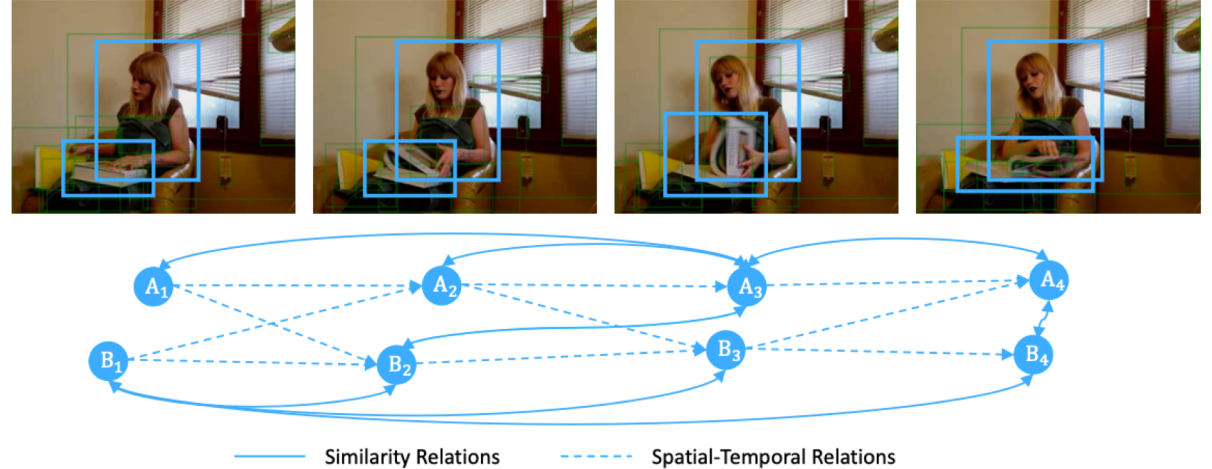


Videos as Space-Time Region Graphs

Jaedong Hwang

Computer Vision Lab.

