# Explicit Modelling of Theory of Mind for Belief Prediction in Nonverbal Social Interactions

**Matteo Bortoletto [a,*], Constantin Ruhdorfer [a], Lei Shi [a] and Andreas Bulling [a]**

[a]University of Stuttgart, Germany

**Abstract.** We propose *MToMnet* – a Theory of Mind (ToM) neural network for predicting beliefs and their dynamics during human social interactions from multimodal input. ToM is key for effective nonverbal human communication and collaboration, yet existing methods for belief modelling have not included explicit ToM modelling or have typically been limited to one or two modalities. MToMnet encodes contextual cues (scene videos and object locations) and integrates them with person-specific cues (human gaze and body language) in a separate *MindNet* for each person. Inspired by prior research on social cognition and computational ToM, we propose three different MToMnet variants: two involving the fusion of latent representations and one involving the re-ranking of classification scores. We evaluate our approach on two challenging real-world datasets, one focusing on belief prediction while the other examining belief dynamics prediction. Our results demonstrate that MToMnet surpasses existing methods by a large margin while at the same time requiring a significantly smaller number of parameters. Taken together, our method opens up a highly promising direction for future work on artificial intelligent systems that can robustly predict human beliefs from their non-verbal behaviour and, as such, more effectively collaborate with humans.

## 1 Introduction

Social interaction and collaboration are essential human skills [14]. To engage in them effectively, humans have developed the ability to predict mental states and beliefs of others by observing their nonverbal behavioural cues, such as gaze or body language – so-called Theory of Mind [36, ToM]. Humans are adept at integrating multiple modalities for this task, including contextual information. Given its importance in human-human interactions, computational ToM has recently emerged as a new frontier in developing intelligent computational agents that can understand and collaborate with humans [18]. Despite a surge of papers on this new task, deep learning methods for predicting other agents' mental states have mainly been studied in constrained artificial environments [3, 38, 32, 15, 40, 39, 30, 31, 8]. Moreover, existing methods typically rely on only one or a few modalities to predict beliefs, such as visual or linguistic cues [28, 42]. Effectively integrating a wider range of modalities for belief prediction remains an open research challenge. Moreover, existing approaches for belief prediction in real-world settings have not explicitly added a ToM mechanism.

In this work, we focus on predicting human beliefs from multimodal inputs and how these beliefs change dynamically in real-world scenarios involving naturalistic dyadic (human-human) interactions. Belief prediction is particularly challenging, and ToM is particularly important when verbal communication is impossible. In these situations, individuals must instead rely on nonverbal cues to convey their intentions and beliefs. Advancing from recent work [13, 12], we propose a multimodal ToM neural network (*MToMnet*) that leverages person-specific nonverbal communication cues (gaze, pose) and contextual cues (video frames, object bounding boxes) to predict beliefs and how they change over time.

MToMnet encodes contextual cues using shared feature extractors and person-specific cues using two independent *MindNet*s – LSTM-based sub-networks that allow our model to encode individual traits. Without a clear theoretical framework for integrating ToM into neural networks, we draw inspiration from computational ToM and social cognition research and study three different variants of MToMnet that add explicit ToM modelling. The Decision-Based MToMnet (DB-MToMnet) adopts a decision-based strategy inspired by recent advancements in referential games [28]. Here, belief prediction for one individual is used to re-rank the predictions for the other. The other two approaches employ a model-based strategy, leveraging MindNets' latent representations. The Implicit Communication MToMnet (IC-MToMnet) enables the communication between MindNets via late fusion of internal representations. The Common Ground MToMnet (CG-MToMnet) is inspired by research in social cognition, in particular by the idea that human communication involves a shared, inter-subjective *common ground* [43]. We use MToMnet latent representations to create such common ground based on this insight.

We evaluate these MToMnet variants on two challenging multimodal real-world datasets that target complementary objectives. The Benchmark for Human Belief Prediction in Object-context Scenarios [12, BOSS] consists of videos of two people tasked to collaborate only using nonverbal communication. BOSS facilitates the evaluation of models' *belief prediction* capabilities, i.e., the ability to correctly predict the belief of both people for each video frame. In contrast, the Triadic Belief Dynamics dataset [13, TBD] focuses on communication events that emerge during in-the-wild social interactions between two people. TBD enables the evaluation of models' ability to predict changes in the *belief dynamics* of a person causally constructed by these events.

We report extensive experiments on both datasets, demonstrating that our approach significantly outperforms state-of-the-art methods while only using a fraction of the parameters. Our results emphasise the importance of explicit Theory of Mind (ToM) modelling for achieving these performance improvements. Moreover, analyses

* Corresponding Author. Email: matteo.bortoletto@vis.uni-stuttgart.de.

of MToMnet's latent representations underline the effectiveness of encoding person-specific cues with independent MindNets. Further post-hoc analyses on TBD show that MToMnet can predict *false beliefs dynamics* – beliefs that do not align with reality – more accurately than previous approaches.

Overall, this work makes four contributions:

- We introduce a multimodal ToM neural network (MToMnet) that combines nonverbal communication cues and visual inputs for belief and belief dynamics prediction.
- We propose three approaches to computationally modelling Theory of Mind, inspired by recent work on computational ToM and social cognition: decision-based, implicit communication, and common ground.
- We demonstrate that explicit ToM modelling allows us to achieve substantial performance gains for two tasks – belief prediction and belief dynamics prediction – at a significantly lower computational cost.
- We report analyses highlighting the effectiveness of modelling person-specific beliefs using independent MindNets and the efficacy of explicit ToM modelling for capturing false beliefs.

## 2 Related Work

### 2.1 Belief Prediction

Predicting beliefs is a challenging task, even in constrained artificial settings, such as grid-world environments [38, 15, 39, 30, 31, 8], 3D worlds consisting of basic geometric shapes [40] or virtual reality environments [37]. First, datasets that take a step towards mental state modelling in real-world settings have been proposed. These datasets consist of videos of human social interactions that have been annotated with rich social cues, such as gaze or body pose. The Benchmark for Human Belief Prediction in Object-context Scenarios (BOSS) focuses on *belief prediction* in dyadic collaborative interactions, i.e. the task of predicting beliefs of two people collaborating with each other [12]. Similarly, the Triadic Belief Dynamics dataset (abbreviated TBD here) focuses on the prediction of *belief dynamics*, i.e. predicting if and how someone's belief changes during social nonverbal interactions [13]. As such, both datasets complement each other in terms of the tasks they evaluate and the scenarios they cover, namely in-the-wild everyday activities (TBD) and collaborative scenarios (BOSS). In this work, we use both datasets to evaluate our method and show that employing a triadic structure and explicitly modelling ToM achieves better performance than existing methods for predicting both beliefs and belief dynamics.

### 2.2 Machine Theory of Mind

Theory of Mind (ToM) has been studied in cognitive science and psychology for decades, but our understanding of how humans develop this essential ability is still severely limited. Mirroring efforts to understand ToM in humans, an increasing number of works in the computational sciences have investigated means to equip artificial intelligent (AI) systems with similar capabilities. Previously proposed models that aim to implement a machine ToM have been based on partially observable Markov decision processes (POMDP) [11, 19], Bayesian methods [2, 27, 13, 29] and deep learning methods [38, 4, 47, 12, 28, 8, 6, 7]. Specifically for predicting beliefs of agents that engage in nonverbal communication, Duan et al. [12] and Fan et al. [13] follow different approaches. Duan et al. [12]

have used deep learning methods based on a ResNet [21] feature extractor for video frames and linear feature extractors for gaze, pose, bounding boxes and object-context relations. In contrast, Fan et al. [13] has used a triadic hierarchical energy-based model to track beliefs dynamics and compared it to neural network baselines that use only RGB frames, histogram of oriented gradients [10, HOG] or handcrafted features. Current deep learning approaches either handle input modalities shallowly or restrict themselves to a limited set of modalities. Energy-based models rely on more upfront engineering work involving the use of handcrafted features. Moreover, none of these approaches models ToM explicitly in their formulation. In this work, we show how explicitly modelling ToM in the neural network architecture can lead to substantial improvements compared to previous approaches.

