

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 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 $v \in \{0, 1\}^M$, in terms of binary hidden variables $h \in \{0, 1\}^J$ and weights $\mathbf{W} \in \mathbb{R}^{M \times J}$ with $p(v_m = 1 | \boldsymbol{h}) = \sigma(\boldsymbol{w}_m^\top \boldsymbol{h} + c_m) \qquad p(h_j = 1) = \sigma(b_j)$ (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 th time step of which is denoted $v_t \in \{0,1\}^M$. The TSBN describes the joint probability as $p_{\theta}(\mathbf{V}, \mathbf{H}) = p(\mathbf{h}_1) p(\mathbf{v}_1 | \mathbf{h}_1) \cdot \pi_{t=2}^T p(\mathbf{h}_t | \mathbf{h}_{t-1}, \mathbf{v}_{t-1}) \cdot p(\mathbf{v}_t | \mathbf{h}_t, \mathbf{v}_{t-1})$ (2)

Each conditional distribution is expressed as

$$p(h_{jt} = 1 | \boldsymbol{h}_{t-1}, \boldsymbol{v}_{t-1}) = \sigma(\boldsymbol{w}_{1j}^{\top} \boldsymbol{h}_{t-1} + \boldsymbol{w}_{3j}^{\top} \boldsymbol{v}_{t-1})$$
$$p(v_{mt} = 1 | \boldsymbol{h}_{t}, \boldsymbol{v}_{t-1}) = \sigma(\boldsymbol{w}_{2m}^{\top} \boldsymbol{h}_{t} + \boldsymbol{w}_{4m}^{\top} \boldsymbol{v}_{t-1})$$

• 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: $p(\boldsymbol{v}_t | \boldsymbol{h}_t, \boldsymbol{v}_{t-1}) = \mathcal{N}(\boldsymbol{\mu}_t, \mathsf{diag}(\boldsymbol{\sigma}_t))$ $\mu_{mt} = w_{2m}^{\top} h_t + w_{4m}^{\top} v_{t-1} + c_m \qquad \log \sigma_{mt}^2 = (w_{2m}')^{\top} h_t + b_t + b_$

- Modeling counts: $p(m{v}_t|m{h}_t,m{v}_{t-1}) = \Pi_{m=1}^M y_{mt}^{v_{mt}}$, where

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

• Going deep: Adding stochastic or deterministic hidden layers.

