arXiv:2506.21429v1 [cs.LG] 26 Jun 2025

# Deception Detection in Dyadic Exchanges Using Multimodal Machine Learning: A Study on a Swedish Cohort

Franco Rugolon[1], Thomas Jack Samuels[2], Stephan Hau[2], and Lennart Högman[2]

[1]Department of Computer and Systems Sciences

[2]Department of Psychology

Stockholm University

Stockholm

Sweden

**Author Note**

Franco Rugolon ⓘD https://orcid.org/0000-0002-7693-0576 Thomas Jack Samuels
ⓘD https://orcid.org/0009-0008-7181-8666

Correspondence concerning this article should be addressed to Thomas Jack
Samuels, Department of Psychology, Stockholm University, Albanovägen 12, Stockholm,
Sweden 106 91. E-mail: thomas.samuels@su.se

## Abstract

This study investigates the efficacy of using multimodal machine learning techniques to detect deception in dyadic interactions, focusing on the integration of data from both the deceiver and the deceived. We compare early and late fusion approaches, utilizing audio and video data—specifically, Action Units and gaze information—across all possible combinations of modalities and participants. Our dataset, newly collected from Swedish native speakers engaged in truth or lie scenarios on emotionally relevant topics, serves as the basis for our analysis.

The results demonstrate that incorporating both speech and facial information yields superior performance compared to single-modality approaches. Moreover, including data from both participants significantly enhances deception detection accuracy, with the best performance (71%) achieved using a late fusion strategy applied to both modalities and participants. These findings align with psychological theories suggesting differential control of facial and vocal expressions during initial interactions. As the first study of its kind on a Scandinavian cohort, this research lays the groundwork for future investigations into dyadic interactions, particularly within psychotherapy settings.

*Keywords:* Multimodal, Machine Learning, Deception, Game Theory, Dyadic Interaction, Deception Detection, Nonverbal

# Deception Detection in Dyadic Exchanges Using Multimodal Machine Learning: A Study on a Swedish Cohort

## An Introduction to Deception

Deception has perhaps been most notably defined by Vrij (2000) as "a successful or unsuccessful deliberate attempt, without forewarning, to create in another a belief which the communicator considers to be untrue." In this sense, deception can be conceptualized as a deliberate act involving either the direct manipulation of relevant information or the omission of pertinent detail (Fuller et al., 2012). The prevalence of these various forms of deception in everyday social interactions, which take place across a range of different human societies (Camara et al., 2024), seems to have been mirrored by increasing academic interest in the topic over the years (Talwar & Crossman, 2022). The finding that individuals engage in deception behavior more than once a day on average (DePaulo et al., 1996), highlights the apparent ubiquity and significance of deception in modern society, further reinforcing the need for greater understanding of this important social phenomenon.

Despite the apparent frequency of deceptive behaviors, detecting deception remains a significant challenge for most individuals, with accuracy rates reported at around 54%, only slightly above chance level (Bond & DePaulo, 2006). This difficulty is underscored by the well-documented 'veracity effect', which suggests that lies are detected less accurately than truths (Levine et al., 1999). Furthermore, related cognitive biases, such as 'truth bias', compound the difficulty of accurate deception detection, with individuals generally tending to overestimate the truthfulness of others in social interactions (Zloteanu & Vuorre, 2024). Such phenomena emphasize the complexity of deception detection from a human perspective, with inter-individual differences in the propensity to deceive potentially further complicating the task (Eskritt et al., 2022).

## Theories of Deception

Although detecting deception poses challenges for many individuals, the belief that nonverbal cues can serve as useful indicators of truthfulness or deception appears to be

particularly pervasive across various cultures and human societies (Bogaard & Meijer, 2022). This belief is rooted in the widespread assumption that nonverbal behaviors are harder to regulate, as they are believed to emerge spontaneously from underlying emotional states (Bogaard & Meijer, 2022). This concept of 'leakage', introduced by Ekman and Friesen (1969), suggests that nonverbal cues emerge as a result of an inability to suppress the oftentimes complex and intense emotional experiences associated with deception. In this framework, nonverbal cues, such as specific facial expressions, bodily poses and other related movements, can be considered as unintentional manifestations of the guilt, anxiety, shame, or even relief, which can occur during deceptive acts. Despite its long-standing influence, the leakage hypothesis (Ekman & Friesen, 1969) has come under significant scrutiny in recent years, with several notable researchers moving away from this theoretical framework (Vrij et al., 2019) in favor of alternative explanations, such as cognitive theories of deception.

In contrast to the leakage hypothesis, cognitive theories of deception primarily focus on the cognitive demands of deception, with less emphasis on the role underlying emotions in deception situations (Vrij et al., 2019). In this sense, nonverbal cues are thought to arise, not as leaked emotional responses, but rather as indicators of cognitive load. Deception is conceptualized as being cognitively taxing as individuals attempt to construct coherent fictional narratives, while simultaneously repressing their knowledge of true events (Van Der Zee et al., 2021). This interplay between maintaining an awareness of the known truth, whilst attempting to replace it with a fabricated untruth, is thought to contribute to cognitive load, which may manifest in subtle nonverbal cues, such as pupil dilation (Van Der Zee et al., 2021). Moreover, cognitive capacity is further tested, as deceivers attempt to retain a mental representation of both their thoughts and those of the individual(s) they intend to deceive (Zhou et al., 2023). Within this framework, cognitive theories acknowledge the role of receivers as active participants in deception situations, with related perspectives, such as Interpersonal Deception Theory (IDT), further

underscoring the importance of the interactive component in deception (Buller & Burgoon, 1996).

Interpersonal Deception Theory (IDT) (Buller & Burgoon, 1996) outlines that deception is a fundamentally interactive process, involving the dynamic exchange of both verbal and nonverbal information, between senders and receivers (Dunbar et al., 2020). Importantly, this theory denotes that senders of deceptive messages actively monitor the reactions of their recipients and adjust their communication based on their perceived degree of believability (Fuller et al., 2012), thereby framing deception as an inherently interpersonal phenomenon (Eskritt et al., 2022).

Related theories of deception, such as self-presentational theory (DePaulo, 1992), also focus on the indelibly social nature of deception behavior, with deceivers continually adjusting and refining their social presentation to appear truthful. Like strategic theories of deception, a focus is on how signals (such as nonverbal cues) can be used to strategically convey messages which bolster the facade of truthfulness (Vrij et al., 2019).

**Nonverbal Indicators of Deception**

One of the fundamental principles underpinning and uniting these different theories of deception is the critical role nonverbal indicators play in enabling deceivers to influence their interaction partners' responses. Unsurprisingly, research into nonverbal indicators of deception has steadily increased over the years (Colasanti et al., 2024), with a significant focus on facial, head and bodily movements, as well as prosodic features, including various nonverbal attributes of vocal communication (Burgoon et al., 2021).

For instance, previous research has provided evidence to suggest that specific facial movements, as categorized using the Facial Action Coding System (FACS) (Ekman & Friesen, 1978), such as chin raiser (AU17) and lip stretcher (AU20), were associated with deception (Shen et al., 2021; Zhou et al., 2023), whilst lip corner puller (AU12) has been previously found in receivers responding to deceptive senders (Sen et al., 2018). Other cues, relating to the eyes, such as pupil dilation (Shen et al., 2021) and gaze fixation

(Colasanti et al., 2024), have also been associated with deception behavior, whilst specific bodily movements, such as shrugging and fidgeting, have similarly been linked with deception in prior research (Matsumoto & Hwang, 2021).

Prosodic features have also been extensively studied, with research indicating that voices with 'rising intonation, less intensity at the beginning of each syllable, and slower speech rate' were associated with doubt or lying (Goupil et al., 2021). These findings align with previous research suggesting that auditory information may be even more useful than visual information in deception detection (Ahmed Khan et al., 2024).

In addition to individual nonverbal cues, considerable research has also examined how the coordination of various nonverbal cues over time, referred to as nonverbal or interactional synchrony, also appears during deceptive interactions. This synchrony has been conceptualized as a social strategy, used by deceivers to enhance their credibility. Indeed, Duran and Fusaroli (2017), found an increased rate of head movement coordination in deceptive conversations relative to honest conversations. It has been hypothesized that the reciprocation of nonverbal coordination by a receiver results in an inability to correctly identify and appraise patterns of nonverbal behavior potentially useful for deception detection (Van Der Zee et al., 2021). As such, deceivers may consciously or non-consciously control the exhibition of nonverbal cues to foster a semblance of credibility, something supported by previous research undertaken by Burgoon and Buller (1994), which demonstrated that deceivers appeared to display more restrained and tense movements relative to their truthful counterparts.

