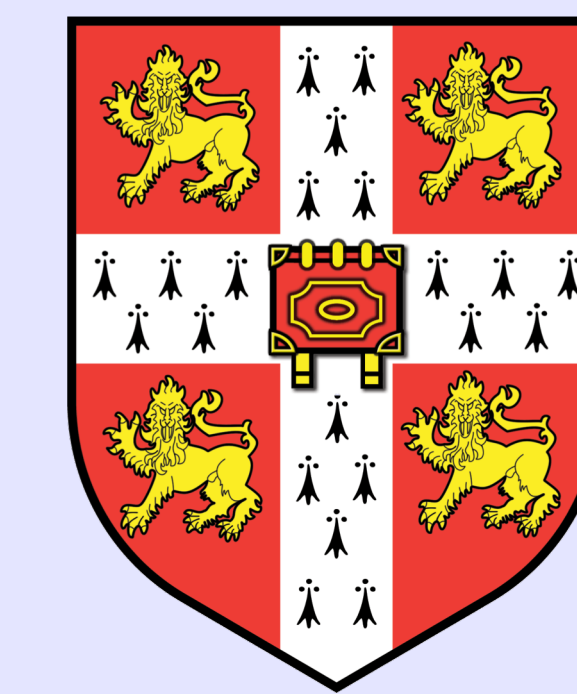


LIMITED-MEMORY BFGS OPTIMIZATION OF RECURRENT NEURAL NETWORK LANGUAGE MODELS FOR SPEECH RECOGNITION



Xunying Liu¹, Shansong Liu¹, Jinze Sha², Jianwei Yu¹, Zhiyuan Xu², Xie Chen² & Helen Meng¹

{xyliu, sslu, jwyu, hmmeng}@se.cuhk.edu.hk, {js2294, zyx22, xc257}@cam.ac.uk

¹The Chinese University of Hong Kong, ²Cambridge University Engineering Department

Introduction

Problem statement and objectives

- Faster and more stable training for deep neural networks (DNNs)
- Investigating 2nd order optimization techniques
- Applied to recurrent neural network language model (RNNLM)

