

# **Glimpse Clouds: Human Activity Recognition from Unstructured Feature Points**

Fabien Baradel, Christian Wolf, Julien Mille,  
Graham W. Taylor

- Authorship
- Background
- Proposed Method
- Experiments
- Conclusion

# ResNet50

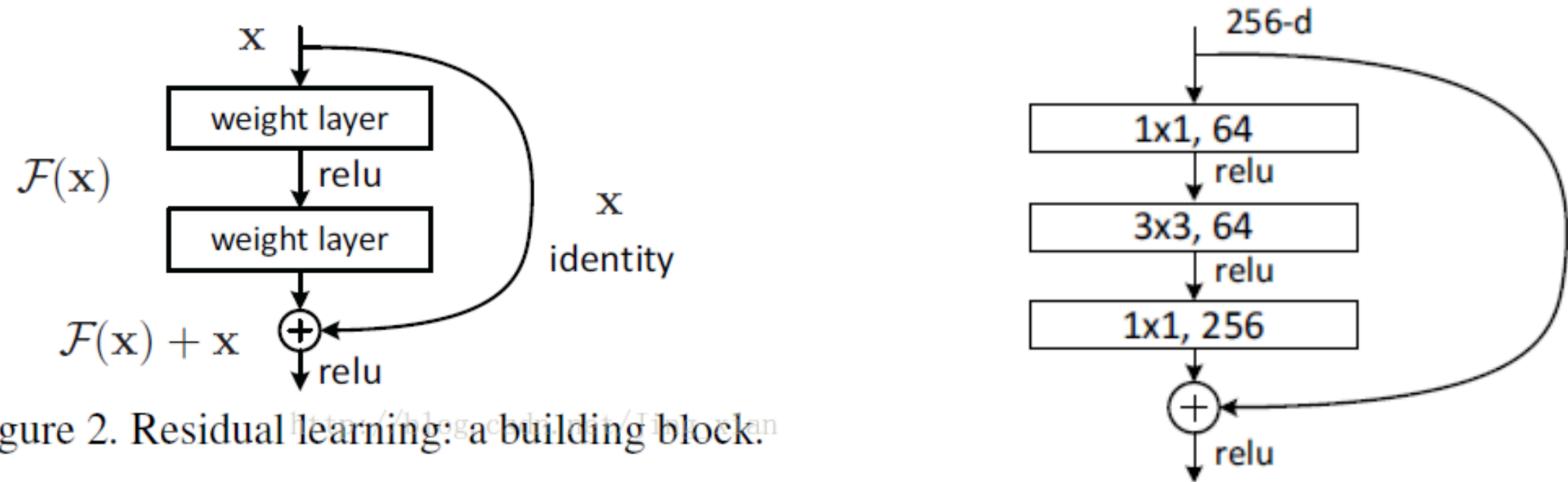


Figure 2. Residual learning: a building block.

- Image classification
- Residual learning & bottleneck
- K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition

# Attention and External Memory

- Soft attention and hard attention
  - Soft: weights each part of the observation dynamically
  - Hard: takes hard decisions when choosing parts of the input data
- External memory for extra information

