# Attend and Interact: Higher-Order Object Interactions for Video Understanding

Chih-Yao Ma[1], Asim Kadav[2], Iain Melvin[2], Zsolt Kira[3], Ghassan AlRegib[1], and Hans Peter Graf[2]

[1]Georgia Institute of Technology, [2]NEC Laboratories America, [3]Georgia Tech Research Institute

CVPR 2018 · SALT LAKE CITY · JUNE 18-22
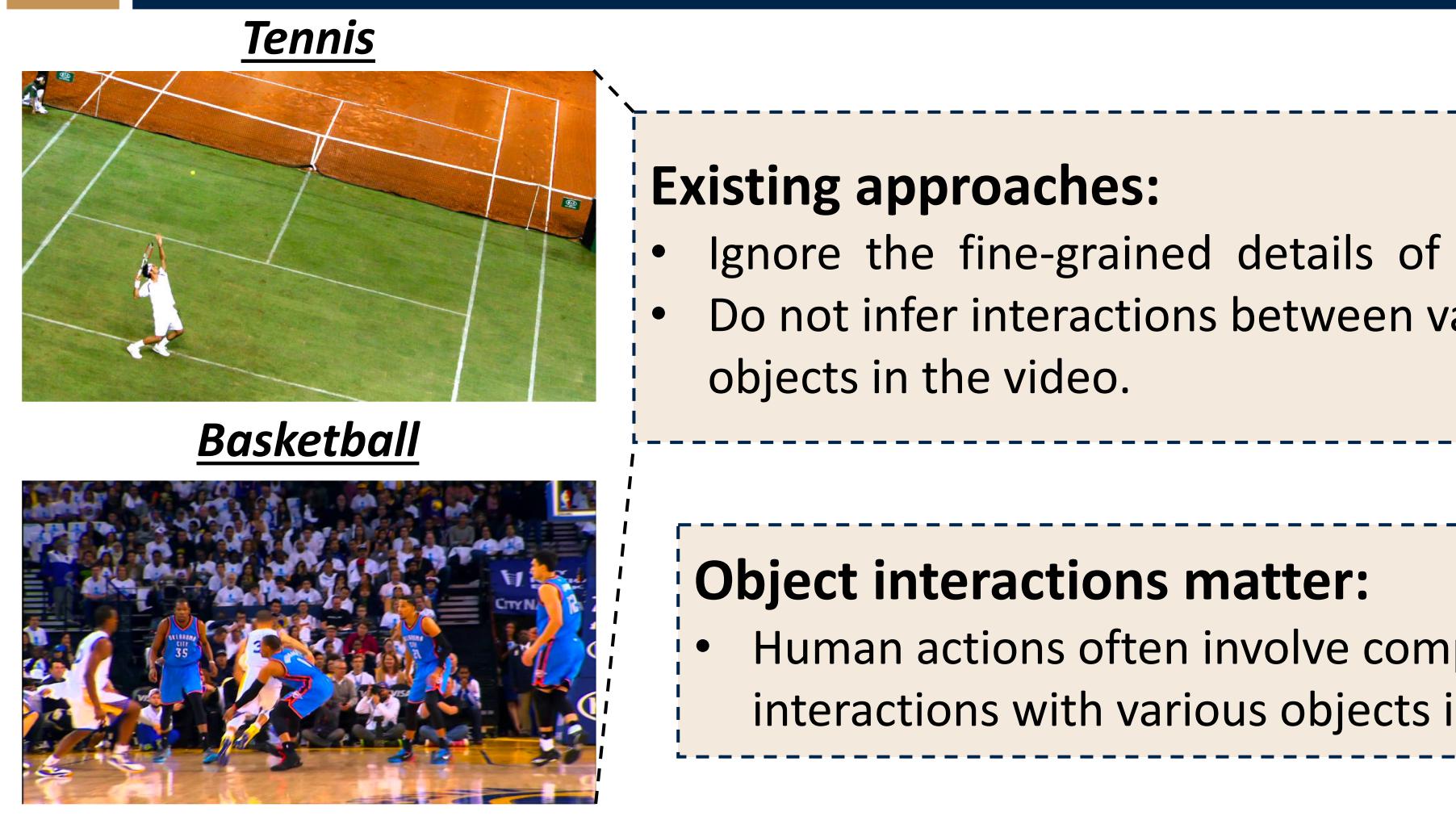
## 1 BACKGROUND: VIDEO UNDERSTANDING

