

Copyright is owned by the Author of the thesis. Permission is given for a copy to be downloaded by an individual for the purpose of research and private study only. The thesis may not be reproduced elsewhere without the permission of the Author.

# English-Persian Phrase-based Statistical Machine Translation: Enhanced Models, Search and Training

---

A THESIS PRESENTED IN FULFILMENT OF THE REQUIREMENTS  
FOR THE DEGREE OF

DOCTOR OF PHILOSOPHY  
IN  
COMPUTER ENGINEERING

at Massey University  
Albany (Auckland), New Zealand

Mahsa Mohaghegh

2012

Copyright ©2012 by Mahsa Mohaghegh. Some Rights Reserved.  
You are free to copy, distribute and transmit the work as well as adapt the work,  
provided it is used for non-commercial purposes and it is cited clearly and correctly.

---

## ABSTRACT

Machine translation (MT), as applied to natural language processing, has undergone substantial development over the past sixty years. While there are a number of different approaches to MT, there has been increasing interest in statistical machine translation (SMT) as the preferred approach to MT. Advances in computational power, together with the exploration of new methods and algorithms have enabled a general improvement in the output quality in a number of systems for various language pairs using this approach. However, there is a significant lack of research work in the area of English/Persian SMT, mainly due to the scarcity of data for this language pair, and the shortage of fundamental resources such as large-scale bilingual corpora. Several research studies have been published on work in the area of machine translation involving the Persian language; however, results producing fluent, usable output are rare.

This thesis shows how SMT was implemented with this language pair for the first time, and how we created a cutting-edge hybrid SMT system capable of delivering high-quality translation output.

We present the development of what is currently the largest English/Persian parallel corpus, constructed using a web crawler to source usable online data, together with the concatenation of existing parallel corpora. As yet another contribution of the research, we propose an improved hybrid corpus alignment method involving sentence length-based and word correspondence-based models to align words, phrases and sentences in the corpus. We also show the impact that the corpus domain can have on the translation output, and the necessity to consider domains of both bilingual and monolingual corpora where they are included in the training and language models.

Two open-source toolkits, Moses and Joshua, were modified to work with the Persian language, and their behaviour and performance results were compared to determine which performed better when implemented with the Persian language.

We present our work in designing, testing, and implementing a novel, three-level Transfer-based automatic post-editing (APE) component based on grammatical rules, which operates by analysing, parsing, and POS-tagging the output, and implements functions as transformers which perform corrections to the text, from lexical

---

transformation to complex syntactical rearrangement. We show that rule-based approaches to the task of post-editing are superior to the commonly-used statistical models, since they incorporate linguistic knowledge, and are strong in terms of syntax, morphology, and structural semantics – qualities which are very desirable when performing grammatical correction and syntactical restructuring.

We implement independent manual evaluation as well as standard automatic techniques, in order to assess more accurately the translation output. This evaluation shows that the use of the APE component is able to improve translation output significantly, that is, by at least 25%, resulting in high-quality translation output.

Our system performs well by using a combination of the capabilities of two main MT approaches – SMT and RBMT – in different areas of the system as a whole. SMT provides the main system with consistent, mathematical-based translation, and the Transfer-based algorithm in the APE component operates with comprehensive linguistic rules in order to improve incorrect sentences, and fine-tune translation output. This results in a robust, state-of-the-art system, which noticeably exceeds other currently available solutions for this language pair.

---

## ACKNOWLEDGEMENTS

*In the name of God who owns soul and wisdom. These are the best attributes of God.*

*~ Ferdowsi (935 – 1020)*

Firstly, I would like to express my gratitude to my advisors Professor Dr.Hossein Sarrafzadeh, and Dr. Rezaul Hasan. This thesis would not have been possible without their support, advice, and valuable ideas and suggestions.

I would also like to thank my committee member Dr. Tom Moir for his constant support and valuable advice and critical comments. His advice and patience was much appreciated.

I wish to extend my heartfelt thanks to my parents, whose continuous encouragement lightens my path into higher education.

At this point, I would like to express my everlasting gratitude to my best friend for his patience, support, and love. This thesis would not have been possible without Mike's continuous encouragement and patience. Thank you for reminding me that there are things more important than this work, and for standing next to me during this time.