## 3 Method

Our multimodal Theory of Mind neural network (MToMnet) combines nonverbal human communication cues (gaze and pose) with contextual cues (e.g. RGB video frames and object bounding boxes) to predict the beliefs of two observed human agents. In stark contrast to previous approaches [13, 12], our method leverages shared feature extractors and two *MindNets* that individually model each person's beliefs (see Figure 1). This design choice is motivated by research in social cognition suggesting that triadic (human-human-context) joint attentional engagement is necessary for effective cooperation [43]. This triadic engagement is reflected by our choice of two MindNets that encode nonverbal cues of each human and shared feature extractors that encode contextual cues available to both humans.

### 3.1 Base MToMnet

Our base MToMnet consists of two separate MindNets and a set of shared feature extractors. Each MindNet encodes individual cues from one person (e.g. human gaze and body language), combines them with contextual features (e.g. scene videos and object locations) coming from the shared features extractors, and adds temporal information. Let $\boldsymbol{x}_C \in \mathcal{C}$ be a set of contextual cues and $\boldsymbol{x}_I \in \mathcal{I}$ a set of individual cues for a specific person in the scene. Contextual cues are encoded by shared features extractors and concatenated

$$\boldsymbol{x}_{ctx} = \|_{C \in \mathcal{C}} \texttt{SharedFeatExtr}_C(\boldsymbol{x}_C) \tag{1}$$

where $\|$ denotes concatenation. Similarly, individual cues for a particular person are encoded and concatenated:

$$\boldsymbol{x}_{ind} = \|_{I \in \mathcal{I}} \texttt{MindNetFeatExtr}_I(\boldsymbol{x}_I) \tag{2}$$

Individual and contextual features are subsequently concatenated and used as input to a normalisation layer [1, LN], followed by a bidirectional LSTM [17] to model temporal information. The final LSTM hidden state is passed on to one or more fully connected (FC) classification heads that yield a probability distribution over classes:

$$P(y_i|\boldsymbol{x}) \propto \exp(\texttt{FC}(\texttt{LSTM}(\texttt{LN}(\boldsymbol{x})))) \tag{3}$$

where $\{y_i\}_{i=0}^{Y}$ are dataset-specific classes and $\boldsymbol{x} = \boldsymbol{x}_{ctx} \ \| \ \boldsymbol{x}_{ind}$. The final belief predictions are obtained using the argmax of $P(y_i|\boldsymbol{x})$:

$$\tilde{b} = \underset{i}{\operatorname{argmax}} \, P(y_i|\boldsymbol{x}) \tag{4}$$

In the following, we refer to the inputs, latent representations, and outputs of the two person-specific MindNets with the indices $\{1, 2\}$. We investigate three different variants of explicitly adding Theory of Mind to this base model architecture.
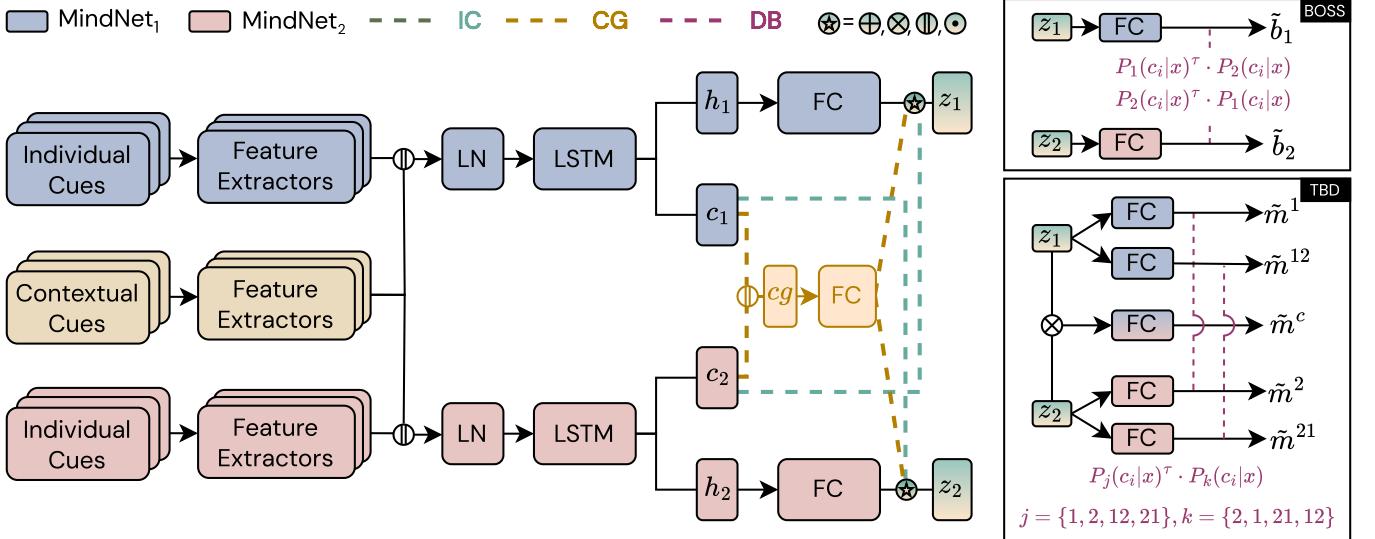
**Figure 1**: Our multimodal Theory of Mind neural network *MToMnet* consists of two separate *MindNet*s – one for each person – to encode individual cues (e.g. human gaze and body language) and integrate them with contextual cues (e.g. scene videos and object locations). We propose three different MToMnet variants: Decision-Based MToMnet (DB-MToMnet) combines class probabilities of the two MindNets to re-rank predictions. Implicit Communication MToMnet (IC-MToMnet) adds a communication mechanism between the two MindNets that exchange their internal LSTM cell state. Common Ground MToMnet (CG-MToMnet) forms a common ground representation by concatenating the two MindNets' cell states. We also study different aggregation operations: $\otimes$ = element-wise multiplication, $\oplus$ = element-wise sum, $\|$ = concatenation, and $\odot$ = cross-attention.

## 3.2 MToMnet Variants

We study three different variants of MToMnet that *draw inspiration* from prior research on computational ToM and social cognition to add explicit ToM modelling: Decision-Based (DB-MToMnet), Implicit Communication (IC-MToMnet), and Common Ground (CG-MToMnet). Our goal is to explore whether an internal ToM mechanism can, similar to humans, also benefit computational agents. We study different operations to combine neural representations for each variant of MToMnet and identify the best for our tasks.

**Decision-Based Theory of Mind (DB-MToMnet).** Inspired by previous work on referential games, where a speaker agent uses an internal listener model to re-rank potential utterances [28], DB-MToMnet incorporates a ToM mechanism to re-rank class label predictions. More specifically, we combine $P(y_i|\boldsymbol{x}_1)$ and $P(y_i|\boldsymbol{x}_2)$ within the MToMnet using a "ToM weight" hyper-parameter $\tau$, and we take the argmax of this score as the final belief prediction:

$$\tilde{b}_1 = \operatorname{argmax}(P(y_i|\boldsymbol{x}_1)^\tau \cdot P(y_i|\boldsymbol{x}_2)) \quad (5)$$

$$\tilde{b}_2 = \operatorname{argmax}(P(y_i|\boldsymbol{x}_2)^\tau \cdot P(y_i|\boldsymbol{x}_1)) \quad (6)$$

In contrast to Liu et al. [28], we apply the weight $\tau$ to the original probability distribution, e.g. to $P(y_i|\boldsymbol{x}_1)$ for MindNet$_1$, and not to the other probability distribution. As such, the hyper-parameter $\tau$ controls the extent to which the prediction from one MindNet impacts that of the other: the larger $\tau$, the smaller the impact.

