



Commensality

*the practice of sharing food and eating
together in a social group*
Ochs and
Shohet, 2006

- important social activity
- one of the most frequent and common human experiences
- time for celebration, making business, and creation of new social bounds
- food has a social and emotional meaning
- several impacts of commensality were observed



Aim

We aim to show the **feasibility** of using machine learning to **recognize a variety of commensal activities** from video data collected in a naturalistic setting

Interaction around the table is peculiar:

- partners shift their attention between the conversation and the food

Commensal activities include:

- actions related to food consumption (e.g., chewing, intake, taking, drinking)
- social signals (e.g., smiling, gazing, passing the food)

The previous works focus on a single activity recognition and/or different contexts

Data collection

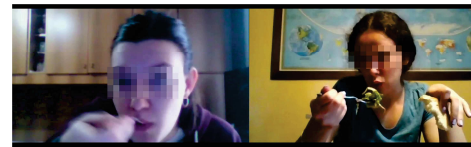
Recordings of dyads eating in a video call:

- participants know each other well
- they eat at their home
- closeup framing including face and upper body
- 18 subjects

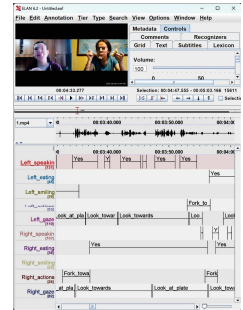


Videos:

- synchronized view
- 96 minutes of synchronized recordings
- average meal duration 9m 23s



Annotation



Five activities annotated:

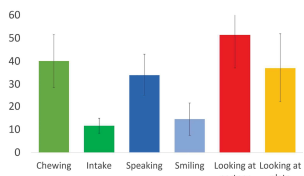
- speaking
- food/drink in-taking
- chewing
- smiling/laughing
- gaze

One expert:

- annotated all videos
- used audio and video for annotation
- annotated each activity independently

Different activities may be performed simultaneously

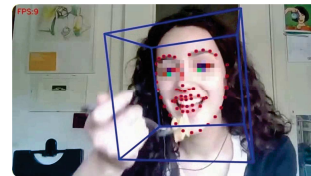
Annotation results



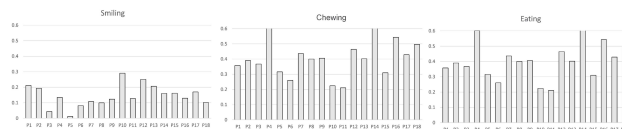
- classes are imbalanced
- social activities are frequent
- interactions are smooth and rich
- nine social metrics computed, e.g.:
 - avg. overlapping speech is 1.6%
 - avg. duration of mutual gaze is 27.85%
- differences between participants and dyads

Feature extraction

- 17 Action Units extracted with OpenFace



- three segment lengths: 10, 25, 50 frames
- only data segments with a single label
- six statistics computed for each AU: min, max, mean, standard deviation, skewness, kurtosis



	Length	Accuracy	F-score	Weigh. F-score		Length	Accuracy	F-score	Weigh. F-score		Length	Accuracy	F-score	Weigh. F-score
RF	10	70.09	59.36	68.29	RF	10	60.07	48.65	58.12	RF	10	68.85	63.57	69.18
SVM	10	69.1	59.38	67.62	SVM	10	62.06	50.8	60.53	SVM	10	65.71	59.32	65.82
RF	25	71.02	60.04	69.34	RF	25	64.2	52.9	62.09	RF	25	70.01	64.81	70.49
SVM	25	71.76	62.50	70.59	SVM	25	65.45	54.6	64.22	SVM	25	68.03	61.66	68.12
RF	50	73.81	58.65	71.91	RF	50	67.62	50.73	65.54	RF	50	72.14	63.54	72.61
SVM	50	73.64	63.50	73.07	SVM	50	65.93	52.12	65.3	SVM	50	71.62	63.27	71.72

Results using cross-validation

Results for leave-one-subject-out

Results on balanced dataset



Classification

- best accuracy 73.6% (SVM, 50 frames)
- better results for longer segments
- differences between SVM and RF are minor
- accuracy drops with leave-one-subject-out to 67.6%
- when balancing the dataset, the results are better only for RF (4.6% on average)

	Speaking	Intaking	Chewing	Smiling	Total
10 frames	7594	1920	5273	1741	16528
25 frames	2636	622	1978	621	5857
50 frames	987	158	882	253	2280

Discussion and next steps

The first attempt to classify commensal activities from a video:

- standard video processing and machine learning techniques
- low-resolution videos of naturalistic interactions

Future extensions:

- other modalities, e.g., hands, gaze
- temporal features
- multi-label classification
- fine-grained activity recognition