

# Deep Temporal Sigmoid Belief Networks for Sequence Modeling

Zhe Gan, Chunyuan Li, Ricardo Henao, David Carlson and Lawrence Carin

Duke University, Durham NC 27708, USA



#### INTRODUCTION

Problem of interest: Developing deep generative models for sequential data.

Main idea:

- Constructing a hierarchy of Temporal Sigmoid Belief Networks (TSBNs).
- TSBN is defined as a sequential stack of Sigmoid Belief Networks (SBNs).

#### Contributions:

- A generalization of Hidden Markov Models (HMMs) and Linear Dynamical Systems (LDS).
- A probabilistic construction of Recurrent Neural Networks (RNNs).
- Closely related to Temporal Restricted Boltzmann Machine (TRBM), but our model has a fully directed generative process.
- Can be utilized to model various data, e.g., binary, real-valued and counts.

**Challenge:** Designing scalable learning and inference algorithms. **Solution:** 

# Stochastic Variational Inference (SVI).

Design a recognition model for fast inference.

### MODEL FORMULATION

**Sigmoid Belief Network:** An SBN models a binary visible vector  $\mathbf{v} \in \{0, 1\}^M$ , in terms of binary hidden variables  $\mathbf{h} \in \{0, 1\}^J$  and weights  $\mathbf{W} \in \mathbb{R}^{M \times J}$  with

$$p(v_m = 1 | \boldsymbol{h}) = \sigma(\boldsymbol{w}_m^{\mathsf{T}} \boldsymbol{h} + c_m) \qquad p(h_j = 1) = \sigma(b_j)$$
 (1)

• SBN is closely related to RBM, which is a Markov random field with the same bipartite structure as the SBN.

**Temporal SBN:** Assume a length-T binary visible sequence, the tth time step of which is denoted  $\mathbf{v}_t \in \{0,1\}^M$ . The TSBN describes the joint probability as

$$p_{\theta}(\mathbf{V}, \mathbf{H}) = p(\mathbf{h}_1)p(\mathbf{v}_1|\mathbf{h}_1) \cdot \prod_{t=2}^{T} p(\mathbf{h}_t|\mathbf{h}_{t-1}, \mathbf{v}_{t-1}) \cdot p(\mathbf{v}_t|\mathbf{h}_t, \mathbf{v}_{t-1})$$
(2)

Each conditional distribution is expressed as

$$p(h_{jt} = 1 | \boldsymbol{h}_{t-1}, \boldsymbol{v}_{t-1}) = \sigma(\boldsymbol{w}_{1j}^{\mathsf{T}} \boldsymbol{h}_{t-1} + \boldsymbol{w}_{3j}^{\mathsf{T}} \boldsymbol{v}_{t-1} + b_j)$$
(3)

$$p(v_{mt} = 1 | \boldsymbol{h}_t, \boldsymbol{v}_{t-1}) = \sigma(\boldsymbol{w}_{2m}^{\mathsf{T}} \boldsymbol{h}_t + \boldsymbol{w}_{4m}^{\mathsf{T}} \boldsymbol{v}_{t-1} + c_m)$$
 (4)

- TSBN can be viewed as a HMM with an exponentially large state space and a highly structured transition matrix.
- TSBN allows for fast sampling of "fantasy" data from the inferred model.

#### **Extensions:**

- Modeling real-valued data:  $p(\boldsymbol{v}_t|\boldsymbol{h}_t,\boldsymbol{v}_{t-1}) = \mathcal{N}(\boldsymbol{\mu}_t,\operatorname{diag}(\boldsymbol{\sigma}_t^2))$ , where  $\mu_{mt} = \boldsymbol{w}_{2m}^{\top}\boldsymbol{h}_t + \boldsymbol{w}_{4m}^{\top}\boldsymbol{v}_{t-1} + c_m$   $\log \sigma_{mt}^2 = (\boldsymbol{w}_{2m}')^{\top}\boldsymbol{h}_t + (\boldsymbol{w}_{4m}')^{\top}\boldsymbol{v}_{t-1} + c_m'$
- Modeling counts:  $p(oldsymbol{v}_t|oldsymbol{h}_t,oldsymbol{v}_{t-1})=\pi_{m=1}^M\,y_{mt}^{v_{mt}}$ , where

$$y_{mt} = \frac{\exp(\boldsymbol{w}_{2m}^{\mathsf{T}}\boldsymbol{h}_t + \boldsymbol{w}_{4m}^{\mathsf{T}}\boldsymbol{v}_{t-1} + c_m)}{\boldsymbol{\Sigma}_{m'=1}^{M} \exp(\boldsymbol{w}_{2m'}^{\mathsf{T}}\boldsymbol{h}_t + \boldsymbol{w}_{4m'}^{\mathsf{T}}\boldsymbol{v}_{t-1} + c_{m'})}$$

• Going deep: Adding stochastic or deterministic hidden layers.

