

CycDA: Unsupervised Cycle Domain Adaptation to Learn from Image to Video

Wei Lin^{1,4} Anna Kukleva² Kunyang Sun^{1,5} Horst Possegger¹ Hilde Kuehne³ Horst Bischof¹
¹Graz University of Technology ²Max-Planck-Institute for Informatics ³Goethe University Frankfurt
⁴Christian Doppler Laboratory for Semantic 3D Computer Vision ⁵Southeast University

Motivation

Image-to-video adaptation

- Web images are a labeling-free data source for action recognition
- Webly-labeled images to unlabeled videos

But ...

- **Spatial domain shift** between web images and video frames
- **Modality gap** between images and videos

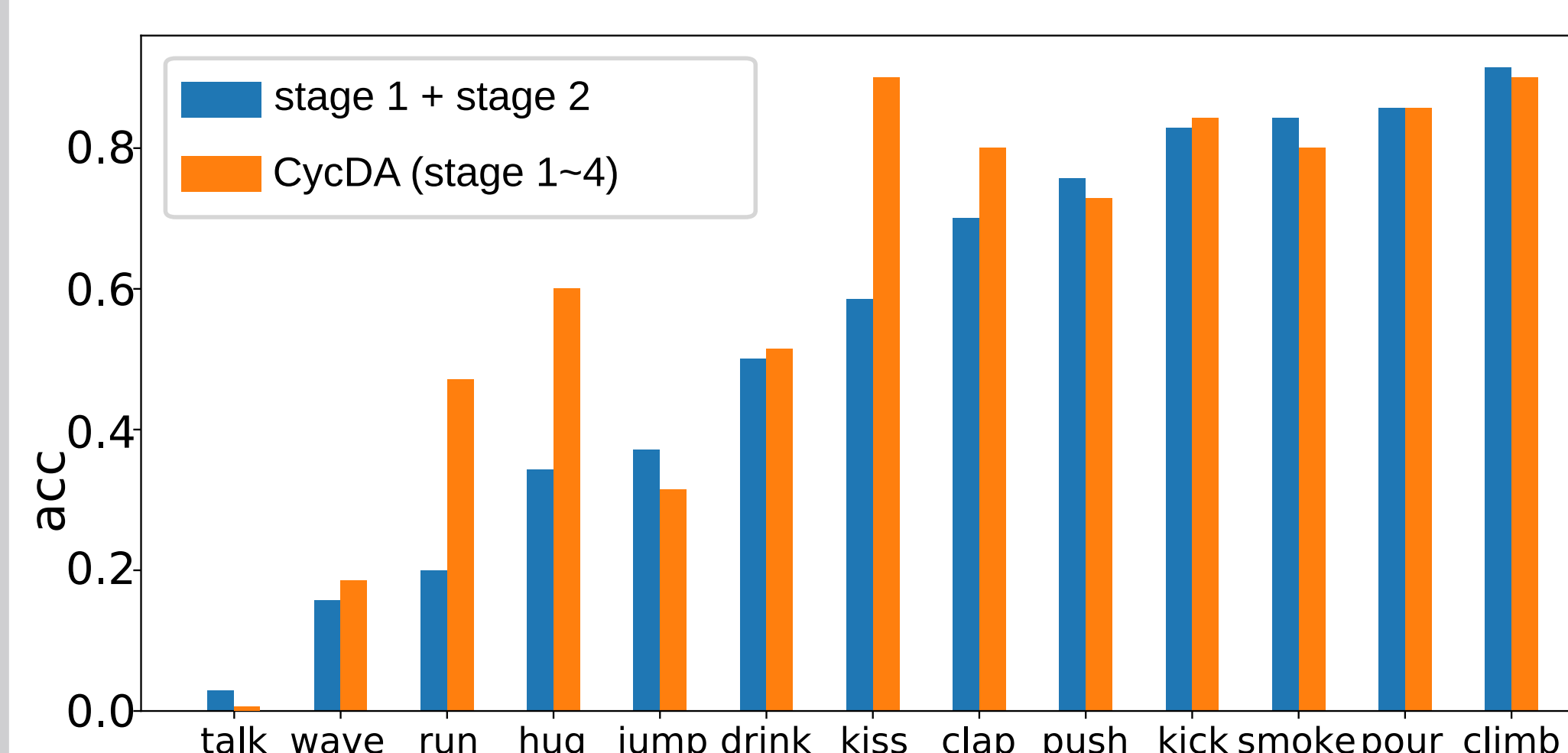


Overview

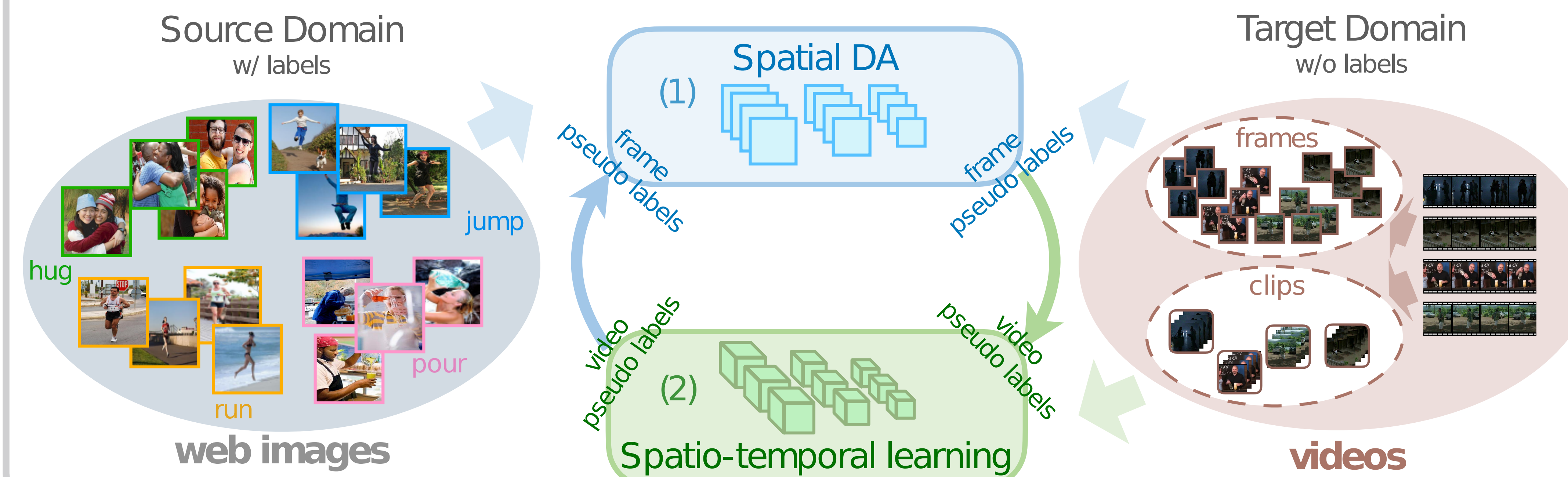
CycDA

