

# Model Evaluation

## Model Statistics:

<i>Model</i>	<i>Correct Answers</i>	<i>Partially Correct Answers</i>	<i>Wrong Answers</i>
<i>deepset/tinyroberta-squad2</i>	26 (65.0%)	5 (12.5%)	9 (22.5%)
<i>deepset/roberta-base-squad2-covid</i>	26 (65.0%)	5 (12.5%)	9 (22.5%)

## Response Times Statistics:

<i>Model</i>	<i>Average</i>	<i>Median</i>	<i>Min</i>	<i>Max</i>
<i>deepset/tinyroberta-squad2</i>	7.2	7.1	6.9	7.5
<i>deepset/roberta-base-squad2-covid</i>	9.5	9.4	9.0	10.3

*\*\*All the response times are in seconds*

## Model answers confidence scores:

For deepset/tinyroberta-squad2 model:

<i>Range</i>	<i>Correct</i>	<i>Partially Correct</i>	<i>Wrong</i>
<i>0.0 - 0.1</i>	2 (7.7%)	0 (0.0%)	2 (22.2%)
<i>0.1 - 0.2</i>	2 (7.7%)	0 (0.0%)	0 (0.0%)
<i>0.2 - 0.3</i>	1 (3.8%)	1 (20.0%)	1 (11.1%)
<i>0.3 - 0.4</i>	0 (0.0%)	1 (20.0%)	0 (0.0%)
<i>0.4 - 0.5</i>	4 (15.4%)	1 (20.0%)	2 (22.2%)
<i>0.5 - 0.6</i>	2 (7.7%)	1 (20.0%)	1 (11.1%)
<i>0.6 - 0.7</i>	4 (15.4%)	0 (0.0%)	0 (0.0%)
<i>0.7 - 0.8</i>	2 (7.7%)	1 (20.0%)	0 (0.0%)
<i>0.8 - 0.9</i>	3 (11.5%)	0 (0.0%)	2 (22.2%)
<i>0.9 - 1.0</i>	6 (23.1%)	0 (0.0%)	1 (11.1%)

For deepset/roberta-base-squad2-covid model:

<i>Range</i>	<i>Correct</i>	<i>Partially Correct</i>	<i>Wrong</i>
<i>0.0 - 0.1</i>	3 (11.5%)	0 (0.0%)	5 (55.6%)
<i>0.1 - 0.2</i>	0 (0.0%)	3 (60.0%)	0 (0.0%)
<i>0.2 - 0.3</i>	2 (7.7%)	0 (0.0%)	0 (0.0%)
<i>0.3 - 0.4</i>	0 (0.0%)	0 (0.0%)	0 (0.0%)
<i>0.4 - 0.5</i>	1 (3.8%)	0 (0.0%)	1 (11.1%)
<i>0.5 - 0.6</i>	4 (15.4%)	0 (0.0%)	1 (11.1%)
<i>0.6 - 0.7</i>	1 (3.8%)	1 (20.0%)	0 (0.0%)
<i>0.7 - 0.8</i>	3 (11.5%)	1 (20.0%)	1 (11.1%)
<i>0.8 - 0.9</i>	4 (15.4%)	0 (0.0%)	1 (11.1%)
<i>0.9 - 1.0</i>	8 (30.8%)	0 (0.0%)	0 (0.0%)

Question Statistics:

- 22 questions were answered correctly by all models
- 6 questions were answered incorrectly by all models
- 30 questions were answered correctly by at least one model
- 8 questions were answered correctly by exactly one model
- 28 questions were answered correctly or partially correct by all models