

Action Segmentation with Joint Self-Supervised Temporal Domain Adaptation



Min-Hung Chen^{1*}



Baopu Li²



Yingze Bao²



Ghassan AlRegib¹



Zsolt Kira¹

¹Georgia Institute of Technology ²Baidu USA

June 17, 2020

[Paper] <https://arxiv.org/abs/2003.02824>

[Project] <https://minhungchen.netlify.app/project/cdas>

Action Segmentation

Action segmentation = Action Recognition + Temporal Segmentation

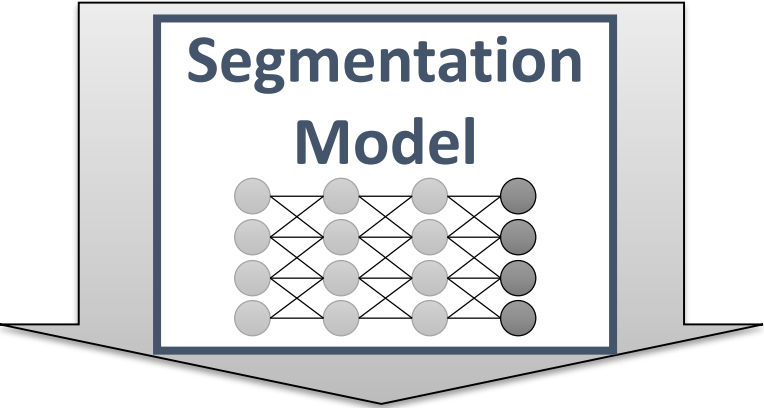
Input Video

Begin



