

# Deep Bayesian Natural Language Processing

Jen-Tzung Chien

Department of Electrical and Computer Engineering  
National Chiao Tung University, Hsinchu, Taiwan  
jtchien@nctu.edu.tw

## 1 Introduction

This *introductory* tutorial addresses the advances in deep Bayesian learning for natural language with ubiquitous applications ranging from speech recognition (Saon and Chien, 2012; Chan et al., 2016) to document summarization (Chang and Chien, 2009), text classification (Blei et al., 2003; Zhang et al., 2015), text segmentation (Chien and Chueh, 2012), information extraction (Narasimhan et al., 2016), image caption generation (Vinyals et al., 2015; Xu et al., 2015), sentence generation (Li et al., 2016), dialogue control (Zhao and Eskenazi, 2016), sentiment classification, recommendation system, question answering (Sukhbaatar et al., 2015) and machine translation (Bahdanau et al., 2014), to name a few. Traditionally, “deep learning” is taken to be a learning process where the inference or optimization is based on the real-valued deterministic model. The “semantic structure” in words, sentences, entities, actions and documents drawn from a large vocabulary may not be well expressed or correctly optimized in mathematical logic or computer programs. The “distribution function” in discrete or continuous latent variable model for natural language may not be properly decomposed or estimated. This tutorial addresses the fundamentals of statistical models and neural networks, and focus on a series of advanced Bayesian models and deep models including hierarchical Dirichlet process (Teh et al., 2006), Chinese restaurant process (Blei et al., 2010), hierarchical Pitman-Yor process (Teh, 2006), Indian buffet process (Ghahramani and Griffiths, 2005), recurrent neural network (Mikolov et al., 2010; Van Den Oord et al., 2016), long short-term memory (Hochreiter and Schmidhuber, 1997; Cho et al., 2014), sequence-to-sequence model (Sutskever et al., 2014), variational auto-encoder (Kingma and Welling, 2014),

generative adversarial network (Goodfellow et al., 2014), attention mechanism (Chorowski et al., 2015; Seo et al., 2016), memory-augmented neural network (Graves et al., 2014; Sukhbaatar et al., 2015), skip neural network (Campos et al., 2018), stochastic neural network (Bengio et al., 2014; Miao et al., 2016), predictive state neural network (Downey et al., 2017) and policy neural network (Mnih et al., 2015; Yu et al., 2017). We present how these models are connected and why they work for a variety of applications on symbolic and complex patterns in natural language. The variational inference and sampling method are formulated to tackle the optimization for complicated models (Rezende et al., 2014). The word and sentence embeddings, clustering and co-clustering are merged with linguistic and semantic constraints. A series of case studies and domain applications are presented to tackle different issues in deep Bayesian processing, learning and understanding. At last, we will point out a number of directions and outlooks for future studies.

## 2 Objective of tutorial

Owing to the current growth in research and related emerging technologies in machine learning and deep learning, it is timely to introduce this tutorial to a large number of researchers and practitioners who are attending ACL 2019 and working on statistical models, deep neural networks, sequential learning and natural language processing and understanding. To the best of our knowledge, there is no similar tutorial presented in previous ACLs. This three-hour tutorial will concentrate on a wide range of theories and applications and systematically present the recent advances in deep Bayesian learning which are impacting the communities of machine learning, natural language processing and human language technology.

### 3 Tutorial outline

- Introduction
  - motivation and background
  - probabilistic models
  - neural networks
  - modern natural language models
- Bayesian Learning
  - inference and optimization
  - variational Bayesian (VB) inference
  - Monte Carlo Markov chain (MCMC)
  - Bayesian nonparametrics (BNP)
  - hierarchical theme and topic model
  - hierarchical Pitman-Yor-Dirichlet proc.
  - nested Indian buffet process
- Deep Learning
  - deep unfolded topic model
  - gated recurrent neural network (RNN)
  - generative adversarial network (GAN)
  - memory-augmented neural network
  - sequence-to-sequence learning
  - convolutional neural network (CNN) (Coffee Break)
  - dilated recurrent neural network
  - attention network using transformer
- Deep Bayesian Processing and Learning
  - Bayesian recurrent neural network
  - variational auto-encoder (VAE)
  - variational recurrent auto-encoder
  - stochastic temporal convolutional net
  - stochastic recurrent neural network
  - regularized recurrent neural network
  - stochastic learning & normalizing flows
  - VAE with VampPrior
  - skip recurrent neural network
  - temporal difference VAE
  - Markov recurrent neural network
  - reinforcement learning & understanding
  - sequence GAN
- Summarization and Future Trend

