

CONTACT INFORMATION (1) 41-35 Elbertson St, Elmhurst, NY 11373, U.S. *Phone:* (1) +1(347)-891-8622 (U.S.)  
 (2) FIT, Tsinghua University, Beijing, 100084, P.R. China. (2) +86-13581700448 (China)  
*E-mail:* fanmiao.cslt.thu@gmail.com  
*WWW:* [http://cslt.riit.tsinghua.edu.cn/mediawiki/index.php/Miao\\_Fan](http://cslt.riit.tsinghua.edu.cn/mediawiki/index.php/Miao_Fan)

RESEARCH INTERESTS **Text Mining & Social Computing**

- *Information Extraction*
- *Social Tag Recommendation*

EDUCATION **New York University, New York, U.S.A. March 2015 - March 2016**

*Joint-supervised Ph.D. Candidate in Computer Science.*

*Junior Research Scientist, Courant Institute of Mathematical Sciences.*

- Research Topic: “Large-scale Entity Relation Extraction based on Low-dimensional Representations” (Proteus Project)<sup>2</sup>.
- Advisor: Prof. Ralph Grishman<sup>3</sup>

**Tsinghua University, Beijing, P.R. China. September 2012 - June 2017 (expected)**

*Ph.D. Candidate in Computer Science.*

- Dissertation Topic: “Large-scale Knowledge Base Completion based on Low-dimensional Representations”.
- Co-advisors: Prof. Thomas Fang Zheng<sup>4</sup> and Prof. Qiang Zhou<sup>5</sup>.
- GPA: **87.3**/100.

**Beijing University of Posts and Telecommunications September 2008 - June 2012**

*B.Eng. in Software Engineering.*

- Dissertation Topic: “Chinese Natural Question Generation Based on Wiki Knowledge”<sup>6</sup>, **Best B.Eng. Dissertation Award.**
- Advisor: Prof. Guoshi Wu<sup>7</sup>.
- GPA: **95.4**/100 (Rank **1**/97).

MY MONOGRAPHS **Miao Fan, Chao Li; *DIY Machine Learning Systems for Kaggle Competitions with Python Programming*, Printing, Tsinghua University Press.**

**Miao Fan, *Distributed Machine Learning and Big Data Analysis with PySpark*, Drafting, Tsinghua University Press.**

PATENTS **Miao Fan, Doo Soon Kim. Detecting Table Region in PDF Documents using Distant Supervision<sup>8</sup>. U.S. Patent submitted.**

**Miao Fan, Qiang Zhou. The Approach to Social Tag Generation and Recommendation based on Monolingual Word Alignments. P.R.China Patent submitted.**

<sup>1</sup><http://scholar.google.com/citations?user=aPlHReAAAAAJ&hl=en>

<sup>2</sup><http://nlp.cs.nyu.edu/index.shtml>

<sup>3</sup><http://scholar.google.com/citations?user=blwKAKUAAAAJ>

<sup>4</sup>[http://scholar.google.com/citations?user=H3MX\\_8IAAAAAJ&hl=en](http://scholar.google.com/citations?user=H3MX_8IAAAAAJ&hl=en)

<sup>5</sup><http://cslt.riit.tsinghua.edu.cn/~qzhou/eng/index.htm>

<sup>6</sup>[http://cslt.riit.tsinghua.edu.cn/mediawiki/images/a/a9/B.ENG\\_Dissertation\\_-Miao\\_Fan-.pdf](http://cslt.riit.tsinghua.edu.cn/mediawiki/images/a/a9/B.ENG_Dissertation_-Miao_Fan-.pdf)