Make milk

Time →



Output Predictions

Begin



background

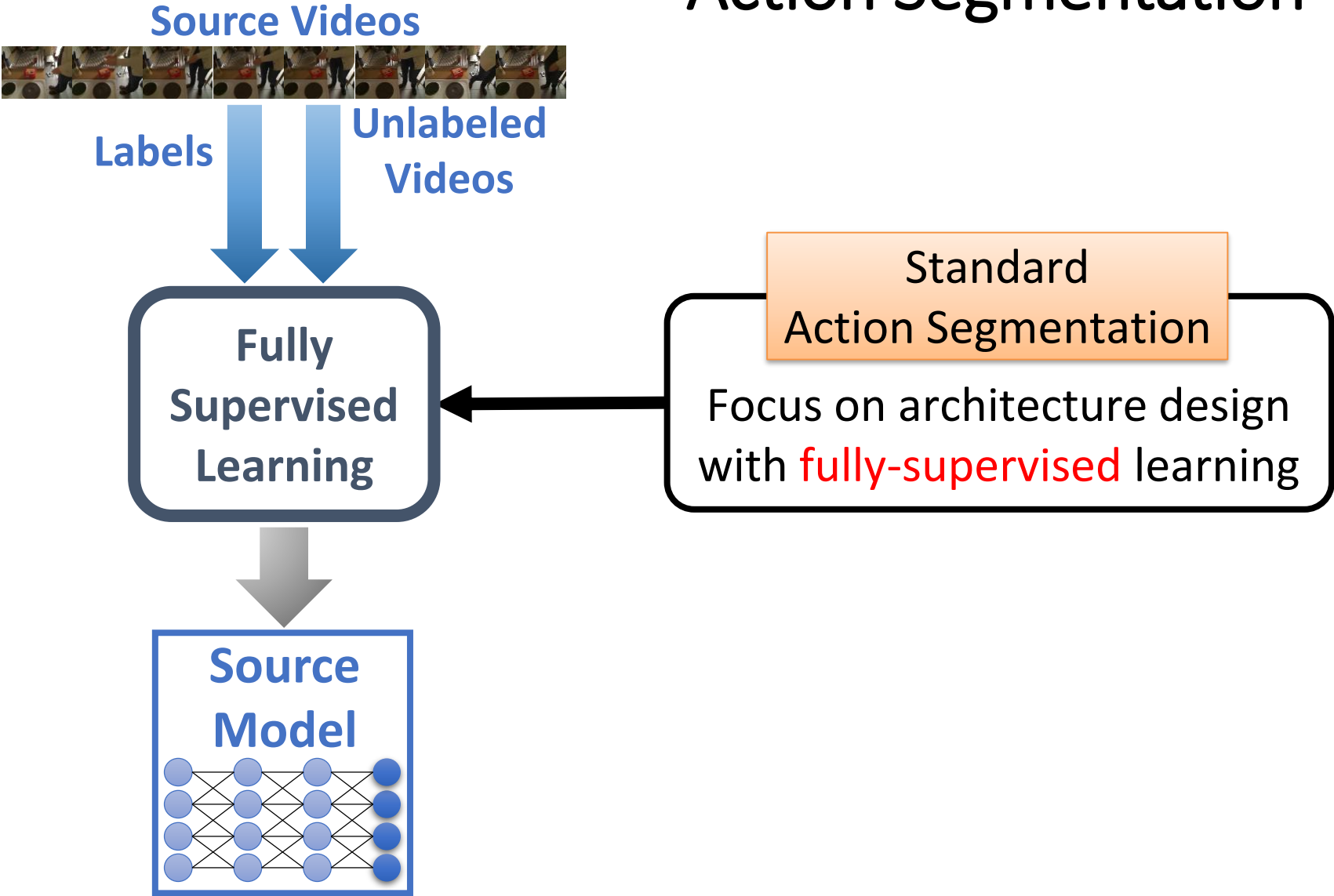
take cup

spoon powder

pour milk

Time →

Action Segmentation



Challenge

Source Videos

Target Videos



Labels

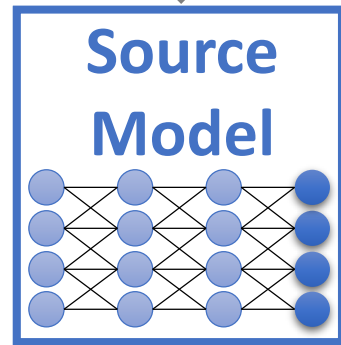
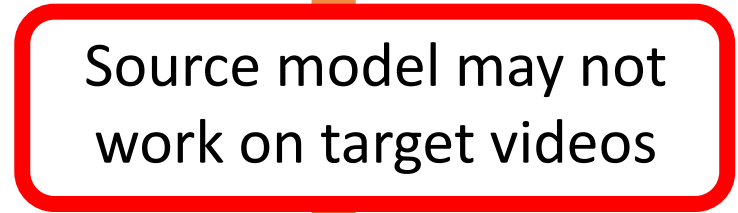
Unlabeled Videos

Unlabeled Videos



Standard Action Segmentation

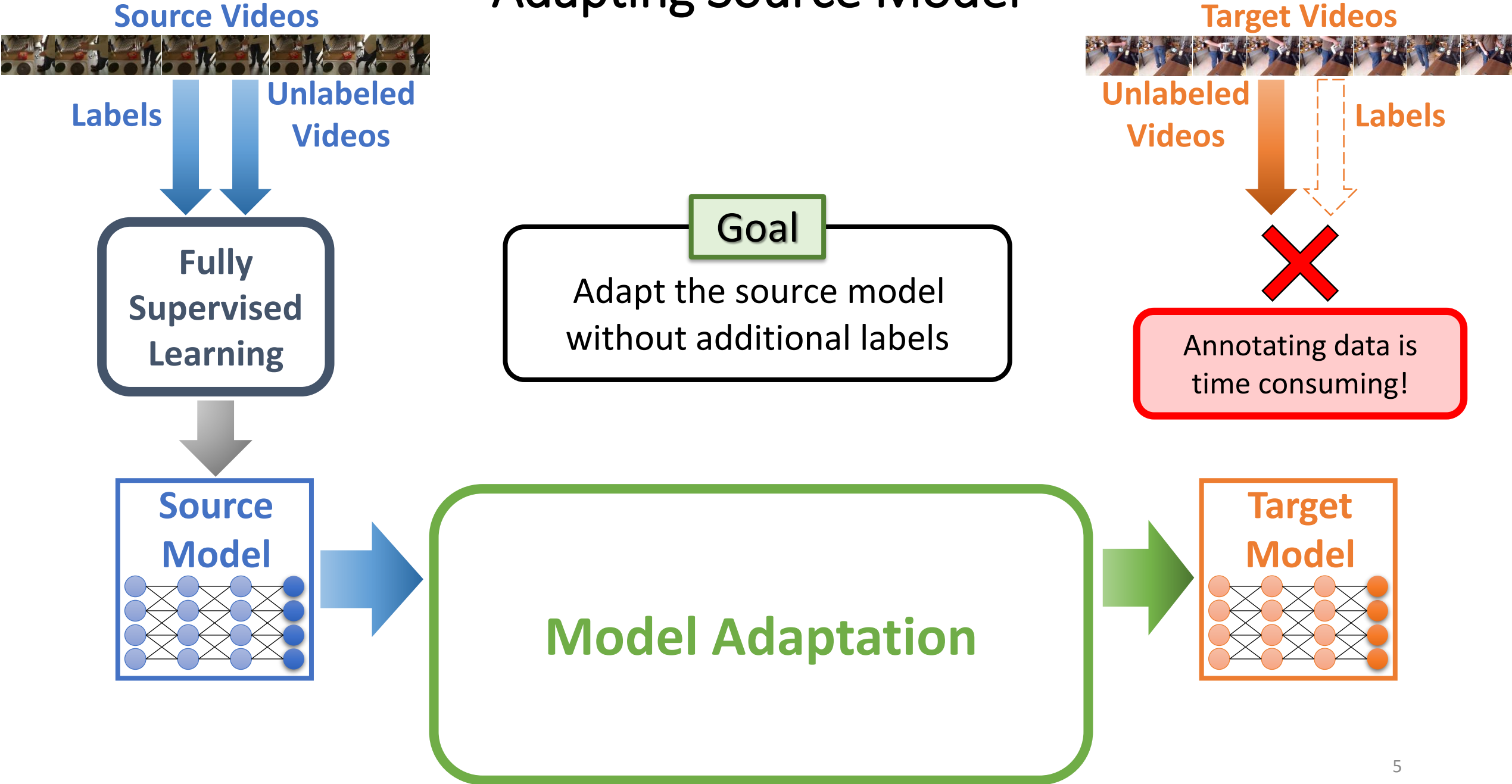
Focus on architecture design with **fully-supervised** learning

