

Model Evaluation

Model Statistics:

Model	Correct Answers	Partially Correct Answers	Wrong Answers
deepset/bert-large-uncased-whole-word-masking-squad2	150 (75.0%)	5 (2.5%)	45 (22.5%)
bert-large-cased-whole-word-masking-finetuned-squad	147 (73.5%)	5 (2.5%)	48 (24.0%)
dmis-lab/biobert-large-cased-v1.1-squad	145 (72.5%)	6 (3.0%)	49 (24.5%)
deepset/roberta-base-squad2	144 (72.0%)	7 (3.5%)	49 (24.5%)
bert-large-uncased-whole-word-masking-finetuned-squad	143 (71.5%)	5 (2.5%)	52 (26.0%)
deepset/minilm-uncased-squad2	138 (69.0%)	7 (3.5%)	55 (27.5%)
rsvp-ai/bertserini-bert-base-squad	138 (69.0%)	5 (2.5%)	57 (28.5%)
distilbert-base-uncased-distilled-squad	134 (67.0%)	4 (2.0%)	62 (31.0%)
distilbert-base-cased-distilled-squad	129 (64.5%)	4 (2.0%)	67 (33.5%)
deepset/bert-base-cased-squad2	125 (62.5%)	3 (1.5%)	72 (36.0%)

Response Times Statistics:

Model	Average	Median	Min	Max
distilbert-base-uncased-distilled-squad	11.0	8.1	7.3	37.7
deepset/minilm-uncased-squad2	11.2	8.3	7.7	31.0
distilbert-base-cased-distilled-squad	12.0	8.5	7.8	99.8
deepset/bert-base-cased-squad2	12.6	9.4	8.2	35.3
rsvp-ai/bertserini-bert-base-squad	12.6	9.5	8.2	37.0
deepset/roberta-base-squad2	13.4	10.6	9.3	31.6
bert-large-cased-whole-word-masking-finetuned-squad	14.8	13.0	9.5	41.9
bert-large-uncased-whole-word-masking-finetuned-squad	15.9	13.2	10.2	59.8
deepset/bert-large-uncased-whole-word-masking-squad2	16.2	13.7	10.1	45.9

<i>dmis-lab/biobert-large-cased-v1.1-squad</i>	16.5	14.2	10.3	37.3
--	------	------	------	------

***All the response times are in seconds*

Model answers confidence scores:

For deepset/roberta-base-squad2 model:

<i>Range</i>	<i>Correct</i>	<i>Partially Correct</i>	<i>Wrong</i>
<i>0.0 - 0.1</i>	7 (4.9%)	1 (14.3%)	12 (24.5%)
<i>0.1 - 0.2</i>	7 (4.9%)	0 (0.0%)	6 (12.2%)
<i>0.2 - 0.3</i>	17 (11.8%)	1 (14.3%)	4 (8.2%)
<i>0.3 - 0.4</i>	12 (8.3%)	1 (14.3%)	15 (30.6%)
<i>0.4 - 0.5</i>	26 (18.1%)	1 (14.3%)	1 (2.0%)
<i>0.5 - 0.6</i>	22 (15.3%)	3 (42.9%)	6 (12.2%)
<i>0.6 - 0.7</i>	12 (8.3%)	0 (0.0%)	0 (0.0%)
<i>0.7 - 0.8</i>	13 (9.0%)	0 (0.0%)	3 (6.1%)
<i>0.8 - 0.9</i>	14 (9.7%)	0 (0.0%)	0 (0.0%)
<i>0.9 - 1.0</i>	14 (9.7%)	0 (0.0%)	2 (4.1%)

For bert-large-uncased-whole-word-masking-finetuned-squad model:

<i>Range</i>	<i>Correct</i>	<i>Partially Correct</i>	<i>Wrong</i>
<i>0.0 - 0.1</i>	1 (0.7%)	0 (0.0%)	8 (15.4%)
<i>0.1 - 0.2</i>	7 (4.9%)	1 (20.0%)	6 (11.5%)
<i>0.2 - 0.3</i>	11 (7.7%)	0 (0.0%)	7 (13.5%)
<i>0.3 - 0.4</i>	16 (11.2%)	0 (0.0%)	5 (9.6%)
<i>0.4 - 0.5</i>	15 (10.5%)	1 (20.0%)	8 (15.4%)
<i>0.5 - 0.6</i>	23 (16.1%)	0 (0.0%)	5 (9.6%)
<i>0.6 - 0.7</i>	19 (13.3%)	2 (40.0%)	3 (5.8%)
<i>0.7 - 0.8</i>	13 (9.1%)	1 (20.0%)	4 (7.7%)