## Graphical Model

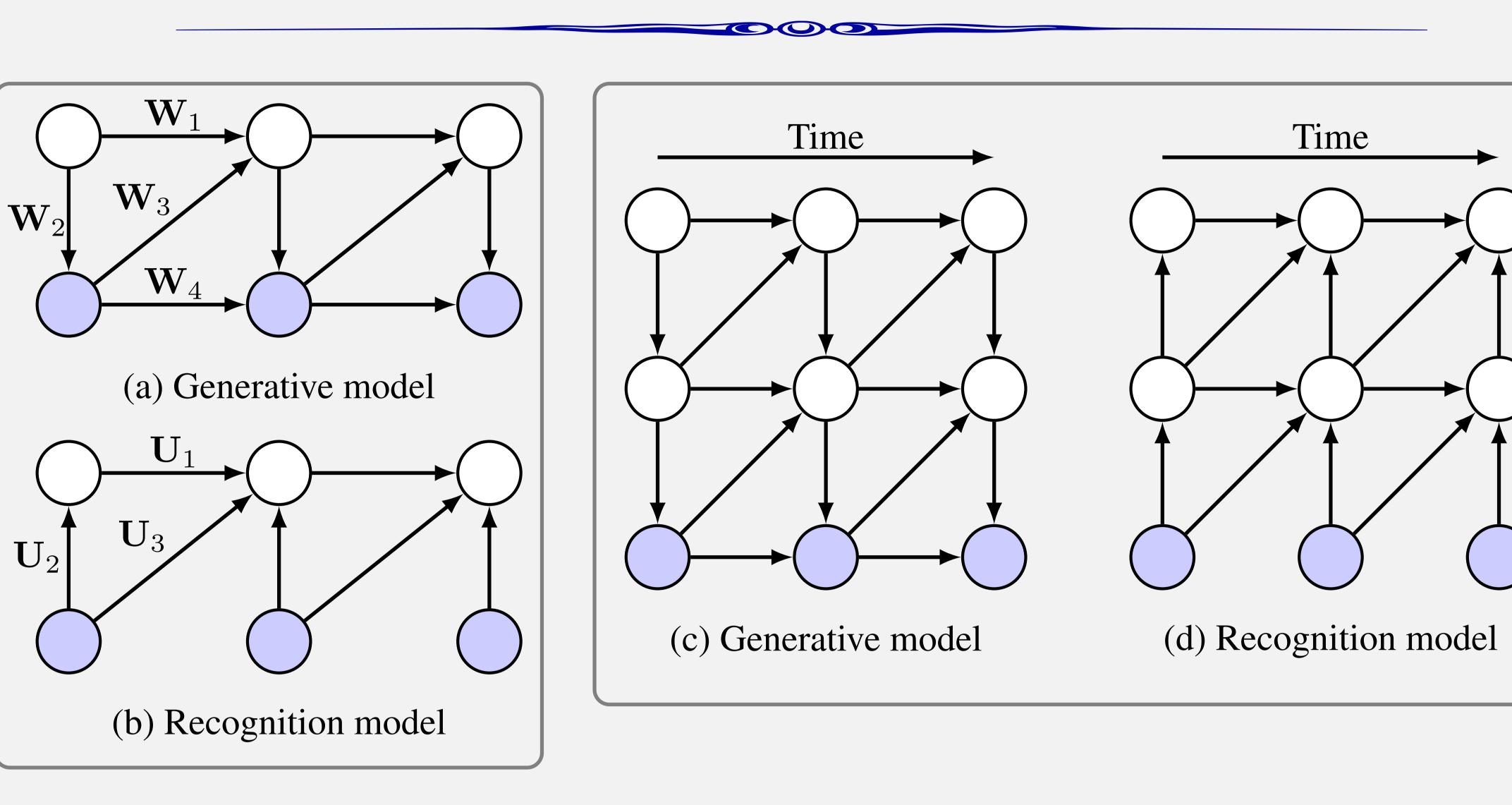


Figure: Graphical model for Deep Temporal Sigmoid Belief Network. (a,b) Generative and recognition model of TSBN. (c,d) Generative and recognition model of a two-layer TSBN.

