotation rate) and 400 frames. We again restrict evaluation to objects present in both views and report \mathcal{J} & \mathcal{F}_m .

On DAVIS, we find that PCC significantly outperforms the prior SOTA approaches for VOS. Similar to the EgoExo setting, PCC shows strong improvements after the first iteration, but performance saturates after the second. On LVOS, although SiamMAE achieves SoTA for temporal gaps of 5 and 10 frames, PCC shows comparatively better

Table 2. SoTA Comparison of Video Object Correspondence Methods on DAVIS-17 [38] Val with a temporal gap of 20 frames

Method	$\mathcal{J} \& \mathcal{F}_m \uparrow$	$\mathcal{J}_m \uparrow$	\mathcal{F}_m
SiamMAE [18]	60.7	58.4	62.9
CrocoV2 + Cont. Pretrain [59]	40.0	37.4	42.5
DINO ViTs/8 [4]	64.5	61.6	67.5
DINO ViTb/8 [4]	66.4	63.7	69.2
DINOv2 + Reg ViTb/14 [8]	62.1	59.6	64.8
PCC Iter 1	64.4	61.3	67.5
PCC Iter 2	69.7	67.0	72.4
PCC Iter 3	70.2	67.8	72.7

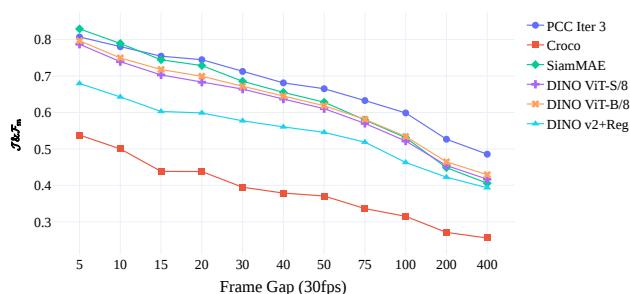


Figure 7. Results on LVOS-V1 Val at different temporal gaps.

532 performance as the frame gap increases, achieving SoTA
533 results for all other temporal gaps. Interestingly, SiamMAE
534 sees a sharper decline as the frame gap widens compared
535 to DINO. We suspect this suggests that at shorter frame
536 gaps, precise object segmentation is most critical, whereas
537 at larger frame gaps, robustness to object semantics be-
538 comes more important.

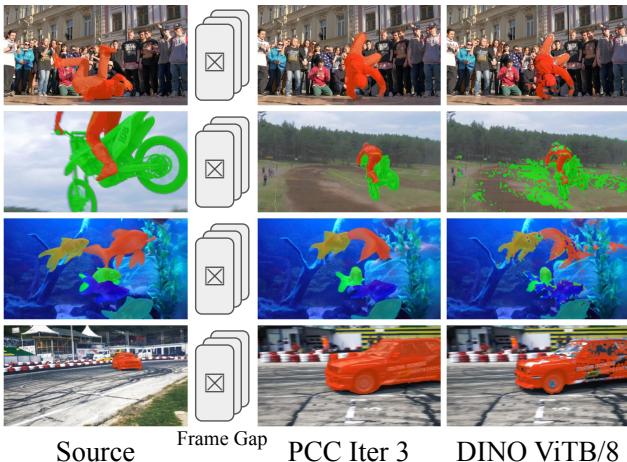


Figure 8. Qualitative results on correspondence across time. Representational approaches degrade quickly with larger frame gaps.

539 We demonstrate qualitative results for object recognition
540 across temporal gaps in 8. We observe that representa-

tional approaches such as DINO struggle to maintain object consistency (bottom row), handle semantically similar distractors with moderate movement (third row), and manage heavy object occlusions (second row). In contrast, PCC maintains high-quality correspondence throughout.

4.4. Ablations

Table 3. **Ablation: Model configuration** Supervised Training, EgoExo4D Validation Set. Grayscale continued training harms supervised training. A deeper decoder improves performance.

