# Transcription Errors: Techniques for WER Reduction in Natural Language Processing

Alright, folks, gather around! Today we're diving deep into WER, I frequently encounter this metric when evaluating the performance of automatic speech recognition (ASR) and machine translation (MT) systems. The tireless workhorse of NLP evaluation, the metric that helps us navigate the ever-evolving landscape of speech recognition and machine translation.

## WER:

So, what exactly is WER? It's the compass that guides us in the fascinating world of speech recognition and machine translation. It helps us gauge the gap between the symphony of words our systems generate and the flawless score, the reference sequence, we're trying to replicate.

## Challenges of Measuring Performance:

Unlike character-by-character comparison, the recognized words may not always have the same number of words as the reference. This is where WER comes in. It leverages the **Levenshtein distance**, which calculates the minimum number of edits (insertions, deletions, substitutions) needed to transform one sequence into another. However, WER provides a high-level view and doesn't delve into the specifics of the errors.

**Calculating WER:** A Multi-Step Process

## Alignment: Bridging the Gap

Dynamic programming algorithms come to our rescue when the recognized and reference sequences have different lengths. These algorithms meticulously compare the sequences and create a correspondence (alignment) between them. This alignment unveils the specific errors (insertions, deletions, substitutions) made by the system.

For instance, if the reference sequence is "the cat sat on the mat" and the recognized sequence is "the cat sat down on a mat," the alignment would reveal one insertion ("down") and one substitution ("a" for "the").

## Error Enumeration: A Meticulous Process

Armed with the alignment, we can meticulously enumerate the number of errors. We count the number of insertions (extra words the system hallucinates), deletions (missing words), and substitutions (misinterpreted words) in the recognized sequence compared to the reference.

In the above example, we would count one insertion and one substitution.

**Normalization: Setting the Standard**

To facilitate fair comparison across diverse tasks and datasets, we normalize the error count by dividing it by the total number of words in the reference sequence. This yields the WER as a percentage, where a lower value signifies superior performance.

For instance, if the reference sequence has 7 words, a total of 2 errors (1 insertion + 1 substitution) would result in a WER of 28.57% (2 / 7 * 100).

Beyond WER: Unveiling the Nuances of Language

While WER is a valuable tool, it has limitations. It doesn't account for:

- **Semantic Dissonance:** A low WER can mask a functionally nonsensical output.

  For example, the system might recognize "the quick brown fox jumps over the lazy dog" perfectly, but substitute "lazy" with "active." Here, the WER would be 0%, despite a significant alteration in meaning. We, as NLP folks, are actively developing techniques that go beyond word-level comparison to assess semantic equivalence. This involves incorporating semantic similarity measures, natural language inference models, and discourse coherence analysis into our evaluation frameworks.

- **Real-World Robustness:** The complexities of human communication pose significant challenges. Background noise, speaker accents, and idiomatic expressions can all trip up speech recognition systems. Our research endeavors to develop algorithms that are robust to these real-world variations. We employ techniques like noise cancellation, speaker diarization, and language adaptation to improve recognition accuracy in noisy environments, with diverse speakers, and across different dialects.

**Normalization: A Level Playing Field for ASR Evaluation**

Normalization is an essential step in ASR evaluation because it ensures that we're comparing apples to apples. Imagine an ASR system that transcribes "they are" while the human reference has "they're." By default, WER would penalize the system for this difference. But for most purposes, humans wouldn't consider this a true error. Here's where normalization comes in.

A normalizer acts as a bridge, transforming the model's output into a representation that aligns better with the human reference. It can handle various discrepancies that we, as NLP researchers, want to disregard during evaluation. For instance, a normalizer can:

- Expand contractions (e.g., "they're" becomes "they are")
- Remove disfluencies (e.g., stutters, filler words like "um")
- Standardize spellings (e.g., "colour" to "color")

By applying these normalizations, we create a level playing field for ASR evaluation. The focus shifts to errors that truly matter, like missing words or incorrect word choices that impact meaning.

For proper scientific comparison of Speech Recognition models, we must ensure that all aspects of the evaluation process are consistent between different models. Therefore, we should also ensure that we are using a consistent normalizer when evaluating different models, just like how we must use a consistent dataset.

Similarly, if models report metrics using a private normalizer, we must run our own evaluations using a public normalizer to compare models. An open-source normalizer such as the Whisper normalizer are excellent choices for this purpose.

It's important to note that WER is not the only metric used to evaluate speech recognition systems. Other metrics include sentence accuracy, character error rate, and speaker diarization accuracy.

Through my exploration of Word Error Rate in Natural Language Processing, I gained a thorough understanding of its underlying mechanisms, calculation methodologies, and the various typologies employed. Notably, WER is a prevalent evaluation, as evidenced by numerous research papers. Having reviewed several such papers, I gained a comprehensive understanding of the specific WER variants employed for ASR model evaluation.