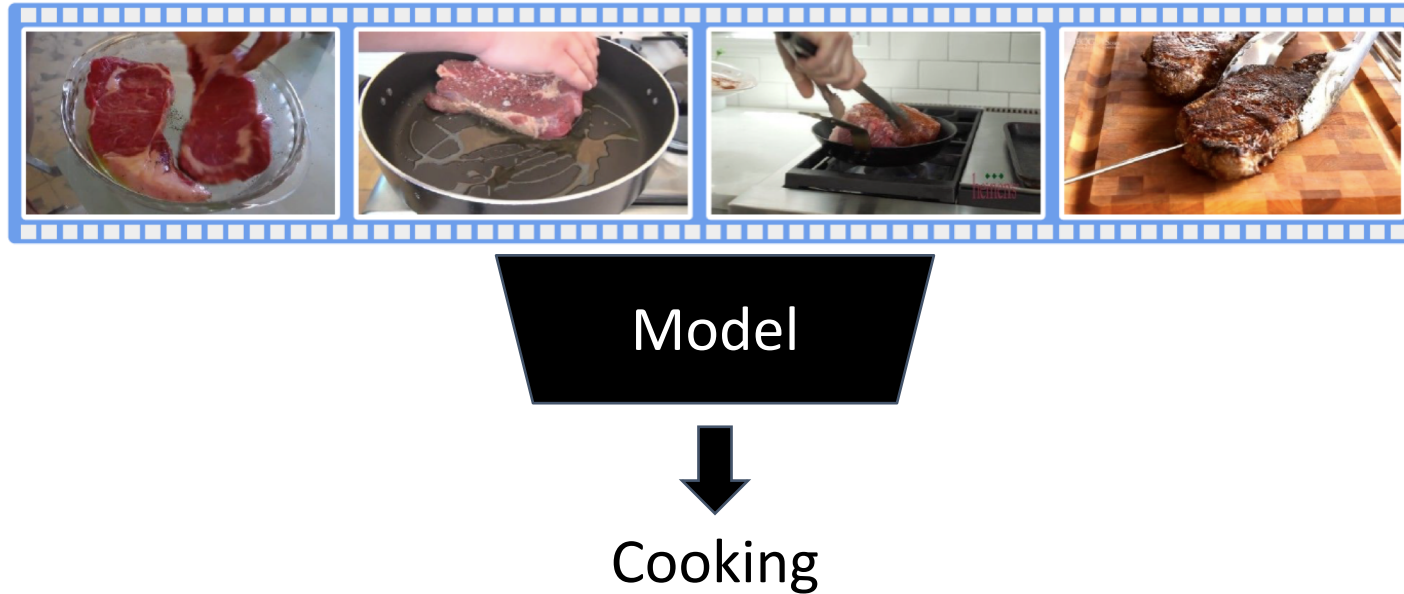
Seoul National University



- TL;DR: Generate and utilize graphs of region proposals to capture relationship between objects.

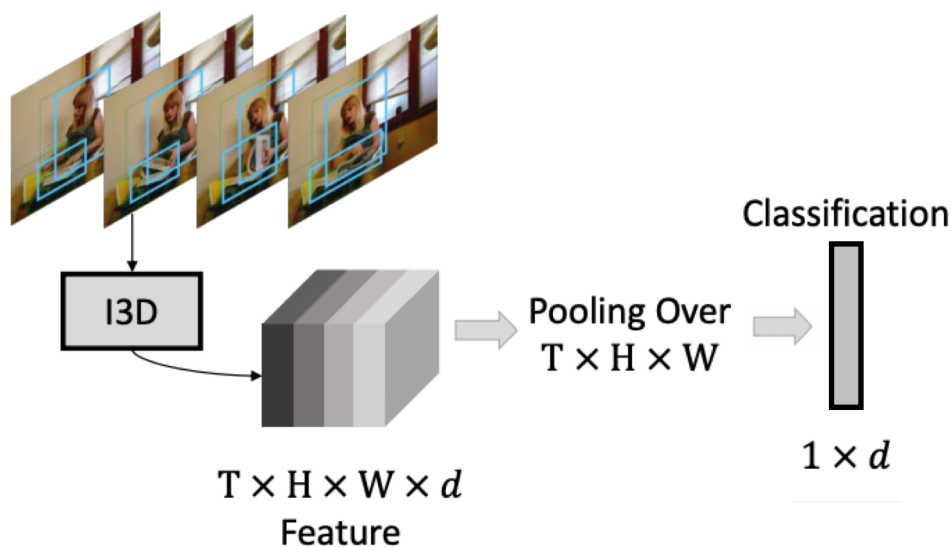
What is Action Recognition?

- Classify which actions in a videos.

