



# Multipar-T: Multiparty-Transformer for Capturing Contingent Behaviors in Group Conversations

**Dong Won Lee , Yubin Kim , Rosalind Picard , Cynthia Breazeal , Hae Won Park**

Massachusetts Institute of Technology

{dongwonl, ybkim95, picard, cynthiab, haewon}@mit.edu,

<https://github.com/dondongwon/Multipar-T>

— IJCAI 2023

2023. 6. 8 • ChongQing



gesis  
Leibniz-Institut  
für Sozialwissenschaften



Reported by JiaWei Cheng



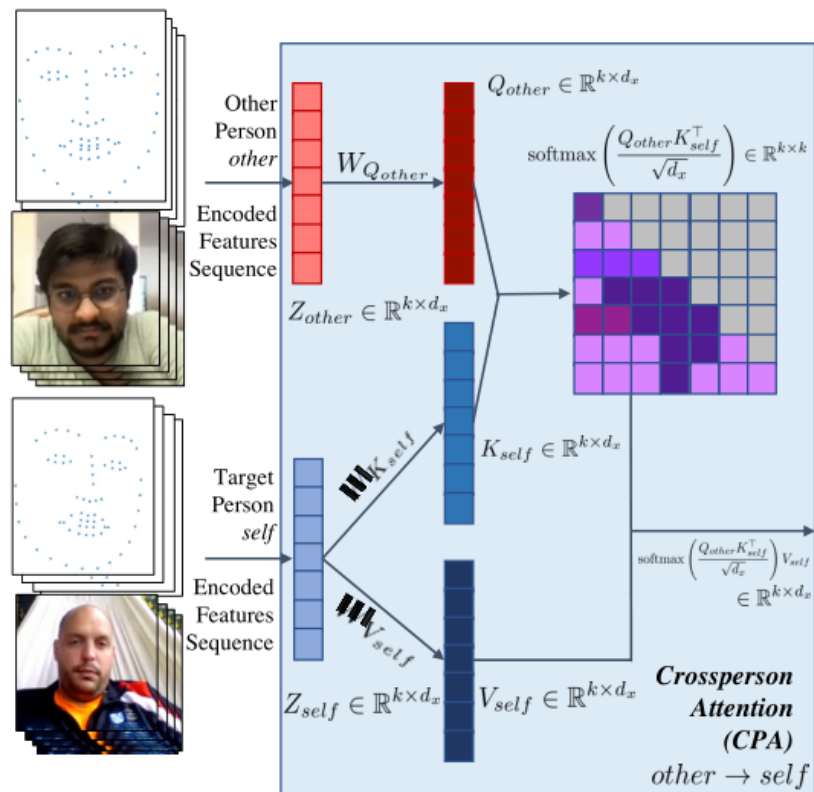
# Motivation

Firstly, the system must perform well in recognizing individual behavioral cues

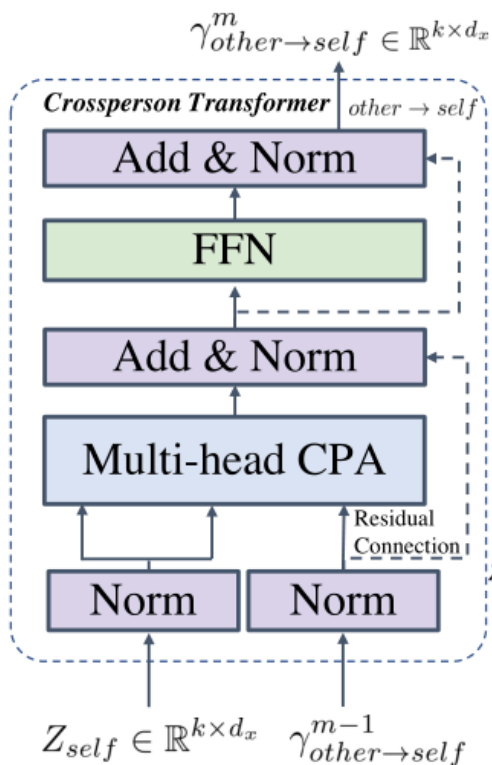
Secondly, it must do so simultaneously, while keeping track of every individual in the group

Finally, it must also recognize the subtle interactions that take place between group members as it can provide more insights into what is being communicated.