<sup>7</sup><http://baike.baidu.com/view/9021485.htm>

<sup>8</sup><http://arxiv.org/abs/1506.08891>

- Miao Fan**, Qiang Zhou, Thomas Fang Zheng, Ralph Grishman. Distributed Representation Learning for Knowledge Bases with Entity Descriptions. *Pattern Recognition Letters*<sup>9</sup> (**IF:1.55**), Elsevier. SCIE Index. Full length article submitted.
- Miao Fan**, Qiang Zhou, Andrew Abel, Thomas Fang Zheng, Ralph Grishman. Probabilistic Belief Embedding for Large-scale Knowledge Population. *Cognitive Computation*<sup>10</sup> (**IF: 1.44**), Springer. SCIE Index. Full length article submitted.
- Miao Fan**, Qiang Zhou, Thomas Fang Zheng. Learning Embedding Representations for Imperfect Knowledge Repository.<sup>11</sup> 2016 IEEE/WIC/ACM International Conference on Web Intelligence (*WI'16*). Full research paper submitted.
- Miao Fan**, Qiang Zhou, Thomas Fang Zheng. Distant Supervision for Entity Linking.<sup>12</sup> The 29th Pacific Asia Conference on Language, Information and Computation (*PACLIC'15*), pp. 79-86. Full paper, oral presentation.
- Miao Fan**, Kai Cao, Yifan He, Ralph Grishman. Jointly Embedding Relations and Mentions for Knowledge Population.<sup>13</sup> The 10th Recent Advances in Natural Language Processing (*RANLP'15*), pp. 186-191. Poster paper.
- Miao Fan**, Qiang Zhou, Thomas Fang Zheng, Ralph Grishman. Large Margin Nearest Neighbor Embedding for Knowledge Representation.<sup>14</sup> The 2015 IEEE/WIC/ACM Web Intelligence Conference (*WI'15*), pp. 53-59, 6-9 December 2015, Singapore. Full paper, oral presentation.
- Miao Fan**, Deli Zhao, Qiang Zhou, Zhiyuan Liu, Thomas Fang Zheng, Edward Y. Chang. Distant Supervision for Relation Extraction with Matrix Completion.<sup>15</sup> The 52th Annual Meeting of the Association for Computational Linguistics (*ACL'14*), pp. 839-849. Full paper, oral presentation.
- Miao Fan**, Qiang Zhou, Emily Chang, Thomas Fang Zheng. Transition-based Knowledge Graph Embedding with Relational Mapping Properties.<sup>16</sup> The 28th Pacific Asia Conference on Language, Information and Computing (*PACLIC'14*), pp. 328-337. Full paper, oral presentation.
- Miao Fan**, Qiang Zhou, Thomas Fang Zheng. Mining the Personal Interests of Microbloggers via Exploiting Wikipedia Knowledge.<sup>17</sup> 15th International Conference on Intelligent Text Processing and Computational Linguistics (*CICLing'14*), pp. 188-200. Full paper, poster presentation.
- Miao Fan**, Qiang Zhou, Thomas Fang Zheng. Content-based Semantic Tag Ranking for Recommendation.<sup>18</sup> The 2012 IEEE/WIC/ACM International Conference on Web Intelligence (*WI'12*), pp. 292-296. Short paper, oral presentation.
- Miao Fan**, Yingnan Xiao, Qiang Zhou. Bringing the Associative Ability to Social Tag Recommendation.<sup>19</sup> *ACL'12 Workshop on Graph-based Methods for Natural Language Processing*, pp. 44-54. Workshop paper, oral presentation.
- Miao Fan**, Guoshi Wu. Aspect Opinion Mining on Customer Reviews.<sup>20</sup> Proceedings of the 2011 International Conference on Informatics, Cybernetics, and Computer Engineering (*ICCE'11*) November 19-20, 2011, Melbourne, Australia. Advances in Intelligent and Soft Computing Volume 112, 2012, pp 27-33.

