

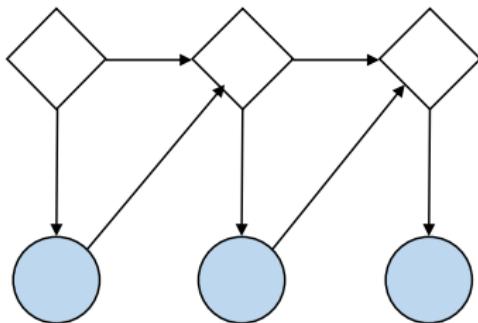
A Review on Recent Sequential Latent Variable Models

Presented by Zhe Gan

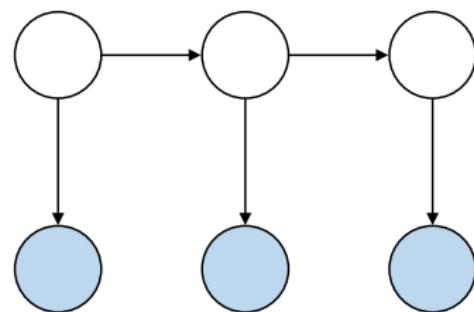
Duke University

October 7th, 2016

Motivation



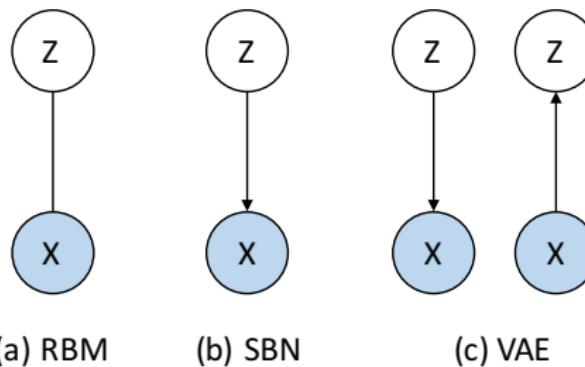
(a) RNN



(b) SSM

- **RNNs** are widely used in deep learning, e.g., LSTM, Seq2Seq Model, attention mechanism;
- **SSMs** are Bayesian probabilistic generative models, e.g., HMM, LDS (KF), EKPF;
- Can we leverage the advantage of both? – incorporate *stochastic* latent variable into RNN.

Deep Generative Models



- **RBM**: *Undirected* graphical model, binary latent variable, inference by CD;
- **SBN**: *Directed* graphical model, binary latent variable, inference by Gibbs, MCEM, NVIL;
- **VAE**: *Directed* graphical model, usually Gaussian latent variable, inference by SGVB.

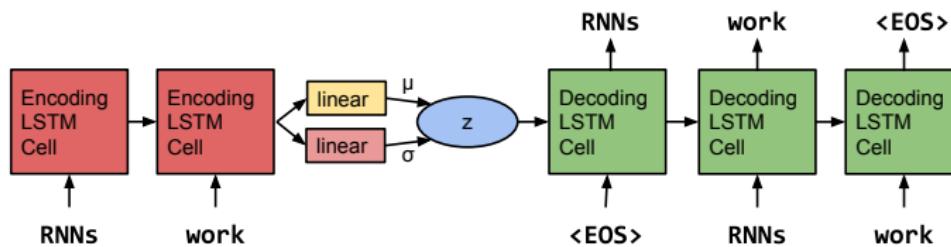
Literature Review

Category	Model	Conferences	Applications
RBM	TRBM [22, 19] RTRBM [20, 5, 17] FCRBM [21, 24] imCRBM [23]	NIPS 2006, AISTATS 2007 NIPS 2009, ICML 2012 & 2014 ICML 2009, NIPS 2011 CVPR 2010	Mocap, Bouncing ball plus Music, Weather Mocap, Facial expression Human pose tracking
SBN	TSBN, FTSBN [10, 18]	NIPS 2015, ICML 2016	plus Dynamical TM
VAE-I	VRAE [8] VRAE [6] NASM [16] VNMT [25]	ICLR 2015 workshop CoNLL 2016 ICML 2016 EMNLP 2016	Music Text modeling QA NMT
VAE-II	STORN [4] VRNN [7] NASMC [11] SRNN [9]	ICLR 2015 workshop NIPS 2015 NIPS 2015 NIPS 2016	Music Speech, Handwriting Music Speech, Music
SSM	DKF [14, 3, 13, 15] SVAE [12]	arXiv 2015 & 2016 NIPS 2016	EHR video of mouse behavior
	Poisson[2, 1, 26]	AISTATS 2015, NIPS 2016, ICDM 2016	Dynamical TM, Music

Overview

- Generating Sentences from a Continuous Space, [CoNLL 2016 \[6\]](#)
- Neural Variational Inference for Text Processing, [ICML 2016 \[16\]](#)
- A Recurrent Latent Variable Model for Sequential Data, [NIPS 2015 \[7\]](#)
- Sequential Neural Models with Stochastic Layers, [NIPS 2016 \[9\]](#)
- Neural Adaptive Sequential Monte Carlo, [NIPS 2015 \[11\]](#)

Generating Sentences from a Continuous Space Model



- **Objective:**

$$\mathcal{L} = -KL(q_{\phi}(z|x) || p(z)) + \mathbb{E}_{q_{\phi}(z|x)}[\log p_{\theta}(x|z)] \quad (1)$$

- Direct implementation fails, the training always tries to yield models that consistently set $q(z|x)$ equal to the prior $p(z)$.
- **Solution:** (i) KL cost annealing; (ii) Word dropout and historyless decoding.

Generating Sentences from a Continuous Space

Experiments

- **Language modeling:** perform roughly the same, or actually a little bit worse than **RNNLM**.

