



DUBLIN CITY UNIVERSITY

SEMESTER 2 EXAMINATIONS 2015/2016

MODULE: CA4012 – Statistical Machine Translation

PROGRAMME(S):
CASE BSc in Computer Applications (Sft.Eng.)

YEAR OF STUDY: 4

EXAMINERS: Prof. Andy Way (Ext: 5074)
Dr. Jinhua Du (Ext: 6716)
Dr. Antonio Toral (Ext: 8712)
Dr. Ian Pitt

TIME ALLOWED: 2 Hours

INSTRUCTIONS: Answer any **four** questions.
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.

Requirements for this paper (Please mark (X) as appropriate)

<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>

Log Tables
Graph Paper
Dictionaries
Statistical Tables
Bible

<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>

Thermodynamic Tables
Actuarial Tables
MCQ Only – Do not publish
Attached Answer Sheet
Exam Paper to be returned with Booklet

QUESTION 1

[TOTAL MARKS: 25]

Q 1(a)

[8 Marks]

How might you argue that “The time for MT is now”? Justify your answer by providing *three* use-cases where you claim that MT is the only solution, i.e. that there is no place for human intervention in the translation pipeline for such use-cases.

Q 1(b)

[6 Marks]

Study these sentences in Abma, an Austronesian language spoken in the South Pacific island of Vanuatu, and their English translations:

1. Mwamni sileng. ⇔ He drinks water.
2. Nutsu mwatbo mwamni sileng. ⇔ The child keeps drinking water.
3. Nutsu mwatbo mwegalgal. ⇔ The child keeps crawling.
4. Mwerava Mabontare mwisib. ⇔ He pulls Mabontare down.
5. Mabontare mwisib ⇔ Mabontare goes down.
6. Mweselkani tela mwesak. ⇔ He carries the axe up.
7. Mwelebbe sileng mwabma. ⇔ He brings water.
8. Mabontare mworob mwesak. ⇔ Mabontare runs up.

Assume the following additional lexical entries:

- sesesrakan ⇔ teacher
- mwegani ⇔ eat

Translate the following sentences into Abma:

- i. The teacher carries the water down.
- ii. The child keeps eating.

Q 1(c)

[5 Marks]

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

Q 1(d)

[6 Marks]

Any statistical approach to MT requires the availability of aligned bilingual corpora which are (i) large, (ii) good-quality, and (iii) representative. Explain why all three requirements are important. What are some of the potential problems if any one of these requirements are not met?

[End of Question 1]

QUESTION 2

[TOTAL MARKS: 25]

Q 2(a)

[15 Marks]

For the source-language sentence (A):

(A) *Kuopion kaupunginvaltuusto hyväksyi liitoksen yksimielisesti maanantaina .*

Assume that the outputs (B)(i) and (B)(ii) were produced by an MT system:

(B)(i) The city council unanimously approved the joint Niiralan on Monday .

(B)(ii) The Kuopion liitoksen City council approved unanimously on Monday .

Assume also that the ‘gold standard’ reference translation is (C):

(C) The city council of Kuopio accepted the annexation unanimously on Monday .

Calculate the BLEU scores of the two candidate translations using maximum n -gram length of 3.

Q 2(b)

[5 Marks]

Why is *standard* BLEU unsuitable for sentence-level evaluation of MT quality, especially when sentences are short? How would you modify BLEU to make it suitable to be used at sentence level?

Q 2(c)

[5 Marks]

Explain the concepts of “fluency” and “adequacy” as they apply to the evaluation of MT output. To support your answer, give example translations (in English) which are:

- fluent and adequate,
- fluent but inadequate,
- disfluent but adequate,
- disfluent and inadequate.

[End of Question 2]

QUESTION 3**[TOTAL MARKS: 25]****Q 3(a)****[10 Marks]**

Given the following sentences:

- $\langle s \rangle$ Denis likes Ada $\langle /s \rangle$
- $\langle s \rangle$ Ada likes Richard $\langle /s \rangle$
- $\langle s \rangle$ Ada hates Java $\langle /s \rangle$

List all the parameters of the unigram and bigram language models trained with these sentences without smoothing.

Q 3(b)**[10 Marks]**

Given the language models you built in Q3(a), but now with add-alpha smoothing where $\alpha=0.3$, calculate the probabilities of the following sentences:

- $\langle s \rangle$ Richard likes Ada $\langle /s \rangle$
- $\langle s \rangle$ Richard hates Ada $\langle /s \rangle$

Q 3(c)**[5 Marks]**

Explain the Markov assumption. Why do you need to take it into account when building n -gram language models? How can language models based on neural networks be non-Markovian?

[End of Question 3]

QUESTION 4**[TOTAL MARKS: 25]****Q 4(a)****[6 Marks]**

Assume the following Chinese—English segment-pairs:

S_1	S_2
yuan	hen yuan
far	far away

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

Assuming each target word is exactly aligned with one source word, list all possible word alignments for the two segment-pairs.

Q 4(b)**[10 Marks]**

For all the word alignments you computed above, state what the following translation probabilities will be after two iterations of the Expectation Maximisation algorithm, and show all the interim steps by which you arrived at these values:

- $t(\text{far}|\text{yuan})$
- $t(\text{away}|\text{yuan})$
- $t(\text{far}|\text{hen})$
- $t(\text{away}|\text{hen})$

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

Explain the term “consistency” as used in phrase extraction.

Q 4(d)**[5 Marks]**

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

	<i>A</i>	<i>B</i>	<i>C</i>
<i>x</i>			
<i>y</i>			
<i>z</i>			

[End of Question 4]

QUESTION 5

[TOTAL MARKS: 25]

Assume the following partial phrase table:

wo	I	0.7	xihuan	like	0.4	kaiche	driving	0.3
wo	me	0.3	xihuan	like to	0.5	kaiche	drive	0.5
			xihuan	likes to	0.1	kaiche	drive a car	0.2

wo xihuan	I like	0.3	xihuan kaiche	like driving	0.3
wo xihuan	I like to	0.5	xihuan kaiche	like to drive	0.5
wo xihuan	I likes to	0.1	xihuan kaiche	like to drive a car	0.1
wo xihuan	me likes to	0.1	xihuan kaiche	likes to drive	0.1

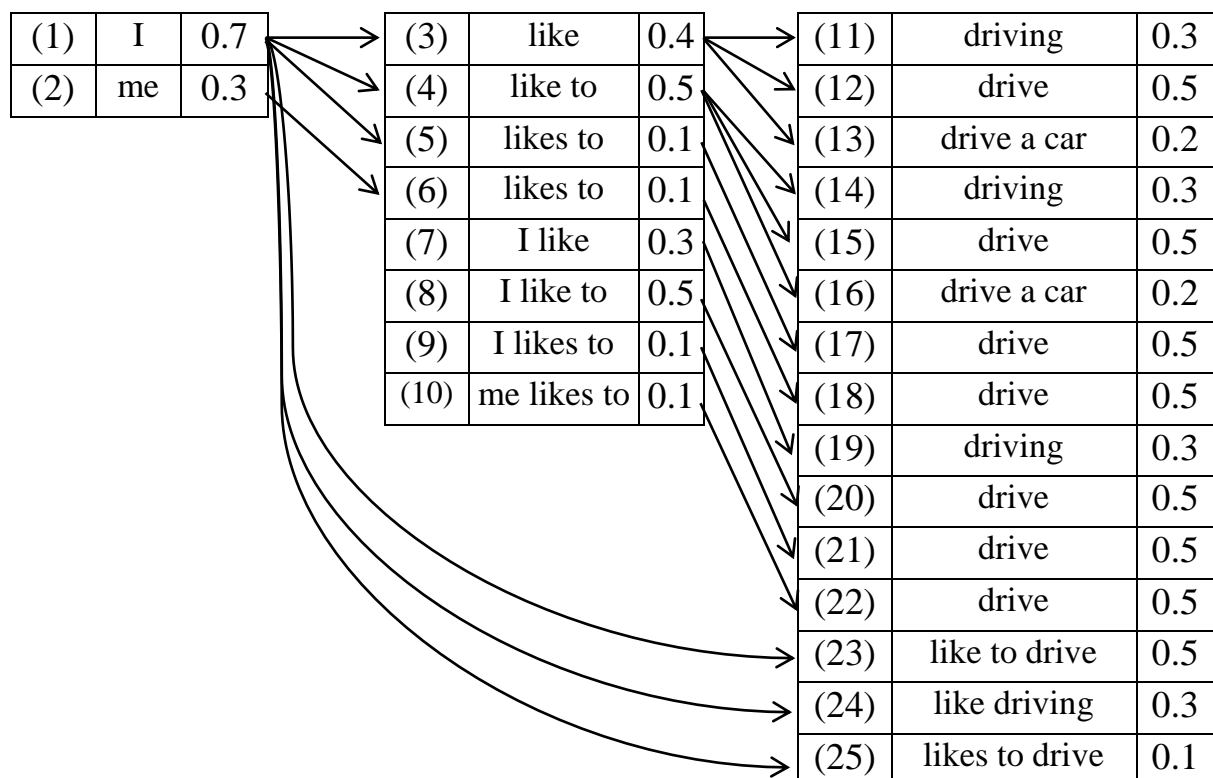
Consider the following input sentence:

wo xihuan kaiche

Assume that:

- Only monotone word order is permitted;
- The language model is ignored.

