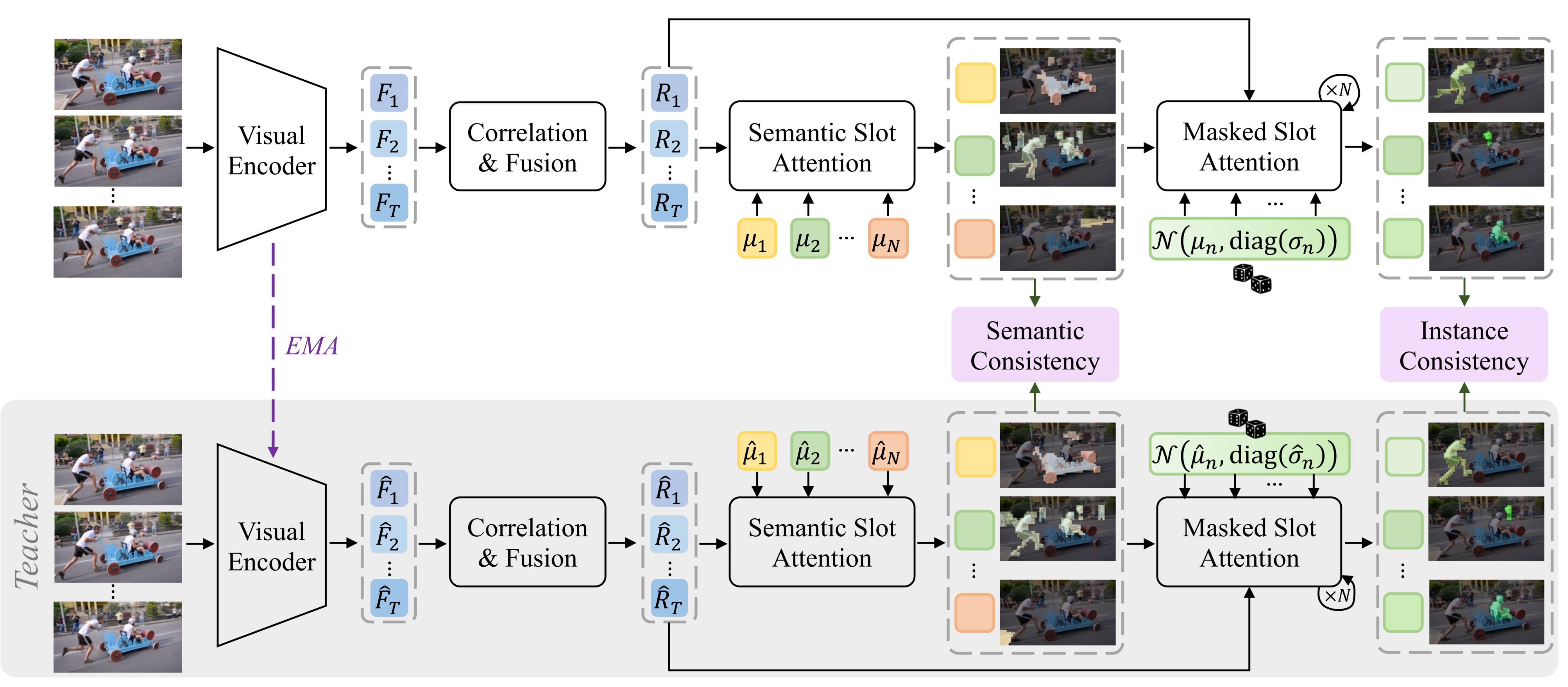


Semantics Meets Temporal Correspondence: Self-supervised Object-centric Learning in Videos ICCV23

Abstract:

- **Task:** Object-centric video representation learning, including unsupervised video object segmentation and spatio-temporal correspondence.
- **Motivation:** Jointly utilize high-level semantics and low-level correspondence to parse the object structures. There exist consistent semantic patterns for similar objects across scenes, while the temporal correspondence cues vary with specific scenes.
- Our solution: We propose a two-stage semantic-aware slot attention for semantic decomposition and instance discrimination. We distill temporally coherent object-centric representations in a fully selfsupervised manner.



