# Natural language processing Mias Ghantous – 213461692 Faisal Omari – 325616894

### Section 1)

#### Part 1:

We have implemented a function called  $make\_list\_of\_sentence$  which takes a sentence and make the tokens as required, how did we do it?

- 1. If we have "'" at the end of the token then this a regular word if the first character is a Hebrew latter and the second to last.
- 2. If we have Hebrew latter at the end of the token and also at the beginning then this is also a regular word
- 3. Anything else is not a Hebrew word that we should include.

And to make the list of all the sentences we implemented  $make\_list$  which uses  $make\_list\_of\_sentence$ 

#### Questions:

1. Min\_count: is the minimum count that would be included in the training, if it was big then we are going to learn on the most frequent words in the language, which leads to overfitting and if it were too small then we would learn on words that we don't use them very often like names of countries which will include noise, we defined it to be in the 1 that means that every word in the corpus will be included but we would have put it bigger.

Window: is the max distance between the current word and the predict word in a sentence, small value leads to over fitting as we are not going to see a lot of repeated combinations, and not enough data on medium frequent combinations of words. And big value leads to a lot of noise as the not frequent combinations would appear in the training, we choose Window = 5 which we think is enough not too big and not too small.

Vector\_size: is the embedding vector size which tells how much does words are similar, too big of a number leads to overfitting because of the number of features, and too small would lead to underfitting because the feature would not be representative enough, we choose 100 which is okay.

2. The problem is that we work with Hebrew and in Hebrew we have a lot of combined words that make one meaning like "בית ספר" and our corpus is not big enough to cope with this and also in Hebrew there is tokens that have many words like "זכשמהבית" which means 'and from the hose' and we created the corpus in a way that all of these words are just one token also we have words like "בצל" which can mean "in the shadow" or "onion" 2 words that have totally different meanings and in the corpus we don't attempt to fix this.

### Section 2)

Part 1. Implementation in function *Section2\_part1*, we put all the words in a list and iterate throw the list.

Our results:

Knesset similar words.txt:

```
א knesset_similar_words.txt ×

1 (0.09953677654266357, (מומשלה, 0.37252214550971985, מושלה), (0.08438344299793243, מומשלה, 0.3935275673866272, מומשלה, 0.4523324966430692394, מומשלה, 0.4523324966430692394, מומשלה, 0.2829371381378174, (מומשלה, 0.5479722619056702, מומשלה, 0.6516565084457397, מומשלה, 0.6516565084573, מומשלה, 0.6516565084457397, מומשלה, 0.6516565084573, מומשלה, 0.6516565084573,
```

Part 2. Implementation in *embeddings\_of\_sentences* which takes the model also the sentences and return the a dictionary from the index of the corpus to the average embedding of each sentence:

Part 3. We search the corpus for sentences that potentially we can use we choose them to be not that long has very common words in the language also sentences that can be said a lot in protocols and also can be said in a different way too, and if we are not satisfied with the result of the sentence we simply change it to another sentence.

our results:

#### Knesset similar sentences.txt:

```
    וארesset_similar_sentences.txt ×
    זה בדיוק אותו דבר : most similar sentence: אבל זה אותו דבר : את זה צריך להביא בחשבון : most similar sentence: ולכן , צריך להביא את זה בחשבון : most similar sentence: אני לא מוכן לקבל את זה : לא , אני מוכן לבדוק את זה : most similar sentence: אני לא מוכן לקבל את זה : most similar sentence: אם כן , רבותי, אנחנו עוברים להצבעה : most similar sentence: אם כן , רבותי, אנחנו עוברים להצבעה : most similar sentence: זה לא דבר שהוא חדש : most similar sentence: זה לא דבר שהוא חדש ? מה הקונספציה שלכם בעניין הזה : most similar sentence: בגלל שאני אומר את דעתי , ודעתי שונה מדעתך : most similar sentence : בגלל שאני אומר את תקיים היום . היום , בכל אופן , לא תהיה הצבעה : most similar sentence : אני לא כל כך הבנתי : most similar sentence : אני לא כל כך הבנתי : most similar sentence : איך קורה דבר כזה : most similar sentence : איך ייתכן דבר כזה : most similar sentence : איך ייתכן דבר כזה : most similar sentence : איך ייתכן דבר כזה : most similar sentence : איך ייתכן דבר כזה : most similar sentence : איך ייתכן דבר כזה : most similar sentence : איך ייתכן דבר כזה : most similar sentence : איך ייתכן דבר כזה : most similar sentence : איך ייתכן דבר כזה : most similar sentence : איך ייתכן דבר כזה : most similar sentence : mo
```