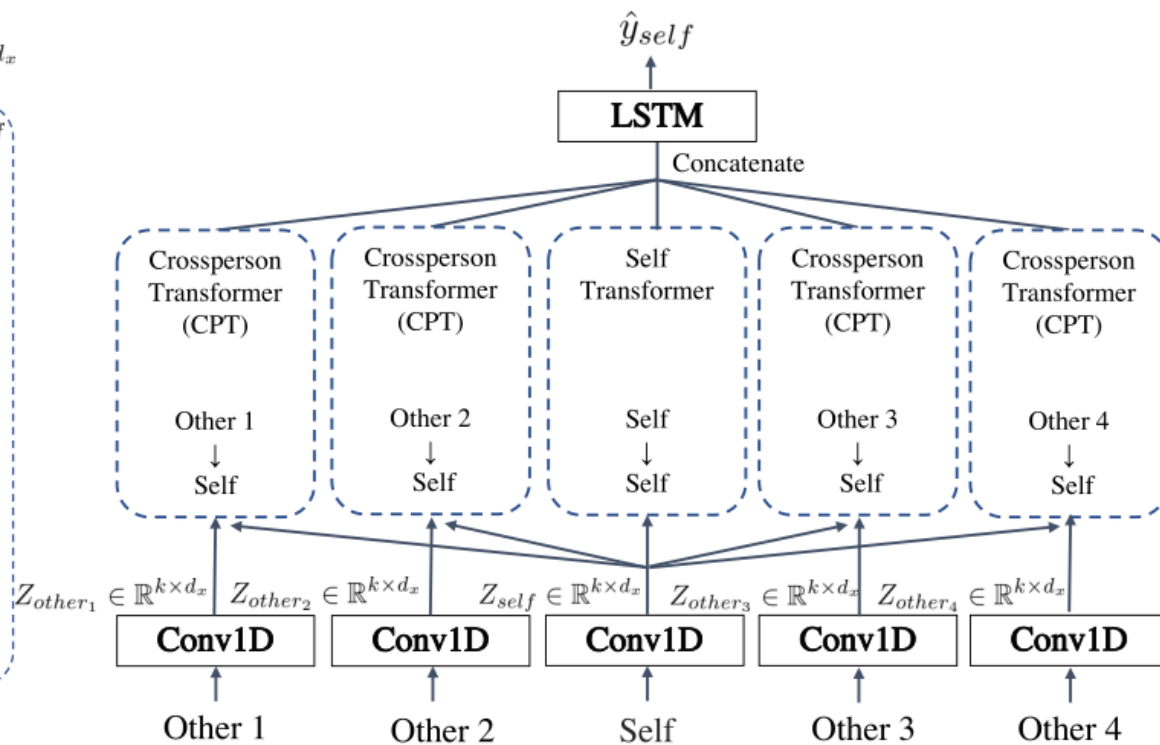
# Overview



(a) Crossperson Attention

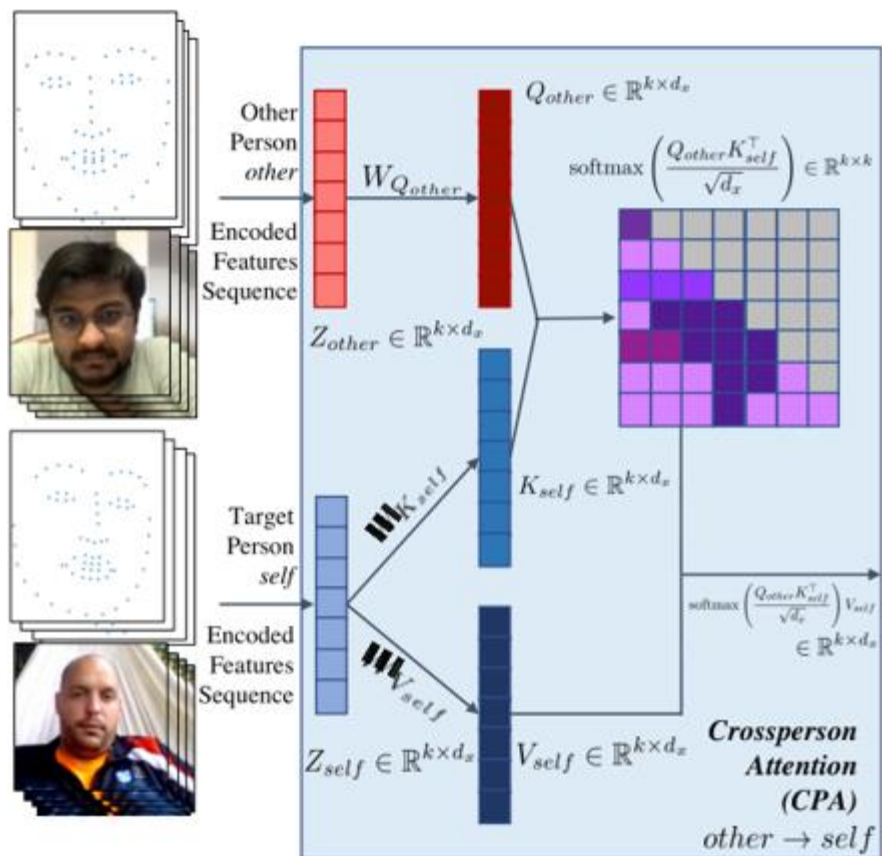


(b) Crossperson Transformer



(c) Multiparty-Transformer (MultiPar-T)

# Method



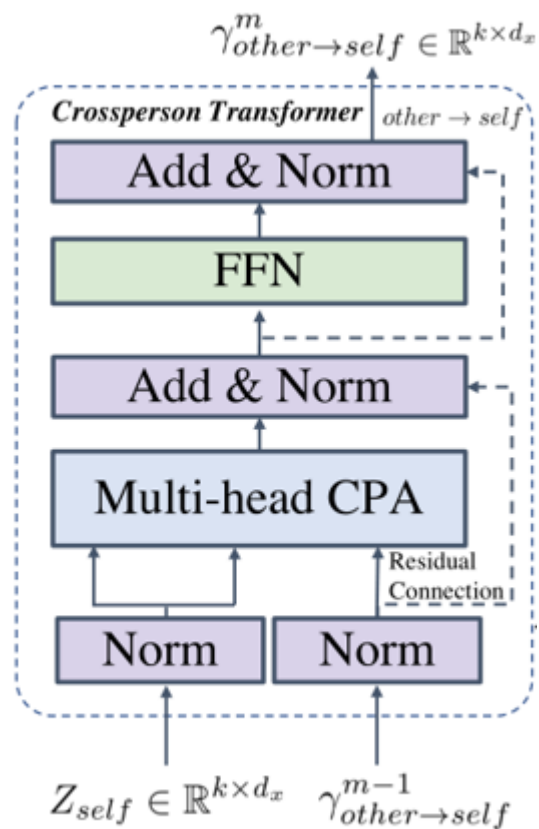
(a) Crossperson Attention

$$\begin{aligned} \text{CPA}_{other \rightarrow self}(Z_{other}, Z_{self}) &= \text{softmax}\left(\frac{Q_{other}K_{self}^T}{\sqrt{d_x}}\right) V_{self} \\ &= \text{softmax}\left(\frac{Z_{other}W_{Q_{other}}(Z_{self}W_{K_{self}})^T}{\sqrt{d_x}}\right) Z_{self}W_{V_{self}}. \end{aligned} \quad (1)$$

$$\begin{aligned} \text{CPA}_{other \rightarrow self}^{multi}(Z_{other}, Z_{self}) \\ = \text{Concat}\left(\text{CPA}_{other \rightarrow self}^1, \dots, \text{CPA}_{other \rightarrow self}^h\right) W^{multi} \end{aligned} \quad (2)$$