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



Q 5(a)**[9 Marks]**

Given this search diagram, calculate the probabilities for all possible hypotheses (search paths). Furthermore, indicate which hypothesis provides the most likely translation for the given input sentence.

Q 5(b)**[6 Marks]**

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

Q 5(c)**[6 Marks]**

Assuming histogram pruning after recombination, where the maximum number of hypotheses in each stack is 4, indicate which hypotheses will be pruned.

Q 5(d)**[4 Marks]**

Provide the fundamental equations of (i) the noisy channel model of SMT, and (ii) the log-linear model of SMT. List *three* frequently used features in log-linear models of SMT.

[End of Question 5]

[END OF EXAM]



DUBLIN CITY UNIVERSITY

AUGUST/RESIT EXAMINATIONS 2015/2016

MODULE: CA4012 – Statistical Machine Translation

PROGRAMME(S):
CASE BSc in Computer Applications (Sft.Eng.)

YEAR OF STUDY: 4

EXAMINERS: Prof. Andy Way (Ext: 5074)
Dr. Jinhua Du (Ext: 6716)
Dr. Antonio Toral (Ext: 8712)
Dr. Ian Pitt

TIME ALLOWED: 2 Hours

INSTRUCTIONS: Answer 4 questions. 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.

Requirements for this paper (Please mark (X) as appropriate)

<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>

Log Tables
Graph Paper
Dictionaries
Statistical Tables
Bible

<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>

Thermodynamic Tables
Actuarial Tables
MCQ Only – Do not publish
Attached Answer Sheet
Exam Paper to be returned with Booklet

QUESTION 1**[TOTAL MARKS: 25]****Q 1(a)****[6 Marks]**

For a language pair of your choice, provide **three** examples of translational phenomena which demonstrate why MT is a difficult problem, no matter what type of system might be built.

Q 1(b)**[8 Marks]**

Provide the fundamental equations of (i) the noisy channel model of SMT, and (ii) the log-linear model of SMT. Name the different components in (i), and describe their basic function. Demonstrate how the two equations might be equivalent.

Q 1(c)**[6 Marks]**

Why do translation companies build customised solutions for their clients? Give **two** reasons why customised engines are likely to produce better output than a freely available web-based system such as Google Translate.

Q 1(d)**[5 Marks]**

Give **two** reasons why human translators are essential cogs in the MT pipeline. In your opinion, to what extent should human translators be fearful of the impact of MT on their profession?

[End of Question 1]

QUESTION 2

[TOTAL MARKS: 25]

Q 2(a)

[10 Marks]

Given two SMT systems that translate between English and Irish (in both directions), their performance can be improved by incorporating linguistic knowledge. Suggest **three** ideas to that end. For each idea:

- state its motivation (i.e. the linguistic problem that it is trying to solve),
- how it could be implemented, and
- for which direction (i.e. English-to-Irish or Irish-to-English).

There should be *at least one* idea for each of these phases: (a) pre-processing, (b) decoding and (c) post-processing. There should be *at least one* idea for each of the following types of linguistic knowledge: (i) morphology and (ii) syntax.

Q 2(b)

[6 Marks]

State the main disadvantage of adding linguistic knowledge to an SMT system. How is this exacerbated for a language like Irish, compared to (say) English?

Q 2(c)

[9 Marks]

State one pre-processing step that is commonly carried out when translating:

1. from Chinese, regardless of what the target language is;
2. between languages that follow different word orders;
3. from a language with a rich inflectional system into English.

[End of Question 2]

QUESTION 3**[TOTAL MARKS: 25]****Q 3(a) [4 Marks]**

What is the main reason to use sentence boundaries in *n*-gram-based ($n > 1$) language models? How are sentence boundaries typically represented in language models?

Q 3(b) [4 Marks]

State how a bigram language model would decompose the sentence "They didn't evaluate their SMT systems ." in order to calculate its probability, both with and without sentence boundaries.

Q 3(c) [3 Marks]

Given the sequence of words $x\ y\ z$, provide the formulae to calculate the probability of z using (i) a bigram language model, and (ii) a trigram language model.

Q 3(d) [7 Marks]

Calculate the probability of the last word in the sequence "the green witch" using bigram and trigram language models. The following sequences occur in the training data the number of times shown:

- "the", 600 times.
- "green", 10 times.
- "the green", 5 times.
- "green witch", twice.
- "the green witch", twice.
- "witch", 5 times.

Q 3(e) [3 Marks]

Why is it a good idea to use "smoothing" in the context of language modelling?

Q 3(f) [4 Marks]

What are the strengths and weaknesses of higher and lower order *n*-gram models?

[End of Question 3]

QUESTION 4**[TOTAL MARKS: 25]****Q 4(a)****[7 Marks]**

Regarding IBM Model 1, what algorithm do we usually use to obtain the word translation probabilities given a parallel corpus? Detail the steps involved in using this algorithm to calculate the IBM Model 1 score, using an example of your choice. .

Q 4(b)**[6 Marks]**

Regarding phrase-based SMT, how do we learn a phrase translation model from a parallel corpus? List the main steps involved, describe the main purpose of each step, and describe what happens at each step and the process involved at each step.

Q 4(c)**[6 Marks]**

How is a “phrase” defined in phrase-based SMT? What are the main differences between a phrase-based translation model and a word-based translation model?

Q 4(d)**[6 Marks]**

Describe the basic rule that we need to follow when we extract phrases. List all phrase pairs with the following word alignment based on the rule you have provided.

	A	B	C	D
M				
X				
Y				
Z				

[End of Question 4]

QUESTION 5

[TOTAL MARKS: 25]

Q 5(a)

[9 Marks]

Assume the following partial phrase table:

<i>ta</i>	<i>she</i>	0.4	<i>shanchang</i>	<i>likes</i>	0.3	<i>paobu</i>	<i>running</i>	0.6
			<i>shanchang</i>	<i>is good at</i>	0.7	<i>paobu</i>	<i>run</i>	0.4
<i>ta shanchang</i>	<i>she likes</i>	0.2	<i>shanchang paobu</i>	<i>likes running</i>	0.3			
<i>ta shanchang</i>	<i>she is good at</i>	0.8	<i>shanchang paobu</i>	<i>is good at running</i>	0.7			

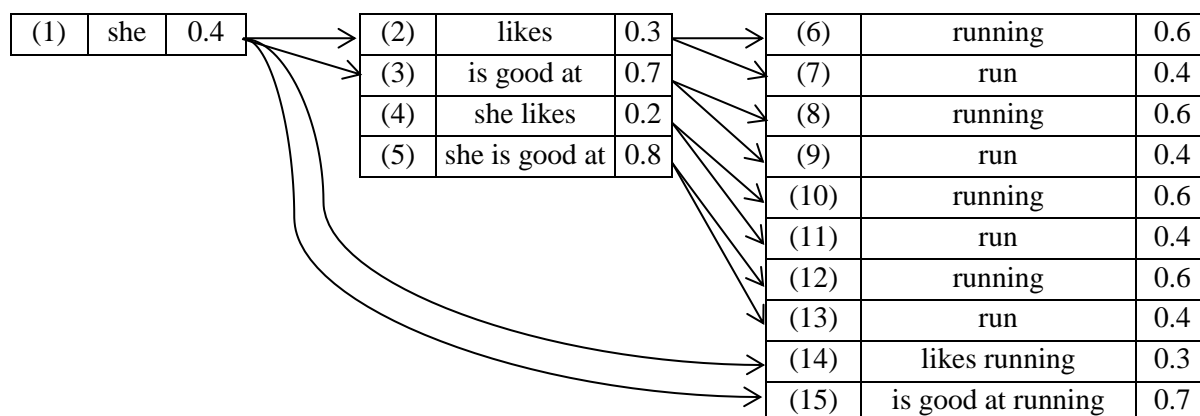
Consider the following input sentence:

ta shanchang paobu

Assume that:

- Only monotone word order is permitted;
- The language model is ignored.