Despite the apparent centrality of nonverbal cues in deception, findings in the field appear to demonstrate greater homogeneity than would have previously been expected. This point is perhaps most famously illustrated by a meta-analysis originally conducted by DePaulo et al. (2003), which found little evidence to support the role of nonverbal cues in deception. These findings suggest that, if such a relationship does exist, it may be difficult to reliably identify and observe (Delmas et al., 2024; Vrij et al., 2019). In this sense, the

field lacks a 'Pinocchio's nose' (Luke, 2019; Vrij, 2004), with some previous research undermined by the use of imprecise measures and potentially inadequate coding systems, leading to mixed findings that offer little support for any single nonverbal feature (Delmas et al., 2024). Such inconsistencies are unsurprising given the influence of broader contextual factors, such as cultural differences, which appear to have a marked impact on the exhibition of nonverbal behaviors (Matsumoto & Hwang, 2021). Additionally, some studies have also failed to account for potential confounding variables, such as the distinction made between low vs high-stakes lies, which may further cloud the picture (Matsumoto & Wilson, 2023). Given these challenges, focusing solely on investigating individual nonverbal cues, or even modalities, may be insufficient (Gupta et al., 2019), Instead, examining 'clusters' or 'constellations' of cues across multiple modalities may offer a more holistic and nuanced understanding of deception (Hartwig & Bond Jr., 2014; Matsumoto & Hwang, 2021).

## What is a modality?

According to Baltrušaitis et al. (2018), a modality is each of the ways in which a phenomenon happens or is experienced. In particular, a specific group of modalities can be associated with the sensory modalities through which humans experience the world, such as sight, hearing, smell, taste, and tact (Baltrušaitis et al., 2018). In this paper, we will focus on visual and vocal signals: visual signals will represent the facial expressions and gaze direction information of the participants, while the vocal signals will represent various parameters associated with the voice of the participants. In addition to modalities, it is important to consider the source of the data, since that will determine, at least in part, the format that the data will have to be analyzed: one of the most relevant distinctions can be made between structured and unstructured data (Zhang et al., 2018). Structured data often originates from forms or sensors, and it refers to information that is organized in a predefined format, or structure: usually it is organized in tables with rows and columns, such as in spreadsheets or databases. Each column has a specific meaning, and each row represents an instance or observation, such as a different subject or point in time. Because

of this organization, structured data is straightforward for computers to store, search, and analyze. In contrast, unstructured data does not follow a clear, predefined format, but can be varies in its contents: it can represent images, text, and other kinds of data. This kind of data is often rich and informative but also more complex to process.

**Multimodal Machine Learning**

To this end, multimodal Machine Learning (ML) can be employed. Multimodal ML is a subfield of ML in which data from multiple modalities is integrated and processed, and different algorithms are applied to it (Baltrušaitis et al., 2018). The goal in Multimodal ML is to develop models that can make predictions based on data from the different modalities, achieving a more holistic representation of real-world phenomena.

The data from each modality can be integrated at different points of the ML process, but the three main strategies are early fusion, late fusion, and joint/intermediate fusion (Baltrušaitis et al., 2018; Xu et al., 2021).

*Early fusion*

In early fusion, the data from the different modalities is integrated before applying any ML technique to it, and for this reason, this fusion technique is also called data-level fusion. Since data from different modalities is usually represented in different ways, each modality is preprocessed independently, with appropriate techniques (e.g., normalization, standardization, tokenization). After this passage, the data from the different modalities is brought to a common representation. The techniques applied in this step are modality-dependent and vary greatly according to the different modalities that are considered in each experiment. Once all the data is in a common representation, it can be concatenated, and a single ML model can be applied to the concatenated data to generate predictions (Gadzicki et al., 2020).

*Late fusion*

Late fusion is also called decision-level fusion: in this process data from the different modalities is preprocessed independently, with appropriate techniques, and a different

model is applied to each modality: for example, one model might analyze audio signals, while another might analyze visual data. These models can be as complex or as simple as needed. Once each model has produced the predictions for the data in each modality, these predictions are used to train a meta-model: an ML model that learns to combine the outputs (the predictions) of several other models(Gadzicki et al., 2020). Meta-models can be seen as a judge that, after listening to the opinions of several experts (the initial models), takes a decision keeping all of these opinions in consideration. In this analogy, the experts do not communicate among them: the initial models are trained only on the data from their specific modality. At the same time, the judge does not see (and would not be able to interpret) the initial data, but only hears the opinions of the experts, and bases their decision only on what the experts tell them. Meta-models, contrarily to the models that are trained on each modality, have to be simple models to reduce the risk of overfitting: they can be as simple as a majority vote or an average of the decisions of the single models. The late fusion predictions are the predictions generated by the meta-model.

### *Joint fusion*

In joint fusion, the data from each modality is preprocessed independently, and it is fed into a Neural Network (NN), where the fusion step will happen in between the data point and the output point (Gadzicki et al., 2020). NNs have the ability to process data from different modalities, depending on the kind of layer that is applied to this data (Kahou et al., 2013). Moreover, NNs have the ability to highlight non-linear relations in the data and to use these relations to create their predictions, if the correct activation function is used (Grossberg, 1988). In joint fusion, data from different modalities is fed to different arms of an NN, which applying the appropriate layer to further process the data and bring it to a common representation. Once the data is in a common representation, the representations of the data from each modality can be fused using different techniques, and a common prediction can be outputted by the model (Baltrušaitis et al., 2018).

Each fusion strategy has different advantages and disadvantages: early fusion

models have the possibility to learn relationships between data from different modalities, but due to the need to have all the data in the same modality (Gadzicki et al., 2020), some of the characteristics of the data will inevitably be lost, while late fusion, having different models for each modality, is able to exploit all of their characteristics, but since the meta-model is only exposed to the decisions of each base model, the cross-correlations between the different modalities are lost (Gadzicki et al., 2020). Joint fusion preserves the ability of the model to learn from each modality independently, and at the same time it preserves, to an extent, the interactions between the different modalities (Baltrušaitis et al., 2018), but the models that originate from this fusion strategy are more complex, needing expert crafting, and the risk of overfitting is generally higher. For this reason, these models are not appropriate for tasks with low amounts of data.

**A Multimodal Understanding of Dyadic Interaction and Deception Detection**

**ML in deception detection**

Deception detection is a traditionally hard problem, which has been faced using both psychological methods, like inquiring about additional details or asking questions to further investigate a statement, or more advanced ones, like polygraph tests (Prome et al., 2024). More recently, ML has been proposed as a tool to help in identifying deception, using behavioral, physiological, and linguistic data (Tang et al., 2018).

ML approaches can be classified depending on the modality (or modalities) that they consider in order to predict whether a subject is being deceptive or not: while in past years the use of a single modality was widespread, in recent years approaches that use multiple modalities are gaining traction (Constâncio et al., 2023): studies using single modalities, such as facial information (Ahmed Khan et al., 2024; Cardaioli et al., 2022; Khan et al., 2021), textual information (Loconte et al., 2023) or body pose information (Poppe et al., 2024) are becoming rarer, with studies using a combination of modalities, such as facial, audio and textual information (Gogate et al., 2017), or facial information and pulse rates (Tsuchiya et al., 2023). Since ML models ability to detect deception is

usually linked to the specific scenario in which the models have been trained, it's important to take these scenarios into consideration: deception can be differentiated according to the importance that it has for the person performing it, in low stakes, medium stakes, and high stakes, where the emotional involvement of the deceiver grows with the stakes (Ahmed Khan et al., 2024). This distinction is important because the ML models do not, usually, detect deception per se, but rather the behavioural changes that originate from a heightened emotional state to infer deception (Constâncio et al., 2023). In this scenario, synthetic datasets, which are commonly used to train ML models in deception detection, and which also usually have the highest number of available samples, usually lack the high stake component that real life datasets can carry.

Multiple studies use databases created in unnatural conditions, in which the subjects do not interact with other humans, to train their models, such as the SMIC dataset (Li et al., 2013), CASME (Yan et al., 2013), CASME II (Yan et al., 2014), SASE-FE (Wan et al., 2017), while other studies used datasets from real life situations, such as the real-life trial dataset (Pérez-Rosas et al., 2015), but rich, multimodal, annotated real-life datasets are rare. Multiple studies (Ahmed Khan et al., 2024; Constâncio et al., 2023; Khan et al., 2021; Prome et al., 2024) cite the importance of micro-expressions (Ekman, 2009), fleeting variations of one's facial expression as possible indicators of deception. To capture such expressions, it's important to consider the temporal component of the data, analyzing the changes in the facial expressions of the subjects, but there are also papers (Cardaioli et al., 2022) which fail to do so, employing techniques such as averaging the activation of the different AUs over a period of time, losing the temporal dimension of the data and the ability to capture rapid changes.