### 4 Target audience

This tutorial will be useful to research students working in natural language processing and researchers who would like to explore machine learning, deep learning and sequential learning. The prerequisite knowledge includes calculus, linear algebra, probability and statistics. This tutorial serves the objectives to introduce novices to major topics within deep Bayesian learning, motivate and explain a topic of emerging importance for natural language understanding, and present a novel synthesis combining distinct lines of machine learning work.

### 5 Description of tutorial content

The presentation of this tutorial is arranged into five parts. First of all, we share the current status of researches and applications on natural language processing, statistical modeling and deep neural network (Bahdanau et al., 2014), and address the key issues in deep Bayesian learning for discrete-valued observation data and latent semantics. Modern natural language models are introduced to address how data analysis is performed from language processing to semantic learning, memory networking, knowledge mining and understanding. Secondly, we address a number of Bayesian models ranging from latent variable model to VB inference (Chien and Chueh, 2011; Chien, 2015b; Chien and Chang, 2014), MCMC sampling and BNP learning (Chien, 2016, 2015a, 2018; Watanabe and Chien, 2015) for hierarchical, thematic and sparse topics from natural language. In the third part, a series of deep models including deep unfolding (Chien and Lee, 2018), RNN (Hochreiter and Schmidhuber, 1997), GAN (Goodfellow et al., 2014), memory network (Weston et al., 2015; Chien and Lin, 2018; Tsou and Chien, 2017), sequence-to-sequence learning (Graves et al., 2006; Gehring et al., 2017), CNN (Kalchbrenner et al., 2014; Xingjian et al., 2015; Dauphin et al., 2017), dilated RNN (Chang et al., 2017) and attention network with transformer (Vaswani et al., 2017; Devlin et al., 2018) are introduced. The coffee break is arranged within this part. Next, the fourth part focuses on a variety of advanced studies which illustrate how deep Bayesian learning is developed to infer the sophisticated recurrent models for natural language understanding. In particular, the Bayesian RNN (Gal and Ghahramani,

2016; Chien and Ku, 2016), VAE (Kingma and Welling, 2014), variational recurrent auto-encoder (Chien and Wang, 2019), neural variational learning (Serban et al., 2017; Chung et al., 2015), stochastic temporal convolutional network (Aksan and Hilliges, 2019), neural discrete representation (Jang et al., 2017; van den Oord et al., 2017), recurrent ladder network (Rasmus et al., 2015; Prémont-Schwarz et al., 2017), stochastic recurrent neural network (Fraccaro et al., 2016; Goyal et al., 2017; Chien and Kuo, 2017), predictive state neural network (Downey et al., 2017), Markov recurrent neural network (Venkatraman et al., 2017; Kuo and Chien, 2018), reinforcement learning (Tegho et al., 2017), sequence GAN (Yu et al., 2017), and temporal difference VAE (Gregor et al., 2019) are introduced in various deep models. Enhancing the prior/posterior representation in variational inference is addressed (Rezende and Mohamed, 2015; Tomczak and Welling, 2018). These sophisticated models open a window to numerous practical tasks such as reading comprehension, sentence generation, dialogue system, question answering and machine translation. Variational inference methods based on normalizing flows (Rezende and Mohamed, 2015) and “variational mixture of posteriors” prior (VampPrior) (Tomczak and Welling, 2018) are addressed. Posterior collapse problem in variational sequential learning is compensated. In the final part, we spotlight on some future directions for deep language understanding which can handle the challenges of big data, heterogeneous condition and dynamic system. In particular, deep learning, structural learning, temporal and spatial modeling, long history representation and stochastic learning are emphasized. Slides of this tutorial are available at (<http://chien.cm.nctu.edu.tw/home/acl-tutorial>).