Model	Standard				Inputless Decoder			
	Train NLL	Train PPL	Test NLL	Test PPL	Train NLL	Train PPL	Test NLL	Test PPL
RNNLM	100 –	95	100 –	116	135 –	600	135 –	> 600
VAE	98 (2)	100	101 (2)	119	120 (15)	300	125 (15)	380

- **Imputing missing words:**

but now , as they parked out front and owen stepped out of the car , he could see _____ .
True: that the transition was complete . **RNNLM:** it , " i said . **VAE:** through the driver 's door .

you kill him and his ___
True: men .

RNNLM: . "

VAE: brother .

not surprising , the mothers dont exactly see eye to eye with me
True: on this matter . **RNNLM:** , i said . **VAE:** , right now .

Generating Sentences from a Continuous Space

Experiments

Sampling from the prior

100% word keep	75% word keep
" no , " he said . " thank you , " he said .	" love you , too . " <i>she put her hand on his shoulder and followed him to the door .</i>
50% word keep	0% word keep
" maybe two or two . " <i>she laughed again , once again , once again , and thought about it for a moment in long silence .</i>	<i>i hear some of of of</i> <i>i was noticed that she was holding the in in in of the the in</i>

Sampling from the posterior:

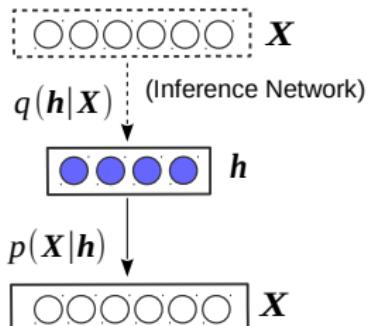
INPUT	we looked out at the setting sun .	i went to the kitchen .	how are you doing ?
MEAN	<i>they were laughing at the same time .</i>	<i>i went to the kitchen .</i>	<i>what are you doing ?</i>
SAMP. 1	<i>ill see you in the early morning .</i>	<i>i went to my apartment .</i>	<i>" are you sure ?</i>
SAMP. 2	<i>i looked up at the blue sky .</i>	<i>i looked around the room .</i>	<i>what are you doing ?</i>
SAMP. 3	<i>it was down on the dance floor .</i>	<i>i turned back to the table .</i>	<i>what are you doing ?</i>

Neural Variational Inference for Text Processing

Document Modeling

Neural Variational Document Model (NVDM)

$$\mathcal{L} = \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{X})} \left[\sum_{i=1}^N \log p_\theta(\mathbf{x}_i|\mathbf{h}) \right] - KL[q_\phi(\mathbf{h}|\mathbf{X})||p(\mathbf{h})] \quad (2)$$



- $\mathbf{h} \in \mathbb{R}^K$, $\mathbf{X} \in \mathbb{R}^{|V|}$, $\mathbf{x}_i \in \mathbb{R}^{|V|}$;
- $p_\theta(\mathbf{x}_i|\mathbf{h})$ is multinomial
- $q_\phi(\mathbf{h}|\mathbf{X})$ is Gaussian

Figure 1. NVDM for document modelling.

Neural Variational Inference for Text Processing

Document Modeling

Model	Dim	20News	RCV1
LDA	50	1091	1437
LDA	200	1058	1142
RSM	50	953	988
docNADE	50	896	742
SBN	50	909	784
fdARN	50	917	724
fdARN	200	—	598
NVDM	50	836	563
NVDM	200	852	550

(a) Perplexity on test dataset.

Word	weapons	medical	companies	define	israel	book
NVDM	guns	medicine	expensive	defined	israeli	books
	weapon	health	industry	definition	arab	reference
	gun	treatment	company	printf	arabs	guide
	militia	disease	market	int	lebanon	writing
NADE	armed	patients	buy	sufficient	lebanese	pages
	weapon	treatment	demand	defined	israeli	reading
	shooting	medecine	commercial	definition	israelis	read
	firearms	patients	agency	refer	arab	books
NVDM	assault	process	company	make	palestinian	relevent
	armed	studies	credit	examples	arabs	collection

(b) The five nearest words in the semantic space.

Space	Religion	Encryption	Sport	Policy
orbit	muslims	rsa	goals	bush
lunar	worship	cryptography	pts	resources
solar	belief	crypto	teams	charles
shuttle	genocide	keys	league	austin
moon	jews	pgp	team	bill
launch	islam	license	players	resolution
fuel	christianity	secure	nhl	mr
nasa	atheists	key	stats	misc
satellite	muslim	escrow	min	piece
japanese	religious	trust	buf	marc

Neural Variational Inference for Text Processing

Question Answering

Neural Answer Selection Model (NASM)

- A question q is associated with a set of answer sentences $\{a_1, \dots, a_n\}$, together with their judgements $\{y_1, \dots, y_n\}$, $y_m = 1$, or 0.
- Each training data point can be treated as a triple (q, a, y) .

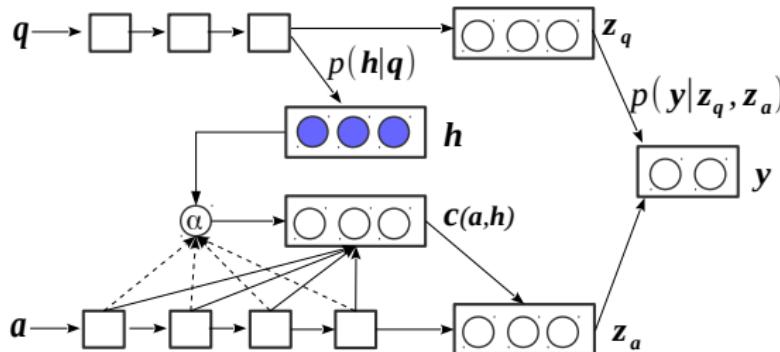


Figure 2. NASM for question answer selection.