Pretraining	Dec. Params.	Ego/Exo Query		
		IoU \uparrow	Loc. Score \downarrow	CA. \uparrow
MAE	50M	36.9/40.2	0.043/0.078	0.57/0.51
MAE	100M	37.7/43.7	0.042/0.060	0.60/0.56
Grayscale	100M	31.0/34.3	0.048/0.089	0.49/0.43

Table 4. **Ablation: Frame Gap** Davis 2017 Validation Set. We ablate the framerate gap on Kinetics-400 during pseudolabel extraction on Iteration 3.

Temporal Gap	$\mathcal{J} \& \mathcal{F}_m$		
	2 Frames	10 Frames	30 Frames
2 seconds	83.9	75.8	65.4
4 seconds	83.5	76.2	66.3
6 seconds	83.7	76.2	67.7

We ablate our model configuration on supervised EgoExo4D correspondence in Table 3. Our results show that pretraining with a grayscale objective on EgoExo4D, rather than a MAE objective, significantly degrades performance. This suggests that the gains from PCC are not due to better representation learning, but rather arise from the full cycle-consistent pipeline.

We investigate the impact of the temporal gap between images for pseudolabeling PCC on Kinetics-400 in Table 4. We find that training with a longer gap improves longer-term correspondence.

5. Conclusion

In this work, we address the task of learning self-supervised visual correspondence across significant camera view shifts and large temporal gaps in video. Our approach, *Predictive Cycle Consistency*, achieves state-of-the-art performance on the EgoExo4D correspondence benchmark and in low-frame-rate scenarios on the DAVIS-17 and LVOS datasets. Furthermore, we show that applying cycle consistency to correspondence predictions *at the object level* is an strikingly effective strategy for generating clean and consistent correspondence pseudolabels. Overall, our work extends the capabilities of self-supervised learning to establish correspondence in highly disjoint and challenging domains, paving the way for improved understanding of dynamic scenes in the natural world.

573

References

- [1] Alexei Baevski, Wei-Ning Hsu, Qiantong Xu, Arun Babu, Jiatao Gu, and Michael Auli. Data2vec: A general framework for self-supervised learning in speech, vision and language. In *International Conference on Machine Learning*, pages 1298–1312. PMLR, 2022. 3
- [2] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. 1
- [3] Steven S. Beauchemin and John L. Barron. The computation of optical flow. *ACM computing surveys (CSUR)*, 27(3):433–466, 1995. 2
- [4] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021. 3, 6, 7, 8
- [5] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15750–15758, 2021. 3
- [6] Ho Kei Cheng and Alexander G Schwing. Xmem: Long-term video object segmentation with an atkinson-shiffrin memory model. In *European Conference on Computer Vision*, pages 640–658. Springer, 2022. 7
- [7] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Alexander Schwing, and Joon-Young Lee. Tracking anything with de-coupled video segmentation. In *ICCV*, 2023. 3, 2
- [8] Timothée Darcret, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers, 2023. 3, 7, 8
- [9] Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 1
- [10] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 3, 4
- [11] Mohamed El Banani, Ignacio Rocco, David Novotny, Andrea Vedaldi, Natalia Neverova, Justin Johnson, and Ben Graham. Self-supervised correspondence estimation via multiview registration. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1216–1225, 2023. 2, 7
- [12] Alexandre Eymaël, Renaud Vandeghen, Anthony Cioppa, Silvio Giancola, Bernard Ghanem, and Marc Van Droogenbroeck. Efficient image pre-training with siamese cropped masked autoencoders. *arXiv:2403.17823*, 2024. 4
- [13] Heng Fan, Liting Lin, Fan Yang, Peng Chu, Ge Deng, Sijia Yu, Hexin Bai, Yong Xu, Chunyuan Liao, and Haibin Ling. Lasot: A high-quality benchmark for large-scale single object tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5374–5383, 2019. 1
- [14] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15180–15190, 2023. 3
- [15] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022. 5
- [16] Kristen Grauman, Andrew Westbury, Lorenzo Torresani, Kris Kitani, Jitendra Malik, Triantafyllos Afouras, Kumar Ashutosh, Vijay Baiyya, Siddhant Bansal, Bikram Boote, et al. Ego-exo4d: Understanding skilled human activity from first-and third-person perspectives. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19383–19400, 2024. 1, 2, 5, 6, 7
- [17] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020. 3
- [18] Agrim Gupta, Jiajun Wu, Jia Deng, and Li Fei-Fei. Siamese masked autoencoders. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, Red Hook, NY, USA, 2024. Curran Associates Inc. 1, 2, 3, 4, 6, 7, 8
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 3, 7
- [20] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 3
- [21] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022. 1, 3, 4
- [22] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1
- [23] Lingyi Hong, Wencho Chen, Zhongying Liu, Wei Zhang, Pinxue Guo, Zhaoyu Chen, and Wenqiang Zhang. Lvos: A benchmark for long-term video object segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 13480–13492, 2023. 2, 6
- [24] Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial intelligence*, 17(1-3):185–203, 1981. 2
- [25] Alexandros Iosifidis, Anastasios Tefas, and Ioannis Pitas. View-invariant action recognition based on artificial neural networks. *IEEE transactions on neural networks and learning systems*, 23(3):412–424, 2012. 3
- [26] Allan Jabri, Andrew Owens, and Alexei Efros. Space-time correspondence as a contrastive random walk. *Advances in neural information processing systems*, 33:19545–19560, 2020. 3

