

LEXICAL PARALLELISM IN TEXT STRUCTURE DETERMINATION  
AND CONTENT ANALYSIS

Yoshiyuki Sakamoto                      Tetsuya Okamoto  
Electrotechnical Laboratory      University of Electcommunications  
Tsukuba, Japan                      Tokyo, Japan

ABSTRACT

In this paper the problem is discussed about the text structure determination and content analysis by lexical parallelism, or the repetition of lexical items. Intersentential relations are determined through the identical, partly identical or lexico-semantic repetition in Japanese scientific texts. Lexical parallelism ratio and lexical parallelism indicator distance are obtained on computer and by hand. And the application of the characteristics to automatic content analysis is discussed.

1. INTRODUCTION

Lexical parallelism, that is, the repetition of lexical items, is an important device for indicating the sentence connections in a text (discourse). The recurrent lexical items, or lexical equivalents need not have the same syntactic function or parts of speech in the two sentences in which they occur. They may be identical in form and in meaning, or they may be related by lexico-semantic relationship, such as synonymy, hyponymy, antonymy. In a special case they may be partly identical both in form and in meaning, as in 超音波 (ultrasonic wave), 音波 (sound wave) and 波 (sound).

Another device for indicating the sentence connections is a syntactic device, such as substitutes, logical connecters, time and place relaters and structural parallelism [1]. For example, in Japanese substitutes--- これ/この (this), ここ (here), われわれ (we/our), それ (it), time relaters--- 次に (next), 以上の (above mentioned), and logical connecters--- および (and), または (or), 第二に (secondly) belong to this device.

Sevbo studied lexical parallelism in normalized text, where substitutes were replaced by their lexical

equivalents and complex sentences were decomposed into successive simple sentences (clauses).

She traced the repetition patterns of lexical items in Subject/Predicate opposition. She assumes here that the syntactic subject or its dependent, direct or indirect, corresponds to "Subject (old information) of elementary thought" and the syntactic predicate or its dependent to "Predicate (new information) of elementary thought" [2].

In Japanese, sentence components occur in any positions before predicate and old information or topic is placed, as a rule, at/near the beginning of a sentence [3]. In the following discussion we analyze the repetition of lexical items in an unnormalized text without regard to their syntactic functions, parts of speech and topic/comment distinctions, assuming that the lexical equivalents at/near the beginning of the sentences function as the keywords in indicating the sentence connections and the contents of a text.

Nouns do not inflect and most verbs and adjectives have the unchanging stems and inflectional suffixes in Japanese. The important concepts and technical terms (noun, verb or

adjective stems) are written in Kanji (Chinese ideographs) or Katakana (square Japanese syllabary). Katakana is used to transcribe foreign technical terms. Hiragana (Japanese cursive syllabary), on the other hand, is used to write post-positional particles and suffixes, denoting case, topic, mood, tense aspect etc. In view of these facts we define lexical items as a word or phrase in Kanji and Katakana.

We have studied lexical parallelisms in a short tale[4], in technical and scientific texts[5,6], based upon Sevbo's approach. The purpose of the present paper is to obtain the characteristics of lexical parallelism in Japanese technical and scientific texts and to explore the possibilities of utilizing these characteristics for automatic content analysis.

Five text samples are used for experiment and discussion. They are the essays on "Ultrasonic amplification"(Text A), "Brain and automaton"(Text B), "Petrochemical industry"(Text C), "Chemical industry in Japan"(Text D) and "Between organism and inanimate matter"(Text E).

## 2. LEXICAL PARALLELISM RATIO

The sentence connection of type  $t$  in position  $w$  is determined between the given  $j$ -th sentence  $S_j$  and the  $i$ -th sentence  $S_i$  ( $i < j$ ), if and only if  $S_i$  is the nearest preceding sentence which contains the lexical item, lexically equivalent to the  $w$ -th lexical item from the beginning of the given sentence  $S_j$  through the type  $t$  repetition ( $t = 1, 2, 3$ ;  $w = 1, 2, 3, 4, 5$ ).

The repetitions of type 1, 2, 3 correspond to the identical, partly identical, lexico-semantic repetitions, respectively.