Neural Variational Inference for Text Processing

Question Answering

- The proposed model employs two different LSTMs to embed question q and a .
- Let $s_q(j)$ and $s_a(i)$ be the state outputs of the two LSTMs, i, j be the positions of the states.
- Objective:

$$\mathcal{L} = \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{q}, \mathbf{a}, \mathbf{y})} [\log p_\theta(\mathbf{y}|\mathbf{q}, \mathbf{a}, \mathbf{h})] - KL[q_\phi(\mathbf{h}|\mathbf{q}, \mathbf{a}, \mathbf{y}) || p_\theta(\mathbf{h}|\mathbf{q})] \quad (3)$$

Neural Variational Inference for Text Processing

Question Answering

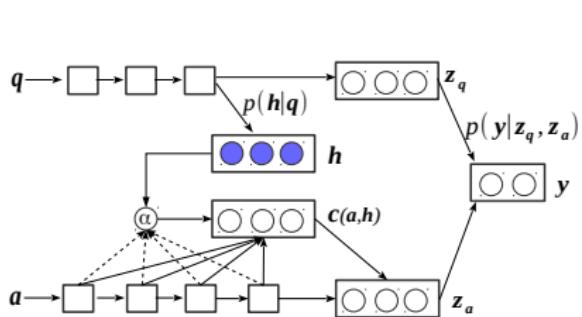


Figure 2. NASM for question answer selection.

(2) Generative Model

$$p_\theta(\mathbf{h}|\mathbf{q}):$$

$$\lambda_\theta = \tanh(\mathbf{W}_1 s_q(|\mathbf{q}|) + \mathbf{b}_1) \quad (39)$$

$$\pi_\theta = \tanh(\mathbf{W}_2 \lambda_\theta + \mathbf{b}_2) \quad (40)$$

$$\mu_\theta = \mathbf{W}_3 \pi_\theta + \mathbf{b}_3 \quad (41)$$

$$\log \sigma_\theta = \mathbf{W}_4 \pi_\theta + \mathbf{b}_4 \quad (42)$$

(1) Inference Network $q_\phi(\mathbf{h}|\mathbf{q}, \mathbf{a}, \mathbf{y})$:

$$s_q(|\mathbf{q}|) = f_q^{\text{LSTM}}(\mathbf{q}) \quad (30)$$

$$s_a(|\mathbf{a}|) = f_a^{\text{LSTM}}(\mathbf{a}) \quad (31)$$

$$s_y = \mathbf{W}_5 \mathbf{y} + \mathbf{b}_5 \quad (32)$$

$$\gamma = s_q(|\mathbf{q}|) || s_a(|\mathbf{a}|) || s_y \quad (33)$$

$$\lambda_\phi = \tanh(\mathbf{W}_6 \gamma + \mathbf{b}_6) \quad (34)$$

$$\pi_\phi = \tanh(\mathbf{W}_7 \lambda_\phi + \mathbf{b}_7) \quad (35)$$

$$\mu_\phi = \mathbf{W}_8 \pi_\phi + \mathbf{b}_8 \quad (36)$$

$$\log \sigma_\phi = \mathbf{W}_9 \pi_\phi + \mathbf{b}_9 \quad (37)$$

$$\mathbf{h} \sim \mathcal{N}(\mu_\phi(\mathbf{q}, \mathbf{a}, \mathbf{y}), \text{diag}(\sigma_\phi^2(\mathbf{q}, \mathbf{a}, \mathbf{y}))) \quad (38)$$

$$p_\theta(\mathbf{y}|\mathbf{q}, \mathbf{a}, \mathbf{h}):$$

$$\mathbf{e}(i) = \mathbf{W}_\alpha^T \tanh(\mathbf{W}_h \mathbf{h} + \mathbf{W}_s s_a(i)) \quad (43)$$

$$\alpha(i) = \frac{\mathbf{e}(i)}{\sum_j \mathbf{e}(j)} \quad (44)$$

$$c(\mathbf{a}, \mathbf{h}) = \sum_i s_a(i) \alpha(i) \quad (45)$$

$$z_a(\mathbf{a}, \mathbf{h}) = \tanh(\mathbf{W}_n c(\mathbf{a}, \mathbf{h}) + \mathbf{W}_m s_a(|\mathbf{a}|)) \quad (46)$$

$$z_q(\mathbf{q}) = s_q(|\mathbf{q}|) \quad (47)$$

$$p_\theta(\mathbf{y} = 1|\mathbf{q}, \mathbf{a}, \mathbf{h}) = \sigma(z_q^T \mathbf{M} z_a + b) \quad (48)$$

Neural Variational Inference for Text Processing

Question Answering: Experiments

Source	Set	Questions	QA Pairs	Judgement
QASent	Train	1,229	53,417	automatic
	Dev	82	1,148	manual
	Test	100	1,517	manual
WikiQA	Train	2,118	20,360	manual
	Dev	296	2,733	manual
	Test	633	6,165	manual

Model	QASent		WikiQA	
	MAP	MRR	MAP	MRR
Published Models				
PV	0.5213	0.6023	0.5110	0.5160
Bigram-CNN	0.5693	0.6613	0.6190	0.6281
Deep CNN	0.5719	0.6621	—	—
PV + Cnt	0.6762	0.7514	0.5976	0.6058
WA	0.7063	0.7740	—	—
LCLR	0.7092	0.7700	0.5993	0.6068
Bigram-CNN + Cnt	0.7113	0.7846	0.6520	0.6652
Deep CNN + Cnt	0.7186	0.7826	—	—
Our Models				
LSTM	0.6436	0.7235	0.6552	0.6747
LSTM + Att	0.6451	0.7316	0.6639	0.6828
NASM	0.6501	0.7324	0.6705	0.6914
LSTM + Cnt	0.7228	0.7986	0.6820	0.6988
LSTM + Att + Cnt	0.7289	0.8072	0.6855	0.7041
NASM + Cnt	0.7339	0.8117	0.6886	0.7069

Q1	how old was sue lyon when she made lolita
A _{NASM}	the actress who played lolita , sue lyon , was fourteen at the time of filming .
A _{LSTM}	the actress who played lolita , sue lyon , was fourteen at the time of filming .
Q2	how much is centavos in mexico
A _{NASM}	the peso is subdivided into 100 centavos , represented by " _UNK_ "
A _{LSTM}	the peso is subdivided into 100 centavos , represented by " _UNK_ "
Q3	what does a liquid oxygen plant look like
A _{NASM}	the blue color of liquid oxygen in a dewar flask
A _{LSTM}	the blue color of liquid oxygen in a dewar flask

Figure 4. A visualisation of attention scores on answer sentences.