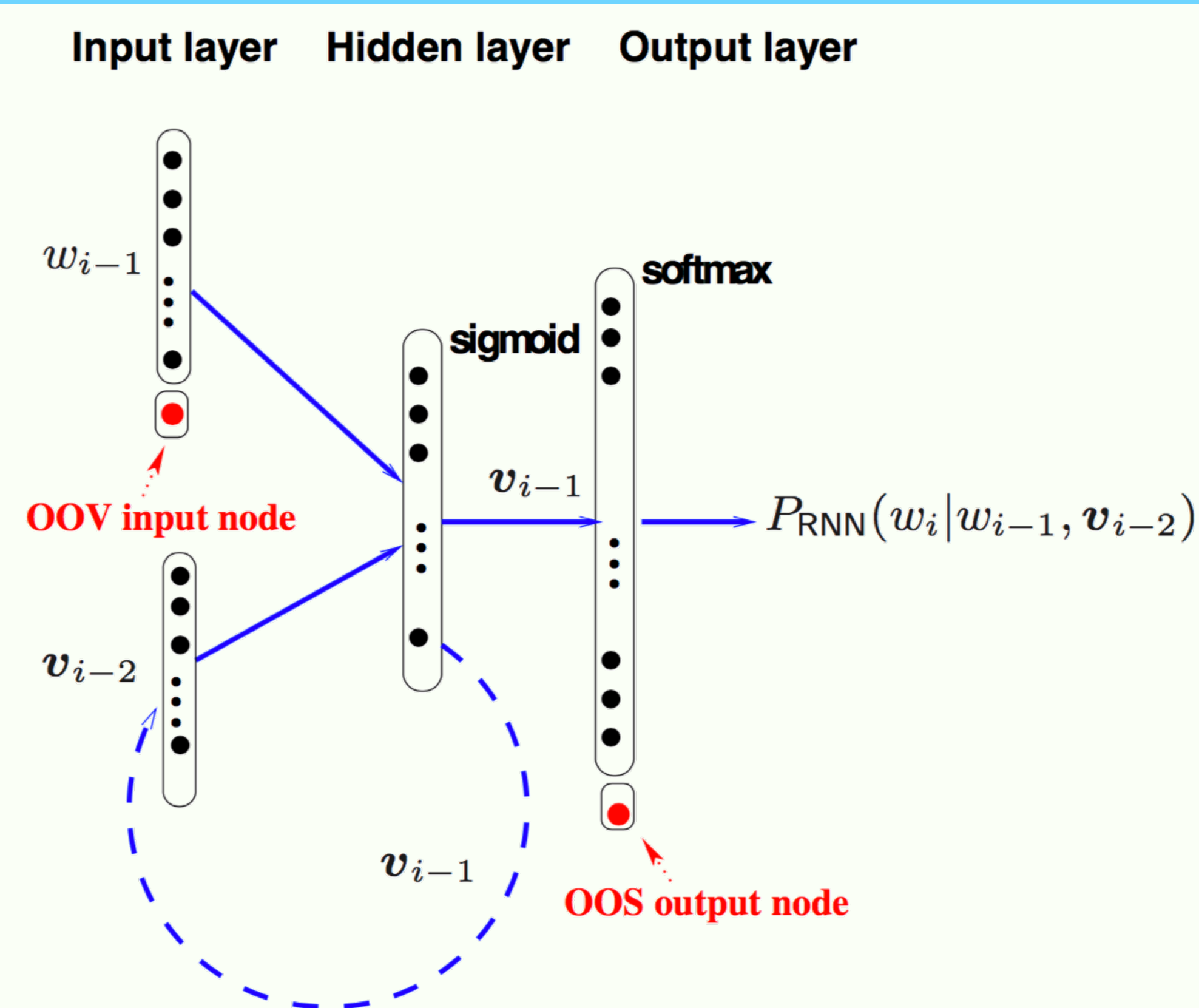
Existing RNNLM training algorithms

- Minimize the cross entropy (CE) using stochastic gradient descent (SGD)
- SGD uses no higher order gradient information, models no correlation between parameters, poorly captures error cost function curvature
- Quadratic approximation to error cost function using Newton methods
- Storing and computing Hessian matrix and its inverse are problematic
- Quasi-Newton methods, e.g. Hessian-free optimization applied to DNN acoustic modeling; iterative conjugate gradient (CG) method
- CG search very expensive for large datasets in Hessian-free optimization

Our approach

- Limited-memory Broyden Fletcher Goldfarb Shannon (L-BFGS) based 2nd order optimization for RNNLM training
- Efficiently approximates the product between inverse Hessian and gradient vector via a recursion over past gradients
- Only require a few vectors representing finite number of past updates of the matrix-vector product, so it's memory efficient

Recurrent Neural Network LMs



RNNLM description

- Vector representation of complete word history $h_1^{i-1} = \langle w_{i-1}, \dots, w_1 \rangle$
- Sigmoid hidden layer activation
- Shortlist output vocabulary plus out-of-shortlist (OOS) output node
- Output probability linearly interpolated with n-gram LMs

$$P(w_i | h_1^{i-1}) = \lambda P_{NG}(w_i | h_1^{i-1}) + (1 - \lambda) P_{RNN}(w_i | h_1^{i-1})$$

RNNLM Training Using SGD

Cross entropy training criterion

$$J^{CE}(\theta) = -\frac{1}{N_w} \sum_{i=1}^{N_w} \ln P_{RNN}(w_i | h_i)$$

- where N_w is the total words of a given sequence
- $P_{RNN}(w_i | h_i) = f_{softmax}(v_{i-1}; \theta)$

SGD training procedure for RNNLM

- Parameter update: $\theta[t+1] = \theta[t] - \eta \frac{\partial J^{CE}(\theta)}{\partial \theta} \Big|_{\theta=\theta[t]}$ θ is the layer wise weight matrix, this is applied to all layers
- Gradient stats for output layer weights: $\frac{\partial J^{CE}(\theta)}{\partial \theta} = -\frac{1}{N_w} \sum_{i=1}^{N_w} v_i \xi_i^T$ ξ_i is the error cost vector $\xi_{i,j} = \delta(w_j | h_i) - P_{RNN}(w_j | h_i)$ $\delta(w_j | h_i) = 0$ or 1 (target prob.)

