# TubeDETR: Spatio-Temporal Video Grounding with Transformers
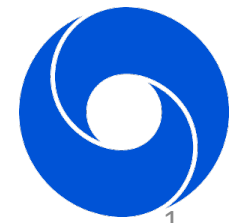
Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, Cordelia Schmid

Project page: https://antoyang.github.io/tubedetr.html

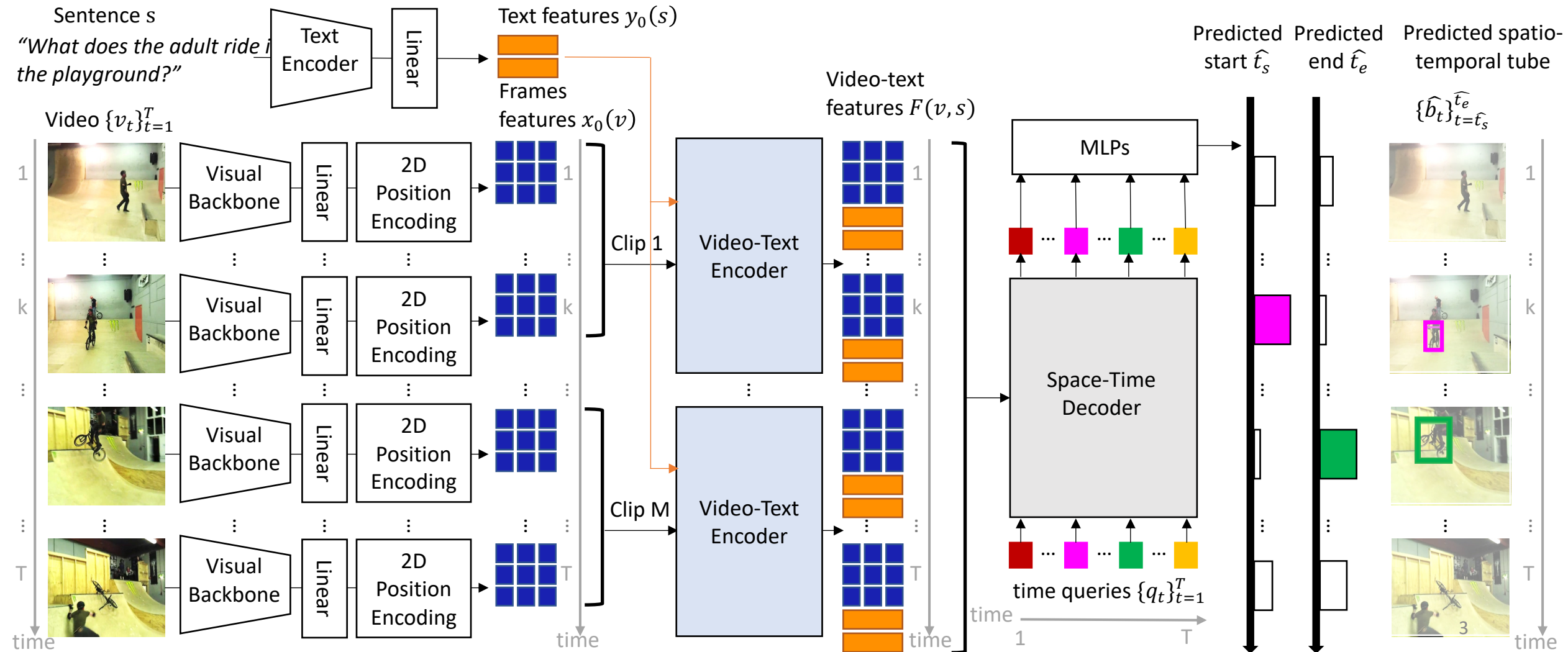Paper: https://arxiv.org/abs/2203.16434
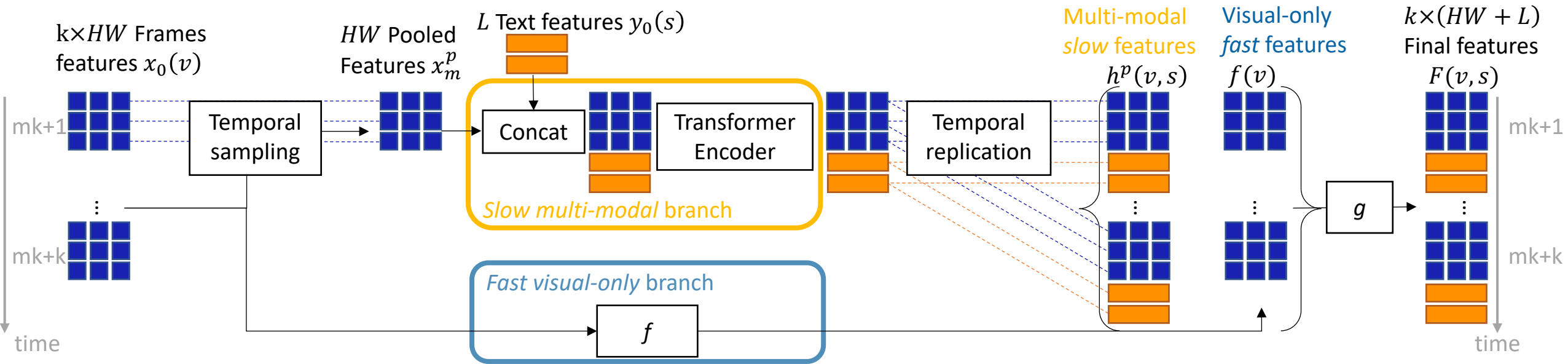
# Spatio-Temporal Video Grounding

