**Captcha Recognition Project Report**

**Introduction:** In this project, I set out to tackle the challenge of recognizing captchas—those little puzzles meant to keep bots at bay. Instead of relying on just one method, I decided to mix things up by handling both image and audio captchas. For images, I used EasyOCR to extract text, and for audio, I fine-tuned Facebook's Wav2Vec2 model to convert speech into text. My goal was to create a robust, end-to-end solution that could deal with the messy, real-world data that captchas usually throw at us.

**Procedures:**

**Getting the Data Together:** First, I gathered all the captcha data. I had a folder full of images (in PNG, JPG, and JPEG formats) and another folder with audio files in WAV format. I used EasyOCR to extract text from the images and saved that info in a CSV file. Then, I made sure each audio file was paired up with its corresponding image text.

**Cleaning and Preprocessing**

**Image Captchas:** With the image captchas, it wasn't too tricky—EasyOCR did the heavy lifting. I just ran it on each image to pull out the text. The results were then stored neatly in a CSV file so I could easily refer back to them later.

**Audio Captchas:** The audio part was a bit more involved. I needed to ensure that all my audio files were in the same format, so I resampled everything to 16kHz. Then, I tackled noise by using a noise reduction technique that used the first half-second of each recording as a guide for what to filter out. I also trimmed any silence at the beginning or end of the audio, normalized the volume to keep things consistent, and even applied a bit of data augmentation (like adding some background noise, stretching the time, or shifting the pitch) to help the model learn under different conditions.

After processing, I saved these cleaned-up audio waveforms as NumPy arrays. To make sure everything was manageable during training, I split my dataset into training and testing sets (about 80/20) and saved them in small, compressed chunks. This not only helped with memory usage but also made loading the data during training much smoother.

**Building the Model Pipeline**

For the audio-to-text task, I used Wav2Vec2. I loaded a pre-trained processor and model and then created a custom data pipeline:

- **Feature Extraction:** I used the processor to turn the raw audio into a format that the model could understand.
- **Text Tokenization:** I tokenized the image text, ensuring everything was in the right format. I also handled padding carefully—replacing it with a special token (-100) so that the loss function wouldn't be thrown off.
- **Custom Data Collator:** Since each audio file is a bit different, I built a custom collator to pad the audio sequences and their labels properly.
- **Custom Trainer:** I even went as far as creating a custom trainer to compute the Connectionist Temporal Classification (CTC) loss correctly, which is crucial for this kind of audio transcription.

I set up the training with a few sensible choices for hyperparameters: a modest batch size, a careful learning rate, and early stopping to prevent overfitting.

**Evaluation**

To see how well everything was working, I used the Word Error Rate (WER) metric. Essentially, WER tells you how many words were misrecognized by the model. After training, I ran the model on the test set, decoded the predictions, and computed the WER. This helped me understand where the model was getting tripped up.

**Results and Observations**

- **Image Captchas:** EasyOCR did a pretty good job overall, though there were a few hiccups with very distorted images.
- **Audio Captchas:** The audio model's performance was encouraging. The noise reduction and augmentation steps made a noticeable difference, though I did see some errors—mostly due to residual background noise or overly aggressive augmentation that sometimes made the speech less clear.
- **WER:** The final Word Error Rate was a solid indicator of the model's performance. While there's always room for improvement, the metric showed that the approach was on the right track.

**Challenges and Limitations:**

I ran into a few challenges along the way:

- **Data Variability:** The captchas were noisy and inconsistent. Using both noise reduction and augmentation helped the model become more robust.
- **Memory Issues:** Working with large datasets can be a pain. By saving the data in compressed chunks, I managed to keep the training process smooth.
- **Handling CTC Loss:** Implementing the custom loss function was tricky, but writing a custom trainer made it manageable. It was all about getting the tokenization and padding just right.

**Potential Improvements and Limitations**

- **Fine-Tuning:** Training a model specifically on captcha data (rather than a generic pre-trained model) might boost accuracy even further.
- **Preprocessing Tweaks:** There's still a chance to refine the noise reduction process—perhaps by trying out different techniques like spectral subtraction.
- **Better Decoding:** I'm considering more advanced decoding strategies, like beam search, which could improve transcription accuracy.
- **Multimodal Approaches:** Integrating the image and audio components more tightly might yield a more comprehensive solution.

**Current Limitations**

- **Data Quality:** No matter how much you preprocess, some captchas are just too messy, and that can lead to errors.
- **Model Size:** Fine-tuning large models like Wav2Vec2 requires a lot of computing power, which might not be practical for every situation.
- **Real-Time Performance:** While this pipeline works great offline, making it fast enough for real-time use would need further optimization.

**Conclusion**

Working on this project was a rewarding challenge. I managed to build a system that can handle both image and audio captchas using some of the best tools available today. Although there are challenges and limitations—like handling very noisy data and managing computational resources—the overall approach proved promising. I believe with some fine-tuning and further experimentation, this solution could be made even more robust.

In the end, this project isn't just about building a model—it's about learning to handle messy real-world data and continuously finding ways to improve the system. I'm excited to see how these ideas can evolve further, potentially leading to even better captcha recognition in the future.