

The Bible is All You Need: Translating the Bible using Transformers

Yorai Shaoul, Dana Rosenfarb, Zachary Metzman, Doron Hazan

Abstract

Hebrew has changed drastically from its use in the Old Testament. In fact, it has changed so much that even native Hebrew speakers have trouble understanding the ancient or Biblical language. Given the immense interest in ancient languages that is shared by linguists, religious figures, and students, we wish to facilitate the understanding of old languages by modern day readers via deep learning methods. To this end, we made use of state-of-the-art translation methods to obtain a deep neural model that translates Biblical Hebrew to Modern Hebrew. In this work we present our collected dataset of a Modern Hebrew Translation of the Bible and two Biblical to Modern Hebrew machine translation models. Additionally, we provide a quantitative and qualitative evaluation of our methods. We report BLEU scores of 82% and 98%, showing the efficacy of our methods in the task of Hebrew translation, and provide a discussion suggesting that adaptation of pretrained language models to new tasks might be superior to end-to-end learning.

The code for this project is available online at https://drive.google.com/drive/folders/1sK_ECSV_buGGQ-xtkKV5m7iXCcLeiT?usp=sharing.

1 Introduction

The Bible is widely considered to be the most read book existing.¹ It is a collection of religious texts or scriptures sacred to Jews, Samaritans, Christians and others. Believers might generally consider the Bible to be a product of divine inspiration, and the Bible is considered as the source containing the fundamentals of the Jewish, Christian and various other religions. The straight link to religion, popularity and reputation of the Bible led it

¹According to the 2015 edition of Guinness World Records



Figure 1: The Hebrew Alphabet and English Pronunciation. (Source: Jewish Museum London.)

to have a profound influence on innumerable humans throughout history.²

Over the years, the Hebrew language has evolved from its ancient roots to a modern and (somewhat) widely used language. As a result, older Hebrew texts (i.e., the Old Testament) take on a different structure and style from their modern translations. Therefore, the Old Testament or the Bible written in Biblical Hebrew is very hard to read and perhaps even uninterpretable by many Hebrew speakers nowadays. Since the Bible and its content have immense influence on so many people, its translation has immense ramifications. Namely, the translation of the Bible has to be precise and interpretable.

However, we found no evidence for previous literature that aimed to establish mapping between Biblical and Modern Hebrew. In this paper, we used existing Natural Language Processing approaches to try complete this task. We looked at various sources that offer Bible translations. Most

²The Hebrew Bible shares most of its content with its ancient Greek translation, the Septuagint, which is the base for the Christian Old Testament. The Christian New Testament is a collection of writings by early Christians, believed to be mostly Jewish disciples of Christ, written in first-century Koine Greek. In this paper we use the Biblical Hebrew Old Testament

of the authors' native language is Hebrew and have experience reading it to make an independent judgement of the quality of the translation offered by the different sources. We concluded that the ultimate source for our project is the public Annotated Bible (מִקְרָא מִבּוֹאָר). This source is publicly available on Wikitext and includes translations that we believe are widely acceptable within the Hebrew-speaking community.

In addition, our work was also motivated by providing potential contributions to these following needs:

- A need to provide a translation framework between ancient and modern languages. This was rarely done in the past.
- A need to shed more light about the Bible and about Biblical Hebrew. By translating the Bible to an interpretable, modern language, we will elucidate the Bible, and make it more accessible to those who desire to understand it.
- A need to deepen our understanding of the differences between Biblical and Modern Hebrew. With the appropriate experiments, we can detect fundamental similarities and differences between the two, such as vocabulary similarity, transformation of words (ancient to modern), etc.
- A need to expand the amount of languages the Bible is translated to (extrapolation). Our models provide framework to translate the Bible to other languages that perhaps do not provide existing translations to the Bible.

Overall, in this work we provide a comprehensive analysis on using three translation models on Biblical to Modern Hebrew, conduct experiments and provide insights that we hope would shed more light on the task of translating the Bible and Biblical Hebrew.

1.1 Issues with Translation

Some of meaning of particular words in the Bible are not actually known with total certainty. Commentators such as Rashi and Ramban interpret the Bible and provide translations. As a result, different biblical translators may translate the Bible not just with different words, but with also fairly different meanings.

2 Related Work

To our knowledge this task has not been researched before. Therefore, we look to general translation literature and literature related to translated of similar languages, partially building on work reviewed by (Zampieri et al., 2020) overviewing the challenges in translation between dialects. In addition, Harrat et al. (2019) provides a detailed overview of translations between dialects of arabic. However, this review is fairly old and does not cite any method that make use of Sequence-to-Sequence models. Laith H. Baniata (2018) uses a more modern approach to translation of arabic dialects by using an encoder-decoder method. The paper used a multi-task learning model to share one decoder among language pairs. There has also been some older literature related to translation between Shakesperian english and modern english (Xu et al., 2012).

While there has not been much literature related to translation between biblical and modern Hebrew, there is a fair amount of literature related to applying NLP methods to Hebrew text. This literature provides general information that is helpful when working with Hebrew text. Tsarfaty et al. (2019) notes some of the challenges with common NLP libraries such as ScaPy in relation to Hebrew.

2.1 History

The Hebrew language was used widely in Israel starting in the 10th century BCE (1000 BC to 901 BC) to 4th century CE. During this period, the language evolved and was even lost for 1400 years. However, the Hebrew used today and the Hebrew from the 10th century BCE are very different. We briefly tell the evolution story of the language.

