



ANALYSIS OF QUANTIZED MODELS

 1 {LHOUAB, JAMESK}@CSE.UST.HK, 2 TESLAZHU@MAIL.USTC.EDU.CN, 2 {FEIGA, TAOQIN, TYLIU}@MICROSOFT.COM

Lu Hou¹, Jinhua Zhu², James T. Kwok¹, Fei Gao³, Tao Qin³, Tie-yan Liu³

¹Hong Kong University of Science and Technology, ²University of Science and Technology of China, ³Microsoft Research, Beijing, China

Research

微软亚洲研究院

Microsoft®

BACKGROUD AND MOTIVATION

Quantized LSTM

- BinaryConnect fails (Hou et al., 2017) on LSTM
- BWN, TWN, LAB, LAT perform much better, but sometimes inferior to the Spectral norm of quantized matrix full-precision network on some tasks (Ardakani et al. 2019)
- SOTA performance when training binarized/ternarized LSTMs with batch normalization (Ardakani et al. 2019)
- Why does batch normalization work for quantized LSTM?
- Does weight/layer normalization also help?

Recurrence of a LSTM

$$\begin{bmatrix} \mathbf{i}_t \\ \mathbf{f}_t \\ \mathbf{a}_t \\ \mathbf{o}_t \end{bmatrix} = \begin{bmatrix} \mathbf{W}_{xi} \mathbf{x}_t + \mathbf{W}_{hi} \mathbf{h}_{t-1} \\ \mathbf{W}_{xf} \mathbf{x}_t + \mathbf{W}_{hf} \mathbf{h}_{t-1} \\ \mathbf{W}_{xa} \mathbf{x}_t + \mathbf{W}_{ha} \mathbf{h}_{t-1} \\ \mathbf{W}_{xo} \mathbf{x}_t + \mathbf{W}_{ho} \mathbf{h}_{t-1} \end{bmatrix} + \begin{bmatrix} \mathbf{b}_i \\ \mathbf{b}_f \\ \mathbf{b}_a \\ \mathbf{b}_o \end{bmatrix},$$

$$\mathbf{c}_t = \sigma(\mathbf{i}_t) \odot \tanh(\mathbf{a}_t) + \sigma(\mathbf{f}_t) \odot \mathbf{c}_{t-1},$$

$$\mathbf{h}_t = \sigma(\mathbf{o}_t) \odot \tanh(\mathbf{c}_t).$$

- storage dominated by $\mathbf{W}_{x*}, \mathbf{W}_{h*}$
- computation dominated by $\mathbf{W}_{x*}\mathbf{x}_t + \mathbf{W}_{h*}\mathbf{h}_{t-1}$
- ullet ightarrow quantize $\mathbf{W}_{x*}, \mathbf{W}_{h*}$

Our Findings

- Quantized LSTM is hard to train due to the exploding gradient problem
- popularly used weight/layer/batch normalization schemes can help stabilize the gradient magnitude in training quantized LSTMs

EXPLODING GRADIENT IN LSTM

Proposition 1
$$\left\| \frac{\partial \xi_m}{\partial \mathbf{h}_{t-1}} \right\| \leq \lambda_1 \left\| \frac{\partial \xi_m}{\partial \mathbf{h}_t} \right\| + \lambda_2 \left\| \frac{\partial \xi_m}{\partial \mathbf{c}_{t+1}} \right\|.$$

- $\lambda_1 = \frac{1}{4} \|\mathbf{W}_{hi}\|_2 + \frac{\gamma_1}{4} \|\mathbf{W}_{hf}\|_2 + \|\mathbf{W}_{ha}\|_2 + \frac{1}{4} \|\mathbf{W}_{ho}\|_2$ $\lambda_2 = \frac{1}{4} \|\mathbf{W}_{hi}\|_2 + \frac{\gamma_1}{4} \|\mathbf{W}_{hf}\|_2 + \|\mathbf{W}_{ha}\|_2$

Corollary 1 When $\lambda_2 = 0$, a necessary condition for exploding gradients in the LSTM is

- empirically, λ_2 is rarely zero
- upper bound of $\|\frac{\partial \xi_m}{\partial \mathbf{h}_{\perp}}\|$ even larger, and gradient explode even more easily