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



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

Q 5(b)

[6 Marks]

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

Q 5(c)

[5 Marks]

Assume histogram pruning after recombination, where the maximum number of hypotheses in each stack is 2. Indicate which hypotheses will be pruned.

Q 5(d)

[5 Marks]

Assume threshold pruning after recombination, where the threshold is 0.5. Indicate which hypotheses will be pruned.

[End of Question 5]

[END OF EXAM]



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): `	Prof. Andy Way Dr. Jinhua Du Dr. Haithem Afli Prof. David Bustard Dr. Ian Pitt	(Ext: 5074) (Ext: 6716) (Ext: 8712)									
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 penyanyi.

The man in front of the house is 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

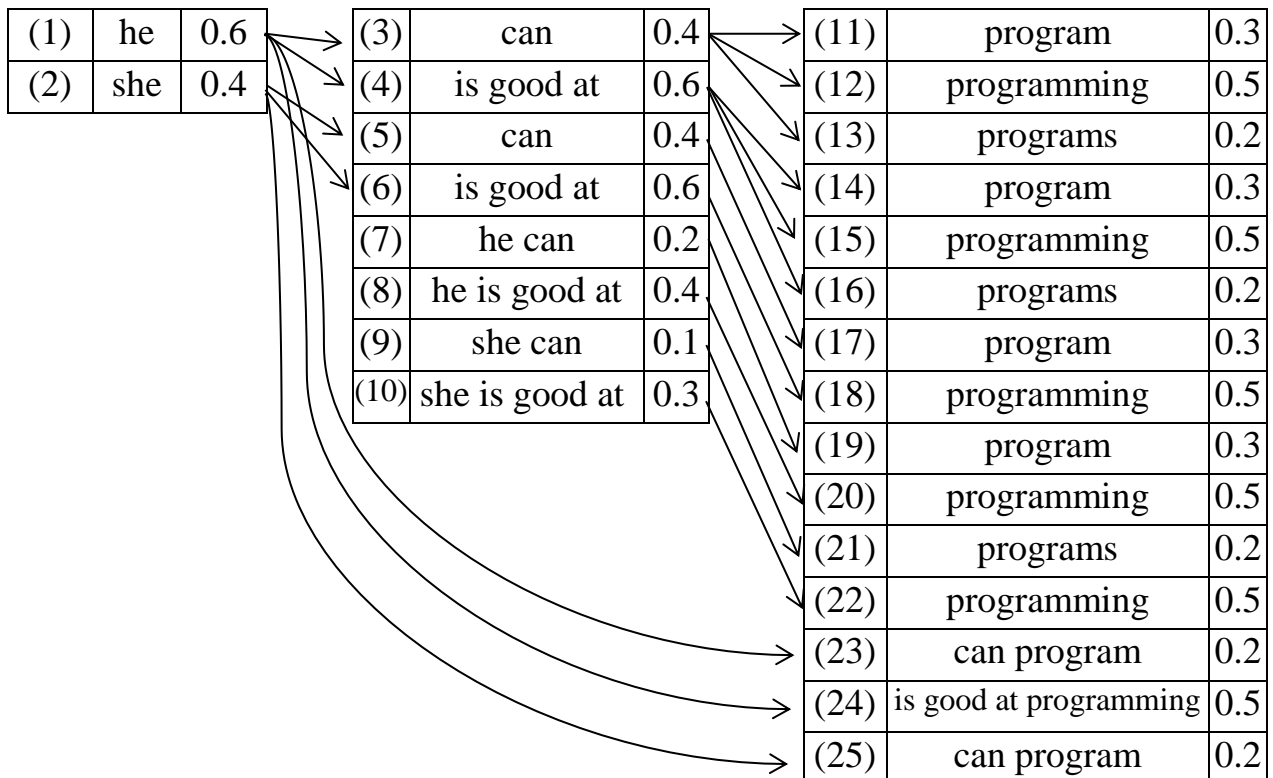
Consider the following input sentence:

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)

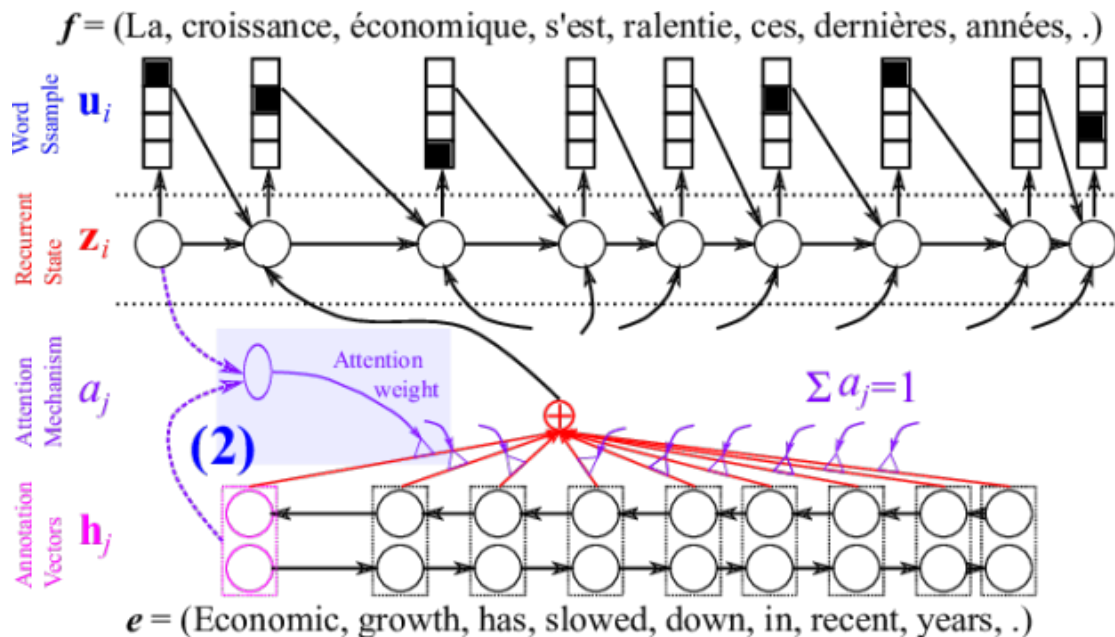
[4 Marks]

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]



DUBLIN CITY UNIVERSITY

SEMESTER 2 EXAMINATIONS 2017/2018

MODULE: CA4012 - Statistical Machine Translation

PROGRAMME(S):

CASE	BSc in Computer Applications (Sft.Eng.)
ECSAO	Study Abroad (Engineering & Computing)
CPSSD	BSc in Computational Problem Solv&SW Dev.

YEAR OF STUDY: 4,O

EXAMINER(S):

Prof. Andrew Way	(Ext: 5074)
Dr. Jinhua Du	(Ext: 6716)
Dr. Mohammed Hassanuzzaman	(Ext: 6912)
Prof. Brendan Tangney	
Dr. Hitesh Tewari	

TIME ALLOWED: 3 Hours

INSTRUCTIONS: Answer 5 questions. You must attempt at least one question from each 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)

[4 Marks]

What are the main advantages of the log-linear model of SMT compared to the noisy channel model of SMT? How might you argue that in some cases, the move from the noisy channel model to the log-linear model could be interpreted as a disadvantage?

Q 1(b)

[6 Marks]

Provide the fundamental equations of these two SMT models. Name the different components in each, and describe their basic function. Demonstrate how the two equations might be equivalent.