*Tennis* *Basketball* *Skiing* *Snowboarding*
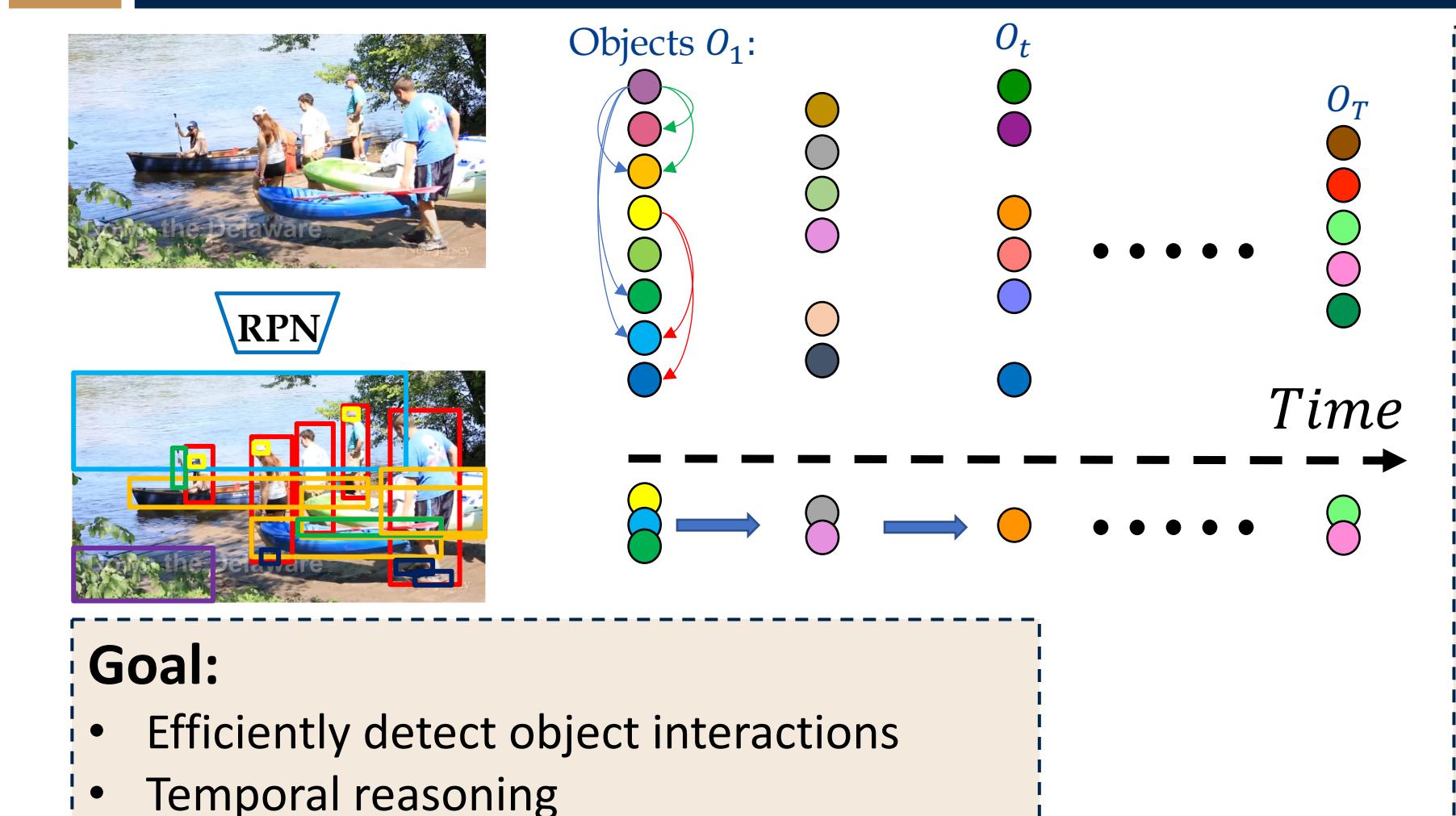
**Existing approaches:**
- Ignore the fine-grained details of the scene.
- Do not infer interactions between various objects in the video.