- Authorship
- Background
- Proposed Method
- Experiments
- Conclusion

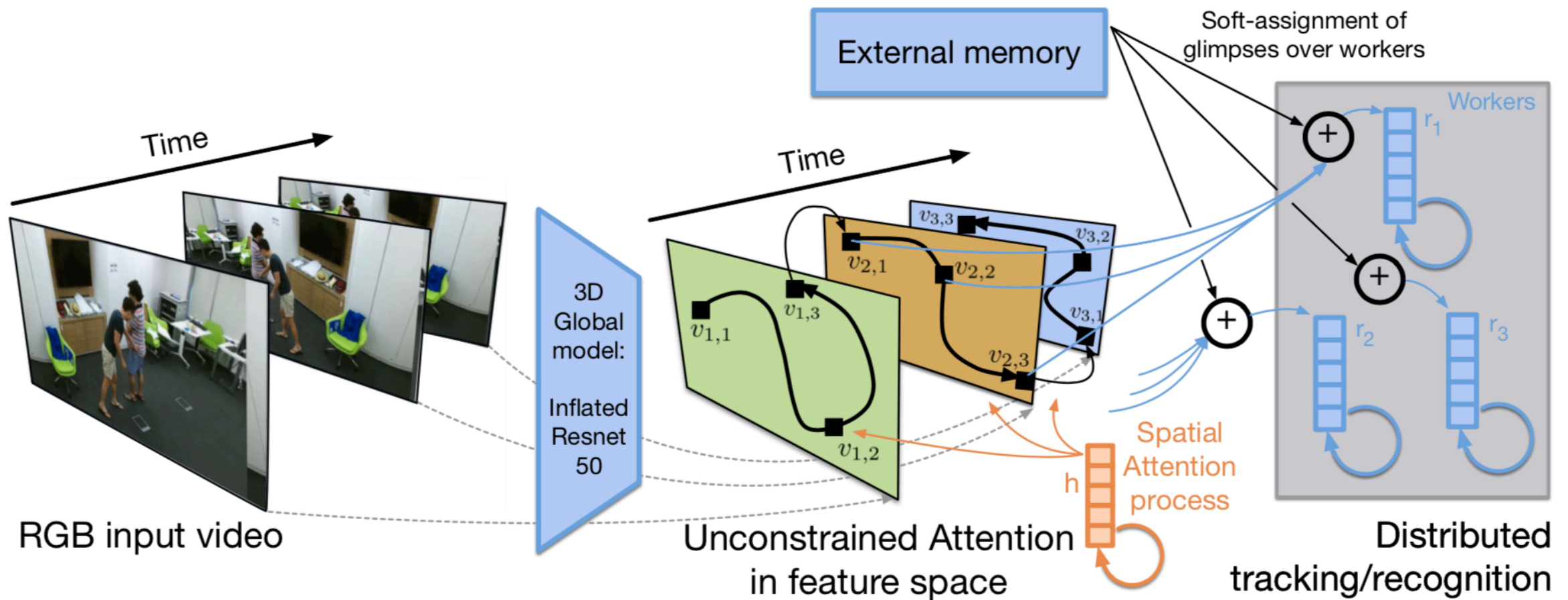


Figure 1. We recognize human activities from unstructured collections of spatio-temporal glimpses with distributed recurrent tracking/recognition and soft-assignment among glimpse points and trackers.

