

## HTX xData Test Q4

The pretrained model and finetuned model are evaluated via generating a transcription, then comparing their results. The word error rate (WER) of the pretrained and finetuned model are as shown below:

| Model      | WER   |
|------------|-------|
| Pretrained | 0.108 |
| Finetuned  | 0.038 |

As seen from the WER, the finetuned model performs similarly to the pretrained model results when it is evaluated against the libspeech dataset, showing that the finetuned model has learnt the appropriate features in the common voice dataset.

The table below shows some instances in which the finetuned model results did not match with the ground truth, whereas the pretrained model matches:

| Finetuned Text  | Ground Truth Text   |
|---|---|
| BUT THE ENGLISHMAN APPEARED NOT TO<br>ATTAUCH ANY IMPORTANCE TO IT              | BUT THE ENGLISHMAN APPEARED NOT TO<br>ATTACH ANY IMPORTANCE TO IT             |
| THE REFERENCES ARE GOOD IN THE VERY<br>NARROW AREA OF HISTOGRAMH METHODS        | THE REFERENCES ARE GOOD IN THE VERY<br>NARROW AREA OF HISTOGRAM METHODS       |
| AT WASN'T MUCH ANYWAY   | IT WASN'T MUCH ANYWAY   |
| A LITTLE LEMONGRASS SHOULD FRESTEN IT<br>UP                                     | A LITTLE LEMONGRASS SHOULD FRESHEN IT<br>UP                                   |
| THE THE BOY AWOKE AS THE SUN ROSE   | THE BOY AWOKE AS THE SUN ROSE   |
| THE BOY WATCHED AS HIS COMPANION<br>WENT TO HIS HORSE AND WITHDREW A<br>CIMITAR | THE BOY WATCHED AS HIS COMPANION WENT<br>TO HIS HORSE AND WITHDREW A SCIMITAR |
| BUT SHE'S GONE  | BUT SHE IS GONE   |
| RED YOU WOULD HAVE TO HAVE BEEN BORN<br>AN ARAB TO UNDERSTAND HE ANSWERED       | YOU WOULD HAVE TO HAVE BEEN BORN AN<br>ARAB TO UNDERSTAND HE ANSWERED         |
| THE HORIZON WAS TINGED WITH RED AND<br>SUDDENLY THE SIGN APPEARED               | THE HORIZON WAS TINGED WITH RED AND<br>SUDDENLY THE SUN APPEARED              |
| GIME ME MY ROBE   | GIVE ME MY ROBE   |

As we can see, it sometimes makes spelling mistakes, such as “ATTACH” is spelled as “ATTAUCH”, “FRESTEN” instead of “FRESHEN”, or “CIMITAR” instead of “SCIMITAR”. One possible way to remedy such errors could be to include more samples where similar words are included. Another way could be to add a spell-checking stage after the transcription is completed to correct for such errors. 84/4076 samples observe such errors.

The table below shows instances where the pretrained model results did not match with the ground truth, whereas the finetuned model matches:

| <b>Pretrained Text</b>  | <b>Ground Truth Text</b>   |
|---|--|
| ARE YOU GOING TO TALK OR WANT YOU   | ARE YOU GOING TO TALK OR AREN'T YOU  |
| HE MOVED ABOUT INVISIBLE BUT EVERY ONE<br>COULD HEAL HIM                        | HE MOVED ABOUT INVISIBLE BUT EVERYONE<br>COULD HEAR HIM                        |
| WE COULD GET TO THE PYRAMIDS BY TO<br>MORROW SAID THE OTHER TAKING THE<br>MONEY | WE COULD GET TO THE PYRAMIDS BY<br>TOMORROW SAID THE OTHER TAKING THE<br>MONEY |
| I'M TRYING TO THINK OF SOMETHING BEFORE<br>THOSE REPORTES GETBACK               | I'M TRYING TO THINK OF SOMETHING BEFORE<br>THOSE REPORTERS GET BACK            |
| WE DON'T HAVE TO GIVE UP OUR CLAP   | WE DON'T HAVE TO GIVE UP OUR CLUB  |
| CAN I GET A WHAT WHAT   | CAN I GET A WOOT WOOT  |
| PEOPLE SAW ME COMING AND WELCOMED<br>ME HE SOUGHED                              | PEOPLE SAW ME COMING AND WELCOMED<br>ME HE THOUGHT                             |
| THE SIMMON BLUE THAT DAY AS IT HAD<br>NEVER BLOWN BEFORE                        | THE SIMUM BLEW THAT DAY AS IT HAD NEVER<br>BLOWN BEFORE                        |
| THE FOLLOWING NIGHT THE BOY APPEARED<br>AT THE ALCHEMIS'S TENT WITH THE HORSE   | THE FOLLOWING NIGHT THE BOY APPEARED<br>AT THE ALCHEMIST'S TENT WITH A HORSE   |
| SAYS FERACHE TO GO THE LINET  | SAYS FOR US TO GO THE LIMIT  |

Like the finetuned model, the pretrained model makes spelling mistakes from the ground truth text, such as “TO MORROW” instead of “TOMORROW”, “REPORTES” instead of “REPORTERS”, or “SOUGHED” instead of “THOUGHT”. In the pretrained model case, it appears to have heard a different word from the ground truth, such as “WHAT” instead of “WOOT”, or “CLAP” instead of “CLUB”. Such errors may need more contextual knowledge to provide the correct transcription. 1145/4076 samples encountered such issues.