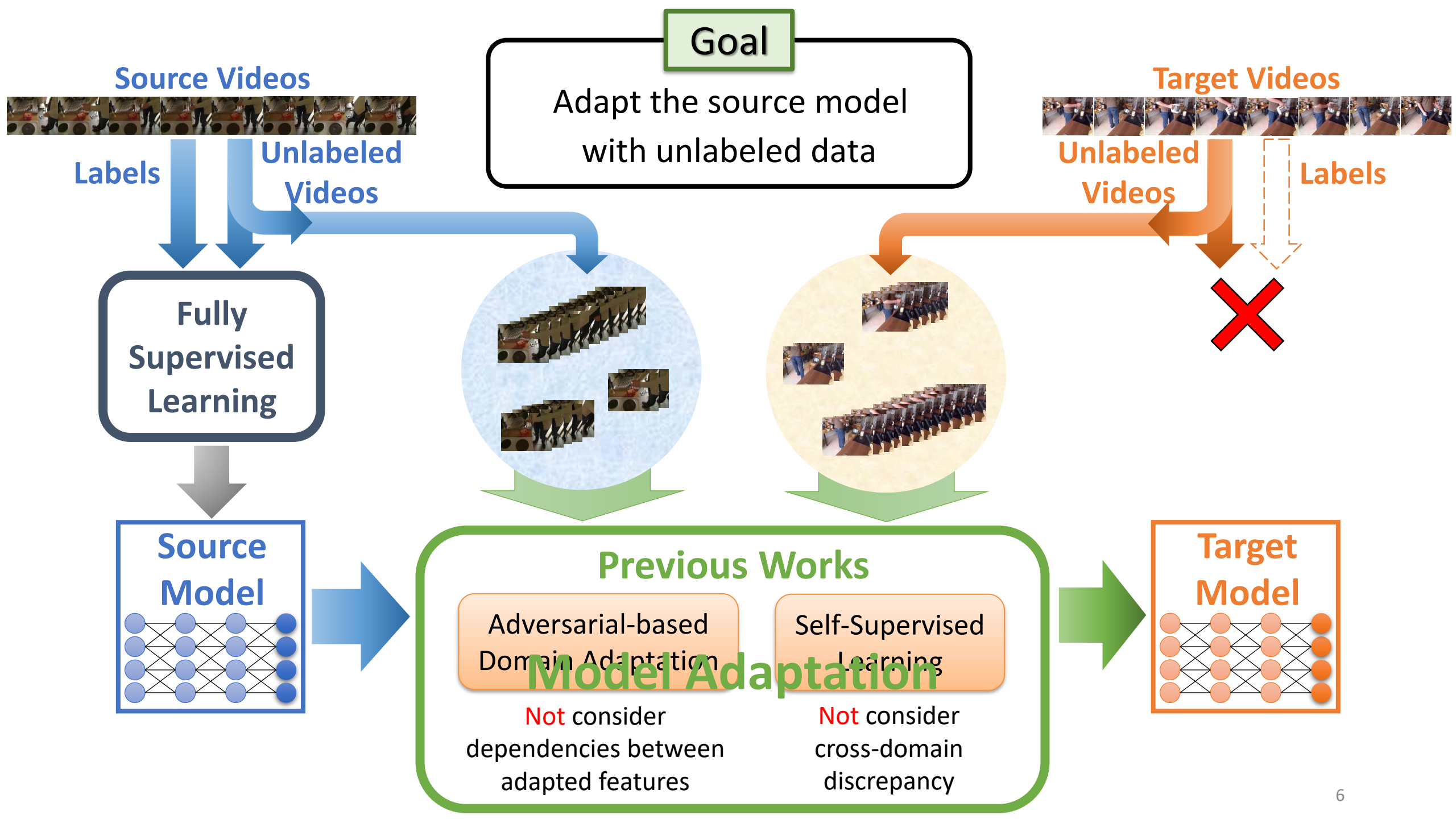


Source vs. Target

Different people perform the same action in different styles

Adapting Source Model





