

Human Social Interaction Modeling Using Temporal Deep Networks

Mohamed R. Amer
SRI International
Mohamed.Amer@sri.com

Behjat Siddique
SRI International
Behjat.Siddique@sri.com

Amir Tamrakar
SRI International
Amir.Tamrakar@sri.com

David A. Salter
SRI International
David.Salter@sri.com

Brian J. Lande
UC Santa Cruz
brianlande@soe.ucsc.edu

Darius Mehri
UC Berkeley
darius_mehri@berkeley.edu

Ajay Divakaran
SRI International
Ajay.Divakaran@sri.com

ABSTRACT

We present a novel approach to computational modeling of social interactions based on modeling of essential social interaction predicates (ESIPs) such as joint attention and entrainment. Based on sound social psychological theory and methodology, we collect a new “Tower Game” dataset consisting of audio-visual capture of dyadic interactions labeled with the ESIPs. We expect this dataset to provide a new avenue for research in computational social interaction modeling. We propose a novel joint Discriminative Conditional Restricted Boltzmann Machine (DCRBM) model that combines a discriminative component with the generative power of CRBMs. Such a combination enables us to uncover actionable constituents of the ESIPs in two steps. First, we train the DCRBM model on the labeled data and get accurate (76%-49% across various ESIPs) detection of the predicates. Second, we exploit the generative capability of DCRBMs to activate the trained model so as to generate the lower-level data corresponding to the specific ESIP that closely matches the actual training data (with mean square error 0.01-0.1 for generating 100 frames). We are thus able to decompose the ESIPs into their constituent actionable behaviors. Such a purely computational determination of how to establish an ESIP such as engagement is unprecedented.

Categories and Subject Descriptors

H.1.2 [Models and Principles]: User/Machine Systems—*Human information processing*; I.2.10 [Artificial Intelligence]: Vision and Scene Understanding

General Terms

Algorithms, Theory, Human Factors

Keywords

Hybrid Models; Deep Learning; DCRBMs; Social Interaction; Computational Social Psychology; Tower Game Dataset;

1. INTRODUCTION

This research brings together multiple disciplines to explore the problem of social interaction modeling. The goal of this work is to leverage research in social psychology, computer vision, signal processing, and machine learning to better understand human social interactions.

As an example application, consider aid-workers or medical personnel deployed in a foreign country. During the course of their deployment, these workers often have to interact with people with whom they share little in common in terms of language, customs and culture. Reducing friction as well as increasing engagement between the workers and the populations they encounter can have an important bearing on the success of their mission. Therefore the ability to impart such professionals, with a general cross-cultural competency which would enable them to smoothly interact with the foreign populations they encounter would be extremely useful. With such an application in mind, we focus on identifying and automatically detecting predicates that facilitate social interactions irrespective of the cultural context. Since our interests lie in aspects of social interactions that reduce conflict and build trust, we focus on social predicates that support rapport: joint attention, temporal synchrony, mimicry, and coordination.

Our orientation to social sensing departs significantly from existing methods [4, 39] that focus on inferring internal or hidden mental states. Instead, inspired by a growing body of research [26, 30, 55], we focus on the process of social interaction. This research argues that social interaction is more than the meeting of two minds, with an additional emphasis on the cognitive, perceptual and motor explanations of the joint and coordinated actions that occur as part of these interactions [37]. Our approach is guided by two key insights. The first is that apart from inferring the mental state of the other, social interactions require individuals to attend each other’s movements, utterances and context to coordinate actions jointly with each other [46]. The second insight is that social interactions involve reciprocal acts, joint behaviors

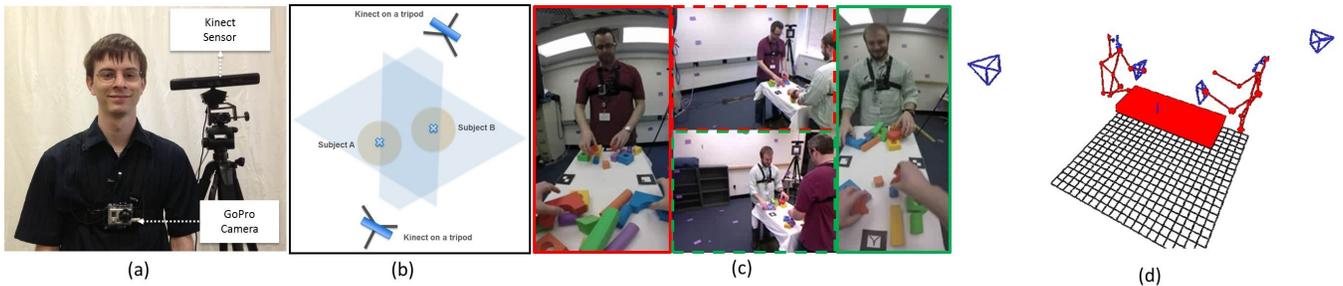


Figure 1: (a) Our capture setup which includes a GoPro camera mounted on each participant’s chest and a Kinect mounted on a tripod. (b) An overhead view of our capture setup involving the two participants. (c) Sample Data Collected: The image outlined in solid red shows the image captured from the GoPro camera mounted on player A (green shirt), while the image outlined in dashed red shows the image captured from the Kinect behind player A and is used to track the upper body of player B (red shirt). Similarly the image outlined in solid green is the image captured from the GoPro mounted on player B and the image outlined in dashed green is the image captured from the Kinect behind player B. (d) A view of our collected data projected in a unified coordinate framework.

along with nested events (e.g. speech, eye gaze, gestures) at various timescales and therefore demand adaptive and cooperative behaviors of their participants [14].

Using the work of [10] as a starting point, which emphasizes the interactive and cooperative aspects of the social interactions, we focus on detecting rhythmic coupling (also known as entrainment and attunement), mimicry (behavioral matching), movement simultaneity, kinematic turn taking patterns, and other measurable features of engaged social interaction. We established that behaviors such as *joint attention* and *entrainment* were the essential predicates of social interaction (ESIPs). With this in mind we focus on developing computational models of social interaction, that utilize multimodal sensing and temporal deep learning models to detect and recognize these ESIPs as well as discover their actionable constituents.

Over the past decade, the fields of computer vision and machine learning have made significant advances. Furthermore, with the availability of complex sensors like Kinect, researchers are able to accurately track full human body poses [47]. This allowed for many different applications in such as activity recognition [42], facial feature tracking [13], and multimodal event detection [22].

The sophistication of our problem requires a machine learning algorithm capable of jointly recognizing, correlating features, and generating multimodal data of dyadic social interactions. Discriminative models focus on maximizing the separation between classes, however, they are often uninterpretable. On the other hand, generative models focus solely on modeling distributions and are often unable to incorporate higher level knowledge. Hybrid models tend to address these problems by combining the advantages of discriminative and generative models. They encode higher level knowledge as well as model the distribution from a discriminative perspective. We propose a novel hybrid model that allows us to recognize classes, correlate features, and generate social interaction data.

This paper proposes new approach to machine learning that answers questions posed by social psychology. Our approach to social sensing is multimodal and attempts to detect the existence of features of social interaction, social interaction itself, and the qualitative and dynamic features of social interaction. We took a multimodal approach be-

cause humans must solve a variety of binding problems to effectively coordinate action. Coordination must span everything from postural sways, eye gazes, head pose, gestures, lexical choice, verbal pitch and intonation, etc.