**Problem definition**

The general task in this experiment was to determine whether the sender was lying or telling the truth during their interaction with a passive receiver. We collected audio and video recordings of participants in the deception experiment that we designed, and we

extracted Facial Action Units (FAC) activation levels and eye movements from each frame of each video recording, and the Geneva minimalistic acoustic parameter set (GeMAPS) (Eyben et al., 2015) from each timepoint in the acoustic recordings. We then aggregated the recordings of all the dyads for each modality into a collection of files containing the features extracted from the video recordings of the sender in all the dyads, and into a separate files collection containing the same features for both participants. We then repeated the process with the features extracted from the audio recordings, to obtain files containing the data extracted from the recordings of the sender, and a different set of files containing the data extracted from the recordings of both participants. We then analyzed these files using three different approaches: the unimodal approach, the early fusion approach, and the late fusion approach. Given the low amount of data available in our dataset, it was not possible to utilize the joint fusion approach, since the risk of overfitting, with sucha complex model, would have been too high.

In the unimodal approach, we considered either the video or audio recordings of the active participant or both participants in all the dyads.

Our objective was to train a classifier using leave-one-out cross-validation (Han et al., 2022) at the dyad level, which, when applied separately to either facial or vocal features extracted from recordings of the sender alone or from both participants, would yield a binary prediction ("lie" or "truth") for each dyad. We then averaged the results of the leave-one-out cross-validation to get the final precision and recall metrics per each class, and the final accuracy metric.

In the early fusion approach, we considered both video and audio recordings of the same dyads involved in the unimodal approach. Since the audio and video recordings had different recording frequencies, our goal was to bring them to the same number of timesteps before learning a classifier with the same procedure used in the unimodal approach. In particular, we experimented with two different methods to bring the extracted features to the same feature space, using either Piecewise Aggregate Approximation (PAA) (Keogh

et al., 2001) or Symbolic Aggregate approXimation (SAX) (Lin et al., 2007) to bring the recordings of both modalities to have the same amount of timesteps. We applied two commonly used timeseries summarization techniques, PAA and SAX. These methods both apply a segmentation technique, calculating the average value of the timeseries for the each of the resulting segments. On top of this segmentation and averaging technique, SAX applies a discretization procedure to reduce the complexity of the values for each timestep.

In the Late fusion approach we used a two-step process, learning the base classifiers with the same procedure as for the unimodal approach, but, in addition, we trained a decision tree on the predictions of the base classifiers for the dyads that were not the target of the classification during the leave-one-out cross-validation procedure, and we used this decision tree to produce the final predictions of "lie" or "truth" for each dyad.

## Method

In our approach, we analyzed both visual (FAC and gaze information extracted using OpenFace) and audio (GEMAPS features extracted using OpenSmile) information from participants who were assigned the task of lying or telling the truth.

### Participants

44 individuals participated in this experimental study, with an age range from 18 to 67 years old (M = 29.52, SD = 11.2). The sample contained 28 men, 14 women and 2 identifying as 'other'. All participants were native Swedish speakers and were at least 18 years old at the time of participation. Additional details relating to the sample for each experimental condition can be found in 1.

Participants were recruited through the use of several mediums, including advertisements posted on digital learning platforms and bulletin boards located throughout Stockholm University, as well as on the research recruitment platform Accindi (accindi.se). All participants received financial compensation in the form of a 100 SEK gift card for a Swedish supermarket chain.

**Procedure**

The study employed a between-subjects design, with participants randomly assigned to experimental dyads, in one of two experimental conditions: truthfulness or deception. Within each dyad participants were assigned the role of either a) sender or b) responder. In the truth-telling condition, person A, the sender, was advised to recall actual experiences, completing seven sentences, such as 'One of my most cherished memories is when...', and 'A major problem I've had was when...'. In the deception condition, the sender was asked to complete the same sentences but with fabricated responses which did not reflect their true experiences, beliefs or opinions. Participants adopting the role of the sender in both conditions were advised to rehearse their statements before attending the experimental session, in an attempt to minimize any differences potentially occurring due to familiarity with the experience or opinion (DePaulo et al., 2003). Participants assigned to the role of responder were tasked to ask follow-up questions after each sentence was read, in an attempt to enter into a meaningful conversation for each statement.

Before the experiment began, the experiment leaders ensured that all participants had been appropriately prepared and provided consent to participate. Participants were then instructed to remove any items which obscured their face or head (e.g. caps or scarves) and remove any oral impediments (e.g. nicotine pouches or food items, such as chewing gum), before being seated in front of their respective cameras. Each participant was recorded with a separate camera which faced directly toward them, as shown in Picture 1.

Participants were then introduced to one another and instructed to engage in conversation for approximately one minute, during which the experimental leaders ensured that the recording equipment was functioning correctly and that the cameras and microphone were capturing the appropriate audiovisual input. Once these preliminary checks were complete, the experiment leader instructed the participants to begin, with the sender reading their first pre-prepared statement. The responder then asked a follow-up question, with conversation unfolding naturally after this point. After two minutes, the

experiment leader instructed the sender to move on to their second statement, and the process was repeated. In total, there were seven rounds of conversation, one per pre-prepared statement, each lasting two minutes, with the total interaction time of being around 15 minutes.

At the end of the interaction, participants were partitioned by a screen and asked to complete a set of questionnaires, rating the quality of their interaction and the perceived truthfulness of their interaction partner. After completing the questionnaires, participants were reunited and debriefed. They were then reimbursed with gift cards and advised to contact the lead researcher with any follow-up questions.

**Materials**

Participants were required to complete an informed consent form before participating in the experiment. The form provided a brief description of the structure of the experiment, as well as relevant information relating to data management and storage.

Post-experimental questionnaires were used to gather additional data. These included a perceived trustworthiness scale, a three-item scale adapted and translated (from English to Swedish) from the scale originally used by Dunbar et al. (2019) ($\alpha = .81$). This scale assessed participants' perceptions of the truthfulness of their interaction partners, with responses recorded on a 7-point Likert scale. Participants also completed a subset of three items taken from the original scale presented by Bernieri et al. (1996) ($\alpha = .83$), evaluating the overall experience of the interaction. The items asked participants to rate how 'engaging', 'cooperative', and 'awkward' they perceived the interaction to be, with the scale translated into Swedish.

**ML methods**

After the experiments were conducted, we collected the video and audio recordings of the sessions, and we used two open source programs, OpenFace (Baltrušaitis et al., 2016) and OpenSmile (Eyben et al., 2010), to automatically extract features from them. OpenFace allows us to extract AU and information about the eye position and gaze angles,

collectively called facial features, while we use OpenSmile with the GeMAPS (Eyben et al., 2015) configuration to extract acoustic features.

The facial features are extracted for each frame of the recording, and are concatenated on a temporal axis to create time series for each feature. In the same way, the acoustic features are extracted from the recordings, extracting a hundred time points per second, and concatenated along the temporal axis to create time series for each feature.

After the concatenation, the time series relative to each modality for each participant were transformed into two-dimensional arrays (features × timesteps). The timeseries of the different participants were then aggregated for each modality, to create a three-dimensional array (participants × features × timesteps) that was modality specific. The features contained in these three-dimensional arrays were then standardized feature-wise and participant-wise using z-score normalization (Han et al., 2022), to reduce the impact of the use of different measuring units on the classification process.

We divided the classification process, as described in the problem formulation subsection, in three phases: unimodal, early fusion, and late fusion. In the unimodal phase, we trained three classifiers, namely Rocket (Dempster et al., 2020), Canonical Interval Forest (Middlehurst et al., 2020) and Z-Time (Lee et al., 2024) on the data from each modality. We repeated the process using data from the sender only and from both participants. To reduce the risk of overfitting, we employed leave-one-out cross-validation, using the data from $n-1$ participants or dyads as our training set and the data from the remaining participant or dyad as our test set, and repeating the process $n$ times, until the data from each participant or dyad was used as test data exactly once. We repeated this process for both modalities, therefore using both data from the facial features and the acoustic features, separately, to predict the binary label to be attributed to the data.