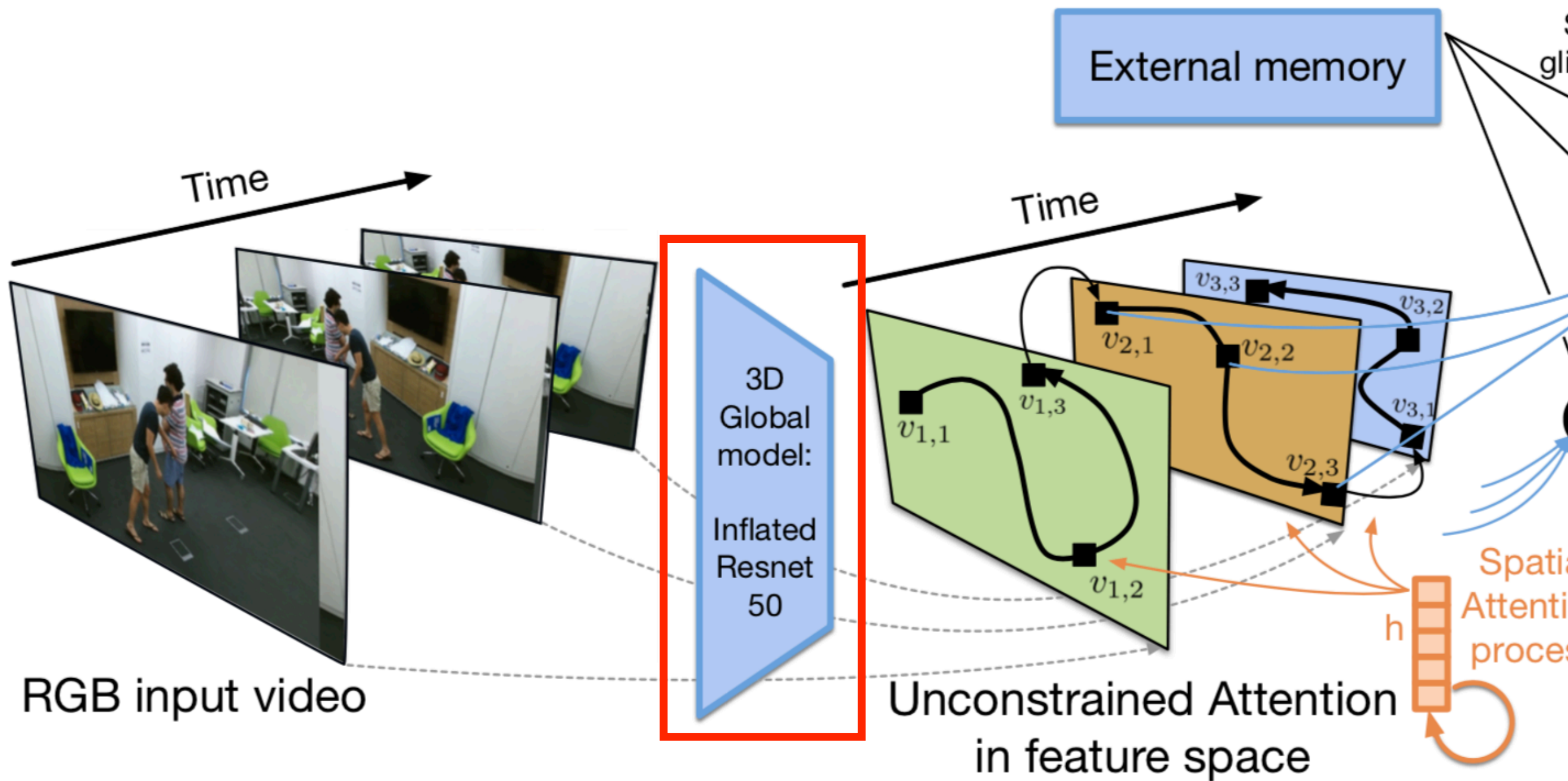