The lexical equivalents in  $S_j$  and  $S_i$  are called lexical parallelism indicators, and  $S_j$  is called a dependent on  $S_i$ .

Lexical parallelism ratio of type  $t$  in position  $w$  is defined as follows;

$$\alpha_w^t = (n / N - 1) * 100$$

where  $n$  is the number of the determined connections in a text;  $N-1$

is the determinable maximum number of the sentence connections in a text,  $N$  being the total number of the sentences in the text;  $t$  is type of lexical repetition and  $w$  is the position, i.e. the sequence number from the beginning of the sentence.

The experiments were carried out to obtain the characteristics of the lexical parallelism in sample texts on computer and by hand.

In computer experiment lexical items, i.e. the sequence in Kanji or Katakana, were identified and segmented by machine character codes without syntactic and morphological analysis. Then the sentence connections of type 1 (identical repetition) are determined in each position and lexical parallelism ratios are obtained (Table 1). On the same samples the optimal sentence connections are determined manually and the lexical parallelism ratios were calculated (Table 2). Except for Text E, the totals of the ratios amount to 72-83% (cf. Table 2) and in computer experiment the ratios of type 1 in the initial position amount to 57-68% (cf. Table 1). And moreover, the initial lexical items ( $w=1$ ) show the maxima in most samples in Table 1 and by far the highest value in all samples in Table 2, and they decrease with increasing  $w$  in Table 2. It is clear from the results that lexical parallelism plays an important role in the intersentential dependency and lexical items at the beginning of the sentences are the most relevant lexical parallelism indicators.

## 3. LEXICAL PARALLELISM INDICATOR DISTANCE

As an example, intersentential dependency determined manually in Text A, which is the essay on "Ultrasonic amplification" with 123 sentences in four paragraphs, is shown in Table 3 and Figure 1. The lexical parallelism indicator distances are shown as well. Lexical parallelism indicator distance is defined as follows:

$$D = \frac{t}{w, j} - i$$

where  $D$  is lexical parallelism indicator distance;  $t$  is type of

lexical repetition;  $w$  is position of the lexical indicator;  $i$  and  $j$  are sequence numbers of the governor sentence and dependent sentence respectively.

The distance is supposed to represent the semantic extent of the lexical parallelism indicators, or better the concepts referred by them.

In Figure 1 a diagonal unit distance line indicates the hypothetical situation, where every sentence depends on the immediately preceding sentence. Data show a tendency to distribute near this line in all samples.

Lexical parallelism indicators show the progress of the author's thought in the text in Table 3. Sevbo pointed out the significance of the indicators with large  $D$  in indicating the contents of paragraphs and texts. The lexical items with large  $D$  are supposed to be the important topics, to which the author of the text returns after commenting on another topics. In the example the items with large  $D(D>10)$  were shown in Figure 2.

These indicators are distributed among paragraphs. For example, the indicator 超音波 (ultrasonic wave) extends over 15 sentences (from 9th to 24th) within paragraph 2, which ranges from 2nd to 40th sentence, and the indicator 進行波管 (traveling-wave tube) extends over 22 sentences (100th-122nd) within paragraph 4 (85th-123rd) as well. The indicator 進行波增幅 (traveling-wave amplification) covers paragraph 3 completely, ranging from the 41st sentence, or the first sentence of the paragraph, through the 67th sentence to 85th sentence, or the first sentence of the next paragraph. In short, these indicators divide the text into the three paragraphs.

In addition, they reflect appropriately the contents of paragraphs in the sample text, as suggested by the fact that they are partly identical with the following paragraph names: "Introduction" (paragraph 1), "What is the ultrasonic wave?" (paragraph 2), "Microwave and traveling-wave tube" (paragraph 3) and "Ultrasonic wave and traveling-wave amplification" (paragraph 4).

These data suggest that the

indicator with large  $D$  may be useful as keywords to the contents of a text.

#### 4. CONCLUSION

Lexical parallelism plays an important role in the intersentential dependency, or text structure and lexical items at the beginning of the sentences are the most relevant lexical parallelism indicators.

The initial lexical parallelism indicators with long lexical parallelism indicator distances reflect the contents of paragraphs and may be useful keywords in information retrieval.