Part 4. Our way is to first try the same word if any of the first 3 is ok then take it Else try a similar word or a word that fits the sentence if any of the first 3 are ok then take it. Else, try more than one word that are related.

For each word we have tried a lot and some of them got nice results and some of them we didn't get useful results.

לחדר: השתמשנו במילה נרדפת "אולם" וכדי לקבל את ה "ל" השתמשנו ב שתבוא וקיבלנו "לאולם" מוכנה: השתמשנו באותה מילה ובמילה שיכולה לחליף אותה כמעת תמיד "יכול" וקבלנו "יכולה".

ההסכם: השתמשנו באותה מילה וקיבלנו ההמשך שיא יכולה להחליף אבל משנה את המשמעות.

טוב: השתמשנו ב "שמש" כדי לנסות לקבל "אור" אבל הגענו למילה פחות טובה "בריא"

פותח: השתמשנו ב "מתחיל" כי היא יכולה להחליף את המילה בלי לשנות את משמעות המשפט וקיבלנו מילה משנה קצת את המשמעות "ממשיך"

שלום: זו המילה הקשה ביותר נסינו לקבל "אהלן" אבל לא הגענו לשום דבר, ואז סמנו מלים שמשתמשים בהם עם חבריים ("תודה", "חברי","רבותי") עם "שלום" וקיבלנו עמיתי, שהיא יכולה להופיע במשפט הזה.

היקר: השתמשנו ב "הטוב" כי היא מילה שיכולה להחליף וקיבלנו "גדול", היא משנה קצת את המשמעות אבל זה מה שיכלנו לעשות.

בשנה: השתמשנו באותה מילה וקיבלנו "השנה".

We can see that in all of the words we have grammatical errors, some of the words changed the meaning of the words a bit.

Our results:

red words sentences.txt:

#### Questions:

1. we assumed that related words would be more similar like:

```
כנסת, ממשלה: קבלה 0.45 שהוא ערך גבוה וזה הגיוני כי כנסת וממשלה מאוד כשורים. ישראל, כנסת: קיבלה 0.11 הקשר של הכנסת הוא יותר גדול. ישראל, חבר: אין כשר ולכן נמוך 0.08 כי שראל, חבר: אין כשר ולכן נמוך 0.08 כי מילה עם עצמה: זה הכי טבעי שכל אחת תקבל 1.0 כי המילה היא הכי קרובה לעצמה ישראל שולחן: אין דמיון ולכן ערך נמוך 0.09
```

but there are also words that have almost no relation, but got high value like and we didn't expect that, also words that are expected them to be more related because we are talking about Knesset corpus, like:

```
חבר, כנסת: יש הרבה הופעות של "חבר כנסת" אבל קבלנו ערך 0.2 שהוא אמור להיות יותר, זה יכול להיות בגלל ש כנסת הופיעה יותר מן "חבר כנסת" .
שולחן, חבר: קיבלנו ערך 0.6 שזה מאוד גדול והם אינם דומים בכלל.
שלום, שולחן: פחות קיבלנו 0.31 צפינו שיהיה.
```

Why this happen? We don't have big enough corpus, if we had one we would have got better results.

- 2. We would have small destine because a lot time, if we have 2 words (antonyms) and have a sentence that use one of the words (let's say the first one), then we can replace the first word with the second word and have a valid sentence.
- 3. We tried a lot of words and here are the results:

```
קר,חם: 0.8986342549324036
קטן,גדול: 0.6460744738578796
קטן,גדול: 0.6482081604003906
רע,טוב: 0.8434505462646484
עצוב,שמח: 0.5723463296890259
עני,עשיר: 0.5578402876853943
אחרון,ראשון: 0.9238438606262207
```

As we can see that all of the antonyms has a score over 0.55 which means that they are very similar, which is as we explained in the previous question, and even got words over 0.84.