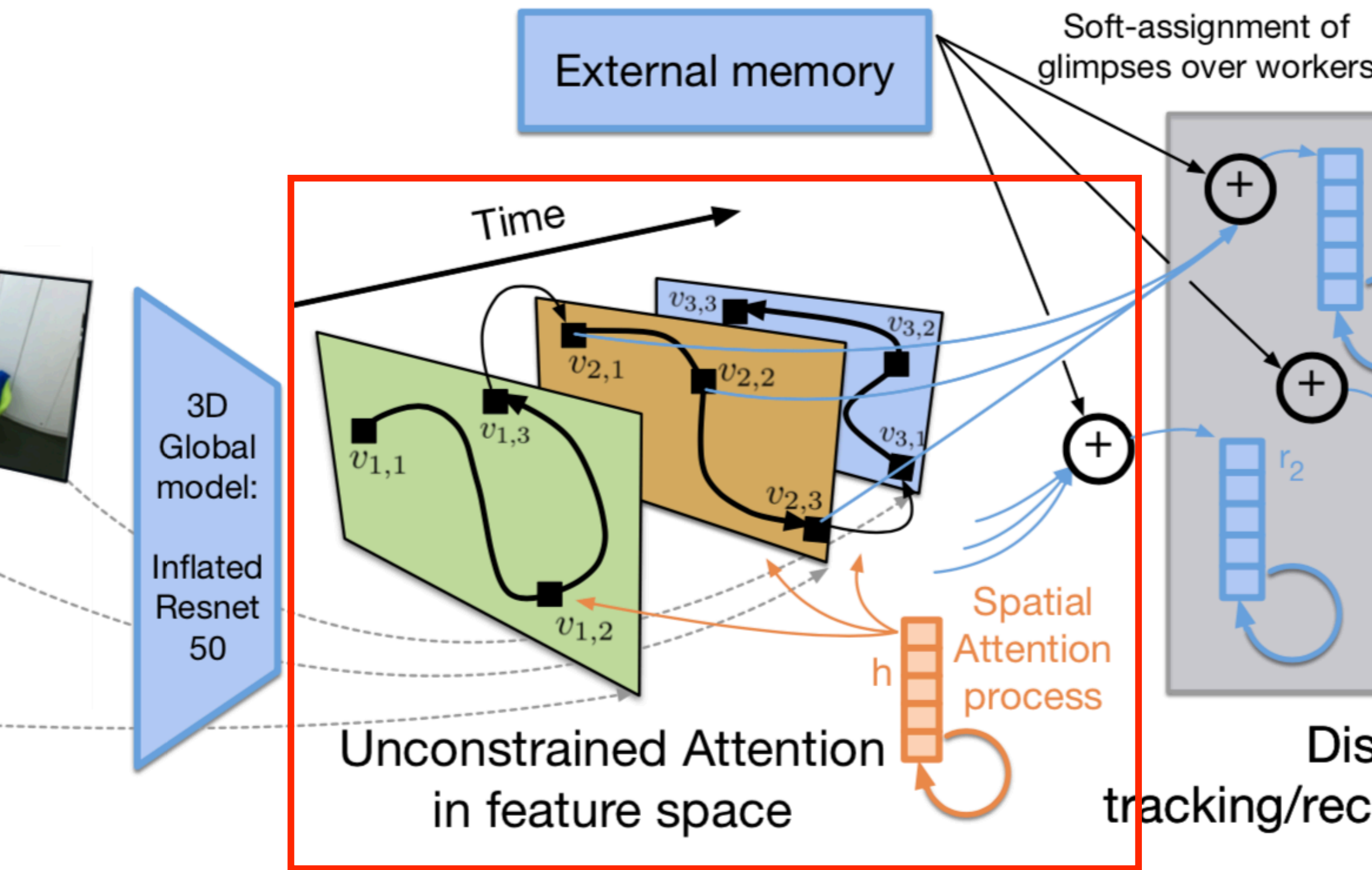