The partly identical repetition and lexico-semantic repetition through the lexical items at/near the beginning of the sentence, firstly, intersentential dependency by syntactic device, secondly, the recognition of topic/comment opposition in the sentence, thirdly, and lastly, the application to automatic keyword or key-sentence extraction in content analysis depend on the future researches.

#### REFERENCE

- [1] Quirk, R., Greenbaum, S., Leech, G., and Svartvik, J., A grammar of contemporary English (Longman, London, 1972).
- [2] Sevbo, I.N., Struktura svjaznovo teksta i avtomatizatsija (Nayka, M., 1969).
- [3] Makino, S., Grammar of repetition (Taisyukan, Tokyo, 1980).
- [4] Okamoto, T., Text structure determination and content analysis by lexical parallelism, Proceeding of the Univ. of Electro-Communications, vol.24(1973), no. 1, 177-190.
- [5] Okamoto, T., Structure analysis of Japanese text, Mathematical linguistics, No.62(1972), 1-11.
- [6] Sakamoto, Y., Okamoto, T., Yatsu, N., Text structure and a model of discourse understanding by lexical parallelism, Proceeding of the 10th annual meeting on information science and technology, (1973), 55-64.

Table 1 Lexical parallelism ratios of type 1 in computer experiment(%)

T \ w	1	2	3	4	5
A	60.4 (75)	61.9 (75)	57.1 (64)	54.2 (58)	56.4 (57)
B	68.2 (71)	64.4 (67)	56.3 (58)	58.4 (59)	57.4 (58)
C	59.4 (41)	45.5 (31)	43.2 (29)	37.5 (24)	32.2 (19)
D	57.2 (71)	61.2 (76)	54.9 (67)	52.5 (60)	56.7 (58)
D	41.1 (37)	53.3 (48)	49.4 (43)	42.1 (35)	50.0 (40)

Note: T - sample texts, w - sequence numbers of indicators, values in() are numbers of determined sentence connections.

Table 2 Lexical parallelism ratios determined by hand(%)

T \ N-1 \ w	1	2	3	4	5
A 122	60.7 (74)	6.6 (8)	3.2 (4)	0.8 (1)	0.8 (1)
B 103	68.9 (71)	9.7 (10)	1.9 (2)	0.9 (1)	0.9 (1)
C 69	50.7 (35)	8.7 (6)	13.0 (9)	2.9 (2)	0 (0)
D 123	54.9 (67)	13.9 (17)	2.4 (3)	1.6 (2)	0 (0)
E 89	29.2 (26)	5.6 (5)	2.2 (2)	1.1 (1)	0 (0)

Note: N-1 --- the determinable maximum number of intersentential relations.