Archaeological findings of Yossi Garfinkel suggest that Hebrew was used over 3,000 years ago. After many exiles, the Jews stopped using Hebrew language as the spoken language. Particularly, the Neo-Babylonian Empire conquering of the ancient Kingdom of Judah in the 6th century BCE, resulted in some loss of spoken Hebrew and the use of Aramaic. Once Cyrus the Great leaders of the Greeks conquered Babylon, Jews began to speak Hebrew once again but alongside Aramaic. In fact, one of the most important books to the Jews, the Gemara, is written in Aramaic. At some point, the common language in Israel was Greek. As Jews were dispersed further from Israel, languages used by Jews began to change even more. During the Middle Ages, one of the greatest Jewish Scholars,

200 Maimonides, actually wrote commentary mostly in
201 Arabic. In Europe, Ashkenazi Jews mainly spoke
202 Yiddish, which actually uses Hebrew characters,
203 and only read Hebrew for religious practices.

204 In the 19th century, Haskalah (Enlightenment)
205 and Zionist movement began to gain popularity as
206 Jews longed to return to their historical home land.
207 As a result, Eliezer Ben-Yehuda began to revive
208 the Hebrew language to create a usable modern
209 language. Since Biblical and Ancient Hebrew are
210 so old, they do not even contain words for mod-
211 ern items. He believed that the revival of the He-
212 brev language could unite all Jews in Israel. Us-
213 ing grammar and methods from many other lan-
214 guages, but specifically other semetic languages,
215 Ben-Yehuda would create what is known as Mod-
216 ern Hebrew. It is the only language that has been
217 successfully revived dead language in history. To-
218 day, approximately 9 million people speak Hebrew
219 around the world. Of these 9 million, approxi-
220 mately 5 million are native Hebrew speakers.³

221 **2.2 Hebrew Language 101**

222 Since Hebrew is spoken by few and differs from En-
223 glish substantially, we provide an overview of the
224 basic characteristics of Hebrew to help the reader
225 better understand the problem at hand. The most
226 noticeable difference between Hebrew and English
227 is the alphabet. The Hebrew alphabet is unique to
228 the language⁴ and contains 22 letters (consonants,
229 no vowels as explained later on), five of which
230 use different forms when placed at the end of a
231 word. Some consonants can be modified by the
232 addition of a dot to represent a vowel, as an aid to
233 children and those learning Hebrew as a foreign lan-
234 guage. Modern Hebrew text, both handed-written
235 and printed, consists of consonants, spaces and
236 western punctuation. Hand-written Hebrew does
237 not join the letters, unlike Arabic and most hand-
238 written English. The Hebrew alphabet is depicted
239 in Figure 1. In addition, Hebrew is written hori-
240 zontally from right to left, like Arabic. Books open to
241 that, for English native speakers, would be the last
242 page. These differences create many challenges
243 when working with deep learning methods that are
244 primarily built for the Latin alphabet.

245 Furthermore, Hebrew has *nikud*, 15 diacritical

246 ³Accoring to the March 18, 2013 Israel Hayom Newsletter

247 ⁴Other languages that were once spoken by Jews such
248 as Aramaic and Yiddish uses the same alphabet. Yiddish is
249 actually based on the German Language. Aramaic is almost a
hybrid dialect.

250 signs used to distinguish between alternative pro-
251 nunciations of letters and to represent vowels - the
252 Hebrew *nikud* system is depicted in Figure 2. Dif-
253 ferent combinations of Hebrew letters and Nikkud
254 can create a very diverse writing system. It is worth
255 mentioning that in the past, there was more of a
256 difference when pronouncing certain letters and
257 certain *nikud* charcters. As time passed by, ceratin
258 differences (such as between Alef and Ayin, or be-
259 tween Qamets and Patach) have disappeared, and
260 today they are pronounced almost the same.

261 Hebrew verb grammar is similar to English in
262 that it has past, present and future tenses, condition-
263 als, imperatives and infinitives. It has the active and
264 passive voice and differentiates between transitive
265 and intransitive. There are some minor differences,
266 however, than can lead to incorrect English verb
267 use. For example, Hebrew does not use the copula
268 to be in the present tense as in English, so begin-
269 ners may correctly say sentences such as "I happy
270 today!"

271 Moreover, Hebrew is a much more inflected lan-
272 guage than English. For example, Hebrew has more
273 verb endings, nouns and pronouns vary in form ac-
274 cording to the preposition that precedes them, and
275 adjectives must agree in number and gender with
276 the nouns they modify. Hebrew has the masculine
277 and feminine genders.

278 Unlike English, many Hebrew "stop words" are
279 connected to the root word. For example, 'and' is
280 written by appending the letter *vav* (ו) in front of a
281 word, 'from' is written by appending the letter *mem*
282 (מ) in front of a word, 'to' is written by appending
283 the letter *lamid* (ל) in front of a word, 'in' is written
284 by appending the letter *bet* (ב) in front of a word,
285 etc.⁵ These features make it difficult for neural
286 networks to learn the langugae and its vocabulary
287 for proper translation.

288 **2.3 Biblical Hebrew Versus Modern Hebrew**

289 The core vocabulary of Modern Hebrew comes
290 from Biblical Hebrew. A large majority of simple
291 nouns and verbs – at least those that existed in the
292 days of the Bible – are the same. But there are
293 also many words in the Hebrew Bible that are not
294 used in Modern Hebrew, and some that are used
295 differently. Moreover, the grammatical structures
296 are different, and so are the tenses, the words' order
297 and the use of *nikud*.

298 ⁵Adding to the challenge of the task at hand, there are no
widely used or maintained Hebrew tokenizers, stemmers, or
stop-word cleaners.