In the early fusion phase, as described in the introduction, the data from different modalities has to be brought to a common space, concatenated, and then a classifier can be applied to the concatenated data. Since both modalities were made of timeseries, albeit

with different dimensionalities ($22 dyads \times 19739 timesteps \times 25 features$ for the facial features of each participant, and $22 dyads \times 66105 timesteps \times 24 features$ for the acoustic features of each participant) we needed to bring the different timeseries to the same number of *timesteps*. To do so, we compared two different techniques: PAA (Keogh et al., 2001) and SAX (Lin et al., 2007). Given an original timeseries $t$ of $N$ timesteps, both techniques allow us to choose a desired number $N'$ of frames, where $N' \leq N$. They then divide the original timeseries in $N'$ fragments of the same duration, to then calculate the average value of each fragment. This generates a new timeseries $t'$ with the desired duration $N'$. In addition to this, SAX also applies a discretization over the values $V$ of the timeseries $t$, dividing the interval between the minimum value $V_{min}$ and the maximum value $V_{max}$ in $V'$ desired values, and creating a dictionary of $V'$ values. The original values $V$ are then mapped to the new $V'$ values, reducing the complexity of the timeseries.

In our early fusion approach, we chose 10000 timesteps for both PAA and SAX, but, without loss of generality, any $N'$ which is less or equal to the duration of both timeseries can be chosen. After applying PAA and SAX, the facial features were represented by a three-dimensional array of $22 dyads \times 10000 timesteps \times 25 features$, or $50 features$ if considering the data of both participants, and the acoustic features by a three-dimensional array of $22 dyads \times 10000 timesteps \times 24 features$, or $48 features$ if considering the data of both participants, allowing us to concatenate the features at each time point. After concatenation, a new three-dimensional array of dimensions $22 dyads \times 10000 timesteps \times F features$ was generated. The actual number of features depended on the use of data from just the sender or both participants for each modality.

Following the same training and testing procedure used in the unimodal phase, we trained a Rocket binary classifier across the entire participant dataset, systematically evaluating all possible combinations of modalities, and using data from either only the sender or both participants for each modality. Furthermore, we used either PAA and SAX to bring the timeseries from different modalities to the same feature space. Again, we used

leave-one-out cross-validation to reduce the risk of overfitting.

In the late fusion phase, as explained in the introduction, the data from each modality is fed to a classifier specific for that modality, and then the predictions of each classifier are either aggregated or fed to a meta-classifier. In our case, we used the probabilistic predictions generated by Rocket on the single modalities, and we used those to train a decision tree and produce the final multimodal predictions.

To do this, we, once more, used leave-one-out cross-validation to reduce the risk of overfitting. To further reduce the risk of overfitting, we kept three as the maximum depth of the tree.

The primary aim of this study was to explore the impact of different modalities, combinations of modalities, and of the dyadic aspect (i.e., using data from the sender, or from the sender and the receiver) on deception detection performance. As such, our focus was on comparative analysis across experimental conditions rather than on maximizing the absolute predictive performance of any single classifier. Consequently, we did not perform hyperparameter tuning for our models. All base classifiers were trained using their default parameter settings to ensure consistency and comparability across conditions, and the decision tree meta-classifier was trained using a reduced maximum depth to reduce the risk of overfitting and to isolate the effects of modality and data origin.

Since there was only a minimal imbalance in the class distribution in the data, we deemed that accuracy was a good performance measure for the classifiers. To give a better idea of the performance of each classifier on each class, we still reported precision and recall per each class for all the classifiers [1].

## Results

Table 2 summarizes the performance of the different classifiers on the various combinations of modalities and participants for this study.

Looking at the results for the unimodal classifiers, we can see that in both single

---

[1] All the code needed to run these experiments is made available on GitHub.

modalities, employing data from both participants either does not change the accuracy of the classifier in two cases, or improves it in the other four cases. For facial expressions, in particular, there is an increase in accuracy using data from both participants, but even the best results are only marginally better than random level. We can also see that sensitivity always improves, while the results for specificity are mixed.

Looking at the results for the early fusion, we can see that applying SAX instead of PAA improves the classification accuracy and the sensitivity in three out of four modality combinations, while using data from the facial expressions of both participants improved the classification accuracy in four out of eight cases, decreasing it in the remaining half, compared to using data from the facial expressions of just the sender, and using data from the voice of both participants improved the classification accuracy in two cases out of eight, decreased it in two other cases, and left it unchanged in the remaining half.

Looking at the results for the late fusion, we can see that this way to perform data fusion generated the two best classification accuracies of the whole experiment, in both cases using data from the facial expressions of the sender. The best accuracy overall was obtained by the decision tree trained on the output of the Rocked classifiers trained on the voices of both participants, and the facial expressions of the sender. We can also see that this classifier shows a perfect precision for class 0 (lying dyads) and a perfect recall for class 1 (truthful dyads), meaning that this classifier was able to spot all the dyads in which the sender was telling the truth, and never misclassified a dyad in which the sender was lying.

Finally, looking at the results for human classification, the receivers showcased a lower accuracy than unimodal methods trained on prosodic data, but a better accuracy than unimodal methods trained on facial expressions. The early fusion approach showed better results than the receivers in three experimental combinations out of eight, in all three cases when using data from both participants in at least a modality and when using SAX for dimensionality reduction. Late fusion showed better results than the receivers in half of the experimental settings.

## Discussion

Overall, the present findings indicate that adopting a multimodal approach to deception detection can be advantageous, as shown by a trend in recent research (Gogate et al., 2017; Tsuchiya et al., 2023) especially when data obtained from different nonverbal modalities are integrated using late fusion. One particularity of this study is that all the data comes from non-verbal modalities, decreasing the need for extremely complex models such as Large Language Models, which would be less interpetable and would require more computational power to run. Additionally, the data for this study can be collected in a relatively non-invasive way, with cameras and microphones that can be located away from the subjects, reducing the likelihood for the recording procedure to influence the behavior of the subjects. While the use of pulse rates, as was done by Tsuchiya et al. (2023) or skin conductivity (Ströfer et al., 2015) can provide information on physiological signals that are more difficult to control for the sender, it can also disrupt the normal flow of the interaction. Since, in this study, we wanted to focus on the dyadic aspect of the interaction, we chose to only use audio and video information. Moreover, even in the video information, we were limited in the type of video information by the positioning of the cameras: as it can be seen in Picture 2, our cameras only captured the face of the participants and, partially, the upper part of the shoulders, rendering the extraction of body positions impossible.

Even with these limitations, Within the current study, the highest accuracy (0.77) was significantly higher than what humans usually are able to achieve.

One particularly notable result from our experiments is the perfect precision achieved for the Lie class in our best-performing model. Precision, in this context, measures the proportion of instances classified as lies that are actually lies. A perfect precision score of 1.00 means that every time the model labeled a statement as a lie, it was correct: there were no false accusations. This is especially important in deception detection tasks, where false positives can carry serious ethical and practical consequences. In forensic, clinical, or even interpersonal settings, wrongly labeling someone as deceptive can

damage trust, reputation, or bring legal outcomes. Therefore, a system that avoids such errors is particularly valuable.

This result was achieved using a Decision Tree model as a meta-model, trained on the outcomes from two Rocket classifiers, which, in turn, were trained on the data from the single modalities, audio and video data. This approach suggests that such late fusion techniques may enable the single modality classifiers to learn unique, complementary information obtained from each modality, which, when combined at a later decision stage, can lead to greater classification accuracy. This advantage seems to outweigh the lack of interactions among the single modalities, which might not be determinant for a good classification result, as this approach does not preserve them. Additionally, the late fusion approach allowed us to maintain the data from each modality in its original format, avoiding the issue of having to conduct dimensionality reduction, which we had to face in the case of early fusion.

In contrast, the use of individual modalities within this study yielded lower accuracy scores. The use of facial expression data alone produced lacklustre returns, with accuracy rates ranging from 0.36 - 0.55, indicating that the use of such information alone may not be sufficient for accurate classification of deception in such contexts. Indeed, data relating to the voice proved to produce marginally better results, although these also fell short of the scores produced by the late fusion multimodal model.

To use early fusion, it was fundamental to reduce the timeseries from different modalities to a common, lower dimensionality. This process was carried out with two different methods, using the same classifier, to compare the impact of the dimensionality reduction techniques. In three out of four cases where the two techniques, PAA and SAX, were compared, the results of the classifier trained on the data generated by SAX were superior than those of the results of the classifier trained on data transformed with PAA. Moreover, the results of the classifiers trained on early fusion data were worse than the results of the classifiers trained using only the voice of both participants, showing that the

direct integration of data from both modalities is not necessarily beneficial during the classification process.