I also want to thank my best friend Dr. Mandana Arzaghi, who pointed out early that there are life and opportunities also outside of Iran. During these last years, geographical distances were against seeing more each of each other, but you were always there for me as a true friend and a good example of perseverance and success. I know you will always be there as a good influence for me.

I would also like to thank Mehdi Mohammadi who has been working at the University of Shikh Bahae for the many fruitful discussions and joint experiments related to Automatic Post-Editing and Hybrid translation approaches.

My PhD studies turned out to be an unforgettable experience, mostly thanks to the support from my colleagues and friends. I was lucky to be part of a great group of scientifically ambitious, intelligent, and industrious researchers at Massey e-Centre.

Lastly, I offer my regards and appreciation to all of those who have supported me in many aspects during the course of the research, my friends, family, colleagues, reviewers, etc, and anyone else not mentioned here.

---

Dedicated to my parents

---

## DECLARATION

The author declares that this is her own work except where due acknowledgement has been given. It is being submitted for a PhD in Engineering to Massey University, New Zealand.

This Thesis describes the research carried out by the author at the School of Engineering, Massey University, New Zealand, from June 2008 to December 2012, supervised by Dr Rezaul Hasan.

## TABLE OF CONTENTS

ABSTRACT.....	i
ACKNOWLEDGEMENTS.....	iii
DECLARATION.....	v
LIST OF FIGURES.....	xi
LIST OF TABLES.....	xii
LIST OF ABBREVIATIONS.....	xv
LIST OF PUBLICATIONS.....	xviii
Chapter 1. Introduction.....	1
1.1 Problem Statement.....	1
1.2 Scope of the Study.....	2
1.3 Research Challenges.....	3
1.4 Contributions to Knowledge.....	4
1.5 Thesis Outline.....	5
Chapter 2. Literature Review.....	7
2.1 Introduction.....	7
2.2 Machine Translation Systems.....	7
2.2.1 Machine Translation Difficulties.....	8
2.2.2 Examples of Use.....	10
2.2.3 Machine Translation Advantages and Disadvantages.....	10
2.3 Machine Translation Approaches.....	11
2.3.1 Statistical Machine Translation.....	13
2.3.2 Advantages of the Statistical Approach for Machine Translation.....	15
2.4 Related Work in Statistical Machine Translation.....	16
2.5 Related Work in Persian MT.....	17
2.6 Existing Machine Translation Tools and Services.....	26
2.7 Online vs. Installable Software.....	29
2.8 Persian Language.....	30
2.9 Characteristics of the Persian Language.....	33
2.10 Persian Alphabet and Pronunciation.....	34
2.11 Persian Corpora.....	35
2.12 Available Persian Text Corpora.....	35
2.12.1 Bijankhan corpus.....	35
2.12.2 Hamshahri Corpus.....	36
2.12.3 Shiraz Corpus.....	37
2.12.4 MULTTEXT-East Framework.....	37



2.12.5	TEP – Tehran English-Persian Corpus (Parallel)	37
2.12.6	PEN: Parallel English-Persian News Corpus	38
2.12.7	TMC – Tehran Monolingual Corpus	38
2.12.8	ELRA	38
2.12.9	Other Corpora	38
2.13	Summary	39
Chapter 3. Statistical Machine Translation and Evaluation Metrics		41
3.1	Introduction	41
3.2	SMT Overview	41
3.3	Bayes Decision Rule	41
3.3.1	Noisy-Channel Model	41
3.3.2	Log-linear Model	44
3.4	Translation Model	45
3.5	Training Model	47
3.6	Parallel Corpus Alignment	47
3.6.1	Word Alignment	47
3.6.2	Phrase Alignment	48
3.7	Language Model	50
3.7.1	Uni-gram Model	50
3.7.2	Bi-gram Model	50
3.7.3	N-grams	51
3.8	Translation and Evaluation for Training Purposes	52
3.9	Decoding Process	52
3.10	Evaluation Metrics	53
3.10.1	BLEU	54
3.10.2	NIST	55
3.10.3	Meteor	56
3.10.4	TER	56
3.11	Open-source Decoding Software	56
3.12	Summary	57
Chapter 4. Initial Tests and Corpus Development		59
4.1	Introduction	59
4.2	Initial Set-up and Testing	59
4.3	Discussion and Analysis of Initial Results	61
4.4	Corpus Development	65
4.5	Alignment	67
4.6	Experiments and Results	68