Our contributions are 3-fold:

- A new problem of computational modeling of essential social interaction predicates (ESIPs). Starting from a socio-psychological framework, we demonstrate the use of multimodal sensors and temporal deep learning models to uncover actionable constituents of ESIPs.
- A new dataset, *Tower Game Dataset*, for analyzing social interaction predicates. The dataset consists of multimodal recordings of two players participating in a tower building game, in the process communicating and collaborating with each other. The dataset has been annotated with ESIPs and will be made publicly available. We believe that it will foster new research in the area of computational social interaction modeling.
- A novel model, Discriminative Conditional Restricted Boltzmann Machine (DCRBM), that introduces a discriminative component to Conditional Restricted Boltzmann Machines (CRBM). The discriminative component enables DCRBMs to directly learn classification models while retaining all the advantages of CRBMs, including their ability to generate missing data. Results on the *Tower Game Dataset* demonstrate that DCRBMs can effectively detect ESIPs as well decompose ESIPs into their constituent actionable behaviors.

Paper organization: In sec. 2 we discuss prior work. In sec. 3 we specify our model, then we explain inference and learning. In sec. 4 we describe our dataset and demonstrate the quantitative results of our approach. In sec. 5 we conclude.

2. RELATED WORK

Social Psychology: The study of social interactions and their associated sociological and psychological implications has received a lot of attention from social science researchers [6, 26, 37]. Early research focused on the “Theory of Mind” according to which individuals ascribe mental states to themselves and others [4], a line of thinking that largely inspired much of the initial work on affective computing. However, more recent work has shown that apart from inferring each

other’s mental states, an important challenge for participants of a social interaction is to pragmatically sustain sequences of action where the action is tightly coupled to one another via multiple channels of observable information (e.g. visible kinematic information, audible speech). In other words, social interactions require dynamically coupled interpersonal motor coordination from their participants [46]. Moreover, detecting coupled behaviors such as kinematic turn taking or simultaneity in movements can help in recognizing engaged social interactions [10].

Affective Computing: refers to the study and development of systems that can automatically detect human affect [8, 39]. Affective computing has long been an active research area due to its utility in a variety of applications that require realistic Human Computer Interaction, such as online tutoring [11] and health screenings [17]. The goal here is to detect the overall mental or emotional state of the person based on external cues. This is typically done based on speech [3], facial expressions [29], gesture/posture [35] and multimodal cues [2, 40, 48]. There has also been work on modeling activities and interactions involving multiple people [23, 43, 44]. However, most of this work deals with short duration task-oriented activities [23, 44] with a focus on their physical aspects. There has been a recent interest in modeling interactions with a focus on the rich and complex social behaviors that they elicit along with their affective impact on the participants [43].

Hybrid Models: consist of a generative model, which usually learns a feature representation of low level input, and a discriminative model for higher level reasoning. Recent work has empirically shown that generative models which learn a rich feature representation tend to outperform discriminative models that rely solely on hand-crafted features [38]. Hybrid models can be divided into three groups, joint methods [12, 24, 25, 32], iterative methods [15, 49], and staged methods [7, 21, 28, 38, 41]. Joint methods optimize a single objective function which consists of both the generative and discriminative energies. Iterative methods consist of a generative and a discriminative model that are trained in an iterative manner, influencing each other. In staged methods, both models are trained separately, with the discriminative model being trained on representations learned by the generative model. Classification is performed after projecting the samples into a fixed-dimensional space induced by the generative model.

Deep Networks: are able to learn rich features in an unsupervised manner, this is what makes deep learning very powerful. They have been successfully applied to many problems [5]. Restricted Boltzmann Machines (RBMs) form the building blocks in deep networks models [20, 45]. In [20, 45], the networks are trained using the Contrastive Divergence (CD) algorithm [19], which demonstrated the ability of deep networks to capture the distributions over the features efficiently and to learn complex representations. RBMs can be stacked together to form deeper networks known as Deep Boltzmann Machines (DBMs), which capture more complex representations. Recently, deep networks based temporal models, capable of modeling a more temporally rich set of problems have been proposed. These include Conditional RBMs (CRBMs) [54] and Temporal RBMs (TRBMs) [18, 51, 52]. CRBMs have been successfully used in both visual and audio domains. They have been used for modeling human motion [54], tracking 3D human pose [53] and phone

recognition [34]. TRBMs have been applied for transferring 2D and 3D point clouds [31], transition based dependency parsing [16], and polyphonic music generation [27].

3. APPROACH

In this section we describe our approach. We first review similar prior work, next we define our model, formulate its inference, and finally show how the model parameters are learned.

3.1 Review of Prior Models

Restricted Boltzmann Machines [45]: An RBM (Fig. 2a) defines a probability distribution $p_{\mathbf{R}}$ as a Gibbs distribution (1), where \mathbf{v} is a vector of visible nodes, \mathbf{h} is a vector of hidden nodes. $E_{\mathbf{R}}$ is the energy function and Z is the partition function which ensures that the distribution is valid. The parameters $\theta_{\mathbf{R}}$ to be learned are \mathbf{a} and \mathbf{b} the biases for \mathbf{v} and \mathbf{h} respectively and the weights W . The RBM architecture is defined as fully connected between layers, with no lateral connections. This architecture implies that \mathbf{v} and \mathbf{h} are factorial given one of the two vectors. This allows for the exact computation of $p(\mathbf{v}|\mathbf{h})$ and $p(\mathbf{h}|\mathbf{v})$.

$$\begin{aligned} p_{\mathbf{R}}(\mathbf{v}, \mathbf{h}; \theta_{\mathbf{R}}) &= \exp[-E_{\mathbf{R}}(\mathbf{v}, \mathbf{h})]/Z(\theta_{\mathbf{R}}), \\ Z(\theta_{\mathbf{R}}) &= \sum_{\mathbf{v}, \mathbf{h}} \exp[-E_{\mathbf{R}}(\mathbf{v}, \mathbf{h})], \quad (1) \\ \theta_{\mathbf{R}} &= \{\mathbf{a}, \mathbf{b}, W\} \end{aligned}$$

In case of binary valued data v_i is defined as a logistic function. In case of real valued data, v_i is defined as a multivariate Gaussian distribution with a unit covariance. A binary valued hidden layer h_j is defined as a logistic function¹. This is done because we want the hidden layer to be a sparse binary code (empirically proven to be better [52, 54]). (2) shows the probability distributions for v

$$\begin{aligned} p(v_i = 1|\mathbf{h}) &= \sigma(a_i + \sum_j h_j w_{ij}), && \text{Binary,} \\ p(v_i|\mathbf{h}) &= \mathcal{N}(a_i + \sum_j h_j w_{ij}, 1), && \text{Real,} \quad (2) \\ p(h_j = 1|\mathbf{v}) &= \sigma(b_j + \sum_i v_i w_{ij}), && \text{Binary.} \end{aligned}$$

The energy function $E_{\mathbf{R}}$ for binary v_i is defined as in (3).

$$E_{\mathbf{R}}(\mathbf{v}, \mathbf{h}) = - \sum_i a_i v_i - \sum_j b_j h_j - \sum_{i,j} v_i w_{i,j} h_j, \quad (3)$$

while, the energy function $E_{\mathbf{R}}$ is slightly modified to allow for the real valued \mathbf{v} as shown in (4).

$$E_{\mathbf{R}}(\mathbf{v}, \mathbf{h}) = - \sum_i \frac{(a_i - v_i)^2}{2} - \sum_j b_j h_j - \sum_{i,j} v_i w_{i,j} h_j \quad (4)$$

Discriminative Restricted Boltzmann Machines [24]: DRBMs are a natural extension of RBMs which have an additional discriminative term for classification. They are based on the model in [24]. DRBM (Fig. 2b) defines a prob-

¹The logistic function $\sigma(\cdot)$ for a variable x is defined as $\sigma(x) = (1 + \exp(-x))^{-1}$.