- Back propagate, e.g. to recurrent layer:

$$\frac{\partial J^{CE}(\zeta)}{\partial \zeta} = -\frac{1}{N_w} \sum_{i=1}^{N_w} v_{i-1} (\xi_i \odot u_i)^T$$

ζ is the recurrent layer weight matrix $u_{i,j} = v_{i,j}(1 - v_{i,j})$ \odot denotes elementwise multiplication

- Back propagation through time (BPTT):

$$\frac{\partial J^{CE}(\zeta)}{\partial \zeta} = -\frac{1}{N_w} \sum_{i=1, \tau=1}^{N_w, N_\tau} v_{i-\tau-1} (\xi_{i-\tau} \odot u_{i-\tau})^T$$

RNNLM Training Using L-BFGS

Newton methods

- Can model the correlation between model parameters using a quadratic approximation to the error cost function

$$J^{CE}(\theta[t] + \Delta\theta) \approx J^{CE}(\theta[t]) + \Delta\theta^T \frac{\partial J^{CE}(\theta)}{\partial \theta} \Big|_{\theta=\theta[t]} + \frac{1}{2} \Delta\theta^T H_t \Delta\theta$$

- Newton direction

$$\Delta\theta = H_t^{-1} \frac{\partial J^{CE}(\theta)}{\partial \theta} \Big|_{\theta=\theta[t]}$$

where the Hessian matrix is computed as $H_{t,i,j} = \frac{\partial^2 J^{CE}(\theta)}{\partial \theta_i \partial \theta_j} \Big|_{\theta=\theta[t]}$

L-BFGS method for RNNLM training

Algorithm 1 For RNNLM output layer weights θ , L-BFGS algorithm approximates inverse Hessian gradient matrix-vector product

- $q_t \leftarrow \frac{\partial J^{CE}(\theta)}{\partial \theta} \Big|_{\theta=\theta[t]}$
- for** $i = t-1, t-2, \dots, t-m$ **do**
- $s_i \leftarrow \theta[i+1] - \theta[i]$, $y_i \leftarrow q_{i+1} - q_i$ (past gradient)
- $\rho_i \leftarrow \frac{1}{y_i^T s_i}$, $\alpha_i \leftarrow \rho_i s_i^T q_t$
- end for**
- $B_t^0 \leftarrow \frac{y_{t-m} s_{t-m}^T}{y_{t-m}^T y_{t-m}}$, $z \leftarrow B_t^0 q_t$
- for** $i = t-m, t-m+1, \dots, t-1$ **do**
- $\beta_i \leftarrow \rho_i y_i^T z$, $z \leftarrow z + (\alpha_i - \beta_i) s_i$
- end for**
- $H_t^{-1} \frac{\partial J^{CE}(\theta)}{\partial \theta} \Big|_{\theta=\theta[t]} \leftarrow z$ (parameter update direction)

Efficient GPU based training parallelization

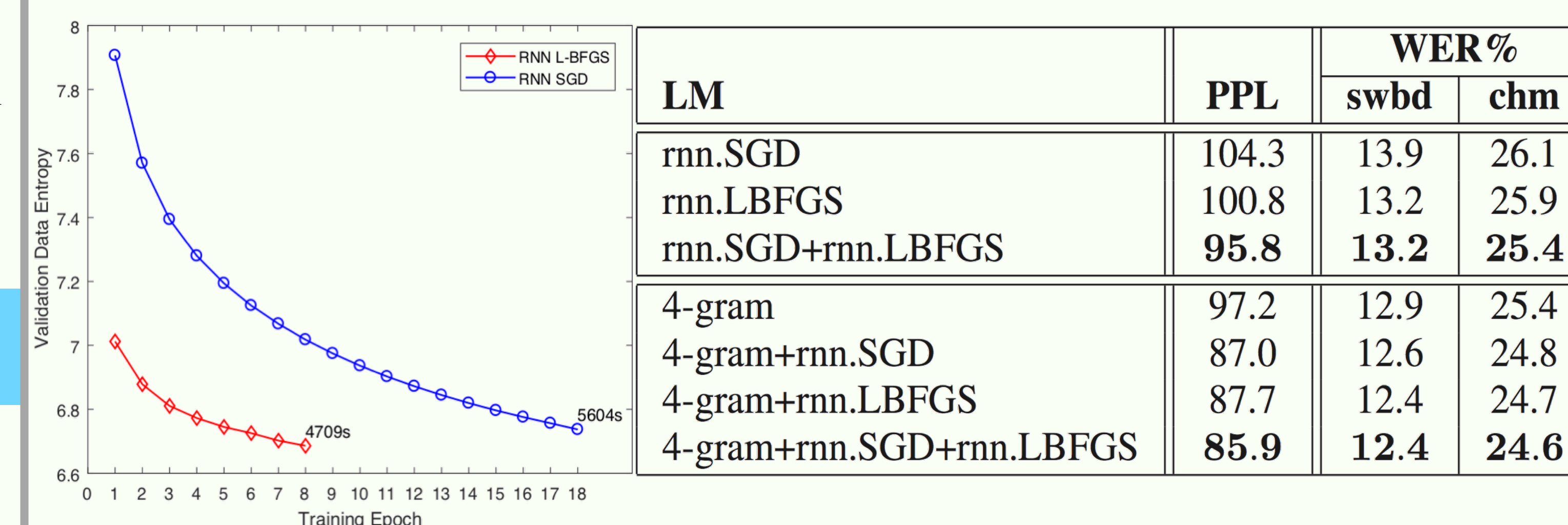
- L-BFGS for RNNLM is implemented as an extension to CUED-RNNLM
- Integrated into an efficient bunch mode GPU parallelization algorithm