- 688 [27] Joel Janai, Fatma Guney, Anurag Ranjan, Michael Black,
689 and Andreas Geiger. Unsupervised learning of multi-frame
690 optical flow with occlusions. In *Proceedings of the European
691 conference on computer vision (ECCV)*, pages 690–
692 706, 2018. 2
- 693 [28] Zhenyu Jiang, Hanwen Jiang, and Yuke Zhu. Doduo: Dense
694 visual correspondence from unsupervised semantic-aware
695 flow. In *arXiv preprint arXiv:2309.15110*, 2023. 1
- 696 [29] Zhouqiang Jiang, Bowen Wang, Tong Xiang, Zhaofeng Niu,
697 Hong Tang, Guangshun Li, and Liangzhi Li. Concatenated
698 masked autoencoders as spatial-temporal learner. *arXiv
699 preprint arXiv:2311.00961*, 2023. 4, 1
- 700 [30] Rico Jonschkowski, Austin Stone, Jonathan T Barron, Ariel
701 Gordon, Kurt Konolige, and Anelia Angelova. What mat-
702 ters in unsupervised optical flow. In *Computer Vision–ECCV
703 2020: 16th European Conference, Glasgow, UK, August 23–
704 28, 2020, Proceedings, Part II 16*, pages 557–572. Springer,
705 2020. 2
- 706 [31] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang,
707 Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola,
708 Tim Green, Trevor Back, Paul Natsev, et al. The kinetics hu-
709 man action video dataset. *arXiv preprint arXiv:1705.06950*,
710 2017. 6, 1
- 711 [32] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao,
712 Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer White-
713 head, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and
714 Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023.
715 3
- 716 [33] Liulei Li, Wenguan Wang, Tianfei Zhou, Jianwu Li, and Yi
717 Yang. Unified mask embedding and correspondence learning
718 for self-supervised video segmentation. In *CVPR*, 2023. 3, 2
- 719 [34] Siyuan Li, Lei Ke, Martin Danelljan, Luigi Piccinelli, Mattia
720 Segu, Luc Van Gool, and Fisher Yu. Matching anything by
721 segmenting anything. *CVPR*, 2024. 3, 2
- 722 [35] Pengpeng Liu, Michael Lyu, Irwin King, and Jia Xu. Self-
723 low: Self-supervised learning of optical flow. In *Proceedings
724 of the IEEE/CVF conference on computer vision and pattern
725 recognition*, pages 4571–4580, 2019. 2
- 726 [36] Maxime Oquab, Timothée Darzet, Theo Moutakanni, Huy V.
727 Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez,
728 Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Rus-
729 sell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-
730 Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nico-
731 las Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou,
732 Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bo-
733 janowski. Dinov2: Learning robust visual features without
734 supervision, 2023. 3, 7, 1
- 735 [37] Gensheng Pei, Tao Chen, Xiruo Jiang, Huafeng Liu, Zeren
736 Sun, and Yazhou Yao. Videomac: Video masked autoen-
737 coders meet convnets. In *CVPR*, 2024. 1, 4, 6
- 738 [38] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Ar-
739 beláez, Alexander Sorkine-Hornung, and Luc Van Gool.
740 The 2017 davis challenge on video object segmentation.
741 *arXiv:1704.00675*, 2017. 2, 6, 7, 8, 1
- 742 [39] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario
743 Amodei, and Ilya Sutskever. Language models are unsuper-
744 vised multitask learners. 2019. 1