The present results may be tentatively interpreted through the frame of existing theories of deception. The relatively poor performance of classifiers trained solely on facial expression datav raises important questions about the robustness of leakage theory's emphasis on the face as the most primary channel through which involuntary emotional cues may 'leak' (Ekman & Friesen, 1969). If such leakage were a consistent or detectable phenomenon, we may expect to see a markedly better performance in terms of classification accuracy for this individual modality. Such results align with previous findings (DePaulo et al., 2003), and may even suggest that such features may be more subject to voluntary control than previously assumed, an idea originally outlined by Burgoon and Buller (1994). From the cognitive perspective, an increase in cognitive load may not be uniformly distributed across modalities, but instead disproportionately affect those modalities which place high demands on temporal and semantic coordination, such as speech (DePaulo et al., 2003). The stronger performance of classifiers trained on prosodic features aligns with this perspective, suggesting this modality may be more susceptible to disruption under cognitive load. Indeed, the improved performance of the late fusion model over the early fusion model may further support the notion that deception cues are not uniformly distributed across modalities, but instead emerge unevenly throughout the dynamic and shifting flow of a nonverbal exchange. This pattern aligns with Interpersonal Deception Theory (IDT) (Buller & Burgoon, 1996), which posits that deception is a multimodal, dynamic process evolving across the course of an interaction.

Such findings appear to be consistent with previous research which has outlined the benefits of using a multimodal approach in the detection of deception, as well as the potential pitfalls of unimodal classification (Camara et al., 2024). For example, previous research on the analysis of single modalities, such as facial expressions, has found evidence to suggest that such displays may be highly idiosyncratic, with previously unexamined

factors, such as personality, potentially playing an important role in explaining underlying variance apparently present in these nonverbal behaviors (Zhou et al., 2023). Furthermore, the examination of single modalities, often presented as single nonverbal channels, may not fully capture the wealth of information being communicated across multiple nonverbal channels simultaneously, leading to an inability for models to correctly classify and thereby detect deception (Matsumoto & Hwang, 2021). Furthermore, the performance comparison between human and ML deception detection confirms the promise of ML in this delicate domain, which might be due to the absence of truth biases showcased, in previous research, by human participants (Zloteanu & Vuorre, 2024). Although deception can be better detected from multiple cues than a single cue, the strongest cue contributes a lot.

However, multimodal deception detection is not a panacea, and the application of such approaches should be carefully considered, especially when the considering the possiblity of Type I errors in not correcting for multiple comparisons (Luke, 2019). The current approach employs a leave-one-out cross-validation in an attempt to reduce overfitting on the current sample. In future, authors may seek to further reduce the likelihood of Type I errors by addressing these issues earlier in the analytical workflow.

This study constitutes an illustration of a method that can be implemented within the burgeoning research field of multimodal deception detection and multichannel nonverbal research in general (Matsumoto & Wilson, 2023). As such current study utilizes a sample of 22 dyads, comprising of 44 individuals who were previously unknown to one another. Whilst previous research has featured samples of a similar size with sufficient statistical power to detect effects (Zhou et al., 2023), we are mindful to not overstate our findings given previous criticism of sample sizes in this field (Luke, 2019). Whilst leave-one-out cross-validation can reduce the likelihood of overfitting in smaller samples, it may still not appropriately address more foundational concerns such as the representativeness of the sample, something which future research could seek to address.

Furthermore, whilst the use of a laboratory setting enabled us to establish a stable

ground-truth for this study, through the use of experimental manipulation (Dunbar et al., 2023), the generalizability of our findings may be limited by the use of an artificial setting which could be viewed as a somewhat unnatural context for the study of human interaction patterns (Burgoon et al., 2021). In addition, although the pairing of strangers helped control for potential familiarity effects in nonverbal interaction patterns (Brandt et al., 1980), it may also have reduced the perceived stakes of engaging in deception, with interactants unlikely to engage in future interactions, thereby lessening the potential consequences of being caught in a lie when compared to real-world scenarios (Dunbar et al., 2019). Such 'low-stakes' have long been assumed to result in fewer nonverbal deceit cues when compared to 'high-stakes' scenarios which could induce stronger emotional responses and result in greater cognitive load (Vrij et al., 2019). However, given that participants were asked to fabricate personal experiences or opinions, it remains an open question whether the potentially emotionally-charged nature of these statements could have, to some degree, heightened the stakes of deception within our study. Additionally, our experimental method seeks to instigate a relatively narrow form of deception, namely falsified experiences or beliefs, and does not account for broader definitions of deception which may include related phenomena such as gas-lighting, feigned ignorance and self-deception. However, adopting this simplified definition of deception did enable us to establish a clear ground truth and verify the effectiveness of the experimental manipulation through post-interaction self-reports.

**Ethical Considerations**

The deployment of ML techniques for deception detection carries a number of important ethical implications, especially when applied in sensitive or high-stakes domains such as legal proceedings, border control, clinical settings, or employment screening (Oravec, 2022). While such systems offer the promise of greater consistency and objectivity compared to human judgment, their use must be carefully evaluated in light of both technical limitations and broader societal concerns.

Under the General Data Protection Regulation (GDPR) in the European Union, individuals are entitled to a "right to explanation" when subjected to decisions made by automated systems that significantly affect them (Goodman & Flaxman, 2017). Since many high-performing ML models, such as those used in our study, are difficult to interpret, their use in legally consequential contexts would require robust explainability techniques. This poses both a technical and legal challenge, as explanations must be meaningful to end-users and stakeholders (e.g., defendants, lawyers, clinicians), and not just technically accurate.

Another potential issue that needs to be considered is that, although our approach outperforms untrained human interlocutors in detecting deception, it still falls slightly short of what trained students were able to achieve in one specific study (77.5% accuracy) (Levine et al., 2014), and it is, in general, far from perfect. This creates a risk of misclassification, particularly false negatives (failing to detect a lie) and false positives (incorrectly labeling truthful behavior as deceptive). The consequences of such errors can be significant, ranging from wrongful suspicion or punishment in legal contexts to breakdowns in trust in therapeutic or interpersonal settings. Even with perfect precision in detecting lies (as achieved in our best model), the broader risk profile must be carefully managed before deployment.

An additional concern relates to the potential for bias in the training data or in the model's interpretation of multimodal signals such as facial expressions, tone of voice, or gaze behavior. These cues can vary significantly across cultures, neurodiverse populations, or individuals with different emotional baselines, which could lead to systematic disadvantages for certain groups. This underscores the need for inclusive datasets, fairness audits, and potentially customized models calibrated for specific populations.

Deception detection systems typically rely on rich and sensitive data, including audio, video, and behavioral cues. The use of such data raises privacy concerns, especially in contexts where individuals may not be fully informed or have not consented to being analyzed. Even in research settings, clear ethical guidelines must be followed to protect

participants' rights, and any operational use must be aligned with principles of informed consent, data minimization, and secure storage.

Lastly, there is a risk that institutional trust in ML systems could lead to over-reliance on these tools, with human decision-makers deferring to algorithmic outputs even when their limitations are known. This can create a false sense of objectivity, where decisions are viewed as "data-driven" despite the system's probabilistic nature and vulnerability to error.

**Future Directions**

Future deception detection efforts using multimodal classification approaches, such as the one outlined here, should prioritize generating and maintaining large open-source multimodal datasets for benchmarking and reproducibility. Many research groups have begun the arduous and time-consuming process of systematically collecting multimodal datasets, which is commendable; however, much larger and more diverse samples are required to ensure that classification models are appropriately trained and tested (Delmas et al., 2024), especially if the models are expected to generalize from a specific task to a more general ability to detect deception in a broad array of different situations.

Ideally, such datasets would also account for the distinction between low-stakes lies, which are more easily collected in experimental settings, and high-stakes lies, which are often the focus of practical deception detection efforts employed in law enforcement and criminal justice contexts (Camara et al., 2024). The concern is that, while low-stakes datasets may potentially be more accessible, models trained on such data may not generalize well to high-stakes deception scenarios, where non-verbal behavioral exhibition may differ meaningfully.

Moreover, future research could also seek to account for individual differences in deception. Previous research has highlighted the seemingly idiosyncratic nature of nonverbal cues in deception contexts, with personality traits, such as extroversion, associated with differing nonverbal behavior exhibition (Zhou et al., 2023). Furthermore,

factors such as childhood development and early social experiences could also be considered, as early exposure to dishonesty during childhood has been linked to a greater tendency to engage in dishonest behavior later in life (Talwar & Crossman, 2022), potentially further complicating the task of detecting deception across individuals. The use of within-participant experimental designs could enable future research to better account for nonverbal behaviors across truthful and deceptive contexts, although the implementation of such designs may be tricky, as participants could be influenced by order effects and may have difficulty switching between conditions, artificially generating cognitive load which may not accurately reflect the experience of deceivers in real life.