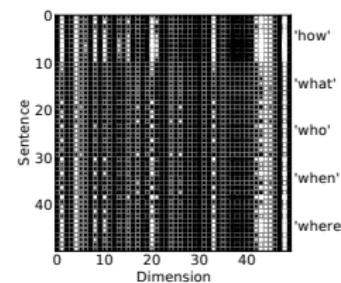
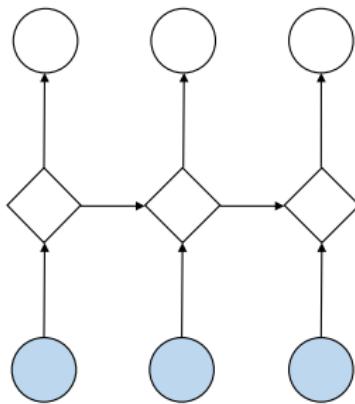


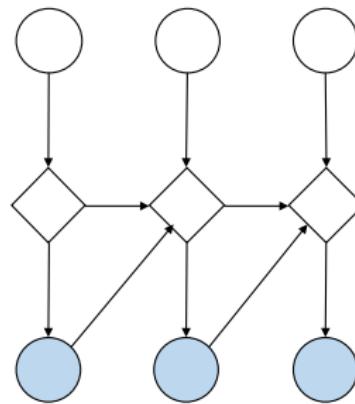
Figure 5. Hinton diagrams of the log standard deviations.

Learning Stochastic Recurrent Networks

ICLR 2015 Workshop



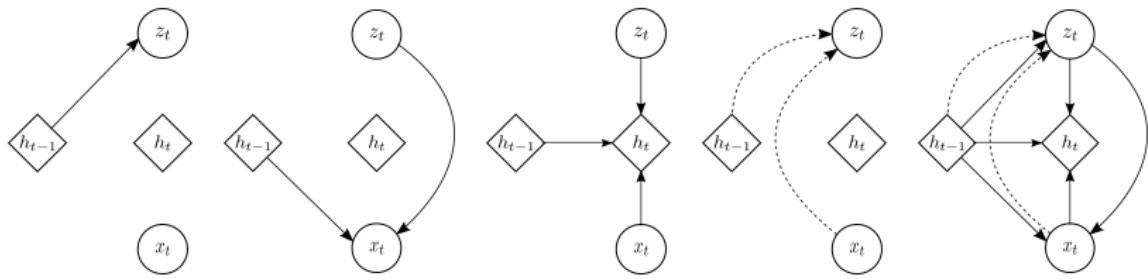
(a) Inference Network



(b) Generative Model

A Recurrent Latent Variable Model for Sequential Data

Model: Generation



(a) Prior

(b) Generation

(c) Recurrence

(d) Inference

(e) Overall

$$p(\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T}) = \prod_{t=1}^T p(x_t | \mathbf{z}_{\leq t}, \mathbf{x}_{<t}) p(z_t | \mathbf{x}_{<t}, \mathbf{z}_{<t})$$

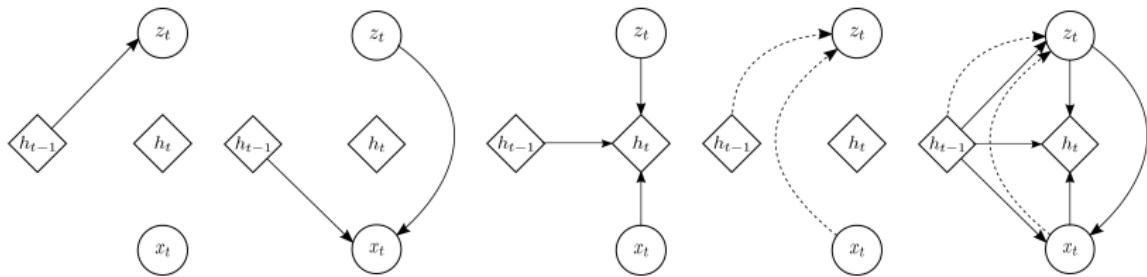
$\mathbf{z}_t \sim \mathcal{N}(\mu_{0,t}, \text{diag}(\sigma_{0,t}^2))$, where $[\mu_{0,t}, \sigma_{0,t}] = \phi_{\tau}^{\text{prior}}(\mathbf{h}_{t-1})$

$\mathbf{x}_t | \mathbf{z}_t \sim \mathcal{N}(\mu_{x,t}, \text{diag}(\sigma_{x,t}^2))$, where $[\mu_{x,t}, \sigma_{x,t}] = \phi_{\tau}^{\text{dec}}(\phi_{\tau}^{\mathbf{z}}(\mathbf{z}_t), \mathbf{h}_{t-1})$

$$\mathbf{h}_t = f_{\theta}(\phi_{\tau}^{\mathbf{x}}(\mathbf{x}_t), \phi_{\tau}^{\mathbf{z}}(\mathbf{z}_t), \mathbf{h}_{t-1})$$