- 40 [40] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
745 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
746 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning
747 transferable visual models from natural language supervi-
748 sion. In *International conference on machine learning*, pages
749 8748–8763. PMLR, 2021. 3
- 750 [41] Tanner Schmidt, Richard Newcombe, and Dieter Fox. Self-
751 supervised visual descriptor learning for dense correspon-
752 dence. *IEEE Robotics and Automation Letters*, 2(2):420–
753 427, 2017. 2
- 754 [42] Jinghuan Shang and Michael S Ryoo. Self-supervised dis-
755 entangled representation learning for third-person imitation
756 learning. In *2021 IEEE/RSJ International Conference on In-
757 telligent Robots and Systems (IROS)*, pages 214–221. IEEE,
758 2021. 1
- 759 [43] Xi Shen, Alexei A Efros, Armand Joulin, and Mathieu
760 Aubry. Learning co-segmentation by segment swapping for
761 retrieval and discovery. In *Proceedings of the IEEE/CVF
762 Conference on Computer Vision and Pattern Recognition*,
763 pages 5082–5092, 2022. 7
- 764 [44] Bradly C Stadie, Pieter Abbeel, and Ilya Sutskever. Third-
765 person imitation learning. *arXiv preprint arXiv:1703.01703*,
766 2017. 1
- 767 [45] Carole H Sudre, Wenqi Li, Tom Vercauteren, Sébastien
768 Ourselin, and M Jorge Cardoso. Generalised dice overlap as
769 a deep learning loss function for highly unbalanced segmen-
770 tations. In *Deep Learning in Medical Image Analysis and
771 Multimodal Learning for Clinical Decision Support: Third
772 International Workshop, DLMIA 2017, and 7th International
773 Workshop, ML-CDS 2017, Held in Conjunction with MIC-
774 CAI 2017, Québec City, QC, Canada, September 14, Pro-
775 ceedings 3*, pages 240–248. Springer, 2017. 5
- 776 [46] Hao Tang, Kevin J Liang, Kristen Grauman, Matt Feiszli,
777 and Weiyao Wang. Egotracks: A long-term egocentric vi-
778 sual object tracking dataset. *Advances in Neural Information
779 Processing Systems*, 36, 2024. 1
- 780 [47] Yansong Tang, Zhenyu Jiang, Zhenda Xie, Yue Cao, Zheng
781 Zhang, Philip HS Torr, and Han Hu. Breaking shortcut:
782 Exploring fully convolutional cycle-consistency for video
783 correspondence learning. *arXiv preprint arXiv:2105.05838*,
784 2021. 2, 3
- 785 [48] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field
786 transforms for optical flow. In *Computer Vision–ECCV
787 2020: 16th European Conference, Glasgow, UK, August 23–
788 28, 2020, Proceedings, Part II 16*, pages 402–419. Springer,
789 2020. 2
- 790 [49] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang.
791 VideoMAE: Masked autoencoders are data-efficient learners
792 for self-supervised video pre-training. In *Advances in Neural
793 Information Processing Systems*, 2022. 4
- 794 [50] Prune Truong, Martin Danelljan, Fisher Yu, and Luc
795 Van Gool. Warp consistency for unsupervised learning of
796 dense correspondences. In *Proceedings of the IEEE/CVF
797 international conference on computer vision*, pages 10346–
798 10356, 2021. 2
- 799 [51] Prune Truong, Martin Danelljan, Fisher Yu, and Luc
800 Van Gool. Probabilistic warp consistency for weakly-
801 supervised semantic correspondences. In *Proceedings of
802*