4. Here are our results again:

```
1 ה בדיוק אותו דבר : most similar sentence ואבל זה אותו דבר : most similar sentence ולכן , צריך להביא את זה בחשבון : most similar sentence ולכן , צריך להביא את זה בחשבון . לא , אני מוכן לבדוק את זה : אני לא מוכן לקבל את זה . אם כן , אנחנו עוברים להצבעה : most similar sentence . אם כן , רבותי, אנחנו עוברים להצבעה . מבחינתך זה דבר שהוא לא נורא : most similar sentence . זה לא דבר שהוא חדש . מה הקונספציה שלכם בעניין הזה : most similar sentence : מה התפקיד שלכם בנושא הזה . אני אומר את דעתי , ודעתי שונה מדעתך : most similar sentence : בגלל שאני אומר את האמת . אני אופן , לא תהיה הצבעה : most similar sentence . בכל מקרה ההצבעה לא תתקיים היום . אני לא כל כך הבנתי : most similar sentence : איך קורה דבר כזה : most similar sentence : איך ייתכן דבר כזה ? איך ייתכן דבר כזה : most similar sentence : איך ייתכן דבר כזה ? איך ייתכן דבר כזה
```

We see that, the sentences 1,2,4,8,9,10 have the same meaning.

In sentence 3: the both are talking about ideas that needs to be checked (the second sentence) and to be agreed up on (the first sentence) and both rejection.

Sentence 5: both of the sides are describing a subject ("דבר").

Sentence 6: the word "נושא" is similar to "נושא" also in the both of the sentences the talker is talking in the same ("אתם") and it's ("אתם"), we can see that by using ("שלכם")

Sentence 7: both of the sentence give the something on the talker point of view.

As we can see we have very nice results and that is because of the way we choose the sentences, for each sentence, the words of the sentence have a high probability to appear together and so they have high probability to appear another time in the corpus, and that's why we can find similar sentences.

# Section 3)

We tired 25, 100,50 and 10 and the got similar accuracy so we choose to be in the middle k=50 and here are the results because we don't want too big or too small of a k:

# For 100:

	or chunk size			f1-score	support	or chunk siz	e 3 s validation: precision		f1-score	support	for chunk siz KNN with cor	ze 5 ss validation:			
	committee	0.61	0.84	0.71	29476	committee	0.69	0.88	0.77	9802		precision	recall	f1-score	support
	plenary	0.74	0.47	0.57	29476	plenary	0.83	0.60	0.70	9802	committee	0.74	0.88	0.80	5872
											plenary	0.85	0.69	0.76	5872
	accuracy			0.65	58952	accuracy			0.74	19604					
	macro avg	0.68	0.65	0.64	58952	macro avg	0.76	0.74	0.73	19604	accuracy			0.78	11744
1	weighted avg	0.68	0.65	0.64	58952	veighted avg	0.76	0.74	0.73	19604	macro avg	0.79	0.78	0.78	11744
											weighted avg	0.79	0.78	0.78	11744
ŀ	(NN with spli					(NN with spli									
		precision	recall	f1-score	support		precision	recall	f1-score	support	KNN with spl:				
												precision	recall	f1-score	support
	committee	0.62	0.84	0.71	2948	committee	0.68	0.87	0.76	981					
	plenary	0.74	0.48	0.58	2948	plenary	0.82	0.59	0.69	980	committee	0.72	0.87	0.79	588
											plenary	0.83	0.66	0.74	587
	accuracy			0.66	5896	accuracy			0.73	1961	2001122011			0.77	1175
	macro ave	0.68	0.66	0.65	5896	macro avg	0.75	0.73	0.72	1961	accuracy	0.78	0.77	0.77	1175
	weighted avg	0.68	0.66	0.65	5896	veighted avg	0.75	0.73	0.72	1961	macro avg weighted avg	0.78	0.77	0.76	1175
ľ	expired usb										weighted avg	0.78	0.//	0.76	11/5