Q 1(c)

[6 Marks]

Until quite recently, the freely available web-based system Google Translate used entirely phrase-based SMT models, but now, certain language pairs are serviced via Neural MT (NMT). How would you build a case for switching over from SMT to NMT? What would the main differences be for Google in providing this new service compared to their previous offering?

Q 1(d)

[4 Marks]

Machine-learning researchers have had remarkable success in building specific models to beat the very best humans in certain fields, e.g. Go, Chess, Jeopardy. Why does MT continue to perform at a level lower than human translators, in general?

[End of Question 1]

QUESTION 2**[TOTAL MARKS: 20]****Q 2(a)****[8 Marks]**

SMT learns from data. What two types of data do we need to build SMT systems? What recommendations regarding training data would you provide to someone intending to build an SMT system so that the best possible performance was achieved?

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

Ideally, the translations output by MT systems should be both adequate and fluent. Explain why. Which components of an SMT system are primarily responsible for ensuring that these two constraints are met?

Q 2(c)**[3 Marks]**

Assume the reference translation for a particular source sentence is “Fruit flies like bananas”. Provide examples of (i) an inadequate but fluent translation; (ii) an adequate disfluent translation; and (iii) an inadequate disfluent translation.

Q 2(d)**[5 Marks]**

Despite the fact that MT has never been used as much as it is today, some translators continue to argue that MT will never be useful to them. Why do you think this is? Give **two** reasons that you believe would persuade them to try MT, and find it useful in their work.

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

SECTION B

QUESTION 3

[TOTAL MARKS: 20]

Q 3(a)

[10 Marks]

Given the following:

1. Source sentence:
Nous avons besoin d'évaluer nos systèmes statistiques de traduction.
2. Machine translation outputs:
We need to evaluate our statistical systems translation.
We evaluate our statistical translation systems.
3. Translation reference:
We need to evaluate our statistical translation systems.

Calculate the BLEU-3 scores of the two candidate translations.

Q 3(b)

[5 Marks]

Why is standard BLEU not suitable to evaluate MT output at sentence level (especially when sentences are short)? How would you modify BLEU to make it suitable to be used at sentence level?

Q 3(c)

[5 Marks]

What are the main benefits of performing automatic evaluation of MT output compared to human evaluation? At the same time, explain why it remains important to conduct human evaluations from time to time.

[End of Question 3]

QUESTION 4**[TOTAL MARKS: 20]****Q 4(a)****[10 Marks]**

Describe Word Error Rate (WER) in your own words. How is it calculated? What is it used for? What are the benefits/disadvantages associated with WER compared to the BLEU MT evaluation metric?

Q 4(b)**[5 Marks]**

Given the following strings:

Reference: They are responsible for the airport security.

MT output: They responsible security are airport.

Calculate the WER Score.

Q 4(c)**[5 Marks]**

Explain the differences between WER and Translation Edit Rate (TER).

[End of Question 4]

[END OF SECTION B]

SECTION C

QUESTION 5

[TOTAL MARKS: 20]

Q 5(a)

[6 Marks]

What is the EM algorithm? How is EM used to compute the word alignment and lexical probabilities in IBM Model 1? Provide **one** specific example of your own choosing which demonstrates how EM is calculated.

Q 5(b)

[5 Marks]

Describe the basic principles of word-based SMT and phrase-based SMT. Name **three** advantages and disadvantages of (i) word-based SMT and (ii) phrase-based SMT models.

Q 5(c)

[5 Marks]

In phrase-based SMT, what rules need to be followed in order to extract parallel phrases from a word-aligned parallel corpus?

Q 5(d)

[4 Marks]

Explain the main differences between higher IBM Models 2—5 and IBM Model 1.

[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>ti</i>	<i>plays</i>	0.6
<i>ti</i>	<i>is playing</i>	0.4

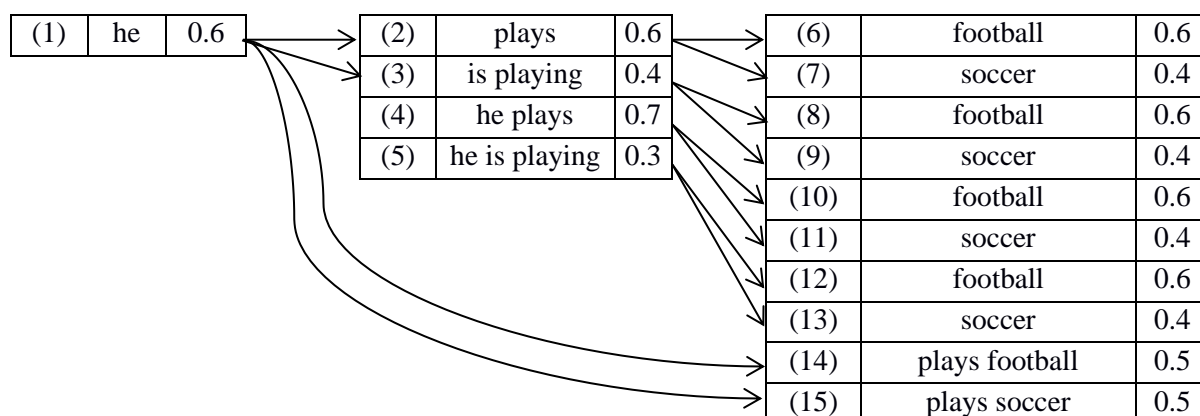
<i>zuqiu</i>	<i>football</i>	0.6
<i>zuqiu</i>	<i>soccer</i>	0.4

<i>ta ti</i>	<i>he plays</i>	0.7
<i>ta ti</i>	<i>he is playing</i>	0.3

<i>ti zuqiu</i>	<i>plays football</i>	0.5
<i>ti zuqiu</i>	<i>plays soccer</i>	0.5

Consider the following input sentence: *ta ti zuqiu*

Assume that (i) only monotone word order is permitted; and (ii) 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). Indicate also which hypothesis provides the optimal translation for the input sentence.

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

Given the search diagram above, indicate which groups of hypotheses can be recombined and indicate which hypothesis should be selected to represent each group.

Q 6(c)**[4 Marks]**

Assuming threshold pruning after recombination, where the threshold is 0.5, indicate which hypotheses will be pruned in the last stack.

[End of Question 6]

QUESTION 7

[TOTAL MARKS: 20]

Q 7(a)

[5 Marks]

What is the activation function in neural machine translation (NMT)? Provide the formulae of **two** commonly used activation functions and specify the output range of these functions.

Q 7(b)

[5 Marks]

What do you understand by the term “gradient”? In your own words, describe the gradient vanishing problem. Provide **one** method to alleviate this problem in neural networks.

Q 7(c)

[5 Marks]

What is the overfitting problem in neural networks? Describe **two** methods to alleviate this problem.

Q 7(d)

[5 Marks]

Draw the basic architecture of a recurrent neural network (RNN), and explain how the context of the input sequence is memorised in an RNN.

[End of Question 7]

[END OF SECTION C]

[END OF EXAM]



DUBLIN CITY UNIVERSITY

AUGUST/RESIT EXAMINATIONS 2017/2018

MODULE: CA4012 - Statistical Machine Translation

PROGRAMME(S):

CASE	BSc in Computer Applications (Sft.Eng.)
CPSSD	BSc in Computational Problem Solv&SW Dev.
ECSAO	Study Abroad (Engineering & Computing)

YEAR OF STUDY: 4,O

EXAMINER(S):

Andrew Way	(Internal)	(Ext: 5074)
Dr. Jinhua Du		(Ext: 6716)
Dr. Mohammed Hassanuzzaman		(Ext: 6912)
Prof. Brendan Tangney	(External)	External
Dr. Hitesh Tewari	(External)	External

TIME ALLOWED: 3 Hours

INSTRUCTIONS: Answer 5 questions. You must attempt at least one question from each 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)

[9 Marks]

Any statistical approach to MT requires the availability of aligned bilingual corpora which are (i) large, (ii) good-quality, and (iii) representative. Explain why all three requirements are important. What are the potential problems if the training corpus is (iv) small, (v) noisy, and (vi) unrepresentative compared to the intended test data?

Q 1(b)

[7 Marks]

Provide the fundamental equations of (i) the noisy channel model of SMT, and (ii) the log-linear model of SMT. With reference to these equations, name the different components in both models, and describe their basic function.