Experiments and Results

Experiment setup

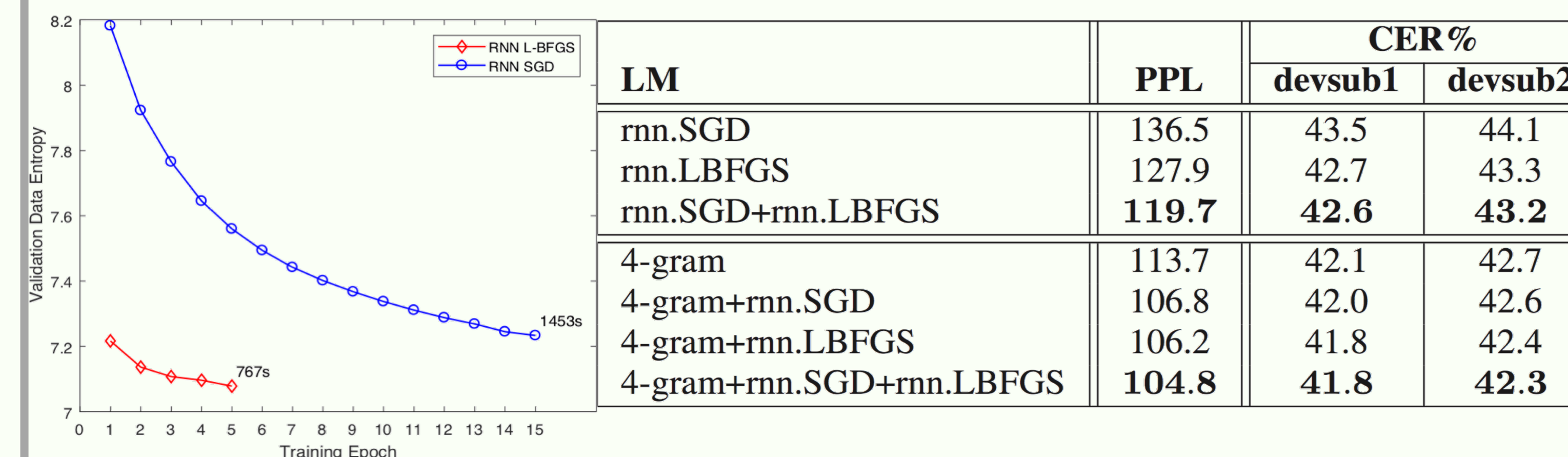
- Datasets: the Switchboard English system (SWBD), 300 hours, 3.6M words of acoustic transcripts and a 30k words lexicon; Babel Cantonese system, 175 hours, 1.1M words of transcripts and a 25k vocabulary
- Acoustic model: MPE trained stacked hybrid DNN-HMM by HTK toolkit
- Target of acoustic DNN: SWBD, 12k tied states; Babel Cantonese, 6k tied states
- RNNLM model structure and training: 512 hidden nodes and Sigmoid activation; SGD training with newbob scheduling plus momentum or L-BFGS method
- Evaluation method: perplexity (PPL) and word error rate (WER)
- GPU card: NVidia K40 GPUs; used to measure speed

Results on Switchboard



- Convergence: 19 epochs using 5604s for SGD; 9 epochs using 4709s for L-BFGS
- 0.7% abs. WER reductions obtained by L-BFGS before interpolation with 4-gram

Results on Babel Cantonese



- Convergence: 16 epochs using 1453s for SGD; 6 epochs using 767s for L-BFGS
- 0.8% abs. WER reductions obtained by L-BFGS before interpolation with 4-gram
- Observed on both tasks, the combination between SGD and L-BFGS is complementary since consistent improvements are obtained

Conclusion

L-BFGS optimization for RNNLM training & Future work

- Successfully applied to RNNLM training
- Consistent improvements over SGD on multiple speech recognition tasks
- Future research on L-BFGS training of advanced forms of NNs