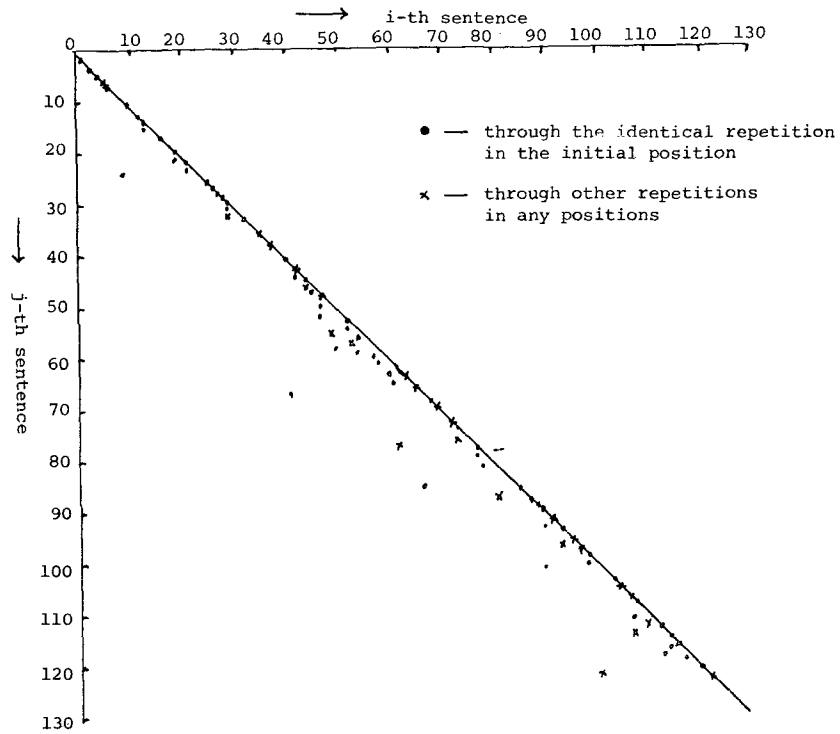


Figure 1 Lexico-semantic intersentential dependency graph in sample text A

Table 3 Lexico-semantic intersentential dependency in sample text A

Indicator	J	I	D	w	t	Indicator	J	I	D	w	t
音の (sound)	1	-	-	-	-	エネルギー (energy)	63	60	3	1	1
音波 (sound)	2	1	1	1	1	位相速度の (phase velocity)	64	63	1	1	2
超音波 (ultrasonic wave)	3	-	-	-	-	の (possessive particle)	65	61	4	1	1
聞 (hear)	4	3	1	1	1	振幅 (amplitude)	66	65	1	2	1
超音波 (ultrasonic wave)	5	4	1	1	1	進行波増幅の (traveling-wave amplification)	67	41	26	1	1
一 (one)	6	5	1	1	2	このこと (this fact)	68	-	-	-	-
第二 (the second)	7	6	1	1	2	電波 (radio wave)	69	68	1	1	1
この場合 (this case)	8	-	-	-	-	波長 (wave)	70	69	1	1	2
ここでは (here)	9	-	-	-	-	電気信号 (electric signal)	71	-	-	-	-
はじめに述べたように (as mentioned before)	10	-	-	-	-	電界 (electric field)	72	-	-	-	-
音 (sound)	11	10	1	1	1	電子 (electron)	73	72	1	5	1
今 (now)	12	-	-	-	-	電子 (electron)	74	73	1	1	1
笛から (whistle)	13	12	1	1	1	これ (this)	75	-	-	-	-
波長より (wavelength)	14	13	1	1	1	外側の (outside)	76	73	3	2	3
波 (wave)	15	13	2	1	1	波 (wave)	77	62	14	3	1
このこと (this fact)	16	-	-	-	-	Ve	78	77	1	1	1
波動的 (wave motion)	17	16	1	1	1	(1)	79	77	2	1	1
われわれ (we)	18	-	-	-	-	これ (this)	80	-	-	-	-
われわれ (our)	19	-	-	-	-	(2)	81	78	3	1	1
波長 (wavelength)	20	19	1	1	1	これ (this)	82	-	-	-	-
われわれ (our)	21	19	2	1	1	マイクロウェーブ通信 (microwave communication)	83	-	-	-	-
目で (eyes)	22	21	1	1	1	物理学 (physics)	84	-	-	-	-
音 (sound)	23	21	2	1	1	進行波増幅 (traveling-wave amplification)	85	67	18	1	1
超音波 (ultrasonic wave)	24	9	15	1	1	二つの (two)	86	85	1	1	1
あなた (you)	25	-	-	-	-	エネルギー源として (energy source)	87	81	6	1	2
コウモリ (bat)	26	25	1	1	1	電気系 (electric system)	88	87	1	1	1
コウモリ (bat)	27	26	1	1	1	圧電現象 (piezoelectric phenomenon)	89	88	1	1	1
コウモリ (bat)	28	27	1	1	1	圧電現象 (piezoelectric phenomenon)	90	89	1	1	1
レーダー (radar)	29	28	1	1	1	あなた (you)	91	-	-	-	-
共に (together with)	30	29	1	2	1	ピックアップ (pick-up)	92	91	1	1	2
レーダー (radar)	31	29	2	1	1	圧電結晶 (piezoelectric crystal)	93	90	3	1	1
音波 (sound wave)	32	29	3	2	2	これ (this)	94	93	1	1	1
超音波 (ultrasonic wave)	33	32	1	1	1	逆効果 (reverse effect)	95	-	-	-	-
例えば (for example)	34	-	-	-	-	リコーパ (recoiler)	96	95	1	1	2
途中 (halfway)	35	-	-	-	-	圧電結晶 (piezoelectric crystal)	97	93	4	3	1
魚群探知機と (fish detector)	36	35	1	1	3	圧電結晶 (in piezoelectric crystal)	98	97	1	1	2
医学的 (medicine)	37	-	-	-	-	電界 (electric field)	99	98	1	1	1
診断には (diagnosis)	38	37	1	3	1	圧電結晶 (piezoelectric crystal)	100	98	2	1	1
新発見 (introduction)	39	-	-	-	-	ロッセル塩 (Rochelle salt)	101	90	11	1	1
既 (discussion)	40	-	-	-	-	このようなとき (such time)	102	-	-	-	-
進行波増幅 (traveling-wave amplification)	41	40	1	1	1	圧電半導体 (piezoelectric semiconductor)	103	-	-	-	-
第一図 (figure 1)	42	-	-	-	-	圧電半導体 (piezoelectric semiconductor)	104	103	1	1	1
振り子 (pendulum)	43	42	1	4	1	CdS (CdS)	105	104	1	2	1
バネ (spring)	44	42	2	1	1	それ (it)	106	-	-	-	-
の (possessive particle)	45	44	1	1	1	CdS結晶 (CdS crystal)	107	106	1	1	2
の (possessive particle)	46	44	2	2	1	増幅器 (amplifier)	108	107	1	1	1
の (possessive particle)	47	45	2	1	1	第三図 (figure 3)	109	-	-	-	-
の (possessive particle)	48	47	1	2	1	装置 (equipment)	110	-	-	-	-
このようなこと (such thing)	49	-	-	-	-	光 (light)	111	107	4	1	1
図 (figure)	50	47	3	1	1	加速電圧 (accelerated voltage)	112	110	1	1	2
このこと (this fact)	51	-	-	-	-	減衰量 (attenuation quantity)	113	112	1	1	1
の (possessive particle)	52	47	5	1	1	結晶 (crystal)	114	107	7	1	2
の (possessive particle)	53	52	1	1	1	電圧 (voltage)	115	114	1	1	1
から (from)	54	52	2	1	1	出力 (output)	116	115	1	1	1
検出 (recognition)	55	49	6	1	2	結晶 (crystal)	117	114	3	1	1
エネルギー (energy)	56	54	2	1	1	減衰 (attenuation)	118	113	5	1	1
振幅 (amplitude)	57	53	4	3	1	加速電圧 (accelerated voltage)	119	117	2	1	1
図 (figure)	58	50	8	1	1	以上の (above mentioned)	120	-	-	-	-
の (possessive particle)	59	54	5	1	1	増幅器 (amplifier)	121	120	1	1	1
エネルギー (energy)	60	57	3	1	1	進行波管 (traveling-wave tube)	122	100	21	2	1
図 (figure)	61	58	3	1	1	誕生 (birth)	123	122	1	1	2
このことから (this fact)	62	-	-	-	-						