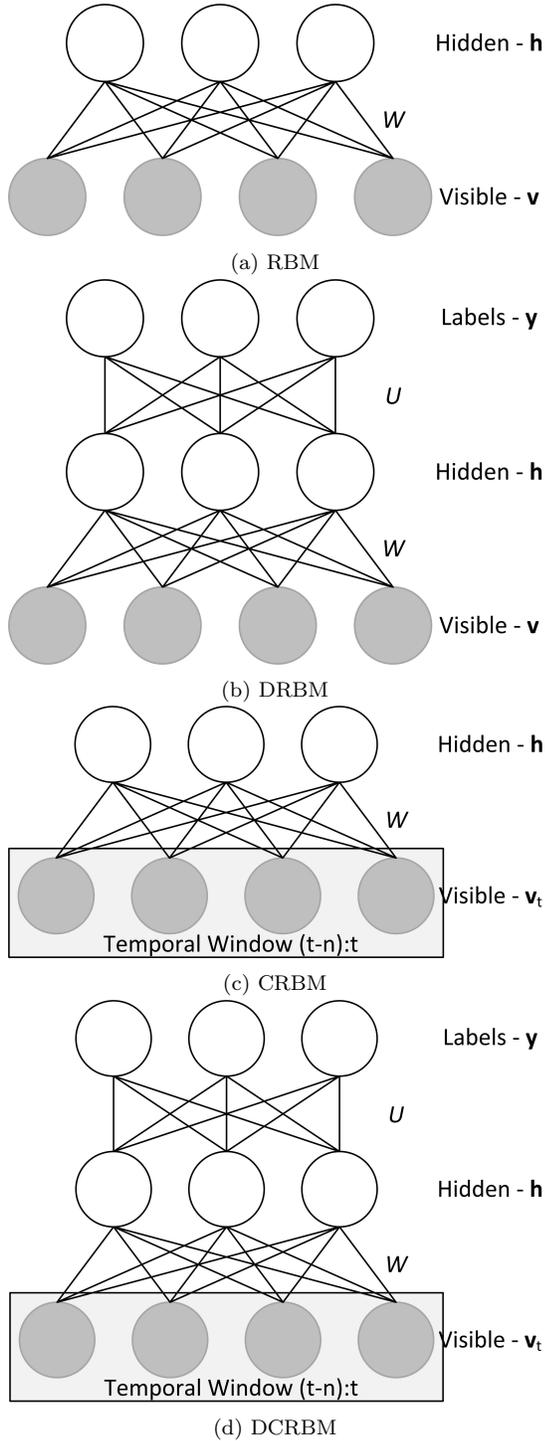


Figure 2: Deep learning models described in sections 3. (a) RBM and (c) CRBM are generative models. (b)DRBM and (d)DCRBM are discriminatively trained hybrid models.

ability distribution p_D as a Gibbs distribution (5).

$$\begin{aligned}
 p_{DR}(\mathbf{y}, \mathbf{v}, \mathbf{h} | \boldsymbol{\theta}_{DR}) &= \exp[-E_{DR}(\mathbf{y}, \mathbf{v}, \mathbf{h})] / Z(\boldsymbol{\theta}_{DR}), \\
 Z(\boldsymbol{\theta}_{DR}) &= \sum_{\mathbf{y}, \mathbf{v}, \mathbf{h}} \exp[-E_{DR}(\mathbf{y}, \mathbf{v}, \mathbf{h})] \\
 \boldsymbol{\theta}_C &= \{\mathbf{a}, \mathbf{b}, \mathbf{s}, W, U\}
 \end{aligned} \tag{5}$$

The probability distribution over the visible layer will follow the same distributions as in (2). The hidden layer \mathbf{h} is defined as a function of the labels y and the visible nodes \mathbf{v} . Also, a new probability distribution for the classifier is defined to relate the label y to the hidden nodes \mathbf{h} as in (6).

$$\begin{aligned}
 p(v_i | \mathbf{h}) &= \mathcal{N}(a_i + \sum_j h_j w_{ij}, 1), \\
 p(h_j = 1 | y_k, \mathbf{v}) &= \sigma(b_j + u_{j,k} + \sum_i v_i w_{ij}), \\
 p(y_k | \mathbf{h}) &= \frac{e^{s_k + \sum_j u_{j,k} h_j}}{\sum_{k^*} e^{s_{k^*} + \sum_j u_{j,k^*} h_j}}
 \end{aligned} \tag{6}$$

The new energy function E_{DR} is defined similar to (7),

$$\begin{aligned}
 E_D(\mathbf{y}, \mathbf{v}, \mathbf{h}) &= -\sum_i (a_i - v_i)^2 / 2 - \sum_j b_j h_j - \sum_k s_k y_k \\
 &\quad - \sum_{i,j} v_i w_{i,j} h_j - \sum_{j,k} h_j u_{j,k} y_k
 \end{aligned} \tag{7}$$

Conditional Restricted Boltzmann Machines [54]: CRBMs are a natural extension of RBMs for modeling short term temporal dependencies. A CRBM (Fig. 2c) is an RBM which takes into account history from the previous time instances $[(t-N), \dots, (t-1)]$ at time (t) . This is done by treating the previous time instances as additional inputs. Doing so does not complicate inference². A CRBM defines a probability distribution p_C as a Gibbs distribution (8).

$$\begin{aligned}
 p_C(\mathbf{v}_t, \mathbf{h}_t | \mathbf{v}_{<t}; \boldsymbol{\theta}_C) &= \exp[-E_C(\mathbf{v}_t, \mathbf{h}_t | \mathbf{v}_{<t})] / Z(\boldsymbol{\theta}_C), \\
 Z(\boldsymbol{\theta}_C) &= \sum_{\mathbf{v}, \mathbf{h}} \exp[-E_C(\mathbf{v}_t, \mathbf{h}_t | \mathbf{v}_{<t})] \\
 \boldsymbol{\theta}_C &= \{\mathbf{a}, \mathbf{b}, A, B, W\}
 \end{aligned} \tag{8}$$

The additional inputs from previous time instances are modeled as directed autoregressive edges from the past N visible nodes and the past M hidden layers, where, N does not have to be equal to M . The concatenated history vector is defined as $\mathbf{v}_{<t}$. The probability distributions are defined in (9).

$$\begin{aligned}
 p(v_i | \mathbf{h}, \mathbf{v}_{<t}) &= \mathcal{N}(a_i + \sum_n A_{n,i} v_{n,<t} + \sum_j h_j w_{ij}, 1), \\
 p(h_j = 1 | \mathbf{v}, \mathbf{v}_{<t}) &= \sigma(b_j + \sum_m B_{m,j} v_{m,<t} + \sum_i v_i w_{ij}).
 \end{aligned} \tag{9}$$

The new energy function $E_C(\mathbf{v}_t, \mathbf{h}_t | \mathbf{v}_{<t})$ in (10) is defined in a manner similar to that of the RBM (4).

$$\begin{aligned}
 E_C(\mathbf{v}_t, \mathbf{h}_t | \mathbf{v}_{<t}) &= -\sum_i \frac{(c_i - v_{i,t})^2}{2} - \sum_j d_j h_{j,t} \\
 &\quad - \sum_{i,j} v_{i,t} w_{i,j} h_{j,t},
 \end{aligned} \tag{10}$$

where,

$$c_i = a_i + \sum_n A_{n,i} v_{n,<t}, \quad d_j = b_j + \sum_m B_{m,j} v_{m,<t}.$$

²Some approximations have been made to facilitate efficient training and inference, more details are available in [54].

