

## Evaluation Results of Text Summarization Challenge 2

Takahiro Fukusima  
Otemon Gakuin University  
fukusima@res.otemon.ac.jp

Manabu Okumura  
Tokyo Institute of Technology  
oku@pi.titech.ac.jp

Hidetsugu Nanba  
Tokyo Institute of Technology<sup>1</sup>  
nanba@pi.titech.ac.jp

### TSC2 Evaluation Results

We show the formal run evaluation results of TSC2 in the following pages. First, the results of ranking evaluation for the two tasks, and then the results evaluation by revision for the two tasks are shown.

#### 1 Ranking evaluation

Here are the results of ranking evaluation for task A (single document summarization) and task B (multi-document summarization). Please note that the best score is 1.0 and the worst is 4.0, and the score for each cell is the average scores for the 30 articles for task A and 30 topics for task B.

System No	Content 20%	Read- ability 20%	Content 40%	Read- ability 40%
1	2.53	2.87	2.60	2.77
2	2.67	2.97	2.50	2.77
3	2.80	2.93	2.90	2.90
4	2.77	2.73	2.80	2.90
5	2.70	2.73	2.60	2.77
6	2.73	2.57	2.63	2.67
7	2.70	2.60	2.50	2.53
8	2.40	2.83	2.60	2.77
9	3.30	3.30	3.20	3.10
10	2.33	2.20	2.10	2.03

**Table 1** Ranking evaluation (task A)

---

<sup>1</sup> Presently with Hiroshima City University

System No	Content Short	Read-ability Short	Content Long	Read-ability Long
1	2.70	3.17	2.50	3.23
2	2.73	2.70	2.77	2.93
3	2.60	2.33	2.97	3.03
4	2.63	2.90	2.80	3.03
5	2.53	3.10	2.73	3.30
6	3.20	3.00	3.47	3.30
7	2.40	2.87	2.63	3.27
8	2.93	2.70	2.53	2.80
9	2.83	2.73	2.53	2.87
10	2.00	2.17	1.83	2.33

**Table 2 Ranking evaluation (task B)**

Next, we show the evaluation scores of human-produced summaries (abstract type 1 and type2) as well as those of baseline system.

	Content 20%	Read-ability 20%	Content 40%	Read-ability 40%
Human (type 1)	1.58	1.61	1.67	1.69
Human (type 2)	1.50	1.57	1.42	1.55
Baseline (Lead)	3.80	3.60	3.83	3.55

**Table 3 Ranking evaluation (task A, human and baseline)**

	Content Short	Read-ability Short	Content Long	Read-ability Long
Human (type 2)	1.65	2.38	1.82	2.38
Baseline (Lead)	2.80	2.20	2.70	2.22
Baseline (Stein)	2.48	2.00	2.50	1.99

**Table 4 Ranking evaluation (task B, human and baseline)**

## 2 Evaluation by revision

Next, we show the results of evaluation by revision for task A and task B. Three types of measures of revision are listed in the following tables. They are 1) average number of revised characters per text, 2) average number of revisions per text and 3) number of summaries which were given up.

[Task A (40%)]

[Average number of revised characters per text]