- **Input text query:** What does the adult ride in the playground?
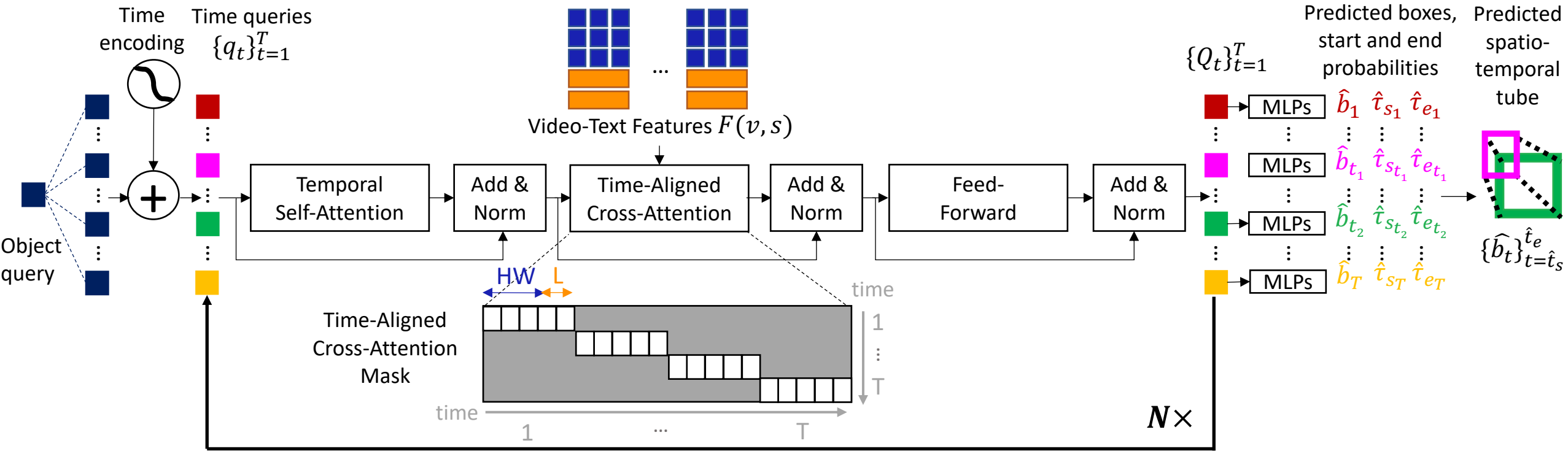- **Output spatio-temporal tube:**

