



DUBLIN CITY UNIVERSITY

SEMESTER 2 EXAMINATIONS 2016/2017

MODULE:	CA4012 – Statistical Machine Translation										
PROGRAMME(S):	<table><tr><td>CASE</td><td>BSc in Computer Applications (Sft.Eng.)</td></tr><tr><td>ECSA</td><td>Study Abroad (Engineering & Computing)</td></tr><tr><td>ECSAO</td><td>Study Abroad (Engineering & Computing)</td></tr><tr><td>CPSSD</td><td>BSc in Computational Problem Solving & Software Development</td></tr></table>			CASE	BSc in Computer Applications (Sft.Eng.)	ECSA	Study Abroad (Engineering & Computing)	ECSAO	Study Abroad (Engineering & Computing)	CPSSD	BSc in Computational Problem Solving & Software Development
CASE	BSc in Computer Applications (Sft.Eng.)										
ECSA	Study Abroad (Engineering & Computing)										
ECSAO	Study Abroad (Engineering & Computing)										
CPSSD	BSc in Computational Problem Solving & Software Development										
YEAR OF STUDY:	4,O,X										
EXAMINER(S):	<table><tr><td>Prof. Andy Way</td><td>(Ext: 5074)</td></tr><tr><td>Dr. Jinhua Du</td><td>(Ext: 6716)</td></tr><tr><td>Dr. Haithem Afli</td><td>(Ext: 8712)</td></tr><tr><td>Prof. David Bustard</td><td></td></tr><tr><td>Dr. Ian Pitt</td><td></td></tr></table>	Prof. Andy Way	(Ext: 5074)	Dr. Jinhua Du	(Ext: 6716)	Dr. Haithem Afli	(Ext: 8712)	Prof. David Bustard		Dr. Ian Pitt	
Prof. Andy Way	(Ext: 5074)										
Dr. Jinhua Du	(Ext: 6716)										
Dr. Haithem Afli	(Ext: 8712)										
Prof. David Bustard											
Dr. Ian Pitt											
TIME ALLOWED:	3 Hours										
INSTRUCTIONS:	Answer five questions. You must attempt <i>at least one</i> question from <i>each</i> of Sections A, B and C. All questions carry equal marks.										

PLEASE DO NOT TURN OVER THIS PAGE UNTIL YOU ARE INSTRUCTED TO DO SO.

The use of programmable or text storing calculators is expressly forbidden.

Please note that where a candidate answers more than the required number of questions, the examiner will mark all questions attempted and then select the highest scoring ones.

There are no additional requirements for this paper.

SECTION A

QUESTION 1

[TOTAL MARKS: 20]

Q 1(a)

[6 Marks]

Why is it good practice to assemble bilingual training data for statistical MT engines which are (i) as large as possible, (ii) of good quality, and (iii) representative? Explain some of the problems which are likely to ensue if any one of these requirements are not met.

Q 1(b)

[10 Marks]

Malay is an Austronesian language spoken by the Malay people and people of other ethnic groups in the Malay Peninsula, southern Thailand, the Philippines and Singapore. Here are eight Malay sentences and their English translations.

Gadis cantik itu tidak kaya.

The beautiful girl is not rich.

Penyanyi itu tidak bahagia.

The singer is not happy.

Kekayaan itu bukan dari teman bapa.

The wealth is not from his father.

Wang bukan kebahagiaan.

Money is not happiness.

Kareta itu tidak berdating dari medan itu.

The car is not coming from the field

Manusia itu depan rumah itu bukan

The man in front of the house is

penyanyi.

not a singer.

Hadiah itu bukan untuk bapa itu.

The gift is not for the father.

Gadis bahagia itu tidak bermenangis.

The happy girl is not crying.

Translate the following into Malay.

- i. The rich girl is not a singer.
- ii. The man is not coming.
- iii. His wealth is not for the girl.
- iv. Beauty is not a gift.
- v. The gift from the singer is not beautiful.

Q 1(c)

[4 Marks]

Describe in your own words how you produced these translations, focusing in particular on the particular types of inferences you made from the parallel data provided. How is this analogous to how SMT works?