Note that A and B are matrices of concatenated vectors of previous time instances of \mathbf{a} and \mathbf{b} .

3.2 Model

Discriminative Conditional Restricted Boltzmann Machines: (DCRBMs) are a natural extension of CRBMs which have an additional discriminative term for classification. They are based on the model in [24], generalized to account for temporal phenomenon using CRBMs. DCRBMs (Fig. 2d) are a simpler version of the Factored Conditional Restricted Boltzmann Machines [54] and Gated Restricted Boltzmann Machines [33]. Both these models incorporate labels in learning representations, however, they use a more complicated potential which involves three way connections into factors. DCRBM defines a probability distribution p_{DC} as a Gibbs distribution (11).

$$\begin{aligned} p_{DC}(\mathbf{y}_t, \mathbf{v}_t, \mathbf{h}_t | \mathbf{v}_{<t}, \boldsymbol{\theta}_{DC}) &= \exp[-E_{DC}(\mathbf{y}_t, \mathbf{v}_t, \mathbf{h}_t | \mathbf{v}_{<t})] / Z(\boldsymbol{\theta}_{DC}), \\ Z(\boldsymbol{\theta}_{DC}) &= \sum_{\mathbf{y}_t, \mathbf{v}_t, \mathbf{h}_t} \exp[-E_{DC}(\mathbf{y}_t, \mathbf{v}_t, \mathbf{h}_t | \mathbf{v}_{<t})], \\ \boldsymbol{\theta}_{DC} &= \{\mathbf{a}, \mathbf{b}, \mathbf{s}, A, B, W, U\}. \end{aligned} \quad (11)$$

The probability distribution over the visible layer will follow the same distributions as in (6). The hidden layer \mathbf{h} is defined as a function of the labels y and the visible nodes \mathbf{v} . A new probability distribution for the classifier is defined to relate the label y to the hidden nodes \mathbf{h} is defined as in (12).

$$\begin{aligned} p(v_{i,t} | \mathbf{h}_t, \mathbf{v}_{<t}) &= \mathcal{N}(a_i + \sum_n A_{n,i} v_{n,<t} + \sum_j h_j w_{ij}, 1), \\ p(h_{j,t} = 1 | y_t, \mathbf{v}_t, \mathbf{v}_{<t}) &= \frac{\sigma(b_j + u_{j,k} + \sum_i v_{i,t} w_{ij} + \sum_m B_{m,j} v_{m,<t})}{1 + \sigma(b_j + u_{j,k} + \sum_i v_{i,t} w_{ij} + \sum_m B_{m,j} v_{m,<t})}, \\ p(y_{k,t} | \mathbf{h}) &= \frac{e^{s_k + \sum_j u_{j,y} h_j}}{\sum_{k^*} e^{s_{k^*} + \sum_j u_{j,k^*} h_j}}. \end{aligned} \quad (12)$$

The new energy function E_{DC} is defined similar to that of the DRBM (7).

$$\begin{aligned} E_{DC}(\mathbf{y}_t, \mathbf{v}_t, \mathbf{h}_t | \mathbf{v}_{<t}) &= -\sum_i (c_i - v_{i,t})^2 / 2 - \sum_j d_{j,k} h_{j,t} \\ &\quad - \sum_k s_k y_{k,t} - \sum_{i,j} v_{i,t} w_{i,j} h_{j,t} - \sum_{j,k} h_{j,t} u_{j,k} y_{k,t} \end{aligned} \quad (13)$$

where,

$$c_i = a_i + \sum_n A_{n,i} v_{n,<t}, \quad d_{j,k} = b_{j,t} + u_{j,k} + \sum_m B_{m,j} v_{m,<t}.$$

Note that A and B are matrices of concatenated vectors of previous time instances of \mathbf{a} and \mathbf{b} .

3.3 Inference and Learning

Inference: For classification we use a bottom up approach, where we maximize the posterior distribution, $p_{DC}(y_{t,k} | \mathbf{v}_t, \mathbf{v}_{<t})$, over all the labels. This is equivalent to activating the hidden layer given the visible layer \mathbf{v}_t , visible layer history $\mathbf{v}_{<t}$, and label $y_{t,k}$ as shown in (14).

$$\begin{aligned} y_t &= \arg \max_k p_{DC}(y_{t,k} | \mathbf{v}_t, \mathbf{v}_{<t}), \quad \text{where,} \\ p_{DC}(y_{t,k} | \mathbf{v}_t, \mathbf{v}_{<t}) &= \frac{e^{s_k} \prod_j (1 + e^{s_k + d_{j,y} + \sum_i v_{i,t} w_{ij}})}{\sum_{k^*} e^{s_{k^*}} \prod_j (1 + e^{s_{k^*} + d_{j,k^*} + \sum_i v_{i,t} w_{ij}})}. \end{aligned} \quad (14)$$

For generation we use a combination of top-down/bottom-up depending on the type of generation by activating the required layers given the available data, as in (12). Figures 4a and 5a show the two cases. The first case (Fig. 4a) deals with partial missing data, where we have partial data for the hidden layer v_t as well as the label y , and our goal is to generate the missing part of the v_t . The second case (Fig. 5a) is when we have a fully missing visible layer v_t and our goal is to generate it given only the class label y . For both cases we assume we have access to some history information.

Learning: Learning our model is done using Contrastive Divergence (CD) [19]. The update equations of the dynamically changing bases $\Delta \mathbf{c}$ and $\Delta \mathbf{d}$ are obtained by first updating ΔA and ΔB as in the case of the real valued CRBM (8) and then combining them with Δa and Δb .

$$\begin{aligned} \Delta w_{i,j} &\propto \langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{recon}, \\ \Delta u_{j,k} &\propto \langle y_k h_j \rangle_{data} - \langle y_k h_j \rangle_{recon}, \\ \Delta a_i &\propto \langle v_i \rangle_{data} - \langle v_i \rangle_{recon}, \\ \Delta b_j &\propto \langle h_j \rangle_{data} - \langle h_j \rangle_{recon}, \\ \Delta s_k &\propto \langle y_k \rangle_{data} - \langle y_k \rangle_{recon}, \\ \Delta A_{k,i,t-n} &\propto v_{k,t-n} \langle v_{i,t} \rangle_{data} - \langle v_{i,t} \rangle_{recon}, \\ \Delta B_{i,j,t-m} &\propto v_{i,t-m} \langle h_{j,t} \rangle_{data} - \langle h_{j,t} \rangle_{recon}, \end{aligned} \quad (15)$$

where $\langle \cdot \rangle_{data}$ is the expectation with respect to the data distribution and $\langle \cdot \rangle_{recon}$ is the expectation with respect to the reconstructed data. The reconstruction is generated by first sampling $p(h_j = 1 | \mathbf{v}, y)$ for all the hidden nodes in parallel. The visible nodes are then generated by sampling $p(v_i | \mathbf{h})$ for all the visible nodes in parallel. Finally, the label nodes are generated using $p(y | \mathbf{h})$ using (12).

4. EXPERIMENTS

In this section, we first discuss existing activity recognition and affective computing datasets. Next we describe the collection and annotation of our *Tower Game Dataset*, which contains recordings of two players building a tower and in the process engaging in a variety of interactive behaviors. Finally, we describe our experimental results on this dataset, demonstrating the effectiveness of our DCRBM model.

4.1 Datasets

Most existing activity recognition benchmarks – e.g., the Weizmann, Trecvid, PETS04, CAVIAR, IXMAS, Hollywood datasets, Olympic Sports and UCF-100 – contain relatively simple and repetitive actions involving a single person [9]. On the other hand, group activity recognition datasets such as UCLA Courtyard, UT-Interactions, Collective Activity datasets, and Volleyball dataset, lack rich social dynamics.