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

Graphical Model

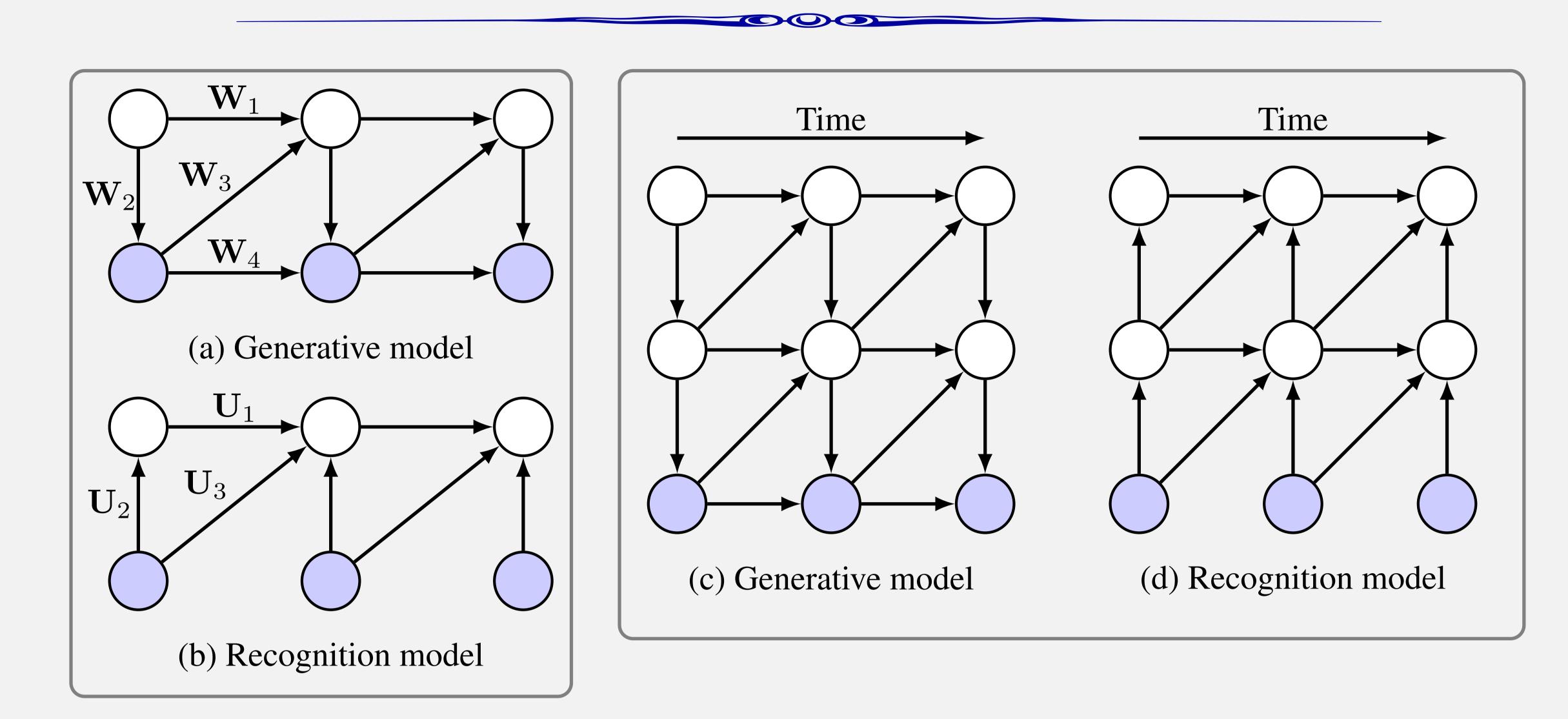


Figure: Graphical model for Deep Temporal Sigmoid Belief Network. (a,b) Generative and recognition model of TSBN. (c,d) Generative and recognition model of a two-layer TSBN.

SCALABLE LEARNING & INFERENCE

Variational Lower Bound Objective $\mathcal{L}(\mathbf{V}, \boldsymbol{\theta}, \boldsymbol{\phi}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{H}|\mathbf{V})}[\log p]$ We construct the approximate posterior $q_{\phi}(\mathbf{H}|\mathbf{V}) = q(\mathbf{h}_1|\mathbf{v}_1) \cdot$ and each conditional distribution is specified as $q(h_{jt} = 1 | h_{t-1}, v_t, v_{t-1}) = \sigma(u_1)$ The recognition model is introduced to achieve fast inference. (NVIL) algorithm

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

- 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.

(3) $+b_{i}$) (4) $+ c_m$)

$$(\boldsymbol{\sigma}_{t}^{2})$$
), where
+ $(\boldsymbol{w}_{4m}^{\prime})^{\top}\boldsymbol{v}_{t-1} + c_{m}^{\prime}$ (5)

(6)m')

$p_{\theta}(\mathbf{V}, \mathbf{H}) - \log q_{\phi}(\mathbf{H} \mathbf{V})]$	(7)
$\mathbf{r} \ q_{\phi}(\mathbf{H} \mathbf{V})$ as a recognition model	
$\prod_{t=2}^{T} q(h_t h_{t-1}, v_t, v_{t-1})$	(8)
sified ac	

$$\boldsymbol{\lambda}_{1j}^{\top}\boldsymbol{h}_{t-1} + \boldsymbol{u}_{2j}^{\top}\boldsymbol{v}_t + \boldsymbol{u}_{3j}^{\top}\boldsymbol{v}_{t-1} + d_j$$
 (9)