## For 50:

for chunk size KNN with cors			f1-score	support	for chunk size KNN with cors	e 3 s validation precision	: recall	f1-score	support		ze 5 ss validation: precision		f1-score	support
committee plenary	0.61 0.73	0.82 0.49	0.70 0.58	29476 29476	committee plenary	0.69 0.82	0.87 0.61	0.77 0.70	9802 9802		0.74 0.85	0.87 0.70	0.80 0.76	5872 5872
accuracy macro avg weighted avg	0.67 0.67	0.65 0.65	0.65 0.64 0.64	58952 58952 58952	accuracy macro avg weighted avg	0.76 0.76	0.74 0.74	0.74 0.73 0.73	19604 19604 19604	macro avg	0.79 0.79	0.78 0.78	0.78 0.78 0.78	11744 11744 11744
KNN with spli	t: precision	recall	f1-score	support	KNN with spli	t: precision	recall	f1-score	support	KNN with spli	it: precision	recall	f1-score	support
committee plenary	0.62 0.73	0.82 0.49	0.70 0.59	2948 2948	committee plenary	0.68 0.80	0.85 0.61	0.76 0.69	981 980		0.73 0.82	0.85 0.69	0.79 0.75	588 587
accuracy macro avg weighted avg	0.67 0.67	0.65 0.65	0.65 0.65 0.65	5896 5896 5896	accuracy macro avg weighted avg	0.74 0.74	0.73 0.73	0.73 0.73 0.73	1961 1961 1961	macro avg	0.78 0.78	0.77 0.77	0.77 0.77 0.77	1175 1175 1175

## For 25:

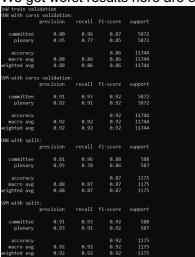
for chunk siz KNN with cors	e 1 s validation: precision		f1-score	support	for chunk siz KNN with cors			f1-score		For chunk siz (NN with cors			f1-score	support
committee plenary	0.62 0.70	0.77 0.53	0.69 0.61	29476 29476	committee plenary	0.70 0.80	0.85 0.63	0.76 0.71	9802 9802	committee plenary	0.75 0.83	0.85 0.72	0.80 0.77	5872 5872
accuracy macro avg weighted avg	0.66 0.66	0.65 0.65	0.65 0.65 0.65	58952 58952 58952	accuracy macro avg weighted avg	0.75 0.75	0.74 0.74	0.74 0.74 0.74	19604 19604 19604	accuracy macro avg veighted avg	0.79 0.79	0.78 0.78	0.78 0.78 0.78	11744 11744 11744
KNN with spli	t: precision	recall	f1-score	support	KNN with spli	t: precision	recall	f1-score	support	(NN with spli	t: precision	recall	f1-score	support
committee plenary	0.63 0.70	0.76 0.54	0.69 0.61	2948 2948	committee plenary	0.70 0.80	0.84 0.63	0.76 0.71	981 980	committee plenary	0.74 0.81	0.84 0.71	0.79 0.76	588 587
accuracy macro avg weighted avg	0.66 0.66	0.65 0.65	0.65 0.65 0.65	5896 5896 5896	accuracy macro avg weighted avg	0.75 0.75	0.74 0.74	0.74 0.74 0.74	1961 1961 1961	accuracy macro avg veighted avg	0.78 0.78	0.77 0.77	0.77 0.77 0.77	1175 1175 1175