**Object interactions matter:**
- Human actions often involve complex interactions with various objects in the scene.
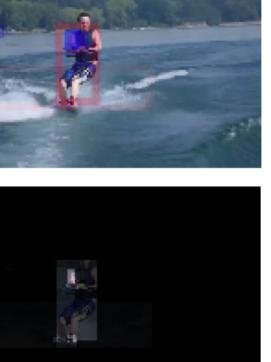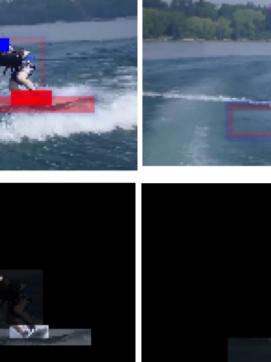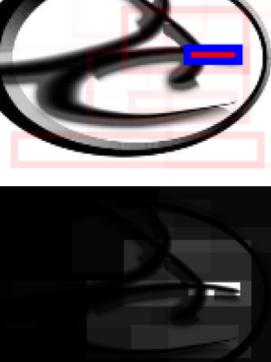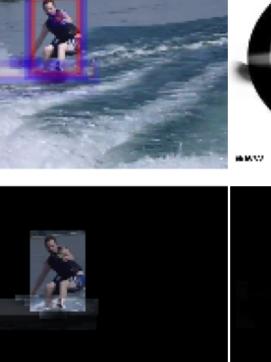
## 2 FINE-GRAINED OBJECT INTERACTIONS

Objects $o_1$ ... $o_t$ ... $o_T$

RPN

Time

**Challenges:**
- Unordered variable-lengths of object sets that span across time
- Video contains hundreds or thousands of frames
- Object-object pairs are too large to fully represented by a finite-capacity neural network

**Goal:**
- Efficiently detect object interactions
- Temporal reasoning

**Higher-order interactions:**
Object Interactions are **not** always between two objects (one-to-one)

## 3 CONTRIBUTION: FROM PAIRWISE TO HIGHER-ORDER INTERACTIONS

**Interactions/relationships:**
$$RN(O) = f_\phi \left( \sum_{i,j} f_\theta(o_i, o_j) \right)$$

Concatenation

**Higher-order interactions:**
- Interactions over groups of inter-related objects
- Covers pair-wise or triplet object relationships as a special case

Dot-product

**Goal:**
- Detect inter-object relationships
- Objects with significant relationships are selected
- Groups of selected object relationships are concatenated.

Higher-Order Interaction

Concatenation [1]:
$$f_\theta(o_i, o_j) = W_\theta^T(o_i \| o_j)$$

Dot-product:
$$f_\theta(o_i, o_j) = \theta(o_i)^T \phi(o_j)$$
$$\rightarrow O^T W_\theta^T W_\phi O$$

inter-relationship

[1] Santoro, Adam, et al. "A simple neural network module for relational reasoning." *NIPS* 2017.

## 4 ACTION RECOGNITION – *SINet*

**Coarse-grained**

Image context

Time

ConvNet $v_{c,1}$ ... $v_{c,T}$ → MLP $g_\phi$ → SDP Attention → $v_c$

**Fine-grained**

RPN

Timestep $t$

Objects (ROIs) $o_{1,t}$, $o_{2,t}$ ... $o_{N,t}$ → $O_1$ ... $O_t$ ... $O_T$ → Recurrent HOI → $h_T = v_{oi,T}$ → $p(y)$

**Coarse-grained:**
- Video frames are encoded via a ConvNet
- Temporal pooling via SDP-Attention

**Fine-grained:**
- Objects (ROIs) are obtained from a RPN
- Progressively detect higher-order interactions via the HOI module

## 5 RECURRENT HIGHER-ORDER INTERACTION (HOI)

Objects $O_t$: $o_{1,t}$ ... $o_{N,t}$   Image context $v_{c,t}$

$K = 3$

MLP $g_{\theta_1}$ / MLP $g_{\theta_2}$ / MLP $g_{\theta_3}$

Attentive Selection $\alpha_1$ → $v_{o,t}^1$, $\alpha_2$ → $v_{o,t}^2$, $\alpha_3$ → $v_{o,t}^3$