Q 1(c)

[4 Marks]

Give **two** reasons why customised engines are likely to produce better output than a freely available web-based system such as Google Translate.

[End of Question 1]

QUESTION 2

[TOTAL MARKS: 20]

Q 2(a)

[5 Marks]

Until the 1980s, all MT systems were rule-based. Now pretty much all MT systems are corpus-based. Why did corpus-based MT replace rule-based MT? Provide **two** examples where you think linguistic rules may still have a place in today's state-of-the-art MT systems.

Q 2(b)

[6 Marks]

The number of translation use-cases where there is no place for a human-in-the-loop are increasing. Provide **three** use-cases where MT is the only solution, i.e. that there is no place for human intervention in the translation pipeline for such use-cases.

Q 2(c)

[5 Marks]

Recently, Microsoft have claimed to have achieved “human parity” for Chinese-to-English neural MT (NMT), while SDL have claimed to have “cracked Russian-to-English NMT”. What do you think of such claims? How would you propose to test their validity?

Q 2(d)

[4 Marks]

Given claims such as those in 2(c) above, some human translators are fearful of the impact of MT on their profession. Give **two** reasons why human translators remain essential cogs in the MT pipeline.

[End of Question 2]

END OF SECTION A

SECTION B

QUESTION 3

[TOTAL MARKS: 20]

Q 3(a)

[6 Marks]

What are the **three** main types of MT evaluation? Briefly describe how you would conduct each of these types of evaluation.

Q 3(b)

[3 Marks]

When deciding whether to adopt MT or not, apart from translation quality, what other criteria need to be taken into account?

Q 3(c)

[6 Marks]

Why was the “brevity penalty” introduced in the BLEU automatic evaluation metric? Is fluency accounted for in BLEU via “*n*-gram precision” or “word precision”? How can the correlation between BLEU and human quality scores be improved?

Q 3(d)

[5 Marks]

How can the MT quality on new text be predicted? Provide **three** features which could be useful in this regard.

[End of Question 3]

QUESTION 4

[TOTAL MARKS: 20]

Q 4(a)

[8 Marks]

Describe in your own words the “ n -gram language model”. Show how a bigram language model would decompose the following sentence to calculate its probability, both with and without sentence boundaries:

“They did n’t evaluate their systems .”

Q 4(b)

[4 Marks]

How can the quality of a language model be evaluated? Why is this useful?

Q 4(c)

[5 Marks]

Why do we need “smoothing” in language modelling? Describe **two** smoothing methods.

Q 4(d)

[3 Marks]

Describe **three** techniques that can be used to manage very large language models.

[End of Question 4]

END OF SECTION B

SECTION C

QUESTION 5

[TOTAL MARKS: 20]

Q 5(a)

[4 Marks]

Assume the following Chinese—English sentence pairs:

S_1	S_2
<i>diannao</i> <i>computer</i>	<i>xin diannao</i> <i>new computer</i>

The source side is Chinese, and the target side is English. 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:

$t(\text{computer}|\text{diannao})$
 $t(\text{new}|\text{diannao})$
 $t(\text{new}|\text{xin})$
 $t(\text{computer}|\text{xin})$

Q 5(c)**[6 Marks]**

Assume the following English—Chinese parallel sentence pair:

<i>good reputation of the school</i> <i>xuexiao de hao shengyu</i>

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

<i>good</i>	<i>reputation</i>	<i>of</i>	<i>the</i>	<i>school</i>

[End of Question 5]

QUESTION 6**[TOTAL MARKS: 20]****Q 6(a)****[5 Marks]**

In the context of phrase-based SMT: (i) what is the definition of the term “phrase”?
(ii) List **three** advantages of using phrases as atomic units in phrase-based SMT compared to word-based SMT.

Q 6(b)**[7 Marks]**

In order to build a phrase-based SMT model from a parallel corpus: (i) what are the **three** basic steps that need to be followed? (ii) What kinds of methods can we use to obtain a symmetrised word alignment from two unidirectional word alignments? (iii) In the process of extracting parallel phrases, what rules need to be followed?

Q 6(c)**[8 Marks]**

After phrase pairs are extracted from the symmetrised word alignment, (i) in your own words, describe the basic idea of how phrase pairs are scored.

(ii) The following table provides four different Chinese translations for the English phrase “Good morning”.

Translation	Counts
zaoshang hao	50
zaoshang	5
zhongwu hao	15
nihao	30

Using the method you described in (i), calculate the probability of each phrase pair.

[End of Question 6]

QUESTION 7

[TOTAL MARKS: 20]

Q 7(a)

[4 Marks]

Why do we typically use “recurrent” neural networks instead of “feed-forward” neural networks in neural machine translation (NMT)? Describe the **single** fundamental difference between these two architectures.

Q 7(b)

[4 Marks]

Explain what the “encoder” and the “decoder” do in an encoder-decoder NMT system.

Q 7(c)

[4 Marks]

A baseline encoder-decoder NMT system can be improved by means of an attention mechanism.

- (i) what is an “attention mechanism”?
- (ii) how can the attention mechanism improve translation quality compared to the baseline encoder-decoder NMT system with no attention?

Q 7(d)

[4 Marks]

In NMT, briefly describe the stochastic gradient descent algorithm. List **two** commonly used optimisers for NMT, and briefly describe how they work.

Q 7(e)

[4 Marks]

In NMT, briefly describe the term “activation function”. List **two** commonly used activation functions in neural networks, and briefly describe how they differ.

[End of Question 7]

END OF SECTION C

[END OF EXAM]

SEMESTER 2 EXAMINATIONS 2018/2019

MODULE: CA4012 - Statistical Machine Translation

PROGRAMME(S):

CASE	BSc in Computer Applications (Sft.Eng.)
CPSSD	BSc in Computational Problem Solv & SW Dev.
ECSA	Study Abroad (Engineering & Computing)
ECSAO	Study Abroad (Engineering & Computing)

YEAR OF STUDY: 4,O,X

EXAMINER(S):

Prof Andrew Way	(Internal)	(Ext: 5074)
Dr Mohammed Hasanuzzaman	(Internal)	(Ext: 6719)
Dr Dimitar Shterionov	(Internal)	(Ext: 6719)
Dr Hitesh Tewari	(External)	
Prof Brendan Tangney	(External)	

TIME ALLOWED: 3 Hours

INSTRUCTIONS: Answer **five** questions. You **must** attempt *at least one* question from *each* 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]

Any statistical approach to MT requires the availability of aligned bilingual corpora which are (i) large, (ii) good-quality, and (iii) representative. Explain why all three requirements are important.

Q 1(b)

[7 Marks]

Provide the fundamental equations of (i) the noisy channel model of SMT, and (ii) the log-linear model of SMT. With reference to these equations, name the different components in both models, and describe their basic function.

Q 1(c)

[4 Marks]

Give **two** reasons why we can expect a customised client-specific MT engine to produce superior quality translations compared to a freely available web-based system such as Google Translate.

Q 1(d)

[3 Marks]

With the switch from statistical to neural MT, should we be concerned that the field will be dominated more than ever by the large providers such as Google Translate and Bing Translator?

[End of Question 1]

QUESTION 2**[TOTAL MARKS: 20]****Q 2(a)****[5 Marks]**

Corpus-based models have outperformed rule-based MT systems for 30 years now. In your opinion, do rules play *no* part in today's state-of-the-art neural MT models, or is the way forward a hybrid combination of rules and data-driven systems? Provide **three** examples to support your argument.

Q 2(b)**[9 Marks]**

In 2016, Google announced that their neural models were “bridging the gap” between MT and human translation quality. In 2018, Microsoft claimed to have achieved “human parity” for Chinese-to-English neural MT. What do you think of such claims? How would you go about testing them? What advice would you give to translators who were concerned at the possibility of losing their job in light of such claims?

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

While in most cases neural models of translation clearly outperform statistical MT (SMT), our ability to explain their outputs is less than it was for SMT. Why do you think this is the case? How would you suggest that neural models could become more inspectable?

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

SECTION B

QUESTION 3

[TOTAL MARKS: 20]

Q 3(a)

[5 Marks]

Explain the Markov assumption. Why do we need to take it into account when building n -gram language models?

Q 3(b)

[5 Marks]

How is the Maximum Likelihood estimate of a *trigram* language model computed? Compute $P(\text{ate}|\text{Bukka})$ from the following unigram and bigram counts.