- 803 *the IEEE/CVF Conference on Computer Vision and Pattern*
804 *Recognition*, pages 8708–8718, 2022. 2
- 805 [52] Carl Vondrick, Abhinav Shrivastava, Alireza Fathi, Sergio
806 Guadarrama, and Kevin Murphy. Tracking emerges by col-
807 orizing videos. In *Proceedings of the European conference*
808 *on computer vision (ECCV)*, pages 391–408, 2018. 1, 2, 4
- 809 [53] Qianqian Wang, Xiaowei Zhou, Bharath Hariharan, and
810 Noah Snavely. Learning feature descriptors using camera
811 pose supervision. In *Computer Vision–ECCV 2020: 16th*
812 *European Conference, Glasgow, UK, August 23–28, 2020,*
813 *Proceedings, Part I 16*, pages 757–774. Springer, 2020. 2
- 814 [54] Xiaolong Wang, Allan Jabri, and Alexei A. Efros. Learn-
815 ing correspondence from the cycle-consistency of time. In
816 *CVPR*, 2019. 1, 2, 3
- 817 [55] Xudong Wang, Rohit Girdhar, Stella X Yu, and Ishan Misra.
818 Cut and learn for unsupervised object detection and instance
819 segmentation. In *Proceedings of the IEEE/CVF Conference*
820 *on Computer Vision and Pattern Recognition*, pages 3124–
821 3134, 2023. 3
- 822 [56] XuDong Wang, Jingfeng Yang, and Trevor Darrell. Segment
823 anything without supervision. *NeurIPS* 2024, 2024. 3
- 824 [57] Yang Wang, Yi Yang, Zhenheng Yang, Liang Zhao, Peng
825 Wang, and Wei Xu. Occlusion aware unsupervised learning
826 of optical flow. In *Proceedings of the IEEE conference on*
827 *computer vision and pattern recognition*, pages 4884–4893,
828 2018. 2
- 829 [58] Daniel Weinland, Mustafa Özysal, and Pascal Fua. Mak-
830 ing action recognition robust to occlusions and viewpoint
831 changes. In *Computer Vision–ECCV 2010: 11th European*
832 *Conference on Computer Vision, Heraklion, Crete, Greece,*
833 *September 5–11, 2010, Proceedings, Part III 11*, pages 635–
834 648. Springer, 2010. 3
- 835 [59] Philippe Weinzaepfel, Thomas Lucas, Vincent Leroy,
836 Yohann Cabon, Vaibhav Arora, Romain Brégier, Gabriela
837 Csurka, Leonid Antsfeld, Boris Chidlovskii, and Jérôme Re-
838 vraud. CroCo v2: Improved Cross-view Completion Pre-
839 training for Stereo Matching and Optical Flow. In *ICCV*,
840 2023. 3, 6, 7, 8, 1
- 841 [60] Weinzaepfel, Philippe and Leroy, Vincent and Lucas,
842 Thomas and Brégier, Romain and Cabon, Yohann and
843 Arora, Vaibhav and Antsfeld, Leonid and Chidlovskii, Boris
844 and Csurka, Gabriela and Revaud Jérôme. CroCo: Self-
845 Supervised Pre-training for 3D Vision Tasks by Cross-View
846 Completion. In *NeurIPS*, 2022. 3, 4, 6
- 847 [61] Jiarui Xu and Xiaolong Wang. Rethinking self-supervised
848 correspondence learning: A video frame-level similarity per-
849 spective. In *Proceedings of the IEEE/CVF International*
850 *Conference on Computer Vision*, pages 10075–10085, 2021.
851 3
- 852 [62] Tinghui Zhou, Philipp Krahenbuhl, Mathieu Aubry, Qixing
853 Huang, and Alexei A Efros. Learning dense correspon-
854 dence via 3d-guided cycle consistency. In *Proceedings of*
855 *the IEEE conference on computer vision and pattern recog-*
856 *nition*, pages 117–126, 2016. 2, 3