QUANTIZATION -> GRADIENT EASIER TO EXPLODE

 λ_1,λ_2 in the upper bound related to the spectral norm of weight matrix

• larger spectral norm \rightarrow more easily to explode

- For any $\mathbf{W} \in \{-1, +1\}^{d \times d}$, $\|\mathbf{W}\|_2 \ge \sqrt{d}$.
- For any $\mathbf{W} \in \{-B_k, \dots, -B_1, B_0, B_1, \dots, B_k\}^{d \times d}$ where $0 = B_0 < B_1 < \dots < B_k$ B_k , $\|\mathbf{W}\|_2 \ge (1-s)B_1\sqrt{d}$, where s is the sparsity of W.

the gradient is more easily to explode

- \bullet when d is large
- for binarization than ternarization

NORMALIZED LSTM

- Apply normalization as $\mathcal{N}(\mathbf{W}_{x*}\mathbf{x}_t)$ and $\mathcal{N}(\mathbf{W}_{h*}\mathbf{h}_{t-1})$
- \bullet \mathcal{N} can be weight, layer or batch normalization

$$egin{bmatrix} \widetilde{\mathbf{i}}_t \ \widetilde{\mathbf{f}}_t \ \widetilde{\mathbf{a}}_t \end{bmatrix} = egin{bmatrix} \mathcal{N}(\mathbf{W}_{xi}\mathbf{x}_t) + \mathcal{N}(\mathbf{W}_{hi}\mathbf{h}_{t-1}) \ \mathcal{N}(\mathbf{W}_{xa}\mathbf{x}_t) + \mathcal{N}(\mathbf{W}_{ha}\mathbf{h}_{t-1}) \ \mathcal{N}(\mathbf{W}_{xa}\mathbf{x}_t) + \mathcal{N}(\mathbf{W}_{ha}\mathbf{h}_{t-1}) \ \mathcal{N}(\mathbf{W}_{xo}\mathbf{x}_t) + \mathcal{N}(\mathbf{W}_{ho}\mathbf{h}_{t-1}) \end{bmatrix} + egin{bmatrix} \mathbf{b}_i \ \mathbf{b}_a \ \mathbf{b}_o \end{bmatrix},$$

 $\mathbf{c}_t = \sigma(\tilde{\mathbf{i}}_t) \odot \tanh(\tilde{\mathbf{a}}_t) + \sigma(\tilde{\mathbf{f}}_t) \odot \mathbf{c}_{t-1},$

 $\mathbf{h}_t = \sigma(\tilde{\mathbf{o}}_t) \odot \tanh(\mathbf{c}_t).$

WEIGHT NORMALIZATION:

decouples length and direction of the weight vector

$$\mathcal{WN}(\mathbf{W}_{j,:}\mathbf{x}) = g_j rac{\mathbf{W}_{j,:}}{\|\mathbf{W}_{j,:}\|}\mathbf{x}$$

• each row $W_{i,:}$ of W (W_{h*} or W_{x*}) is separately normalized as

• $g_* = \max_{1 \le j \le d} g_j$; $\mathbf{D}_* = \operatorname{diag}([\|(\mathbf{W}_{h*})_{1,:}\|, \|(\mathbf{W}_{h*})_{2,:}\|, \dots, \|(\mathbf{W}_{h*})_{d,:}\|]^\top)$

Proposition 2 With weight normalization,

$$\left\| \frac{\partial \xi_{m}}{\partial \mathbf{h}_{t-1}} \right\| \leq \left(\frac{g_{i}}{4} \left\| \mathbf{D}_{i}^{-1} \mathbf{W}_{hi} \right\|_{2} + \frac{\gamma_{1} g_{f}}{4} \left\| \mathbf{D}_{f}^{-1} \mathbf{W}_{hf} \right\|_{2} + g_{a} \left\| \mathbf{D}_{a}^{-1} \mathbf{W}_{ha} \right\|_{2} + \frac{g_{o}}{4} \left\| \mathbf{D}_{o}^{-1} \mathbf{W}_{ho} \right\|_{2} \right) \left\| \frac{\partial \xi_{m}}{\partial \mathbf{h}_{t}} \right\| + \left(\frac{g_{i}}{4} \left\| \mathbf{D}_{i}^{-1} \mathbf{W}_{hi} \right\|_{2} + \frac{\gamma_{1} g_{f}}{4} \left\| \mathbf{D}_{f}^{-1} \mathbf{W}_{hf} \right\|_{2} + g_{a} \left\| \mathbf{D}_{a}^{-1} \mathbf{W}_{ha} \right\|_{2} \right) \left\| \frac{\partial \xi_{m}}{\partial \mathbf{c}_{t+1}} \right\|.$$

• if \mathbf{W}_{h*} is scaled by a factor α , \mathbf{D}_{*} will also be scaled by $\alpha \to \mathbf{D}_{*}^{-1}\mathbf{W}_{h*}$ not affected.