Figure 1. We recognize human activities from unstructured collections of spatio-temporal glimpses by learning global context and soft-assignment among glimpse points and trackers.

- Proceed from the ResNet-50 network and inflate the 2D spatial convolutional kernels into 3D kernels
- Input:  $T * H * W * 3$
- Output:  $T * H' * W' * C'$
- Temporal receptive field is more than one frame





from unstructured collections of spatio-temporal glimpses with distribu

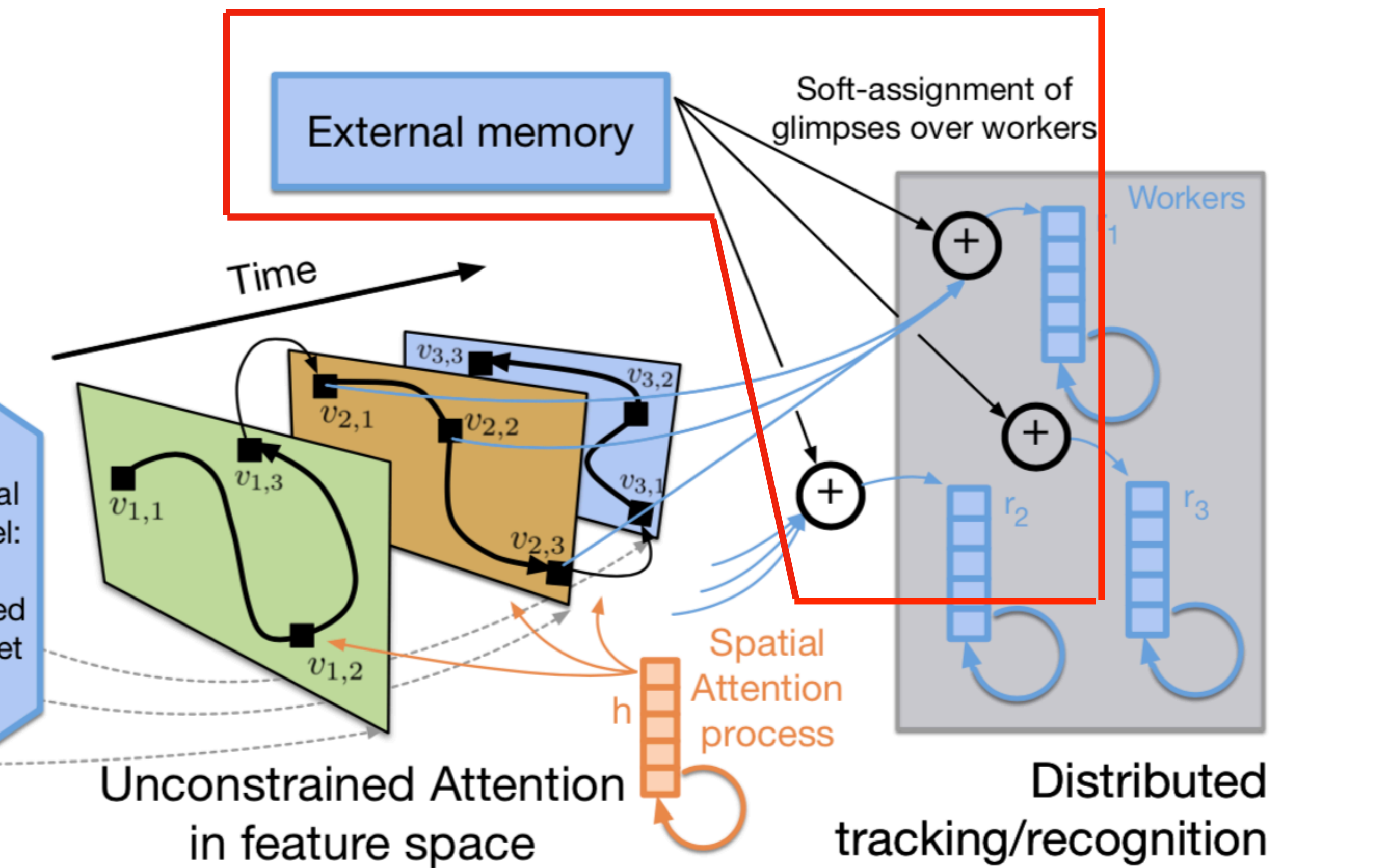
- Attention representation  $\mathbf{l}_{t,g} = [x_g, y_g, s_g^x, s_g^y]_t^\top$
- Attention feature  $\mathbf{z}_{t,g} = \Gamma(\mathbf{Z}_{t,g}) = \frac{1}{H'W'} \sum_m \sum_n \mathbf{Z}_{t,g}(m, n)$
- Feature dimension:  $1 * C$

To be trained

- GRU to predict next:
- P.s. c & r coming soon

$$\mathbf{h}_g = \Omega(\mathbf{h}_{g-1}, [\mathbf{z}_{g-1}, \mathbf{r}_t] | \theta)$$