- solves image-to-video adaptation by decoupling a) **domain alignment** b) **spatio-temporal learning**
- is a cyclic alternation of 4 stages between spatial and spatio-temporal learning
- has extensive evaluations with SOTA results on **image-to-video adaptation** and **mixed-source (image+video) to video adaptation**

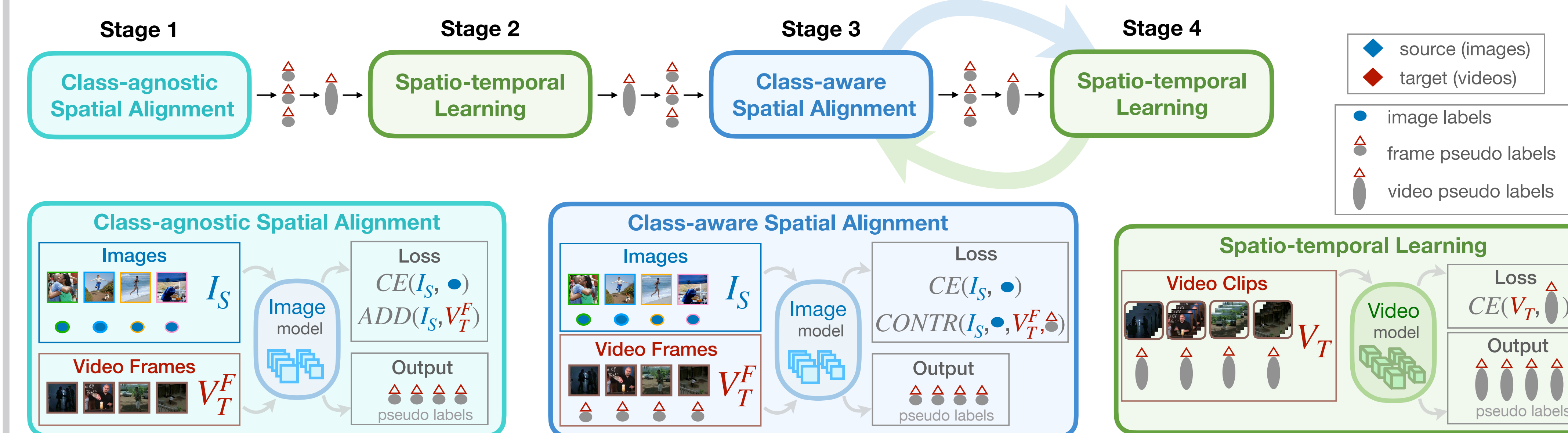
Category-level Pseudo Labels



CycDA



Stages



Stage 1: Class-agnostic domain alignment

- Task: domain alignment without category information. The classes of source and target data could be misaligned.
- Model: image model ϕ^I
- Objective: supervised cross entropy loss and adversarial domain discrimination loss

$$\min_{\theta_E^I, \theta_C^I} \mathcal{L}_{CE}(I_S) + \beta \max_{\theta_E^I} \min_{\theta_D^I} \mathcal{L}_{ADD}(I_S, V_T^F)$$