## 6 Instructor

Jen-Tzung Chien is now with the Department of Electrical and Computer Engineering, National Chiao Tung University, Taiwan, where he is currently the University Chair Professor. He held the visiting researcher position with the IBM T. J. Watson Research Center, Yorktown Heights, NY, in 2010. His research interests include machine learning, deep learning, natural language processing and computer vision. He served as the associate editor of the IEEE Signal Processing Letters in 2008-2011, the guest editor of the

IEEE Transactions on Audio, Speech and Language Processing in 2012, the organization committee member of ICASSP 2009, the area coordinator of Interspeech 2012, EUSIPCO 2017-2019, the program chair of ISCSLP 2018, the general chair of MLSP 2017, and currently serves as an elected member of the IEEE Machine Learning for Signal Processing (MLSP) Technical Committee. He received the Best Paper Award of IEEE Automatic Speech Recognition and Understanding Workshop in 2011 and the AAPM Farrington Daniels Award in 2018. Dr. Chien has published extensively including the books “Bayesian Speech and Language Processing”, Cambridge University Press, in 2015, and “Source Separation and Machine Learning”, Academic Press, in 2018. He has served as the Tutorial Speaker for APSIPA 2013, ISCSLP 2014, Interspeech 2013, 2016, ICASSP 2012, 2015, 2017, COLING 2018, AAAI 2019, KDD 2019, and IJCAI 2019. (<http://chien.cm.nctu.edu.tw/>)

## Acknowledgments

This work was partially supported by the Ministry of Science and Technology, Taiwan, under MOST 108-2634-F-009-003.

## References

- Emre Aksan and Otmar Hilliges. 2019. STCN: stochastic temporal convolutional networks. In *Proc. of International Conference on Learning Representations*, page 2019.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- Yoshua Bengio, Eric Thibodeau-Laufer, Guillaume Alain, and Jason Yosinski. 2014. Deep generative stochastic networks trainable by backprop. In *Proc. of International Conference on Machine Learning*.
- David M. Blei, Thomas L. Griffiths, and Michael I. Jordan. 2010. The nested Chinese restaurant process and Bayesian nonparametric inference of topic hierarchies. *Journal of the ACM*, 57(2). Article 7.
- David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3(3):993–1022.
- Víctor Campos, Brendan Jou, Xavier Giró-i Nieto, Jordi Torres, and Shih-Fu Chang. 2018. Skip RNN: learning to skip state updates in recurrent neural networks. In *Proc. of International Conference on Learning Representations*.

- William Chan, Navdeep Jaitly, Quoc Le, and Oriol Vinyals. 2016. Listen, attend and spell: A neural network for large vocabulary conversational speech recognition. In *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 4960–4964.
- Shiyu Chang, Yang Zhang, Wei Han, Mo Yu, Xiaoxiao Guo, Wei Tan, Xiaodong Cui, Michael Witbrock, Mark A Hasegawa-Johnson, and Thomas S Huang. 2017. Dilated recurrent neural networks. In *Advances in Neural Information Processing Systems*, pages 77–87.
- Ying-Lang Chang and Jen-Tzung Chien. 2009. Latent Dirichlet learning for document summarization. In *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 1689–1692.
- Jen-Tzung Chien. 2015a. Hierarchical Pitman-Yor-Dirichlet language model. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 23(8):1259–1272.
- Jen-Tzung Chien. 2015b. Laplace group sensing for acoustic models. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 23(5):909–922.
- Jen-Tzung Chien. 2016. Hierarchical theme and topic modeling. *IEEE Transactions on Neural Networks and Learning Systems*, 27(3):565–578.
- Jen-Tzung Chien. 2018. Bayesian nonparametric learning for hierarchical and sparse topics. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 26(2):422–435.
- Jen-Tzung Chien and Ying-Lan Chang. 2014. Bayesian sparse topic model. *Journal of Signal Processing Systems*, 74(3):375–389.
- Jen-Tzung Chien and Chuang-Hua Chueh. 2011. Dirichlet class language models for speech recognition. *IEEE Transactions on Audio, Speech, and Language Processing*, 19(3):482–495.
- Jen-Tzung Chien and Chuang-Hua Chueh. 2012. Topic-based hierarchical segmentation. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(1):55–66.
- Jen-Tzung Chien and Yuan-Chu Ku. 2016. Bayesian recurrent neural network for language modeling. *IEEE Transactions on Neural Networks and Learning Systems*, 27(2):361–374.
- Jen-Tzung Chien and Kuan-Ting Kuo. 2017. Variational recurrent neural networks for speech separation. In *Proc. of Annual Conference of International Speech Communication Association*, pages 1193–1197.
- Jen-Tzung Chien and Chao-Hsi Lee. 2018. Deep unfolding for topic models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(2):318–331.
- Jen-Tzung Chien and Ting-An Lin. 2018. Supportive attention in end-to-end memory networks. In *Proc. of IEEE International Workshop on Machine Learning for Signal Processing*, pages 1–6.
- Jen-Tzung Chien and Chun-Wei Wang. 2019. Variational and hierarchical recurrent autoencoder. In *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 3202–3206.
- Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proc. of Conference on Empirical Methods in Natural Language Processing*, pages 1724–1734.
- Jan K. Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio. 2015. Attention-based models for speech recognition. In *Advances in Neural Information Processing Systems*, pages 577–585.
- Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron C Courville, and Yoshua Bengio. 2015. A recurrent latent variable model for sequential data. In *Advances in Neural Information Processing Systems*, pages 2980–2988.
- Yann N. Dauphin, Angela Fan, Michael Auli, and David Grangier. 2017. Language modeling with gated convolutional networks. In *Proc. of International Conference on Machine Learning*, pages 933–941.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Carlton Downey, Ahmed Hefny, Byron Boots, Geoffrey J Gordon, and Boyue Li. 2017. Predictive state recurrent neural networks. In *Advances in Neural Information Processing Systems*, pages 6055–6066.
- Marco Fraccaro, Søren Kaae Sønderby, Ulrich Paquet, and Ole Winther. 2016. Sequential neural models with stochastic layers. In *Advances in Neural Information Processing Systems*, pages 2199–2207.
- Yarin Gal and Zoubin Ghahramani. 2016. A theoretically grounded application of dropout in recurrent neural networks. In *Advances in Neural Information Processing Systems*, pages 1019–1027.
- Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. 2017. Convolutional sequence to sequence learning. In *Proc. of International Conference on Machine Learning*, pages 1243–1252.