LAYER NORMALIZATION

normalizes activities in each layer

• input $\mathbf{x} \in \mathbb{R}^d$ ($\mathbf{W}_{x*}\mathbf{x}_t$ or $\mathbf{W}_{h*}\mathbf{h}_{t-1}$) with mean μ and standard deviation

$$\mathbf{y} = \mathcal{L}\mathcal{N}(\mathbf{x}) = \mathbf{g} \odot \mathbf{z} + \mathbf{b}, \text{where } \mathbf{z} = (\mathbf{x} - \mu \mathbf{1})/\sigma$$

• for $\mathcal{LN}(\mathbf{W}_{h*}\mathbf{h}_{t-1})$: $g_* = g_k, \sigma_* = \sigma_k$, where $k = \arg\max_{1 \leq j \leq d} g_j$

Proposition 3 With layer normalization,

$$\left\| \frac{\partial \xi_{m}}{\partial \mathbf{h}_{t-1}} \right\| \leq \left(\frac{1}{4} \frac{g_{i}}{\sigma_{i}} \|\mathbf{W}_{hi}\|_{2} + \frac{\gamma_{1}}{4} \frac{g_{f}}{\sigma_{f}} \|\mathbf{W}_{hf}\|_{2} + \frac{g_{a}}{\sigma_{a}} \|\mathbf{W}_{ha}\|_{2} + \frac{1}{4} \frac{g_{o}}{\sigma_{o}} \|\mathbf{W}_{ho}\|_{2} \right) \left\| \frac{\partial \xi_{m}}{\partial \mathbf{h}_{t}} \right\| + \left(\frac{1}{4} \frac{g_{i}}{\sigma_{i}} \|\mathbf{W}_{hi}\|_{2} + \frac{\gamma_{1}}{4} \frac{g_{f}}{\sigma_{f}} \|\mathbf{W}_{hf}\|_{2} + \frac{g_{a}}{\sigma_{a}} \|\mathbf{W}_{ha}\|_{2} \right) \left\| \frac{\partial \xi_{m}}{\partial \mathbf{c}_{t+1}} \right\|.$$

• if the elements of \mathbf{W}_{h*} grow twice as large, the corresponding σ_* will be twice as large

BATCH NORMALIZATION

operates on a minibatch (N samples in a batch)

- ullet $\mathbf{H}_t = [\mathbf{h}_t^1, \dots, \mathbf{h}_t^N]^ op \in \mathbb{R}^{N imes d}; \mathbf{X}_t = [\mathbf{x}_t^1, \dots, \mathbf{x}_t^N]^ op \in \mathbb{R}^{N imes r}$
- input $\mathbf{X} \in \mathbb{R}^{N \times d}$ ($\mathbf{X}_t \mathbf{W}_{x*}^{\top}$ or $\mathbf{H}_{t-1} \mathbf{W}_{h*}^{\top}$), with mean μ_j and std σ_j for the jth column

$$\mathbf{y}_{:,j} = \mathcal{BN}(\mathbf{X}_{:,j}) = g_j \frac{\mathbf{X}_{:,j} - \mu_j \mathbf{1}}{\sigma_j} + b_j \mathbf{1}$$

• for $\mathcal{BN}(\mathbf{H}_{t-1}\mathbf{W}_{h*}^{\top})$: $(\sigma_*, g_*) = \arg\max_{1 \leq j \leq d} \frac{g_j}{\sigma_j}$

Proposition 4 With batch normalization,

$$\sum_{k=1}^{N} \left\| \frac{\partial \xi_{m}}{\partial \mathbf{h}_{t-1}^{k}} \right\|^{2} \leq \left(\frac{1}{2} \frac{g_{i}^{2}}{\sigma_{i}^{2}} \|\mathbf{W}_{hi}\|_{2}^{2} + \frac{\gamma_{2}^{2}}{2} \frac{g_{f}^{2}}{\sigma_{f}^{2}} \|\mathbf{W}_{hf}\|_{2}^{2} + 8 \frac{g_{a}^{2}}{\sigma_{a}^{2}} \|\mathbf{W}_{ha}\|_{2}^{2} + \frac{1}{4} \frac{g_{o}^{2}}{\sigma_{o}^{2}} \|\mathbf{W}_{ho}\|_{2}^{2} \right) \sum_{k=1}^{N} \left\| \frac{\partial \xi_{m}}{\partial \mathbf{h}_{t}^{k}} \right\|^{2} + \left(\frac{1}{2} \frac{g_{i}^{2}}{\sigma_{i}^{2}} \|\mathbf{W}_{hi}\|_{2}^{2} + \frac{\gamma_{2}^{2}}{2} \frac{g_{f}^{2}}{\sigma_{f}^{2}} \|\mathbf{W}_{hf}\|_{2}^{2} + 8 \frac{g_{a}^{2}}{\sigma_{a}^{2}} \|\mathbf{W}_{ha}\|_{2}^{2} \right) \sum_{k=1}^{N} \left\| \frac{\partial \xi_{m}}{\partial \mathbf{c}_{t+1}^{k}} \right\|^{2}.$$