- Semantic Consistency: First determine dense patch correspondence across time with optimal transport, then align the semantic distributions. $\mathcal{L}_{sem} = -\sum_{\substack{t,j=1\\t\neq j}}^{r} \sum_{\substack{u,v=1\\t\neq j}}^{r} \pi_{tj}^*[u,v] \widehat{\mathcal{M}}_j[n,v] \log \mathcal{M}_t[n,u]$ **Instance Consistency:** First filter out non-existing semantic centers and invalid instance slots, then encourage temporally coherent object representations.

$$\mathcal{L}_{obj} = -\sum_{\substack{t,j=1\\t\neq j}}^{T}\sum_{n=1}^{N}\sum_{p=1}^{P}\mathcal{I}_{t}[n,p] \left\{ \hat{\mathcal{I}}_{j}[n,\varepsilon(p)] \left| \left| \mathcal{O}_{t}[n,p] - \hat{\mathcal{O}}_{j}[n,\varepsilon(p)] \right| \right|_{2} + \sum_{\substack{q=1\\q\neq\varepsilon(p)}}^{P} \operatorname{relu}\left(\lambda - \left| \left| \mathcal{O}_{t}[n,p] - \hat{\mathcal{O}}_{j}[n,q] \right| \right|_{2} \right) \right\}$$

Rui Qian, Shuangrui Ding, Xian Liu, Dahua Lin

Method:

- Fuse the frame-wise feature and simple temporal correlation map to carry on rich semantics and temporal correspondence cues.
- Formulate a set of learnable Gaussian distributions consisting of mean and deviation vectors. The former represents potential semantic centers, the latter introduces random perturbation to capture temporal correspondence patterns.
- Develop two-stage semantic-aware slot attention, which first uses mean vectors as slot initialization for semantic decomposition, then performs random sampling around each semantic center to identify instances within the semantic area.

Experiment:

Unsupervised video object discovery: Single object segmentation on DAVIS-2016, SegTrack-v2, FBMS-59; Multiple object segmentation on DAVIS-2017-Unsupervised.

Model	RGB	Flow	DAVIS	S
CIS [89]	\checkmark	\checkmark	71.5	6
AMD [57]	\checkmark	×	57.8	5
DINO [12]	\checkmark	×	52.3	4
SIMO [53]	×	\checkmark	67.8	6
MG [88]	×	\checkmark	68.3	5
OCLR [84]	X	\checkmark	72.1	6
GWM [16]	\checkmark	\checkmark	71.2	6
SMTC	\checkmark	X	71.8	6
SMTC[†]	\checkmark	×	70.8	6

Model	Backbone	$\mathcal{J}\&\mathcal{F}$
DINOSAUR [72]	ViT-S/16	21.4
SMTC	ViT-S/16	40.5
SMTC	ResNet-50	39.0
RVOS [76]	ResNet-101	41.2
ProReduce [56]	ResNet-101	68.3

Dense spatio-temporal correspondence: Semi-supervised video object segmentation on DAVIS-2017; Pose tracking on JHMDB; Human part tracking on VIP.

		DAVIS		JHMDB		VIP	
Model	Backbone	$\mathcal{J}\&\mathcal{F}$	${\mathcal J}$	\mathcal{F}	PCK@0.1	PCK@0.2	mIoU
Supervised [40]	ResNet-50	66.0	63.7	68.4	59.2	78.3	39.5
MoCo [38]	ResNet-50	65.4	63.2	67.6	60.4	79.3	36.1
VFS [86]	ResNet-50	68.9	66.5	71.3	60.9	80.7	43.2
DINO [12]	ViT-S/16	61.8	60.2	63.4	45.4	75.2	37.9
TimeCycle [81]	ResNet-50	40.7	41.9	39.4	57.7	78.5	28.9
UVC [55]	ResNet-50	56.3	54.5	58.1	56.5	76.6	34.2
MAST [51]	ResNet-18	65.5	63.3	67.6	-	-	-
CRW [44]	ResNet-18	67.6	64.8	70.2	58.8	80.3	37.6
SFC [43]	ResNet-18+ResNet-50	71.2	68.3	74.0	61.9	83.0	38.4
SMTC	ViT-S/16	67.6	64.1	71.2	53.2	79.6	39.2
SMTC	ResNet-50	73.0	69.4	76.6	62.5	84.1	38.8

Conclusion:

- representation learning.





Unify semantic discrimination and temporal correspondence for object-centric video

Achieve fully self-supervised video object instance identification with semantic structure.

Efficiently transfer object knowledge to more general video understanding tasks.