# TubeDETR Architecture Overview

# Video-Text Encoder

# Space-Time Decoder

# Training

- **Loss:** Combination of spatial localization ($\mathcal{L}_1$ , $gIoU$) and temporal localization ($KL$ , $att$) objectives

$$\mathcal{L} = \lambda_{\mathcal{L}_1} \mathcal{L}_{\mathcal{L}_1}(\hat{b}, b) + \lambda_{gIoU} \mathcal{L}_{gIoU}(\hat{b}, b) + \lambda_{KL} \mathcal{L}_{KL}(\hat{\tau}_s, \hat{\tau}_e, \tau_s, \tau_e) + \lambda_{att} \mathcal{L}_{att}(A)$$

$\lambda_{\bullet}$ : scalar weights of the individual losses

$\hat{b}$ and $b$ : predicted and ground truth boxes

$\hat{\tau}_s$ and $\tau_s$ : predicted and ground truth start probability distribution

$\hat{\tau}_e$ and $\tau_e$ : predicted and ground truth end probability distribution

$A$ : temporal self-attention matrix

- **Initialization:** from MDETR weights pretrained on Visual Genome, COCO and Flickr

# Ablations: Space-Time Decoder

**Temporal** · **Video** · **Spatial**

**IoU** · **IoU** · **IoU**

| | Time Encoding | Self Attention | m_tIoU | m_vIoU | vIoU @0.3 | vIoU @0.5 | m_sIoU |
|---|---|---|---|---|---|---|---|
| 1. | ✗ | - | 23.9 | 12.2 | 15.3 | 6.1 | 47.0 |
| 2. | ✗ | Temporal | 25.2 | 13.0 | 16.9 | 6.5 | 47.3 |
| 3. | ✓ | - | 41.7 | 21.3 | 28.7 | 17.4 | 46.5 |
| 4. | ✓ | Temporal | **45.9** | **24.3** | **33.2** | **22.0** | **47.7** |

Time encoding matters.

Temporal self-attention helps.

Table 1. Effect of the time encoding and the temporal self-attention in our space-time decoder on the VidSTG validation set.

# Ablations: Weights initialization

Temporal    Video    Spatial

IoU    IoU    IoU

| | Pre-Training | Decoder Self-Attention Transfer | m_tIoU | m_vIoU | vIoU @0.3 | vIoU @0.5 | m_sIoU |
|---|---|---|---|---|---|---|---|
| 1. | ✗ | ✗ | 42.8 | 18.8 | 25.1 | 15.6 | 38.5 |
| 2. | ✓ | ✗ | 43.8 | 22.4 | 29.9 | 19.1 | 46.5 |
| 3. | ✓ | Temporal | **45.9** | **24.3** | **33.2** | **22.0** | **47.7** |

Table 2. Effect of the weight initialization for our model on the VidSTG validation set.

MDETR pretraining matters.

Transferring spatial self-attention to temporal self-attention helps.

# Ablations: Video-Text Encoder

- Our encoder is memory-efficient.
- Fast branch matters.