A Recurrent Latent Variable Model for Sequential Data

Model: Inference



(a) Prior

(b) Generation

(c) Recurrence

(d) Inference

(e) Overall

$$q(\mathbf{z}_{\leq T} | \mathbf{x}_{\leq T}) = \prod_{t=1}^T q(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{<t})$$

$\mathbf{z}_t | \mathbf{x}_t \sim \mathcal{N}(\boldsymbol{\mu}_{z,t}, \text{diag}(\boldsymbol{\sigma}_{z,t}^2))$, where $[\boldsymbol{\mu}_{z,t}, \boldsymbol{\sigma}_{z,t}] = \phi_\tau^{enc}(\phi_\tau^x(\mathbf{x}_t), \mathbf{h}_{t-1})$

A Recurrent Latent Variable Model for Sequential Data Experiments

- **Speech Modeling:** train the model directly on raw audio signal, represented as a sequence of 200-d frames.
 - **Blizzard:** 300 hours of English spoken by a single female speaker.
 - **TIMIT:** 6300 English sentences, read by 630 speakers.
 - **Onomatopoeia:** 6738 non-linguistic human-made sounds, such as coughing, screaming, laughing and shouting, recorded from 51 voice actors.
 - **Accent:** English paragraphs read by 2046 different English speakers.
- **Handwriting Generation**
 - **IAM-OnDB:** 13040 handwritten lines written by 500 writers.

A Recurrent Latent Variable Model for Sequential Data Experiments

Table 1: Average log-likelihood on the test (or validation) set of each task.

Models	Speech modelling				Handwriting
	Blizzard	TIMIT	Onomatopoeia	Accent	IAM-OnDB
RNN-Gauss	3539	-1900	-984	-1293	1016
RNN-GMM	7413	26643	18865	3453	1358
VRNN-I-Gauss	≥ 8933 ≈ 9188	≥ 28340 ≈ 29639	≥ 19053 ≈ 19638	≥ 3843 ≈ 4180	≥ 1332 ≈ 1353
VRNN-Gauss	≥ 9223 $\approx \mathbf{9516}$	≥ 28805 $\approx \mathbf{30235}$	≥ 20721 $\approx \mathbf{21332}$	≥ 3952 ≈ 4223	≥ 1337 ≈ 1354
VRNN-GMM	≥ 9107 ≈ 9392	≥ 28982 ≈ 29604	≥ 20849 ≈ 21219	≥ 4140 $\approx \mathbf{4319}$	≥ 1384 $\approx \mathbf{1384}$

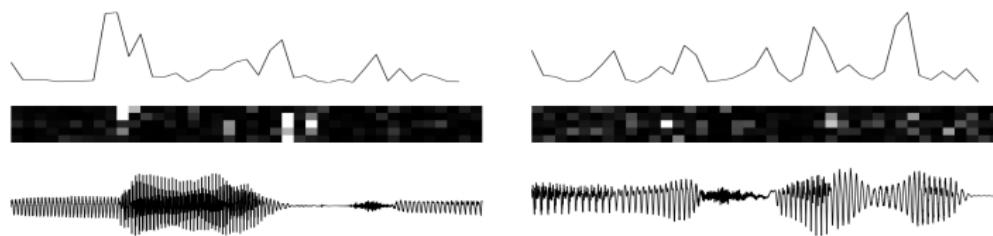


Figure 2: The top row represents the difference δ_t between $\mu_{z,t}$ and $\mu_{z,t-1}$. The middle row shows the dominant KL divergence values in temporal order. The bottom row shows the input waveforms.

A Recurrent Latent Variable Model for Sequential Data Experiments

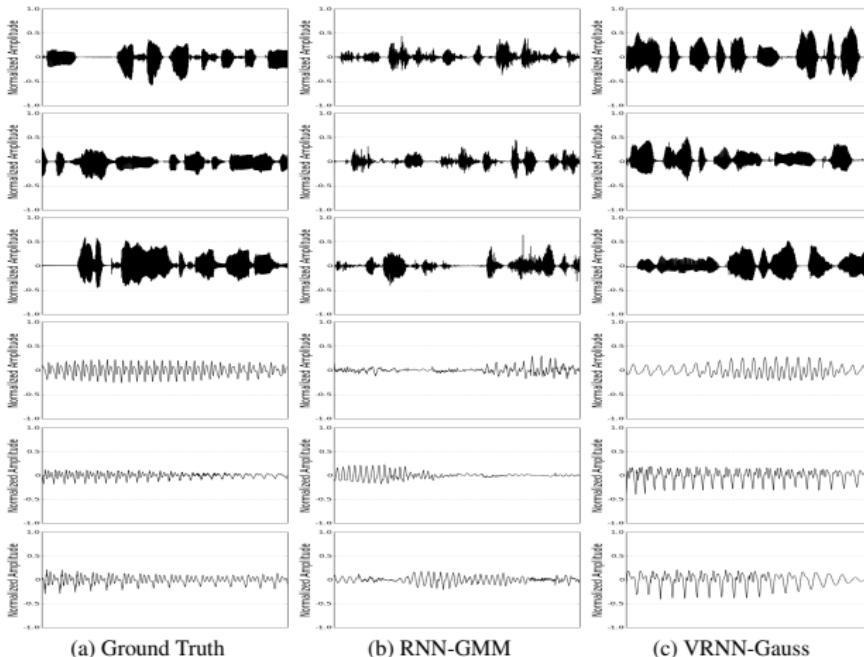


Figure 3: Examples from the training set and generated samples from RNN-GMM and VRNN-Gauss. Top three rows show the global waveforms while the bottom three rows show more zoomed-in waveforms. Samples from (b) RNN-GMM contain high-frequency noise, and samples from (c) VRNN-Gauss have less noise. We exclude RNN-Gauss, because the samples are almost close to pure noise.