Other relevant datasets include the Multimodal Dyadic Behavior (MMDB) dataset [43], which focuses on analyzing dyadic social interactions between adults and children in a developmental context. This dataset was collected in a semi-structured format, where children interact with an adult examiner in a series of pre-planned games. However, due to its narrow focus on analysis of social behaviors to diagnose developmental disorders in children, we believe it is not general enough. Another dataset is the Mimicry database [50] which focuses on studying social interactions between humans with the aim of analyzing mimicry in human-human interactions. This dataset was collected in an unstructured format where the two humans talk to each other about different subjects.

There are a number of issues with the aforementioned datasets, including: (a) unnatural, acted activities in constrained scenes; (b) limited spatial and temporal coverage; (c) poor diversity of activity classes; (d) Lack of rich social interactions; (e) Narrow focus on a single behavior (e.g. mimicry); and (f) Unstructured or uncontrolled collection setup. Hence, we propose our new Tower Game Dataset to address the above issues.

Tower Game Dataset is a simple game of tower building often used in social psychology to elicit different kinds of interactive behaviors from the participants. It is typically played between two people working with a small fixed number of simple toy blocks that can be stacked to form various kinds of towers. We choose these tower games as they force the players to engage and communicate with each other in order to achieve the objectives of the game, thereby evoking behaviors such as *joint-attention* and *entrainment* from the participants. The game, due to its simplicity, allows for total control over the variables of an interaction. Due to the small number of blocks involved, the number of potential moves (actions) is limited. Also since the game involves interacting with physical objects, *joint-attention* is mediated through concrete objects. Furthermore, only two players are involved, ensuring that we can stay in the realm of dyadic interactions.

There are many different variants of the game. We settled on two variants designed to elicit maximum communication between the players, namely, (i) the **architect-builder** variant and (ii) the **distinct-objective** variant. Furthermore, in order to maximize the amount of non-verbal communication, we prohibited the participants from verbally communicating with each other.

The **architect-builder** variant involves one participant playing the role of the architect, who decides the kind of tower to build and how to build it. The second participant is the builder, who has control of all the building blocks and is the only one actually manipulating the blocks. The goal here is for the architect to communicate to the builder how to build the tower so that builder can build the desired tower.

The **distinct-objective** variant is slightly more complicated and is designed to elicit more interaction between the players. In this variant, each player is given half of the building blocks required to build the tower. Each player is also given a particular rule, restricting the configuration of the tower being built, that they are required to enforce. An example rule could be that no two blocks of the same color may be placed such that they are touching each other. To make the play interesting, each player only knows their own rule and is not aware of rule given to the other player. The rules are selected at random from a small rule book. While certain combinations of rules may result in some conflict between the objectives of the two players, this is typically not the case. However, since each player needs to adhere to their rule, it means that they will need to correct an action taken by the other if it conflicts with their rule. In the process, each player also tries to figure out the rule assigned to the other player so that the process of building the tower is more efficient. Also, when the subjects played multiple sessions of this game, the pieces used were changed and the area of the table upon which they could place blocks was reduced in size.

Capture Setup: Our sensors include a pair of Kinect cameras that record color videos, depth video and track skeletons of the players and a pair of GoPro cameras mounted on the chest of each player (Fig. 1(a)). External lapel microphones were attached to the GoPro cameras. However, the audio captured from them was used only for data synchronization purposes. Since the players were not allowed to verbally communicate with each other, very little speech (or paralinguistic) data exists.

In order to ensure optimal data capture from the Kinect cameras (i.e. minimal occlusions and optimal skeleton tracking), they were mounted on tripods facing one another, slightly to the right and back of each of the participants and slightly elevated, ensuring that each camera got an unobstructed view of the other participant. The overhead layout is shown in Fig. 1(b). These videos are of VGA resolutions (640x480) and were captured at 30Hz. The GoPro cameras were set to capture at full HD (1920x1080) resolution and at the widest angle available. They were placed on the harnesses rotated 90 degrees so as to capture the face of the other player as well as the blocks on the table (Fig. 1(c)).

In each session, the subjects play the game by standing at either end of a small rectangular table as shown in Fig. 1(c). The person supervising the data collection enters player information and other meta-data about the game session into a form and then starts recording. He/she then instructs the players to begin their game session. They first manually activate the GoPro cameras to start recording and then clap their hands before starting their sessions. These claps were used to automatically synchronize the GoPro videos with the Kinect videos. The final dataset consists of the following data types for each game session:

1. Two Kinect videos (RGB)
2. Two depth videos (depth encoded in RGB)
3. Two GoPro videos (distortion corrected)
4. Intrinsic and extrinsic calibrations for the two Kinect cameras
5. Intrinsic calibrations and video frame aligned sequences of camera poses for the GoPro cameras
6. Kinect tracked skeletons for the two participants
7. Face and head pose tracking for the two individuals from the GoPro cameras when visible
8. Eye Gaze information (3d vectors) for the two participants whenever available
9. Object positions (2d bounding boxes, not 3d positions) and tracks for all the blocks within each gaming session.

Data Annotation: Since our focus is on *joint attention* and *entrainment*, we annotated 112 videos which were divided into 1213 10-second segments indicating the presence or absence of these two behaviors in each segment. To annotate the videos, we developed an innovative annotation schema drawn from concepts in the social psychology literature [1, 6]. The annotation schema is a series of questions, that could be used as a guideline to assist the annotators. The annotation schema associates high level social interaction predicates with more objectively perceptible measures. For example, *Joint attention* is the shared focus of two individuals on a common subject and it involves eye gaze (on a person and on an object) and body language. Similarly,

entrainment is the alignment in the behavior of two individuals and it involves simultaneous movement, tempo similarity, and coordination. Each measure was rated using a low, medium, high measure for the entire 10 second segment. We hired six undergraduate sociology and psychology students to annotate the videos. The students were given a general introduction to the survey instrument and were then asked to code representative samples of the videos. The videos were annotated after ensuring that all the students as a group were annotating the sample videos accurately and reliably.

The dataset will be released with the acceptance of this paper. We will also publish a fully detailed description of the collection, capture, and annotation.

4.2 Quantitative Results

In this section we describe the set of experiments we conducted to evaluate our proposed model.

Implementation Details: For our experiments, we relied only on the skeleton features. We use the 11 joints from the upper body of the two players since the tower game almost entirely involves only upper body actions.

Using the 11 joints we extracted a set of first order static and dynamic handcrafted skeleton features. The static features are computed per frame. The features consist of, relationships between all pairs of joints of a single actor, as well as the relationships between all pairs of joints of both the actors. The dynamic features are extracted per window (a set of 300 frames). In each window, we compute first and second order dynamics (velocities and accelerations) of each joint, as well as relative velocities and accelerations of pairs of joints per actor, and across actors. The dimensionality of the static and dynamic features is (257400 D). To reduce their dimensionality we use Principle Component Analysis (PCA) (100 D), Bag-of-Words (BoW) (100 and 300 D) [36]. We also extracted Deep Learning features using RBMs and CRBMs (50 dimensions)

For the DRBM and DCRBM we used the raw joint locations normalized with respect to a selected origin point. We used the same dimensionality for both models $D(v) = 66, D(h) = 50$. For DCRBM we empirically evaluated history windows of different sizes, and found that a window of size $n = 15$ works the best.

Results: For the purpose of this paper we focused on the three ECIPs, Coordination, Simultaneous Movement, and Tempo Similarity. As a baseline we used a multi-class Support Vector Machine and the different types of features defined above to classify a certain ECIP.