Bukka sandwich	10
ate Bukka	8
Bukka the	15
Bukka	35
Bukka ate	16
ate the	20
sandwich ate	2
sandwich Bukka	3

Q 3(c)

[4 Marks]

Assume (i) a unigram and (ii) a bigram language model trained on standard English. Based on these two different language models, how would you expect the probability of the phrase “the sandwich Hakka ate” to compare to the probability of the sentence “Hakka ate the sandwich”?

Q 3(d)

[6 Marks]

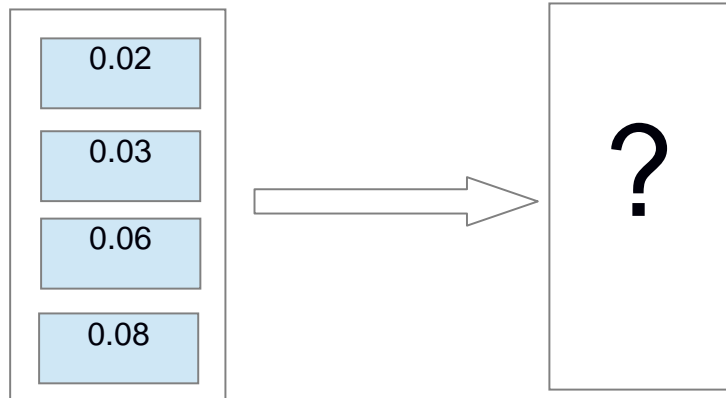
Why do n -gram language models need to be smoothed? Name **three** methods of smoothing, and explain one method in detail.

[End of Question 3]

QUESTION 4
Q 4(a)

[TOTAL MARKS: 20]
[6 Marks]

In decoding, what is pruning, and why it is important? How many of the following hypotheses will be pruned if we prune all hypotheses that are at least 0.4 times worse than the best hypothesis? Show your calculations.



Q 4(b)

[7 Marks]

Assume the following partial phrase table:

se	she	0.4
bhalobase	likes	0.3
bhalobase	likes to	0.5
se bhalobase	he likes	0.3
se bhalobase	he likes to	0.6
khete	eat	0.6
khete	eating	0.7
khete bhalobase	likes eating	0.3
khete bhalobase	likes to eat	0.7

Assume that only monotone word order is permitted and that the language model is ignored.

Draw the search graph constructed during decoding for the sentence “se bhalobase khete”.

Q 4(c)

[7 Marks]

Using the constructed search graph for Q 4(b), calculate all possible hypotheses, and indicate which hypothesis provides the optimal translation for the above sentence.

[End of Question 4]

[END OF SECTION B]

SECTION C

QUESTION 5

[TOTAL MARKS: 20]

Q 5(a)

[4 Marks]

Explain Word Error Rate (WER). How is the WER score computed?

Q 5(b)

[4 Marks]

Given the following two candidate translations (from *MT system 1* and from *MT system 2*), compute their WER scores with respect to the reference provided.

- *Candidate MT system 1*: near the shore the ship sank
- *Candidate MT system 2*: the big ship sank close to the shore
- *Reference*: the large ship sank near the shore

Q 5(c)

[6 Marks]

Why is the BLEU score typically calculated over the entire document, rather than on a sentence-by-sentence level? Give **two** examples of the relevant shortcomings of this metric.

Q 5(d)

[6 Marks]

In BLEU we compute the *n*-gram clipped precision (typically for 1-, 2-, 3- and 4-grams). What is clipped precision and why do we need to perform clipping? To support your explanation, provide **two** examples.

[End of Question 5]

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

The goal of MT is to find a sentence e that is the most likely translation of a source-language sentence f , $p(e|f)$.

Why is the translation model concerned with the translation probabilities of *words* or *phrases*, rather than the sentence as a whole? Use the following corpus to provide **one** example that would support your answer.

Source (Dutch)	Target (English)
ze speelt tennis graag .	she likes to play tennis .
hij eet pizza graag .	he likes eating pizza .
ze eet spaghetti graag .	she likes eating spaghetti .
hij speelt voetbal graag .	he likes to play football .

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

Expectation maximisation (EM) is an algorithm to iteratively estimate translation probabilities and alignments. Explain what the following three formulae compute. How are they employed in the EM algorithm?

1.	$\prod_{j=1}^{l_e} \frac{t(e_j f_{a(j)})}{\sum_{i=0}^{l_f} t(e_j f_i)}$
2.	$\sum_a p(a \mathbf{e}, \mathbf{f}) \sum_{j=1}^{l_e} \delta(e, e_j) \delta(f, f_{a(j)})$
3.	$\frac{\sum_{(\mathbf{e}, \mathbf{f})} c(e f; \mathbf{e}, \mathbf{f})}{\sum_e \sum_{(\mathbf{e}, \mathbf{f})} c(e f; \mathbf{e}, \mathbf{f})}$

Q 6(c)**[4 Marks]**

In your own words, explain IBM model 1, IBM model 2, IBM model 3 and IBM Model 4, focusing in particular on their differences and similarities.

Q 6(d)**[4 Marks]**

Assume the following two sentence pairs and alignments:

Greek: i gata ekatse sto halaki
English: the cat sat on the mat
alignment: 1->1, 2->2, 3->3, 4->4, 5->4, 6->5

Greek: i gata ekatse sto halaki
English: sat on the mat the cat
alignment: 1->3, 2->4, 3->4, 4->5, 5->1, 6->2

Using the above examples, explain and illustrate the deficiencies of IBM Model 1.
How might you propose a solution to this problem?

[End of Question 6]

QUESTION 7

[TOTAL MARKS: 20]

Q 7(a)

[6 Marks]

In neural machine translation, the most commonly used architecture is the encoder-decoder model. Explain the basic principles of the encoder-decoder model.

What type of neural networks are suitable for use in the encoder-decoder model, and why?

Q 7(b)

[7 Marks]

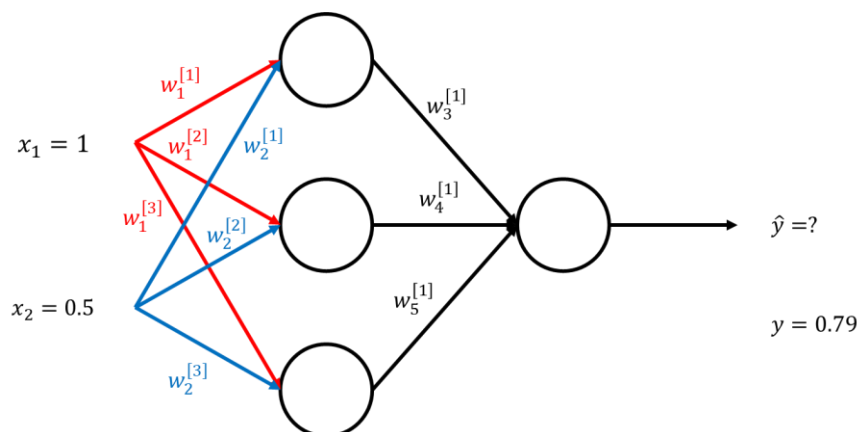
In training neural networks, a common problem is overfitting. Explain what overfitting is in your own words. Provide **two** ways of dealing with this problem, together with a description of how they work.

Q 7(c)

[7 Marks]

Below is given a simple neural network with two inputs, one output and one hidden layer. All the neurons in the hidden layer and the output layer use the **sigmoid** activation function and **all biases are equal to 1**.

Given the input $x_1=1$, $x_2=0.5$ and the expected output of $y=0.79$, perform **one** pass of a forward and backward propagation on this network. All initial weights are uniformly distributed (as is noted in the table). What are the missing values in the table below the network, corresponding to the **weights**, **network output** and **error** after the forward and the backward steps?