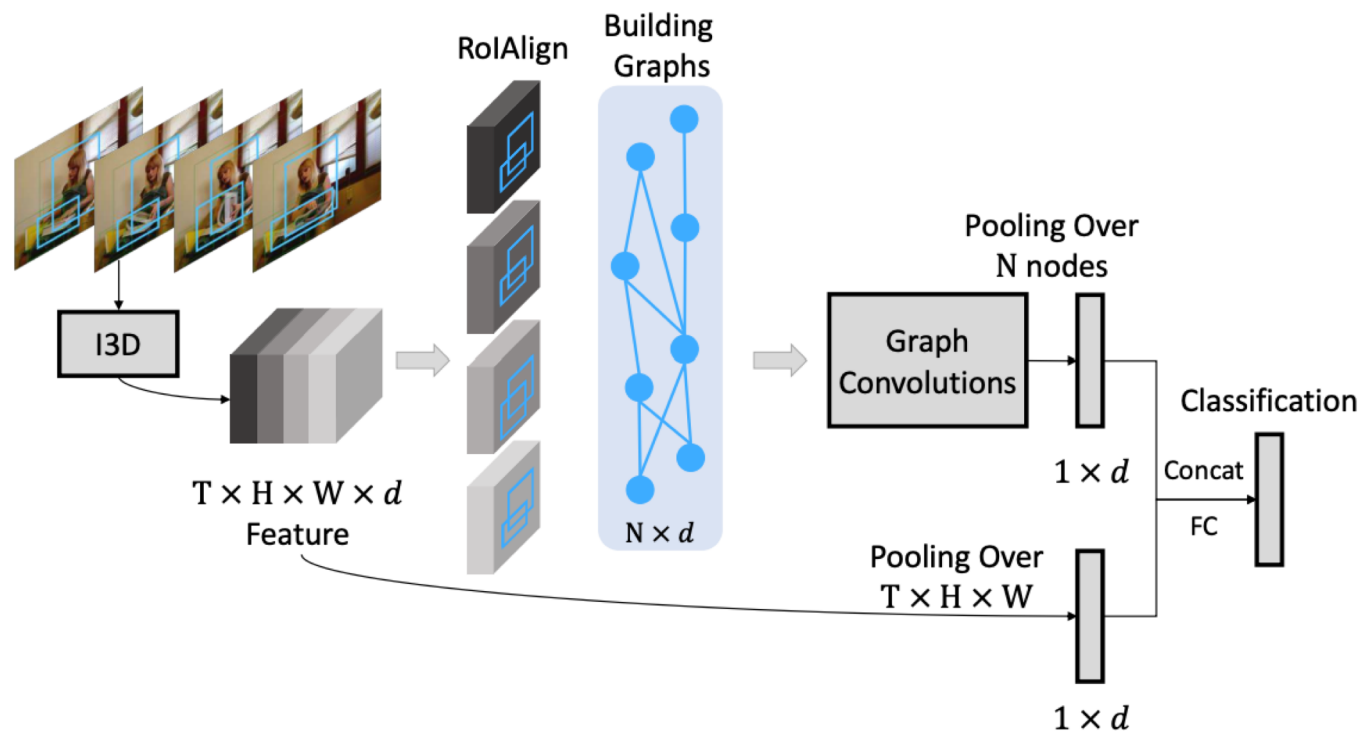


- Depending on temporal direction, label can be changed
 - Door open vs. Door close