We divided our evaluation into two tasks. The first task is the *Classification Task*. We use the raw features of the two players and our goal is to predict the level (strength) of the three ECIPs. Each ECIP can be *low*, *medium* or *high*, hence random classification accuracy is 33%. The data is split into two sets, a training set consisting of 70% of the instances, and a test set consisting of the remaining 30%. We performed a 5 fold cross validation to guarantee unbiased results. Figure 3 shows our average classification accuracy on the Tower Game Dataset using different features and baselines combinations as well as the results from our DCRBM model. The evaluation is done with respect to the six annotators $\{A_1, A_2, \dots, A_6\}$ as well as the mean annotation. We can see that the DCRBM model outperforms all the

Simultaneous Movement							
Classifier/Annotator	A1	A2	A3	A4	A5	A6	All
SVM + Raw Skeleton	45.20	37.45	32.61	38.39	37.06	50.47	39.52
SVM + PCA (100D)	42.83	21.26	37.36	39.35	31.76	62.73	47.84
SVM + BoW (100D)	33.17	36.55	35.46	36.90	39.25	50.02	44.27
SVM + BoW (300D)	38.48	33.46	40.79	41.60	41.23	50.08	42.84
SVM + RBM (v=66, h=50)	43.52	31.06	42.17	38.73	41.90	56.36	43.17
D-RBM (v=66, h=50)	49.17	34.83	44.02	40.17	45.47	66.55	44.02
SVM + CRBM (v=66, h=50, n=15)	48.54	33.37	43.46	34.08	44.39	66.04	42.45
D-CRBM (v=66, h=50, n=15)	55.15	40.79	48.61	45.39	50.02	70.48	49.17

(a) Simultaneous Movement

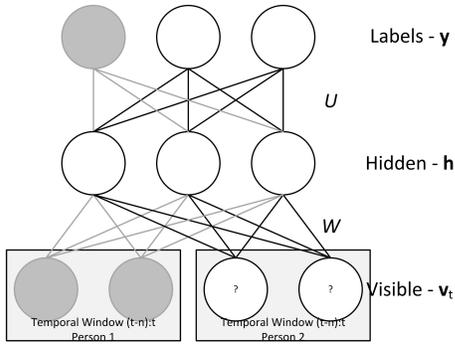
Coordination							
Classifier/Annotator	A1	A2	A3	A4	A5	A6	All
SVM + Raw Skeleton	59.00	54.56	46.91	73.31	79.21	50.57	52.21
SVM + PCA (100D)	58.75	48.01	52.25	79.68	83.03	44.61	58.16
SVM + BoW (100D)	52.15	63.57	40.06	58.03	73.25	47.38	55.76
SVM + BoW (300D)	60.30	64.94	35.09	59.53	74.88	51.32	47.54
SVM + RBM (v=66, h=50)	61.57	50.01	41.46	68.94	82.15	50.38	58.94
D-RBM (v=66, h=50)	71.57	60.15	48.93	79.38	87.57	53.03	59.75
SVM + CRBM (v=66, h=50, n=15)	70.52	47.35	43.09	78.15	87.86	53.74	59.38
D-CRBM (v=66, h=50, n=15)	85.71	68.22	57.15	82.53	89.30	55.86	62.08

(b) Coordination

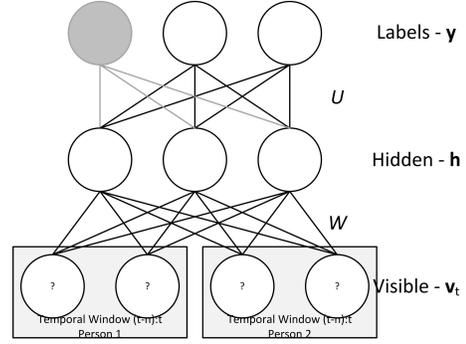
Tempo Similarity							
Classifier/Annotator	A1	A2	A3	A4	A5	A6	All
SVM + Raw Skeleton	68.65	44.55	53.50	82.35	72.00	71.65	59.29
SVM + PCA (100D)	64.86	50.48	51.77	86.24	82.21	81.32	72.81
SVM + BoW (100D)	65.10	45.60	40.49	58.03	73.25	47.38	65.64
SVM + BoW (300D)	57.04	41.29	46.06	77.21	58.05	69.64	54.36
SVM + RBM (v=66, h=50)	66.38	50.77	44.21	82.65	75.32	73.00	68.04
D-RBM (v=66, h=50)	69.55	51.24	51.77	87.03	77.86	80.24	71.32
SVM + CRBM (v=66, h=50, n=15)	68.06	51.03	45.65	86.10	77.01	77.38	71.71
D-CRBM (v=66, h=50, n=15)	71.86	55.77	54.03	88.76	85.49	83.01	76.52

(c) Tempo Similarity

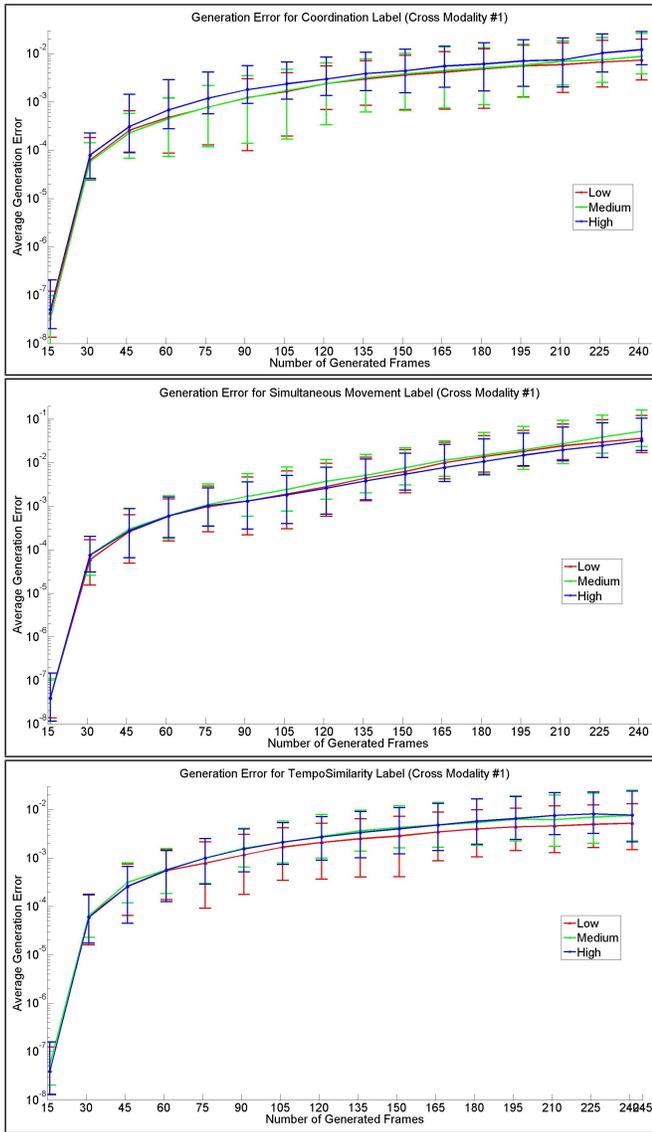
Figure 3: Average classification results on ESIPs. It is clear that the DCRBM outperforms all other baselines on the three ESIPs.