4.6.1	Overview of earlier English-Persian experiments .....	68
4.6.2	Further experiments in the English-Persian Translation Direction.....	68
4.6.3	Experiments in the Persian-English Translation Direction .....	73
4.7	Summary .....	75
Chapter 5.	Hierarchical Phrase-Based Translation Model .....	76
5.1	Introduction .....	76
5.2	Hierarchical Phrase-Based Overview.....	76
5.3	Thrax .....	78
5.4	Moses vs. Joshua.....	79
5.4.1	Syntax Models .....	79
5.4.2	String-to-tree Models.....	79
5.4.3	Text Rule Table Format .....	79
5.5	Data Preparation.....	81
5.6	Experiment Results and Evaluation .....	82
5.6.1	System configuration .....	82
5.6.2	Results.....	82
5.7	Joshua 4.0.....	87
5.7.1	Experiments and Results.....	88
5.8	Multiple References .....	91
5.9	Summary .....	92
Chapter 6.	APE - Automatic Post-Editing System: Background .....	93
6.1	Introduction .....	93
6.2	Motivation for an APE Approach .....	93
6.3	Related Work.....	95
6.4	Other Hybrid Approaches .....	104
6.5	Proposed APE Approach.....	105
6.6	Description of our APE Approach .....	107
6.7	Persian Dependency Treebank .....	109
6.8	Corpus Study for POS-Tagging Experiments .....	109
6.8.1	Related Work .....	109
6.8.2	POS-tagging for Persian language – difficulties.....	111
6.9	Parsing Approaches.....	112
6.9.1	Link Grammar Parser.....	112
6.9.2	Data-Driven Dependency Parsing.....	113
6.9.3	MaltParser .....	114
6.10	Initial Steps for an RBMT-APE Approach.....	114
6.10.1	MLETagger.....	114

6.10.2	Tagger class .....	114
6.10.3	MLETagger class .....	115
6.10.4	CoNLL class .....	116
6.11	POS-Tagger .....	116
6.12	Summary .....	120
Chapter 7.	APE Method Development, Experiments and Results .....	121
7.1	Introduction .....	121
7.2	Program class .....	121
7.2.1	ParserDataLine class .....	122
7.2.2	DataPreparation class .....	122
7.3	MSTParser Details .....	123
7.3.1	Class Parser .....	123
7.3.2	Training Parser .....	123
7.3.3	Parsing inputs .....	124
7.4	Transformers .....	126
7.4.1	OOV Remover class .....	127
7.4.2	Dictionary class .....	127
7.4.3	TransferEngineclass .....	128
7.4.4	NumberPreserverclass .....	128
7.4.5	IncompleteDependentTransformerclass .....	128
7.4.6	IncompleteEndedPREMTransformerclass .....	130
7.4.7	AdjectiveArrangementTransformerclass .....	131
7.4.8	NoSubjectSentenceTransformerclass .....	132
7.4.9	PluralNounsTransformer class .....	133
7.4.10	VerbArrangementTransformerclass .....	133
7.4.11	Transliteratorclass .....	135
7.4.12	ConjoinedTokenTransfer class .....	136
7.4.13	Syntactic Valency Lexicon .....	137
7.4.14	VerbValency class .....	139
7.4.15	MissingVerb Transformer class .....	139
7.4.16	MozafOfAlefEndedTokenTransformer class .....	141
7.5	Experiments and Results .....	143
7.5.1	Baseline SMT .....	143
7.5.2	Automatic Evaluation .....	143
7.5.3	Manual Evaluation .....	145
7.6	Summary .....	146
Chapter 8.	Discussion and Conclusions .....	147

