# **Dyadformer**

#### A Multi-modal Transformer for Long-Range Modeling of Dyadic Interactions

David Curto, Albert Clapés, Javier Selva, Sorina Smeureanu, Julio C. S. Jacques Junior, David Gallardo-Pujol, Georgina Guilera, David Leiva, Thomas B. Moeslund, Sergio Escalera, and Cristina Palmero

Presenter: Albert Clapés (postdoc AAU)

Slides by: Javier Selva







### **Motivation**

- Growing interest in **social interaction understanding**.
- **Dyadic context** is crucial to understand both **individual behavior** and **interaction** characteristics.
- Evaluation on UDIVA v0.5!
- **Lacking** in competing methods:
  - Long-range temporal information for self-reported personality prediction.
  - Joint modeling of both interlocutors.
  - Use of **metadata** is uncommon.



Source: https://pixabay.com/photos/people-girls-women-students-2557396

#### Contributions

- Purely data-driven method to model (and infer) self-reported personality in <u>dyadic interactions.</u>
- It is using larger temporal context w.r.t. previous works (**+30 seconds**).
- Leverage cross-attention to **jointly model both subjects** in the interaction.
- New state-of-the-art results on UDIVA v0.5.

### **Self-attention and Cross-attention**

(+)

&

- **Transformer**-based architecture -
- Q, K, V **Attention** -





Self-attention

# **The Dyadformer**

Input: sequences of video & audio clip-level features.

# Embeddings: pretrained R(2+1)D and VGGish.

- **Add** both **positional** and **metadata** information.
- Weight sharing across streams & intra-encoder layers.



#### Input

- Pre-segmented ~3-second clips of aligned audio-video.
- T-clips sequences as input (T ∈ {6,12}) resulting in up to ~30 seconds.
  - 67.5/10.9/5.6 hours (Train / Validation / Test)
- Metadata:
  - **Participant** (age, gender, cultural background, session index of participant, and pre-session mood/fatigue).
  - **Session** (task order within the session and its difficulty).
  - **Dyadic information** (relationship between interactants).



Self-reportedOCEANvaluesasz-scores(standardized)

values) clipped the range [-3,3]:

- **O**pen-mindedness
- Conscientiousness
- Extraversion
- Agreeableness
- **N**egative emotionality

for **both subjects** simultaneously.

#### Dataset: UDIVA v0.5







- Audiovisual recordings of face-to-face dyadic interactions.
- **134 participants** (17-75 y.o., 55.2% male)
- Participant and session **metadata**
- 145 interaction sessions (4 tasks each)

- **Multilingual**: Spanish (73.1%), Catalan (17.25%) and English (9.65%)
- Includes self-reported OCEAN personality scores.
- Other labels.

#### **Evaluation**

Both **loss** and **evaluation metric** are computed at sequence level and based on **MSE**<sub>seq</sub>.

We also report  $MSE_{part}$  – the MSE at participant level – by aggregating (i.e. median) their OCEAN predictions before computing the error.

**Pearson Correlation** is measure correctness of predictions' distribution.

**Experiments**:

- **Ablation:** #layers, types of cross-attention, and temporal receptive field.
- Comparison to SOTA

## Ablation



encoder

 $\bigcirc \rightarrow \bullet \bullet \bullet (M)$ 

Audio



# Ablation

#### Ablated:

- L = #layers in the encoders (sharing weights).
- **T** = #clips in the audiovisual sequences.
- Arch. = model architectures.

#### Findings:

- The more **temporal context** the better (T = 12).
- Cross-modal & cross-subject (DF<sub>xm,xs</sub>) improve prediction.



T = #clips in the sequences; L = #layers in the encoders

## Ablation results per OCEAN trait and task

Different tasks enact different traits.

Models using **crosssubject interactions** show **better** Pearson **Correlation**.





### **SOTA COMPARISION**

put	Feature extraction & STE	Query preproces	ssor			()) B	roadcast concatenat	on	
- T						€ E: ⊙ R	E Expansion of dimensions Residual addition		
Face chunk	Arch.	0	С	Е	Α	Ν	Avg.		
	LEam [2]	0.744	0.794	0.886	0.653	1.012	0.818		
Sec.	DF <sub>xm,xs</sub>	0.646	0.664	0.699	0.614	0.989	0.722		
Local Context chunk	R(2+1)D * shared wegtes * shared wegtes * ST		Local Context features				OCEAN		
-	Arch.	0	C	E	A	N	Avg.	]	
Extended Contex	LEam [2]	-0.084	0.404	0.386	0.219	0.439	0.273		
C chunk	DF <sub>xm,xs</sub>	0.401	0.517	0.490	0.392	0.350	0.430	]	
Extended metadata Audio signal	$\rightarrow \overbrace{VGGish}^{\text{Extends}} \rightarrow \stackrel{\text{Stendsmetadat}}{\stackrel{\text{Stendsmetadat}}{\stackrel{\text{Total Stendsmetadat}}{\stackrel{\text{Total Stendsmetadat}}{\stackrel$	$ \begin{array}{c} \operatorname{rd} \\ \operatorname{s} \\ \\ \operatorname{s} \\ \end{array} \rightarrow \begin{array}{c} \operatorname{\mathfrak{S}} \\ \operatorname{\mathfrak{S}} \end{array} \rightarrow \begin{array}{c} \operatorname{\mathfrak{S}} \\ \operatorname{\mathfrak{S}} \end{array} = \begin{array}{c} \\ \operatorname{\mathfrak{S}} \\ \operatorname{\mathfrak{S}} \end{array} $	$ \xrightarrow{1 \times 1 \times 100} \rightarrow \mathbb{C} $		Keys Values	Dropout LayerNorm FC + ReU		Updated query	

- In [2], we worked with **3second chunks** (<u>here, ~30</u> <u>seconds</u>).
- Multi-modal fusion was done by concatenating (<u>here, self-</u> /cross-attention).
- Regresses personality of the target subject (<u>here, both</u> <u>simultaneously</u>).



- Successful data-driven architecture for predicting self-reported personality on observed dyadic interactions.
- **Multi-modal fusion** via cross-attention, i.e. cross-modal attention.
- Explicitly modeling **both subjects jointly**, i.e. cross-subject attention.
- Relatively long temporal receptive fields (+30 secs).
- SOTA on UDIVA v0.5.

#### **Future work**

- Exploiting UDIVA v0.5 **dialogue transcripts**.
- Designing **self-supervised pre-text tasks**.
- **Transfer perceived personality** prediction model and/or **multi-task**.
- Perceived personality as a proxy for selfreported personality, Then, learn a mapping from perceived -> self-reported.
- **Contrasting interaction behaviors** of the same person with different interlocutors.

•		
• •		
$\bigcirc$		

#### **Questions**?

### **Open to collaboration** $\bigcirc$

mailto: alcl@create.aau.dk

#### Other interests:

- Video-related analyses and/or multi-modal (incl. video) -> sequential models (e.g., Transformers).
- XAI and active supervision in vision through language.
- Learning from small/noisy data (i.e., self-supervision, self-learning, knowledge distillation, etc).
- **Applications**: Affective, assistive, edge, and personalized AI.