## Conclusion

Perhaps there is indeed no universal 'Pinocchio's nose'—a single, unmistakable cue signaling deception. Challenges such as the inherent granularity of high-resolution multimodal data, which can introduce significant amounts of noise, as well as individual and cultural differences in expressive behaviors and communication styles, may inherently limit the reliability of deception detection systems (Vrij et al., 2019). Despite these complexities, this study represents a robust advancement toward more reliable automated detection of deception: by leveraging a novel dataset collected from an underrepresented cohort (native Swedish speakers), this research has systematically compared unimodal and multimodal methods, demonstrating the advantages of multimodal analysis and highlighting the increased predictive power achieved by integrating data from both participants in dyadic interactions. Through this comprehensive approach, the study underscores the importance and effectiveness of capturing dyadic dynamics rather than relying solely on isolated individual cues. Thus, while perfect detection of deception may remain elusive, this work decisively advances our understanding and methodological capabilities in deception research.

## References

Ahmed Khan, H. U. D., Bajwa, U. I., Ratyal, N. I., Zhang, F., & Anwar, M. W. (2024). Deception detection in videos using the facial action coding system. *Multimedia Tools and Applications*, 1–15.

Baltrušaitis, T., Ahuja, C., & Morency, L.-P. (2018). Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, *41*(2), 423–443.

Baltrušaitis, T., Robinson, P., & Morency, L.-P. (2016). Openface: An open source facial behavior analysis toolkit. *2016 IEEE winter conference on applications of computer vision (WACV)*, 1–10.

Bernieri, F. J., Gillis, J. S., Davis, J. M., & Grahe, J. E. (1996). Dyad rapport and the accuracy of its judgment across situations: A lens model analysis [Place: US Publisher: American Psychological Association]. *Journal of Personality and Social Psychology*, *71*(1), 110–129. https://doi.org/10.1037/0022-3514.71.1.110

Bogaard, G., & Meijer, E. H. (2022). No evidence that instructions to ignore nonverbal cues improve deception detection accuracy. *Applied Cognitive Psychology*, *36*(3), 636–647. https://doi.org/10.1002/acp.3950

Bond, C. F., Jr., & DePaulo, B. M. (2006). Accuracy of Deception Judgments: Appendix A. *Personality and Social Psychology Review*, *10*(3). https://doi.org/10.1207/s15327957pspr1003_2A

Brandt, D. R., Miller, G. R., & Hocking, J. E. (1980). The truth-deception attribution: Effects of familiarity on the ability of observers to detect deception. *Human Communication Research*, *6*(2), 99–110. https://doi.org/10.1111/j.1468-2958.1980.tb00130.x

Buller, D. B., & Burgoon, J. K. (1996). Interpersonal Deception Theory. *Communication Theory*, *6*(3), 203–242. https://doi.org/10.1111/j.1468-2885.1996.tb00127.x

Burgoon, J. K., & Buller, D. B. (1994). Interpersonal deception: III. Effects of deceit on
    perceived communication and nonverbal behavior dynamics. *Journal of Nonverbal
    Behavior*, *18*(2), 155–184. https://doi.org/10.1007/BF02170076

Burgoon, J. K., Wang, X., Chen, X., Pentland, S. J., & Dunbar, N. E. (2021). Nonverbal
    behaviors "speak" relational messages of dominance, trust, and composure.
    *Frontiers in Psychology*, *12*. https://doi.org/10.3389/fpsyg.2021.624177

Camara, M. K., Postal, A., Maul, T. H., & Paetzold, G. H. (2024). Can lies be faked?
    Comparing low-stakes and high-stakes deception video datasets from a Machine
    Learning perspective. *Expert Systems with Applications*, *249*, 123684.
    https://doi.org/10.1016/j.eswa.2024.123684

Cardaioli, M., Miolla, A., Conti, M., Sartori, G., Monaro, M., Scarpazza, C., & Navarin, N.
    (2022). Face the truth: Interpretable emotion genuineness detection. *2022
    International Joint Conference on Neural Networks (IJCNN)*, 01–08.

Colasanti, M., Melis, G., Monaro, M., Ricci, E., Bosco, F., Rossi, M., Biondi, S.,
    Verrocchio, M. C., Di Domenico, A., Mazza, C., & Roma, P. (2024). Did You
    Commit a Crime There? Investigating the Visual Exploration Patterns of Guilty,
    Innocent, Honest, and Dishonest Subjects When Viewing a Complex Mock Crime
    Scene. *Journal of Nonverbal Behavior*, *48*(1), 47–71.
    https://doi.org/10.1007/s10919-023-00438-5

Constâncio, A. S., Tsunoda, D. F., Silva, H. d. F. N., Silveira, J. M. d., & Carvalho, D. R.
    (2023). Deception detection with machine learning: A systematic review and
    statistical analysis. *Plos one*, *18*(2), e0281323.

Delmas, H., Denault, V., Burgoon, J. K., & Dunbar, N. E. (2024). A Review of Automatic
    Lie Detection from Facial Features. *Journal of Nonverbal Behavior*, *48*(1), 93–136.
    https://doi.org/10.1007/s10919-024-00451-2

Dempster, A., Petitjean, F., & Webb, G. I. (2020). Rocket: Exceptionally fast and accurate time series classification using random convolutional kernels. *Data Mining and Knowledge Discovery*, *34*(5), 1454–1495.

DePaulo, B. M. (1992). Nonverbal behavior and self-presentation [Place: US Publisher: American Psychological Association]. *Psychological Bulletin*, *111*(2), 203–243. https://doi.org/10.1037/0033-2909.111.2.203

DePaulo, B. M., Kashy, D. A., Kirkendol, S. E., Wyer, M. M., & Epstein, J. A. (1996). Lying in everyday life. *Journal of Personality and Social Psychology*, *70*(5), 979–995. https://doi.org/10.1037/0022-3514.70.5.979

DePaulo, B. M., Lindsay, J. J., Malone, B. E., Muhlenbruck, L., Charlton, K., & Cooper, H. (2003). Cues to deception. *Psychological Bulletin*, *129*(1), 74–118. https://doi.org/10.1037/0033-2909.129.1.74

Dunbar, N. E., Burgoon, J. K., Chen, X., Wang, X., Ge, S., Huang, Q., & Nunamaker, J. (2023). Detecting ulterior motives from verbal cues in group deliberations. *Frontiers in Psychology*, *14*. https://doi.org/10.3389/fpsyg.2023.1166225

Dunbar, N. E., Giles, H., Bernhold, Q., Adams, A., Giles, M., Zamanzadeh, N., Gangi, K., Coveleski, S., & Fujiwara, K. (2019). Strategic synchrony and rhythmic similarity in lies about ingroup affiliation. *Journal of Nonverbal Behavior*, *44*(1), 153–172. https://doi.org/10.1007/s10919-019-00321-2

Dunbar, N. E., Giles, H., Bernhold, Q., Adams, A., Giles, M., Zamanzadeh, N., Gangi, K., Coveleski, S., & Fujiwara, K. (2020). Strategic Synchrony and Rhythmic Similarity in Lies About Ingroup Affiliation. *Journal of Nonverbal Behavior*, *44*(1), 153–172. https://doi.org/10.1007/s10919-019-00321-2

Duran, N. D., & Fusaroli, R. (2017). Conversing with a devil's advocate: Interpersonal coordination in deception and disagreement [Publisher: Public Library of Science]. *PLOS ONE*, *12*(6), e0178140. https://doi.org/10.1371/journal.pone.0178140

Ekman, P. (2009). *Telling lies: Clues to deceit in the marketplace, politics, and marriage (revised edition)*. WW Norton & Company.

Ekman, P., & Friesen, W. V. (1969). Nonverbal leakage and clues to deception†. *Psychiatry*, *32*(1), 88–106. https://doi.org/10.1080/00332747.1969.11023575

Ekman, P., & Friesen, W. V. (1978). Facial Action Coding System. https://doi.org/10.1037/t27734-000

Eskritt, M., Fraser, B., & Bosacki, S. (2022). Did You Just Lie to Me? Deception Detection in Face to Face versus Computer Mediated Communication. *The Journal of Social Psychology*, *162*(5), 566–579. https://doi.org/10.1080/00224545.2021.1933884

Eyben, F., Scherer, K. R., Schuller, B. W., Sundberg, J., André, E., Busso, C., Devillers, L. Y., Epps, J., Laukka, P., Narayanan, S. S., et al. (2015). The geneva minimalistic acoustic parameter set (gemaps) for voice research and affective computing. *IEEE transactions on affective computing*, *7*(2), 190–202.