<sup>9</sup><http://www.journals.elsevier.com/pattern-recognition-letters/>

<sup>10</sup><http://link.springer.com/journal/12559>

<sup>11</sup><http://arxiv.org/pdf/1503.08155v1.pdf>

<sup>12</sup><http://aclweb.org/anthology/Y15-1010>

<sup>13</sup><http://arxiv.org/pdf/1504.01683v1.pdf>

<sup>14</sup><http://arxiv.org/pdf/1504.01684v1.pdf>

<sup>15</sup><http://arxiv.org/abs/1411.4455>

<sup>16</sup><http://pan.baidu.com/s/1pJqLPIb>

<sup>17</sup><http://pan.baidu.com/s/1bntmkpx>

<sup>18</sup><http://pan.baidu.com/s/1kT41lyf>

<sup>19</sup><http://pan.baidu.com/s/1o6wn7lg>

<sup>20</sup><http://pan.baidu.com/s/1jGj47dG>

- Enterprises Information Laboratory**, Beijing University of Posts and Telecommunications.  
*Team leader* **April, 2009 - May, 2011**  
Studying on developing *Feature-Opinion Recommender System*<sup>21</sup> funded by the National Innovative Experimental Program.  
Advisor: Prof. Guoshi Wu.
- Natural Language Processing and Computational Social Science Lab**, Tsinghua University.  
*Group member* **May, 2011 - December, 2011**  
Studying on social tagging at Sina micro-blog platform.  
Advisor: Prof. Maosong Sun<sup>22</sup> and Dr. Zhiyuan Liu<sup>23</sup>.
- The Association for Computational Linguistics.**  
*Member* **June, 2012 - Present**
- Center for Speech and Language Technology**, Tsinghua University.  
*Ph. D. candidate* **September, 2012 - Present**  
Exploring large scale knowledge extraction approaches which can be applied to build entity-based search engine or Goolge Knowledge Graph. Advisor: Prof. Thomas Fang Zheng and Prof. Qiang Zhou.
- DolphinNLP Group & Social Card Project.**  
*Founder* **June, 2012 - December, 2013**
- ACM Transactions on Intelligent Systems and Technology**<sup>24</sup> (IF<sup>25</sup>: 9.39)  
*Reviewer* **November, 2014**
- 2014 International Doctoral Forum**<sup>26</sup>  
*Technical Committee Track Chair*<sup>27</sup> **December 5-7, 2014**
- C++ Programming**, Tsinghua University.  
*Teacher Assistant of Instructor Chao Li*<sup>28</sup> **September, 2014 - January, 2015**
- Natural Language Processing**<sup>29</sup>, New York University.  
*Teacher Assistant of Instructor Adam Meyers*<sup>30</sup> **September, 2015 - December, 2015**  
Special talks of “Statistical NLP: A Machine Learning Perspective”<sup>31</sup>.
- The Eighth International Conference on Creative Content Technologies (CONTENT**

<sup>21</sup><http://pan.baidu.com/s/1jGj47dG>

<sup>22</sup><http://scholar.google.com/citations?user=zIgTOHMAAAAJ&hl=en>

<sup>23</sup><http://scholar.google.com/citations?user=dT0v5u0AAAAJ>

<sup>24</sup><http://tist.acm.org/editors.html>

<sup>25</sup><http://www.spinellis.gr/blog/20140808/>

<sup>26</sup><http://www.nwpu-aslp.org/phdforum/2014/index.html>

<sup>27</sup><http://www.nwpu-aslp.org/phdforum/2014/content/orgcommittee.html>

<sup>28</sup><http://dbgroup.cs.tsinghua.edu.cn/lichao/>

<sup>29</sup><http://cs.nyu.edu/courses/fall15/CSCI-UA.0480-006/>

<sup>30</sup><http://nlp.cs.nyu.edu/people/meyers.html>

<sup>31</sup><http://1drv.ms/1Saijf0>

2016)<sup>32</sup>

Technical Program Committee Member<sup>33</sup>

March 20-24, 2016

**Hulu Inc.**, Beijing, P.R. China.

*Occupation:* Research intern

May, 2013 - September, 2013

- *Brief Intro:* Researcher and computer model developer for extracting relation instances from movie plots.
- *Programming Language:* Python.
- *Supervisor:* Post Dr. Tao Xiong.