# For 10:

for chunk siz	e 1													
KNN with cors	s validation				for chunk siz					for chunk siz				
	precision	recall	f1-score	support	KNN with cors					(NN with cors	s validation			
						precision	recall	f1-score	support		precision	recall	f1-score	support
committee	0.62	0.79	0.69	29476										
plenary	0.71	0.51	0.59	29476	committee	0.68	0.85	0.76	9802	committee	0.74	0.86	0.80	5872
					plenary	0.81	0.60	0.69	9802	plenary	0.84	0.70	0.76	5872
accuracy			0.65	58952	200110201			0.73	19604					
macro avg	0.66	0.65	0.64	58952	accuracy	0.74	0.73			accuracy			0.78	11744
weighted avg	0.66	0.65	0.64	58952	macro avg	0.74	0.73	0.73	19604	macro avg	0.79	0.78	0.78	11744
merginees and					weighted avg	0.74	0.73	0.73	19604	veighted avg	0.79	0.78	0.78	11744
KNN with spli		100 111 111 111												
	precision	recall	f1-score	support	KNN with spli					(NN with spli				
				-append		precision	recall	f1-score	support		precision	recall	f1-score	support
committee	0.61	0.78	0.69	2948	committee	0.68	0.85	0.76	981					
plenary	0.70	0.51	0.59	2948		0.80	0.61	0.69	980	committee	0.73	0.85	0.78	588
, , ,					plenary	0.80	0.01	0.09	980	plenary	0.82	0.68	0.74	587
accuracy			0.64	5896	accuracy			0.73	1961					4475
macro avg	0.66	0.64	0.64	5896		0.74	0.73	0.73	1961	accuracy			0.77	1175
weighted avg	0.66	0.64	0.64	5896	macro avg					macro avg	0.77	0.77	0.76	1175
weighted avg	0.00	0.04	0.04	3890	weighted avg	0.74	0.73	0.72	1961	weighted avg	0.77	0.77	0.76	1175

#### Questions:

1. We got worst results here are our results from the past home work:



for chunk siz	e 5										
KNN with corss validation:											
	precision	recall	f1-score	support							
committee	0.74	0.86	0.80	5872							
plenary	0.84	0.70	0.76	5872							
,											
accuracy			0.78	11744							
macro avg	0.79	0.78	0.78	11744							
weighted avg	0.79	0.78	0.78	11744							
weighted avg	0.79	0.78	0.78	11/44							
MM											
KNN with spli											
	precision	recall	f1-score	support							
committee	0.73	0.85	0.78	588							
plenary	0.82	0.68	0.74	587							
accuracy			0.77	1175							
macro avg	0.77	0.77	0.76	1175							
weighted avg	0.77	0.77	0.76	1175							
weighted avg	0.77	0.77	0.70	11/3							

- It could be because that the current feature vector is not as representative of the class as the previous feature vector, and that is because previously we had feature for each word, but now 100 features are going to represent all the words.
- 3. We always get that size 1 is worse than 3, and 3 worse than 5, why?

if we have small chunk size then we would have bigger feature matrix which is worst for run time, also each chunk would have features that doesn't represent the types good enough, because the "feature power" would split among the smaller chunks rather than be powerful in one chunk, which make the classification task harder.

## Section 4)

1. Form grammar perspective we see that most the sentences have good grammar And most of them also make since like:

but here an example on a sentence that didn't make since:

Original sentence: מסיבות שונות , [\*] 1,275 התיקים שנפתחו בשנת [\*] - טרם הוגשו כתבי אישום . DictaBERT sentence: מסיבות שונות , בכל 1,275 התיקים שנפתחו בשנת שנפתחו - טרם הוגשו כתבי אישום . DictaBERT tokens: בכל, שנפתחו

Original sentence: בסופו של דבר , גם אני הייתי [\*] אליפויות העולם רצוף שבע שנים , וכן יצאנו עם דגל ישראל ואני יודעת שחוץ מאימא שלי , אף [\*] לא יודע שבאמת אנחנו מייצגים את ישראל :DictaBERT sentence: בסופו של דבר , גם אני הייתי בכל אליפויות העולם רצוף שבע שנים , וכן יצאנו עם דגל ישראל ואני יודעת שחוץ מאימא שלי , אף אחד לא יודע שבאמת אנחנו מייצגים את ישראל ב-DictaBERT tokens: ". בכל אחד ". בכל אחד בסופו של דבר , גם אני הייתי בכל אליפויות העולם רצוף שבע שנים , וכן יצאנו עם דגל ישראל ואני יודעת שחוץ מאימא שלי , אף אחד לא יודע שבאמת אנחנו מייצגים את ישראל. אחד ב-. במל אחד ב-. במופר של מייצגים את ישראל ב-. ".