• if the elements of W_{h*} grow twice as large, the corresponding σ_* will be twice as large

EXPERIMENTS

Character-level Language Modeling

• Rite Por Character (RPC) and size (in KR) of 1-lawer I STM

	guanti	normali	Ize (in KB) of 1-				ToxtQ	
preci-	quanti-	normali-	War and Peace		Penn Treebank		Text8	
sion	zation	zation	BPC	size	BPC	size	BPC	size
full		_	1.72	4800	1.45	4504	1.46	63375
		weight	1.73	4816	1.45	4520	1.48	63438
	_	layer	1.69	4832	1.43	4536	1.45	63500
		batch (shared)	1.72	4864	1.45	4568	1.46	63625
		batch (separate)	1.72	8032	1.45	7736	1.46	86000
1-bit	SBN	batch (separate)	1.78	3794	1.60	3785	1.54	27464
		_	4.24	158	2.51	149	N/A	2011
	Binary-	weight	1.74	174	1.50	165	1.50	2073
	Connect	layer	1.69	190	1.49	181	1.47	2136
		batch (shared)	1.72	222	1.51	213	1.47	2261
		batch (separate)	1.72	3390	1.50	3381	1.48	24636
2-bit	STN	batch (separate)	1.72	3944	1.60	3521	1.51	15303
		_	6.35	308	5.84	289	N/A	3990
	Ter-	weight	1.72	324	1.42	305	1.42	4053
	Connect	layer	1.67	340	1.43	321	1.44	4115
		batch (shared)	1.70	372	1.44	353	1.44	4240
		batch (separate)	1.71	3540	1.45	3521	1.44	26615

Word-level Language Modeling

• Test Perplexity and size (in KB) of 1-layer LSTM with d hidden units

			a = 500		a = 000	
precision	quantlization	normalization	PPL	size	PPL	size
full		-	91.5	2817	87.6	13213
	SBN	batch (separate)	92.2	852	87.2	2068
	_	-	8247.4	93	1244.2	423
		weight	87.6	102	84.8	443
	BinaryConnect	layer	89.4	111	82.3	463
		batch (shared)	92.4	130	84.8	504
1-bit		batch(separate)	91.9	767	85.6	1885
	alternating LSTM	-	103.1	180	_	-
	STN	batch (separate)	90.7	940	86.1	2481
		-	113.8	180	113.8	835
		weight	86.5	190	84.9	856
	TerConnect	layer	88.2	199	82.5	876
2-bit		batch (shared)	90.6	218	85.8	917
2-DI		batch (separate)	91.6	855	86.5	2298

Observations

- Vanilla BinaryConnect and TerConnect fail, but normalized versions work
- Normalized quantized LSTM is comparable to the full-precision baseline
- Applying weight/layer/batch (shared) normalization perform similarly or better than SBN and STN (Ardakani et al., 2019), while being much smaller