<table>
<tr><td colspan="8">(a) VidSTG</td><td colspan="8">(b) HC-STVG2.0</td></tr>
<tr><td>Fast</td><td>Res.</td><td>Temp. Stride</td><td>m_tIoU</td><td>m_vIoU</td><td>vIoU@0.3</td><td>vIoU@0.5</td><td>m_sIoU</td><td>Mem. (GB)</td><td>Fast</td><td>Res.</td><td>Temp. Stride</td><td>m_tIoU</td><td>m_vIoU</td><td>vIoU@0.3</td><td>vIoU@0.5</td><td>m_sIoU</td><td>Mem. (GB)</td></tr>
<tr><td>1. —</td><td>224</td><td>1</td><td>46.5</td><td>25.2</td><td>34.1</td><td>23.0</td><td>49.1</td><td>23.9</td><td>1. —</td><td>224</td><td>1</td><td>52.8</td><td>35.0</td><td>55.3</td><td>28.3</td><td>63.9</td><td>14.3</td></tr>
<tr><td>2. ✓</td><td>224</td><td>2</td><td>46.0</td><td>25.0</td><td>34.3</td><td>22.9</td><td>49.0</td><td>16.2</td><td>2. ✓</td><td>224</td><td>2</td><td>53.7</td><td>35.8</td><td>56.7</td><td>29.6</td><td>64.3</td><td>10.2</td></tr>
<tr><td>3. ✓</td><td>224</td><td>5</td><td>45.9</td><td>24.3</td><td>33.2</td><td>22.0</td><td>47.7</td><td>11.8</td><td>3. ✓</td><td>224</td><td>5</td><td>53.2</td><td>35.0</td><td>54.5</td><td>29.0</td><td>63.2</td><td>8.0</td></tr>
<tr><td>4. ✓</td><td>288</td><td>2</td><td>46.4</td><td>25.9</td><td>35.0</td><td>23.9</td><td>50.5</td><td>23.7</td><td>4. ✓</td><td>288</td><td>2</td><td>**53.9**</td><td>**36.4**</td><td>58.1</td><td>**30.7**</td><td>**65.4**</td><td>13.9</td></tr>
<tr><td>5. ✓</td><td>320</td><td>3</td><td>46.4</td><td>25.9</td><td>35.7</td><td>23.7</td><td>**50.7**</td><td>23.6</td><td>5. ✓</td><td>320</td><td>3</td><td>53.6</td><td>36.2</td><td>57.5</td><td>30.4</td><td>65.2</td><td>13.8</td></tr>
<tr><td>6. ✓</td><td>352</td><td>4</td><td>**46.9**</td><td>**26.2**</td><td>**36.1**</td><td>**24.1**</td><td>**50.7**</td><td>24.4</td><td>6. ✓</td><td>352</td><td>4</td><td>**53.9**</td><td>**36.4**</td><td>**58.8**</td><td>30.6</td><td>64.9</td><td>14.3</td></tr>
<tr><td>7. ✗</td><td>352</td><td>4</td><td>46.6</td><td>24.8</td><td>34.0</td><td>21.6</td><td>48.3</td><td>18.1</td><td>7. ✗</td><td>352</td><td>4</td><td>53.1</td><td>34.7</td><td>55.9</td><td>27.4</td><td>63.0</td><td>11.3</td></tr>
<tr><td>8. ✓</td><td>384</td><td>5</td><td>46.8</td><td>26.0</td><td>35.5</td><td>24.0</td><td>50.4</td><td>26.1</td><td>8. ✓</td><td>384</td><td>5</td><td>53.6</td><td>36.3</td><td>57.5</td><td>30.4</td><td>65.3</td><td>15.2</td></tr>
</table>

Table 3. Comparison of performance-memory trade-off with various temporal strides $k$, spatial resolutions (Res.), with or without the fast branch in our video-text encoder, on the VidSTG validation set (left, Table 3a) and the HC-STVG2.0 validation set (right, Table 3b).

# Comparison with state of the art

- State-of-the-art results on: VidSTG and HC-STVG.

| Method | Pretraining Data | VidSTG | | | | | | | | HC-STVG1 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Declarative Sentences | | | | Interrogative Sentences | | | | | | |
| | | m_tIoU | m_vIoU | vIoU@0.3 | vIoU@0.5 | m_tIoU | m_vIoU | vIoU@0.3 | vIoU@0.5 | m_vIoU | vIoU@0.3 | vIoU@0.5 |
| 1. STGRN [102] | Visual Genome | **48.5** | 19.8 | 25.8 | 14.6 | **47.0** | 18.3 | 21.1 | 12.8 | — | — | — |
| 2. STGVT [72] | Visual Genome + Conceptual Captions | — | 21.6 | 29.8 | 18.9 | — | — | — | — | 18.2 | 26.8 | 9.5 |
| 3. STVGBert [68] | ImageNet + Visual Genome + Conceptual Captions | — | 24.0 | 30.9 | 18.4 | — | 22.5 | 26.0 | 16.0 | 20.4 | 29.4 | 11.3 |
| 4. TubeDETR (Ours) | ImageNet | 43.1 | 22.0 | 29.7 | 18.1 | 42.3 | 19.6 | 26.1 | 14.9 | 21.2 | 31.6 | 12.2 |
| 5. TubeDETR (Ours) | ImageNet + Visual Genome + Flickr + COCO | 48.1 | **30.4** | **42.5** | **28.2** | 46.9 | **25.7** | **35.7** | **23.2** | 32.4 | 49.8 | 23.5 |

Table 4. Comparison to the state of the art on the VidSTG test set and the HC-STVG1 test set.

# Qualitative results

- **Interactive Demo:** http://stvg.paris.inria.fr/

- **Query:** What is beneath the adult in the snow?