Parameter Learning: We apply the Neural Variational Inference and Learning

$$\log q_{\phi}(\mathbf{H}|\mathbf{V})) \times \nabla_{\phi} \log q_{\phi}(\mathbf{H}|\mathbf{V})] \quad (11)$$

Qualitative Evaluation

0	(U	6	0	5)	0		0
6	(*)	010	0	0				0
20)			· (8	2		0	
)	0	6	8	0	0			C
0	2	0	•	0		2	6	0
0	9	G	0	16	()	0	20
0	6	20	C		9		6	G
0				0	0	(0	0
0	3		(0	6	9	0
90	0	0	0	5	9 6	10	(1)	1

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.

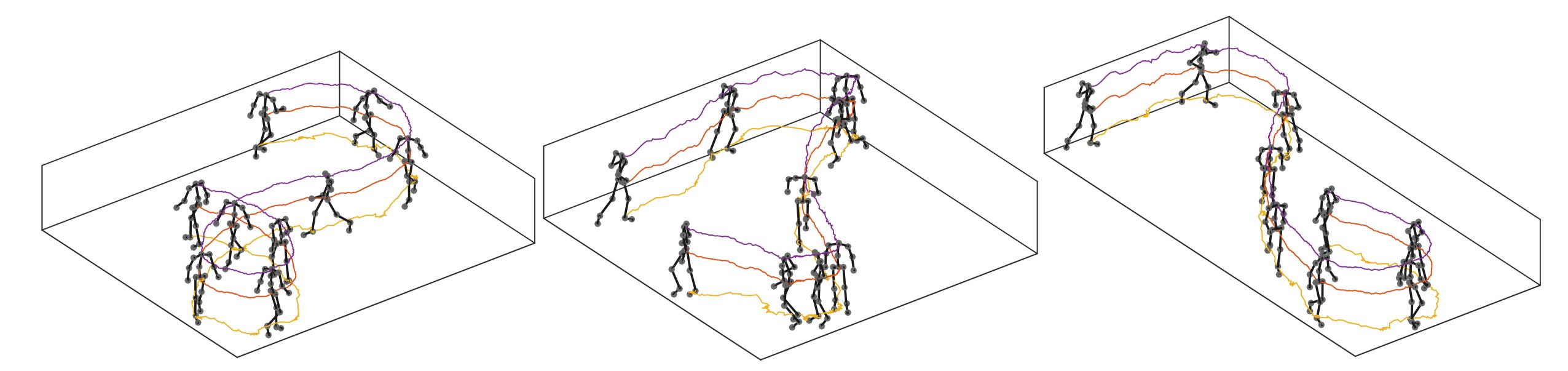


Figure: Motion trajectories generated from the TSBN trained on the motion capture dataset. (Left) Walking. (Middle) Running-running-walking. (Right) Running-walking.

