

Dynamic Regions Graph Neural Networks for Spatio-Temporal Reasoning

Andrei Nicolicioiu* 1 Marius Leordeanu^{1 2 3} Iulia Duta* 1

¹Bitdefender, ²University Politehnica of Bucharest, ³Institute of Mathematics of the Romanian Academy

Contributions

- DyReG can discover salient regions that correlate with objects locations, while being unsupervised at the object-level.
- Our model creates localised graph nodes that improves the relational processing and obtains superior results on video classification tasks.
- Localising the nodes, as a form of hard attention leads to a more **explainable** model.

Our approach:

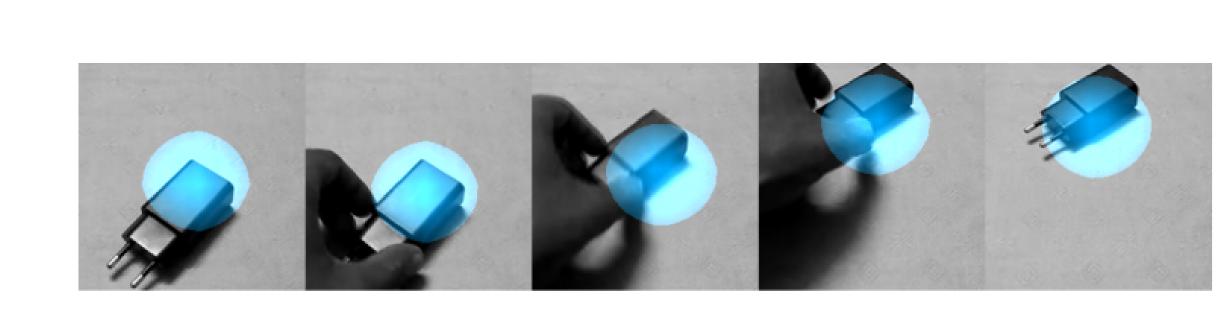
- generate location and size of N regions.
- create nodes by extracting features from these regions using bilinear interpolation.
- process the nodes with a spatio-temporal GNN and project each node into its initial location.

Results

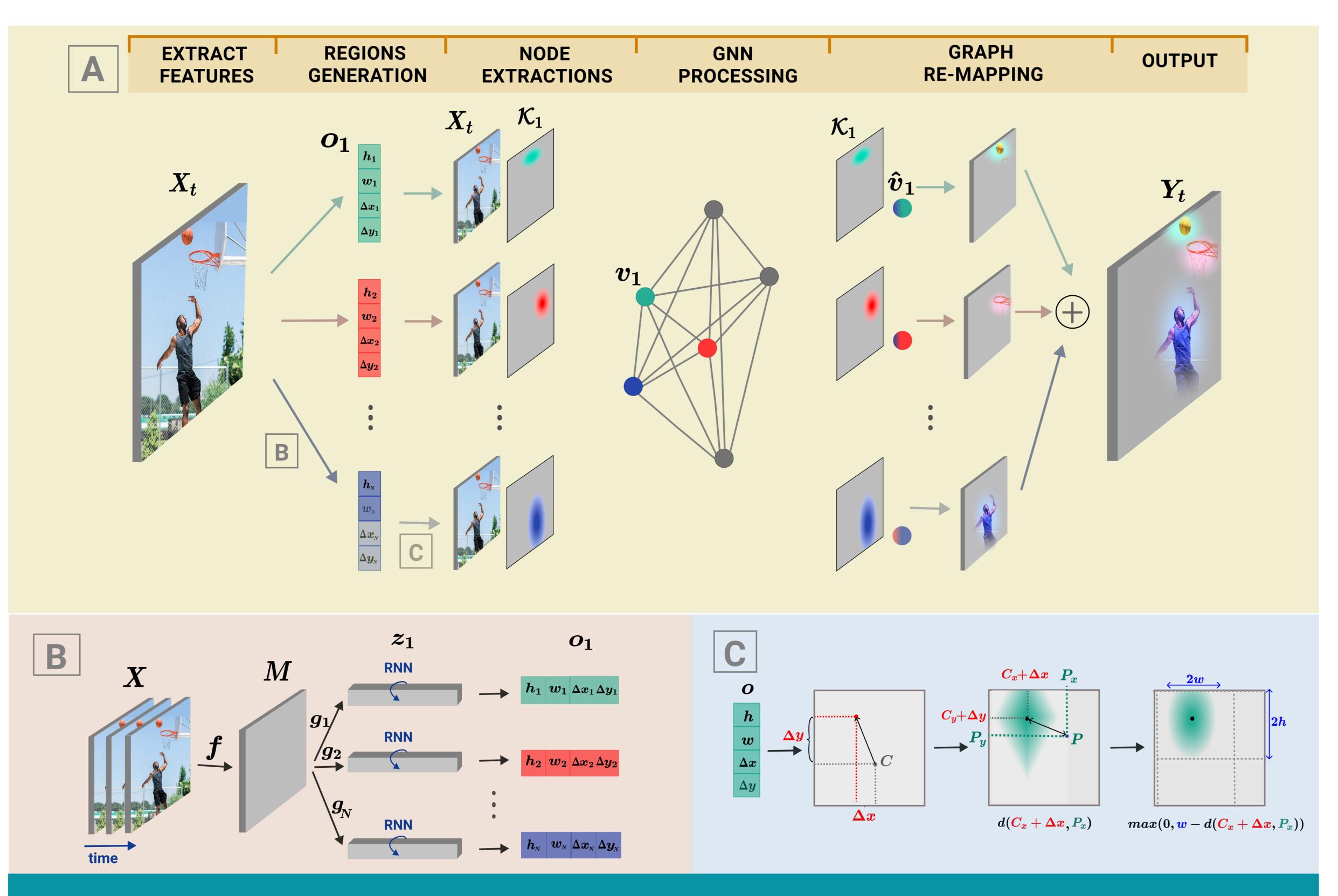
Something-Something is a real world dataset involving human-object interactions.