Self-Supervised Cross-View Correspondence with Predictive Cycle Consistency

Supplementary Material

857

6. Implementation Details

858

In this section, we detail the specific setups for the following components:

859

1. Our Grayscale Colorization model (Section 3.2),
2. Our ViT Correspondence models (Section 3.5),
3. Our implementation of predictive approaches SiamMAE [18] and CroCoV2 [59], and
4. Our parameterization of SAM (Segment Anything Model) for tracking.

860

6.1. Hyperparameters

861

6.1.1 Grayscale and Correspondence Model

862

We implement our Grayscale Colorization model and Correspondence Model using the CroCoV2 [59] base architecture. Starting from the CroCoV2 Base-Decoder checkpoint, we continue pretraining with either the Grayscale Colorization objective or the original Cross-View MAE objective from CroCoV2 on datasets such as EgoExo4D [16] or Kinetics-400 [31]. The hyperparameters for continued pretraining are listed in Table 5.

863

Table 5. Hyperparameters for Grayscale Colorization and Correspondence Model Continual Training.

	Grayscale Colorization (Sec. 3.2)	Correspondence Model (Sec. 3.5)
Encoder Layers	12	12
Encoder Embed Dim	768	768
Decoder Layers	12	12
Decoder Embed Dim	768	768
MLP Dim	3072	3072
Learning rate	1.5×10^{-4}	1.5×10^{-4}
Adam β_1 / β_2	0.9 / 0.98	0.9 / 0.98
Weight decay	0.01	0.01
Learning rate schedule	Linear Decay	Linear Decay
Dropout	0.1	0.1
Warmup updates	2,000	2,000
Batch size	256	256
Updates	10,000	10,000
Training Objective	Colorization (RBG MSE Loss)	MAE
Kinetics-400 Time Gap	4-48 Frames	4-48 Frames

864

This continued pretraining results in the Grayscale Colorization model that we use to initialize PCC, extracting correspondence with the technique in Sec. 3.3. For our final PCC Correspondence Model, we further train using PCC pseudolabels as described in Section 3.5. Table 6 outlines the hyperparameters used for this additional training.

865

6.1.2 Baseline Implementation

866

To ensure fairness during evaluation, we continually pretrain CroCoV2 [59] and SiamMAE [18] on EgoExo4D [16] before measuring correspondence. The settings for continually pretraining CroCoV2 are outlined in Section 6.1.

Table 6. PCC Correspondence Model Hyperparameters. For each domain (EgoExo4D or Kinetics-400) we initialize our PCC Correspondence Model parameters with a continually pretrained MAE (Table 5)

	EgoExo4D [16] Correspondence	Kinetics-400 [16] Correspondence
Encoder Layers	12	12
Encoder Embed Dim	768	768
Decoder Layers	12	12
Decoder Embed Dim	768	768
MLP Dim	3072	3072
Learning rate	1.5×10^{-4}	1.5×10^{-4}
Adam β_1 / β_2	0.9 / 0.98	0.9 / 0.98
Weight decay	0.01	0.01
Learning rate schedule	Linear Decay	Linear Decay
Dropout	0.1	0.1
Warmup updates	2,000	2,000
Batch size	256	256
Updates	10,000	10,000
Training Objective	DICE + BCE	DICE + BCE
Kinetics-400 Time Gap	-	60 Frame Gap (2 sec)
EgoExo Parameters	50/50 Ego→Exo/Exo→Ego	-
Image size	240x240 (Ego) 240x416 (Exo)	224x224

Since the SiamMAE [18] code and checkpoints are not publicly available, we reimplement their approach by adapting the published CAT-MAE [29] codebase and checkpoints. We continually train SiamMAE using the CAT-MAE hyperparameters on Kinetics-400 for 60,000 steps with a batch size of 256. To validate our reimplementation, we evaluate our model on the DAVIS-2017 validation set [38], achieving a $\mathcal{J} & \mathcal{F}_m$ score of 70.6, closely matching the original SiamMAE score of 71.4. For EgoExo4D, we continually pretrain this checkpoint for an additional 60000 steps at a batch size of 256, otherwise using the same settings.

We exclude DINO [36] models from continual pretraining on EgoExo4D due to a lack of diversity of data for image augmentation (EgoExo4D only has 123 unique sites used for data collection). Additionally, models employing exponentially moving average teachers require extensive tuning of the moving average temperature, making continual pretraining more challenging.