8.1	Research Contributions .....	149
8.2	Directions for Future Work.....	150
	References.....	152
	Appendix I: .....	162
	Persian Alphabet .....	162
	Persian Numerals .....	163
	Appendix II: .....	164
	Language Model Example .....	164
	Test Set, Output, Reference, Score Example:.....	166
	Appendix III:.....	168
	APE Diagram 1 .....	168
	APE Diagram 2 .....	169
	Example of MLETagger on Output and Reference Set:.....	170

## LIST OF FIGURES

Figure 2-1: Indo-European Languages .....	30
Figure 2-2: The Top-Affluence Languages of the World.....	31
Figure 2-3: Hamshahri Corpus Version 1 sample.....	36
Figure 3-1: Example of Persian-English alignment (1) .....	46
Figure 3-2: Example of Persian-English alignment (2) .....	46
Figure 3-3: Example of Persian-English alignment (3) .....	48
Figure 3-4: English – Persian Bi-text grid .....	49
Figure 4-1: BLEU scores for various tests.....	61
Figure 4-2: BLEU scores vs. language model sentences for each system configuration .....	63
Figure 4-3: NIST scores vs. language model sentences for each system configuration .....	64
Figure 4-4: Domain percentages for NSPEC corpus .....	67
Figure 5-1: BLEU Scores Pe-En Joshua vs. Moses.....	83
Figure 5-2: NIST scores Pe-En Joshua vs. Moses .....	84
Figure 5-3: BLEU scores English-Persian.....	85
Figure 5-4: NIST scores English-Persian .....	86
Figure 5-5: NPEC Corpus composition .....	88
Figure 6-1: Syntactic Selection.....	104
Figure 6-2: Stochastic Selection .....	105
Figure 6-3: SMT-fed RBMT.....	105
Figure 6-4: Hybrid Architecture .....	105
Figure 6-5: High-level Diagram of the Rule-based APE.....	106
Figure 6-6: Output text parsed with MSTParser.....	108
Figure 6-7: Reference text parsed with MSTParser.....	108
Figure 6-8: Dependency parsing example .....	113
Figure 6-9: POS-Tagging Approaches.....	118
Figure 7-1: BLEU score before and after APE.....	144
Figure 7-2: NIST score before and after APE .....	144
Figure 7-3: Manual evaluation comparison .....	146

## LIST OF TABLES

Table 3-1: English-Persian Probability Example.....	46
Table 3-2: Phrase alignment examples .....	49
Table 4-1: Training model and Persian language model sizes .....	60
Table 4-2: BLEU scores for test with different sized models.....	60
Table 4-3: Results for 817-sentence training model.....	62
Table 4-4: Results for 1011-sentence training model.....	62
Table 4-5: Results for 2343-sentence training model.....	63
Table 4-6: Bilingual corpora used in the training model .....	69
Table 4-7: Bilingual corpora after hybrid alignment method .....	69
Table 4-8: Monolingual corpora used to train the language model .....	70
Table 4-9: Evaluation metric scores with Hamshahri-based language model.....	71
Table 4-10: Evaluation metric scores with BBC News-based language model .....	71
Table 4-11: Evaluation metric scores with IRNA-based language model.....	71
Table 4-12: Evaluation metric score comparison between Google Translate and System 5 with IRNA-based language model.....	72
Table 4-13: Monolingual corpora used to train the language model.....	73
Table 4-14: Evaluation metric scores with News Commentary-based language model .....	73
Table 4-15: Evaluation metric scores with Europarl v4-based language model .....	74
Table 4-16: Evaluation metric score comparison between Google Translate and System 5 with News Commentary-based language model.....	75
Table 5-1: Monolingual corpora composition .....	81
Table 5-2: Parallel corpora composition.....	83
Table 5-3: BLEU scores Pe-En Joshua vs. Moses.....	83
Table 5-4: NIST scores Pe-En Joshua vs. Moses .....	84
Table 5-5: BLEU scores English-Persian .....	85
Table 5-6: NIST scores En-Pe Joshua vs. Moses .....	85
Table 5-7: Baseline System Components .....	89
Table 5-8: Statistics of eight test sets used in automatic and manual evaluation .....	89
Table 5-9: Difference of BLEU and NIST Score after using Joshua 4.0 on eight test sets.....	90
Table 5-10: Multi-BLEU Joshua 4.0 on eight test sets .....	90