[End of Question 1]

QUESTION 2**[TOTAL MARKS: 20]****Q 2(a)****[4 Marks]**

Machine translation is being used by millions of people on a daily basis. Provide *two* use-cases which freely available services such as Google Translate are good for, and provide reasons behind your selection.

Q 2(b)**[6 Marks]**

What is the market need for MT companies which build customised engines for their clients? Why can those clients expect better translation quality compared to using (say) Google Translate?

Q 2(c)**[6 Marks]**

In a post-editing MT pipeline, what is the role of human translators? How important is their role compared to the MT system? How can their feedback improve the quality of the MT engines?

Q 2(d)**[4 Marks]**

Why do you think some human translators are still reluctant to embrace MT technology?

[End of Question 2]**[END OF SECTION A]**

SECTION B

QUESTION 3

[TOTAL MARKS: 20]

Q 3(a)

[9 Marks]

Assume two SMT systems that translate from English→Arabic (a morphological rich language) and from Arabic→English.

Suggest *three* ways in which their performance may be improved by incorporating linguistic knowledge.

For each idea, (i) state its motivation (i.e. what linguistic problem it is trying to solve), (ii) how you would implement it, and (iii) for which direction (i.e. English-to-Arabic or Arabic-to-English) it is deployed.

NB, There should be *at least one* idea for each of these phases: pre-processing, decoding and post-processing. There should be *at least one* idea for each of these types of linguistic knowledge: morphology and syntax.

Q 3(b)

[5 Marks]

Even though adding linguistic information into an SMT system may be beneficial, provide *two* reasons why it may be less than straightforward to do in practice.

Q 3(c)

[6 Marks]

Describe *one* different pre-processing step that is commonly carried out when translating in the following three cases:

- i. From Arabic, regardless of the target language.
- ii. Between languages that follow different word orders.
- iii. From a language with a rich inflectional system into English.

[End of Question 3]

QUESTION 4**[TOTAL MARKS: 20]****Q 4(a)****[5 Marks]**

Show how a bigram language model would decompose the following sentence to calculate its probability, both with and without sentence boundaries.

“They didn 't evaluate their systems .”

Q 4(b)**[7 Marks]**

Calculate the probability of the last token (.) in the sequence “ the green apple .” using bigram and trigram language models. The following sequences occur in the training data the number of times shown:

“green witch”, twice.

“the green witch”, twice.

“witch”, 5 times.

“ . ”, 20 times.

“apple”, 8 times.

“green apple .”, 6 times.

“ green apple”, 7 times.

“apple .”, 5 times.

“the”, 600 times.

“green”, 10 times.

“the green”, 5 times.

Q 4(c)**[4 Marks]**

How we can deal with unseen n -grams in the context of statistical language modelling? State *one* method which can be applied for this task.

Q 4(d)**[2 Marks]**

How can we evaluate a language model? Why is this useful?

Q 4(e)**[2 Marks]**

Define the term ‘continuous space language model’.

[End of Question 4]

[END OF SECTION B]

SECTION C

QUESTION 5

[TOTAL MARKS: 20]

Q 5(a)

[4 Marks]

Assume the following English—Chinese sentence pairs:

S_1	S_2
<i>car</i> <i>che</i>	<i>my car</i> <i>wode che</i>

The source side is English, and the target side is Chinese. In this question, the *NULL* token is ignored.

Assuming that only one-to-one alignment is allowed, list all possible word alignments for the two sentence pairs.

Q 5(b)

[10 Marks]

Considering all the word alignments you computed in (a), (i) state what the following translation probabilities will be after two iterations of the Expectation Maximisation algorithm, and (ii) show all the steps followed to arrive at these values:

$t(\text{che}|\text{car})$
 $t(\text{wode}|\text{car})$
 $t(\text{che}|\text{my})$
 $t(\text{wode}|\text{my})$

Q 5(c)

[6 Marks]

List all phrase pairs that are consistent with the following word alignment:

	A	B	C	D
W				
X				
Y				
Z				

[End of Question 5]