Temporal Domain Permutation

Previous Works

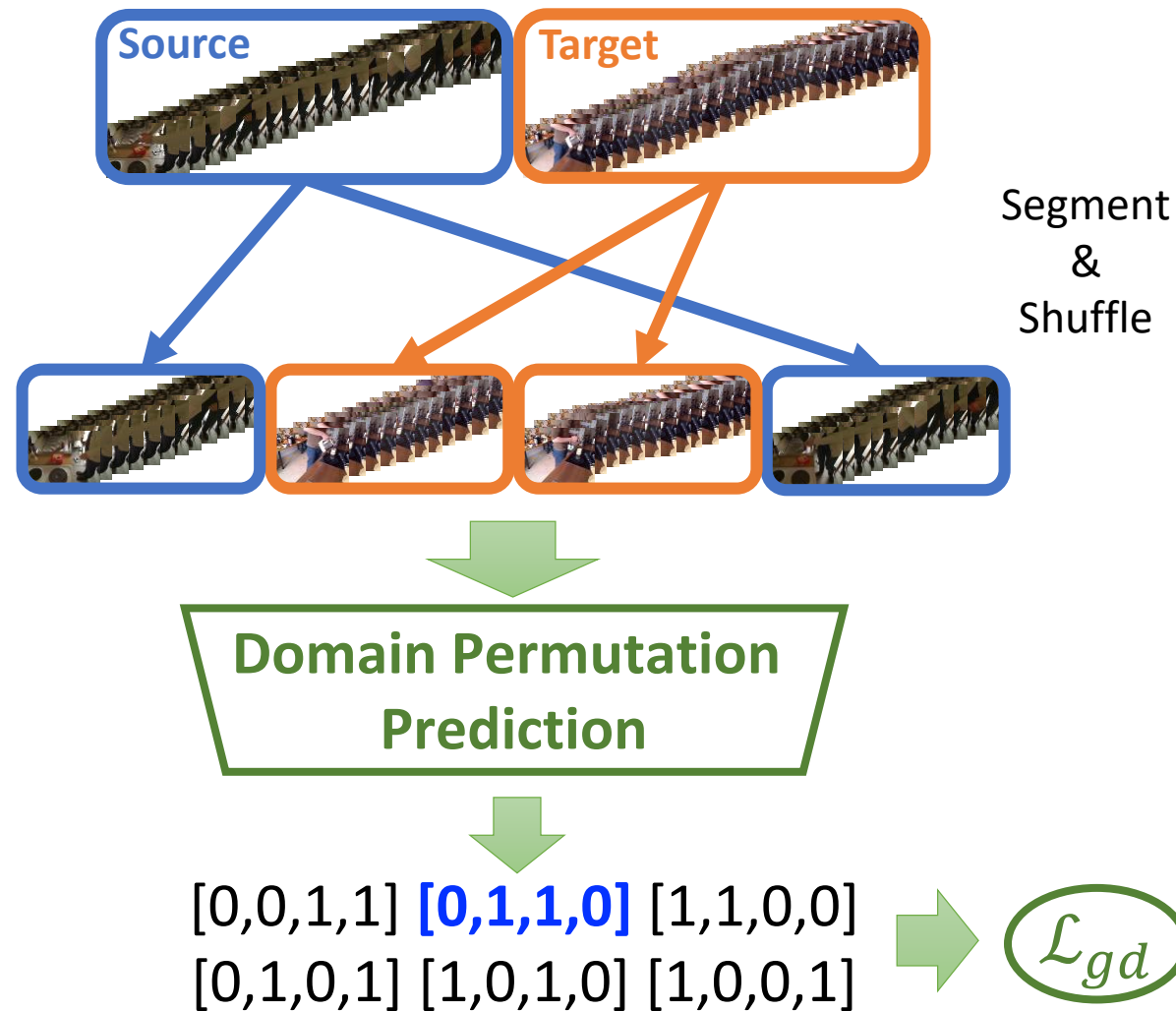
- Predict temporal orders
- Predict binary domains