Table 5-11: Multiple-reference BLEU/NIST scores for Joshua 1.3-based system output .....	91
Table 5-12: Multiple-reference Multi-BLEU scores for Joshua 1.3-based system output .....	91
Table 5-13: Multiple-reference BLEU/NIST scores for Joshua 4.0-based system output .....	91
Table 5-14: Multiple-reference Multi-BLEU scores for Joshua 4.0-based system output .....	92
Table 6-1: Tag Names.....	110
Table 6-2: Examples of pos-tagging Persian output.....	117
Table 7-1: DataPreparation class .....	122
Table 7-2: Parsing Inputs .....	125
Table 7-3: POS-Tagger: Parts of speech categorised .....	125
Table 7-4: IncompleteDependentTransformerclass – Before .....	129
Table 7-5: IncompleteDependentTransformerclass – After .....	129
Table 7-6: IncompleteEndedPREMTransformerclass- Before .....	130
Table 7-7: IncompleteEndedPREMTransformerclass- After .....	130
Table 7-8: AdjectiveArrangementTransformerclass- Before .....	131
Table 7-9: AdjectiveArrangementTransformerclass- After.....	132
Table 7-10: NoSubjectSentenceTransformer class - Before.....	132
Table 7-11: NoSubjectSentenceTransformer class - After .....	132
Table 7-12: PluralNounsTransformer class - Before .....	133
Table 7-13: PluralNounsTransformer class - After .....	133
Table 7-14: VerbArrangementTransformer class -Before.....	134
Table 7-15: VerbArrangementTransformer class -After .....	134
Table 7-16: En-Fa Transliteration (1).....	136
Table 7-17: En-Fa Transliteration (2).....	136
Table 7-18: ConjoinedTokenTransfer class - Before.....	137
Table 7-19: ConjoinedTokenTransfer class - After .....	137
Table 7-20: Syntactic Valency Lexicon.....	139
Table 7-21: MissingVerb Transformer class- Before .....	140
Table 7-22: MissingVerb Transformer class - After.....	141
Table 7-23: MozafOfAlefEndedTokenTransformer class - Before.....	142
Table 7-24: MozafOfAlefEndedTokenTransformer class - After .....	142

Table 7-25: Scores of APE based on SMT Joshua version 4.0 .....	143
Table 7-26: Scores of two human evaluators for 153 test sentences .....	145
Table 7-27: Mutual score for both human evaluator I and evaluator II.....	145



## LIST OF ABBREVIATIONS

ACL	Complement Clause of Adjective
ADV	Adverb
ADVC	Adverbial Complement of Verb
AJCONJ	Conjunction of Adjective
AJPP	Prepositional Complement of Adjective
AJUCL	Adjunct Clause
AOL	America Online
APE	Automatic Post Editing
APOSTMOD	Adjective Post-Modifier
APP	Apposition
APREMOD	Adjective Pre-Modifier
ASR	Automatic Speech Recognition
AVCONJ	Conjunction of Adverb
BLEU	Bilingual Evaluation Understudy
CAT	Computer Assisted Translation
CFGs	Synchronous Context-Free Grammars
CNW	Canada Newswire
COMPPP	Comparative Preposition
C-STAR	International Consortium for Research on Speech Translation
DARPA	Defence Advance Research Project Agency
DG	Dependency Grammar
EBMT	Example-Based Machine Translation
EGIU	English Grammar in Use
ELRA	European Language Resource Association
EM	Expectation–Maximization Algorithm
En	English
ENC	Enclitic Non-Verbal Element
Fa	Farsi
FAMT	Fully Automatic Machine Translation
FLDB	Farsi Linguistic Database
FST	Finite State Transducer
FTD	US Air Force’s Foreign Technology Division
GLP	Gross Language Product
HAMT	Human-Assisted Machine Translation
Hiero	Hierarchical
HMT	Hybrid Machine Translation
IBM	International Business Machines Corporation
IEEE	Institute of Electrical and Electronics Engineers, Inc.
IR	Information Retrieval
IRNA	Iranian News Agency
IWSLT	International Workshop on Spoken Language Translation
LG	Link Grammar
LV	Linking Verb
LVP	Light Verb Particle
MAHT	Machine-Assisted Human Translation
MAP	Maximum A-Posteriori
MERT	Minimum Error Rate Training
MESU	Measure
METEOR	Metric for Evaluation of Translation with Explicit Ordering
MLE	Maximum Likelihood Estimation
MOS	Mosnad