Mark	Name	Sound	Hebrew	Trans.	Class	Type
אָ	Qamets	"a" as in aqua	אָ	a	Long	A-Type
ָ	Patach	"a" as in aqua	ָ	a	Short	
ְ	Chateph Patach	"a" as in aqua	ְָ	a	Reduced	
ֱָ	Qamets Hey	"a" as in aqua	ֱָהָא	ah	Long	
ֵָ	Tsere	"ei" as in eight "e" as in they	ֵָרִי	ei / e	Long	E-Type
ֶָ	Segol	"e" as in red	ֶָנוֹלָ	e	Short	
ַָָ	Chateph Segol	"e" as in red	ַָָנוֹלָ	e	Reduced	
ָָָ	Tsere Yod	"ei" as in eight	ָָָוָדָ	ei	Long	
ְָָָ	Segol Yod	"ey" as in obey	ְָָָוָדָ	ey	Long	I-Type
ִָָָ	Chireq	"ee" as in green	ִָָָרָקָ	i	Short	
ְִָָָ	Chireq Yod	"ee" as in green	ְִָָָוָדָ	i	Long	
ָָָּ	Cholem	"o" as in yellow	ָָָּלָםָ	o	Long	O-Type
ְָָָּ	Chateph Qamets	"o" as in yellow	ְָָָּלָםָ	o	Reduced	
ָָָֻ	Qamets Chatuph	"o" as in yellow	ָָָֻלָםָ	o	Short	
ָָָׁ	Cholem Vav	"o" as in yellow	ָָָׁוָוָ	o	Long	
ָָָּ	Qibbutz	"u" as in blue	ָָָּבָרָןָ	u	Short	U-Type
ְָָָּ	Ehureq	"u" as in blue	ְָָָּבָרָןָ	u	Long	
ָָָֻ	Sheva'	Vocal: short "e" Silent: no sound	ָָָֻבָרָןָ	e or '	(vocal) Short	

Figure 2: The Hebrew Nikkud and English Pronunciation. (Source: Hebrew4Christian.)

3 Computational Resources

Our models were trained using the free version of Google Colab. Google Colab provides GPUs based on supply and demand, so there is no guarantee that models are trained on the same GPU. At times one can be assigned a Nvidia P100, T4, P4, or K80. Normally, a the worse performing GPUs are assigned for training, since we do not have the payed version. In addition, sessions can be terminated at any point and are limited to 12 hours per day, blocking users if excessive use is detected. As a result, we made the best of our resources. We wanted to complete more experiments with additional models, but were unfortunately prohibited. In addition, we attempted to use pretrained models and prebuilt models so that we do not exhaust our limited computation time. To implement our models we make using of PyTorch and HuggingFace.

4 Dataset

4.1 Collection and Cleaning

Despite the accurate translation, there are many challenges with preprocessing and cleaning the text to use pre-existing models. Since there are nearly no common Hebrew text cleaner Python libraries we carried out our dataset construction and cleaning ourselves by solving the following issues:

1. We need to remove non-translation characters such as editors commentary, numbers, etc.
2. We need to handle strange characters such as multiple attached spaces and replace unknown characters with better representations such as replacing “\u2009” with whitespace.
3. In some cases, we need to reverse the order of sentences since Hebrew goes from right to left.
4. We need to remove *nikud*: transforming יִשְׂרָאֵל to יִשְׂרָאֵל. *Nikud* is also an area for which we experiment because we are curious to see if any of our models can represent and learn *nikud* in addition to semantics and structure.

4.2 Data Statistics

Our dataset, which includes Biblical Hebrew sentences and their Modern translations, includes a total of 9222 examples that were collected from 50 Bible sections. It is summarized in Table 1. As previously discussed in Section 2.2, one major difference between Biblical Hebrew and Modern Hebrew is the use of *nikud*, which is the vowel symbol system used to specify pronuciations of words. For example, the word “Israel” is traditionally written as יִשְׂרָאֵל in contrast to today’s writing which strips the *nikud* to obtain יִשְׂרָאֵל. Therefore, we provide two datasets, one of plain Hebrew and the other of Hebrew with *nikud*.

An expected difference between the two datasets is the number of unique characters each contains. This difference is caused by the *nikud* symbols as they are counted as unique characters in addition to conventional letters. For example, יִשְׂרָאֵל has 10 characters while יִשְׂרָאֵל has only 5. We note that this property leads to single words being represented in multiple ways (via variation in their *nikud*), and results in a far greater number of unique words appearing in the *nikud* dataset. For example, the number of unique words in the no-*nikud* Train Set Source is approximately 40% of the *nikud* counterpart.

We also observe that the fraction of source words that are only seen once in the datasets is greater in the *nikud* dataset. On average, 50% of the words in the plain dataset are seen once while 60% of the *nikud* dataset are singular. The fact that the one-appearance word fraction is approximately equal between the Target datasets sheds light on the lesser use of *nikud* in modern Hebrew. Today it is more common to be using the same *nikud* symbols re-

	Test	Train	Val	Test <i>nikud</i>	Train <i>nikud</i>	Val <i>nikud</i>	
Source Max Sentences Length	43	41	32	43	41	32	450
Target Max Sentences Length	42	41	36	42	41	36	451
Source Sentences	922	7379	921	922	7379	921	452
Target Sentences	922	7379	921	922	7379	921	453
Source Unique Word Count	4380 (66%)	13762 (40%)	4363 (64 %)	7765 (81%)	32915 (50%)	7800 (78%)	454
Target Unique Word Count	4960 (65%)	15872 (40%)	4976 (64 %)	5270 (67%)	17554 (41%)	5276 (65%)	455
Source Word Count	12301	100855	12664	12464	102309	12835	455
Target Word Count	13692	111740	14107	13970	114063	14375	456
Source Character Count	58588	480209	60376	111155	910549	114587	456
Target Character Count	67444	550744	69495	113611	926910	116970	457

Table 1: Statistics for our Biblical Hebrew to Modern Hebrew translation dataset. We report statistics for both the clean dataset (no *nikud*) and the raw dataset (with *nikud*). In parentheses next to Unique Word Counts, we add the fraction of these words that appear only once within their set. We note that the *nikud* sets generally have a higher single-word fractions and more unique words.

gardless of the context of words. This insight might suggest that the *nikud* source dataset contains more contextual information in the words themselves and might be better for learning. Section 4.3 further examines the similarity between the source and the target sentences.