Quantitative Evaluation

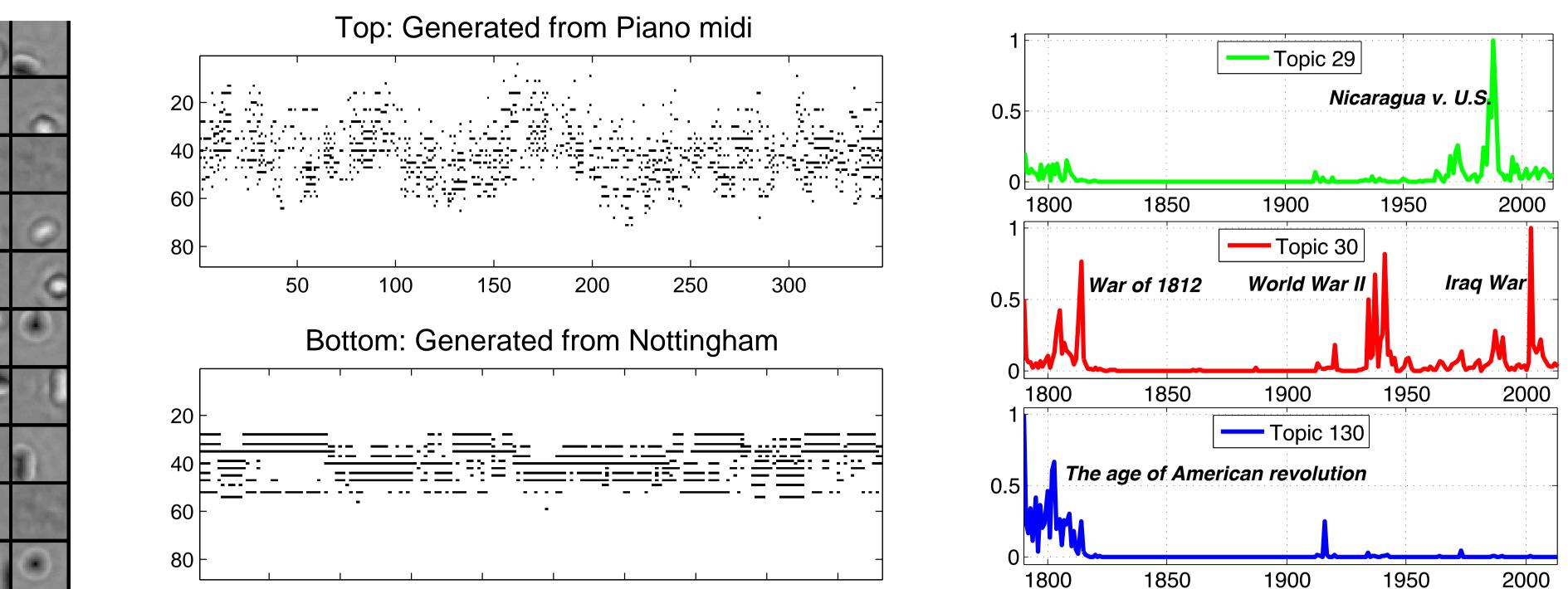
Table: Prediction error	for th	ne bouncing balls.	Table: Prediction	on error for the	e motion capture.
Model Dim	Order	Pred. Err.	Model	Walking	Running
DTSBN-s 100-100	2	2.79 ± 0.39	DTSBN-s	4.40 ± 0.28	2.56 ± 0.40
DTSBN-d 100-100	2	2.99 ± 0.42	DTSBN-d	4.62 ± 0.01	2.84 ± 0.01
TSBN 100	4	3.07 ± 0.40	TSBN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

Table: Log-likelihood for the music dataset.					c dataset.	Table: Prediction precision for STU.								
	Model	Piano.	Nott. Mı	use.	JSB.	Model	Dim		MP			F	P	
	TSBN	-7.98	-3.67 -6	.81	-7.48	HMSBN	25	0.32	27 ± 0	.002	0.3	353:	± 0.0)70
	RNN-NADE	-7.05	-2.31 -5	.60	-5.56	DHMSBN-s	25-25	0.29	99 ± 0	.001	0.3	378:	± 0.0)06
	RTRBM	-7.36	-2.62 -6	.35	-6.35	GP-DPFA	100	0.22	23±0	.001	0.3	189:	± 0.0)03
	RNN	-8.37	-4.46 -8	.13	-8.71	DRFM	25	0.2	17 ± 0	.003	0.	177:	± 0.0)10

Dynamic Topic Modeling

Table: Top
Topic #29
family
budget
Nicaragua
free
future
freedom
excellence
drugs





p 8 most probable words associated with the STU topics.

Topic $#30$	Topic #130	Topic $\#64$	Topic $\#70$	Topic $\#74$
officer	government	generations	Iraqi	Philippines
civilized	country	generation	Qaida	islands
warfare	public	recognize	lraq	axis
enemy	law	brave	Iraqis	Nazis
whilst	present	crime	AI	Japanese
gained	citizens	race	Saddam	Germans
lake	united	balanced	ballistic	mines
safety	house	streets	terrorists	sailors

ACKNOWLEDGEMENTS

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