Intuition

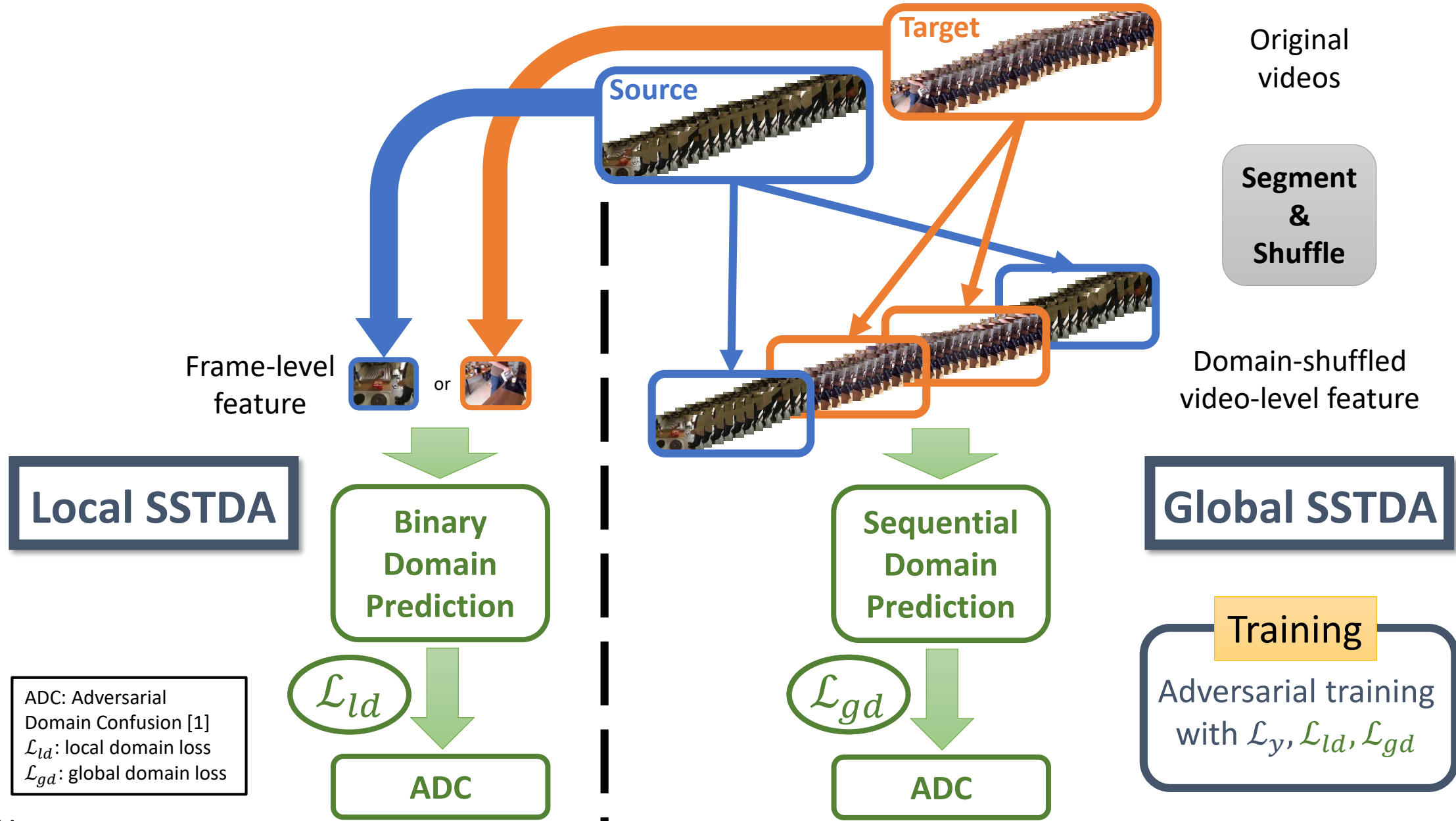
DA for classification \rightarrow domain classification
DA for segmentation \rightarrow domain segmentation