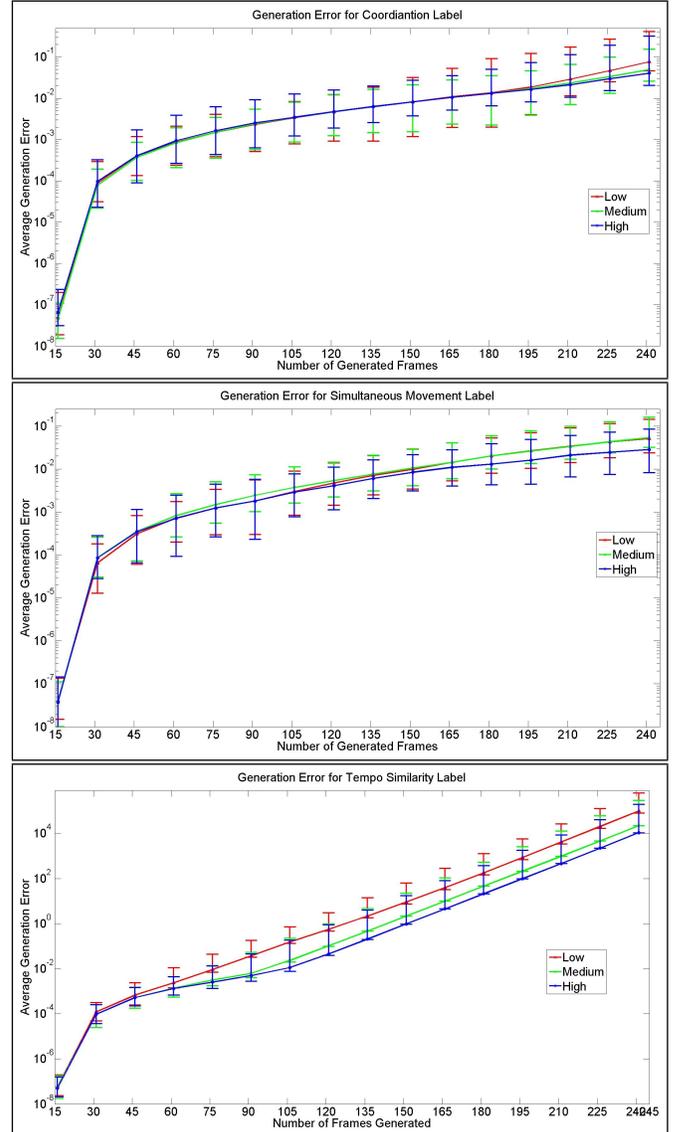
(a) Generating partial visible layer data



(a) Generating full visible layer data



(b) Average generation error (missing partial visible)



(b) Average generation error (missing full visible)

Figure 4: Average generation error for the partial visible layer by varying the generated window size for the three different ESIPs.

Figure 5: Average generation error for the full visible layer by varying the generated window size for the three different ESIPs.

other models for each of the three measures across all annotators, thereby demonstrating its effectiveness on detecting these entrainment measures. Furthermore, the DCRBM model outperforms the PCA and BoW based features which are derived from the high dimensional handcrafted features, demonstrating its ability to learn a rich representation starting from the raw skeleton features. Finally, the performance of the DCRBM model indicates that the joint learning and inference of DCRBMs is superior to the staged approach of the SVM + CRBM model.

The second task is the *Generation Task*, where we are given the class label and our goal is to generate the data (i.e. the raw features) for that label. This task allows us to visualize what the classifier has learned. For generation, we initialize the model using 15 frames for each person, and then generate sequences of lengths varying from 16 to 300 frames. We measure the mean error between the ground-truth data and the generated data for each class label over 50 video instances. For this experiment, we evaluated generated sequences of varying length using a normalized mean squared error metric defined in (16).

$$\text{Generation Error} = \left(\frac{\|\mathbf{v}_{\text{Generated}} - \mathbf{v}_{\text{Groundtruth}}\|}{\|\mathbf{v}_{\text{Groundtruth}}\|} \right)^2 \quad (16)$$

Generation is done in two different settings. In the first setting, given partial visible player data (one player’s features) as well as the class label, the goal is to generate the other player’s data. Figure 4 shows our average generation error using our DCRBM model for generating the partial visible layer. In the second setting, given only the class label, the goal is to generate the entire visible layer data (i.e. the raw features for both the players). Figure 5 shows our average generation error on using our DCRBM model for generating the full visible layer. We can see that the generation is relatively low (< 0.1) in all cases (except for Tempo Similarity³ when generating the entire visible layer data) demonstrating the effectiveness of DCRBM model for generating data. Also, the error is similar across different levels (strengths) for each measure indicating that the model is relatively stable. Finally, the error increases with the length of the generated sequence, which is expected as the possibility of variation in the ground-truth sequences increases with length.

Therefore, the classification task shows that DCRBMs can effectively detect the constituents of *entrainment* (an ESIP). Similarly, the generation task shows that DCRBMs can effectively generate raw skeleton data of the actors while modeling the different strengths of each constituent (measure).

5. CONCLUSIONS AND FUTURE WORK

We presented a novel approach to computational modeling of social interactions based on modeling of essential social interaction predicates (ESIPs) such as joint attention and entrainment. Our data collection was guided by social psychological theory and methodology. We introduce a new “Tower Game” dataset consisting of audio-visual capture of dyadic interactions labeled with the ESIPs, that should spark new research in computational social interaction modeling. We

³Tempo Similarity measures the similarity in the rate of the motion of the two players, and when data from both the players is missing generating their raw features based on whether their rate of motion is similar is extremely under constrained

proposed a novel joint Discriminative Conditional Restricted Boltzmann Machine (DCRBM) model that enabled us to uncover actionable constituents of the ESIPs in two steps. First, we trained the DCRBM model and second, used it to generate lower-level data corresponding to ESIP’s with high accuracy.

Such purely computational decomposition of ESIPs into actionable behavioral constituents is unprecedented and powerful, and offers rich possibilities for further research. First, we can substantially advance the understanding of ESIPs by uncovering mid-level predicates using the hidden layers of the DCRBM thus going beyond the current low-level feature generation to a multi-level understanding of the semantics of ESIPs. Second, we would like to extend our framework to multimodal streams that also include gaze, facial behaviors, head pose and audio so as to get a full understanding of actionable behaviors that make up the ESIPs. For instance, we may find out that coordinating gaze and gestural behavior is the most effective in establishing rapport, or perhaps not. Third, such a comprehensive multimodal and semantic model would capture the overall “rules of engagement” in a social interaction. Such a model would therefore lend itself to monitoring and training applications such as automatic assessment of the efficacy of an interaction in terms of establishment of rapport-engagement and generation of “interaction-realistic” avatar behaviors in a virtual reality environment that convey realism in terms of interaction dynamics rather than through photo or audio realism, and thus achieve immersion and engagement, as well as more efficacious human-robot interaction. We have thus laid the foundation of a computational approach that enables us to move from “folklore” based methods of establishing ESIPs to methods that are systematically arrived at through computational analysis of data from scientific observations.

Acknowledgments

This work is supported by DARPA W911NF-12-C-0001. The views, opinions, and/or conclusions contained in this paper are those of the author and should not be interpreted as representing the official views or policies, either expressed or implied of the DARPA or the DoD.

