



CVPR JUNE 19-24 2022 NEW ORLEANS LOUISIANA

DirecFormer: A Directed Attention in Transformer Approach to Robust Action Recognition

Thanh-Dat Truong¹, Quoc-Huy Bui², Chi Nhan Duong³,
Han-Seok Seo⁴, Son Lam Phung⁵, Xin Li⁶, Khoa Luu¹

¹CVIU Lab, University of Arkansas ²NextG, FPT Software

³Concordia University ⁴Dep. of Food Science, University of Arkansas

⁵University of Wollongong ⁶West Virginia University

This material is based upon work supported by the National Science Foundation under [Award No. QIA-1946391](#). Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.





<https://vizle.offnote.co>

Contact us: vizle@offnote.co

This document was generated automatically by **Vizle**

Your **Personal Video Reader Assistant**

Learn from Videos **Faster** and **Smarter**

VIZLE **PRO / BIZ**

PDF, PPT ~~Watermarks~~

- Convert *entire* videos
- *Customize* to retain all essential content
- Include Spoken *Transcripts*
- Customer support

Visit <https://vizle.offnote.co/pricing> to learn more

VIZLE **FREE PLAN**

PDF only ~~Watermarks~~

- Convert videos *partially*
- Slides may be *skipped**
- Usage restrictions
- No Customer support

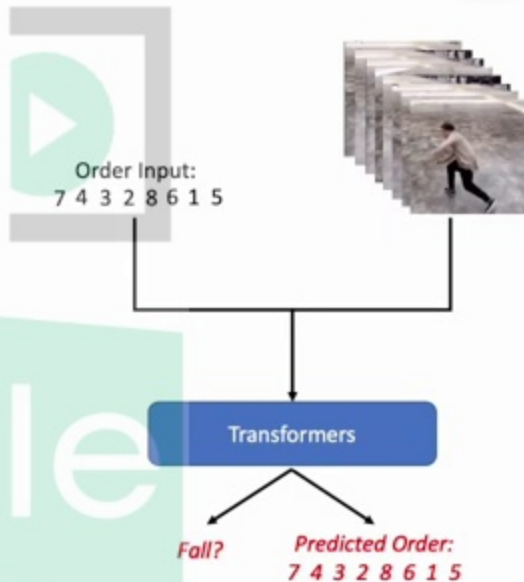
Visit <https://vizle.offnote.co> to try free

Login to Vizle to unlock more slides*



Vizle Motivation

- The goal of action recognition is to predict an action given a set of *video frames in the correct order*.
- Given a set of *video frames shuffled in a random order* and different from the original one, will it be classified as *the same label as the original recognition result*?
- Are action recognition models able to *correct the incorrectly-ordered frames to the right ones* and provide *an accurate prediction*?

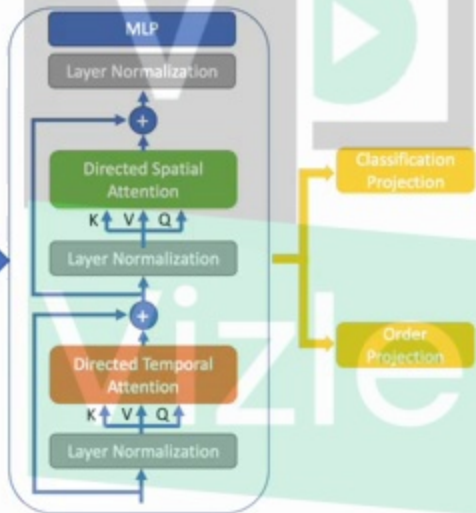




The Proposed Approach

Input
Video

Patch
Representation

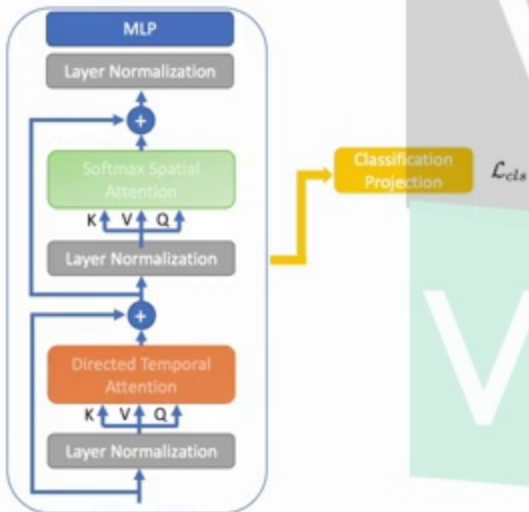




Ablation Study

Table 1. Ablation Study On **Jester**. $X - Y$ denotes for the attention types of temporal and spatial dimension, respectively. X (and Y) could be either S : Softmax or C : Cosine.

Models	Attention Time-Space	\mathcal{L}_{ord}	\mathcal{L}_{self}	Top 1	Top 5
DirecFormer	$S - C$			94.52	99.26
DirecFormer	$S - C$	✓		94.65	99.25
DirecFormer	$C - S$			95.52	99.20






Experimental Results

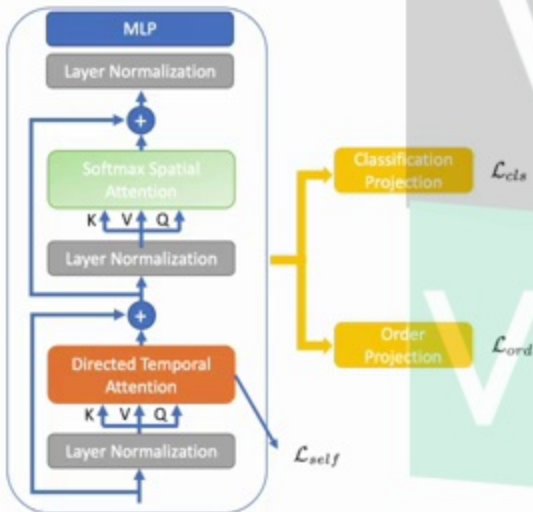


Table 4. **Comparison with the SOTA methods on Kinetics 400.** $X - Y$ denotes for the attention types of temporal and spatial dimension, respectively. X (and Y) could be either S : Softmax or C : Cosine.

Models	Attention Time-Space	Top 1	Top 5
DiracFormer	$S - C$	80.16	94.55
DiracFormer	$C - S$	81.69	94.62



<https://vizle.offnote.co>

Contact us: vizle@offnote.co

This document was generated automatically by **Vizle**

Your **Personal Video Reader Assistant**

Learn from Videos **Faster** and **Smarter**

VIZLE **PRO / BIZ**

PDF, PPT ~~Watermarks~~

- Convert *entire* videos
- *Customize* to retain all essential content
- Include Spoken *Transcripts*
- Customer support

Visit <https://vizle.offnote.co/pricing> to learn more

VIZLE **FREE PLAN**

PDF only ~~Watermarks~~

- Convert videos *partially*
- Slides may be *skipped**
- Usage restrictions
- No Customer support

Visit <https://vizle.offnote.co> to try free

Login to Vizle to unlock more slides*