Our Method

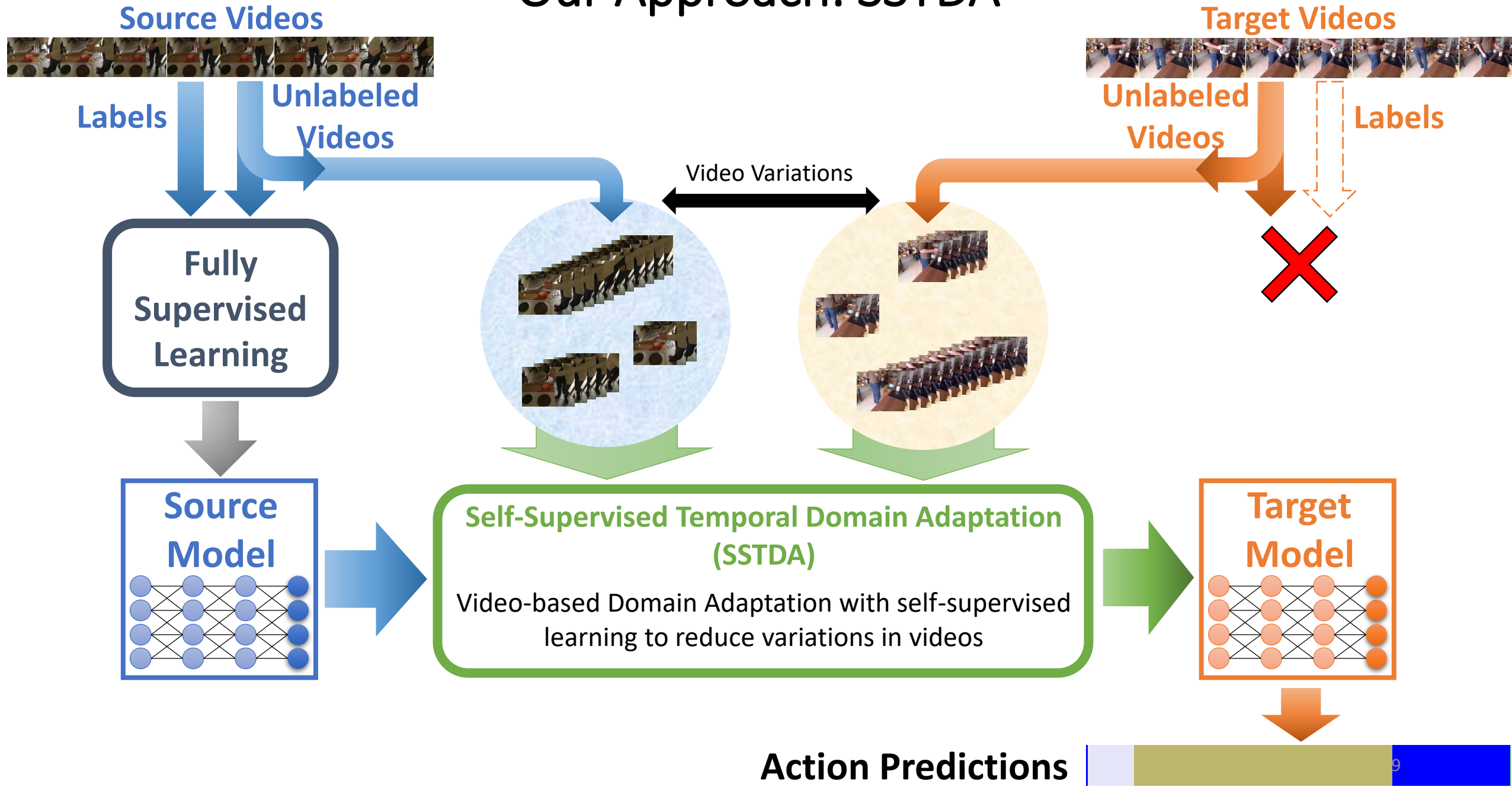
Predict **temporal permutations of domains**