Original sentence: [\*] רחבות , שקשורות בעיקר לצרכים של השקיפות בכל מה שקשור בתהליך קבלת התרומות האלה [\*] Original sentence: בעצם הגענו להסכמות מאוד רחבות , שקשורות בעיקר לצרכים של השקיפות בכל מה שקשור בתהליך קבלת התרומות האלה ": DictaBERT sentence DictaBERT tokens: ", מאוד

The first sentence the token "שנפתחו" is a token a good token from grammar stand point also has no meaning.

The second and third,

we can see that the model predicted '\"' without there to be any closing one which is not what grammar say.

2. Most of the results in the previous Homework were punctuation mark but now their actual words, and even when the previous model got words our current model gives more accurate predictions like:

Original sentence: אבל הנושא הוא [\*] אם אתה בעד [\*] ההתנתקות או נגד . Committee sentence:אבל הנושא הוא צריך אם אתה בעד . ההתנתקות או נגד . Committee tokens: ,צריך,

Probability of committee sentence in committee corpus: -83.088

Probability of committee sentence in plenary corpus: -75.356

This sentence is more likely to appear in corpus: plenary

. אבל הנושא הוא נושא אם אתה בעד או ההתנתקות או נגד Plenary sentence: . אבל הנושא הוא נושא אם אתה בעד או

צושא,או :Plenary tokens

Probability of plenary sentence in plenary corpus: -70.322 Probability of plenary sentence in committee corpus: -90.4

This sentence is more likely to appear in corpus: plenary

. אבל הנושא הוא [\*] אם אתה בעד [\*] ההתנתקות או נגד

DictaBERT sentence: אבל הנושא הוא לא אם אתה בעד תוכנית ההתנתקות או נגד

לא, תוכנית :DictaBERT tokens

### Other example:

. עזבי את [\*] אני לא מכיר את חוק [\*] והבנייה , 84 אני לא מכיר את חוק .

. עזבי את זה 84 , אני לא מכיר את חוק חובת והבנייה:Committee sentence

זה,חובת :Committee tokens

Probability of committee sentence in committee corpus: -83.109

Probability of committee sentence in plenary corpus: -86.292

This sentence is more likely to appear in corpus: committee

. עזבי את זה 84 , אני לא מכיר את חוק ההסדרים והבנייה :Plenary sentence

זה,ההסדרים :Plenary tokens

Probability of plenary sentence in plenary corpus: -83.718

Probability of plenary sentence in committee corpus: -83.758

This sentence is more likely to appear in corpus: plenary

. עזבי את [\*] 84 , אני לא מכיר את חוק [\*] והבנייה :Original sentence

. עזבי את סעיף 84, אני לא מכיר את חוק התכנון והבנייה , 84 אני לא מכיר את חוק התכנון והבנייה .

DictaBERT tokens: סעיף, התכנון

As we see the current model is better than the previous one.

3. As we can see above the model didn't work good on some sentences, reasons may be this happens:

We see that in the first sentence the token must be a number (year number) it could be that he does not predict a numbers well because there are endless possibilities.

As for the second and third, if he can see all the sentence then he would know there is no another ' \" ' which make him predict another token.