4.3 Base Line and Translation Similarity

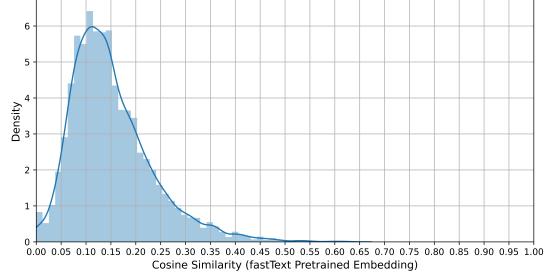
To provide a quantitative sense of the differences of Modern and Biblical Hebrew, we provide two measures of similarity for the source and target: cosine similarity between Hebrew embeddings from pre-trained fastText model for Hebrew and the BLEU score for non-translated sources and targets. The latter acts as a baseline BLEU score for our translation task—we want to perform better than this value.

There is a pretrained fastText 300 (Mikolov et al., 2018) Hebrew word embedder from Facebook that is implemented in Pytorch. We use this pretrained model to create feature representations for our source and target sentences. With these feature representations, we can compare the two sentences with a number of metrics. Specifically, we use Cosine Similarity defined as

$$\text{CosineSimilarity}(u, v) = 1 - \frac{u \cdot v}{\|u\|_2 \|v\|_2} \quad (1)$$

where u and v are vector embeddings. Figure 3 depicts the distribution of the cosine similarity measures for each source and target sentences. We see that the embeddings are not drastically different from each other as would be the case for two different languages. The mean cosine similarity is 0.1516 with a minimum of just about zero. One should note that the source translates similar texts so that the Biblical Hebrew is more understandable

Figure 3: Distribution of Cosine Similarity of sentences with pretrained Hebrew fastText Embedding

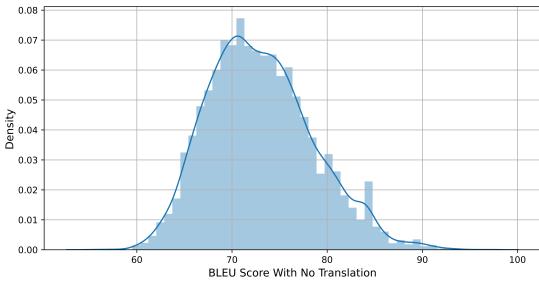


to those who speak Modern Hebrew. If a sentence in the Bible does not need translating because the Modern Hebrew translations is the same, translators may still translate the sentence with a different wording. The translator may be reluctant to translate a sentence exact as written in biblical text even if that is a modern translation.

Even though the cosine similarities of the embeddings are fairly low, the BLEU scores are higher, indicating that this translation task may not be straightforward not because of the difference in words between the sources and targets but because of nuanced word choice. Figure 4 depicts the distribution of BLEU scores for the source and the target sentences. The maximum BLEU score is 97.40, which is extremely high and caused by translations taking on the same form as their sources. The mean of the distribution of BLEU scores is 73.03. The pretrained Hebrew fastText model does poorly in capturing the biblical Hebrew embedding, and the BLEU scores shed light on similar word choices in the translations to the sources.

500
501
502
503
504
505
506
507
508
509
510
511
512
513
514

Figure 4: Distribution of BLEU Scores Between
Source and Target



5 Models

515
516
517
518
519
520
521
522
523
524
525

There are a number of models used for experiments and comparison in the translation task described in this paper. These models are mainly encoder-decoder models, but we apply a number of variations. The first model we apply is a typical Seq2Seq model with attention. The second model is a fine-tuned pretrained Hebrew-to-Hebrew model with 12 encoders and 12 decoders from Facebook’s multilingual model (Fan et al., 2020). In addition, we attempt to use a Seq2Seq model with a copy mechanism with little success.

5.1 End-To-End Transformer Model (ETET)

526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542

In sequence-to-sequence problems (for which a neural net transforms a given sequence of elements, such as the sequence of words in a sentence, into another sequence), the initial solutions are based on the use of RNNs in an encoder-decoder architecture. These architectures are limited when working with long sequences, as their ability to retain information from the first elements is lost when new elements are incorporated into the sequence. In the encoder, the hidden state for each step is associated with a specific word from the input sentence. Therefore, if the decoder only accesses the last hidden state of the encoder, it will lose relevant information about the first elements of the sequence. To deal with this limitation, the attention mechanism is created.

543
544
545
546
547
548
549

The attention mechanism is based on the idea that instead of paying attention to the last state of the encoder, in each step of the decoder we look at all the states of the encoder, accessing information about all the elements of the input sequence, calculating a weighted sum of all the past encoder states. This allows the decoder to assign greater weight

550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565

or importance to certain elements of the input for each element of the output.

556
557
558
559
560
561
562
563
564
565

But this approach continues to have an important limitation—each sequence must be treated one element at a time, which is very time consuming and computationally inefficient.

566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595

The transformer model extracts features for each word, using a self-attention mechanism to figure out how important all the other words in the sentence are with respect to that original word. As no recurrent units are used for this process, it is very parallelizable and efficient. Our model is Transformers based, and is based on Vaswani et al. (2017) paper “Attention is all you need” and on Klein et al. (2017) (Harvard’s NLP) implementation of it.

566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595

As seen in Figure 5, there is an encoder model on the left side and the decoder on the right one. The encoder maps an input sequence of symbol representations (x_1, \dots, x_n) to a sequence of continuous representations $z = (z_1, \dots, z_n)$. Given z , the decoder can then generate an output sequence (y_1, \dots, y_n) of symbols, one element at a time. At each time step, the model consumes the previously generated symbols as an additional input to generate the next. The Transformer follows this overall architecture using stacked self-attention and pointwise, fully connected layers for both the encoder and decoder.