- Zoubin Ghahramani and Thomas L. Griffiths. 2005. Infinite latent feature models and the Indian buffet process. In *Advances in Neural Information Processing Systems*, volume 18, pages 475–482.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2672–2680.
- Anirudh Goyal, Alessandro Sordoni, Marc-Alexandre Côté, Nan Ke, and Yoshua Bengio. 2017. Z-forcing: Training stochastic recurrent networks. In *Advances in Neural Information Processing Systems 30*, pages 6713–6723.
- Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. 2006. Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks. In *Proc. of International Conference on Machine Learning*, pages 369–376.
- Alex Graves, Greg Wayne, and Ivo Danihelka. 2014. Neural Turing machines. *arXiv preprint arXiv:1410.5401*.
- Karol Gregor, George Papamakarios, Frederic Besse, Lars Buesing, and Theophane Weber. 2019. Temporal difference variational auto-encoder. In *Proc. of International Conference on Learning Representations*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Computation*, 9(8):1735–1780.
- Eric Jang, Shixiang Gu, and Ben Poole. 2017. Categorical reparameterization with Gumbel-softmax. In *Proc. of International Conference on Learning Representations*.
- Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. 2014. A convolutional neural network for modelling sentences. In *Proc. of Annual Meeting of the Association for Computational Linguistics*, pages 655–665.
- Diederik P Kingma and Max Welling. 2014. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*.
- Che-Yu Kuo and Jen-Tzung Chien. 2018. Markov recurrent neural networks. In *Proc. of IEEE International Workshop on Machine Learning for Signal Processing*, pages 1–6.
- Jiwei Li, Will Monroe, Alan Ritter, and Dan Jurafsky. 2016. Deep reinforcement learning for dialogue generation. In *Proc. of Conference on Empirical Methods in Natural Language Processing*, pages 1192–1202.
- Yishu Miao, Lei Yu, and Phil Blunsom. 2016. Neural variational inference for text processing. In *Proc. of International Conference on Machine Learning*, pages 1727–1736.
- Tomáš Mikolov, Martin Karafiát, Lukáš Burget, Jan Černocký, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model. In *Proc. of Annual Conference of International Speech Communication Association*, pages 1045–1048.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. 2015. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533.
- Karthik Narasimhan, Adam Yala, and Regina Barzilay. 2016. Improving information extraction by acquiring external evidence with reinforcement learning. In *Proc. of Conference on Empirical Methods in Natural Language Processing*, pages 2355–2365.
- Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. 2017. Neural discrete representation learning. In *Advances in Neural Information Processing Systems*, pages 6309–6318.
- Isabeau Prémont-Schwarz, Alexander Ilin, Tele Hao, Antti Rasmus, Rinu Boney, and Harri Valpola. 2017. Recurrent ladder networks. In *Advances in Neural Information Processing Systems*, pages 6011–6021.
- Antti Rasmus, Mathias Berglund, Mikko Honkala, Harri Valpola, and Tapani Raiko. 2015. Semi-supervised learning with ladder networks. In *Advances in Neural Information Processing Systems*, pages 3546–3554.
- Danilo Jimenez Rezende and Shakir Mohamed. 2015. Variational inference with normalizing flows. In *Proc. of International Conference on Machine Learning*, pages 1530–1538.
- Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. 2014. Stochastic backpropagation and approximate inference in deep generative models. In *Proc. of International Conference on Machine Learning*, pages 1278–1286.
- George Saon and Jen-Tzung Chien. 2012. Large-vocabulary continuous speech recognition systems: A look at some recent advances. *IEEE Signal Processing Magazine*, 29(6):18–33.
- Minjoon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. 2016. Bidirectional attention flow for machine comprehension. *arXiv preprint arXiv:1611.01603*.
- Iulian V. Serban, Alessandro Sordoni, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron Courville, and Yoshua Bengio. 2017. A hierarchical latent variable encoder-decoder model for generating dialogues. In *Proc. of AAAI Conference on Artificial Intelligence*, pages 3295–3301.