Images of the Output:

### dictabert results.txt

```
set_dictabert_results.txt ×
  original sentence: ציור (*) ביצירת מקומות עבודה רבים יותר , בתשלום שכר (*) בעד העבודה :
• ש צורך גם ביצירת מקומות עבודה רבים יותר , בתשלום שכר גבות בעד העבודה -
• שצורך גם ביצירת מקומות עבודה רבים יותר , בתשלום שכר גבות בעד
 Original sentence: אבל הנושא הוא [*] אם אתה בעד [*] ההתנתקות או נגד.
DictaBERT sentence: אבל הנושא הוא לא אם אתה בעד תוכנית ההתנתקות או נגד.
DictaBERT tokens: לא, תוכנית
 Original sentence: לפון המשיח מתנגדת (•) החוק הזאת , (פ) למפרון הבשיח בדרך DictaBERT sentence: מ- לפון המשטח מתנגדת להצעת החוק הזאת , כדי לחביא לפפרון הבשיח בדרך אורת ה
- לכן , המתשלה מתנגדת להצעת החוק הזאת , כדי לחביא לפפרון הבשיח בדרך אורת בהסרט הלבוע DictaBERT total מידעת, בלי, אורת בהסרט ה
  Original sentence: לא הוראת שעה (*) [*] א DictaBERT sentence: ההצמדה היא כבר לא הוראת שנה
היא, כבר לא הוראת שנה
  Original sentence: עובי את (*) אני לא מכיר את חוק (*) והבניה: BictaBERT sentence.
עובי את סעיף 84 , אני לא מכיר את חוק התכנון והבניה: DictaBERT tokens: סעיף, התכנון
  Original sentence: מיליארד [*] עלות החוק (*] .
DictaBERT sentence: שלות החוק היא מיליארד
DictaBERT tokens: היא, שקל
  Original sentence: מסיבות שונות , [*] – טרם הונשו כתבי אישו (1,275 (*) , מסיבות שונות , ב775 (1,275 בשנת 1,275 בשנת 1,275 בשנת שונתת במלל 1,275 המיקים שנפתחו בשנת שופתחו – טרם הונשו כתבי אישו 1,275 המיקים שנפתחו בשנת שופתחו – טרם הונשו כתבי אישו במלל שופתחו במלל, שנפתחו
  Original sentence: (*), על תשובתך המקיפה (ד'), (ד'), (ד'), על תשובתך המקיפה ב-
DictaBERT sentence: תודה לד', מיכל , על תשובתך המקיפה .
תודה, מיכל
  Original sentence: (#) אחר של המטבש המיד בוא נראה את האו (#) אומר רק דבר אחד : בוא נראה את הצד האחר של המטבט .
ט. אני אומר רק דבר אחד : בוא נראה את הצד האחר של המטבט ב. DictaBERT tokens: אני, הצר
  Original sentence: אני מנסה להבין את הפער הזה , איך יכול (*) שמצד אחד הכל (*) , הכל תחת רגולציה ומצד שני , (*) ומזהמים ברישיון :
אני מנסה להבין את הפער הזה , איך יכול להיות שמצד אחד הכל בסדר , הכל תחת רגולציה ומצד שני , הכל ומזהמים ברישיון :
DictaBERT tokens: להיות, בסדר, הכל
  Original sentence: א שליפויות תעולם רצוף שבע שנים, וכן יצאנו עם דבל ישראל ואני יודעה שמוץ מאימא שלי , אף (*) לא יודע שבאתת אנהנו מייצנים את ישראל בי בי שראל ואני יודעה שמוץ מאימא שלי , אף אחד לא יודע שבאתת אנהנו מייצנים את ישראל בי שראל ואני יודעה שחוץ מאימא שלי , אף אחד לא יודע שבאתת אנהנו מייצנים את ישראל בי שראל ואני יודעה שחוץ מאימא שלי , אף אחד לא יודע שבאתת אנהנו מייצנים את ישראל בי שראל ואני יודעה שחוץ מאימא שלי , אף אחד לא יודע שבאתת אנהנו מייצנים את ישראל בי שראל בי שראל ואני יודעה שחוץ מאימא שלי .