$$\mathbf{l}_g = W_l^\top [\mathbf{h}_g, \mathbf{c}_t]$$



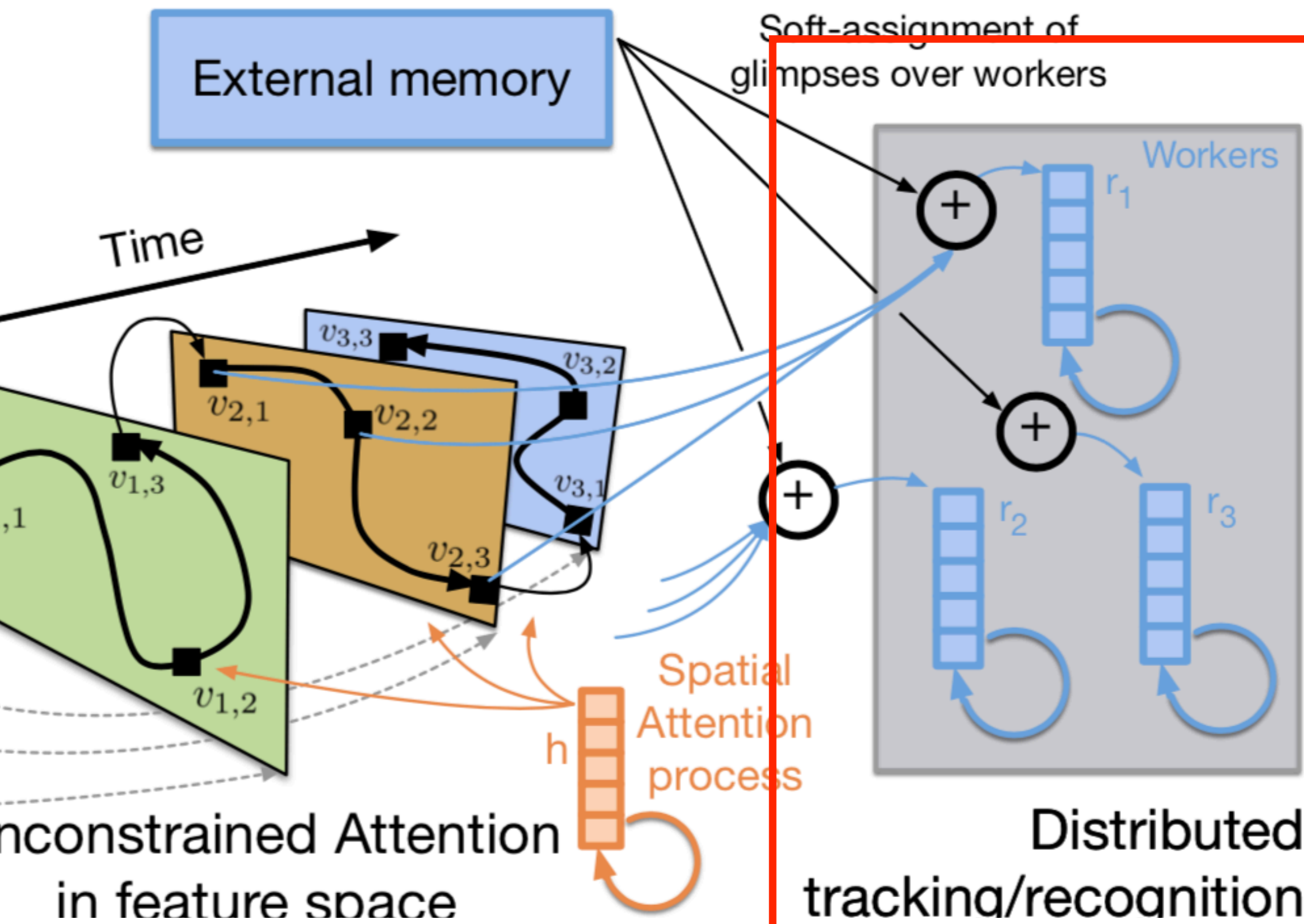
structured collections of spatio-temporal glimpses with distributed recurrent tra

To be trained

- “what” and “where” features:  $\mathbf{v}_{t,g} = \mathbf{z}_{t,g} \otimes \Lambda(\mathbf{l}_{t,g} | \theta_{\Lambda})$
- “workers”: GRUs, take the “important” features
- Input to the workers, linear combination:  $\tilde{\mathbf{v}}_{t,c} = \mathbf{V}_t \mathbf{p}_{t,c}$
- Higher weights to glimpses assigned to this worker in the past
- Using external memory to get the combination weight

To be trained

- Definition of distance:  $\phi(x, y) = \sqrt{(x - y)^\top D(x - y)}$
- External memory bank:  $M = \{m_k\}$
- Each worker's own bank:  $W_c = \{w_{c,k}\}$
- Distance pre-trained as similarity
- p sums to 1:  $p_{t,c,g} = \sigma_\alpha \left( \sum_k e^{-t^{m_k}} \times w_{c,k} [1 - \phi(v_{t,g}, m_k)] \right)$
- Memory bank delete the oldest one when eviction





To be trained

$$\mathbf{r}_{t,c} = \Psi_c(\mathbf{r}_{t-1,c}, \tilde{\mathbf{v}}_{t,c} | \theta_{\Psi_c})$$
$$\mathbf{r}_t = \sum_c \mathbf{r}_{t,c}$$

- “workers”: GRUs

- Recognition:

$$\mathbf{q}_c = W_c \cdot \mathbf{r}_c$$

$$\hat{\mathbf{y}} = \text{softmax} \left( \sum_c^C \mathbf{q}_c \right)$$

- $\mathbf{c}_t$  obtained by global average pooling over the spatial domain of the penultimate feature maps of a given timestep
- Pose estimation:  $\mathbf{y}_t^p = W_p^\top \mathbf{c}_t$
- Recall(P16):  $\mathbf{h}_g = \Omega(\mathbf{h}_{g-1}, [\mathbf{z}_{g-1}, \mathbf{r}_t] \mid \theta)$   
 $\mathbf{l}_g = W_l^\top [\mathbf{h}_g, \mathbf{c}_t]$



