

### Model Evaluation Results:

The data collected in different phases of the ablation study is evaluated using RoBERTa, trained on SNLI. The 4 phases are: regular crowdsourcing, traffic signal feedback, AutoFix, and Full System availability.

Data	ROBERTA (S2)	S22
Crowdsourcing	1	1.000
Traffic Signal	1	0.856
AutoFix	1	0.803
Full System	1	0.748
Full Data (SNLI Dev set)	0.89	--

### DQI Evaluation Results:

#### Parameters:

T1 : Number of unique vocabulary words/Size of Dataset  
T2 : Sentence Length Deviation  
C1 : DQI component C1: Vocabulary  
W1 : Standard Deviation of Word Distribution  
W2 : Proportion of words with sufficient inductive bias  
Bi T1 : Standard Deviation of Bigram Distribution  
Bi T2 : Proportion of bigrams with sufficient inductive bias  
Word : Raw count of Words  
S : Raw count of Sentences  
Adj : Raw count of Adjectives  
Adv : Raw count of Adverbs  
N : Raw count of Nouns  
V : Raw count of Verbs  
Bi : Raw count of Bigrams  
Tri : Raw count of Trigrams  
C3 : DQI Component 3 : Inter-Sample Semantic Textual Similarity  
C4 : DQI Component 4: Inter-Sample Word Similarity  
WO : Average premise-hypothesis overlap across samples  
STS : Average premise-hypothesis similarity across samples

**Results for Original SNLI Samples (Sentence 1 and Sentence 2):**

Data	T 1	T 2	C 1	W1	W2	Bi T1	Bi T2	Word	S	Adj	Adv	N	V	Bi	Tri	C3
Crowdsourcing	6. 0 7	5. 5 5	1 1. 6 2	144 .93	0.5 7	62 1.0 1	0. 88	419	138	99	9	20 1	88	99 4	114 9	202
Traffic Signal	6. 4 3	5. 3 3	1 1. 7 6	170 .09	0.7 1	66 2.5 5	0. 86	444	138	88	17	21 3	102	10 27	116 7	202
AutoFix	5. 8 1	5. 0 8	1 0. 8 9	131 .92	0.6	56 4.5 4	0. 89	401	138	83	12	20 4	80	93 8	104 1	202
Full	6. 6 2	5. 8 5	1 2. 4 8	203 .01	0.6 9	92 2.0 6	0. 96	457	138	90	8	24 5	89	10 55	115 6	202

**Results for Original SNLI Samples (Sentence 1 and Sentence 2)- Here the last 12 columns represent values for the entailment label:**

Data	C 4	W O	S T S	W1	W2	Bi T1	Bi T2	Word	S	Adj	Adv	N	V	Bi	Tri
Crowdsourcing	0. 0 0 8	0. 0 7 4	0. 0 0 1	108 .08	0.3 6	36 2.2 7	0. 83	163	46	35	4	69	36	34 7	369
Traffic Signal	0. 0 0 8	0. 0 7 1	0. 0 0 1	108 .54	0.4 3	40 6.4 5	0. 84	167	46	32	1	83	41	34 6	363
AutoFix	0. 0 0 7	0. 0 8 1	0. 0 0 1	98. 14	0.3 4	47 7.1 0	0. 89	164	46	37	6	73	37	35 0	360
Full	0. 0 0 8	0. 0 7 4	0. 0 0 1	134 .56	0.3 5	60 1.6 5	0. 93	169	46	28	3	91	35	36 6	370



**Results for Original SNLI Samples (Sentence 1 and Sentence 2)- Here the columns represent values for the neutral label:**

Data	W1	W2	Bi T1	Bi T2	Word	S	Adj	Adv	N	V	Bi	Tri
Crowdsourcing	126 .35	0.3 6	53 1.0 5	0. 90	181	46	43	1	89	38	39 0	407
Traffic Signal	210 .84	0.7 8	66 8.4 2	0. 94	224	46	45	16	93	54	45 5	465
AutoFix	136 .15	0.4 8	36 4.0 3	0. 81	172	46	35	3	87	35	33 9	350
Full	162 .47	0.4 0	74 9.4 5	0. 97	178	46	40	3	85	39	38 5	381

**Results for Original SNLI Samples (Sentence 1 and Sentence 2)- Here the columns represent values for the neutral label:**