	Deletion		Insertion		Replacement			
	(unimportant)	(readability)	(important)	(readability)	(C) deletion	(C) Insertion	(R) deletion	(R) Insertion
F0101	44.6 (9.8 %)	0.3 (0.1 %)	44.6 (9.8 %)	1.6 (0.3 %)	3.7 (0.8 %)	5.3 (1.2 %)	3.6 (0.8 %)	2.2 (0.5 %)
F0102	38.8 (8.4 %)	2.2 (0.5 %)	43.8 (9.5 %)	1.5 (0.3 %)	6.0 (1.3 %)	4.3 (0.9 %)	5.1 (1.1 %)	2.3 (0.5 %)
F0103	59.8 (13.1 %)	2.6 (0.6 %)	77.3 (17.0 %)	3.3 (0.7 %)	13.2 (2.9 %)	9.8 (2.2 %)	4.5 (1.0 %)	1.1 (0.2 %)
F0104	59.9 (12.9 %)	2.0 (0.4 %)	79.7 (17.2 %)	3.4 (0.7 %)	12.1 (2.6 %)	7.4 (1.6 %)	2.4 (0.5 %)	1.3 (0.3 %)
F0105	75.1 (16.2 %)	0.9 (0.2 %)	75.7 (16.3 %)	1.8 (0.4 %)	9.0 (1.9 %)	7.8 (1.7 %)	3.0 (0.6 %)	1.5 (0.3 %)
F0106	61.5 (13.0 %)	2.0 (0.4 %)	71.1 (15.0 %)	2.1 (0.5 %)	5.8 (1.2 %)	5.7 (1.2 %)	2.6 (0.5 %)	1.5 (0.3 %)
F0107	57.2 (12.6 %)	2.7 (0.6 %)	77.2 (16.9 %)	6.3 (1.4 %)	4.1 (0.9 %)	2.2 (0.5 %)	3.4 (0.7 %)	1.7 (0.4 %)
F0108	68.7 (14.7 %)	1.9 (0.4 %)	76.8 (16.4 %)	0.4 (0.1 %)	10.0 (2.1 %)	5.9 (1.3 %)	5.3 (1.1 %)	1.9 (0.4 %)
Ld	75.5 (16.4 %)	1.2 (0.3 %)	78.5 (17.1 %)	0.3 (0.1 %)	3.9 (0.9 %)	15.3 (3.3 %)	0.3 (0.1 %)	0.2 (0.0 %)
Free	11.1 (2.3 %)	1.2 (0.3 %)	19.6 (4.1 %)	1.5 (0.3 %)	0.4 (0.1 %)	0.9 (0.2 %)	1.3 (0.3 %)	0.5 (0.1 %)
part	11.2 (2.4 %)	1.7 (0.4 %)	15.3 (3.3 %)	1.3 (0.3 %)	0.5 (0.1 %)	0.5 (0.1 %)	1.9 (0.4 %)	0.9 (0.2 %)
ALL	49.7 (10.7 %)	1.7 (0.4 %)	58.8 (12.7 %)	2.2 (0.5 %)	6.3 (1.4 %)	5.5 (1.2 %)	3.1 (0.7 %)	1.4 (0.3 %)

(Please note that ‘C’ stands for content and ‘R’ for readability in the table.)

[Average number of revisions per text]

	Deletion		Insertion		Replacement	
	(unimportant)	(readability)	(important)	(readability)	(content)	(readability)
F0101	2.0	0.1	1.5	0.4	0.5	0.7
F0102	1.6	0.4	1.5	0.4	0.4	0.8
F0103	2.3	0.2	2.4	0.2	0.4	0.5
F0104	2.4	0.4	2.7	0.5	0.4	0.5
F0105	2.0	0.3	1.7	0.1	0.7	0.7
F0106	2.8	0.2	2.3	0.4	0.3	0.6
F0107	2.5	0.6	1.8	0.2	0.1	0.5
F0108	2.0	0.4	2.4	0.1	0.4	0.6
ld	2.9	0.1	0.7	0.1	0.4	0.1
free	0.4	0.4	1.2	0.4	0.1	0.3
part	0.7	0.6	0.9	0.3	0.1	0.4
ALL	1.9	0.3	1.8	0.3	0.3	0.5

[number of given-up summaries]

	number of given-up summaries
F0101	2
F0102	2
F0103	1
F0104	0
F0105	3
F0106	2
F0107	1
F0108	2
Id	15
free	0
part	0
ALL	2.5

[Task A (20%)]

[Average number of revised characters per text]