MOZ	Ezafe Dependent
MPEC	Modern Persian-English corpus
MST	Parser Maximum Spanning Tree Parser
MT	Machine Translation
NADV	Adverb of Noun
NCL	Clause of Noun
NCONJ	Conjunction of Noun
NE	Non-Verbal Element of Infinitive
NEZ	Ezafe Complement of Adjective
NIST	National Institute of Standards and Technology
NLP	Natural Language Processing
NPEC	News Persian English Corpus
NPOSTMOD	Post-Modifier of Noun
NPP	Preposition Of Noun
NPREMOD	Pre-Modifier of Noun
NPRT	Particle of Infinitive
NSPEC	News Subtitle Persian-English Corpus
NVE	Non-Verbal Element
ODJ	Object
ODJ2	Second Object
OOV	Out of Vocabulary
PAHO	Pan American Health Organization
PARCL	Participle Clause
PART	Interrogative Particle
PB	Phrase Based
PCONJ	Conjunction of Preposition
PCTS	Parallel Corpus Test Set
PeEn-SMT	Persian-English Statistical Machine Translation
PEN	Parallel English-Persian News Corpus
POS	Part of Speech
POSDEP	Post-Dependent
PPL	Perplexity Threshold
PREDEP	Pre-Dependent
PREM	Pre-Modifier
PRO	Predicate
PROG	Progressive Auxiliary
PUNC	Punctuation Mark
RBMT	Rule-Based Machine Translation
RHS	Right Hand Side
ROOT	Root
SAMT	Syntax Augmented Machine Translation
SBJ	Subject
SCFG	Stochastic Context-Free Grammar
SDL	Scalable Enterprise Translation Server
SDL	SDL Language Weaver
SMT	Statistical Machine Translation
SOV	Subject-Object-Verb
SRILM	Sri Language Model
SVO	Subject-Verb-Object
TAM	Tamiz
TEP	Tehran English-Persian Corpus
TER	Translation Error Rate Te
TER	Translation Error Rate
TMC	Tehran Monolingual Corpus
U	Unicode

UN	United Nations
UTIRE	University of Tehran Information Retrieval Evaluation System
VCL	Complement Clause of Verb
VCONJ	Conjunction of Verb
VPP	Prepositional Complement of Verb
VPRT	Verb Particle
WER	Word Error Rate
WSD	Word Sense Disambiguation