	Input 1	Input 2	Weights from input 1 to hidden layer			Weights from input 2 to hidden layer			Weights from hidden layer to output layer			Network Output (\hat{y})	Expected Output (y)	Error
			$w_1^{[1]}$	$w_1^{[2]}$	$w_1^{[3]}$	$w_2^{[1]}$	$w_2^{[2]}$	$w_2^{[3]}$	$w_3^{[1]}$	$w_4^{[1]}$	$w_5^{[1]}$			
Forward	$x_1 = 1$	$x_2 = 0.5$	0.3	0.3	0.3	0.3	0.3	0.3	0.3	0.3	0.3		0.79	
Backward	$x_1 = 1$	$x_2 = 0.5$											0.79	

[End of Question 7]

[END OF SECTION C]

[END OF EXAM]

AUGUST/RESIT EXAMINATIONS 2018/2019

MODULE: CA4012 - Statistical Machine Translation

PROGRAMME(S):

CASE	BSc in Computer Applications (Sft.Eng.)
CPSSD	BSc in Computational Problem Solv & SW Dev.
ECSA	Study Abroad (Engineering & Computing)
ECSAO	Study Abroad (Engineering & Computing)

YEAR OF STUDY: 4,O,X

EXAMINER(S):

Prof Andrew Way	(Internal)	(Ext: 5074)
Dr Mohammed Hasanuzzaman	(Internal)	(Ext: 6719)
Dr Dimitar Shterionov	(Internal)	(Ext: 6719)
Dr Hitesh Tewari	(External)	
Prof Brendan Tangney	(External)	

TIME ALLOWED: 3 Hours

INSTRUCTIONS: Answer **five** questions. You **must** attempt *at least one* question from *each* 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]

For a language pair of your choice, provide **three** examples of translational phenomena which demonstrate why MT is a difficult problem, no matter what type of system might be built.

Q 1(b)

[4 Marks]

What are the main advantages of the log-linear model of SMT compared to the noisy channel model of SMT? How might you argue that in some cases, the move from the noisy channel model to the log-linear model could be interpreted as a disadvantage?

Q 1(c)

[6 Marks]

Until recently, the freely available web-based system Google Translate used entirely phrase-based SMT models, but now its engines are built using Neural MT (NMT). From the perspective of MT, what factors would Google have taken into account in coming to this decision? What would the main differences be for Google in providing this new service compared to their previous offering?

Q 1(d)

[4 Marks]

With the advent of AI approaches, some researchers consider MT to be a 'solved problem'. Do you agree with that claim? Would you expect NMT to beat SMT on all language pairs, in all domains?

[End of Question 1]

QUESTION 2**[TOTAL MARKS: 20]****Q 2(a)****[8 Marks]**

SMT learns from data. What two types of data do we need to build SMT systems? What recommendations regarding training data would you provide to someone intending to build an SMT system so that the best possible performance was achieved?

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

Why did developers of SMT systems switch from word-based to phrase-based models?

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

Ideally, the translations output by MT systems should be both adequate and fluent. Explain why. Which components of an SMT system are primarily responsible for ensuring that these two constraints are met?

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

MT has never been used as much as it is today, but some translators continue to argue that MT will never be useful to them. Why do you think this is? Give **two** reasons that you believe would persuade them to try MT, which could prove useful in their work.

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

SECTION B

QUESTION 3

[TOTAL MARKS: 20]

Q 3(a)

[8 Marks]

Explain the n -gram language model. Show how a bigram language model would decompose the following sentence in order to calculate its probability, both with and without sentence boundaries.

“I strongly believe that Australia is going to win ICCWC2019.”

Q 3(b)

[3 Marks]

How can we measure the quality of a language model? Why is this useful?

Q 3(c)

[4 Marks]

Describe **four** techniques to manage very large language models.

Q 3(d)

[5 Marks]

Assume three separate language models based on unigrams, bigrams and trigrams trained on standard English. Based on each of these language models, how would you expect the probability of the sentence ‘the doughnut Homer ate’ to compare to the probability of ‘Homer ate the doughnut’?

[End of Question 3]

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

What is the main task of decoding? Explain **two** types of errors encountered in decoding. How can the quality of the decoding process be evaluated?

Q 4(b)**[9 Marks]**

In decoding, what is hypothesis recombination? Why it is important? Explain **two** types of pruning strategies, using examples of your choice.

Q 4(c)**[6 Marks]**

Assume the following partial phrase table:

<i>ta</i>	<i>she</i>	0.4
<i>shanchang</i>	<i>likes</i>	0.3
<i>shanchang</i>	<i>is good at</i>	0.7
<i>paobu</i>	<i>running</i>	0.6
<i>paobu</i>	<i>run</i>	0.4
<i>ta shanchang</i>	<i>she likes</i>	0.2
<i>ta shanchang</i>	<i>she is good at</i>	0.8
<i>shanchang paobu</i>	<i>likes running</i>	0.3
<i>shanchang paobu</i>	<i>is good at running</i>	0.7

Assume that:

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

Draw the search diagram for the following input sentence:

ta shanchang paobu

[End of Question 4]

[END OF SECTION B]

SECTION C

QUESTION 5

[TOTAL MARKS: 20]

Q 5(a)

[4 Marks]

Explain **two** approaches to human evaluation of MT output. Identify **two** drawbacks to human evaluation.

Q 5(b)

[6 Marks]

Name **three** issues that the METEOR evaluation metric addresses which BLEU, WER and TER do not. Provide examples to support your answer.

Q 5(c)

[6 Marks]

Given the following two translations (from MT system 1 and MT system 2), compute their BLEU score (max. n -gram size 3) with respect to the reference.

Candidate MT system 1: near the shore the ship sank

Candidate MT system 2: the big ship sank close to the shore

Reference: the large ship sank near the shore

Q 5(d)

[4 Marks]

Identify **two** drawbacks of automatic evaluation metrics, and **two** reasons why they are useful.

[End of Question 5]

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

Training an MT system is typically based on a corpus of parallel sentences. Why do we need to compute word alignments? What alignment patterns are you aware of, and why is word alignment a challenging task?

Q 6(b)**[7 Marks]**

Explain the four steps of the expectation maximisation algorithm (as defined for translation modelling). Focus on what counts or probabilities are computed at each step and how are they connected.

Q 6(c)**[7 Marks]**

What are the differences between the Normal and the Simplified IBM model 1? For the Simplified IBM model 1, write down the formula and explain what it takes into account in terms of word alignment.

[End of Question 6]

QUESTION 7

[TOTAL MARKS: 20]

Q 7(a)

[7 Marks]

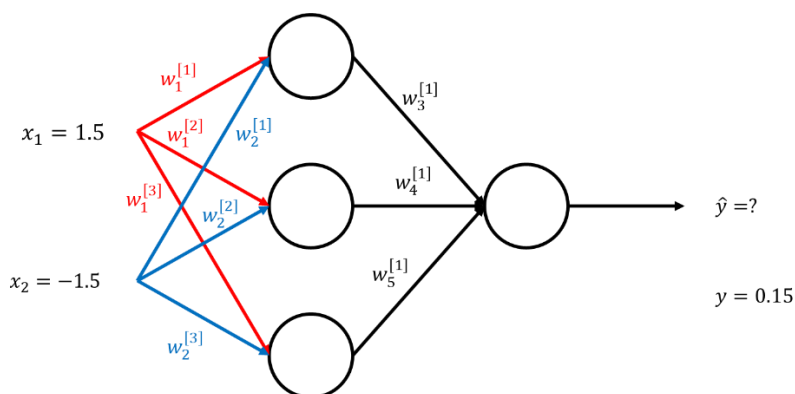
What are the differences between a multi-layer perceptron and a Recurrent Neural Network? Which one would you use in an encoder-decoder NMT architecture, and why?

Q 7(b)

[6 Marks]

Below we provide a simple neural network with two inputs, one output layer, and one hidden layer. All biases are equal to 1. All initial weights are uniformly distributed. Their values are noted in the table below.

Given the inputs $x_1 = 1.5$, $x_2 = -1.5$ and the expected output of $y = 0.15$, choose two activation functions and for each perform one forward pass through this network. Explain the differences between the chosen activation functions. Fill in the blanks in the table below the network, corresponding to the weights, network output and error after the forward and the backward steps. Use the quadratic cost function in your calculations.