Eyben, F., Wöllmer, M., & Schuller, B. (2010). Opensmile: The munich versatile and fast open-source audio feature extractor. *Proceedings of the 18th ACM international conference on Multimedia*, 1459–1462.

Fuller, C. M., Marett, K., & Twitchell, D. P. (2012). An Examination of Deception in Virtual Teams: Effects of Deception on Task Performance, Mutuality, and Trust. *IEEE Transactions on Professional Communication*, *55*(1), 20–35. https://doi.org/10.1109/TPC.2011.2172731

Gadzicki, K., Khamsehashari, R., & Zetzsche, C. (2020). Early vs late fusion in multimodal convolutional neural networks. *2020 IEEE 23rd international conference on information fusion (FUSION)*, 1–6.

Gogate, M., Adeel, A., & Hussain, A. (2017). Deep learning driven multimodal fusion for automated deception detection. *2017 IEEE symposium series on computational intelligence (SSCI)*, 1–6.

Goodman, B., & Flaxman, S. (2017). European union regulations on algorithmic decision-making and a "right to explanation". *AI Magazine*, *38*(3), 50–57. https://doi.org/10.1609/aimag.v38i3.2741

Goupil, L., Ponsot, E., Richardson, D., Reyes, G., & Aucouturier, J.-J. (2021). Listeners' perceptions of the certainty and honesty of a speaker are associated with a common prosodic signature. *Nature Communications*, *12*(1), 861. https://doi.org/10.1038/s41467-020-20649-4

Grossberg, S. (1988). Nonlinear neural networks: Principles, mechanisms, and architectures. *Neural networks*, *1*(1), 17–61.

Gupta, V., Agarwal, M., Arora, M., Chakraborty, T., Singh, R., & Vatsa, M. (2019). Bag-of-Lies: A Multimodal Dataset for Deception Detection. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 83–90. https://doi.org/10.1109/CVPRW.2019.00016

Han, J., Pei, J., & Tong, H. (2022). *Data mining: Concepts and techniques.* Morgan kaufmann.

Hartwig, M., & Bond Jr., C. F. (2014). Lie Detection from Multiple Cues: A Meta-analysis [_eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/acp.3052]. *Applied Cognitive Psychology*, *28*(5), 661–676. https://doi.org/10.1002/acp.3052

Kahou, S. E., Pal, C., Bouthillier, X., Froumenty, P., Gülçehre, Ç., Memisevic, R., Vincent, P., Courville, A., Bengio, Y., Ferrari, R. C., et al. (2013). Combining modality specific deep neural networks for emotion recognition in video. *Proceedings of the 15th ACM on International conference on multimodal interaction*, 543–550.

Keogh, E., Chakrabarti, K., Pazzani, M., & Mehrotra, S. (2001). Dimensionality reduction for fast similarity search in large time series databases. *Knowledge and information Systems*, *3*, 263–286.

Khan, W., Crockett, K., O'Shea, J., Hussain, A., & Khan, B. M. (2021). Deception in the eyes of deceiver: A computer vision and machine learning based automated deception detection. *Expert Systems with Applications*, *169*, 114341.

Lee, Z., Lindgren, T., & Papapetrou, P. (2024). Z-time: Efficient and effective interpretable multivariate time series classification. *Data mining and knowledge discovery*, *38*(1), 206–236.

Levine, T. R., Blair, J. P., & Clare, D. D. (2014). Diagnostic utility: Experimental demonstrations and replications of powerful question effects in high-stakes deception detection. *Human Communication Research*, *40*(2), 262–289.

Levine, T. R., Park, H. S., & McCornack, S. A. (1999). Accuracy in detecting truths and lies: Documenting the "veracity effect". *Communication Monographs*, *66*(2), 125–144. https://doi.org/10.1080/03637759909376468

Li, X., Pfister, T., Huang, X., Zhao, G., & Pietikäinen, M. (2013). A spontaneous micro-expression database: Inducement, collection and baseline. *2013 10th IEEE International Conference and Workshops on Automatic face and gesture recognition (fg)*, 1–6.

Lin, J., Keogh, E., Wei, L., & Lonardi, S. (2007). Experiencing sax: A novel symbolic representation of time series. *Data Mining and knowledge discovery*, *15*, 107–144.

Loconte, R., Russo, R., Capuozzo, P., Pietrini, P., & Sartori, G. (2023). Verbal lie detection using large language models. *Scientific Reports*, *13*(1), 22849.

Luke, T. J. (2019). Lessons from pinocchio: Cues to deception may be highly exaggerated. *Perspectives on Psychological Science*, *14*(4), 646–671. https://doi.org/10.1177/1745691619838258

Matsumoto, D., & Hwang, H. C. (2021). Clusters of nonverbal behavior differentiate truths and lies about future malicious intent in checkpoint screening interviews. *Psychiatry, Psychology and Law*, *28*(4), 463–478. https://doi.org/10.1080/13218719.2020.1794999

Matsumoto, D., & Wilson, M. (2023). Behavioral indicators of deception and associated mental states: Scientific myths and realities. *Journal of Nonverbal Behavior*, *48*(1), 11–23. https://doi.org/10.1007/s10919-023-00441-w

Middlehurst, M., Large, J., & Bagnall, A. (2020). The canonical interval forest (cif) classifier for time series classification. *2020 IEEE international conference on big data (big data)*, 188–195.

Oravec, J. A. (2022). The emergence of "truth machines"?: Artificial intelligence approaches to lie detection. *Ethics and Information Technology*, *24*(1), 6. https://doi.org/10.1007/s10676-022-09621-6

Pérez-Rosas, V., Abouelenien, M., Mihalcea, R., & Burzo, M. (2015). Deception detection using real-life trial data. *Proceedings of the 2015 ACM on international conference on multimodal interaction*, 59–66.

Poppe, R., van der Zee, S., Taylor, P. J., Anderson, R. J., & Veltkamp, R. C. (2024). Mining bodily cues to deception. *Journal of Nonverbal Behavior*, *48*(1), 137–159.

Prome, S. A., Ragavan, N. A., Islam, M. R., Asirvatham, D., & Jegathesan, A. J. (2024). Deception detection using ml and dl techniques: A systematic review. *Natural Language Processing Journal*, 100057.

Sen, T., Hasan, M. K., Teicher, Z., & Hoque, M. E. (2018). Automated Dyadic Data Recorder (ADDR) Framework and Analysis of Facial Cues in Deceptive Communication. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, *1*(4), 1–22. https://doi.org/10.1145/3161178

Shen, X., Fan, G., Niu, C., & Chen, Z. (2021). Catching a Liar Through Facial Expression of Fear. *Frontiers in Psychology*, *12*, 675097. https://doi.org/10.3389/fpsyg.2021.675097

Ströfer, S., Noordzij, M. L., Ufkes, E. G., & Giebels, E. (2015). Deceptive Intentions: Can Cues to Deception Be Measured before a Lie Is Even Stated? [Publisher: Public

Library of Science]. *PLOS ONE*, *10*(5), e0125237.

https://doi.org/10.1371/journal.pone.0125237

Talwar, V., & Crossman, A. (2022). Liar, liar . . . sometimes: Understanding

social-environmental influences on the development of lying. *Current Opinion in*

*Psychology*, *47*, 101374. https://doi.org/10.1016/j.copsyc.2022.101374

Tang, H., Lu, X., Cui, Z., Feng, C., Lin, Q., Cui, X., Su, S., & Liu, C. (2018). Resting-state

functional connectivity and deception: Exploring individualized deceptive

propensity by machine learning. *Neuroscience*, *395*, 101–112.

Tsuchiya, K., Hatano, R., & Nishiyama, H. (2023). Detecting deception using machine

learning with facial expressions and pulse rate. *Artificial Life and Robotics*, *28*(3),

509–519.

Van Der Zee, S., Taylor, P., Wong, R., Dixon, J., & Menacere, T. (2021). A liar and a

copycat: Nonverbal coordination increases with lie difficulty. *Royal Society Open*

*Science*, *8*(1), 200839. https://doi.org/10.1098/rsos.200839

Vrij, A. (2000). *Detecting lies and deceit: The psychology of lying and the implications for*

*professional practice*. John Wiley.