6. REFERENCES

- [1] L. Adamson and et al. Rating parent-child interactions: Joint engagement, communication dynamics and shared topics in autism, down syndrome, and typical development. *JADD*, 2012.
- [2] M. Amer, B. Siddiquie, S. Khan, A. Divakaran, and H. Sawhney. Multimodal fusion using dynamic hybrid models. In *WACV*, 2014.
- [3] M. Amer, B. Siddiquie, C. Richey, and A. Divakaran. Emotion detection in speech using deep networks. In *ICASSP*, 2014.
- [4] S. Baron-Cohen. Mindblindness: An essay on autism and theory of mind. In *MIT*, 1997.
- [5] Y. Bengio. Learning deep architectures for ai. In *FTML*, 2009.
- [6] F. J. Bernieri. Coordinated movement and rapport in teacher-student interactions. *Journal of Non-Verbal Behavior*, 1988.
- [7] A. Bosch, A. Zisserman, and M. Xavier. Scene classification using a hybrid generative/discriminative

- approach. In *TPAMI*, 2008.
- [8] R. Calvo and S. D’Mello. Affect detection: An interdisciplinary review of models, methods, and their applications. In *IEEE Transactions on Affective Computing*, 2010.
- [9] J. M. Chaquet, E. J. Carmona, and A. Fernández-Caballero. A survey of video datasets for human action and activity recognition. *CVIU*, 117(6):633 – 659, 2013.
- [10] H. De Jaegher, E. Di Paolo, and S. Gallagher. Can social interaction constitute social cognition? *Trends in Cognitive Science*, 2010.
- [11] S. D’Mello, R. W. Picard, and A. Graesser. Toward an affect-sensitive autotutor. In *IEEE Intelligent Systems*, 2007.
- [12] G. Druck and A. McCallum. High-performance semi-supervised learning using discriminatively constrained generative models. In *ICML*, 2010.
- [13] G. Fanelli, M. Dantone, J. Gall, A. Fossati, and L. V. Gool. Random forests for real time 3d face analysis. In *IJCV*, 2012.
- [14] V. Fantasia, H. D. Jaegher, and A. Fasulo. We can work it out: an enactive look at cooperation. *Frontiers in Psychology*, 2014.
- [15] A. Fujino, N. Ueda, and K. Saito. Semi-supervised learning for a hybrid generative/discriminative classifier based on the maximum entropy principle. In *TPAMI*, 2008.
- [16] N. Garg and J. Henderson. Temporal restricted boltzmann machines for dependency parsing. In *ACL*, 2011.
- [17] S. Ghosh, M. Chatterjee, and L.-P. Morency. A multimodal context-based approach for distress assessment. In *ICMI*, 2014.
- [18] C. Hausler and A. Susemihl. Temporal autoencoding restricted boltzmann machine. In *CoRR*, 2012.
- [19] G. E. Hinton. Training products of experts by minimizing contrastive divergence. In *NC*, 2002.
- [20] G. E. Hinton, S. Osindero, and Y. W. Teh. A fast learning algorithm for deep belief nets. In *NC*, 2006.
- [21] T. Jebara and et. al. Probability product kernels. In *MLR*, 2004.
- [22] H. S. Koppula and A. Saxena. Learning spatio-temporal structure from rgb-d videos for human activity detection and anticipation. In *ICML*, 2013.
- [23] T. Lan, Y. Wang, W. Yang, S. Robinovitch, and G. Mori. Discriminative latent models for recognizing contextual group activities. In *PAMI*, 2012.
- [24] H. Larochelle and Y. Bengio. Classification using discriminative restricted boltzmann machines. In *ICML*, 2008.
- [25] J. Lasserre, C. Bishop, and T. Minka. Principled hybrids of generative and discriminative models. In *CVPR*, 2006.
- [26] S. C. Levinson. On the human “interaction engine”. In *Roots of Human Sociality Culture, Cognition and Interaction*. Berg, 2006.
- [27] N. B. Lewandowski, Y. Bengio, and P. Vincent. Modeling temporal dependencies in high-dimensional sequences: Application to polyphonic music generation and transcription. In *ICML*, 2012.
- [28] X. Li, T. Lee, and Y. Liu. Hybrid generative-discriminative classification using posterior divergence. In *CVPR*, 2011.
- [29] Y. li Tian, T. Kanade, and J. F. Cohn. Recognizing action units for facial expression analysis. In *IEEE PAMI*. 2001.
- [30] M. M. Louwse, R. Dale, E. G. Bard, and P. Jeuniaux. Behavior matching in multimodal communication is synchronized. *Cognitive Science*, 2012.
- [31] L. S. M. D. Zeiler, G. W. Taylor, I. Matthews, and R. Fergus. Facial expression transfer with input-output temporal restricted boltzmann machines. In *NIPS*, 2011.
- [32] A. Mccallum, C. Pal, G. Druck, and X. Wang. Multi-conditional learning: Generative/discriminative training for clustering and classification. In *AAAI*, 2006.
- [33] R. Memisevic and G. E. Hinton. Unsupervised learning of image transformations. In *CVPR*, 2007.
- [34] A. R. Mohamed and G. E. Hinton. Phone recognition using restricted boltzmann machines. In *ICASSP*, 2009.
- [35] S. Mota and R. W. Picard. Automated posture analysis for detecting learner’s interest level. In *CVPRW*, 2003.
- [36] J. Niebles, H. Wang, and L. Fei-Fei. Unsupervised learning of human action categories using spatial-temporal words. *IJCV*, 79(3):299–318, 2008.
- [37] E. D. Paolo and H. D. Jaegher. The interactive brain hypothesis. *Frontiers in Human Neuroscience*, 2012.
- [38] A. Perina and et al. Free energy score spaces: Using generative information in discriminative classifiers. In *TPAMI*, 2012.
- [39] R. W. Picard. *Affective Computing*. MIT Press, 1995.
- [40] G. Ramirez, T. Baltrusaitis, and L. P. Morency. Modeling latent discriminative dynamic of multi-dimensional affective signals. In *ACII*, 2011.
- [41] M. A. Ranzato and et. al. On deep generative models with applications to recognition. In *CVPR*, 2011.
- [42] M. Raptis, D. Kirovski, and H. Hoppe. Real-time classification of dance gestures from skeleton animation. In *SCA*, 2011.
- [43] J. M. Rehg and et al. Decoding children’s social behavior. *CVPR*, 2013.
- [44] M. Ryoo and J. Aggarwal. Semantic representation and recognition of continued and recursive human activities. In *IJCV*, 2009.
- [45] R. Salakhutdinov and G. E. Hinton. Reducing the dimensionality of data with neural networks. In *Science*, 2006.
- [46] N. Sebanz and G. Knoblich. Prediction in joint action: What, when, and where. *Topics in Cognitive Science*, 2009.
- [47] J. Shotton and et al. Real-time human pose recognition in parts from a single depth image. In *CVPR*, 2011.
- [48] B. Siddiquie, S. Khan, A. Divakaran, and H. Sawhney. Affect analysis in natural human interactions using joint hidden conditional random fields. In *ICME*, 2013.
- [49] C. Sminchisescu, A. Kanaujia, and D. Metaxas.

- Learning joint top-down and bottom-up processes for 3d visual inference. In *CVPR*, 2006.
- [50] X. Sun, J. Lichtenauer, M. F. Valstar, A. Nijholt, and M. Pantic. A multimodal database for mimicry analysis. In *ACII*, 2011.
- [51] I. Sutskever, G. Hinton, and G. Taylor. The recurrent temporal restricted boltzmann machine. In *NIPS*, 2008.
- [52] I. Sutskever and G. E. Hinton. Learning multilevel distributed representations for high-dimensional sequences. In *AISTATS*, 2007.
- [53] G. W. Taylor and et. al. Dynamical binary latent variable models for 3d human pose tracking. In *CVPR*, 2010.
- [54] G. W. Taylor, G. E. Hinton, and S. T. Roweis. Two distributed-state models for generating high-dimensional time series. In *Journal of Machine Learning Research*, 2011.
- [55] M. Tomasello. *The Cultural Origins of Human Cognition*. Harvard University Press, 2001.