**Implicit Communication Theory of Mind (IC-MToMnet).** Conceptually similar to DB-MToMnet, IC-MToMnet enables communication between the two MindNets via internal representations instead of exchanging ranking scores. Specifically, given the LSTM outputs

$$\boldsymbol{h}_1, \boldsymbol{c}_1 = \text{LSTM}_1(\boldsymbol{x}_1) \qquad \boldsymbol{h}_2, \boldsymbol{c}_2 = \text{LSTM}_2(\boldsymbol{x}_2) \quad (7)$$

where $\boldsymbol{h}$ and $\boldsymbol{c}$ indicate the LSTM hidden state and cell state, we aggregate one MindNet's hidden state with the other MindNet's cell

state, and vice versa. As we use bidirectional LSTMs, we use a fully connected layer to project the hidden state to the cell state dimension:

$$\boldsymbol{z}_1 = \text{FC}(\boldsymbol{h}_1) \star \boldsymbol{c}_2 \qquad \boldsymbol{z}_2 = \text{FC}(\boldsymbol{h}_2) \star \boldsymbol{c}_1 \quad (8)$$

where $\star$ can be one of the following aggregation operations: addition, multiplication, concatenation, or cross-attention [45]. $\boldsymbol{z}_1$ and $\boldsymbol{z}_2$ are used to obtain the predictions in the final classification layers.

**Common Ground Theory of Mind (CG-MToMnet).** This final variant is inspired by the idea that the human communicative context is not limited to the surrounding environment but involves a wider, shared and inter-subjective context known as common ground [9, 43]. Tomasello [43] refers to "common ground" as *shared experience between individuals* that is critical for all human communication. As such, common ground represents a broad concept that may include perception, attention, and knowledge. In this work, we build such common ground by combining the LSTM cell state $\boldsymbol{c}$ of each MindNet. We chose the LSTM cell state as it represents the memory of the network, storing information over time. In practice, we concatenate the cell state of the two LSTMs to form a *common ground representation*:

$$\boldsymbol{cg} = \text{FC}(\boldsymbol{c}_1 \| \boldsymbol{c}_2) \quad (9)$$

The $\boldsymbol{cg}$ tensor is then aggregated with the LSTM hidden states to obtain $\boldsymbol{z}_1$ and $\boldsymbol{z}_2$:

$$\boldsymbol{z}_1 = \text{FC}(\boldsymbol{h}_1) \star \boldsymbol{cg} \qquad \boldsymbol{z}_2 = \text{FC}(\boldsymbol{h}_2) \star \boldsymbol{cg} \quad (10)$$

where $\star$ has the same meaning as before. $\boldsymbol{z}_1$ and $\boldsymbol{z}_2$ are used to obtain the predictions in the final classification layers.

## 4 Experiments

### 4.1 Datasets

**BOSS.** The Benchmark for Human Belief Prediction in Object-context Scenarios (BOSS) is a real-world dataset consisting of videos

**Figure 2**: Examples from BOSS [12] and TBD [13]. BOSS includes third-person video frames, bounding boxes (top left), 3D gaze (top centre) and body pose (top right). TBD includes third-person (bottom left) and first-person (bottom centre) video frames, 2D gaze (bottom centre, pink dot) and body pose (bottom right).

of two humans performing collaborative tasks without verbal communication [12]. The dataset consists of 900 third-person videos recorded with 10 participants in 15 different object-context situations. In each video, each person stands in front of a table, one with contextual objects and the other with objects that can be selected based on the context (see Figure 2, top). The person standing in front of the contextual objects table receives a contextual task and must non-verbally communicate which object on the other table should be selected. In Figure 2 (top), the task is *"Circle the words on the magazine's cover"*, for which the correct object to be selected by the participant on the right is the marker.

The computational task is to predict each person's beliefs for each video frame correctly. This translates into a classification problem where the possible classes $\{c_i\}_{i=0}^N$ are the different objects, with $N = 27$. The dataset is annotated with gaze estimation, pose estimation, bounding boxes, ground truth beliefs, and an object-context relations matrix (OCR) that describes which objects are more likely to appear given a certain context. For instance, a hammer is more likely to be present when there are nails rather than when there is a pizza. Given that we noticed that the original bounding box annotations were highly inaccurate during preliminary experiments, we opted for re-extracting them using YOLOv5 [24]. [1]

**Triadic Belief Dynamics Dataset (TBD).** Fan et al. [13] have collected a dataset covering nonverbal communication in rich social interactions. Participants were not provided a detailed script but only with the type of nonverbal communication they could use. The dataset consists of 88 videos recorded with 12 people in seven different scenarios. It includes first- and third-person video frames and gaze, pose, and bounding box annotations (see Figure 2, bottom). We opted for also evaluating on TBD because it differs from and complements BOSS in two distinct ways. First, people were asked to perform three types of nonverbal communication – no communication, attention following, and joint attention – that, while highly relevant for belief prediction, do not directly involve collaboration. Second, TBD differs from BOSS in that the task is not to predict a person's belief about objects in the scene itself but its dynamics, i.e. if and how this belief changes over time. Concretely, given video clips of five frames, a model has to classify the belief dynamics for a selected object in the scene into four classes: *occur*, *disappear*, *update*, and *null*. Importantly, TBD involves predicting belief dynamics not only for first-order ($m^1, m^2$) but also second-order beliefs ($m^{12}, m^{21}$) – i.e. beliefs over another person's beliefs – and a common mind ($m^c$) that corresponds to a common ground between the two participants.

---

[1] Link to code and improved annotations in the Appendix.

## 4.2  Implementation Details

**MToMnet.** BOSS and TBD share most input modalities, with a few exceptions. For BOSS, the contextual cues consist of third-person RGB videos, object bounding boxes, and the OCR matrix. Individual cues are derived from 3D gaze and pose. TBD contains third-person RGB videos and object bounding boxes as contextual cues and first-person RGB videos, pose, and 2D gaze as individual cues. Given these differences, we implemented dataset-specific feature extractors. All MToMnet variants encoded RGB video frames using a three-layer CNN with 16, 32, and 64 filters, respectively. Each convolutional layer was followed by ReLU activation and max pooling. The OCR matrix and poses were processed by a graph convolutional network [26]. The OCR matrix naturally represents an adjacency matrix for our graphs, and we used its normalised values as node features. For the poses, we defined the adjacency matrix based on connections of body joints and used the 3D joint coordinates as node features. Object bounding boxes are fed into a fully connected layer. Each MindNet uses a one-layer bidirectional LSTM, preceded by layer normalisation. All layers in our MToMnet models have hidden dimension 64 and are followed by GELU activation [22] and dropout [41] with $p = 0.1$. For DB-MToMnet we set the "ToM weight" to $\tau = 2$. For BOSS, each MindNet outputs a belief prediction. For TBD, classification layers for $m^1$ and $m^{12}$ take $z_1$ as input, whereas classification layers for $m^2$ and $m^{21}$ take $z_2$ as input. Since $m^c$ represents the "common mind" between both agents [13], $z_1$ and $z_2$ are first aggregated by performing element-wise multiplication and then fed into the classification layer. Additional details on the architecture are provided in Appendix.

**Training.** We trained all models for 300 epochs using three distinct random seeds. Cross-entropy was employed as the loss function, and model checkpoints were saved based on the highest validation accuracy for BOSS and the highest macro F1 score for TBD. These metrics were chosen for easier comparison with the original works [12, 13]. We used the Adam optimiser [25] with a learning rate of $5 \cdot 10^{-4}$. Additional details are provided in Appendix.