QUESTION 6**[TOTAL MARKS: 20]****Q 6(a)****[10 Marks]**

Assume the following partial phrase table:

<i>ta</i>	<i>he</i>	0.6
<i>ta</i>	<i>she</i>	0.4

<i>henhui</i>	<i>can</i>	0.4
<i>henhui</i>	<i>is good at</i>	0.6

<i>biancheng</i>	<i>program</i>	0.3
<i>biancheng</i>	<i>programming</i>	0.5
<i>biancheng</i>	<i>programs</i>	0.2

<i>ta henhui</i>	<i>he can</i>	0.2
<i>ta henhui</i>	<i>she can</i>	0.1
<i>ta henhui</i>	<i>he is good at</i>	0.4
<i>ta henhui</i>	<i>she is good at</i>	0.3

<i>henhui biancheng</i>	<i>can program</i>	0.2
<i>henhui biancheng</i>	<i>is good at programming</i>	0.5
<i>henhui biancheng</i>	<i>can programming</i>	0.1
<i>henhui biancheng</i>	<i>programs</i>	0.2

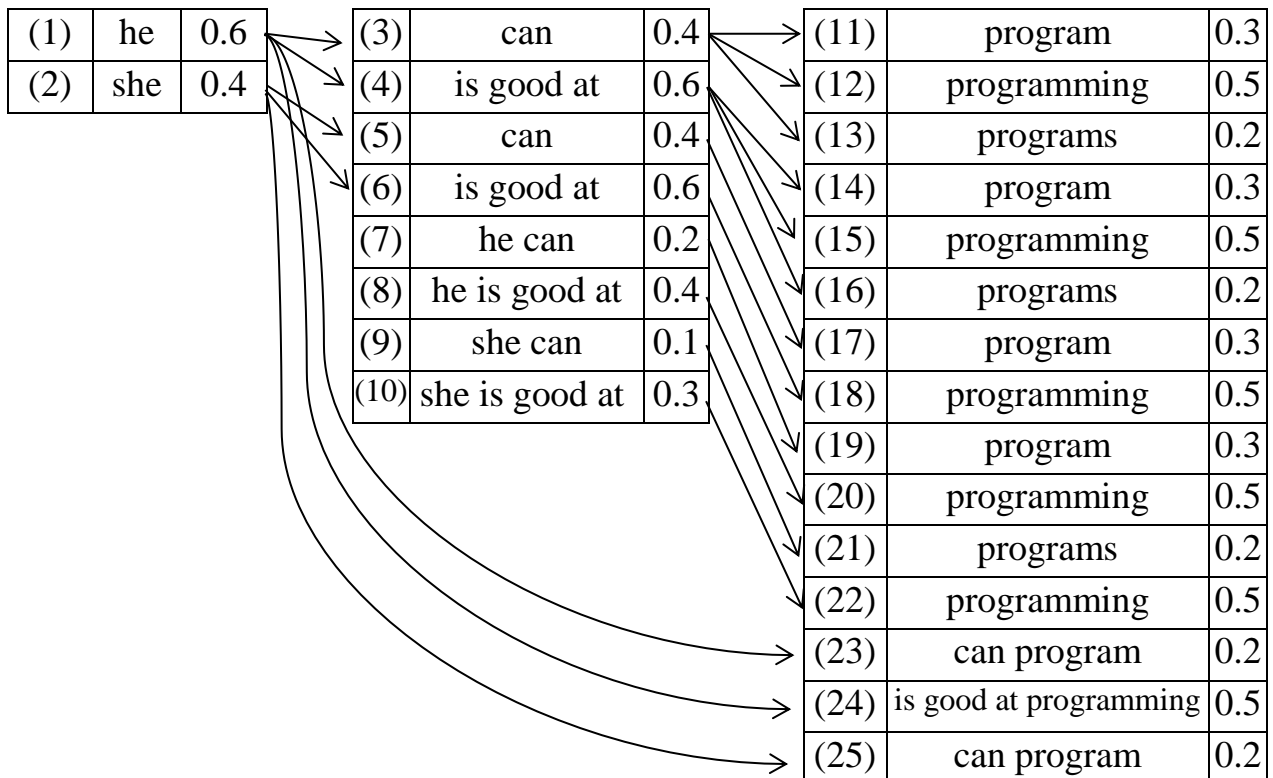
Consider the following input sentence:

ta henhui biancheng

Assuming that:

- only monotone word order is permitted;
- the language model is ignored.