Stage 2 & Stage 4: Spatio-temporal learning

- Task: spatio-temporal learning with pseudo label target data
 - Pseudo labeling: use image model ϕ^I to predict pseudo labels on target frames, and temporally aggregate frame-level labels into video-level labels
 - Model: video model ϕ^V
 - Objective: supervised cross entropy loss
- $$\min_{\theta_E^V, \theta_C^V} \hat{\mathcal{L}}_{CE}(\hat{V}_T)$$

Stage 3: Class-aware domain alignment

- Task: Domain alignment on the category level
- Pseudo labeling: use the video model ϕ^V to generate video-level labels, and disseminate into frame-level labels
- Model: image model ϕ^I
- Objective: cross-domain contrastive learning

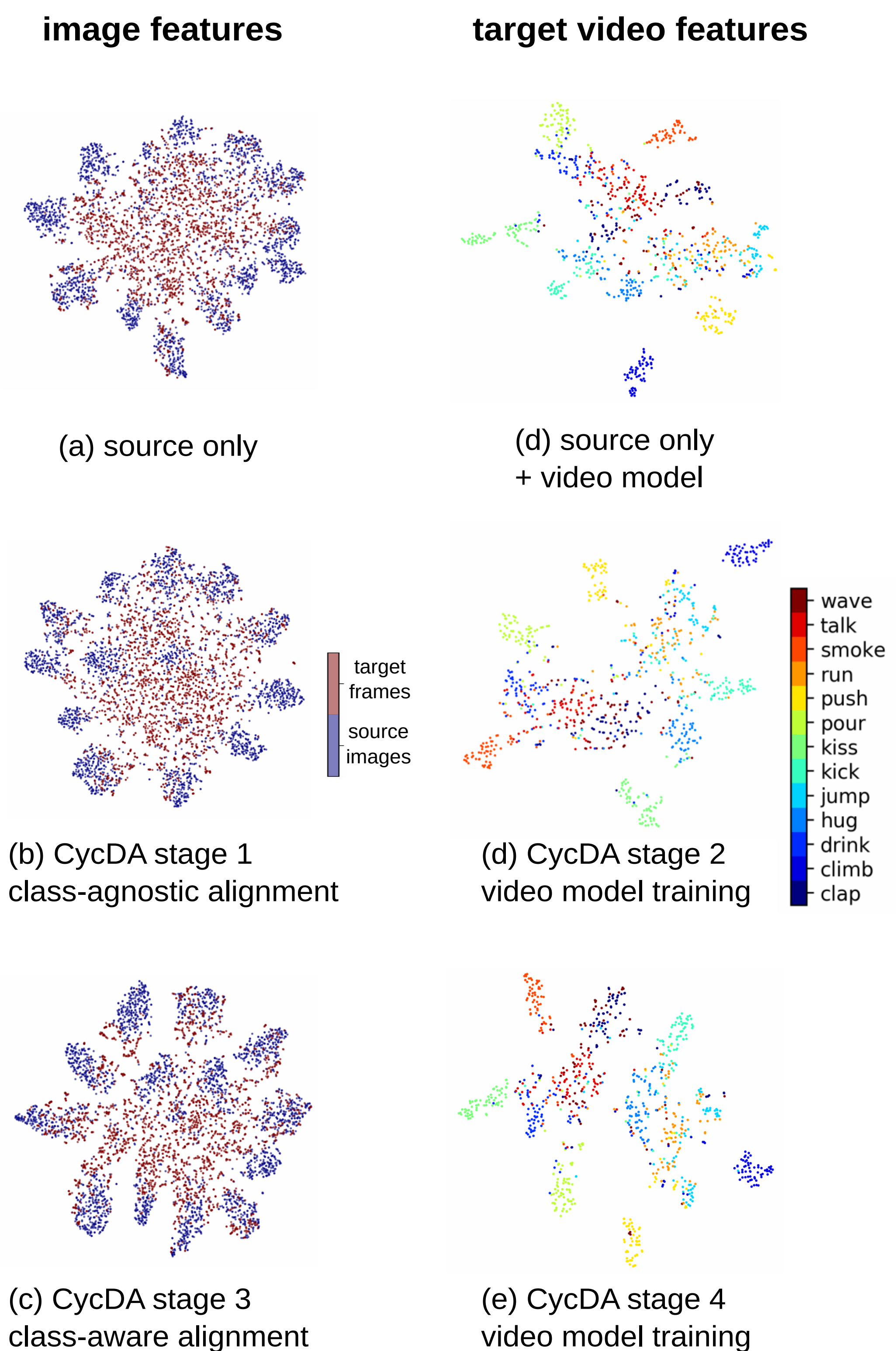
$$\mathcal{L}_{CONTR}(I_S, V_T^F) = - \sum_{z_j^F \in Z_T^F} \log \frac{h(z_j^F, z_{j+}^I)}{h(z_j^F, z_{j+}^I) + h(z_j^F, z_{j-}^I)}$$

Pseudo labeled target sample $z_j^F \in Z_T^F$
 Positive sample from source $z_{j+}^I \in I_S$
 Negative sample from source $z_{j-}^I \in I_S$

Cycling of the Stages

Stage 3 and stage 4 can be performed iteratively.

Feature Distribution



Nearest Neighbor Search