	Deletion		Insertion		Replacement			
	(unimportant)	(readability)	(important)	(readability)	(C) deletion	(C) Insertion	(R) deletion	(R) Insertion
F0101	30.2 (12.8 %)	3.6 (1.5 %)	35.0 (14.8 %)	0.6 (0.3 %)	3.7 (1.6 %)	7.4 (3.1 %)	1.3 (0.6 %)	1.0 (0.4 %)
F0102	25.3 (10.7 %)	1.7 (0.7 %)	28.8 (12.1 %)	0.1 (0.0 %)	4.5 (1.9 %)	7.9 (3.3 %)	2.2 (0.9 %)	1.8 (0.8 %)
F0103	22.2 (10.6 %)	0.4 (0.2 %)	40.6 (19.4 %)	0.0 (0.0 %)	11.2 (5.4 %)	7.0 (3.3 %)	0.8 (0.4 %)	0.2 (0.1 %)
F0104	14.5 (6.8 %)	0.2 (0.1 %)	27.5 (12.8 %)	0.3 (0.1 %)	4.2 (2.0 %)	1.7 (0.8 %)	1.1 (0.5 %)	0.3 (0.1 %)
F0105	27.7 (12.4 %)	4.7 (2.1 %)	35.7 (16.0 %)	0.0 (0.0 %)	3.4 (1.5 %)	3.3 (1.5 %)	1.4 (0.6 %)	0.4 (0.2 %)
F0106	47.7 (20.2 %)	0.6 (0.3 %)	55.5 (23.5 %)	0.2 (0.1 %)	4.2 (1.8 %)	3.8 (1.6 %)	0.8 (0.3 %)	0.4 (0.2 %)
F0107	24.2 (11.0 %)	3.1 (1.4 %)	37.1 (16.9 %)	1.1 (0.5 %)	5.7 (2.6 %)	6.3 (2.9 %)	0.7 (0.3 %)	0.9 (0.4 %)
F0108	31.0 (13.6 %)	1.1 (0.5 %)	35.6 (15.6 %)	1.0 (0.4 %)	8.1 (3.5 %)	4.0 (1.7 %)	3.3 (1.4 %)	1.4 (0.6 %)
Id	45.3 (19.2 %)	1.1 (0.5 %)	52.7 (22.3 %)	0.0 (0.0 %)	0.0 (0.0 %)	0.0 (0.0 %)	0.0 (0.0 %)	0.0 (0.0 %)
free	8.9 (3.6 %)	3.1 (1.2 %)	17.4 (7.1 %)	0.3 (0.1 %)	1.8 (0.7 %)	1.4 (0.6 %)	0.3 (0.1 %)	0.2 (0.1 %)
part	9.7 (3.8 %)	0.6 (0.2 %)	11.7 (4.6 %)	0.6 (0.2 %)	0.3 (0.1 %)	0.9 (0.3 %)	1.2 (0.5 %)	0.7 (0.3 %)
ALL	23.9 (10.2 %)	1.9 (0.8 %)	31.2 (13.3 %)	0.4 (0.2 %)	4.0 (1.7 %)	4.1 (1.7 %)	1.2 (0.5 %)	0.7 (0.3 %)

(Please note that ‘C’ stands for content and ‘R’ for readability in the table.)

[Average number of revisions per text]

	Deletion		Insertion		Replacement	
	(unimportant)	(readability)	(important)	(readability)	(content)	(readability)
F0101	1.4	0.4	1.3	0.2	0.5	0.3
F0102	1.2	0.4	1.0	0.0	0.4	0.5
F0103	0.8	0.1	1.2	0.0	0.2	0.1
F0104	0.8	0.1	1.2	0.1	0.1	0.2
F0105	1.2	0.1	0.7	0.0	0.4	0.2
F0106	2.1	0.2	1.7	0.1	0.1	0.2
F0107	0.8	0.6	0.9	0.1	0.2	0.1
F0108	1.4	0.1	1.1	0.1	0.2	0.6
ld	1.9	0.1	1.3	0.0	0.0	0.0
free	0.6	0.4	1.1	0.1	0.2	0.1
part	0.7	0.3	1.1	0.1	0.1	0.2
ALL	1.1	0.3	1.1	0.1	0.2	0.3

[number of given-up summaries]

	number of given-up summaries
0101	5
F0102	5
F0103	16
F0104	16
F0105	14
F0106	10
F0107	17
F0108	14
ld	23
free	1
part	1
ALL	11.1

[Task B (LONG)]

[Average number of revised characters per text]

