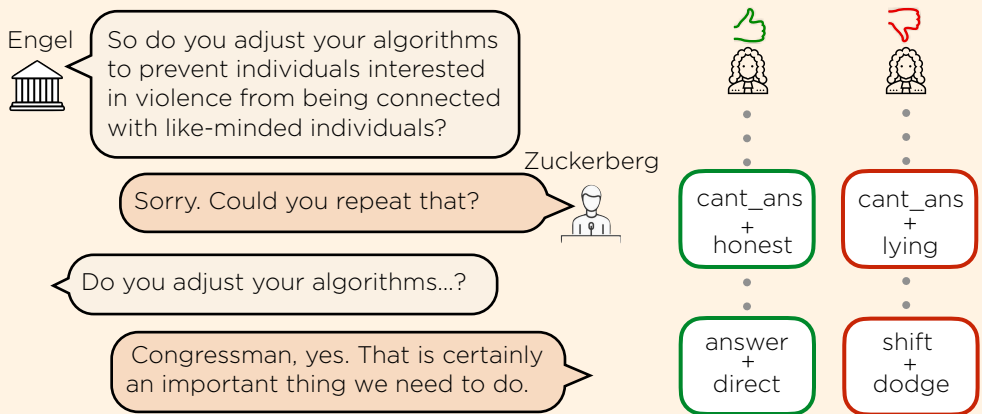
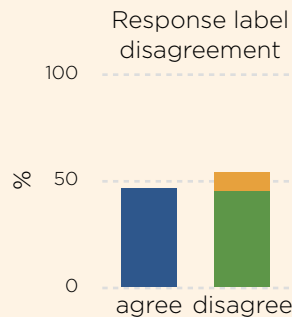




Dataset has *multiple* and *subjective* labels for interpreting each response:



Disagreements are frequent, with usually opposing intents:

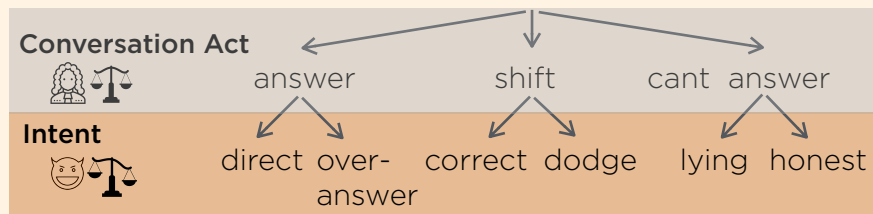


IAA Krippendorff's  $\alpha$ : 0.49  
 Conversation Act: 0.65  
 Intent: 0.38

Most frequent disagreements:

- ans+direct vs. shift+dodge
- shift+correct vs. shift+dodge
- cant\_ans+honest vs. cant\_ans+lying

Hierarchical taxonomy for response label:



Taking into account annotator bias improves prediction of all labels:

Hierarchical: separate classifiers for CA and intent, with forced consistency during training

+Annotator: add annotator sentiment towards the witness

Multi-label classification of all response labels