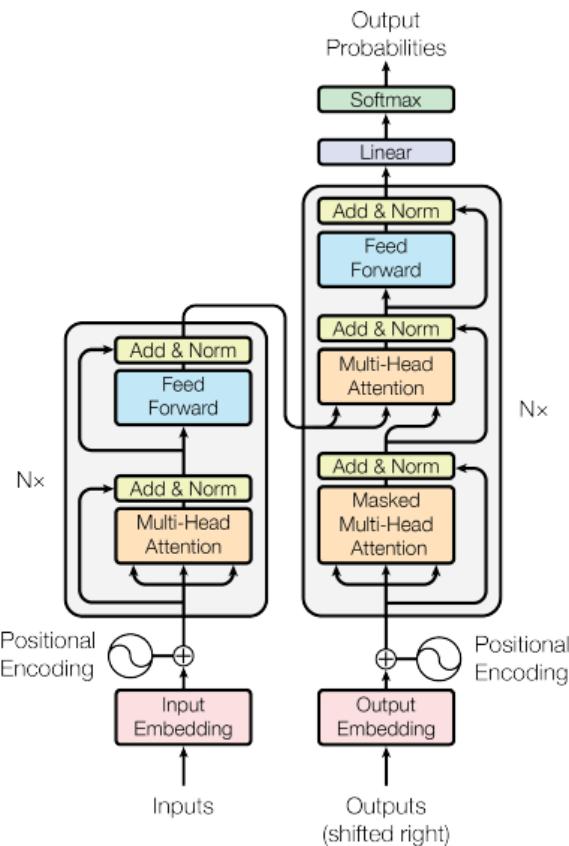
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
599

The encoder is composed of a stack of 6 identical layers. Each layer has two sub-layers. The first sub-layer is a multi-head self-attention mechanism, and the second is a simple, position-wise fully connected feed-forward network. Residual connections are employed around each of the two sub-layers followed by layer normalization. The decoder is also composed of a stack of 6 identical layers. In addition to the two sub-layers in each encoder layer, the decoder inserts a third sub-layer, which performs multi-head attention over the output of the encoder stack. Similar to the encoder, we employ residual connections around each of the sub-layers, followed by layer normalization. Figure 6 presents the attention mask that shows the position each target word (row) is allowed to look at (column). Words are blocked for attending to future words during training.

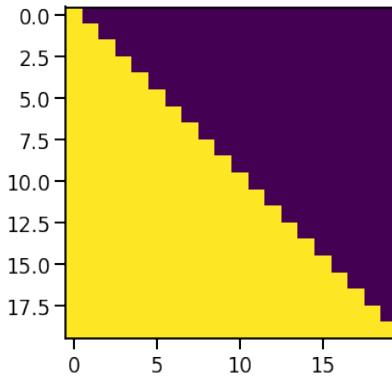
596
597
598
599
599

As explained before, the attention maps a query and a set of key-value pairs to an output, where the query, keys, values, and output are all vectors. The output is computed as a weighted sum of the

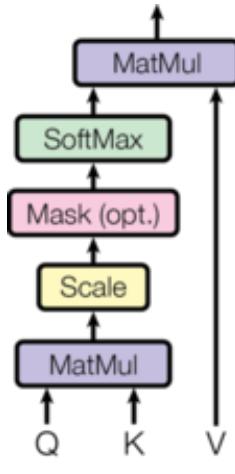
600
601
602 Figure 5: Seq2Seq Attention Model (Source: Vaswani
603 et al. (2017))
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635



636
637
638
639
640
641
642
643
644
645
646
647
648
649 Figure 6: Attention Mask (Source: Vaswani et al.
(2017))



650
651
652 Figure 7: This Figure depicts the scaled dot-product
653 attention. (Source: Vaswani et al. (2017))
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699



values, where the weight of each value is computed by a compatibility function of the query with the corresponding key.

The particular attention used here is a scaled dot-product attention and is depicted in Figure 7. The input consists of queries and keys of dimension d_k , and values of dimension d_v . Then dot products of the query with all keys are computed, each is divided by $\sqrt{d_k}$, and a softmax function is applied to obtain the weights on the values. The use of multi-head attention allows the model to jointly attend to information from different representation subspaces at different positions. Due to the reduced dimension of each head, the total computational cost is similar to that of single-head attention with full dimensionality.

In addition to attention sub-layers, each of the layers in the encoder and decoder contains a fully connected feed-forward network, which is applied to each position separately and identically. This consists of two linear transformations with a ReLU activation in between.

Since tokenization for Hebrew is not ideal, a custom tokenizing function based on “Pytorch’s NLP From Scratch” (Robertson, 2021) was created. Each word was represented as a one-hot vector, and in order to keep track of one-hot encoding a mapping of words to indexes was created and a word counter to use later to replace rare words.

Learned embeddings—the usual learned linear transformation and softmax function—are used to convert the input tokens and output tokens to vectors of the dimension of the model. In addition, they are used to convert the decoder output to pre-

700 dictated next-token probabilities.
701

702 Since the model contains no recurrence and no
703 convolution, we must provide the model information
704 about the relative or absolute position of the
705 tokens in the sequence for the model to make use
706 of the order of the sequence. To address that issue,
707 “positional encodings” were added to the input
708 embeddings which are located at the bottoms of
709 the encoder and decoder stacks. There are many
710 choices of positional encodings, we chose to use
711 sine and cosine functions of different frequencies.

712 5.2 Fine-Tuned Multilingual Transformer 713 Model (FTMT)

714 Recently, Facebook has designed a state-of-the-art
715 multilingual translation model which focuses on
716 non-english centric translations (Fan et al., 2020).
717 The team created “a true Many-to-Many multilingual
718 translation model that can translate directly
719 between any pair of 100 languages” (Fan et al.,
720 2020). This model uses attention and transform-
721 ers with 12 encoder and 12 decoder layers for a
722 total of 1.2 billion parameters. Since this model is
723 designed for multilingual translaiton, it requires a
724 special prefix language ID token for both the source
725 and target text. Additional details on the model con-
726 struction and implementation can be found in Fan
727 et al. (2020).