- Sainbayar Sukhbaatar, Arthur Szlam, Jason Weston, and Rob Fergus. 2015. End-to-end memory networks. In *Advances in Neural Information Processing Systems*, pages 2440–2448.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, pages 3104–3112. .
- Christopher Tegho, Paweł Budzianowski, and Milica Gašić. 2017. Uncertainty estimates for efficient neural network-based dialogue policy optimisation. *arXiv preprint arXiv:1711.11486*.
- Yee Whye Teh. 2006. A hierarchical Bayesian language model based on Pitman-Yor processes. In *Proc. of International Conference on Computational Linguistics and Annual Meeting of the Association for Computational Linguistics*, pages 985–992.
- Yee Whye Teh, Michael I. Jordan, Matthew J. Beal, and David M. Blei. 2006. Hierarchical Dirichlet process. *Journal of the American Statistical Association*, 101(476):1566–1581.
- Jakub Tomczak and Max Welling. 2018. VAE with a VampPrior. In *International Conference on Artificial Intelligence and Statistics*, pages 1214–1223.
- Kai-Wei Tsou and Jen-Tzung Chien. 2017. Memory augmented neural network for source separation. In *Proc. of IEEE International Workshop on Machine Learning for Signal Processing*, pages 1–6.
- Aaron Van Den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. 2016. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008.
- Arun Venkatraman, Nicholas Rhinehart, Wen Sun, Lerrel Pinto, Martial Hebert, Byron Boots, Kris Kitani, and J Bagnell. 2017. Predictive-state decoders: Encoding the future into recurrent networks. In *Advances in Neural Information Processing Systems*, pages 1172–1183.
- Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition*, pages 3156–3164.
- Shinji Watanabe and Jen-Tzung Chien. 2015. *Bayesian Speech and Language Processing*. Cambridge University Press.
- Jason Weston, Sumit Chopra, and Antoine Bordes. 2015. Memory networks. In *Proc. of International Conference on Learning Representation*.
- Shi Xingjian, Zhoung Chen, Hao Wang, Dit-Yan Yeung, Wai-kin Wong, and Wang-chun Woo. 2015. Convolutional LSTM network: A machine learning approach for precipitation nowcasting. In *Advances in Neural Information Processing Systems*, pages 802–810.
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *Proc. of International Conference on Machine Learning*, pages 2048–2057.
- Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. 2017. SeqGAN: Sequence generative adversarial nets with policy gradient. In *Proc. of AAAI Conference on Artificial Intelligence*, volume 31, pages 2852–2858.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In *Advances in Neural Information Processing Systems*, pages 649–657.
- Tiancheng Zhao and Maxine Eskenazi. 2016. Towards end-to-end learning for dialog state tracking and management using deep reinforcement learning. *arXiv preprint arXiv:1606.02560*.