0.8 - 0.9	13 (9.1%)	0 (0.0%)	2 (3.8%)
0.9 - 1.0	25 (17.5%)	0 (0.0%)	4 (7.7%)

For distilbert-base-cased-distilled-squad model:

<i>Range</i>	<i>Correct</i>	<i>Partially Correct</i>	<i>Wrong</i>
0.0 - 0.1	2 (1.6%)	0 (0.0%)	6 (9.0%)
0.1 - 0.2	5 (3.9%)	2 (50.0%)	15 (22.4%)
0.2 - 0.3	12 (9.3%)	1 (25.0%)	4 (6.0%)
0.3 - 0.4	15 (11.6%)	0 (0.0%)	9 (13.4%)
0.4 - 0.5	14 (10.9%)	0 (0.0%)	6 (9.0%)
0.5 - 0.6	22 (17.1%)	0 (0.0%)	8 (11.9%)
0.6 - 0.7	10 (7.8%)	0 (0.0%)	4 (6.0%)
0.7 - 0.8	15 (11.6%)	0 (0.0%)	4 (6.0%)
0.8 - 0.9	9 (7.0%)	0 (0.0%)	4 (6.0%)
0.9 - 1.0	25 (19.4%)	1 (25.0%)	7 (10.4%)

For deepset/bert-large-uncased-whole-word-masking-squad2 model:

<i>Range</i>	<i>Correct</i>	<i>Partially Correct</i>	<i>Wrong</i>
0.0 - 0.1	8 (5.3%)	0 (0.0%)	5 (11.1%)
0.1 - 0.2	8 (5.3%)	1 (20.0%)	7 (15.6%)
0.2 - 0.3	10 (6.7%)	0 (0.0%)	6 (13.3%)
0.3 - 0.4	15 (10.0%)	2 (40.0%)	4 (8.9%)
0.4 - 0.5	24 (16.0%)	0 (0.0%)	3 (6.7%)
0.5 - 0.6	17 (11.3%)	1 (20.0%)	5 (11.1%)
0.6 - 0.7	14 (9.3%)	0 (0.0%)	5 (11.1%)
0.7 - 0.8	19 (12.7%)	0 (0.0%)	6 (13.3%)
0.8 - 0.9	16 (10.7%)	0 (0.0%)	0 (0.0%)
0.9 - 1.0	19 (12.7%)	1 (20.0%)	4 (8.9%)

For distilbert-base-uncased-distilled-squad model:

Range	Correct	Partially Correct	Wrong
0.0 - 0.1	4 (3.0%)	1 (25.0%)	13 (21.0%)
0.1 - 0.2	4 (3.0%)	0 (0.0%)	8 (12.9%)
0.2 - 0.3	12 (9.0%)	1 (25.0%)	8 (12.9%)
0.3 - 0.4	13 (9.7%)	1 (25.0%)	6 (9.7%)
0.4 - 0.5	19 (14.2%)	0 (0.0%)	7 (11.3%)
0.5 - 0.6	12 (9.0%)	0 (0.0%)	5 (8.1%)
0.6 - 0.7	17 (12.7%)	0 (0.0%)	3 (4.8%)
0.7 - 0.8	19 (14.2%)	0 (0.0%)	3 (4.8%)
0.8 - 0.9	9 (6.7%)	1 (25.0%)	4 (6.5%)
0.9 - 1.0	25 (18.7%)	0 (0.0%)	5 (8.1%)

For rsvp-ai/bertserini-bert-base-squad model:

Range	Correct	Partially Correct	Wrong
0.0 - 0.1	4 (2.9%)	0 (0.0%)	10 (17.5%)
0.1 - 0.2	8 (5.8%)	3 (60.0%)	8 (14.0%)
0.2 - 0.3	9 (6.5%)	0 (0.0%)	14 (24.6%)
0.3 - 0.4	18 (13.0%)	0 (0.0%)	5 (8.8%)
0.4 - 0.5	19 (13.8%)	0 (0.0%)	4 (7.0%)
0.5 - 0.6	16 (11.6%)	1 (20.0%)	4 (7.0%)
0.6 - 0.7	13 (9.4%)	0 (0.0%)	5 (8.8%)
0.7 - 0.8	15 (10.9%)	1 (20.0%)	1 (1.8%)
0.8 - 0.9	14 (10.1%)	0 (0.0%)	4 (7.0%)
0.9 - 1.0	22 (15.9%)	0 (0.0%)	2 (3.5%)

