

Model Evaluation

Model Statistics:

<i>Model</i>	<i>Correct Answers</i>	<i>Partially Correct Answers</i>	<i>Wrong Answers</i>
<i>deepset-xlm-roberta-large-squad2</i>	141 (70.5%)	22 (11.0%)	37 (18.5%)

Response Times Statistics:

<i>Model</i>	<i>Average</i>	<i>Median</i>	<i>Min</i>	<i>Max</i>
<i>deepset-xlm-roberta-large-squad2</i>	15.6	15.4	13.3	26.3

***All the response times are in seconds*

Model answers confidence scores:

For deepset-xlm-roberta-large-squad2 model:

<i>Range</i>	<i>Correct</i>	<i>Partially Correct</i>	<i>Wrong</i>
<i>0.0 - 0.1</i>	9 (6.4%)	3 (13.6%)	11 (29.7%)
<i>0.1 - 0.2</i>	2 (1.4%)	3 (13.6%)	10 (27.0%)
<i>0.2 - 0.3</i>	7 (5.0%)	1 (4.5%)	2 (5.4%)
<i>0.3 - 0.4</i>	14 (9.9%)	3 (13.6%)	3 (8.1%)
<i>0.4 - 0.5</i>	15 (10.6%)	7 (31.8%)	3 (8.1%)
<i>0.5 - 0.6</i>	15 (10.6%)	1 (4.5%)	3 (8.1%)
<i>0.6 - 0.7</i>	18 (12.8%)	0 (0.0%)	3 (8.1%)
<i>0.7 - 0.8</i>	18 (12.8%)	1 (4.5%)	1 (2.7%)
<i>0.8 - 0.9</i>	15 (10.6%)	3 (13.6%)	0 (0.0%)
<i>0.9 - 1.0</i>	28 (19.9%)	0 (0.0%)	1 (2.7%)

Question Statistics:

- 141 questions were answered correctly by all models
- 37 questions were answered incorrectly by all models
- 141 questions were answered correctly by at least one model
- 141 questions were answered correctly by exactly one model
- 163 questions were answered correctly or partially correct by all models