$$Z_p = \text{Conv1D}(X_p) + \text{PE}(X_p) \quad (3)$$

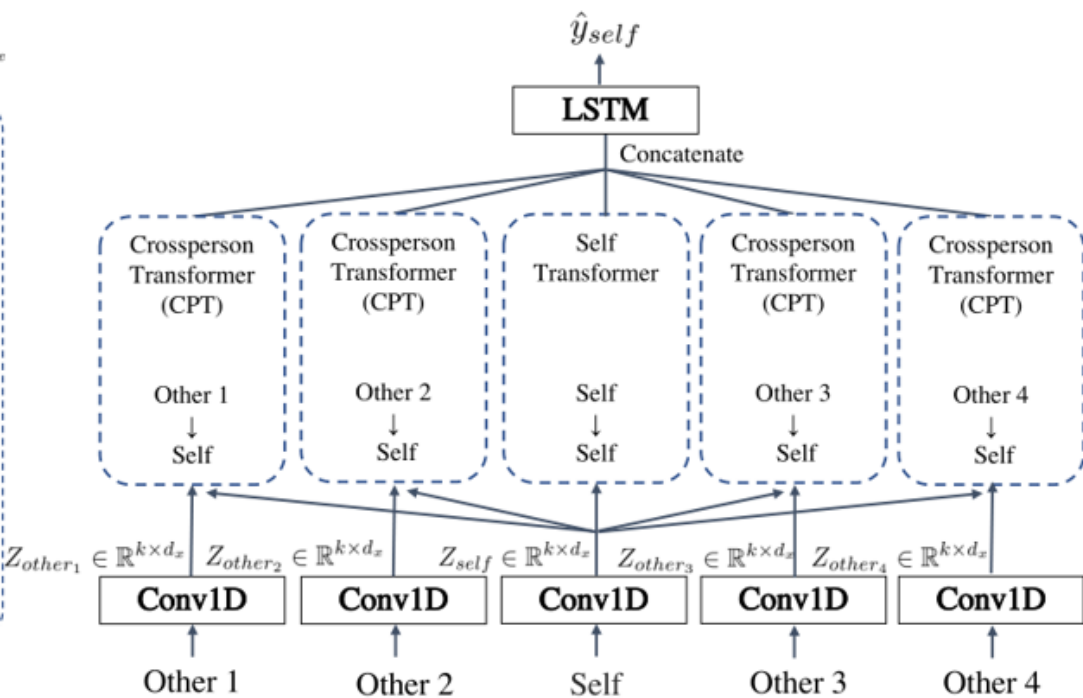
# Method



(b) Crossperson Transformer

$$\begin{aligned} \gamma_{other \rightarrow self}^m &= \text{CPT}_{other \rightarrow self}^m(\gamma_{other \rightarrow self}^{m-1}, Z_{self}) \\ \hat{\gamma}_{other \rightarrow self}^m &= \text{CPA}_{other \rightarrow self}^{m, multi}(\text{Norm}(\gamma_{other \rightarrow self}^{m-1}), \text{Norm}(Z_{self})) \\ &\quad + \text{Norm}(\gamma_{other \rightarrow self}^{m-1}) \\ \gamma_{other \rightarrow self}^m &= \text{Norm}(\text{FFN}(\hat{\gamma}_{other \rightarrow self}^m) + \hat{\gamma}_{other \rightarrow self}^m) \end{aligned} \quad (4)$$

# Method



(c) Multiparty-Transformer (MultiPar-T)

$$\zeta_{self}^n, hidden^n = \text{LSTM}([\gamma_{1 \rightarrow self}^M \parallel \dots \parallel \gamma_{P \rightarrow self}^M], hidden^{n-1})$$

$$\hat{Y}_{self} = \text{FFN}(\zeta_{self}^k) \quad \text{for } n \in [k] \quad (5)$$

$$L_{Focal} = -\frac{1}{N} \sum_i \sum_c (1 - \hat{Y}_{ic})^\alpha y_{ic} \log(\hat{Y}_{ic}) \quad (6)$$