	Model	Top 1	Top 5
non-Graph	GST [1] TSM [2] SmallBig [3] STM [4]	62.6 63.4 63.8 64.2	87.9 88.5 88.9 89.8
Graph	TRG [5]	59.8	87.4
	DyReG - r4 DyReG - r3-4-5	64.3 64.8	88.9 89.4

Accuracy on Something-Something-V2 dataset.



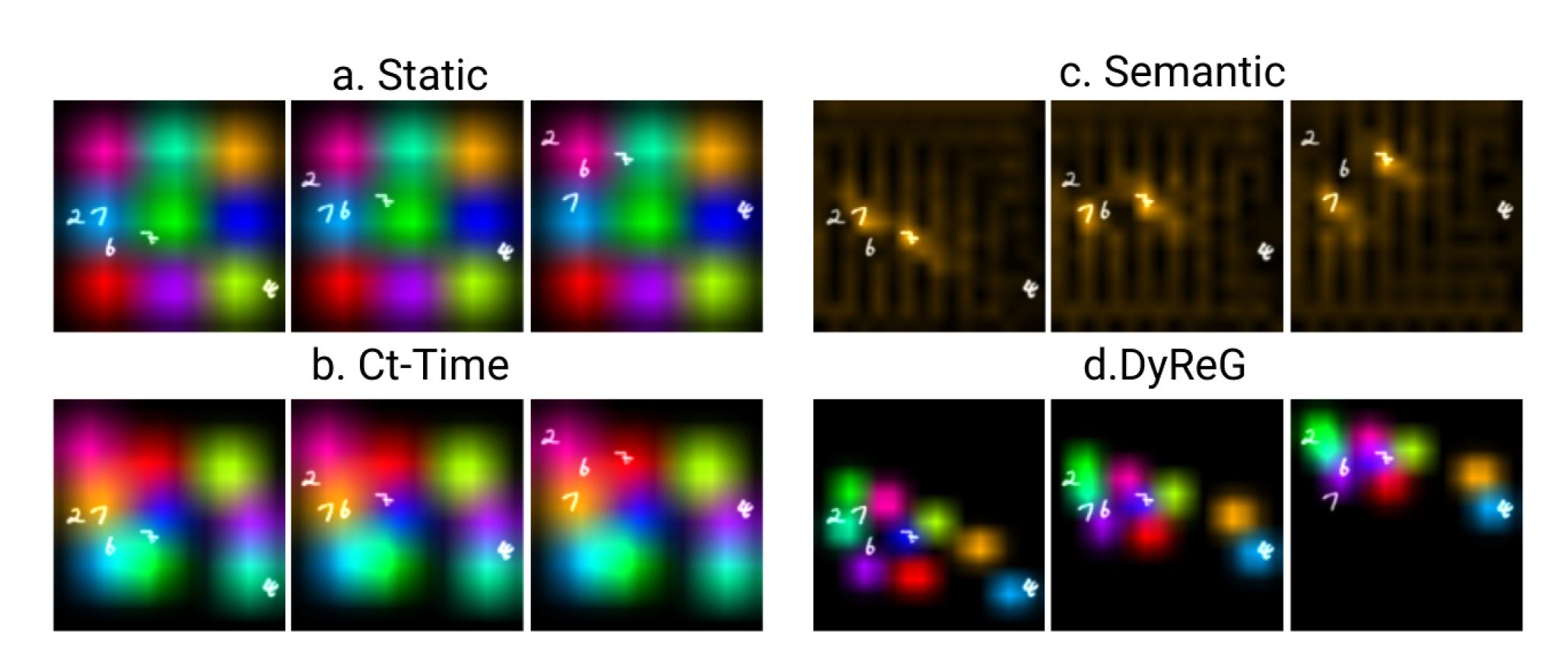
Visualisation of a single predicted kernel on Something-Something-V2.