For our baselines, we adapt the K-Nearest-Neighbor implementation from [54]. While originally designed for multiple video frames, we modify it to treat all evaluation scenarios as two-frame videos. The algorithm inherently supports different resolutions for the first frame query and subsequent frames, accommodating the differing aspect ratios of Ego-view and Exo-view images. As detailed in Section 4.2, we resample all Ego and Exo videos to have a minimum resolution of 480p and perform a grid search to optimize the parameter k and the temperature. For EgoExo4D evaluation, we omit the neighborhood size parameter, as there is no spatial continuity between Ego and Exo views.

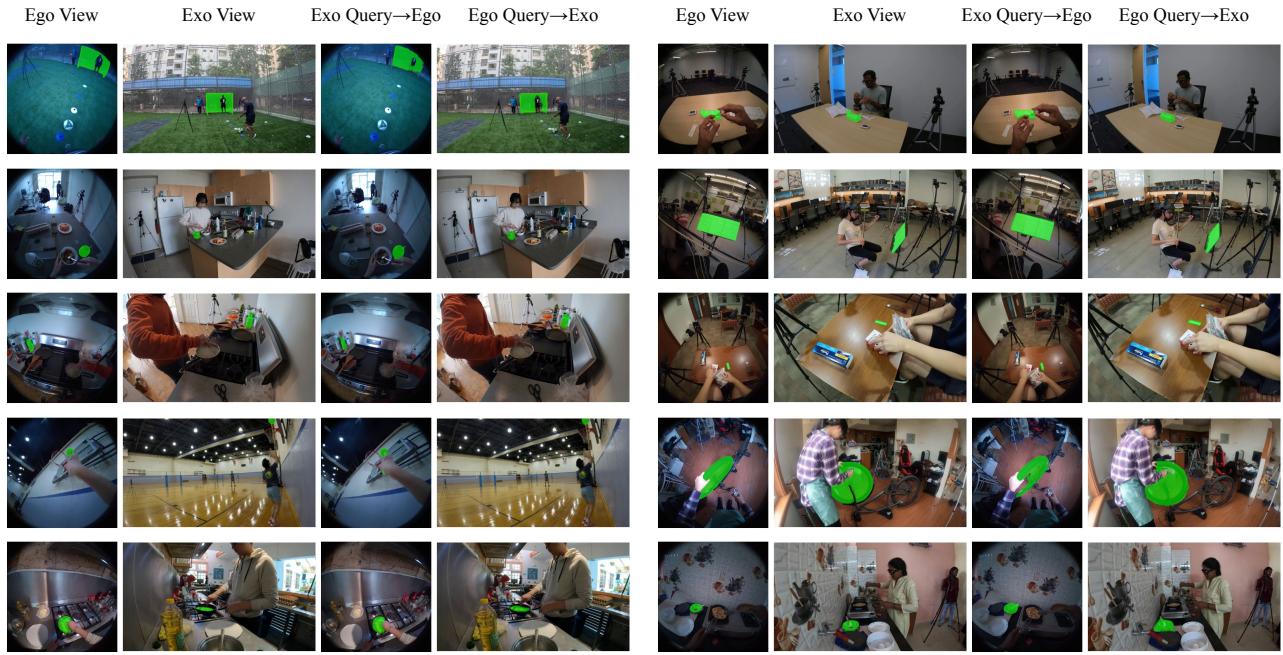


Figure 9. Additional Qualitative Results on the EgoExo4D Correspondence Benchmark.

918 However, for DAVIS-17 and LVOS evaluations, we inde-
 919 pendently grid search the neighborhood size parameter for
 920 each temporal distance.

921 6.1.3 SAM Configuration

922 To extract image segmentations from raw images in
 923 EgoExo4D, we use SAM with the standard point-grid
 924 prompting configuration, as demonstrated in [7, 33]. We
 925 note that this is different from the configuration of MASA
 926 [34], which uses bounding boxes extracted from an off-the-
 927 shelf object detection model using textual object descrip-
 928 tions. Because SAM is traditionally run on third-person
 929 videos, we gridsearch the Predicted IoU Threshold (0.88)
 930 and the Stability Score Threshold to (0.94) to have the high-
 931 est IoU with ground truth object segmentation masks from
 932 the EgoExo4D validation set.