## LIST OF PUBLICATIONS

1. **Mahsa Mohaghegh**, Abdolhossein Sarrafzadeh, Mehdi Mohammadi, *GRAFIX: Automated Rule-Based Post Editing System to Improve English-Persian SMT* (Short paper- COLING 2012, the 24th International Conference on Computational Linguistics. Mumbai, India, December 2012)  
<http://aclweb.org/anthology-new/C/C12/C12-2085.pdf>
2. **Mahsa Mohaghegh**, Abdolhossein Sarrafzadeh, *A Hierarchical Phrase-Based Model for English-Persian Statistical Machine Translation* (Full Paper - Innovations 12, 8th International Conference on Innovations in Information Technology. AL AIN ,UAE, March 2012)  
[http://ieeexplore.ieee.org/xpls/abs\\_all.jsp?arnumber=6207733](http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=6207733)
3. **Mahsa Mohaghegh**, *Advancements in English-Persian Hierarchical Statistical Machine Translation*. (Short paper - NZCSRSC New Zealand Computer Science Research Student Conference. Otago, April 2012)  
<https://sites.google.com/a/nzcsrsc.ac.nz/nzcsrsc2012/programme/posters>
4. **Mahsa Mohaghegh**, Abdolhossein Sarrafzadeh, Tom Moir, *Improving Persian-English Statistical Machine Translation: Experiments in Domain Adaptation* (Full paper – *In Proceedings of the 2nd Workshop on South and Southeast Asian Natural Language Processing (WSSANLP)*, IJCNLP 2011, pages 9–15,Chiang Mai, Thailand, November 2011)  
<http://www.aclweb.org/anthology/W/W11/W11-3002.pdf>
5. **Mahsa Mohaghegh**, Abdolhossein Sarrafzadeh, *An Overview of the Challenges and Progress in PeEn-SMT: First Large Scale Persian-English SMT System* (Full Paper - Innovations 11, 7th International Conference on Innovations in Information Technology. Abu Dhabi, April 2011)  
[http://ieeexplore.ieee.org/xpls/abs\\_all.jsp?arnumber=5893841](http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=5893841)
6. **Mahsa Mohaghegh**, Abdolhossein Sarrafzadeh, *The Impact of Domain for Language Model in the PeEn-SMT: First Large Scale Persian-English SMT System*. (Short paper-NZCSRSC New Zealand Computer Science Research Student Conference. Palmerston North, April 2011)  
<https://sites.google.com/a/maori.geek.nz/nzcsrsc2011/papers/paper-abstracts#TOC-Mahsa-Mohaghegh-The-Impact-of-Domain-for-Language-Model-in-the-PeEn-SMT:-First-Large-Scale-Persian-English-SMT-System->
7. **Mahsa Mohaghegh**, Abdolhossein Sarrafzadeh, *Multilingual Information Service System for Tourists* (Poster – *NZBio Conference*, Auckland, March 2011)
8. **Mahsa Mohaghegh**, Abdolhossein Sarrafzadeh, Tom Moir. *Improved Language Modelling for English-Persian Statistical Machine Translation*. (Full Paper - *In Proceedings of SSST-4 Workshop at COLING-2010*, Beijing, China, August 2010)

<http://www.aclweb.org/anthology-new/W/W10/W10-3810.pdf>

9. **Mahsa Mohaghegh**, Abdolhossein Sarrafzadeh. *Performance Evaluation of Statistical English-Persian Machine Translation*. (Full Paper – JADT2010 – 10<sup>th</sup> International Conference on Statistical Analysis of Textual Data. Sapienza, University of Rome, June 2010)

[http://www.ledonline.it/ledonline/JADT-2010/allegati/JADT-2010-1091-1100\\_114-Sarrafzadeh.pdf](http://www.ledonline.it/ledonline/JADT-2010/allegati/JADT-2010-1091-1100_114-Sarrafzadeh.pdf)

10. **Mahsa Mohaghegh**. *A Statistical Approach to English-Persian Machine Translation*. (Short Paper – NZCSRSC – New Zealand Computer Science Research Student Conference. Wellington, April 2010)

<http://ecs.victoria.ac.nz/Events/NZCSRSC2010/Papers>

11. **Mahsa Mohaghegh**, Tom Moir, Abdolhossein Sarrafzadeh. *A Statistical Approach to English-Persian Machine Translation*. (Poster – NZBio Conference, Auckland, March 2010)

[http://www.academia.edu/238736/A\\_Statistical\\_Approach\\_to\\_English-Persian\\_Machine\\_Translation](http://www.academia.edu/238736/A_Statistical_Approach_to_English-Persian_Machine_Translation)

12. **Mahsa Mohaghegh**, Abdolhossein Sarrafzadeh. *Analysis of the Effect of Data Variation in a Statistical English-Persian Machine Translation*. (Full Paper - Innovations '09 – 6<sup>th</sup> International Conference on Innovations in Information Technology. Dubai, December 2009.

<http://dl.acm.org/citation.cfm?id=1802285>

<http://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=05413782>

13. **Mahsa Mohaghegh**, Tom Moir. *First English-Persian Machine Translation*. (Poster – NZPGC – New Zealand Postgraduate Conference. Wellington, November 2009)

[http://www.academia.edu/238733/the\\_first\\_english-persian\\_statistical\\_machine\\_translation](http://www.academia.edu/238733/the_first_english-persian_statistical_machine_translation)