# Experiments

Model	All Engagement Classes			High Dis-Eng.	Low Dis-Eng.	Low Eng.	High Eng.
	Accuracy	Weighted F1	Macro F1	F1	F1	F1	F1
ConvLSTM [Del Duchetto <i>et al.</i> , 2020]	$0.859 \pm 0.01$	$0.857 \pm 0.02$	$0.699 \pm 0.05$	0.741	$0.459 \pm 0.22$	$0.699 \pm 0.12$	$0.907 \pm 0.01$
OCTCNN-LSTM [Steinert <i>et al.</i> , 2020]	$0.769 \pm 0.08$	$0.695 \pm 0.14$	$0.410 \pm 0.10$	0.588	$0.119 \pm 0.17$	$0.233 \pm 0.33$	$0.864 \pm 0.05$
TEMMA [Chen <i>et al.</i> , 2020a]	$0.823 \pm 0.02$	$0.822 \pm 0.02$	$0.561 \pm 0.11$	0.286	$0.254 \pm 0.19$	$0.621 \pm 0.13$	$0.885 \pm 0.01$
EnsModel [Thong Huynh <i>et al.</i> , 2019]	$0.760 \pm 0.07$	$0.675 \pm 0.12$	$0.302 \pm 0.03$	0	$0.000 \pm 0.00$	$0.160 \pm 0.23$	$0.860 \pm 0.05$
BOOT [Wang <i>et al.</i> , 2019]	$0.817 \pm 0.03$	$0.822 \pm 0.03$	$0.636 \pm 0.09$	0.714	$0.320 \pm 0.24$	$0.658 \pm 0.12$	$0.873 \pm 0.02$
HTMIL [Ma <i>et al.</i> , 2021]	$0.820 \pm 0.02$	$0.818 \pm 0.02$	$0.460 \pm 0.05$	0	$0.000 \pm 0.00$	$0.633 \pm 0.12$	$0.880 \pm 0.02$
GAT [Zhang <i>et al.</i> , 2022]	$0.739 \pm 0.06$	$0.631 \pm 0.08$	$0.261 \pm 0.03$	0	$0.000 \pm 0.00$	$0.006 \pm 0.01$	$0.848 \pm 0.04$
MuT [Tsai <i>et al.</i> , 2019]	$0.847 \pm 0.02$	$0.845 \pm 0.02$	$0.624 \pm 0.12$	0.625	$0.310 \pm 0.25$	$0.665 \pm 0.12$	$0.901 \pm 0.01$
Multipar-T (Ours)	<b><math>0.888 \pm 0.03</math></b>	<b><math>0.887 \pm 0.03</math></b>	<b><math>0.751 \pm 0.05</math></b>	<b>0.800</b>	<b><math>0.559 \pm 0.07</math></b>	<b><math>0.759 \pm 0.11</math></b>	<b><math>0.927 \pm 0.02</math></b>

Table 1: Results and standard deviations for engagement recognition models for 3 seeds (std dev for High Dis-Eng. not reported due to 2 seeds not having corresponding labels). Despite high accuracy and weighted-F1 scores, many previous baselines fail at infrequent disengagement classes. Multipar-T outperforms other approaches across all metrics.

# Experiments

Attention Direction	Ablation	All Classes			High Dis-Eng.	Low Dis-Eng.	Low Eng.	High Eng.
		Accuracy	Weighted F1	Macro F1	Binary F1	Binary F1	Binary F1	Binary F1
	Multipar-T <i>w/o</i> Crossperson Transformer	0.847 + 0.0154	0.844 + 0.14	0.661 + 0.018	0.588	0.433 + 0.1	0.66 + 0.12	0.901 + 0.01
$CPA_{self \rightarrow other}$	Multipar-T <i>w/o</i> Self Transformer	0.847 + 0.0167	0.845 + 0.021	0.624 + 0.12	0.625	0.31 + 0.25	0.665 + 0.12	0.901 + 0.01
	Multipar-T	0.865 + 0.03	0.862 + 0.036	0.735 + 0.02	<b>0.769</b>	<b>0.587 + 0.12</b>	0.698 + 0.15	0.912 + 0.02
$CPA_{other \rightarrow self}$	Multipar-T <i>w/o</i> Self Transformer	<b>0.883 + 0.02</b>	<b>0.884 + 0.024</b>	<b>0.75 + 0.04</b>	<b>0.769</b>	<b>0.555 + 0.11</b>	<b>0.762 + 0.08</b>	<b>0.923 + 0.02</b>
	Multipar-T	<b>0.883 + 0.02</b>	<b>0.885 + 0.02</b>	<b>0.75 + 0.06</b>	0.714	0.557 + 0.19	<b>0.766 + 0.08</b>	<b>0.923 + 0.02</b>