#### SCALABLE LEARNING & INFERENCE

#### Variational Lower Bound Objective

$$\mathcal{L}(\mathbf{V}, \boldsymbol{\theta}, \boldsymbol{\phi}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{H}|\mathbf{V})}[\log p_{\boldsymbol{\theta}}(\mathbf{V}, \mathbf{H}) - \log q_{\boldsymbol{\phi}}(\mathbf{H}|\mathbf{V})]$$
 (5)

We construct the approximate posterior  $q_{\phi}(\mathbf{H}|\mathbf{V})$  as a recognition model.

$$q_{\phi}(\mathbf{H}|\mathbf{V}) = q(\mathbf{h}_1|\mathbf{v}_1) \cdot \prod_{t=2}^{I} q(\mathbf{h}_t|\mathbf{h}_{t-1}, \mathbf{v}_t, \mathbf{v}_{t-1})$$
(6)

and each conditional distribution is specified as

$$q(h_{jt} = 1 | \boldsymbol{h}_{t-1}, \boldsymbol{v}_t, \boldsymbol{v}_{t-1}) = \sigma(\boldsymbol{u}_{1j}^{\mathsf{T}} \boldsymbol{h}_{t-1} + \boldsymbol{u}_{2j}^{\mathsf{T}} \boldsymbol{v}_t + \boldsymbol{u}_{3j}^{\mathsf{T}} \boldsymbol{v}_{t-1} + d_j)$$
(7)

The recognition model is introduced to achieve fast inference.

Parameter Learning: We apply the Neural Variational Inference and Learning (NVIL) algorithm.

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{V}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{H}|\mathbf{V})} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{V}, \mathbf{H})]$$
(8)

$$\nabla_{\phi} \mathcal{L}(\mathbf{V}) = \mathbb{E}_{q_{\phi}(\mathbf{H}|\mathbf{V})}[(\log p_{\theta}(\mathbf{V}, \mathbf{H}) - \log q_{\phi}(\mathbf{H}|\mathbf{V})) \times \nabla_{\phi} \log q_{\phi}(\mathbf{H}|\mathbf{V})]$$
 (9)

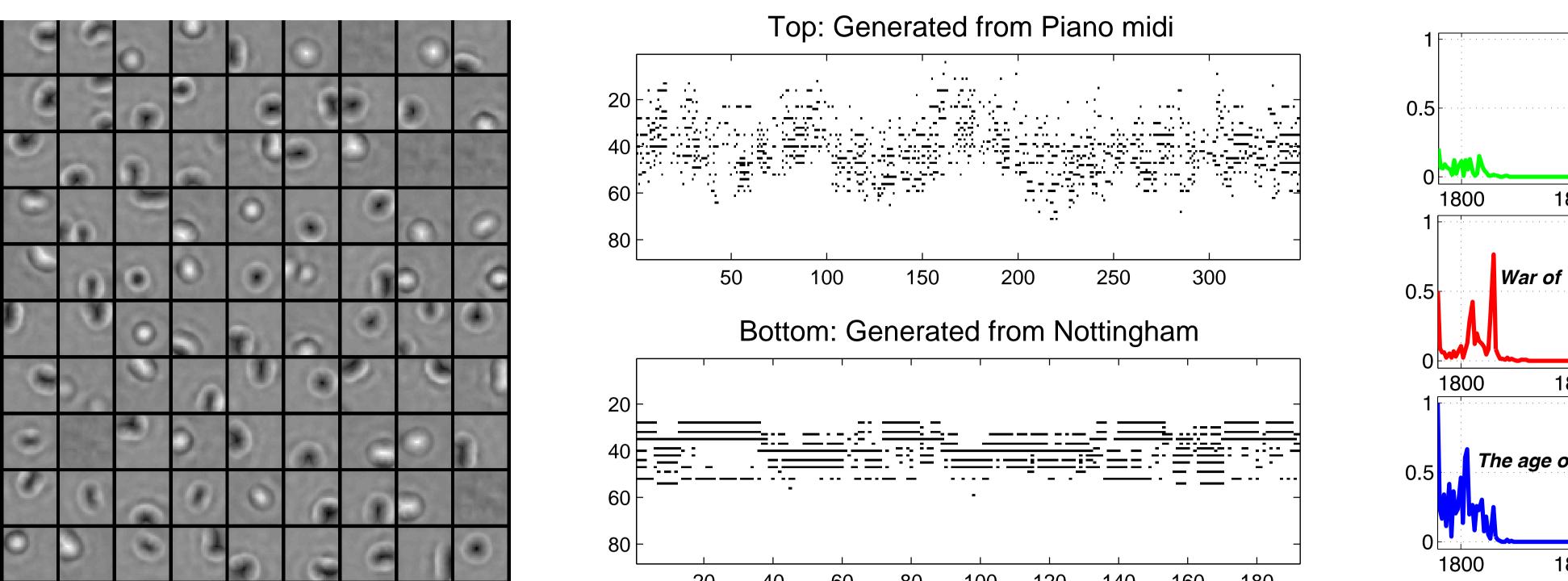
- Use Monte Carlo methods to approximate expectations.
- Variance reduction: (i) centering the learning signal by subtracting the baseline; (ii) variance normalization.
- Use RMSprop for optimization.