Data	W1	W2	Bi T1	Bi T2	Word	S	Adj	Adv	N	V	Bi	Tri
Crowdsourcing	144 .37	0.3 9	49 2.8 2	0. 87	194	46	45	3	10 3	35	40 1	410
Traffic Signal	131 .17	0.4 7	36 9.0 8	0. 81	168	46	32	5	85	36	34 2	369
AutoFix	106 .91	<b>0.3 3</b>	39 0.2 13	0. 87	170	46	27	5	85	36	35 3	372
Full	170 .83	0.5 7	61 8.1 9	1	212	46	53	4	11 3	34	41 2	429

**Results for VAIDA's Produced Samples (Sentence 1 and Sentence 22):**

Data	T 1	T 2	C 1	W1	W2	Bi T1	Bi T2	Word	S	Adj	Adv	N	V	Bi	Tri	C3
Crowdsourcing	6. 2 4 6	5. 5 8	1 1. 8 3	156 .32	0.6 2	68 7.9 1	0. 89	431	138	89	6	22 3	91	10 27	114 5	202
Traffic Signal	7. 3 9	5. 2 2	1 2. 6 2	248 .08	0.8 1	91 9.5 3	0. 88	503	136	109	14	11 7	237	11 17	121 9	202
AutoFix	7. 0 8	5. 0 4	1 2. 1 3	237 .01	0.7 2	88 6.3 6	0. 91	489	138	89	20	24 7	104	10 30	107 1	202
Full	7. 3 4	5. 9 5	1 3. 3 0	263 .46	0.7 8	11 49. 68	1	507	138	111	12	24 7	113	11 07	117 3	202

**Results for VAIDA's Produced Samples (Sentence 1 and Sentence 22)- Here the last 12 columns represent values for the entailment label:**

Data	C 4	W O	S T S	W1	W2	Bi T1	Bi T2	Word	S	Adj	Adv	N	V	Bi	Tri
Crowdsourcing	0. 01	0. 08	0. 0 0 1	114. 09	0.3 9	44 6.8 8	0.9 2	169	46	37	2	82	38	35 5	367
Traffic Signal	0. 0 1	0. 07	0. 00 1	162. 89	0.5 6	70 8.3 0	0.8 8	200	46	45	8	79	54	418	425
AutoFix	0. 0 1	0. 0 8	0. 0 0 1	168. 45	0.5 0	76 0.2 5	0.9 3	201	46	37	11	96	43	39 6	395
Full	0.	0.	0.	181	0.5	84	0.9	204	46	36	2	11	37	403	403

	0 1	09	00 1	.25		3.3 4	8					5			
--	--------	----	---------	-----	--	----------	---	--	--	--	--	---	--	--	--

**Results for VAIDA's Produced Samples (Sentence 1 and Sentence 22)- Here the last 12 columns represent values for the neutral label:**

Data	W1	W2	Bi T1	Bi T2	Word	S	Adj	Adv	N	V	Bi	Tri
Crowdsourcing	128. 41	0.4 6	52 2.3 7	0.9 0	173	46	37	2	81	43	361	369
Traffic Signal	251 .87	0.6 3	90 7.0 4	0.9 8	223	46	48	13	95	54	439	442
AutoFix	217 .46	0.7 2	52 8.8 8	0. 83	198	46	43	7	93	42	35 5	339
Full	212. 53	0.5 7	11 32. 03	0.9 8	190	46	35	6	88	48	391	367

**Results for VAIDA's Produced Samples (Sentence 1 and Sentence 22)- Here the last 12 columns represent values for the contradiction label:**

Data	W1	W2	Bi T1	Bi T2	Word	S	Adj	Adv	N	V	Bi	Tri
Crowdsourcing	156 .21	0.5 6	48 5.5 4	0. 89	202	46	37	1	11 5	36	427	441
Traffic Signal	238. 90	0.6 8	53 8.9 9	0.8 8	209	46	44	6	10 2	44	392	398
AutoFix	184. 01	0.5 7	60 2.6 7	0.8 8	189	46	36	6	93	40	367	362

Full	214. 29	0.7 1	74 5.7 3	0.9 3	223	46	54	6	11 4	41	417	422
------	------------	----------	----------------	----------	-----	----	----	---	---------	----	-----	-----