Self-Supervised Temporal Domain Adaptation (SSTDA)



Our Approach: SSTDA



Experimental Results



50Salads [1]

Source-only: results from directly running the official released code of MS-TCN [2]

50Salads	F1@10	F1@25	F1@50	Edit score
Source-only [2]	75.4	73.4	65.2	68.9
Local SSTDA	79.2	77.8	70.3	72.0
SSTDA	83.0	81.5	73.8	75.8
SSTDA (65%)	77.7	75.0	66.2	69.3

Effectively exploit unlabeled target videos for action segmentation

Comparison: Unlabeled Target Videos

50Salads	F1@10	F1@25	F1@50	Edit score
Source-only	75.4	73.4	65.2	68.9
VCOP [1]	75.8	73.8	65.9	68.4
DANN [2]	79.2	77.8	70.3	72.0
JAN [3]	80.9	79.4	72.4	73.5
MADA [4]	79.6	77.4	70.0	72.4
MSTN [5]	79.3	77.6	71.5	72.1
MCD [6]	78.2	75.5	67.1	70.8
SWD [7]	78.2	76.2	67.4	71.6
SSTDA	83.0	81.5	73.8	75.8

Jointly adapt domains with **multiple temporal scales** can better address discrepancy problems for videos

Visualization: 50Salads



Ground Truth



Expectation



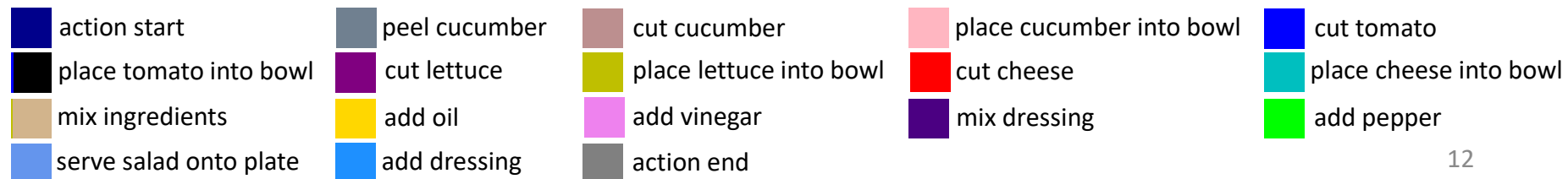
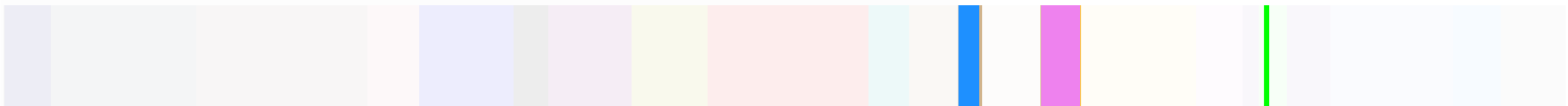
MS-TCN [1]



Local SSTDA

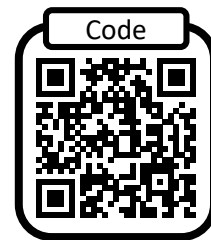
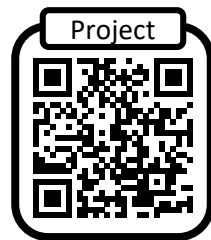


SSTDA



Summary

- Goal: adapt action segmentation models using unlabeled videos
- Approach: **Self-Supervised Temporal Domain Adaptation (SSTDA)**
 - Perform domain adaptation for multiple temporal scales
 - Learn feature representations with domain-invariant temporal dynamics
- Outperform other self-supervised methods and image-based DA methods
- Improve action segmentation by large margins using unlabeled target videos



Poster: #93 @ Session 2.4
Date: June 17 (Wed.)
Q&A Time: 16 - 18 & 04 - 06

[Paper] <https://arxiv.org/abs/2003.02824>

[Project] <https://minhungchen.netlify.app/project/cdas>

[Code] <https://github.com/cmhungsteve/SSTDA>