Table 2: Ablation results for Self Transformer and Crossperson Transformer mechanisms. Attending to *other*'s and own *self* behaviors boosts performance. We refer the readers to Figure 1. Multipar-T *w/o* Crossperson Transformer refers to the ablation of all pairwise Crossperson Transformers with only the Self Transformer remaining. Multipar-T *w/o* Self Transformer refers the ablation of the Self Transformer and utilizing the pairwise Crossperson Transformers. Results with different directions of Crossperson Attention are displayed, where  $CPA_{other \rightarrow self}$  performs well generally and  $CPA_{self \rightarrow other}$  performs well for disengaged instances.



# Experiments

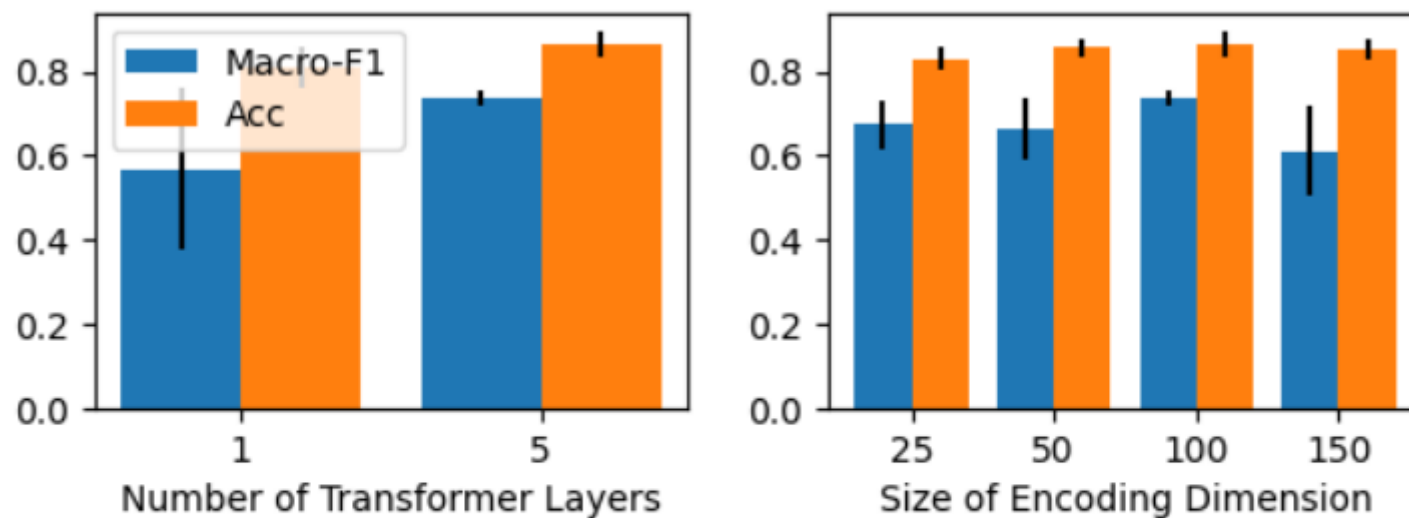


Figure 3: Macro-F1 and Accuracy scores for important hyperparameters for Multipar-T. (Left) Multi-layered transformers and (Right) encoding dimension of  $d_x = 100$  boosts performance.

