The model faced difficulty to identify named entities correctly and identified dates, numbers and abbreviations wrongly. Model predicted numbers and dates as named entities. Also, it was confused with acronyms and struggle to understand the context of some phrases related to time. One of the observations is "Sentence: Germany imported 47,600 sheep from Britain last year , nearly half of total imports .   FP: 47,600,  FP: nearly half,  FP: last year ". To improve the performance, a filter was implemented which contains a list of labels (MONEY, PERCENT, TIME, DATE, CARDINAL, ORDINAL, QUANTITY). Any identified entities falling into these categories are excluded from the final predictions. This approach reduced false positives.

The word "context" refers to the surrounding textual information. For the data given, a context is defined as the 5 tokens to the left and the 5 tokens to the right of the mention, used for disambiguating between different entities. This helps to figure out which thing we're talking about by looking at the words around it. It helps predict the right thing based on how it's used in a sentence. Another type of context could be considering syntactic relationships which involve analyzing the grammatical structure of the sentence.

From the sentences, we want to know the information about Ruth Bader Ginsburg, her profession as a justice, and when did she serve in the Supreme Court. Pronouns like "She" make it tricky because we need to figure out who or what they're talking about. If we misunderstand, we might get the facts wrong.