728 With our aim creating a Hebrew-to-Hebrew task,
729 we choose to use this model and take advantage
730 of the state-of-the-art decoding and encoding mod-
731 els and pretraining. A significant factor that lead
732 to our choice is the uniqueness of this pretrained
733 Hebrew encoding *and* decoding model. Tokenization
734 that support Hebrew inputs are uncommon
735 and not available for many other models. This
736 pretrained model, which is provided by Facebook
737 and implemented by HuggingFace, allowed us to
738 achieve promising results with minimal domain-
739 specific tranining.⁶ This model is pretrained with
740 7.5 billion training sentences from 100 languages.
741 Given we have under 10,000 Hebrew sentences, we
742 fine-tune this model to be able to translate from
743 Biblical Hebrew to Modern Hebrew.

744 To train this model, we use an AdamW optimizer
745 (Loshchilov and Hutter, 2017) with a learning rate
746 of 0.00001. A number of different learning rates
747 were tested and ranked by the performance on our
748 validation set. We found that with this learning rate

749 ⁶https://huggingface.co/transformers/model_doc/m2m_100.html

750 and optimizer, the model achieved very promis-
751 ing performance on the validation set after just 4
752 epochs. The loss function used is a negative log
753 likelihood loss implemented in HuggingFace as
754 NLLossBackward. In total, we train this model
755 for only 10 epochs because of our limited compu-
756 tational resources described in Section 3.

757 5.3 CopyNet

758 In simple terms, a Copy Mechanism is a mech-
759 anism on the decoder of a Seq2Seq model used
760 for translation that copies words from the source
761 sentence and uses them in the output sentence. In-
762 tuitively, we believed that the CopyNet architecture
763 could achieve promising performance since many
764 of the words in the Modern Hebrew and Biblical
765 Hebrew are the same (see Fig. 4 for an illustration).
766 To implement this model, we followed the imple-
767 mentation described in Gu et al. (2016) and even
768 attempted to modify a number of implementa-
769 tions of CopyNet found throughout the web. While
770 performance was able to improve from the mean
771 BLEU of 0.000265, it never reached a score above
772 0.1. A number of different training procedures
773 were used, but the model seemed to perform poorly.
774 One reason for this was the inadequate GPU mem-
775 ory to hold the copy vocabulary size needed for
776 Hebrew which has many variation of words as de-
777 scribed in Section 2.2. Because of this lackluster
778 performance, we decided on focusing on the first
779 two methods that could be trained with our limited
780 resources.

781 6 Experimental Evaluation

782 To determine the efficacy of our methods for trans-
783 lating between Biblical Hebrew and Modern He-
784 brew we conducted three classes of experiments.
785 In this section we discuss our experimental evalua-
786 tion of the training process of our two models on
787 plain and *nikud* Hebrew inputs, present quantitative
788 translation results, and finally qualitatively discuss
789 the quality of translations of free-input text.

791 6.1 Experiments

792 The first and most basic experiment we conduct
793 is an observation of the evolution of the training
794 losses reported when training our models. When
795 fitting the parameters of our models to fit the train-
796 ing samples, we interpret a decreasing trend in the
797 losses as successful learning. We train and eval-
798 uate our models on the data specified in Table 1,

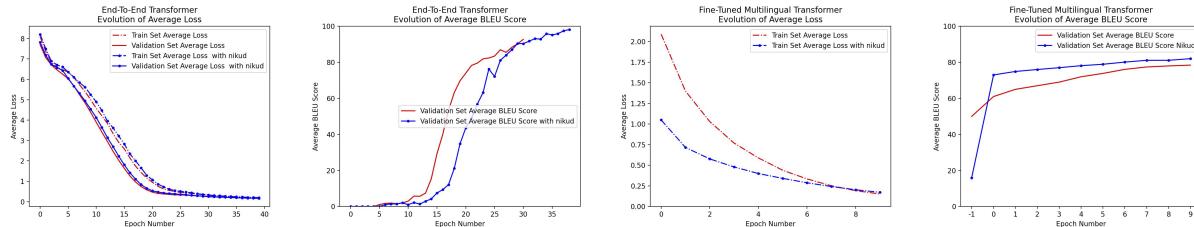


Figure 8: Experimental results. We illustrate our models’ ability to learn Hebrew-to-Hebrew translations from labeled data and provide a quantitative evaluation of the quality of our translations. Our End-To-End Transformer (ETET) and Fine-Tuned Multilingual Transformer (FTMT) models both show a decreasing loss evolution during training and an increasing associated BLEU score. Given that the BLEU scores reported are generated from previously unseen test data, these properties suggest that our models can learn effective translations between Biblical and Modern Hebrew texts. Values at epoch -1 signifies the BLEU score achieved by FTMT before any fine-tuning took place.

which includes our Hebrew-to-Hebrew dataset and a variant of that dataset that includes *nikud*.

To complement the analysis of the learning in our models, we further wish to understand if the learning that took place has been successful for the task at hand (Hebrew-Hebrew translation). To this end, we use the BLEU score metric to quantify the quality of the outputted translations. We compare our models’ translations of Biblical Hebrew sentences with reference translations to recover the BLEU scores. Finally, to verify that our learning pipeline can also generalize to novel inputs and is robust to adversarial queries, we test our models’ performance qualitatively on free-text inputs. We evaluate the outputs of our models that are generated for inputs of Hebrew sentences from different modalities. By using from the Bible, from poetry, and simply arbitrary sentences, we wish to understand qualitatively if our models are robust.

The following sections detail the results of our evaluation. Section 6.1.1 discusses the training-loss evolution of our models, Section 6.1.2 presents our BLEU score results, and Section 6.1.3 covers our observed qualitative results from translating free-text inputs.

