

1 Analysis: "Lost in the Middle" Phenomenon

2 To analyze the "lost in the middle" phenomenon in our LLM-based annotations, we conducted
3 an additional experiment based on the reviewer's suggestion. For a given user, we divided the
4 data into time segments and calculated inter-annotator agreement (IAA) using Cohen's Kappa
5 scores across different models and settings. The data was segmented based on the submission_id,
6 author_key_name, and stance_id_timestamp. For each group (i.e., each combination of submission_id
7 and author_key_name), the timestamps were divided into equal segments. The number of entries for
8 each group was divided by the desired number of segments (3), and the division was done as evenly
9 as possible, with each segment containing a roughly equal number of time-stamped entries. Figure
10 1 in the rebuttal PDF reports the comparison statistics of IAA scores for the stance detection task
11 across initial, middle, and later time stamps.

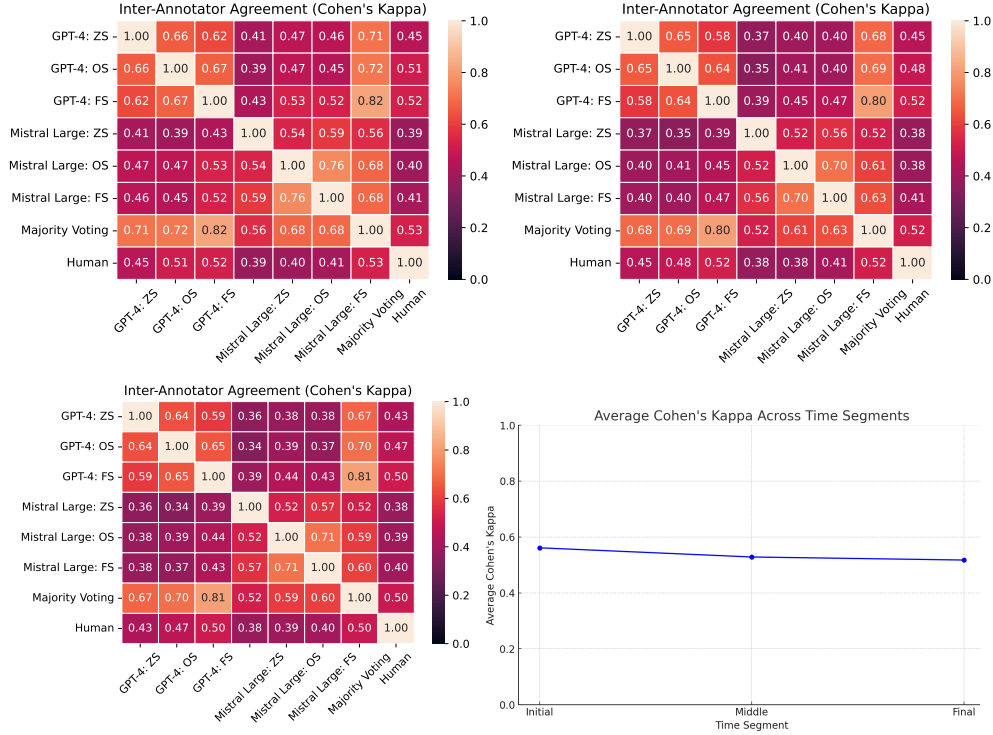


Figure 1: The inter-annotator agreement (IAA) on the USDC test dataset was measured using Cohen's Kappa score across three segments: initial, middle, and later time stamps. The top two rows represent the initial and middle time stamps, while the bottom left corresponds to the later time stamp. The bottom right reports the average Kappa score across all time segments.

12 Analysis: "Recency Bias" Phenomenon

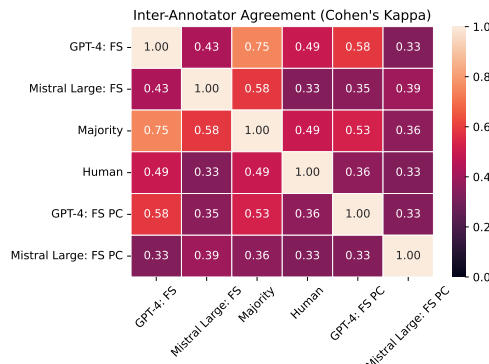


Figure 2: Inter-annotator agreement (IAA) on the test dataset was calculated for both the full conversations and the prior context for a given user. In this context, "GPT-4 FS PC" and "Mistral Large: FS PC" refer to the annotations based on prior context.