**Baselines.** We compare our approach with the original models [13, 12]. For BOSS, Duan et al. [12] have proposed four models based on a ResNet34 backbone for encoding video frames and fully connected layers for other modalities like gaze, pose, bounding boxes, and OCR. The models are CNN, CNN+GRU, CNN+LSTM, and CNN+Conv1D, each incorporating different types of recurrent or convolutional layers before passing the concatenated latent representations to two classification layers. To ensure a fair comparison, we re-trained these models for 300 epochs, as the original models were only trained for five epochs.

Fan et al. [13] have evaluated similar approaches on TBD. Their first model (CNN) uses a ResNet50 to extract frame features, followed by fully connected layers for belief dynamics classification. The CNN+HOG-LSTM model incorporates histograms of oriented gradient [10, HOG] features of the frame patch gazed at by the participants along with full frame features. The CNN+HOG+Mem model concatenates the history of predicted belief dynamics with frame and HOG features, while the Feats+Memory model combines hand-crafted features with the history of predicted belief dynamics using a multi-layer perceptron. Fan et al. [13] achieved state-of-the-art performance on TBD by using a hierarchical graphical model (denoted here as HGM) trained using a beam-search algorithm on handcrafted events derived from raw pixels.
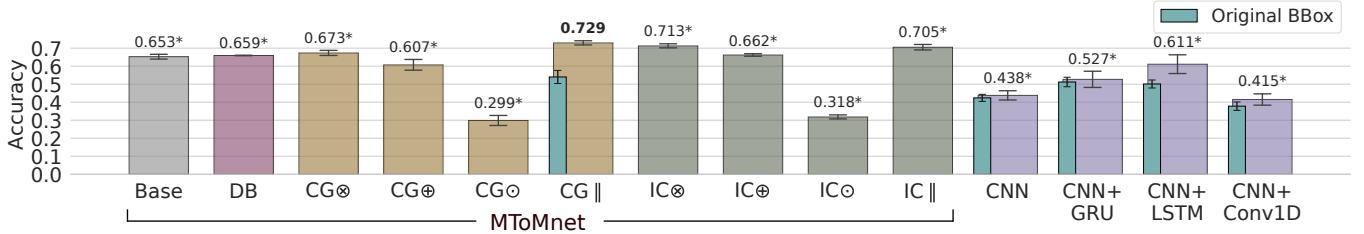
**Figure 3**: Accuracy for belief prediction on BOSS for our MToMnet models and baselines [12] using all input modalities. Scores significantly different from CG∥-MToMnet according to a paired t-test ($p < 0.05$) are marked with a *.

**Table 1**: Macro F1 scores for previous approaches [13] and MToMnet variants on TBD. Scores significantly different from CG∥-MToMnet according to a paired t-test ($p < 0.05$) are marked with a *.

| Model | $m^1$ | $m^2$ | $m^{12}$ | $m^{21}$ | $m^c$ | Average |
|---|---|---|---|---|---|---|
| Chance | 0.103 | 0.104 | 0.102 | 0.101 | 0.100 | 0.102 |
| CNN | 0.171 | 0.167 | 0.169 | 0.174 | 0.250 | 0.186 |
| CNN+HOG-LSTM | 0.167 | 0.132 | 0.205 | 0.182 | 0.250 | 0.187 |
| CNN+HOG+Mem | 0.285 | 0.285 | 0.246 | 0.250 | 0.155 | 0.244 |
| Feats+Mem | 0.274 | 0.288 | 0.230 | 0.227 | 0.191 | 0.242 |
| HGM | 0.431 | 0.443 | 0.351 | 0.349 | 0.299 | 0.375 |
| -MToMnet | | | | | | |
| Base | $0.276 \pm 0.055$* | $0.473 \pm 0.085$* | $0.313 \pm 0.014$* | $0.473 \pm 0.010$ | $0.566 \pm 0.095$* | $0.420 \pm 0.023$ |
| DB | $0.442 \pm 0.067$* | $0.313 \pm 0.051$* | $0.425 \pm 0.006$* | $0.385 \pm 0.013$* | $0.459 \pm 0.010$* | $0.405 \pm 0.013$ |
| CG∥ | $\mathbf{0.477 \pm 0.078}$ | $0.460 \pm 0.065$ | $0.452 \pm 0.010$ | $0.467 \pm 0.013$ | $\mathbf{0.583 \pm 0.080}$ | $\mathbf{0.488 \pm 0.022}$ |
| CG⊕ | $0.461 \pm 0.059$ | $0.472 \pm 0.076$ | $0.459 \pm 0.009$ | $0.468 \pm 0.012$* | $0.544 \pm 0.091$ | $0.481 \pm 0.022$ |
| CG⊗ | $0.462 \pm 0.062$* | $\mathbf{0.478 \pm 0.068}$ | $0.451 \pm 0.010$ | $0.469 \pm 0.014$* | $0.559 \pm 0.103$* | $0.484 \pm 0.023$ |
| CG⊙ | $0.461 \pm 0.057$ | $0.464 \pm 0.062$* | $\mathbf{0.469 \pm 0.015}$ | $0.475 \pm 0.009$ | $0.564 \pm 0.099$ | $0.486 \pm 0.022$ |
| IC∥ | $0.462 \pm 0.067$* | $0.463 \pm 0.067$ | $0.459 \pm 0.002$ | $0.471 \pm 0.014$ | $0.556 \pm 0.093$* | $0.482 \pm 0.022$ |
| IC⊕ | $0.465 \pm 0.074$* | $0.474 \pm 0.073$* | $0.466 \pm 0.025$* | $0.473 \pm 0.015$ | $0.561 \pm 0.098$ | $\mathbf{0.488 \pm 0.025}$ |
| IC⊗ | $0.417 \pm 0.066$* | $0.423 \pm 0.072$ | $0.464 \pm 0.001$ | $\mathbf{0.486 \pm 0.007}$ | $0.450 \pm 0.007$* | $0.448 \pm 0.014$ |
| IC⊙ | $0.455 \pm 0.072$* | $0.257 \pm 0.002$* | $0.462 \pm 0.032$* | $0.274 \pm 0.016$* | $0.554 \pm 0.084$* | $0.401 \pm 0.019$ |

## 4.3 Model Performance

Results for the different MToMnet variants and the baselines on BOSS are shown in Figure 3. We use the following notation to indicate the different aggregation operations: ⊗ = element-wise multiplication, ⊕ = element-wise sum, ∥ = concatenation, ⊙ = self-attention. As can be seen from the figure, already the Base-MToMnet (i.e. without explicit ToM modelling) outperforms all baselines (0.653 accuracy), despite requiring less than 3% of their parameters (∼ 450k vs 21M). Second, incorporating explicit ToM modelling (DB, CG, IC) yields further performance improvements, but the choice of aggregation is critical. CG∥-MToMnet exhibits the highest overall performance amongst all models (0.729), followed by IC⊗-MToMnet (0.713), and IC∥-MToMnet (0.705). In contrast, the DB-MToMnet (0.659) only achieves a marginal improvement over the Base-MToMnet (0.653). Figure 3 also shows results for CG∥-MToMnet and baselines obtained using the original bounding box annotations. While CG∥-MToMnet still outperforms all the baselines, the baselines do not benefit from the improved bounding boxes. This is likely attributed to the shallow feature aggregation in the baseline models. The paired t-test revealed a significant difference ($p < 0.05$) between CG∥-MToMnet and the other models. To further validate our architectural choice, we evaluated a single MindNet model. Our methods outperform this model and achieve an accuracy of 0.61, except for CG⊙ and IC⊙. We report modality ablation studies in Appendix.