**Microsoft Research Asia, Machine Learning Group**<sup>34</sup>, Beijing, P.R. China.

*Occupation:* Research intern

April, 2014 - July, 2014

- *Brief Intro:* Studying the link prediction in large-scale incomplete knowledge base, i.e. Freebase and WordNet, based on the low-dimensional embedding representation of entities and relationships without extra free texts. Developing parallel algorithms for large-scale knowledge embedding under the Microsoft Cosmos System<sup>35</sup>.
- *Programming Language:* C++, C#, Microsoft-SCOPE (Cosmos).
- *Supervisor:* Researcher Jianwen Zhang<sup>36</sup>.

**Baidu Inc., Natural Language Group**, Beijing, P.R. China.

*Occupation:* Research intern

November, 2014 - January, 2015

- *Brief Intro:* R&D for Baidu Chatbot (Smart Search) System<sup>37</sup>. Modeling/designing multi-round dialogue algorithms/system for chatbot.
- *Programming Language:* C++.
- *Supervisor:* Senior Researcher Shiqi Zhao<sup>38</sup> and Dr. Rui Yan<sup>39</sup>.

**Bosch Research**, Palo Alto, CA. U.S.

*Occupation:* Research intern authorized by NYU

June, 2015 - September, 2015

- *Brief Intro:* Table detection and extraction from PDF files, such as electronic manuals and academic articles<sup>40</sup>. We also design the end-to-end system (PDFExtraction) for Bosch Research.
- *Programming Language:* Java.
- *Supervisor:* Senior Research Engineer, Doo Soon Kim<sup>41</sup>.

**DailyCast Lab**, Beijing, P.R.China.

*Occupation:* Leading Researcher and Developer

May, 2016 - Present

<sup>32</sup><http://www.iaria.org/conferences2016/CONTENT16.html>

<sup>33</sup><http://www.iaria.org/conferences2016/ComCONTENT16.html>

<sup>34</sup><http://research.microsoft.com/en-us/groups/ml/>

<sup>35</sup><http://blogs.msdn.com/b/seliot/archive/2010/11/05/cosmos-petabytes-perfectly-processed-perfunctorily.aspx>

<sup>36</sup>

<http://research.microsoft.com/en-us/people/jiazhan/>

<sup>37</sup><http://baike.baidu.com/view/14785371.htm?fr=aladdin>

<sup>38</sup><http://ir.hit.edu.cn/~zhaosq/>

<sup>39</sup><https://sites.google.com/site/ruiyan516/>

<sup>40</sup>We collect all free published articles (9,466) from ACL Anthology: <http://aclweb.org/anthology/>

<sup>41</sup><https://sites.google.com/site/2soonk/>

- *Brief Intro*: Leading the research and development of short-video recommender systems in Dailycast<sup>42</sup> App.
- *Programming Language*: Python and Java.

#### HONORS AND AWARDS

HUAGUANG 2nd Class Ph.D. Scholarship (1,000 CNY), 2015.  
 HUAWEI Ph.D. Fellowship (10,000 CNY), 2014.  
 RIIT Graduate Research Award (3,000 CNY), Tsinghua University, 2014.  
 Excellent Graduate Award in Tsinghua University, 2013.  
 National Scholarship for Undergraduates (8,000 CNY), 2011.  
 The IBM Chinese Excellent Undergraduate Scholarship (4,000 CNY), 2011.  
 Excellent Undergraduate Award in Beijing, China, 2011.  
 Excellent Academic Paper Award of 11th Creative Award in Beijing University of Posts and Telecommunications, 2011.  
 Great Leader Award of National Innovative Experimental Program, China, 2011.  
 MCM/ICM (Mathematical Contest in Modeling/ Interdisciplinary Contest in Modeling) **Meritorious Winner** (First Prize), 2011.  
 “Tang Jun&Sun Chunlan” Enterprise Scholarship (5,000 CNY), 2011.

#### PARTICIPATED CONTESTS

##### **Kaggle Competitions**

**April, 2015 - Present**

- *Profile*: <https://www.kaggle.com/michaelfan>
- *Rank*: 2,768th / 465,867 (Top 1%)
- *Awards*: Top25%(Bronze) \* 2
- *Tools*: I use *NLTK*<sup>43</sup>, *Scikit-learn*<sup>44</sup>, *Pandas*<sup>45</sup>, *Xgboost*<sup>46</sup>, *TensorFlow*<sup>47</sup> and *Spark*<sup>48</sup> to process the data, build new models and make predictions.