6.1.1 Training Translators

As illustrated in Fig. 8, we observe that both of our models have successfully learned to reduce their training loss when trained on Biblical Hebrew sentences alongside their Modern Hebrew translations as supervision. The loss evolution that appear when training our End-To-End Transformer model (ETET) shows a very similar rate of loss change over time for the Train set and the Validation set when training on plain Hebrew and *nikud* Hebrew. It could be that ETET showcased a similar loss

evolution for both variants of Hebrew because its randomly initialized state was sufficiently unbiased to treat the different datasets as realtively similar.

On the other hand, our Fine-Tuned Multilingual Transformer (FTMT) model did behave differently when trained on plain and *nikud* Hebrew inputs. Perhaps counter-intuitively, FTMT showed a higher initial loss when trained on plain Hebrew despite being pretrained on this type of language.

6.1.2 Translation Quality

Through evaluating the BLEU scores of the translated sentences with reference translations we conclude that both of our models have learned to successfully perform Biblical Hebrew translation. Our model ETET achieved BLEU scores of 92% and 98% for plain and *nikud* Hebrew repsectively. Such high scores could point to lack of diversity in our dataset. Given that our target translations are all from the same source it could be that our Test and Train sets have sufficiently similar statistical occurrence of sentences and translations to facilitate near-perfect translation.

We gain insights into the effects of pretraining on large amounts of data through the translation quality evaluation for the FTMT model. We first note that the initial performance score FTMT achieves without *any* Fine-Tuning is 50% for plain Hebrew and 15% for *nikud*. This an expected result given that FTMT was pretrained on plain Hebrew. We additionally see the main two benefits of pretraining: fast adaptation to new data and robustness to overfitting. The first is seen in the dramatic jump in the BLEU scores for both the *nikud* Hebrew trained model and the second is seen in the monotonic increase of the scores in the following epochs for both training modalities. It seems as if the pretrained

900		גואשין גיא אקלים לא טהור וαι סורי	סאית בא אלות תא מאיין	סאית בא אלות תא מאיין	סִירְקָוְלָרְאַנְטָרְבָּה	הה גונול או רור או מגה	גַּעֲמֵי רַעֲמֵי, גַּעֲמֵי צַדְקֵי, וְיַם.	הה ים רומי, גני ים צדוקים.	950
901		In the beginning God created the heaven and the earth.	In the beginning God created the heaven and the earth.	The biggest mountain is the tallest mountain.	The biggest mountain is the tallest mountain.	You are beautiful my wife, you are beautiful your eyes like doves.	You are beautiful my wife, you are beautiful your eyes like doves.		951
902									952
903									953
904	FTMT	במהלך נסחנות השילוח נאכון תרילוח באנט תרילוח ...	במהלך נסחנות השילוח נאכון תרילוח באנט תרילוח ...	אתם באה אלותם שמשם תא מאיין ואה	אתם באה אלותם שמשם תא מאיין ואה	הה גונול או רור או מגה	הה גונול או רור או מגה	הה גונול או רור או מגה	954
905		The kettle come the climbing in their climbing (repeats).		At the beginning God created the heaven the earth and the	The Cohen (repeats).	The big mountain is the tallest mountain.	אֶתְהַדְּרָה שְׁלִיחָה וְלִזְנָה אֶתְהַדְּרָה שְׁלִיחָה וְלִזְנָה שְׁלִיחָה אֶתְהַדְּרָה	אֶתְהַדְּרָה שְׁלִיחָה וְלִזְנָה אֶתְהַדְּרָה שְׁלִיחָה וְלִזְנָה שְׁלִיחָה אֶתְהַדְּרָה	955
906							רִיחָה שְׁלִיחָה Beautifull Wine Sheddin Ideas The Wine Sheddin...		956
907	FTMT Nikud	תאָגָת אֱלֹהִים וְאֶתְהַדְּרָה שְׁלִיחָה וְלִזְנָה	תאָגָת אֱלֹהִים וְאֶתְהַדְּרָה שְׁלִיחָה וְלִזְנָה	בְּקַבְּשָׂה נִמְצָא אֶלְמָן כְּנָפָת	בְּקַבְּשָׂה נִמְצָא אֶלְמָן כְּנָפָת	הַרְמָנוֹת הַרְמָנוֹת הַרְמָנוֹת	הַרְמָנוֹת הַרְמָנוֹת הַרְמָנוֹת	רִיחָה שְׁלִיחָה My wife you are beautiful your eyes like doves.	957
908		Agreement God created the heavens and the earth		Please in them there were gods and the name		The husband son husband is.	רִיחָה שְׁלִיחָה My wife you are beautiful your eyes like doves.	רִיחָה שְׁלִיחָה Their place companies companies...	958
909									959
910	ETET	No output.		To parts he'll divide it from the skin of the land and the land and the land...	No output.	No output.	No output.	No output.	960
911									961
912	ETET Nikud	לְכַפֵּר לְכַפֵּר לְכַפֵּר לְכַפֵּר לְכַפֵּר לְכַפֵּר	No output.	No output.	No output.	No output.	No output.	No output.	962
913		You have sinned for kettle and butcher its family and on himself will forgive the land will be given.							963
914									964
915									965
916									966
917									967
918									968
919									969
920									970
921									971
922									972
923									973
924									974
925									975
926									976
927									977
928									978
929									979
930									980
931									981
932									982
933									983
934									984
935									985
936									986
937									987
938									988
939									989
940									990
941									991
942									992
943									993
944									994
945									995
946									996
947									997
948									998
949									999

Figure 9: Three sample sentences (with *nikud* and without) translated with out two models. English translations are added next to the Hebrew texts. We observe that the Fine-Tuned model FTMT produces more robust translations when compared to the End-to-End model ETET, which is limited in its vocabulary and training. Unsurprisingly, free-text that include *nikud* are generally translated better by models that were trained on text with *nikud*.

parameters were positioned sufficiently well in the parameter space to allow fast refinements instead of a convergence to an unfavorable minima.