Evaluation scores for belief dynamics prediction on TBD are shown in Table 1, where we also report paired t-test results ($p < 0.05$) between CG∥-MToMnet and other MToMnet variants. As for
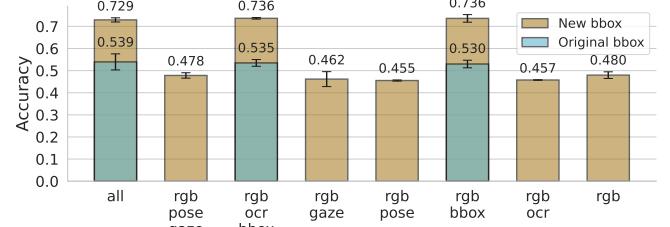


**Figure 4**: Modality ablation study for BOSS.

BOSS, our Base-MToMnet already achieves better performance than the best baseline (HGM) – with the only exception of $m^1$ (0.276 vs. 0.431) and $m^{12}$ (0.313 vs. 0.351). Adding explicit ToM modelling leads to further performance gains for all three variants, up to 30% on HGM. Considering the average performance across different aggregation types, the best model is again the one inspired by social cognition, CG-MToMnet. In particular, the single best performing models are CG∥-MToMnet and IC⊕-MToMnet, on par (0.488). Remarkably, all our MToMnet variants achieved their highest F1 scores on $m^c$, except for IC⊗-MToMnet. In contrast, baselines typically find classifying belief dynamics for $m^c$ one of the most challenging tasks, achieving lower scores. Our CG∥-MToMnet (0.583) substantially outperformed the best baseline, HGM (0.299), by a substantial margin of improvement.

These results highlight the effectiveness of our proposed MToMnet architecture and underline the significance of explicit ToM modelling in achieving superior performance with significantly reduced computational costs.
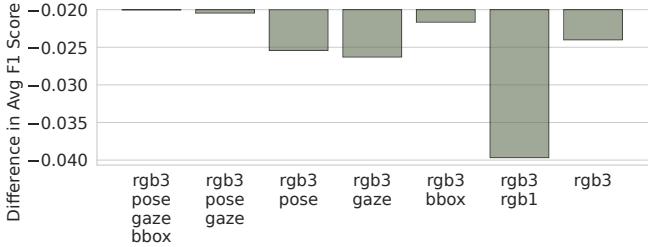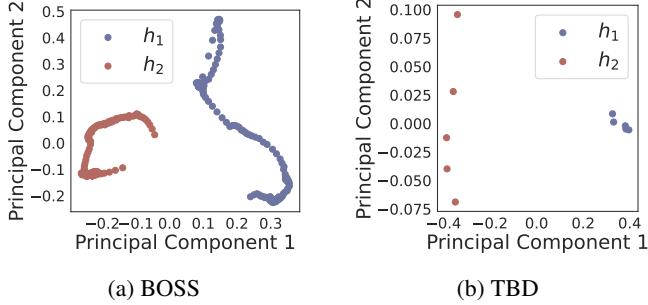
**Figure 5**: Modality ablation study for TBD.



(a) BOSS

(b) TBD

**Figure 6**: Examples from PCA results for $h_1$ and $h_2$ from CG‖-MToMnet, taken from BOSS and TBD test set.

**Modality Ablation Study.** Figure 4 shows the accuracy achieved by ablated versions of our best-performing model, CG‖-MToMnet, on the BOSS dataset. The results highlight the significant impact of including bounding boxes as input modality (rgb+ocr+bbox and rgb+bbox). This result aligns with previous research [12], emphasising the crucial role of knowing the objects present in the scene for the task. When excluding bounding boxes, the accuracy scores decreased. We suspected that the newly introduced bounding box annotations might be a major contributing factor to this outcome. Therefore, for the sake of completeness, we conducted additional experiments where CG‖-MToMnet was trained and evaluated using the original bounding box data. The results, depicted in Figure 4 (teal), demonstrate that when utilising the original bounding boxes, the performance gap with other ablated versions of the model drastically decreases. Nevertheless, our model performs better than the baselines even when using original bounding box annotations.

The differences in averaged F1 scores across $m^i$, $i = \{1, 2, 12, 21\}$ on TBD between the complete CG‖-MToMnet and its ablated versions are reported in Figure 5. Our full model achieves an F1 score of $0.488$. Ablating modalities generally result in performance degradation, with the worst-performing version experiencing an 8.9% decline, attaining an F1 score of $0.449$ (rgb1+rgb3). Integrating multiple behavioural cues, such as pose and gaze, contributes positively to performance. Specifically, the combination of third-person RGB frames with gaze and pose achieved an F1 score of $0.487$, outperforming models where third-person RGB frames were combined with only gaze or pose.

### 4.4 Modelling Person-Specific Features

To assess the efficacy of encoding individual cues and predicting individuals' beliefs using independent MindNets, we compared the feature representations from the two LSTMs, $h_1$ and $h_2$, by employing Principal Component Analysis [35, PCA]. As the examples in Figure 6 show, the LSTMs hidden states $h_1$ and $h_2$ from CG‖-MToMnet show a clear disentanglement of principal components both on BOSS and TBD. We found the same behaviour for all our

**Table 2**: Label counts for TBD train and test set.

| Mind | Train/test count | | | |
|---|---|---|---|---|
| | Occur | Disappear | Update | Null |
| $m^1$ (all) | 156/48 | 0/0 | 2176/731 | 1750/582 |
| $m^1$ (false belief) | 0/1 | 0/0 | 5/0 | 158/53 |
| $m^2$ (all) | 146/46 | 9/0 | 2579/860 | 1348/455 |
| $m^2$ (false belief) | 0/1 | 0/0 | 1/1 | 118/47 |
| $m^{12}$ (all) | 75/25 | 16/5 | 916/298 | 3075/1033 |
| $m^{12}$ (false belief) | 0/0 | 0/0 | 0/0 | 19/4 |
| $m^{21}$ (all) | 81/27 | 11/1 | 821/260 | 3169/1073 |
| $m^{21}$ (false belief) | 0/0 | 0/0 | 0/0 | 52/20 |

MToMnet variants and provide examples for all of them in Appendix. This finding suggests that each MindNet represents person-specific cues and fuses them with contextual cues differently for different persons, highlighting the value of such architectural choice.

### 4.5 False Belief Dynamics Prediction

TBD involves predicting belief dynamics for both first- and second-order beliefs, i.e., self-beliefs ($m^1, m^2$) and beliefs held over another person's beliefs ($m^{12}, m^{21}$). Human beliefs, however, are not always aligned with reality: individuals may hold a first- or second-order *false belief* [49]. We show an example in Figure 7. Thinking her partner had left the room, a participant moved an apple from the backpack to behind her laptop. She believes her partner believes the apple is still in the backpack. However, she is unaware that her partner is actually standing by the door and saw her move the apple to the table. This situation exemplifies a second-order (*she believes he thinks*) false belief. Despite not being used in [13], the TBD dataset includes false belief annotations, allowing us to perform post-hoc analyses on models' capabilities to predict such false belief dynamics. That is, we do not train models specifically to recognise false beliefs but evaluate models' belief dynamics predictions that, according to the provided annotations, correspond to false beliefs.

Figure 8 summarises the accuracy for false belief dynamics predictions on TBD, considering first-order, second-order, and both false belief types. We report accuracy as we were interested in detecting whether a model predicted correctly (positive class) or not (negative class) a (a certain type of) false belief. For first-order false belief and joint first- and second-order false belief dynamics prediction, all IC- and CG-MToMnet variants outperform Base, except for IC⊙. In particular, our CG-MToMnet variants achieve the highest accuracy, improving over Base-MToMnet by a large margin on first-order false beliefs (up to $0.78$ vs. $0.49$) and on joint first- and second-order false beliefs (up to $0.80$ vs. $0.57$).