$h_{t-1}$ → LSTM Cell → $h_t$

$repeat(W_{h_k}h_{t-1} + W_{c_k}v_{c,t}) + g_{\theta_k}(O_t)$
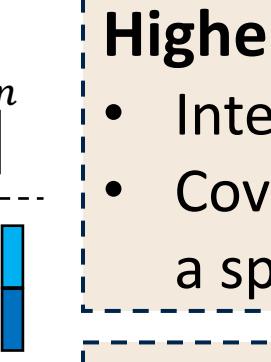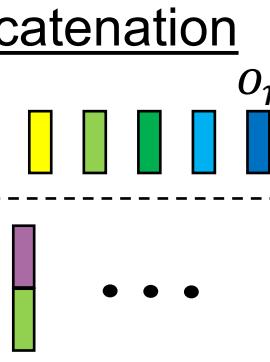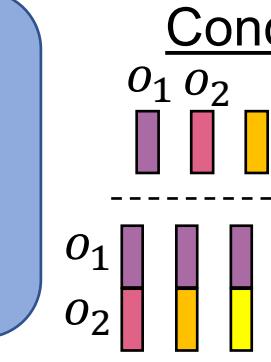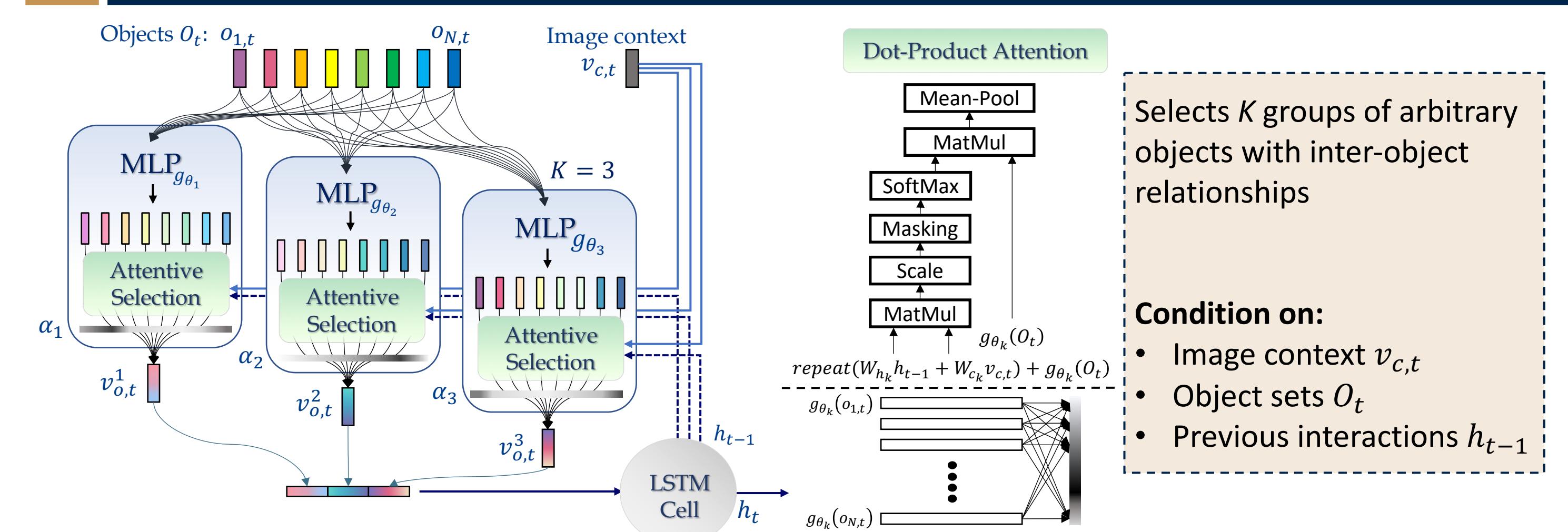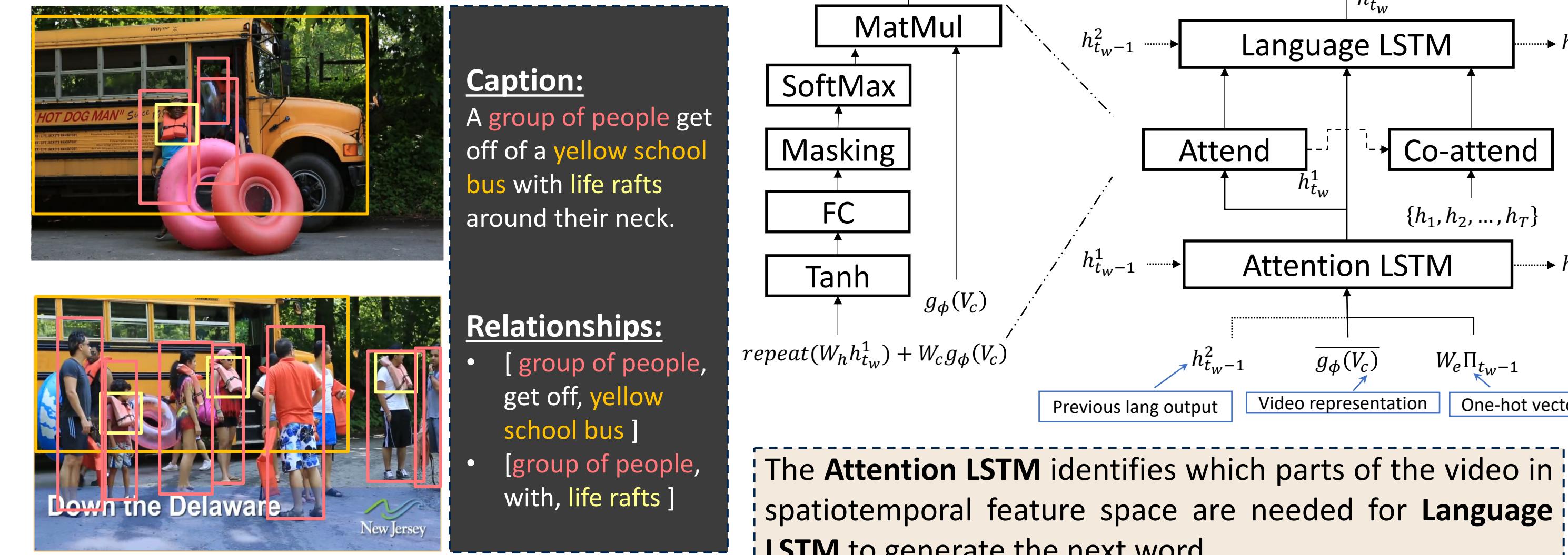
$g_{\theta_k}(o_{1,t})$ ... $g_{\theta_k}(o_{N,t})$