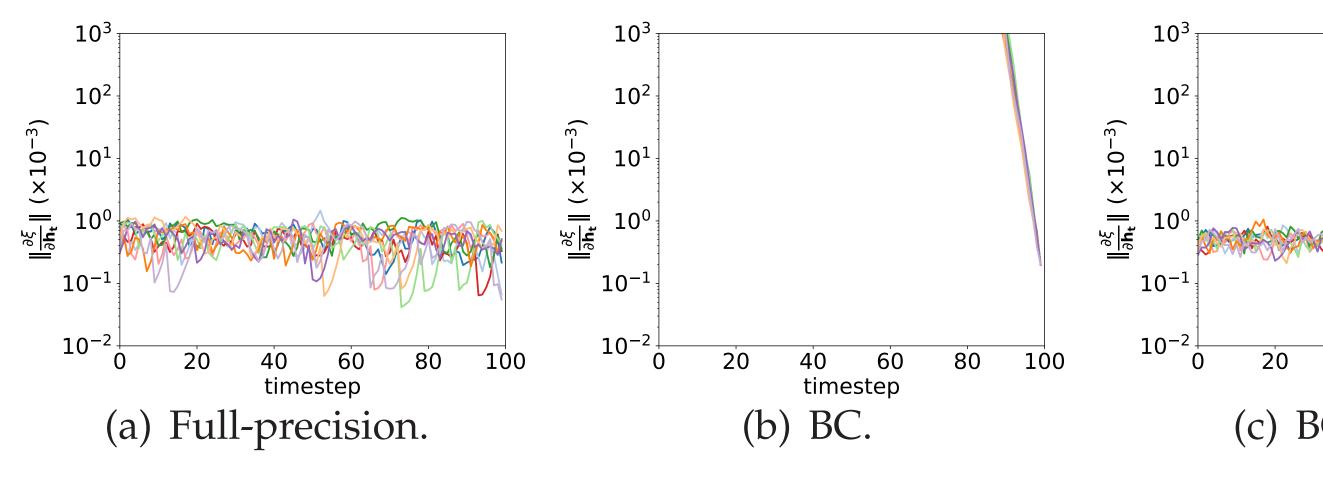
More experiments in the paper!

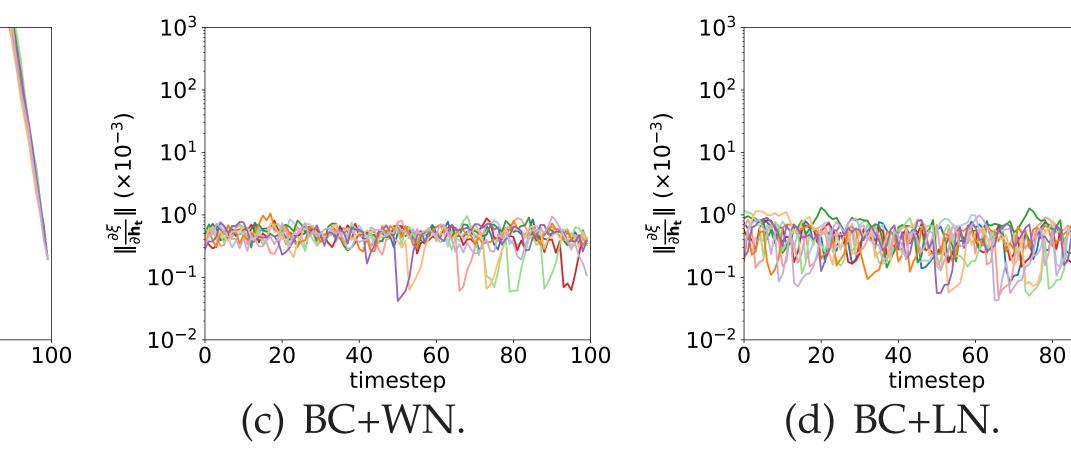
CONCLUSION

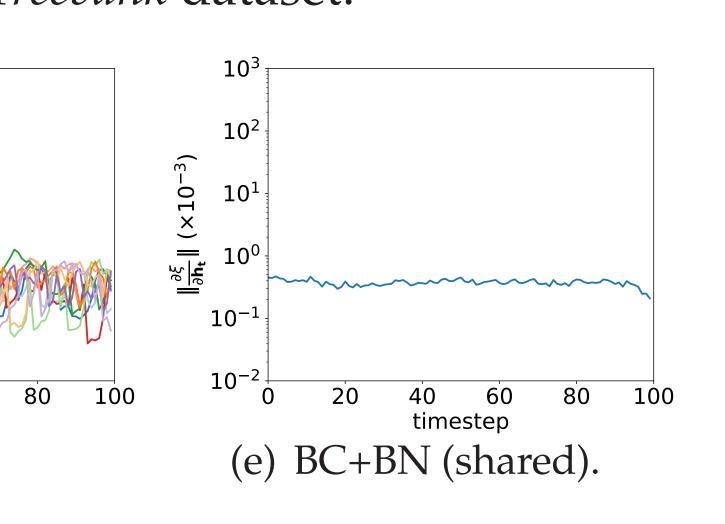
- WHY: Quantization tends to increase spectral norm of weights in LSTM, making the exploding gradient problem much more severe than its full-precision counterpart.
- HOW: By using normalization, backpropagation of $\left\| \frac{\partial \xi_m}{\partial \mathbf{h}_t} \right\|$ in the quantized LSTM is not affected by the possibly large scaling of the weight matrix caused by quantization, and the exploding gradient problem can be alleviated.
- CODE: https://github.com/houlu369/Normalized-Quantized-LSTN