Discover salient regions, biased towards instances, that are useful for visual relational processing.

Model	Accuracy	Description
Static Nodes Ct-Time Nodes	81.48 86.77	Optimize regions across dataset Keep regions fixed in time
Semantic Nodes	82.41	Attend to all the input positions
DyReG Nodes	95.09	Full model with dynamic regions

Ablation on MultiSyncMNIST of different types of node extraction.



The goal of MultiSyncMNIST is to detect a group of digits that move synchronously.

Node Region Generation

- Produce the most salient N = 9 regions using a global processing.
 - 1. Each node is modeled by a function with its own set of parameters.
 - 2. A recurrent function is used to achieve consistency across time.
 - 3. Produce the location and size of each region.

$$M_{t} = f(X_{t}) \in \mathbb{R}^{H' \times W' \times C'}$$

$$\mathbf{\hat{m}}_{i,t} = g_{i}(M_{t}) \in \mathbb{R}^{C'}, \forall i \in \overline{1, N}$$

$$\mathbf{z}_{i,t} = \text{GRU}(\mathbf{z}_{i,t-1}, \mathbf{\hat{m}}_{i,t}) \in \mathbb{R}^{C'}, \forall i \in \overline{1, N}$$

$$\mathbf{o}_{i,t} = (\Delta x_{i,t}, \Delta y_{i,t}, w_{i,t}, h_{i,t}) = \alpha \odot W \mathbf{z}_{i,t} \in \mathbb{R}^{4}$$

Node Features Extraction

- Learn to generate region's parameters using the video-level supervision.
 - Make the node feature extraction differentiable w.r.t. region's parameters.
- Extract node features using an interpolation kernel.
- The kernel decreases with the distance to the center and is non-zero up to a maximal distance of w_i .

$$\mathcal{K}^{(i)}(p_x, p_y) = k_x^{(i)}(p_x)k_y^{(i)}(p_y) \in \mathbb{R}$$
$$k_x^{(i)}(p_x) = \max(0, w_i - d(\Delta x_i, p_x))$$

Graph Processing

- Process the nodes with a spatio-temporal GNN similar to our previous work [6].
 - At each time step, send messages between nodes.
 - Across time, update each node independently using a RNN.

$$\mathbf{v}_{i,t} = \sum_{j=1}^{N} a(\mathbf{v}_{j,t}, \mathbf{v}_{i,t}) \text{MLP}(\mathbf{v}_{j,t}, \mathbf{v}_{i,t}) \in \mathbb{R}^{C}$$

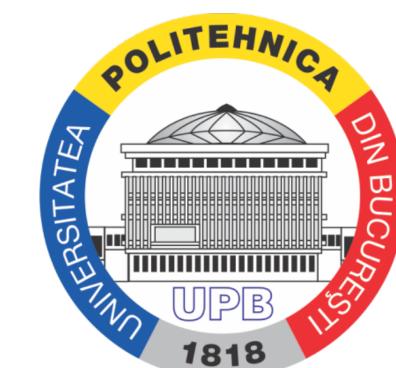
$$\mathbf{\hat{v}}_{i,t+1} = \text{GRU}(\mathbf{\hat{v}}_{i,t}, \mathbf{v}_{i,t}) \in \mathbb{R}^{C}$$