#### EXPERIMENTS

#### Datasets:

- Bouncing balls: Synthetic videos of 3 bouncing balls, binary valued.
- Motion capture: Walking & running sequences collected by CMU & MIT.
- **Polyphonic music:** A collection of 88-dim binary sequences, that span the whole range of piano from A0 to C8.
- State of the Union (STU): transcripts of 225 US STU addresses, from 1790 to 2014. Vocab size is 2375.

#### Qualitative Evaluation



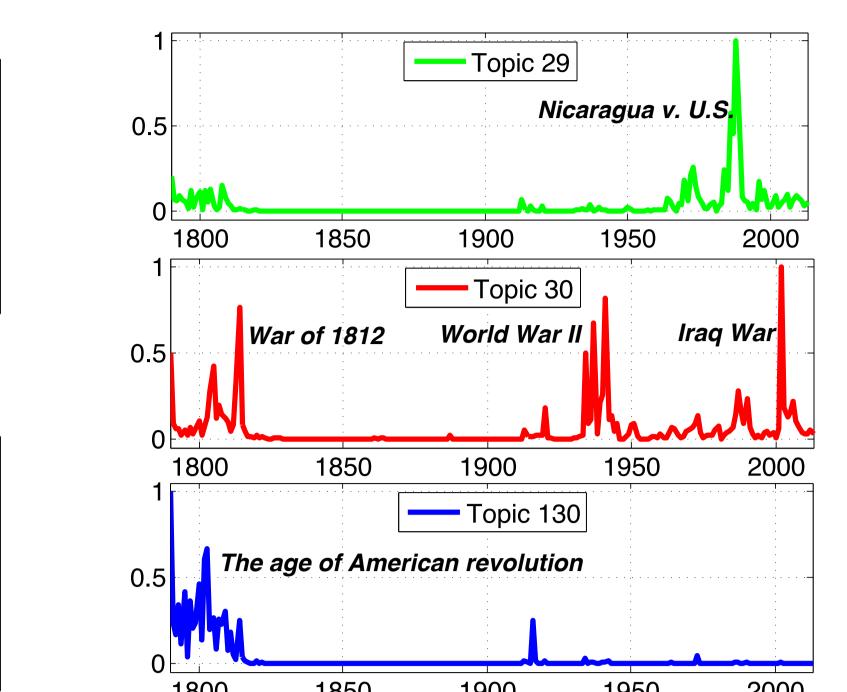


Figure: (Left) Dictionaries learned on the videos of bouncing balls. (Middle) Samples generated from TSBN trained on the polyphonic music. Each column is a sample vector of notes. (Right) Time evolving from 1790 to 2014 for three selected topics learned from the STU dataset.