	Deletion		Insertion		Replacement											
	(unimportant)	(readability)	(important)	(readability)	(C) deletion	(C) Insertion	(R) deletion	(R) Insertion								
F0201	117.3	(16.6 %)	4.3	(0.6 %)	112.5	(15.9 %)	4.8	(0.7 %)	10.9	(1.5 %)	17.7	(2.5 %)	2.3	(0.3 %)	1.9	(0.3 %)
F0202	112.2	(15.9 %)	6.8	(1.0 %)	135.3	(19.2 %)	3.9	(0.6 %)	7.0	(1.0 %)	8.4	(1.2 %)	3.6	(0.5 %)	1.6	(0.2 %)
F0203	118.1	(16.7 %)	5.4	(0.8 %)	135.3	(19.1 %)	8.7	(1.2 %)	12.4	(1.8 %)	12.1	(1.7 %)	1.7	(0.2 %)	1.0	(0.1 %)
F0204	131.9	(18.4 %)	10.1	(1.4 %)	126.6	(17.7 %)	4.5	(0.6 %)	12.2	(1.7 %)	13.2	(1.8 %)	1.3	(0.2 %)	1.6	(0.2 %)
F0205	152.9	(21.5 %)	6.4	(0.9 %)	128.7	(18.1 %)	3.3	(0.5 %)	20.4	(2.9 %)	32.6	(4.6 %)	3.9	(0.5 %)	2.8	(0.4 %)
F0206	112.4	(22.5 %)	12.1	(2.4 %)	209.0	(41.9 %)	22.8	(4.6 %)	8.7	(1.7 %)	15.3	(3.1 %)	1.6	(0.3 %)	2.7	(0.5 %)
F0207	117.6	(16.1 %)	17.5	(2.4 %)	132.8	(18.1 %)	17.0	(2.3 %)	8.5	(1.2 %)	16.1	(2.2 %)	6.4	(0.9 %)	2.4	(0.3 %)
F0208	94.0	(13.4 %)	9.8	(1.4 %)	139.8	(20.0 %)	3.0	(0.4 %)	7.1	(1.0 %)	9.4	(1.4 %)	6.7	(1.0 %)	2.0	(0.3 %)
F0209	148.4	(21.9 %)	14.7	(2.2 %)	173.5	(25.6 %)	8.5	(1.3 %)	4.1	(0.6 %)	5.5	(0.8 %)	2.2	(0.3 %)	1.3	(0.2 %)
human	69.3	(11.4 %)	17.8	(2.9 %)	72.5	(12.0 %)	19.8	(3.3 %)	9.6	(1.6 %)	9.5	(1.6 %)	2.5	(0.4 %)	3.9	(0.6 %)
Id	102.2	(17.3 %)	17.2	(2.9 %)	132.6	(22.5 %)	6.6	(1.1 %)	9.9	(1.7 %)	11.2	(1.9 %)	2.6	(0.4 %)	1.6	(0.3 %)
stein	89.7	(15.9 %)	12.9	(2.3 %)	100.1	(17.7 %)	1.0	(0.2 %)	6.9	(1.2 %)	12.7	(2.3 %)	1.2	(0.2 %)	0.8	(0.1 %)
ALL	112.3	(16.9 %)	11.2	(1.7 %)	129.6	(19.5 %)	8.3	(1.3 %)	9.7	(1.5 %)	13.2	(2.0 %)	3.0	(0.5 %)	1.9	(0.3 %)

(Please note that 'C' stands for content and 'R' for readability in the table.)

[Average number of revisions per text]

	Deletion		Insertion		Replacement	
	(unimportant)	(readability)	(important)	(readability)	(content)	(readability)
F0201	3.8	0.7	7.2	1.4	1.1	0.9
F0202	5.2	0.6	3.5	0.4	0.7	0.5
F0203	5.1	0.6	3.8	0.5	0.9	0.6
F0204	4.2	0.6	3.4	0.7	1.4	0.7
F0205	8.1	0.6	5.4	1.7	3.0	1.3
F0206	3.2	0.2	4.7	0.7	0.8	0.6
F0207	7.0	1.1	4.1	1.1	1.1	1.1
F0208	4.8	0.7	4.0	0.4	0.8	0.9
F0209	4.6	0.5	3.9	0.5	0.5	0.5
Human	3.0	0.9	3.4	7.8	1.0	1.2
Id	5.7	0.9	2.9	0.4	0.7	0.5
stein	4.0	0.5	2.2	0.3	0.8	0.5
ALL	4.9	0.7	4.0	1.4	1.0	0.8

[number of given-up summaries]

	number of given-up summaries
F0201	0
F0202	0
F0203	0
F0204	2
F0205	12
F0206	14
F0207	2
F0208	5
F0209	8
human	2
ld	3
stein	3
ALL	4.3

[Task B (SHORT)]

[Average number of revised characters per text]