6.1.3 Free Text Translation

We report qualitative results on translating free-text inputs using all four variants of our models (ETET, ETET-*nikud*, FTMT, FTMT-*nikud*). Example sentences taken from Biblical and poetry sources alongside Modren Hebrew inputs are included in Fig. 9. Sensibly, the models that were trained on plain Hebrew struggle with recovering good translations for *nikud* inputs and vice versa. Given that *nikud* is encoded as extra unicode characters in words this result is expected.

Focusing on ETET, we observe that a large fraction of the inputs failed to produce an output. The reason for these failures is the small vocabulary size of ETET that is only extracted from our dataset and leads to novel input words being unable to be tokenized. This is a severe drawback of our End-To-End model.

On the other hand, we observed that FTMT (both trained on plain language and *nikud*) has successfully handled novel inputs and produced reasonable outputs. Such robustness can be attributed to the extensive pretraining this model has undergone and to the newly introduced Hebrew tokenization that

was released with the M2M100 model. This result suggests that our FTMT model specialized to perform well on the translation task (achieving a BLEU score of 82%) without losing the ability to handle arbitrary inputs.

7 Discussion

Through evaluating our two models on the task of Ancient Hebrew to Modern Hebrew translation we have observed the important role large datasets play in the construction of effective language models. In our experiments we have shown that learning a language model end-to-end solely from a small dataset yields a model achieving impressive translation BLEU scores that fails to adapt to samples outside the dataset. Our ETET model overfit to the nature of the data in the dataset and suffered from a small vocabulary.

Furthermore, the role of pretraining has become increasingly clear in the development of our FTMT model. The adaptation of an established language model to a specific task was facilitated for the FTMT model’s precomputed initial parameters. Since both our transformer models are very large, it is intuitive to believe that small amounts of training data would not be able to excite all of their parameters sufficiently to induce meaningful learning.

1000 Reliance on past training is important.

1002 8 Conclusion

1003 In this work we have presented two methods to
1004 successfully perform Biblical Hebrew to Modern
1005 Hebrew translation, achieving BLEU scores of
1006 82% and 98%. We have presented the nuances
1007 of this translation problem alongside the origins
1008 and uniqueness of the Hebrew language and pro-
1009 vided insights into potential factors contributing to
1010 the successes and drawbacks of our methods.

1011 In future works, we will like to further exper-
1012 iment with a Copy Mechanism to see if it can work.
1013 In addition, our corpus only includes the Torah,
1014 not the prophets or scrolls. With these additional
1015 texts, we can hopefully have more trianing data and
1016 improve the models performance. In addition, a
1017 website that allows for translation of biblical He-
1018 brew in a child-friendly format will allow students
1019 all over Israel to use our reseach to improve thier
1020 learning.

1022 References

1023 Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi
1024 Ma, Ahmed El-Kishky, Siddharth Goyal, Man-
1025 deep Baines, Onur Celebi, Guillaume Wenzek,
1026 Vishrav Chaudhary, Naman Goyal, Tom Birch, Vi-
1027 taliy Liptchinsky, Sergey Edunov, Edouard Grave,
1028 Michael Auli, and Armand Joulin. 2020. **Bey-**
1029 **ond english-centric multilingual machine transla-**
1030 **tion.** *CoRR*, abs/2010.11125.

1031 Jiatao Gu, Zhengdong Lu, Hang Li, and Victor
1032 O. K. Li. 2016. **Incorporating copying mech-**
1033 **anism in sequence-to-sequence learning.** *CoRR*,
1034 abs/1603.06393.

1035 Salima Harrat, Karima Meftouh, and Kamel Smaili.
1036 2019. **Machine translation for arabic dialects**
1037 **(survey).** *Information Processing Management*,
1038 56(2):262–273. Advance Arabic Natural Language
1039 Processing (ANLP) and its Applications.

1040 Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senel-
1041 lart, and Alexander M. Rush. 2017. **Opennmt:**
1042 **Open-source toolkit for neural machine translation.**

1043 Seong-Bae Park Laith H. Baniata, Seyoung Park. 2018.
1044 A neural machine translation model for arabic di-
1045 alects that utilizes multitask learning (mtl). *Compu-*
1046 *tational Intelligence and Neuroscience*, 2018.

1047 Ilya Loshchilov and Frank Hutter. 2017. **Fixing**
1048 **weight decay regularization in adam.** *CoRR*,
1049 abs/1711.05101.

Tomas Mikolov, Edouard Grave, Piotr Bojanowski,
Christian Puhrsch, and Armand Joulin. 2018. **Ad-**
1050 **vances in pre-training distributed word representa-**
1051 **tions.** In *Proceedings of the International Confer-*
1052 *ence on Language Resources and Evaluation (LREC*
1053 *2018)*.

Sean Robertson. 2021. **Nlp from scratch: Translation**
1055 **with a sequence to sequence network and attention.**

Reut Tsarfaty, Amit Seker, Shoval Sadde, and Stav
Klein. 2019. **What’s wrong with hebrew nlp? and**
1056 **how to make it right.** *CoRR*, abs/1908.05453.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz
Kaiser, and Illia Polosukhin. 2017. **Attention is all**
1061 **you need.**

Wei Xu, Alan Ritter, Bill Dolan, Ralph Grishman, and
Colin Cherry. 2012. **Paraphrasing for style.** In *Pro-*
1064 *ceedings of COLING 2012*, pages 2899–2914, Mum-
1065 *bai, India. The COLING 2012 Organizing Commit-*
1066 *tee.*

Marcos Zampieri, Preslav Nakov, and Yves Scherrer.
2020. **Natural language processing for similar lan-**
1069 **guages, varieties, and dialects: A survey.** *Natural*
1070 *Language Engineering*, 26(6):595–612.

1072 Acknowledgments

1074 *Thank you for a wonderful semester!*