A Recurrent Latent Variable Model for Sequential Data Experiments

one time defend Batman. The
the door opened an
-o the doctor. You c
was just thinking how he
or should one assume that
2.1. Methods of construction

(a) Ground Truth

(part of) Train 157 down after
the role as in it happens when
silvery freedom comic, i think
Gta, local television culture
follows in unselfishness in
Ryder, of other with children.

(b) RNN-Gauss

I feel like we have probably
as much as I can do.
Please to feel, as Daffy, our
was much V.T., my chest
red twice to feel so I am
a Mabat lab. hopefully

(c) RNN-GMM

freight container and denote
was not for me approach to me at
. And very significant in the history of
junkie. It clearly perceives when in an
e being hysteric - Wagner, RW. no
men of the first forage, could!!

(d) VRNN-GMM

Figure 4: Handwriting samples: (a) training examples and unconditionally generated handwriting from (b) RNN-Gauss, (c) RNN-GMM and (d) VRNN-GMM. The VRNN-GMM retains the writing style from beginning to end while RNN-Gauss and RNN-GMM tend to change the writing style during the generation process. This is possibly because the sequential latent random variables can guide the model to generate each sample with a consistent writing style.

Sequential Neural Models with Stochastic Layers

Model

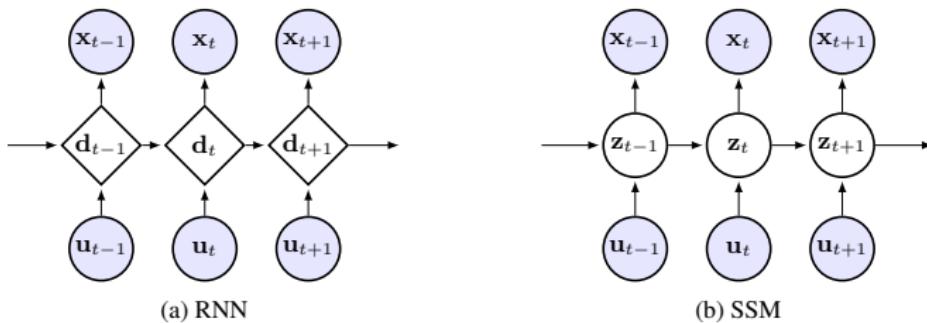


Figure 1: Graphical models to generate $\mathbf{x}_{1:T}$ with a recurrent neural network (RNN) and a state space model (SSM). Diamond-shaped units are used for deterministic states, while circles are used for stochastic ones. For sequence generation, like in a language model, one can use $\mathbf{u}_t = \mathbf{x}_{t-1}$.

Sequential Neural Models with Stochastic Layers Model

Like SSM, \mathbf{z}_t directly depends on \mathbf{z}_{t-1} . The true posterior of \mathbf{z}_t depends on $\mathbf{z}_{t-1}, \mathbf{d}_{t:T}, \mathbf{x}_{t:T}$.

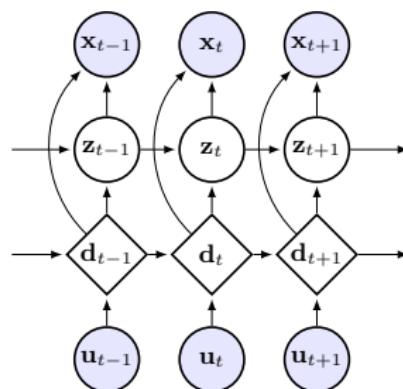
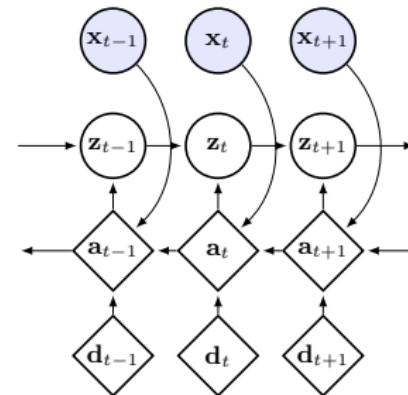
(a) Generative model p_θ (b) Inference network q_ϕ

Figure 2: A SRNN as a generative model p_θ for a sequence $\mathbf{x}_{1:T}$. Posterior inference of $\mathbf{z}_{1:T}$ and $\mathbf{d}_{1:T}$ is done through an inference network q_ϕ , which uses a backwards-recurrent state \mathbf{a}_t to approximate the nonlinear dependence of \mathbf{z}_t on future observations $\mathbf{x}_{t:T}$ and states $\mathbf{d}_{t:T}$; see Equation (7).

Sequential Neural Models with Stochastic Layers

Experiments

Models	Blizzard	TIMIT
SRNN (smooth+Res _q)	≥ 11991	≥ 60550
SRNN (smooth)	≥ 10991	≥ 59269
SRNN (filt+Res _q)	≥ 10572	≥ 52126
SRNN (filt)	≥ 10846	≥ 50524
VRNN-GMM	≥ 9107 ≈ 9392	≥ 28982 ≈ 29604
VRNN-Gauss	≥ 9223 ≈ 9516	≥ 28805 ≈ 30235
VRNN-I-Gauss	≥ 8933 ≈ 9188	≥ 28340 ≈ 29639
RNN-GMM	7413	26643
RNN-Gauss	3539	-1900

Table 1: Average log-likelihood per sequence on the test sets. For TIMIT the average test set length is 3.1s, while the Blizzard sequences are all 0.5s long. The non-SRNN results are reported as in [8]. Smooth: g_{ϕ_a} is a GRU running backwards; filt: g_{ϕ_a} is a feed-forward network; Res_q: parameterization with residual in (12).