# Loss Function

- 3 parts:  $\mathcal{L} = \mathcal{L}_D(\hat{\mathbf{y}}, \mathbf{y}) + \mathcal{L}_P(\hat{\mathbf{y}}^p, \mathbf{y}^p) + \mathcal{L}_G(\mathbf{l}, \mathbf{y}^p)$

- 1 classification, cross-entropy loss

- 2 pose estimation

$$\mathcal{L}_{G_1}^t(\mathbf{l}, \mathbf{y}^p) = \frac{1}{1 + \sum_{g_1}^G \sum_{g_2}^G \|\mathbf{l}_{t,g_1}, \mathbf{l}_{t,g_2}\|}$$

- 3 glimpse

$$\mathcal{L}_{G_2}^t(\mathbf{l}, \mathbf{y}^p) = \sum_g^G \min_j \|\mathbf{l}_{t,g}, \mathbf{y}_j^p\|$$

$$\mathbf{l}_{t,g} = [x_{t,g}, y_{t,g}, s_{t,g}^x, s_{t,g}^y]^T$$

$$\mathcal{L}_G(\mathbf{l}, \mathbf{y}^p) = \sum_t^T (\mathcal{L}_{G_1}^t(\mathbf{l}, \mathbf{y}^p) + \mathcal{L}_{G_2}^t(\mathbf{l}, \mathbf{y}^p))$$

- Authorship
- Background
- Proposed Method
- Experiments
- Conclusion

# Two Datasets

- NTU RGB+D Dataset
- N-UCLA: Northwestern-UCLA Multiview Action 3D Dataset
  - Use the model trained on NTU as a pretrained model and fine-tune it on N-UCLA

Methods	Pose	RGB	CS	CV	Avg
Lie Group [51]	✓	-	50.1	52.8	51.5
Skeleton Quads [13]	✓	-	38.6	41.4	40.0
Dynamic Skeletons [18]	✓	-	60.2	65.2	62.7
HBRNN [11]	✓	-	59.1	64.0	61.6
Deep LSTM [43]	✓	-	60.7	67.3	64.0
Part-aware LSTM [43]	✓	-	62.9	70.3	66.6
ST-LSTM + TrustG. [33]	✓	-	69.2	77.7	73.5
STA-LSTM [46]	✓	-	73.2	81.2	77.2
Ensemble TS-LSTM [29]	✓	-	74.6	81.3	78.0
GCA-LSTM [34]	✓	-	74.4	82.8	78.6
JTM [53]	✓	-	76.3	81.1	78.7
MTLN [23]	✓	-	79.6	84.8	82.2
VA-LSTM [59]	✓	-	79.4	87.6	83.5
View-invariant [35]	✓	-	80.0	87.2	83.6
DSSCA - SSLM [44]	✓	✓	74.9	-	-
STA-Hands [5]	X	X	82.5	88.6	85.6
Hands Attention [6]	✓	✓	84.8	90.6	87.7
C3D†	-	✓	63.5	70.3	66.9
Resnet50+LSTM†	-	✓	71.3	80.2	75.8
<b>Glimpse Clouds</b>	-	✓	<b>86.6</b>	<b>93.2</b>	<b>89.9</b>

Table 2. Results on the NTU RGB+D dataset with Cross-Subject and Cross-View settings (accuracies in %); († indicates method has been re-implemented).