	Input 1	Input 2	Weights from input 1 to hidden layer			Weights from input 2 to hidden layer			Weights from hidden layer to output layer			Network Output (\hat{y})	Expected Output (y)	Error
			$w_1^{[1]}$	$w_1^{[2]}$	$w_1^{[3]}$	$w_2^{[1]}$	$w_2^{[2]}$	$w_2^{[3]}$	$w_3^{[1]}$	$w_4^{[1]}$	$w_5^{[1]}$			
Forward (activation function 1)	$x_1 = 1.5$	$x_2 = -1.5$	0.3	0.3	0.3	0.3	0.3	0.3	0.3	0.3	0.3		$y = 0.15$	
Forward (activation function 2)	$x_1 = 1.5$	$x_2 = -1.5$											$y = 0.15$	

Q 7(c)

[7 Marks]

Explain the gradient descent algorithm for training neural networks. What is the difference between gradient descent and stochastic gradient descent? What is the role of the learning rate?

[End of Question 7]

[END OF SECTION C]

[END OF EXAM]

SEMESTER 2 EXAMINATIONS 2021/2022

MODULE: CA4012 - Machine Translation

PROGRAMME(S):

CASE BSc in Computer Applications (Sft.Eng.)

YEAR OF STUDY: 4

EXAMINER(S): Prof. Andy Way (ext. 5074)
Dr. Maja Popović (ext. 6736)
Dr. Guodong Xie
Dr. Sheila Castilho
Dr. Pintu Lohar

TIME ALLOWED: 3 Hours

INSTRUCTIONS: Answer 5 questions. You must attempt *at least one* question from *each* 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]

Any statistical approach to MT requires the availability of aligned bilingual corpora which are (i) large, (ii) good-quality, and (iii) representative. Explain why all three requirements are important.

Q 1(b)

[7 Marks]

Provide the fundamental equations of (i) the noisy channel model of SMT, and (ii) the log-linear model of SMT. With reference to these equations, name the different components in both models, and describe their basic function.

Q 1(c)

[4 Marks]

What are the main advantages of the log-linear model of SMT compared to the noisy channel model of SMT? How might you argue that in some cases, the move from the noisy channel model to the log-linear model could be interpreted as a disadvantage?

Q 1(d)

[3 Marks]

Do you agree with the claim by some researchers that with the advent of AI approaches, MT should now be considered to be a 'solved problem'?

[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]**

Give *three* ways in which professional human translators are of critical importance in the development and testing of MT systems.

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

Despite the obvious advances in quality seen from the move to neural models, why do you think some human translators are still reluctant to embrace MT technology? How might you try to persuade them that MT can be a useful tool in the translator's armoury?

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

SECTION B

QUESTION 3

[TOTAL MARKS: 20]

Q 3(a)

[9 Marks]

Calculate the WER score for the source sentence (SRC), reference sentence (REF) and the MT output hypothesis (HYP) below:

SRC: Kuopion kaupunginvaltuusto hyväksyi liitoksen yksimielisesti maanantaina

REF: The union was unanimously approved by the Kuopio City Council on Monday

HYP: The city council unanimously approved the joint Niiralan on Monday

Q 3(b)

[6 Marks]

Explain the concepts of “adequacy”, “fluency” and “comprehensibility” and state why each of them is important for the evaluation of MT output.

Q 3(c)

[5 Marks]

What are advantages and disadvantages of automatic evaluation of MT systems?

[End of Question 3]

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

Explain the Markov assumption. Why do we need to take it into account when building n -gram language models?

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

Consider the following sentences:

- *I love the cat .*
- *I love the dog .*
- *They see the dog .*
- *They talk to the girl .*
- *They talk to the dog .*
- *I talk to the cat .*

Using a unigram language model, calculate the probability of the sentence "*I see the cat .*"

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

Explain the advantages and disadvantages of lower and higher order n -gram models.

Q 4(d)**[5 Marks]**

How can the problem of unseen n -grams be overcome in count-based n -gram language modelling? Describe one possible method which can be applied to resolve this issue.

[End of Question 4]

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

Explain the differences between the IBM models of word alignment. How does the “Hidden Markov model (HMM)” differ from “IBM model 2”?

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

What is the main disadvantage of using sentence-level translation probabilities in SMT? Derive the fundamental equation of translation probability when word alignment is introduced.

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

Name and define *four* types of word alignments based on mapping cardinality between source and target words. Which of those alignments are supported by the IBM-models?

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

Calculate IBM-1 lexical probabilities for the following parallel training corpus:

English	French
the house	la maison
the dog	le chien
the cat	le chat

[End of Question 5]**[End of SECTION B]**

SECTION C

QUESTION 6

[TOTAL MARKS: 20]

Q 6(a)

[10 Marks]

Select the following statements which represent the characteristics of neural MT. (2 marks for each correct choice, -2 marks for each incorrect choice but the total score cannot be negative).

- a) it can be regarded as using only local context
- b) it can be regarded as using global context
- c) it can be regarded as being guided by neural language model
- d) the neural language model is one of its many components
- e) it includes a coverage constraint mechanism that enables avoiding repetitions and omissions in translation
- f) sometimes a part of a sentence can be translated more than once
- g) there are no obvious word alignments
- h) it is possible to provide a model introspection and understand well the source of any translation errors
- i) because of the sub-word technology, it sometimes generates misspelled or non-existing words

Q 6(b)

[5 Marks]

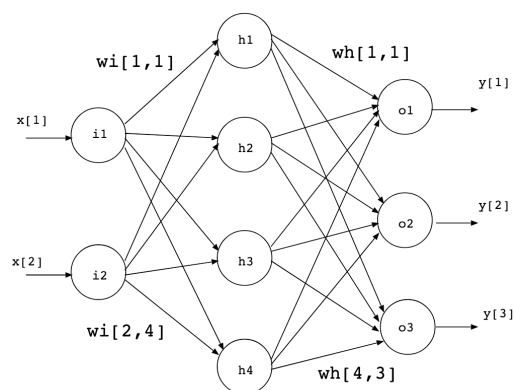
If a neural network with three artificial neurons in the output layer is used for classifying images of animals, and the outputs of the network represent the labels y_1 = “cat”, y_2 = “dog” and y_3 = “mouse”, which activation function would you choose for the neurons in the output layer? Explain your reasoning.

Q 6(c)

[5 Marks]

For the neural network on the right, the following weights connecting the input layer with the hidden layer are given:

- $w_{i11} = 0.7$
- $w_{i12} = 1.4$
- $w_{i13} = 3.2$
- $w_{i14} = -0.3$
- $w_{i21} = -2.3$
- $w_{i22} = 1.2$
- $w_{i23} = -4.7$
- $w_{i24} = 2.3$



If the input values coming into the network are $x_1=3$ and $x_2=2$, what will be the input to the activation function of the fourth neuron in the hidden layer? (h_4)

[End of Question 6]

QUESTION 7**[TOTAL MARKS: 20]****Q 7(a)****[6 Marks]**

Explain why feed-forward neural networks are not appropriate for MT, and why recurrent neural networks are.

Q 7(b)**[8 Marks]**

The following Spanish sentence:

Me gustan los gatos

should be translated into English as:

I like cats

If we have an NMT system which decodes using greedy search, what would be the translation of the entire sentence if the following probabilities of partial hypotheses are generated by the model:

first target word:

they	0.28
I	0.16
we	0.10
like	0.02
it	0.01

first two target words:

they like	0.15
they are	0.10
I like	0.11
I want	0.09
I think	0.03
we like	0.05
we will	0.02
like me	0.002
like you	0.001
it will	0.001

first three target words (complete hypothesis):

they are nice	0.002
they are here	0.001

they are mine	0.003
they like cats	0.08
they like dogs	0.06
they like me	0.01
I like cats	0.09
I like dogs	0.05
I think we	0.002
I think that	0.0001
we like cats	0.01
we like mice	0.005
we will go	0.0012
we will see	0.0001
like me ,	0.0012
like you ,	0.0007
it will be	0.0002
it will go	0.0001

Q 7(c)

[6 Marks]

Explain what is meant by the term “vanishing gradient”. What network architectures are prone to it, under which circumstances?

[End of Question 7]

QUESTION 8**[TOTAL MARKS: 20]****Q 8(a)****[5 Marks]**

Name and describe *five* differences between RNN and transformer architectures used for NMT.

Q 8(b)**[5 Marks]**

Explain the term “attention” in the context of neural networks. Which types of attention are used in the transformer architecture?

Q 8(c)**[5 Marks]**

What are “word representations”? Where do they come from?

Q 8(d)**[5 Marks]**

What is the dimension of the input layer of a recurrent neural language model? Explain your answer.

[End of Question 8]***[End of SECTION C]******[END OF EXAM]***