Videos as Space-Time Region Graphs

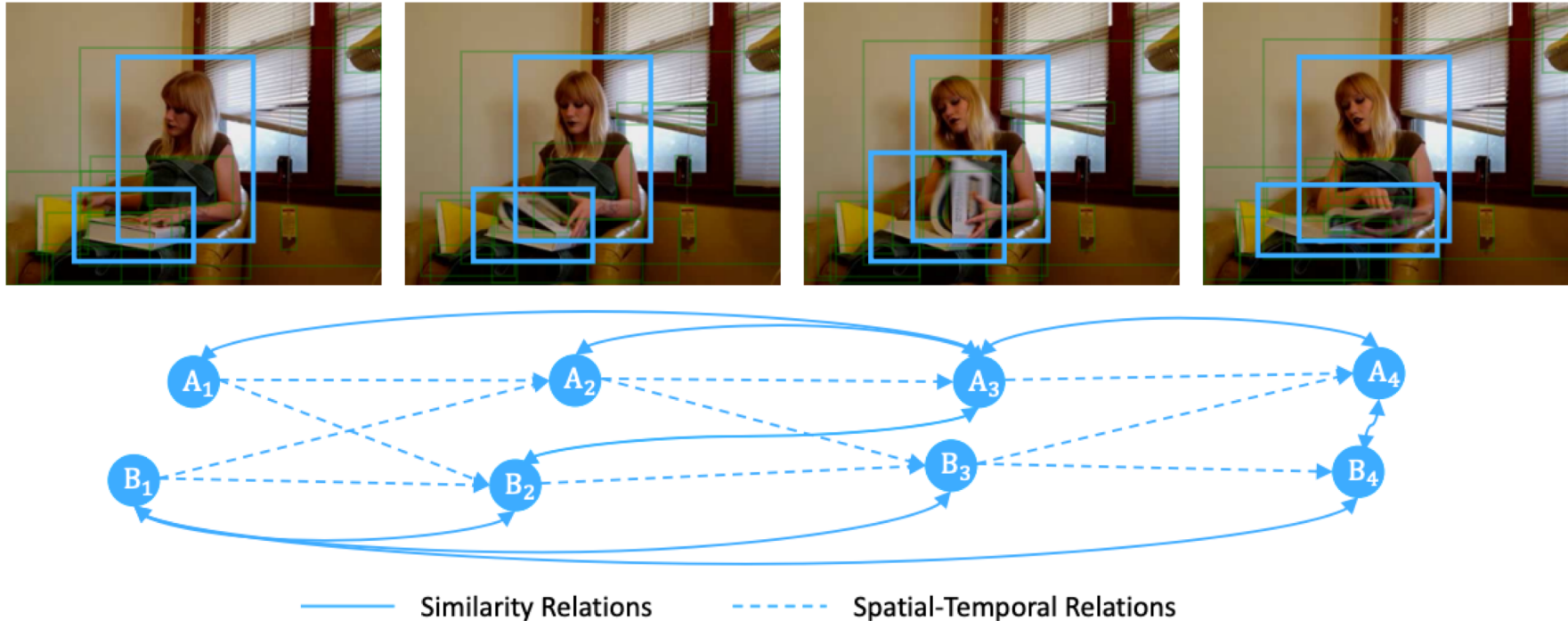


Previous architecture



Proposed architecture

Two Graphs for region proposals



- Similarity graph – fully connected, similarity
- Spatiotemporal graph – partially connected, overlap

Building Graphs – Similarity Graph

- Fully connected directed graph.
- Transition matrix, \mathbf{G}_{ij}^{sim} is defined as correlation btw embedded features.

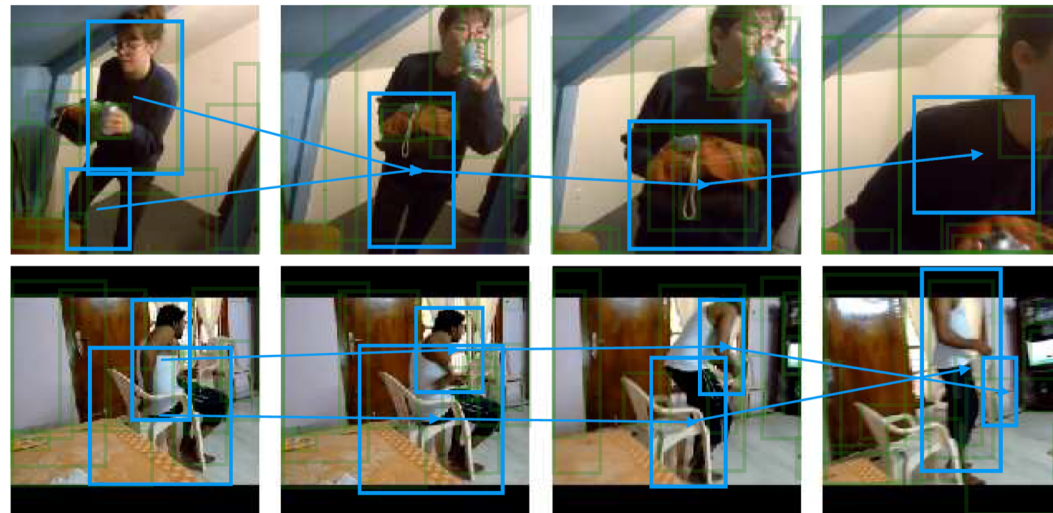
$$F(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi'(\mathbf{x}_j)$$

where $\phi(\mathbf{x}) = \mathbf{W}\mathbf{x}$, $\phi'(\mathbf{x}) = \mathbf{W}'\mathbf{x}$

$$\mathbf{G}_{ij}^{sim} = \frac{\exp F(\mathbf{x}_i, \mathbf{x}_j)}{\sum_{j=1}^N \exp F(\mathbf{x}_i, \mathbf{x}_j)}$$