Building upon this foundation, the discussion will now shift towards exploring how WER is leveraged for evaluation like on specific models such as 'facebookresearch/libri-light', 'wav2vec2.0Large-10h-LV-60k', 'wav2vec2.0with Libri-Light', 'wav2vec 2.0'etc.

Here are the different models listed and their WERs on the Libri-Light test-clean dataset:

- wav2vec 2.0 Large-10h-LV-60k (WER: 2.5%)
- wav2vec 2.0 with Libri-Light (WER: 1.8%)
- wav2vec 2.0 (WER: 4.1%)

**'facebookresearch/libri-light'**

The facebookresearch/libri-light model specifically uses character-level WER. This leverages Character-Level WER as a key metric to assess the performance of ASR systems, particularly when dealing with limited or no supervision. Unlike traditional word-level WER, character-level WER offers a finer-grained analysis, providing a more nuanced understanding of the recognition accuracy.

**Character-Level WER**

Character-level WER measures the disparity between the reference text (ground truth) and the system's output (decoded text) by analyzing individual characters. It calculates the number of insertions, deletions, and substitutions required to transform the decoded text into the reference text. A lower character-level WER signifies a more accurate speech recognition system.

character-level WER becomes particularly valuable. Here's why:

**Captures Granular Errors:** Word-level WER might overlook minor errors like typos or misspellings within a word. Character-level WER identifies these discrepancies, providing a more comprehensive picture of the recognition fidelity.

**Improved Error Localization:** By pinpointing errors at the character level, researchers can identify specific weaknesses in the model and target their improvement efforts more effectively.

Character-level WER is calculated using the Levenshtein distance algorithm. This algorithm determines the minimum number of single-character edits (insertions, deletions, or substitutions) needed to transform one string (decoded text) into another (reference text). The resulting distance is then normalized by the length of the reference text and expressed as a percentage.

A character-level WER of 0% signifies perfect recognition, where every character in the decoded text matches the reference text. Conversely, a higher WER indicates a greater number of errors. The specific threshold for acceptable WER depends on the application and the level of noise or complexity in the speech data.

**Character-Level WER:**

- Analyzes errors at the level of **individual characters**.
- Counts the number of insertions, deletions, and substitutions of **individual characters** needed to make the recognized text match the reference text.
- Highly sensitive to **minor errors** within a word.
- Provides a **more detailed picture** of the system's recognition accuracy at the character level.

**Word Error Rate (WER):**

- Focuses on **whole words** being correct or incorrect.
- Counts the number of insertions, deletions, and substitutions of **entire words** needed to make the recognized text match the reference text.
- Less sensitive to minor errors within a word (e.g., typos, misspellings).
- Offers a **broader overview** of the system's performance in capturing word sequences.

**'wav2vec2.0with Libri-Light'**

**Whisper** is a text normalization technique specifically designed for ASR applications. It aims to pre-process text data to create a consistent format for the model, improving its ability to recognize speech. Here are some key techniques employed by Whisper:

- **Number Normalization:** Numbers can be represented in various ways (e.g., "twenty-one", "21", "four hundred twenty-two"). Whisper can convert these variations to a standard format (e.g., all digits) for better recognition. This is particularly helpful for ASR systems that struggle to differentiate between spoken numbers and written numbers with punctuation or spelled-out forms.
- **Lowercasing:** Speech recognition systems often struggle with capitalization inconsistencies. Whisper converts all text to lowercase to make it case-independent. This is because capitalization can vary depending on the speaker's emphasis, formality of speech, or even typos during transcription. By converting everything to lowercase, the ASR model can focus on the sequence of letters rather than capitalization cues.
- **Clock Time Normalization:** Time expressions like "3:15 PM" or "quarter past three" are standardized into a common format (e.g., "0315"). This is important because ASR models might struggle to interpret spoken time formats that rely on relative terms like "quarter past" or "half past." Converting everything to a 24-hour clock format with minutes (e.g., "HHMM") ensures consistency and avoids confusion for the model.
- **Short Form Expansion:** Abbreviations and contractions ("e.g.", "don't") are expanded to their full forms ("for example", "do not") for easier recognition. ASR models are trained on a specific vocabulary, and encountering abbreviations or contractions outside of their training data can lead to errors. Expanding these forms helps the model match the spoken word to the corresponding full term in its vocabulary.
- **Punctuation Removal:** Punctuation marks can sometimes confuse ASR models, especially those trained on unpunctuated speech data. Removing commas, periods, question marks, and other punctuation can improve recognition accuracy. Punctuation can introduce pauses or breaks in speech that might not be relevant to the meaning of the sentence, and removing them allows the model to focus on the core words.
- **Removal of Disfluencies:** Speech disfluencies like "um," "uh," and stutters can be removed to create a cleaner representation of the intended message. Disfluencies are hesitations, sound fillers, or repetitions that people use in natural conversation but don't carry meaning. By removing them, the ASR model can concentrate on the words that convey the actual content of the speech.
- **Dialect Normalization:** For multilingual systems, dialect-specific variations in pronunciation and word choice can be normalized to a standard form. This is particularly useful for ASR models that need to handle a diverse range of speakers with different dialects or accents. By normalizing pronunciations and vocabulary to a common reference, the model becomes more adaptable and less prone to errors due to dialectal variations.