For deepset/minilm-uncased-squad2 model:

Range	Correct	Partially Correct	Wrong
0.0 - 0.1	4 (2.9%)	1 (14.3%)	9 (16.4%)

0.1 - 0.2	6 (4.3%)	2 (28.6%)	7 (12.7%)
0.2 - 0.3	12 (8.7%)	2 (28.6%)	3 (5.5%)
0.3 - 0.4	9 (6.5%)	1 (14.3%)	9 (16.4%)
0.4 - 0.5	20 (14.5%)	0 (0.0%)	3 (5.5%)
0.5 - 0.6	21 (15.2%)	0 (0.0%)	6 (10.9%)
0.6 - 0.7	18 (13.0%)	0 (0.0%)	8 (14.5%)
0.7 - 0.8	15 (10.9%)	0 (0.0%)	5 (9.1%)
0.8 - 0.9	10 (7.2%)	0 (0.0%)	3 (5.5%)
0.9 - 1.0	23 (16.7%)	1 (14.3%)	2 (3.6%)

For dmis-lab/biobert-large-cased-v1.1-squad model:

Range	Correct	Partially Correct	Wrong
0.0 - 0.1	3 (2.1%)	1 (16.7%)	7 (14.3%)
0.1 - 0.2	9 (6.2%)	0 (0.0%)	3 (6.1%)
0.2 - 0.3	11 (7.6%)	0 (0.0%)	6 (12.2%)
0.3 - 0.4	15 (10.3%)	3 (50.0%)	6 (12.2%)
0.4 - 0.5	11 (7.6%)	1 (16.7%)	6 (12.2%)
0.5 - 0.6	15 (10.3%)	1 (16.7%)	7 (14.3%)
0.6 - 0.7	19 (13.1%)	0 (0.0%)	3 (6.1%)
0.7 - 0.8	18 (12.4%)	0 (0.0%)	3 (6.1%)
0.8 - 0.9	16 (11.0%)	0 (0.0%)	4 (8.2%)
0.9 - 1.0	28 (19.3%)	0 (0.0%)	4 (8.2%)

For deepset/bert-base-cased-squad2 model:

Range	Correct	Partially Correct	Wrong
0.0 - 0.1	4 (3.2%)	1 (33.3%)	22 (30.6%)
0.1 - 0.2	6 (4.8%)	0 (0.0%)	10 (13.9%)
0.2 - 0.3	11 (8.8%)	1 (33.3%)	11 (15.3%)

0.3 - 0.4	14 (11.2%)	0 (0.0%)	6 (8.3%)
0.4 - 0.5	11 (8.8%)	0 (0.0%)	8 (11.1%)
0.5 - 0.6	18 (14.4%)	1 (33.3%)	3 (4.2%)
0.6 - 0.7	16 (12.8%)	0 (0.0%)	4 (5.6%)
0.7 - 0.8	8 (6.4%)	0 (0.0%)	1 (1.4%)
0.8 - 0.9	10 (8.0%)	0 (0.0%)	1 (1.4%)
0.9 - 1.0	27 (21.6%)	0 (0.0%)	6 (8.3%)

For bert-large-cased-whole-word-masking-finetuned-squad model:

<i>Range</i>	<i>Correct</i>	<i>Partially Correct</i>	<i>Wrong</i>
0.0 - 0.1	2 (1.4%)	0 (0.0%)	9 (18.8%)
0.1 - 0.2	5 (3.4%)	0 (0.0%)	6 (12.5%)
0.2 - 0.3	15 (10.2%)	0 (0.0%)	6 (12.5%)
0.3 - 0.4	15 (10.2%)	1 (20.0%)	5 (10.4%)
0.4 - 0.5	17 (11.6%)	1 (20.0%)	7 (14.6%)
0.5 - 0.6	19 (12.9%)	2 (40.0%)	2 (4.2%)
0.6 - 0.7	21 (14.3%)	0 (0.0%)	2 (4.2%)
0.7 - 0.8	15 (10.2%)	0 (0.0%)	6 (12.5%)
0.8 - 0.9	16 (10.9%)	1 (20.0%)	2 (4.2%)
0.9 - 1.0	22 (15.0%)	0 (0.0%)	3 (6.2%)

Question Statistics:

- 84 questions were answered correctly by all models
- 21 questions were answered incorrectly by all models
- 174 questions were answered correctly by at least one model
- 5 questions were answered correctly by exactly one model
- 88 questions were answered correctly or partially correct by all models