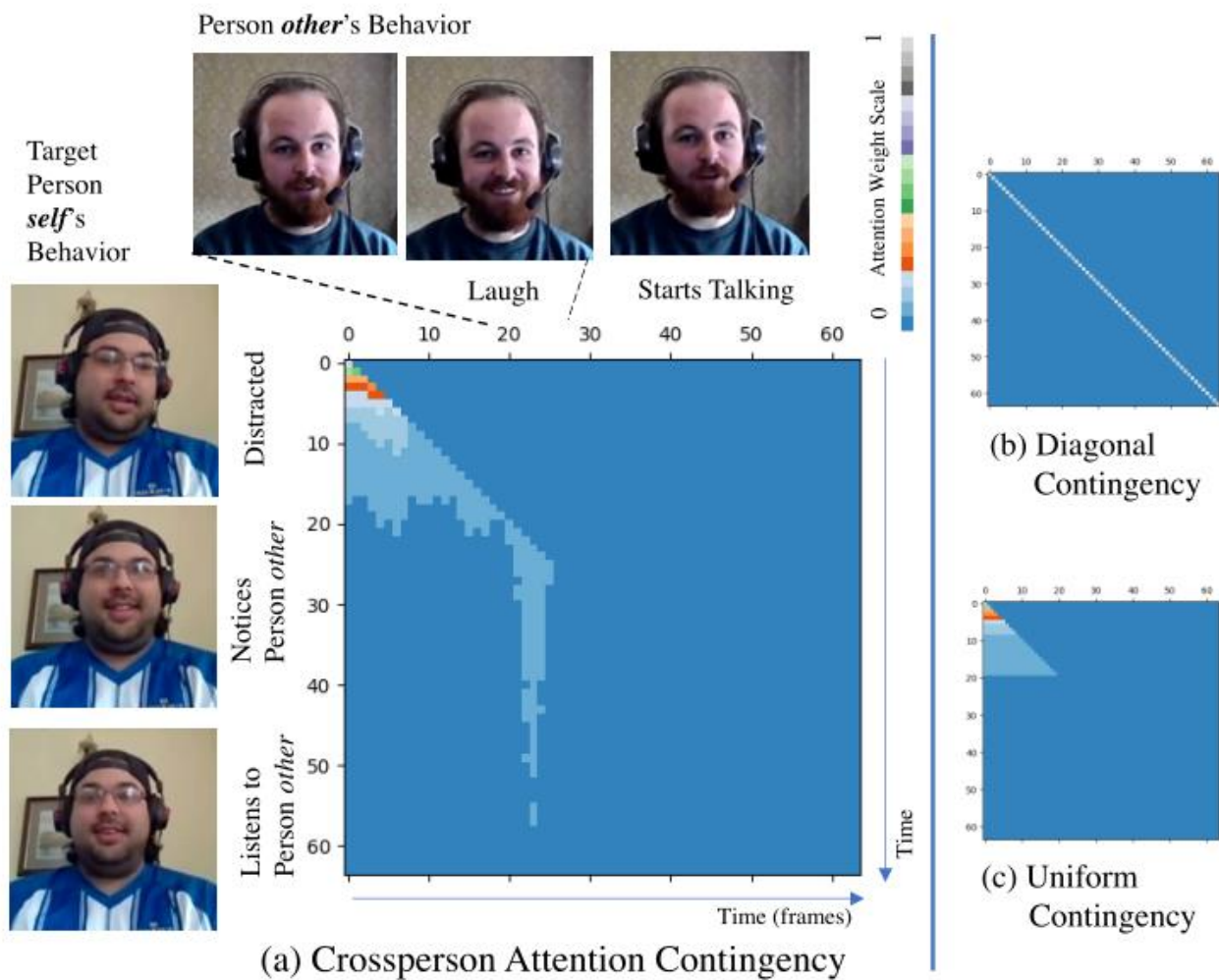


# Experiments

Model	Accuracy	Weighted F1	Macro F1
I3D [Wang <i>et al.</i> , 2018]	$0.751 \pm 0.07$	$0.658 \pm 0.08$	$0.254 \pm 0.05$
TimeSformer [Bertasius <i>et al.</i> , 2021]	$0.806 \pm 0.03$	$0.752 \pm 0.05$	$0.337 \pm 0.14$
SlowFast [Feichtenhofer <i>et al.</i> , 2019]	$0.718 \pm 0.11$	$0.628 \pm 0.12$	$0.232 \pm 0.02$
Multipar-T (Ours)	<b><math>0.828 \pm 0.02</math></b>	<b><math>0.823 \pm 0.02</math></b>	<b><math>0.466 \pm 0.06</math></b>

Table 3: Raw video-based action recognition models and Multipar-T trained with less computationally heavy training set-up. Results and standard deviation are reported for 3 seeds. We see the limitations of training end-to-end raw video-based models.

# Experiments





# Thanks!