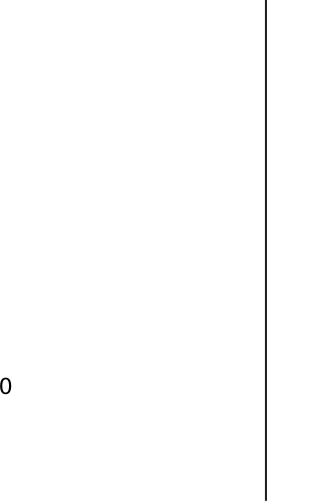
OBSERVATIONS

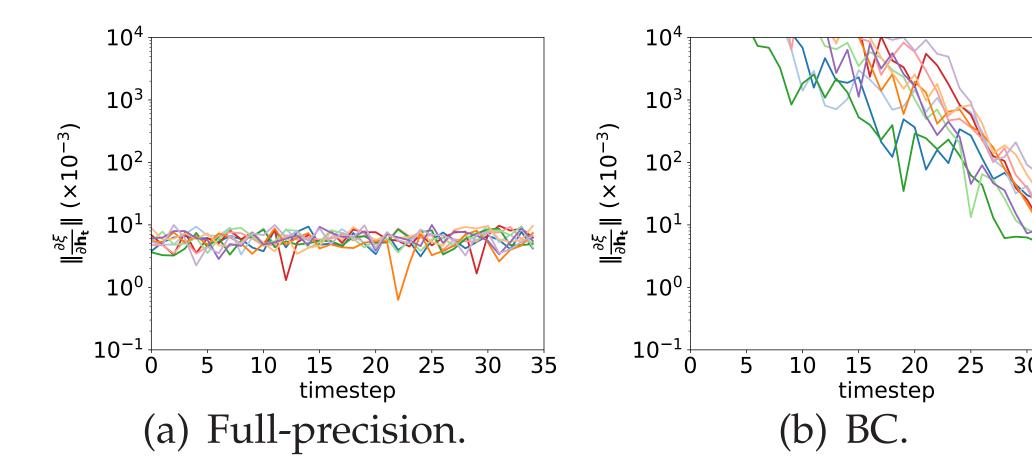
Figure 1: Gradient norms of character-level language modeling on Penn Treebank dataset.

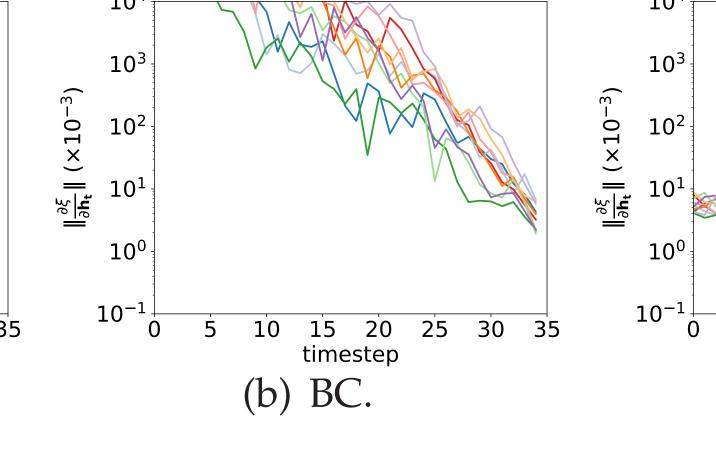












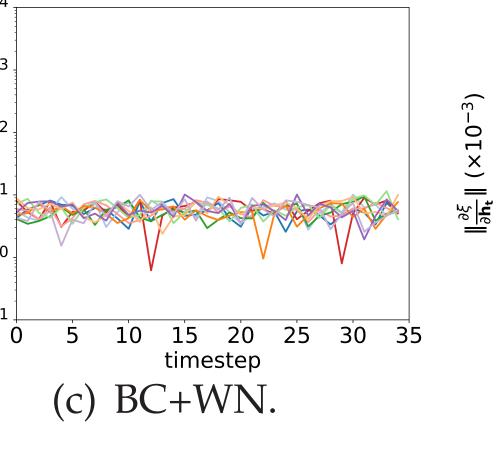


Figure 2: Gradient norms of word-level language modeling on Penn Treebank dataset.