  Original sentence: (*) אבל אומרים לו , חמשת ע (*) אבל (*) מאוד מאוד ארוך וכראי שתפעה את הבדיקה DictaBERT sentence: אבל אומרים לו , חשמע , החור אצל הרופא מאוד מאוד ארוך וכראי שתעשה את הבדיקה הזאת DictaBERT (אבל אומר האוד ארוך וכראי שתעשה התור, הרופא הזאת EATLER (אבל אות הור אצל האוד את הור אות אות האוד את החור את הרובא הזאת בהצראה בהור את החור את הרובא הזאת בהצראה בהור את החור את הרובא הזאת הרובא הדיקה הור את הרובא הדיקה הור את הרובא הדיקה הור את הרובא הדיקה הור את הרובא הרוב
 Original sentence: [*] באנגלית [*] .
DictaBERT sentence: עיתוני העולם הם ברובם באנגלית
Original sentence: אותו : [*] , אתה קיצעת ולא נתת משכורת למורה שלי , [*] אותו : [*] אותו ? DictaBERT sentence: אז היא שאלה אותו : מה , אתה קיצעת ולא נתת משכורת למורה שלי DictaBERT tokens: שאלה, מה
 Original sentence: [*] רחבות , שקשורות בעיקר לצרכים של השקיפות בכל מה שקשור בתהליך קבלת התרומות האלה [*] DictaBERT sentence: ". האלה " בעים הגענו להסכמות מאוד רחבות , שקשורות בעיקר לצרכים של השקיפות בכל מה שקשור בתהליך קבלת התרומות האלה " .
. בעצם הגענו להסכמות מאוד רחבות , שקשורות בעיקר לצרכים של השקיפות בכל מה שקשור בתהליך קבלת התרומות האלה " . מ
  Original sentence: (☀) לא מופיע : Original sentence אוד מחברי (☀) בבית ורואה שאף אחד מחברי (☀) לא מופיע : ObctaBERT sentence . אני חושכת שזה מאוד לא נכון לנהל כזה דיון שרואים שנאמת הציבור יושב בבית ורואה שאף אחד מחברי הנוסת לא חופיע : DictaBERT tokens .

DictaBERT tokens: מנהל, יושם, הכנסת
  Original sentence: בחינה במל-אביב, וגם אין מספיק מודעות לספורט הזה , במדינה בארץ, שהוא נמצא בפארק דרום במל-אביב, וגם אין מספיק מודעות לספורט הזה , במדינה: DictaBERT sentence 
ציפרי הרב , זה די קשה להחשיך , לשעות סקי ולהגיע להישנים , כי גם יש לנו אתר אחד , היחיד בארץ , שהוא נמצא בפארק דרום במל-אביב , וגם אין מספיק מודעות לספורט הזה , במדינה 
DictaBERT tokens: היחיד אור , אתר ולעשרי, ,, אתר
  original sentence: מל ידי נ'ניסר לסלו מהמזרחי ובמשך 8 שנים ניהלנו את קבוצות המוקד מלוס (≨) לכיינ'ינג בעניין מסרים מנצחים (≨) להקשר באשר לישראל DictaBERT sentence: .
הוא נוסד בשנת 2002 של ידי ג'ניסר לסלו מהמזרחי ובמשך 8 שנים ניהלנו את קבוצות החוקד מלוס קבוצות לביינ'ינג בעניין מסרים מנצחים ' להקשר באשר לישראל DictaBERT tokens: ,
  Original sentence: [*] את נושא (*), לדוגמה (*).
DictaBERT sentence: בחר את נושא השיחה, לדוגמה
בחר, השיחה
  Original sentence: [*] (*) מראש ?
DictaBERT sentence: מה לא לעצור ולחשוב מראש
DictaBERT tokens: מה, לא, ולחשוב
   Original sentence: זה לא [*] בשום [*]
DictaBERT sentence: לא פוגע בשום דבר
DictaBERT tokens: פוגע, דבר
  Original sentence: אנחנו צריכים [*] טובי המשפטנים.
אנחנו צריכים את טובי המשפטנים .
DictaBERT tokens: א
  Original sentence: [*] התייחטות ממוקדת (*)
DictaBERT sentence: זה תחום שצריך לקבל התייחטות מחוקדת.
DictaBERT tokens: זה, לקבל
  Original sentence: [*] הוא מאוד מובדר.
• לא אתן , נושא זה הוא מאוד מובדר .
• אני לא אתן , נושא זה הוא מאוד מובדר .
• אני, זה DictaBERT tokens: אני, זה
```