**Dot-Product Attention**

Mean-Pool / MatMul / SoftMax / Masking / Scale / MatMul

Selects $K$ groups of arbitrary objects with inter-object relationships

**Condition on:**
- Image context $v_{c,t}$
- Object sets $O_t$
- Previous interactions $h_{t-1}$

## 6 VIDEO CAPTIONING – *SINet-Caption*

**Caption:**
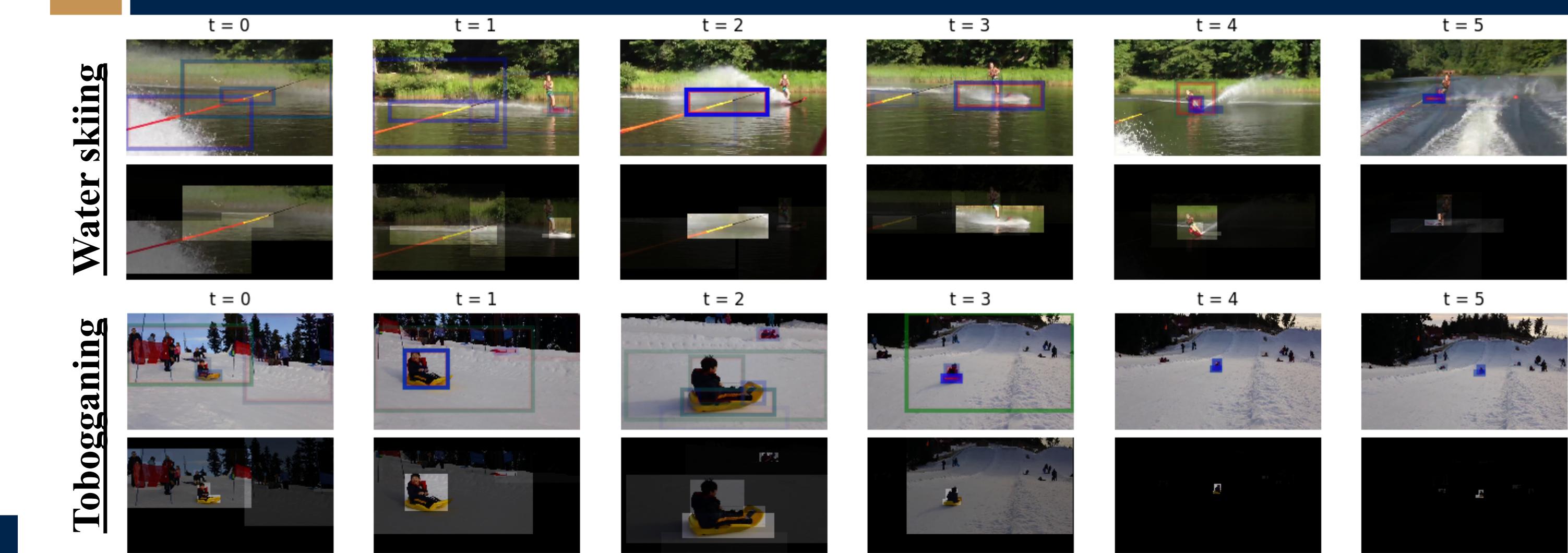A group of people get off of a yellow school bus with life rafts around their neck.