No te: 1) English equivalents are shown in ( ); 2) underlined Hiragana sequences are postpositional particles, denoting topic, case, contrast, etc.; 3) hyphen means that J-th sentence was not connected with any preceding sentence by lexical equivalence.

Symbols: J, I -- sequence numbers of the dependent sentence and governor sentence respectively; D -- lexical parallelism indicator distance; w -- sequence number of the lexical indicator from the beginning sentence; t -- type of lexical repetition, 1, 2, 3 -- identical, partial, lexico-semantic respectively.

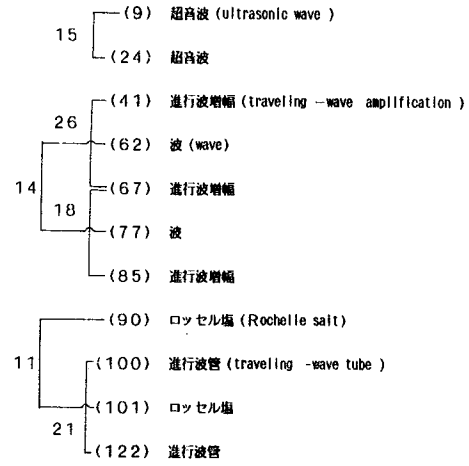


Fig. 2 Distribution of long distance indicator ( $D > 10$ )

Note : numbers in ( ) correspond to the sequence numbers of the sentences, the numbers on the lines to the distances.