The table below shows instances where both finetuned model and pretrained model did not match with ground truth:

| Ground Truth Text   | Finetuned Text   | Pretrained Text   |
|---|--|---|
| AT OTHER TIMES AT A<br>CRUCIAL MOMENT I MAKE IT<br>EASIER FOR THINGS TO<br>HAPPEN             | AT OTHER TIMES AT A CRUCIAL<br>MOMENT YOU MAKE IT EASIER<br>FOR THINGS TO HAPPEN               | AT OTHER TIMES AT A<br>CRUISING MOMENT YOU<br>MAKE IT EASIER FOR THINGS<br>TO HAPPEN      |
| AND THEN THEY WANT THE<br>PERSON TO CHANGE  | AND THEN THEY WANTED THE<br>PERSON TO CHANGE   | AND THEN THEY WANTED<br>THE PERSON TO CHANGE  |
| THE BASKETBALL BOUNCED<br>OFF HIS SHIELD OF TITANIUM  | THE BASKETBALL BOUNCED<br>ON THE SHIELDS OF TITANIUM   | THE BASKETBALL BOUNCED<br>OF HIS SHIELD OF TITANIUM                                       |
| GOTTA BE GENTLE TO SUIT ME  | THET GOTTA BE GENTLE TO<br>SEE ME  | THET GOINTO BE GENTLE TO<br>SIDME   |
| SUPPOSE THERE WAS A<br>SHANE YORK AND HE WALKED<br>INTO THIS OFFICE                           | SUPPOSE THERE WAS A<br>SHANEYOK AND HE WALKED<br>INTO THE LITE OFFICE                          | SUPPOSE THERE WAS A<br>SHAME YOG AND HE<br>WALKED INTO POLISE OFFICE                      |
| I NEED YOU TO BE<br>SPONTANEOUS HE ASKED ME<br>OUT TO DIN DIN                                 | I NEED YOU TO BE<br>SPONTANEOUS HE HAS KAD<br>ME OUT YOUR DIN DIN                              | I NEED YOU TO PEA<br>SPONTENEIOUS HE HAS<br>CATTED ME OUT YOUR<br>DINGDING                |
| THE BOY LOOKED AROUND<br>FOR THE OVENS AND OTHER<br>APPARATUS USED IN<br>ALCHEMY BUT SAW NONE | THE BOY LOOKED AROUND<br>FOR THE OVENS AND OTHER<br>APPARATUS USED IN ALCHEMY<br>BUT SAW NOUNG | THE BOY LOOKED AROUND<br>FOR THE OVENS AND OTHER<br>APPARATUS USED IN ALCEMY<br>TUT SONAN |
| THEY'RE FORMING CLUBS   | THEY 'RE FORMING CLUBS   | THEY ARE FORMING CLUBS  |
| THIS MORNING I FOUND A<br>CALCULATOR TAPED TO MY<br>WII                                       | THIS MORNING I FOUND A<br>CALCULATOR TAPED TO MY<br>WIIT                                       | THIS MORNING I FOUND A<br>CALCULATOR TAKE TO MY<br>WEAT                                   |
| ON TOP OF ALL THAT THE<br>WEEDS KEEP GROWING AND<br>THE GARBAGE HAS TO BE<br>TAKEN OUT        | ON TOP OF ALL THAT THE<br>WIED'S KEEP GROWING AND<br>THE GARBIT HAS TO BE TAKEN<br>OUT         | ON TOP OF ALL THAT THE<br>WEEDS KEEP GROWING AND<br>THE GABBIT HAS TO BE<br>TAKEN OUT     |

From the table, the errors seen are commonly spelling mistakes of words that are less commonly used but do retain their phonetics. For example, the ground truth “BUT SAW NONE” is transcribed as “ BUT SAW NOUNG” for the finetuned model, and “TUT SONAN” in the pretrained model. 690/4076 samples have such error.

Overall, this ML model performs well on the Common Voice dataset, even without much finetuning. While there are some spelling mistakes, most of the transcribed text is largely similar to the ground truth. An increase in the number of training epochs would likely improve the accuracy of the finetuned model. As the finetuned model performs better than the pretrained model, it is likely that the finetuned model learnt from the accent, pace, pitch etc. variability of the common voice dataset. Some experiments to improve the accuracy may include increasing the size of the model, especially if the

variability of the dataset is large. Another experiment could be to replace the CTC decoder with the n-gram language model, which will help with alleviating the spelling mistakes observed in the output.