	Deletion		Insertion		Replacement			
	(unimportant)	(readability)	(important)	(readability)	(C) deletion	(C) Insertion	(R) deletion	(R) Insertion
F0201	71.4 (20.6 %)	3.1 (0.9 %)	87.2 (25.1 %)	2.3 (0.7 %)	8.4 (2.4 %)	9.2 (2.7 %)	2.4 (0.7 %)	1.5 (0.4 %)
F0202	59.9 (18.5 %)	3.2 (1.0 %)	95.9 (29.7 %)	1.8 (0.6 %)	7.0 (2.2 %)	7.5 (2.3 %)	0.8 (0.2 %)	0.6 (0.2 %)
F0203	60.5 (18.5 %)	4.6 (1.4 %)	93.1 (28.5 %)	1.1 (0.3 %)	7.6 (2.3 %)	9.3 (2.8 %)	1.0 (0.3 %)	0.5 (0.2 %)
F0204	71.0 (21.1 %)	9.4 (2.8 %)	81.7 (24.2 %)	4.8 (1.4 %)	12.2 (3.6 %)	15.1 (4.5 %)	1.3 (0.4 %)	1.2 (0.3 %)
F0205	87.7 (24.4 %)	1.3 (0.4 %)	58.1 (16.2 %)	5.2 (1.4 %)	14.2 (4.0 %)	32.7 (9.1 %)	3.4 (1.0 %)	3.7 (1.0 %)
F0206	49.4 (18.5 %)	2.8 (1.0 %)	111.2 (41.7 %)	1.4 (0.5 %)	3.9 (1.5 %)	5.4 (2.0 %)	1.6 (0.6 %)	1.1 (0.4 %)
F0207	61.3 (18.2 %)	4.3 (1.3 %)	81.9 (24.3 %)	4.5 (1.3 %)	3.0 (0.9 %)	2.0 (0.6 %)	2.0 (0.6 %)	2.4 (0.7 %)
F0208	57.5 (17.4 %)	6.0 (1.8 %)	86.3 (26.2 %)	5.2 (1.6 %)	0.7 (0.2 %)	2.1 (0.6 %)	1.3 (0.4 %)	0.4 (0.1 %)
F0209	56.7 (17.9 %)	5.4 (1.7 %)	96.6 (30.4 %)	0.8 (0.2 %)	3.0 (0.9 %)	5.0 (1.6 %)	1.3 (0.4 %)	0.6 (0.2 %)
human	35.3 (11.9 %)	4.9 (1.6 %)	30.9 (10.4 %)	8.6 (2.9 %)	7.9 (2.7 %)	10.8 (3.7 %)	2.8 (1.0 %)	1.7 (0.6 %)
ld	47.1 (17.5 %)	14.0 (5.2 %)	89.2 (33.1 %)	0.2 (0.1 %)	3.6 (1.3 %)	8.5 (3.2 %)	2.8 (1.0 %)	0.8 (0.3 %)
stein	51.4 (18.9 %)	2.1 (0.8 %)	77.5 (28.5 %)	1.6 (0.6 %)	4.9 (1.8 %)	7.2 (2.6 %)	3.0 (1.1 %)	0.9 (0.3 %)
ALL	58.7 (18.6 %)	5.2 (1.6 %)	81.8 (25.9 %)	3.1 (1.0 %)	6.4 (2.0 %)	9.2 (2.9 %)	2.0 (0.6 %)	1.2 (0.4 %)

(Please note that ‘C’ stands for content and ‘R’ for readability in the table.)

[Average number of revisions per text]

	(unimportant)	Deletion (readability)	(important)	Insertion (readability)	(content)	Replacement (readability)
F0201	3.5	0.5	4.3	0.8	1.1	0.7
F0202	3.5	0.4	2.4	0.2	0.7	0.2
F0203	3.6	0.3	2.8	0.2	0.5	0.4
F0204	2.7	0.5	2.3	0.2	1.2	0.7
F0205	5.5	0.4	2.5	0.8	2.0	0.7
F0206	2.0	0.4	3.4	0.6	0.4	0.4
F0207	3.5	0.4	2.7	0.3	0.6	0.6
F0208	2.4	0.5	2.3	0.4	0.2	0.3
F0209	2.5	0.5	2.2	0.2	0.3	0.4
human	1.9	0.8	2.4	2.0	0.9	0.7
ld	2.8	0.7	2.4	0.2	0.5	0.4
stein	3.0	0.3	1.8	0.2	0.4	0.3
ALL	3.1	0.5	2.6	0.5	0.7	0.5

[number of given-up summaries]

	number of given-up summaries
F0201	2
F0202	2
F0203	2
F0204	4
F0205	12
F0206	14
F0207	3
F0208	9
F0209	9
human	3
ld	5
stein	3
ALL	5.7