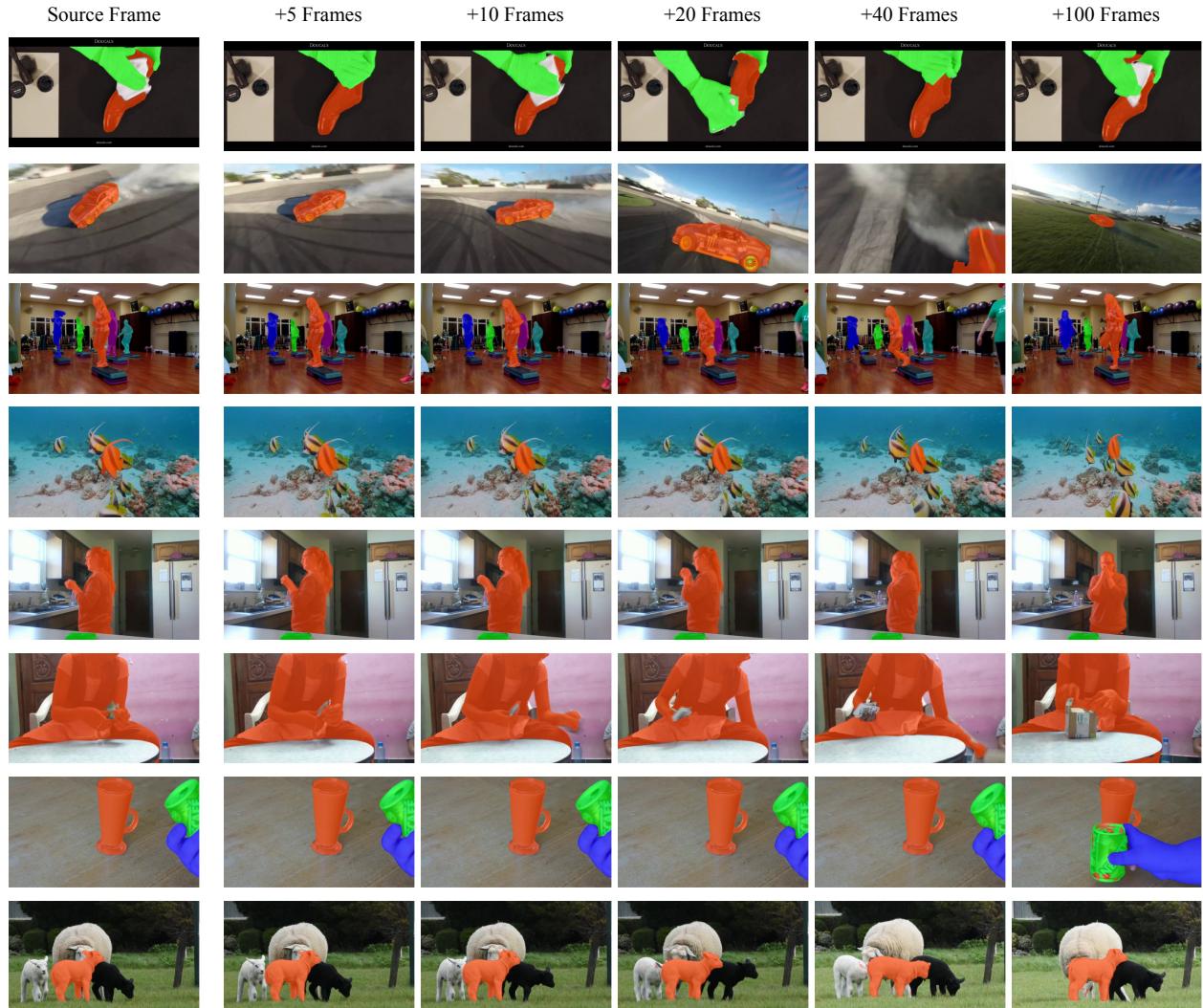


Figure 10. Additional Qualitative Results on LVOS with various frame gaps.