Methods	Spatial Attention	Soft Workers	$L_D$	$L_P$	$L_G$	CS	CV	Avg
Global model	-	-	✓	-	-	84.5	91.5	88.0
Global model	-	-	✓	✓	-	85.5	92.1	88.8
Global model+ $\sum$ Glimpses + GRU	-	-	✓	✓	-	85.8	92.4	89.1
Glimpse Clouds	✓	✓	✓	-	-	85.7	92.5	89.1
Glimpse Clouds	✓	✓	✓	✓	-	86.4	93.0	89.7
Glimpse Clouds	✓	✓	✓	-	✓	86.1	92.9	89.5
Glimpse Clouds	✓	✓	✓	✓	✓	<b>86.6</b>	<b>93.2</b>	<b>89.9</b>
Glimpse Clouds + Global model	✓	✓	✓	✓	✓	86.6	93.2	89.9

Table 3. Results on NTU: ablation study

Glimpse	Type of attention	CS	CV	Avg
3D tubes	Attention	85.8	92.7	89.2
Seq. 2D	Random sampling	80.3	87.8	84.0
Seq. 2D	Saliency	86.2	92.9	89.5
Seq. 2D	<b>Attention</b>	<b>86.6</b>	<b>93.2</b>	<b>89.9</b>

Table 4. Results on the NTU: different attention and alternative strategies.

Predict glimpse location  
through a location network.  
No past information

- Other experiments determines other parameters like the worker number or the glimpse number

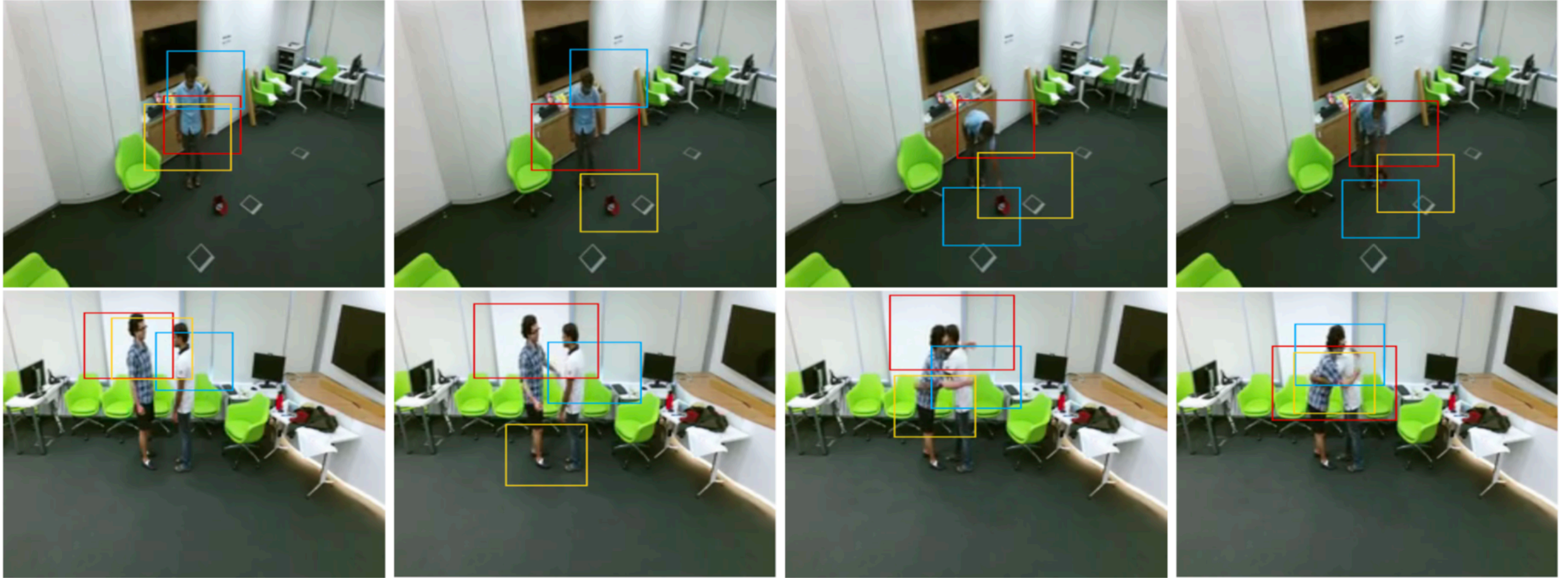


Figure 3. An illustration of the glimpse distribution for several sequences of the NTU dataset. Here we set 3 glimpses per frame ( $G=3$ , Red: first glimpse, Blue: second glimpse, Yellow: third one).



- Authorship
- Background
- Proposed Method
- Experiments
- Conclusion

- Only RGB when testing, though pose information when training
- Attention process, soft-assigned
- A little bit complex but code released