**Benefits of Whisper Normalization:**

- **Improved WER:** Consistent text format leads to more accurate speech recognition by reducing ambiguity in the input data. The ASR model can focus on recognizing the core phonemes and words in the speech, rather than getting sidetracked by variations in how numbers, times, and abbreviations are expressed.
- **Reduced Training Time:** Normalized data allows the model to learn patterns more efficiently. By removing inconsistencies and converting everything to a standard format, the model needs less training data to achieve the same level of accuracy. This translates to faster development cycles and lower computational costs.
- **Increased Generalizability:** The model becomes less sensitive to variations in spoken language. Accents, dialects, and informal speech patterns can all introduce inconsistencies in how people express themselves. Whisper normalization helps to mitigate these variations, making the ASR model more robust and adaptable to different speaking styles. This is particularly important for real-world applications where the system needs to handle a wide range of speakers and speaking environments.

The wav2vec 2.0 model with Libri-Light pre-training data is a popular choice for ASR tasks. It leverages WER to evaluate its performance. By applying Whisper normalization, NLP practitioners can significantly enhance the performance of their ASR systems. This pre-processing step helps bridge the gap between the spoken word and its textual representation, leading to more accurate and robust speech recognition

## The Landscape of WER in Research Papers

The research papers I examined employed WER as the evaluation metric. This preference stems from WER's ability to offer a standardized measure of performance, allowing for straightforward comparison of different ASR models on a single scale. However, it's crucial to acknowledge that the specific WER type (CER or WER) might not be explicitly stated in the papers. Examining the research focus can often provide clues. For instance, papers emphasizing phoneme recognition accuracy are more likely to utilize CER.

While WER enjoys widespread adoption, it's not without its shortcomings. One key limitation is its indifference to the severity of errors. A single substitution that drastically alters the sentence meaning carries the same weight as a minor substitution with minimal impact. Additionally, WER doesn't account for semantic accuracy. Sentences with identical word sequences can convey different meanings.

## A Broader Evaluation Spectrum

To gain a more comprehensive understanding of ASR system performance, it's beneficial to consider metrics beyond WER. Here are some complementary options:

- **Sentence Error Rate (SER):** This metric extends the concept of WER to the sentence level. It analyzes insertions, deletions, and substitutions of entire sentences within the reference and hypothesized transcripts. While WER focuses on individual word accuracy, SER provides insights into how well the ASR system preserves the overall structure and meaning of the spoken content. For instance, an ASR system might produce a series of correctly recognized words that are nonetheless reordered in the output compared to the original sentence. WER would not penalize such an error, whereas SER would capture this discrepancy.
- **Mean Opinion Score (MOS):** This subjective metric relies on human listeners to evaluate the perceived quality of the speech recognition output. Listeners are typically presented with ASR outputs and asked to rate their naturalness, clarity, and fidelity to the original speech on a predefined scale. While WER and SER provide quantitative measures of accuracy, MOS offers valuable insights into the human perception of ASR performance. A high WER or SER might not necessarily translate to poor speech recognition quality if the errors are not noticeable to human listeners. Conversely, an ASR system with a low WER or SER could still produce unnatural-sounding output that is difficult for humans to understand.

By employing a combination of these metrics, NLP folks can develop a richer understanding of ASR system strengths and weaknesses, ultimately leading to the creation of more accurate and nuanced speech recognition models

Reference: https://cdn.openai.com/papers/whisper.pdf, https://github.com/sabeswari-kante/WER-Word-Error-Rate/blob/main/Basic_WER_huggingface.ipynb, https://github.com/openai/whisper/blob/main/notebooks/LibriSpeech.ipynb, https://www.clari.com/blog/word-error-rate/#:~:text=Put%20simply%2C%20WER%20is%20the,are%20a%20little%20more%20nuanced., https://medium.com/nlplanet/two-minutes-nlp-intro-to-word-error-rate-wer-for-speech-to-text-fc17a98003ea, https://en.wikipedia.org/wiki/Word_error_rate, https://github.com/sabeswari-kante/WER-Word-Error-Rate/blob/main/Oct2023-WER.pdf.