Then we have the following search diagram (partial search space):



Given the search diagram, calculate the probabilities for all possible hypotheses (search paths). In addition, indicate which hypothesis provides the optimal translation for the input sentence.

Q 6(b)

[6 Marks]

Given the search diagram, indicate (i) which group of hypotheses can be recombined, and (ii) which hypothesis should be selected to represent each group.

Q 6(c)

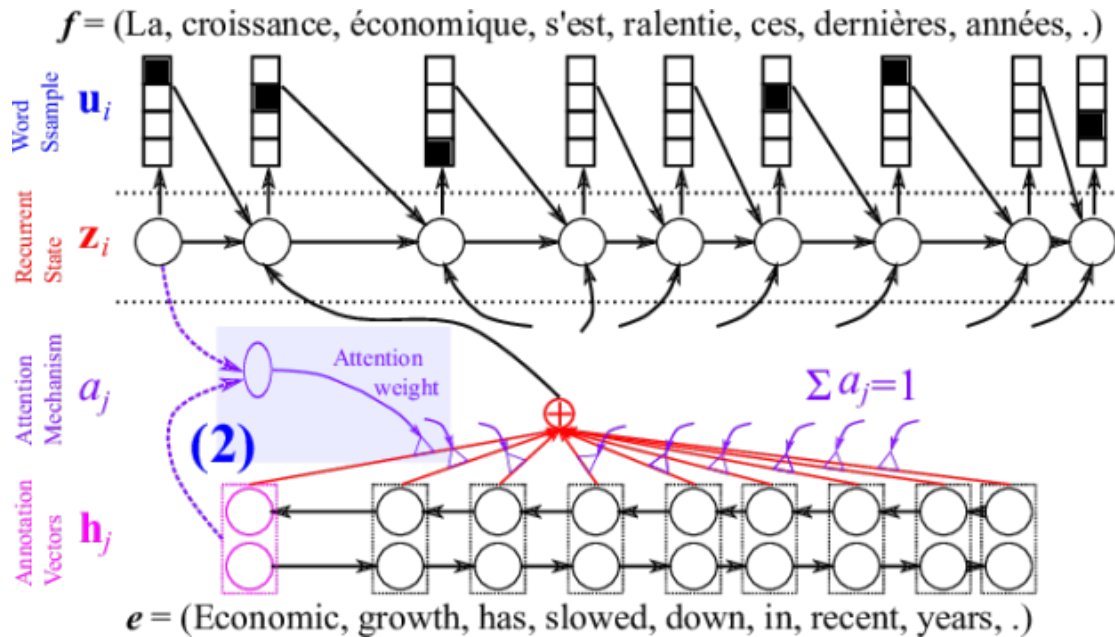
[4 Marks]

Assume histogram pruning after recombination, where the maximum number of hypotheses in each stack is 3. Indicate which hypotheses will be pruned. Explain your answer.

[End of Question 6]

QUESTION 7**[TOTAL MARKS: 20]**

The figure below shows the architecture of the mainstream neural MT model, an encoder-decoder framework with an attention mechanism.



Q 7(a) **[5 Marks]**

Given the figure above, explain how the English sentence e is translated into the French sentence f .

Q 7(b) **[5 Marks]**

Define the role of the “Attention Mechanism”. How are attention weights learned?

Q 7(c) **[5 Marks]**

In NMT, the decoder predicts a target word at each time-step by means of a conditional probability. Write out the equation predicting the i th target word y_i given the input string x , and explain what each variable in the equation indicates.

Q 7(d) **[5 Marks]**

What is the role of the activation function in neural networks? Write out the formula of the commonly used Sigmoid function $f(x)$, draw its curve, and calculate its derivative.

[End of Question 7]

[END OF SECTION C]

[END OF EXAM]