Graph Re-Mapping

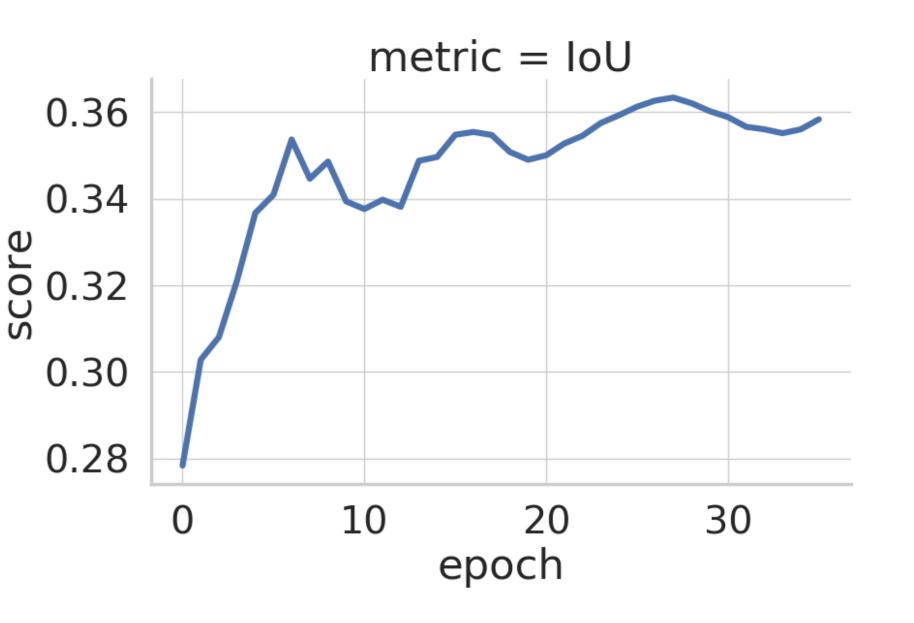
• The features of each updated node are sent to their initial region in the input, as indicated by their corresponding kernel.

$$\mathbf{y}_{p_x,p_y,t} = \sum_{i=1}^{N} \mathcal{K}_t^{(i)}(p_x,p_y) \mathbf{\hat{v}}_{i,t} \in \mathbb{R}^C$$





Object-centric representation



IoU between our DyReG regions and groundtruth boxes on Something-Something-V2.

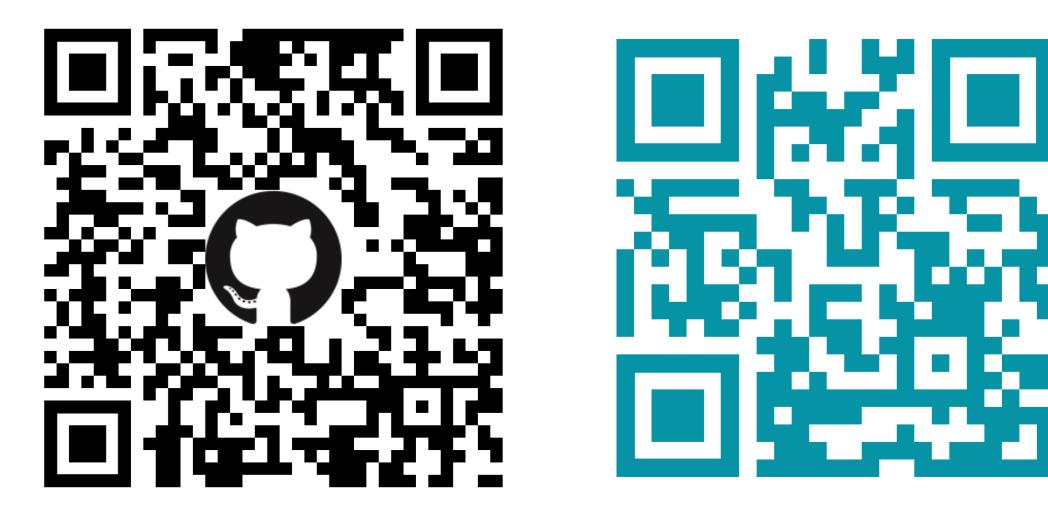
Model	IoU
Grid baseline	26.8
Center baseline	27.8
Detected boxes	47.1
DyReG	39.0

IoU for our best DyReG model and other baselines on Something-Something-V2.

Compute IoU between ground-truth objects and predicted node regions.

- we do not enforce this metric in any way and learn without object-level supervision.
- score increases during training showing that regions correlate with objects.

Code and team homepage:





References

- [1] Luo and Yuille ICCV 2019,
- [2] Lin et al. ICCV 2019,
- [3] Li et al. CVPR 2020,
- [4] Jiang et al. ICCV 2019,
- [5] Zhang et al. TIP 2020,
- [6] Nicolicioiu et al. NeurIPS 2019