Figure 8 also shows that predicting second-order false belief dynamics – generally more difficult – is easier on TBD. To understand why, we compared the distribution of all labels with those corresponding to false beliefs, as shown in Table 2. In training and test sets, most false belief labels for all minds $m^i$, $i = \{1, 2, 12, 21\}$, are *null*. For false beliefs associated with $m^{12}$ and $m^{21}$, *null* is the only possible label. Thus, most false beliefs in the dataset correspond to situations where individuals assume that nothing changed although something did (*occur*, *disappear*, or *update*) – akin to the Sally-Anne test [5]. However, when considering the overall label distribution, the most frequent label for $m^1$ and $m^2$ is *update*, thus leading to biases towards predicting the *update* class during training. This ultimately leads to lower accuracy in predicting first-order false beliefs, where *null* is the most prevalent label. However, our MToMnet variants, especially CG-MToMnet, can overcome this bias.

**Figure 7**: Example of *second-order* false belief from TBD [13]. Thinking her partner had left the room, a participant moved an apple from the backpack to behind her laptop. *She believes her partner believes* the apple is still in the backpack. However, she is unaware that her partner is actually standing by the door and saw her move the apple to the table.
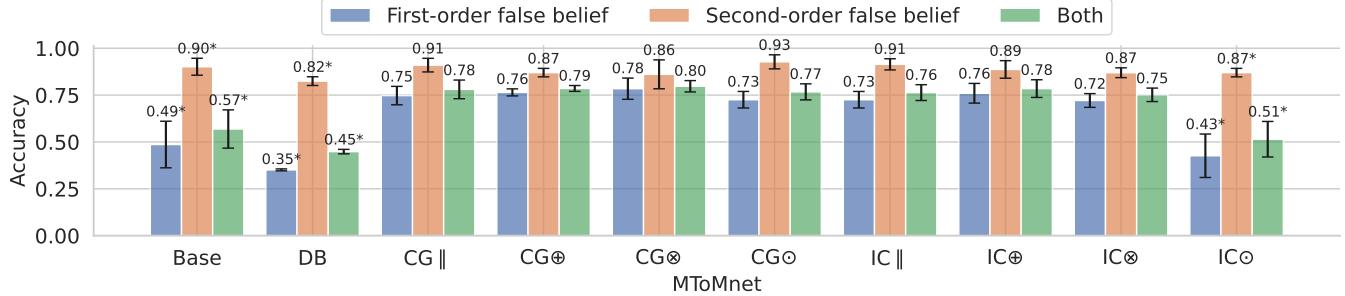


**Figure 8**: Accuracy of false belief prediction on TBD for different variants of our MToMnet (Base = no explicit ToM modelling, DB = decision-based ToM, CG = common ground ToM, IC = implicit communication ToM). Scores significantly different from CG∥-MToMnet according to a paired t-test ($p < 0.05$) are marked with a *.

## 5   Discussion

**Social Cognition Is All You Need.**   In this work, we presented MToMnet – a Theory of Mind neural network for predicting beliefs and their dynamics during nonverbal human social interactions. Our best performing MToMnet variant (**CG-MToMnet**) **achieved new state-of-the-art results** for belief (dynamics) prediction on both BOSS (Figure 3) and TBD (Table 1), as well as for false belief dynamics prediction on TBD (Figure 8). At the same time, **MToMnet requires considerably fewer parameters than previous state-of-the-art methods** – approximately 460k vs. 21M for BOSS and 1M for TBD (see Appendix) – and, as a result, shows faster training (less than 3 hours vs. 17 to 20 hours for baselines on BOSS [12]). These findings are important as they underline the significant potential of adopting concepts developed in the cognitive sciences when designing computational agents. This also represents a paradigm shift compared to the recent trend of improving performance mainly through upscaling model complexity. Our analyses of the latent representations showed that the two MindNets are effective in capturing individuals' information in distinct ways for modelling beliefs (Figure 6), further supporting the architectural design choices of MToMnet. Remarkably, **CG-MToMnet can be easily adapted to interactions involving more than two interacting agents**, thanks to its shared common ground representation across all MindNets. This adaptability is crucial for future work exploring scenarios with several human and computational agents. Extending DB- and IC-MToMnet to involve more than two agents is more challenging. For DB-MToMnet, it would be crucial to find the right set of $\tau$ such that probabilities do not vanish. One interesting idea for future work is to make $\tau$ learnable. IC∥-MToMnet faces the challenge of aggregating numerous hidden states, which could result in a quite large vector. A broader limitation arises when dealing with dynamic environments, where the number of individuals in a scene can vary. Together with extending models to dynamic environments, another promising direction is to improve the integration of different input features, a facet not explored in this work. For example, on BOSS, the OCR matrix could be used to define an additional term in the loss function to enforce object-context relations.

**Explicit Modelling of ToM Improves Performance.**   A key contribution of our work is the explicit modelling of ToM using multimodal individual and contextual cues. Our experiments demonstrated that MToMnet not only achieves new state-of-the-art performance on the two most common benchmark datasets but also outperforms existing baselines by a large margin. Instrumental to these improvements is the explicit ToM modelling that allowed us to improve our Base-MToMnet by 19% on BOSS (CG∥-MToMnet, Figure 3), and up to 60% on TBD (CG∥-MToMnet, Table 1, Average). This particularly shows for the challenging task of false belief prediction on TBD for which explicit ToM modelling leads to significant improvements of up to 60% over the Base-MToMnet model (Figure 8). This finding is particularly significant considering the large body of work and long-standing efforts on false belief prediction [16, 3, 38, 44].

**Limitations of Current Benchmarks.**   Our ablation studies highlight a fundamental limitation inherent in current benchmarks due to the quality and consistency of the data, negatively affecting model performance. In this work, we found inaccuracies in bounding box annotations in the BOSS dataset and re-extracted them. This led to a substantial enhancement in performance, highlighting the importance of precise bounding box information. Rectifying similar issues in gaze data proved unfeasible, as participant faces in the BOSS dataset were deliberately obscured for privacy reasons. Despite data limitations, our models outperformed the baselines with both original and revised annotations (Figure 4).

## 6   Conclusion

In this work, we proposed MToMnet, a Theory of Mind network that predicts beliefs and their dynamics during human social interactions from multimodal input. Building on social cognition and ToM research, we designed three MToMnet variations: one decision-based and two model-based. Across two real-world datasets, MToMnet outperformed existing methods in both belief prediction and belief dynamics prediction, despite having fewer parameters. These results advance the state-of-the-art in belief prediction thus facilitating better collaboration with humans.

# Acknowledgements

# References

[1] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.

[2] C. L. Baker, R. Saxe, and J. B. Tenenbaum. Action understanding as inverse planning. *Cognition*, 113(3):329–349, 2009.

[3] C. L. Baker, J. Jara-Ettinger, R. Saxe, and J. B. Tenenbaum. Rational quantitative attribution of beliefs, desires and percepts in human mentalizing. *Nature Human Behaviour*, 1(4):0064, 2017.

[4] C.-P. Bara, S. CH-Wang, and J. Chai. MindCraft: Theory of mind modeling for situated dialogue in collaborative tasks. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1112–1125. Association for Computational Linguistics, 2021. doi: 10.18653/v1/2021.emnlp-main.85.

[5] S. Baron-Cohen, A. M. Leslie, and U. Frith. Does the autistic child have a "theory of mind"? *Cognition*, 21(1):37–46, 1985.

[6] M. Bortoletto, C. Ruhdorfer, A. Abdessaied, L. Shi, and A. Bulling. Limits of theory of mind modelling in dialogue-based collaborative plan acquisition. In *Proc. 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1–16, 2024.

[7] M. Bortoletto, C. Ruhdorfer, L. Shi, and A. Bulling. Benchmarking mental state representations in language models. In *ICML 2024 Workshop on Mechanistic Interpretability*, 2024. URL https://openreview.net/forum?id=yEwEVoH9Be.

[8] M. Bortoletto, L. Shi, and A. Bulling. Neural reasoning about agents' goals, preferences, and actions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 456–464, 2024.

[9] H. H. Clark. *Using language*. Cambridge University Press, 1996.

[10] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, volume 1, pages 886–893. IEEE, 2005.

