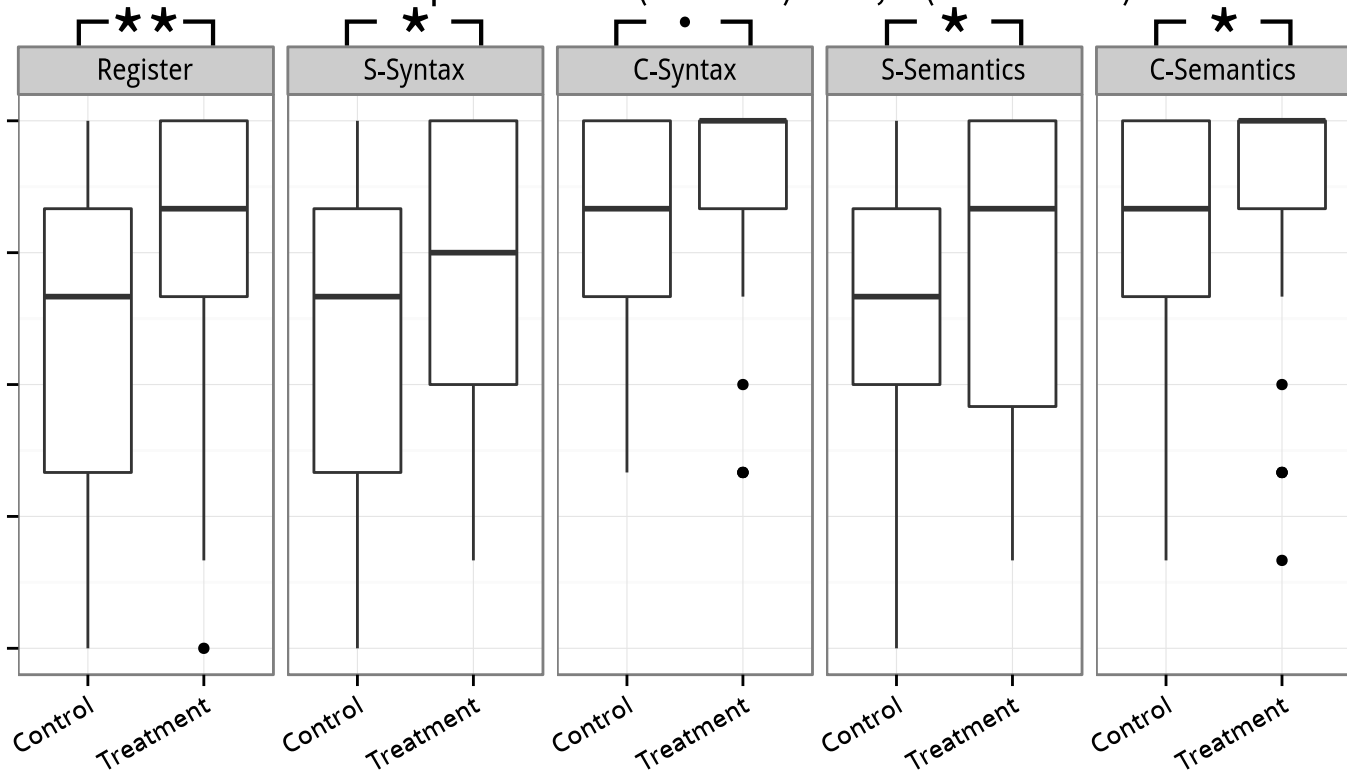


Welch two sample t -test: $N(\text{Control}) = 66$, $N(\text{Treatment}) = 36$

Average annotator evaluation score (3-point scale)



Condition