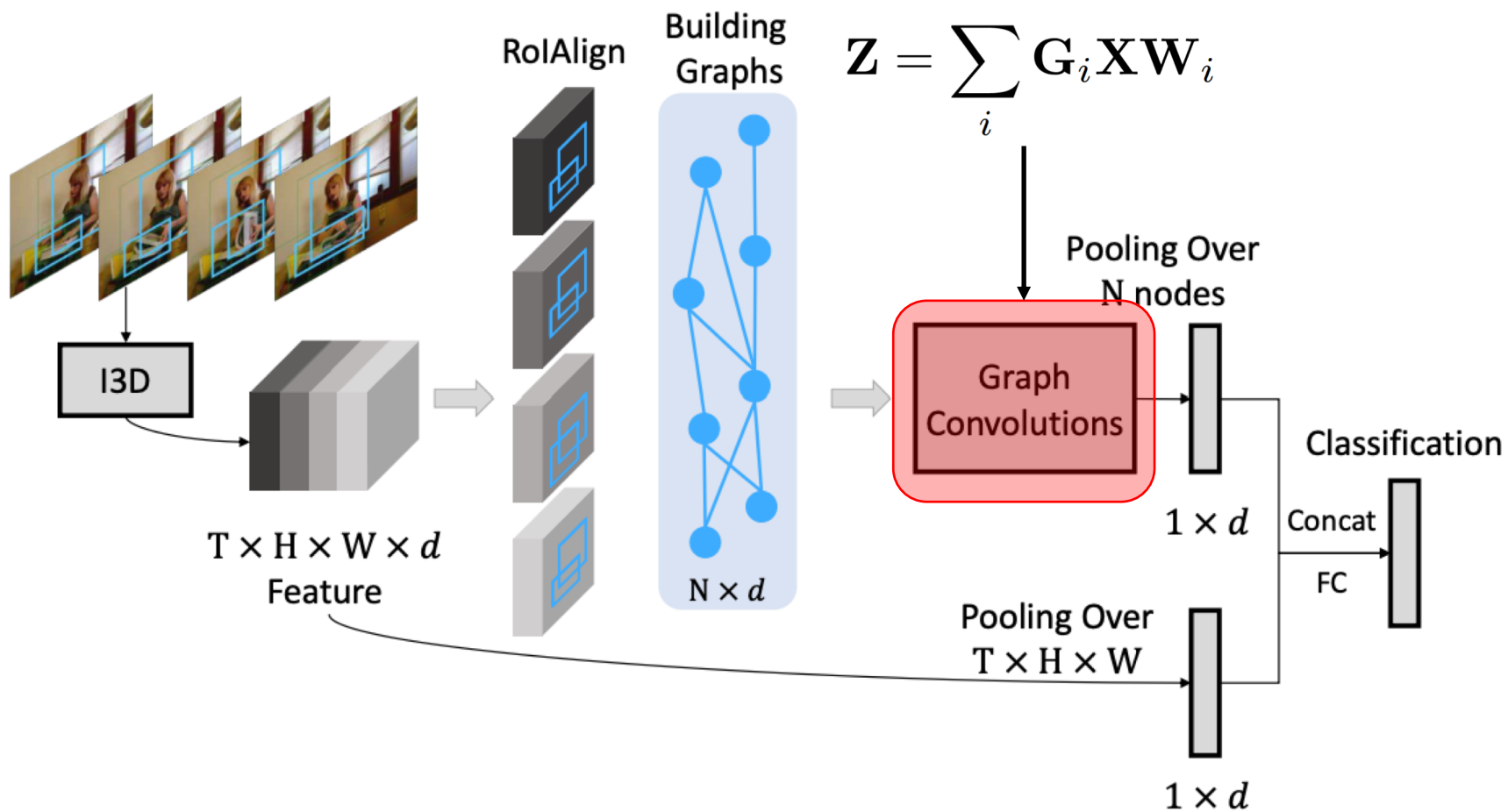
Building Graphs – Spatiotemporal graph

- Use spatio-temporal relation between region proposals.
- Two directed graph (forward graph, backward graph).
- If a proposal at t has an overlap with a proposal at $t + 1$, connect them (vice versa in backward graph).



$$G_{ij}^{front} = \frac{\sigma_{ij}}{\sum_{j=1}^N \sigma_{ij}},$$

σ_{ij} is an IoU of \mathbf{X}_i and \mathbf{X}_j



Result

- Dataset
 - Charades: 157 classes, multi-action, 8k train, 1.8k validation.
 - Something-Something: 174 classes, single-action 86k train, 12k validation, 11k test.

model, R50, I3D	mAP
baseline	31.8
Proposal+AvgPool	32.1
Spatial-Temporal GCN	34.2
Similarity GCN	35.0
Joint GCN	36.2



Conclusion

- first uses a Graph Convolutional Network for reasoning with multiple relation edges on video action recognition.
- presents a novel graph representation with variant relationships between objects in a long range video.