[11] P. Doshi, X. Qu, A. Goodie, and D. Young. Modeling recursive reasoning by humans using empirically informed interactive pomdps. In *Proceedings of the International Conference on Autonomous Agents and Multiagent Systems*, pages 1223–1230, 2010.

[12] J. Duan, S. Yu, N. Tan, L. Yi, and C. Tan. Boss: A benchmark for human belief prediction in object-context scenarios. *arXiv preprint arXiv:2206.10665*, 2022.

[13] L. Fan, S. Qiu, Z. Zheng, T. Gao, S.-C. Zhu, and Y. Zhu. Learning triadic belief dynamics in nonverbal communication from videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7312–7321, 2021.

[14] K. Floyd. *Interpersonal communication*. McGraw-Hill, 2011.

[15] K. Gandhi, G. Stojnic, B. M. Lake, and M. R. Dillon. Baby intuitions benchmark (bib): Discerning the goals, preferences, and actions of others. *Advances in Neural Information Processing Systems*, 34:9963–9976, 2021.

[16] N. D. Goodman, C. L. Baker, E. B. Bonawitz, V. K. Mansinghka, A. Gopnik, H. Wellman, L. Schulz, and J. B. Tenenbaum. Intuitive theories of mind: A rational approach to false belief. In *Proceedings of the twenty-eighth annual conference of the cognitive science society*, volume 6. Cognitive Science Society Vancouver, 2006.

[17] A. Graves and J. Schmidhuber. Framewise phoneme classification with bidirectional lstm networks. In *Proceedings of the IEEE International Joint Conference on Neural Networks*, volume 4. IEEE, 2005.

[18] N. Gurney and D. V. Pynadath. Robots with theory of mind for humans: A survey. In *Proceedings of the IEEE International Conference on Robot and Human Interactive Communication*. IEEE, 2022.

[19] Y. Han and P. Gmytrasiewicz. Learning others' intentional models in multi-agent settings using interactive pomdps. *Advances in Neural Information Processing Systems*, 31, 2018.

[20] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard,

T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, Sept. 2020. doi: 10.1038/s41586-020-2649-2. URL https://doi.org/10.1038/s41586-020-2649-2.

[21] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.

[22] D. Hendrycks and K. Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.

[23] J. D. Hunter. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):90–95, 2007. doi: 10.1109/MCSE.2007.55.

[24] G. Jocher. Yolov5 by ultralytics, 2020. URL https://github.com/ultralytics/yolov5.

[25] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In Y. Bengio and Y. LeCun, editors, *International Conference on Learning Representations*, 2015. URL http://arxiv.org/abs/1412.6980.

[26] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*. OpenReview.net, 2017. URL https://openreview.net/forum?id=SJU4ayYgl.

[27] J. J. Lee, F. Sha, and C. Breazeal. A bayesian theory of mind approach to nonverbal communication. In *Proceedings of the ACM/IEEE International Conference on Human-Robot Interaction*, pages 487–496. IEEE, 2019.

[28] A. Liu, H. Zhu, E. Liu, Y. Bisk, and G. Neubig. Computational language acquisition with theory of mind. In *International Conference on Learning Representations*, 2023.

[29] A. Netanyahu, T. Shu, B. Katz, A. Barbu, and J. B. Tenenbaum. Phase: Physically-grounded abstract social events for machine social perception. In *Proceedings of the aaai conference on artificial intelligence*, volume 35, pages 845–853, 2021.

[30] D. Nguyen, P. Nguyen, H. Le, K. Do, S. Venkatesh, and T. Tran. Learning theory of mind via dynamic traits attribution. In *Proceedings of the International Conference on Autonomous Agents and Multiagent Systems*, pages 954–962, 2022.

[31] D. Nguyen, P. Nguyen, H. Le, K. Do, S. Venkatesh, and T. Tran. Memory-augmented theory of mind network. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(10):11630–11637, Jun. 2023. doi: 10.1609/aaai.v37i10.26374.

[32] T. N. Nguyen and C. Gonzalez. Cognitive machine theory of mind. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, 2020.

[33] T. pandas development team. pandas-dev/pandas: Pandas, Feb. 2020. URL https://doi.org/10.5281/zenodo.3509134.

[34] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems*, 32, 2019.

[35] K. Pearson. Liii. on lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin philosophical magazine and journal of science*, 2(11):559–572, 1901.

[36] D. Premack and G. Woodruff. Does the chimpanzee have a theory of mind? *Behavioral and Brain Sciences*, 1(4):515–526, 1978.

[37] X. Puig, T. Shu, S. Li, Z. Wang, Y.-H. Liao, J. B. Tenenbaum, S. Fidler, and A. Torralba. Watch-and-help: A challenge for social perception and human-ai collaboration. In *International Conference on Learning Representations*, 2020.

[38] N. Rabinowitz, F. Perbet, F. Song, C. Zhang, S. A. Eslami, and M. Botvinick. Machine theory of mind. In *International Conference on Machine Learning*, pages 4218–4227. PMLR, 2018.

[39] M. Sclar, G. Neubig, and Y. Bisk. Symmetric machine theory of mind. In *International Conference on Machine Learning*, pages 19450–19466. PMLR, 2022.

[40] T. Shu, A. Bhandwaldar, C. Gan, K. Smith, S. Liu, D. Gutfreund, E. Spelke, J. Tenenbaum, and T. Ullman. Agent: A benchmark for core psychological reasoning. In *International Conference on Machine Learning*, pages 9614–9625. PMLR, 2021.

[41] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.

[42] E. Takmaz, N. Brandizzi, M. Giulianelli, S. Pezzelle, and R. Fernandez. Speaking the language of your listener: Audience-aware adaptation via plug-and-play theory of mind. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 4198–4217, Toronto, Canada, July 2023. Association for Computational Linguistics. URL https://aclanthology.org/2023.findings-acl.258.

[43] M. Tomasello. *Origins of human communication*. MIT Press, 2010.

[44] T. Ullman. Large language models fail on trivial alterations to theory-of-mind tasks. *arXiv preprint arXiv:2302.08399*, 2023.

[45] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 2017.

[46] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, İ. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.

[47] Y. Wang, F. Zhong, J. Xu, and Y. Wang. ToM2C: Target-oriented Multi-agent Communication and Cooperation with Theory of Mind. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=2t7CkQXNpuq.

[48] Wes McKinney. Data Structures for Statistical Computing in Python. In Stéfan van der Walt and Jarrod Millman, editors, *Proceedings of the 9th Python in Science Conference*, pages 56 – 61, 2010. doi: 10.25080/Majora-92bf1922-00a.

[49] H. Wimmer and J. Perner. Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children's understanding of deception. *Cognition*, 13(1):103–128, 1983.

## A    Appendix

### A.1    CNN Architecture

Our convolutional neural network (CNN) feature extractor consists of three convolutional layer blocks, followed by ReLU activation and max pooling. The convolutional layers have 16, 32 and 64 filters, kernel size three and stride one. For the max pooling layers we employ stride two. The last max pooling layer is global. This complete CNN stack is followed by a projection layer with size 64.

### A.2    MToMnet for TBD

The Triadic Belief Dynamics (TBD) dataset poses the task of predicting belief dynamics for five minds $m^i$, $i = \{1, 2, 12, 21, c\}$, into four possible classes: *appear*, *disappear*, *update* and *null*. In our MToMnet models, we use a separate fully connected classification layer for each mind (cf. Figure 1 in the main paper). Specifically, classification layers for $m^1$ and $m^{12}$ take $z_1$ as input, whereas classification layers for $m^2$ and $m^{21}$ take $z_2$ as input. Since $m^c$ represents the common ground between both minds, $z_1$ and $z_2$ are first aggregated by performing element-wise multiplication and then fed into the classification layer. Our Base-MToMnet, by design, does not incorporate explicit ToM modelling. As a result, $z_1$ and $z_2$ are not computed, and the predictions for $m^i$ are generated by feeding the fully connected layers with $h_1$ and $h_2$.