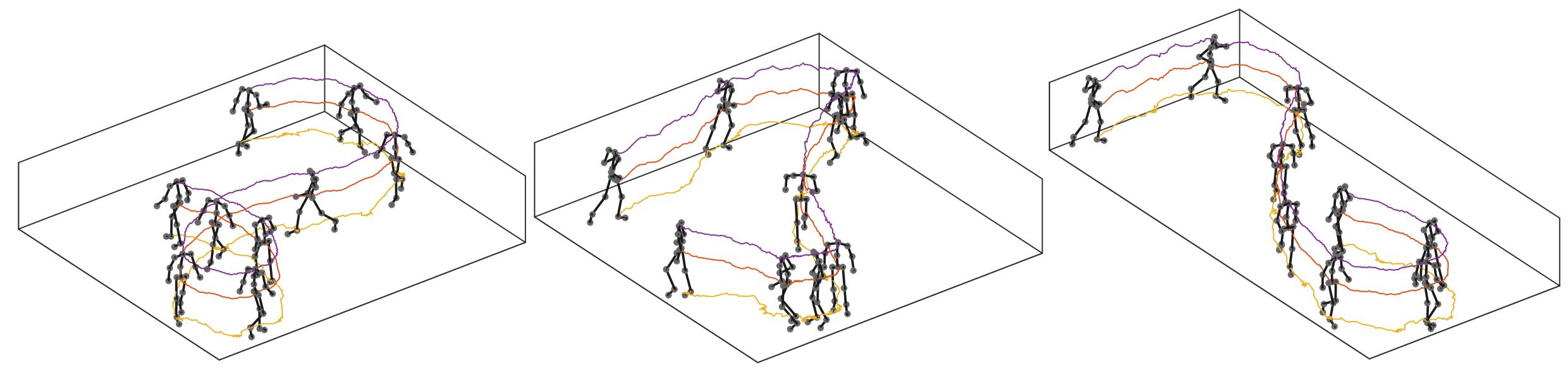


Figure: Motion trajectories generated from the TSBN trained on the motion capture dataset. (Left) Walking. (Middle) Running-running-walking. (Right) Running-walking.

#### Quantitative Evaluation

Гab	ole: Predict	tion error	for th	e bounc	ing balls.	Table: Prediction	on error for the	e motion capt
	Model	Dim	Order	Pred.	Err.	Model	Walking	Running
	DTSBN-s	100-100	2	2.79 ±	0.39	DTSBN-s	$4.40 \pm 0.28$	$2.56 \pm 0.40$
	DTSBN-d	100-100	2	$2.99 \pm$	0.42	DTSBN-d	$4.62 \pm 0.01$	$2.84 \pm 0.01$
	TSBN	100	4	$3.07 \pm$	0.40	TSBN	$5.12 \pm 0.50$	$4.85 \pm 1.26$
	TSBN	100	1	9.48 ±	0.38	HMSBN	$10.77 \pm 1.15$	$7.39 \pm 0.47$
	RTRBM	3750	1	3.88 ±	0.33	ss-SRTRBM	$8.13 \pm 0.06$	$5.88 \pm 0.05$
	SRTRBM	3750	1	$3.31 \pm$	0.33	g-RTRBM	$14.41 \pm 0.38$	$10.91 \pm 0.27$

	Table: Log-likelihood for the music data									
	Model	Piano.	Nott.	Muse.	JSB.					
	TSBN	-7.98	-3.67	-6.81	-7.48					
	RNN-NADE	-7.05	-2.31	-5.60	-5.56					
	RTRBM	-7.36	-2.62	-6.35	-6.35					
	RNN	-8.37	-4.46	-8.13	-8.71					

Table: Prediction precision for STU.

Model Dim MP PP

HMSBN 25  $0.327\pm0.002~0.353\pm0.070$ DHMSBN-s 25-25  $0.299\pm0.001~0.378\pm0.006$ GP-DPFA 100  $0.223\pm0.001~0.189\pm0.003$ DRFM 25  $0.217\pm0.003~0.177\pm0.010$ 

#### Dynamic Topic Modeling

Table: Top 8 most probable words associated with the STU topics.

Topic #29 Topic #30 Topic #130 Topic #64 Topic #70 Topic #74

family officer government generations Iraqi Philippines budget civilized country generation Qaida islands

Nicaragua warfare public recognize Iraq axis

free enemy law brave Iraqis Nazis

future whilst present crime Al Japanese

freedom gained citizens race Saddam Germans excellence lake united balanced ballistic mines

drugs safety house streets terrorists sailors

# ACKNOWLEDGEMENTS

This research was supported by ARO, DARPA, DOE, NGA and ONR.