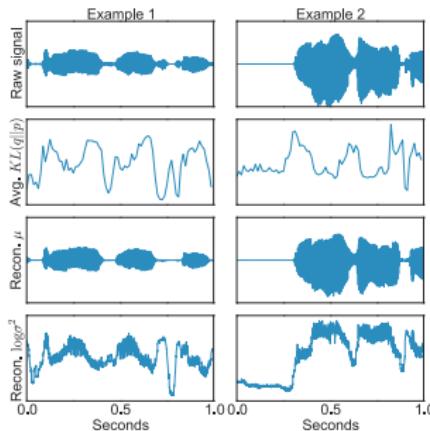


Figure 3: Visualization of the average KL term and reconstructions of the output mean and log-variance for two examples from the Blizzard test set.

Models	Nottingham	JSB chorales	MuseData	Piano-midi.de
SRNN (smooth+Res _q)	≥ -2.94	≥ -4.74	≥ -6.28	≥ -8.20
TSBN	≤ -3.67	≤ -7.48	≤ -6.81	≤ -7.98
NASMC	≈ -2.72	≈ -3.99	≈ -6.89	≈ -7.61
STORN	≈ -2.85	≈ -6.91	≈ -6.16	≈ -7.13
RNN-NADE	≈ -2.31	≈ -5.19	≈ -5.60	≈ -7.05
RNN	≈ -4.46	≈ -8.71	≈ -8.13	≈ -8.37

Table 2: Average log-likelihood on the test sets. The TSBN results are from [13], NASMC from [16], STORN from [3], RNN-NADE and RNN from [4].

Neural Adaptive Sequential Monte Carlo

Methods: Sequential Monte Carlo

Consider a probabilistic model

$$p(\mathbf{z}_{1:T}, \mathbf{x}_{1:T}) = p(\mathbf{z}_1)p(\mathbf{x}_1|\mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t|\mathbf{z}_{1:t-1})p(\mathbf{x}_t|\mathbf{z}_{1:t}, \mathbf{x}_{1:t-1})$$

The goal of SMC is to approximate the posterior as

$$p(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) \approx \sum_{n=1}^N \tilde{w}_t^{(n)} \delta(\mathbf{z}_{1:T} - \mathbf{z}_{1:T}^{(n)}) \quad (4)$$

through a weighted set of N sampled trajectories drawn from a proposal distribution $\{\mathbf{z}_{1:T}^{(n)}\}_{n=1:N} \sim q(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$,

$$q(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) = q(\mathbf{z}_1|\mathbf{x}_1) \prod_{t=2}^T q(\mathbf{z}_t|\mathbf{z}_{1:t-1}, \mathbf{x}_{1:t}) \quad (5)$$

Neural Adaptive Sequential Monte Carlo Methods

How to get normalized importance weights:

$$w(\mathbf{z}_{1:T}^{(n)}) = \frac{p(\mathbf{z}_{1:T}^{(n)}, \mathbf{x}_{1:T})}{q(\mathbf{z}_{1:T}^{(n)} | \mathbf{x}_{1:T})}, \tilde{w}(\mathbf{z}_{1:T}^{(n)}) = \frac{w(\mathbf{z}_{1:T}^{(n)})}{\sum_n w(\mathbf{z}_{1:T}^{(n)})} \propto \tilde{w}(\mathbf{z}_{1:T-1}^{(n)}) \frac{p(\mathbf{z}_T^{(n)} | \mathbf{z}_{1:T-1}^{(n)}) p(\mathbf{x}_T | \mathbf{z}_{1:T}^{(n)}, \mathbf{x}_{1:T-1})}{q(\mathbf{z}_T^{(n)} | \mathbf{z}_{1:T-1}^{(n)}, \mathbf{x}_{1:T})}$$

using Sequential Importance Sampling (**SIS**), or Sequential Importance Resampling (**SIR**).

Adapting proposals by descending the inclusive KL divergence

$$-\frac{\partial}{\partial \phi} \text{KL}[p_\theta(\mathbf{z}_{1:T} | \mathbf{x}_{1:T}) || q_\phi(\mathbf{z}_{1:T} | \mathbf{x}_{1:T})] \approx \sum_t \sum_n \tilde{w}_t^{(n)} \frac{\partial}{\partial \phi} \log q_\phi(\mathbf{z}_t^{(n)} | \mathbf{x}_{1:t}, \mathbf{z}_{1:t-1}^{A_{t-1}^{(n)}}).$$

Model parameter learning by maximum likelihood

$$\frac{\partial}{\partial \theta} \log[p_\theta(\mathbf{x}_{1:T})] \approx \sum_t \sum_n \tilde{w}_t^{(n)} \frac{\partial}{\partial \theta} \log p_\theta(\mathbf{x}_t, \mathbf{z}_t^{(n)} | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1}^{A_{t-1}^{(n)}}).$$

Neural Adaptive Sequential Monte Carlo Methods

Algorithm 1 Stochastic Gradient Adaptive SMC (batch inference and learning variants)

Require: proposal: q_ϕ , model: p_θ , observations: $X = \{\mathbf{x}_{1:T_j}\}_{j=1:M}$, number of particles: N
repeat

$\{\mathbf{x}_{1:T_j}^{(j)}\}_{j=1:m} \leftarrow \text{NextMiniBatch}(X)$

$\{\mathbf{z}_{1:t}^{(i,j)}, \tilde{w}_t^{(i,j)}\}_{i=1:N, j=1:m, t=1:T_j} \leftarrow \text{SMC}(\theta, \phi, N, \{\mathbf{x}_{1:T_j}^{(j)}\}_{j=1:m})$

$\Delta\phi = \sum_j \sum_{t=1}^{T_j} \sum_i \tilde{w}_t^{(i,j)} \frac{\partial}{\partial \phi} \log q_\phi(\mathbf{z}_t^{(i,j)} | \mathbf{x}_{1:t}^{(j)}, \mathbf{z}_{1:t-1}^{A_{t-1}^{(i,j)}})$