### A.3    Model Parameter Counts

Table 3 shows the number of parameter for each deep learning baseline we compare to and for each of our MToMnet variant. As can be seen from the table, our MToMnet has a significantly smaller number of parameters when compared to the baselines. This discrepancy becomes especially noticeable in the case of BOSS, as our MToMnet variants have less than 3% of the parameters compared to the baselines (CNN, CNN+LSTM, CNN+GRU, CNN+Conv1D) – approximately 460k versus 21M.

Baselines for TBD are based on a ResNet50, which has approximately 23 million parameters. However, most of the ResNet50 layers are frozen, reducing the count of trainable parameters.

| Model | Number of parameters | |
|---|---|---|
| | **BOSS** | **TBD** |
| CNN | 21,411,270 | – |
| CNN+LSTM | 24,166,550 | – |
| CNN+GRU | 23,473,302 | – |
| CNN+Conv1D | 23,544,470 | – |
| CNN | – | 162,457 |
| CNN+HOG-LSTM | – | 1,254,997 |
| CNN+HOG+Mem | – | 164,057 |
| Feats+Mem | – | 985,024 |
| -MToMnet | | |
| Base | 452,374 | 465,716 |
| DB | 452,374 | 465,716 |
| CG$\|$ | 460,886 | 474,228 |
| CG$\oplus$ | 460,886 | 474,228 |
| CG$\otimes$ | 460,886 | 474,228 |
| CG$\odot$ | 493,654 | 506,996 |
| IC$\|$ | 452,630 | 465,972 |
| IC$\oplus$ | 452,630 | 465,972 |
| IC$\otimes$ | 452,630 | 465,972 |
| IC$\odot$ | 485,398 | 498,740 |

**Table 3**: Number of Parameters for models evaluated on BOSS and TBD.

Our MToMnet has approximately one third of the parameters in CNN+HOG-LSTM and one half of Feats+Mem's parameters. CNN and CNN+HOG+Mem have less trainable parameters than MToMnet, but they also achieve much lower F1 scores (see Table 1, main paper).

### A.4    Code

Our code and bounding box annotations are available at https://git.hcics.simtech.uni-stuttgart.de/public-projects/mtomnet.

### A.5    Infrastructure & Tools

We ran our experiments on a server running Ubuntu 22.04, equipped with NVIDIA Tesla V100-SXM2 GPUs with 32GB of memory and Intel Xeon Platinum 8260 CPUs. We trained our models using PyTorch [34] with 1, 42 and 123 as random seeds. All models were trained on a single GPU card, using a batch size of four for BOSS and 64 for TBD. We used the Adam optimiser [25] with $\beta_1 = 0.9$, $\beta_2 = 0.99$, $\epsilon = 10^{-8}$, and a learning rate of $\eta = 5 \cdot 10^{-4}$. We did not perform any exhaustive hyper-parameter tuning as we noticed small variability for different configurations.

We performed our data analysis using NumPy [20], Pandas [33, 48], and SciPy [46]. Figures were made using Matplotlib [23].

### A.6    Modelling Person-Specific Features – Additional PCA Examples

Figure 9 shows additional examples of PCA on the latent representations $h_1$ and $h_2$ for each of our MToMnet variants. As the plots show, there is a clear disentanglement of principal components both on BOSS and TBD, for all MToMnet variants.

### A.7    Ethical Impact

While our work lays the foundation and remains distant from specific applications or immediate societal impact, we acknowledge that assertions related to modelling and predicting mental states can carry substantial ethical implications. Mishandling such information has
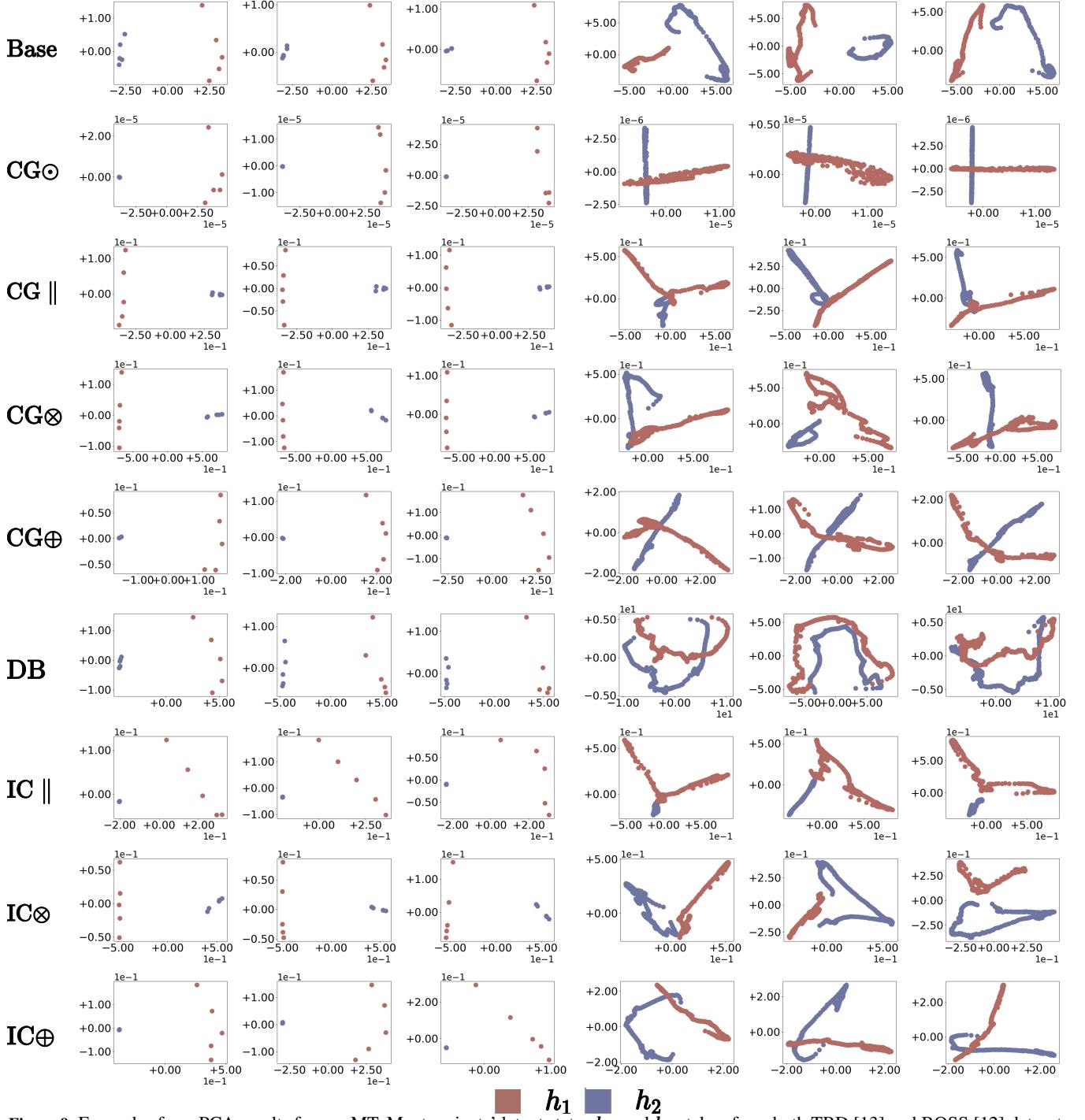
**Figure 9**: Examples from PCA results for our MToMnet variants' latent states $h_1$ and $h_2$, taken from both TBD [13] and BOSS [12] datasets.

the potential to result in the inappropriate use of personal data, reinforce biases or misunderstand complex psychological traits, potentially leading to unintended consequences on individuals' well-being.