Vrij, A. (2004). Why professionals fail to catch liars and how they can improve [_eprint:

https://onlinelibrary.wiley.com/doi/pdf/10.1348/1355325041719356]. *Legal and*

*Criminological Psychology*, *9*(2), 159–181.

https://doi.org/10.1348/1355325041719356

Vrij, A., Hartwig, M., & Granhag, P. A. (2019). Reading Lies: Nonverbal Communication

and Deception. *Annual Review of Psychology*, *70*(1), 295–317.

https://doi.org/10.1146/annurev-psych-010418-103135

Wan, J., Escalera, S., Anbarjafari, G., Jair Escalante, H., Baró, X., Guyon, I., Madadi, M.,

Allik, J., Gorbova, J., Lin, C., et al. (2017). Results and analysis of chalearn lap

multi-modal isolated and continuous gesture recognition, and real versus fake

expressed emotions challenges. *Proceedings of the IEEE international conference on computer vision workshops*, 3189–3197.

Xu, Z., So, D. R., & Dai, A. M. (2021). Mufasa: Multimodal fusion architecture search for electronic health records. *Proceedings of the AAAI Conference on Artificial Intelligence*, *35*(12), 10532–10540.

Yan, W.-J., Li, X., Wang, S.-J., Zhao, G., Liu, Y.-J., Chen, Y.-H., & Fu, X. (2014). Casme ii: An improved spontaneous micro-expression database and the baseline evaluation. *PloS one*, *9*(1), e86041.

Yan, W.-J., Wu, Q., Liu, Y.-J., Wang, S.-J., & Fu, X. (2013). Casme database: A dataset of spontaneous micro-expressions collected from neutralized faces. *2013 10th IEEE international conference and workshops on automatic face and gesture recognition (FG)*, 1–7.

Zhang, L., Xie, Y., Xidao, L., & Zhang, X. (2018). Multi-source heterogeneous data fusion. *2018 International conference on artificial intelligence and big data (ICAIBD)*, 47–51.

Zhou, X., Jenkins, R., & Zhu, L. (2023). An Honest Joker reveals stereotypical beliefs about the face of deception. *Scientific Reports*, *13*(1), 16649. https://doi.org/10.1038/s41598-023-43716-4

Zloteanu, M., & Vuorre, M. (2024). A Tutorial for Deception Detection Analysis or: How I Learned to Stop Aggregating Veracity Judgments and Embraced Signal Detection Theory Mixed Models. *Journal of Nonverbal Behavior*, *48*(1), 161–185. https://doi.org/10.1007/s10919-024-00456-x

**Table 1**

*Summary statistics for the age, gender and education level of the participants in the study.*

|  |  | Overall | Lie | Truth |
| --- | --- | --- | --- | --- |
| n |  | 44 | 20 | 24 |
| Age, mean (SD) |  | 29.1 (11.1) | 30.4 (8.8) | 28.0 (12.9) |
| Age Groups, n (%) | 18-30 | 30 (68.2) | 11 (55.0) | 19 (79.2) |
|  | 31-40 | 8 (18.2) | 7 (35.0) | 1 (4.2) |
|  | 41-50 | 4 (9.1) | 2 (10.0) | 2 (8.3) |
|  | 51-60 | 1 (2.3) |  | 1 (4.2) |
|  | 61-70 | 1 (2.3) |  | 1 (4.2) |
| Gender, n (%) | Female | 28 (63.6) | 13 (65.0) | 15 (62.5) |
|  | Male | 14 (31.8) | 7 (35.0) | 7 (29.2) |
|  | Other | 2 (4.5) |  | 2 (8.3) |
| Education Level, n (%) | Middle school | 3 (6.8) | 1 (5.0) | 2 (8.3) |
|  | High School | 16 (36.4) | 7 (35.0) | 9 (37.5) |
|  | Vocational training | 3 (6.8) | 2 (10.0) | 1 (4.2) |
|  | Bachelor Degree | 14 (31.8) | 7 (35.0) | 7 (29.2) |
|  | Post graduate qualification | 8 (18.2) | 3 (15.0) | 5 (20.8) |

**Table 2**

*A table displaying the performance of the different classifiers on both unimodal and multimodal data. The best result for each metric is highlighted in **bold**, and the accuracy results that are at least human-level are highlighted with underlined text.*

| Modality | Participants | Method | Performance | | | | |
|---|---|---|---|---|---|---|---|
| | | | Precision Class Lie | Precision Class Truth | Recall Class Lie | Recall Class Truth | Accuracy |
| Voice | Sender | Rocket | 0.58 | 0.70 | **0.70** | 0.58 | 0.64 |
| | | CIF | 0.62 | 0.64 | 0.50 | 0.75 | 0.64 |
| | | Z-time | 0.60 | 0.67 | 0.60 | 0.67 | 0.64 |
| | Both | Rocket | 0.67 | 0.69 | 0.60 | 0.75 | 0.68 |
| | | CIF | 0.67 | 0.62 | 0.40 | 0.83 | 0.64 |
| | | Z-time | 0.67 | 0.69 | 0.60 | 0.75 | 0.68 |
| Facial Expressions | Sender | Rocket | 0.25 | 0.43 | 0.20 | 0.50 | 0.36 |
| | | CIF | 0.29 | 0.47 | 0.40 | 0.58 | 0.41 |
| | | Z-time | 0.29 | 0.47 | 0.20 | 0.58 | 0.41 |
| | Both | Rocket | 0.50 | 0.57 | 0.40 | 0.67 | 0.55 |
| | | CIF | 0.43 | 0.53 | 0.30 | 0.67 | 0.50 |
| | | Z-time | 0.20 | 0.47 | 0.10 | 0.67 | 0.41 |
| Multimodal Early Fusion | Voice: Sender | Rocket with PAA | 0.50 | 0.56 | 0.30 | 0.75 | 0.55 |
| | Facial Expressions: Sender | Rocket with SAX | 0.44 | 0.54 | 0.40 | 0.58 | 0.50 |
| | Voice: Sender | Rocket with PAA | 0.50 | 0.58 | 0.50 | 0.58 | 0.55 |
| | Facial Expressions : Both | Rocket with SAX | 0.57 | 0.60 | 0.40 | 0.75 | 0.59 |
| | Voice: Both | Rocket with PAA | 0.29 | 0.47 | 0.20 | 0.58 | 0.41 |
| | Facial Expressions: Sender | Rocket with SAX | 0.62 | 0.64 | 0.50 | 0.75 | 0.64 |
| | Voice: Both | Rocket with PAA | 0.50 | 0.58 | 0.50 | 0.58 | 0.55 |
| | Facial Expressions: Both | Rocket with SAX | 0.56 | 0.62 | 0.50 | 0.67 | 0.59 |
| Multimodal Late Fusion | Voice: Sender Facial Expressions: Sender | Decision Tree | 0.75 | **0.71** | 0.60 | 0.83 | 0.73 |
| | Voice: Sender Facial Expressions: Both | Decision Tree | 0.40 | 0.50 | 0.40 | 0.50 | 0.45 |
| | Voice: Both Facial Expressions: Sender | Decision Tree | **1.00** | **0.71** | 0.50 | **1.00** | **0.77** |
| | Voice: Both Facial Expressions: Both | Decision Tree | 0.43 | 0.53 | 0.30 | 0.67 | 0.50 |
| Receiver | Voice: Sender Facial Expressions: Sender | Personal experience | **1.00** | 0.57 | 0.10 | **1.00** | 0.59 |

**Figure 1**

*A photograph showing the experimental setup, with the two cameras in the center and two*
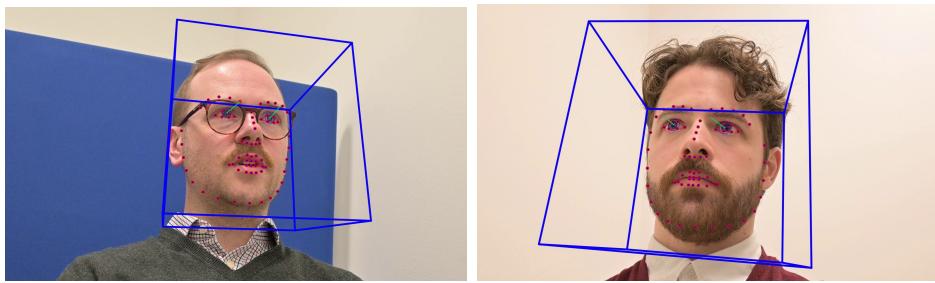
*participants taking part in the experimental interaction.*

**Figure 2**

*A screenshot showcasing the extraction process for facial features from the video recordings.*