#### FAVORITE (READ) BOOKS

Yoshua Bengio, Ian Goodfellow and Aaron Courville; *Deep Learning*<sup>49</sup>, MIT Press book in preparation. I also contribute feedback on several chapters<sup>50</sup>.

Toby Segaran; *Programming Collective Intelligence: Building Smart Web 2.0 Applications*; O’Reilly, August 2007, first edition.

Howard B. Bandy; *Quantitative Technical Analysis: An Integrated Approach to Trading System Development and Trading Management*; Blue Owl Press, January 2015, first edition.

Steven Bird, Ewan Klein and Edward Loper; *Natural Language Processing with Python*; O’Reilly, June 2009, first edition.

Rishi K. Narang; *Inside the Black Box: A Simple Guide to Quantitative and High-Frequency Trading*, John Wiley & Sons, Inc, 2013, second edition.

<sup>42</sup><https://play.google.com/store/apps/details?id=com.jiandaola.dailycast>

<sup>43</sup>Natural Language Toolkit: <http://www.nltk.org/>

<sup>44</sup>Machine learning library in Python:<http://scikit-learn.org/stable/index.html>

<sup>45</sup>Python data analysis library: <http://pandas.pydata.org/>

<sup>46</sup>Extreme Gradient Boosting: <https://xgboost.readthedocs.org/en/latest/>

<sup>47</sup>Google Deep Learning Framework: <http://www.tensorflow.org/>

<sup>48</sup>Lightning-fast cluster computing: <http://spark.apache.org/>

<sup>49</sup><http://www.iro.umontreal.ca/~bengioy/dlbook/>

<sup>50</sup><http://goodfeli.github.io/dlbook/contents/acknowledgements.html>

Kevin P. Murphy; *Machine Learning: A Probabilistic Perspective*<sup>51</sup>, MIT Press, 2012, first edition.  
Christopher M. Bishop; *Pattern Recognition and Machine Learning*<sup>52</sup>, Springer, 2007, first edition.  
Yoshua Bengio; *Learning Deep Architectures for AI*<sup>53</sup>, Foundations and Trends in Machine Learning, Vol.2 No.1 (2009) 1-127.  
Wes McKinney; *Python for Data Analysis*; O'Reilly, October 2012, first edition.  
Gavin Hackeling; *Mastering Machine Learning with Scikit-learn*<sup>54</sup>; Packt Publishing, October 2014, first edition.  
Femi Anthony; *Mastering pandas*<sup>55</sup>; Packt Publishing, June 2015, first edition.  
Luis P. Coelho, Willi Richert; *Building Machine Learning Systems with Python*<sup>56</sup>; Packt Publishing, March 2015, second edition.  
Nick Pentreath; *Machine Learning with Spark*<sup>57</sup>; Packt Publishing, February 2015, first edition.

## HOBBIES

Swimming, playing table tennis and the piano.

---

<sup>51</sup><http://www.amazon.com/gp/product/B00AF1AYTQ/>

<sup>52</sup><http://www.amazon.com/gp/product/0387310738/>

<sup>53</sup>[http://www.iro.umontreal.ca/~bengioy/papers/ftml\\_book.pdf](http://www.iro.umontreal.ca/~bengioy/papers/ftml_book.pdf)

<sup>54</sup><https://www.packtpub.com/packtlib/book/Big-Data-and-Business-Intelligence/9781783988365>

<sup>55</sup><https://www.packtpub.com/packtlib/book/Big-Data-and-Business-Intelligence/9781783981960>

<sup>56</sup><https://www.packtpub.com/packtlib/book/Big-Data-and-Business-Intelligence/9781784392772>

<sup>57</sup><https://www.packtpub.com/packtlib/book/Big-Data-and-Business-Intelligence/9781783288519>