$\Delta\theta = \sum_j \sum_{t=1}^{T_j} \sum_i \tilde{w}_t^{(i,j)} \frac{\partial}{\partial \theta} \log p_\theta(\mathbf{x}_t^{(j)}, \mathbf{z}_t^{(i,j)} | \mathbf{x}_{1:t-1}^{(j)}, \mathbf{z}_{1:t-1}^{A_{t-1}^{(i,j)}}) \quad (\text{optional})$

$\phi \leftarrow \text{Optimize}(\phi, \Delta\phi)$

$\theta \leftarrow \text{Optimize}(\theta, \Delta\theta) \quad (\text{optional})$

until convergence

Dataset	LV-RNN (NASMC)	LV-RNN (Bootstrap)	STORN (SGVB)	FD-RNN	sRNN	RNN-NADE
Piano-midi-de	7.61	7.50	7.13	7.39	7.58	7.03
Nottingham	2.72	3.33	2.85	3.09	3.43	2.31
MuseData	6.89	7.21	6.16	6.75	6.99	5.60
JSBChorales	3.99	4.26	6.91	8.01	8.58	5.19

Takeaways

- **Modeling:**
 - ① RNN + stochastic latent variable
 - ② applications in mocap, music, text, speech, handwriting, etc.
 - ③ mixture RNN model?
- **Training:**
 - ① SGVB
 - ② NASMC
 - ③ combine SMC + VAE?

References I

-  H. Wallach A. Schein and M. Zhou.
Poisson gamma dynamical systems.
In *NIPS*, 2016.
-  Ayan Acharya, Joydeep Ghosh, and Mingyuan Zhou.
Nonparametric bayesian factor analysis for dynamic count matrices.
In *AISTATS*, 2015.
-  Evan Archer, Il Memming Park, Lars Buesing, John Cunningham, and Liam Paninski.
Black box variational inference for state space models.
arXiv preprint arXiv:1511.07367, 2015.
-  Justin Bayer and Christian Osendorfer.
Learning stochastic recurrent networks.
In *ICLR workshop*, 2015.
-  Nicolas Boulanger-Lewandowski, Yoshua Bengio, and Pascal Vincent.
Modeling temporal dependencies in high-dimensional sequences: Application to polyphonic music generation and transcription.
In *ICML*, 2012.
-  Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio.
Generating sentences from a continuous space.
In *CoNLL*, 2016.
-  Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron C Courville, and Yoshua Bengio.
A recurrent latent variable model for sequential data.
In *NIPS*, 2015.

References II

-  Otto Fabius and Joost R van Amersfoort.
Variational recurrent auto-encoders.
In *ICLR workshop*, 2015.
-  Marco Fraccaro, Søren Kaae Sønderby, Ulrich Paquet, and Ole Winther.
Sequential neural models with stochastic layers.
In *NIPS*, 2016.
-  Zhe Gan, Chunyuan Li, Ricardo Henao, David E Carlson, and Lawrence Carin.
Deep temporal sigmoid belief networks for sequence modeling.
In *NIPS*, 2015.
-  Shixiang Gu, Zoubin Ghahramani, and Richard E Turner.
Neural adaptive sequential monte carlo.
In *NIPS*, 2015.
-  Matthew J Johnson, David Duvenaud, Alexander B Wiltschko, Sandeep R Datta, and Ryan P Adams.
Structured vaes: Composing probabilistic graphical models and variational autoencoders.
In *NIPS*, 2016.
-  Maximilian Karl, Maximilian Soelch, Justin Bayer, and Patrick van der Smagt.
Deep variational bayes filters: Unsupervised learning of state space models from raw data.
arXiv preprint arXiv:1605.06432, 2016.
-  Rahul G Krishnan, Uri Shalit, and David Sontag.
Deep kalman filters.
arXiv preprint arXiv:1511.05121, 2015.

References III

-  **Rahul G Krishnan, Uri Shalit, and David Sontag.**
Structured inference networks for nonlinear state space models.
arXiv preprint arXiv:1609.09869, 2016.
-  **Yishu Miao, Lei Yu, and Phil Blunsom.**
Neural variational inference for text processing.
In *ICML*, 2016.
-  **Roni Mittelman, Benjamin Kuipers, Silvio Savarese, and Honglak Lee.**
Structured recurrent temporal restricted boltzmann machines.
In *ICML*, 2014.
-  **Jiaming Song, Zhe Gan, and Lawrence Carin.**
Factored temporal sigmoid belief networks for sequence learning.
In *ICML*, 2016.
-  **Ilya Sutskever and Geoffrey E Hinton.**
Learning multilevel distributed representations for high-dimensional sequences.
In *AISTATS*, 2007.
-  **Ilya Sutskever, Geoffrey E Hinton, and Graham W Taylor.**
The recurrent temporal restricted boltzmann machine.
In *NIPS*, 2009.
-  **Graham W Taylor and Geoffrey E Hinton.**
Factored conditional restricted boltzmann machines for modeling motion style.
In *ICML*, 2009.

References IV

-  Graham W Taylor, Geoffrey E Hinton, and Sam T Roweis.
Modeling human motion using binary latent variables.
In *NIPS*, 2006.
-  Graham W Taylor, Leonid Sigal, David J Fleet, and Geoffrey E Hinton.
Dynamical binary latent variable models for 3d human pose tracking.
In *CVPR*, 2010.
-  Matthew D Zeiler, Graham W Taylor, Leonid Sigal, Iain Matthews, and Rob Fergus.
Facial expression transfer with input-output temporal restricted boltzmann machines.
In *NIPS*, 2011.
-  Biao Zhang, Deyi Xiong, and Jinsong Su.
Variational neural machine translation.
In *EMNLP*, 2016.
-  Yizhe Zhang, Yue Zhao, Lawrence David, Ricardo Henao, and Lawrence Carin.
Dynamic poisson factor analysis.
In *ICDM*, 2016.