**Relationships:**
[ group of people, get off, yellow school bus ]
[ group of people, with, life rafts ]

Down the Delaware · New Jersey

MatMul / SoftMax / Masking / FC / Tanh $g_\phi(V_c)$

$h_{t_w-1}^2$ → Language LSTM → $h_{t_w}^2$

Attend / Co-attend   $\{h_1, h_2, ..., h_T\}$

$h_{t_w-1}^1$ → Attention LSTM → $h_{t_w}^1$

$repeat(W_h h_{t_w}^1) + W_c g_\phi(V_c)$

$h_{t_w-1}^2$ — Previous lang output; $\overline{g_\phi(V_c)}$ — Video representation; $W_e \Pi_{t_w-1}$ — One-hot vector

The **Attention LSTM** identifies which parts of the video in spatiotemporal feature space are needed for **Language LSTM** to generate the next word.

## 7 EXPERIMENT – KINETICS & ACTIVITYNET CAPTIONS

| Method | Top-1 | Top-5 | FLOP($e^9$) |
|---|---|---|---|
| **Prior Arts** | | | |
| I3D (25 FPS) (test) | 71.1 | 89.3 | |
| TSN (Inception-ResNet-v2) (2.5 FPS) | 73.0 | 90.9 | |
| **Ours (1 FPS)** | | | |
| Img feat + LSTM (baseline) | 70.6 | 89.1 | |
| Img feat + temporal SDP-Attn | 71.1 | 89.6 | |
| Obj feat (mean-pooling) | 72.2 | 90.2 | |
| Obj pairs (mean-pooling) | 73.4 | 90.8 | 18.3 |
| Img + obj feat (mean-pooling) | 73.1 | 91.1 | |
| SINet (K = 1) | 73.9 | 91.3 | **2.7** |
| SINet (K = 2) | 74.2 | 91.5 | 5.3 |
| SINet (K = 3) | **74.2** | **91.7** | 8.0 |

| Method | B@4 | R | M | C |
|---|---|---|---|---|
| **Test set** | | | | |
| LSTM-YT (C3D) | 1.24 | - | 6.56 | 14.86 |
| S2VT (C3D) | 2.62 | - | 7.85 | 20.97 |
| H-RNN | 2.53 | - | 8.02 | 20.18 |
| S2VT + Full context | 3.98 | - | 9.46 | 24.56 |
| LSTM-A + policy gradient + retrieval (ResNet + P3D ResNet) | - | 12.84 | - | - |
| **Validation set (Avg. 1st and 2nd)** | | | | |
| LSTM-A + policy gradient + retrieval (ResNet + P3D ResNet) | **3.13** | 14.29 | 8.73 | 14.75 |
| SINet-Caption – img (ResNeXt) | 1.84 | 20.46 | 9.56 | 43.12 |
| SINet-Caption – obj (ResNeXt) | 1.92 | 20.67 | 9.56 | 44.02 |
| SINet-Caption – img + obj – no co-attn | 2.03 | 21.08 | 9.79 | 44.81 |
| SINet-Caption – img + obj – co-attn | 1.98 | **21.25** | **9.84** | **44.84** |

## 8 QUALITATIVE RESULTS

**Water skiing** — t = 0, t = 1, t = 2, t = 3, t = 4, t = 5

**Tobboganing** — t = 0, t = 1, t = 2, t = 3, t = 4, t = 5

t = 141 | the | t = 143 | man | t = 137 | is | t = 135 | then | t = 174 | shown | t = 117 | on | t = 137 | the | t = 145 | water | t = 143 | skiing

**Distinguish interactions when actions with common objects presented – _horse_:**
- People are _riding_ horses.
- A woman is _brushing_ a horse.
- People are playing _polo_ on a field.
- The man ties _up_ the calf.

riding | brushing | (playing) polo | (ties) up