

DEVELOPMENT OF BASIC PRACTICAL TECHNIQUES FOR JAPANESE LETTER
STRING PROCESSING - AUTOMATIC KEYWORD EXTRACTION AND AUTOMATIC
READING

K. Araki, K. Hinatsu, K. Itayama, T. Sahara, Y. Sakagami and
F. Takano

The Japan Information Center of Science and Technology (JICST)
2-5-2 Nagatacho Chiyodaku Tokyo 100 Japan

Japanese is a peculiar language among the thousands of languages in the world. There exist only two of the same class; Japanese and Korean. Japanese is written both in Chinese characters (ideograph) and in Kana (Katakana and Hiragana - phonetic symbols) in mixture without any space. Moreover, Chinese characters in Japanese have, in most cases, several readings and play several roles depending on the context and letter string characteristics. So for written Japanese, it was very difficult to segment letter string and extract adequate terms from sentence and to give them correct readings automatically, which are indispensable for terminology, automatic reading, automatic indexing, key-boarding from on-line terminals otherwise more than 2,000 character key-board is necessary.

The authors invented efficient algorithm and developed computer programmes and dictionaries for successful solution of the problems above for the first time in Japan.

The system consists of two subsystems called K-KACS (Kanji-Kana Automatic Conversion System) and JAKAS (Japanese Keyword Automatic Selection).

Some Chinese characters act both as suffix, preffix or preposition and as parts of meaningful words. We comprehensive-

ly collected such characters (about 500) and those terms in which the characters are included not as fixes or prepositions but as important part (about 8000 words). Letter string which is matched with dictionary term is passed but the letter remained and coincides with the special character itself is cut.

In case of long letter string without such special letter, sentence is cut by those terms of dictionary which are thought to be definite within reasonable amount. That is:

dog liver nucleus DNase
大肝臓細胞核DN/T-t
indefinite type of word. definite type of word.

Equally, among the variety of readings - in some cases more than 8 - some are special and definite and others are indefinite but obey to rules. We collected these special readings (about 25,000) for about 2,000 Chinese characters and developed algorithm and programme to select the correct reading for each Chinese character with the precision higher than 99.94 %.

As the dictionary is small enough and logic is simple, implementation and maintenance are relatively easy and the speed is very high.

JICST adopted this system for its information file production and services of more than 400,000 citations per year and save costs.

By the development of the techniques, processing of Japanese has become to be able to cope with western languages. We were awarded for the work The Prize of Learning of Japan Association of Information and Documentation in 1980, and have applied patent (Japan Patent Kokai Showa 55 (1980) - 102074).

Information File of JICST

000167 81/11/05

- ☐ A010 731455001 P 10842810425810506 1101447
 ☐ B010 1B B 03030 G 534.2-7/-8
 ☐ B020 2D C 06020 V 551.483
 ☐ B510 1E 05072 G 08061 P 0303
 ☐ C010 a105 0100EN 800978514 E 759 B D Z U S A 80 Oceans Oceans
 ☐ C020 20 ④ 1980
 ☐ C030 ④ 127-131
 ☐ D010 I Univ. New Hampshire Univ. New Hampshire
 ☐ E010 ICOX P COX PG
 ☐ E020 IHARVEY P HARVEY PG
 ☐ E030 IRENTIS P RENTIS PG
 ☐ E040 ISIVAPRASAD K SIVAPRASAD KE
 ☐ E050 IYILDIZ A YILDIZ A
 ☐ E060 IYILDIZ M YILDIZ M
 ☐ F010 Sound propagation in a shallow water region overlying a viscoelastic halfspace
 ☐ G010 粘弹性半空間上の浅海領域の音波伝播
 ☐ H010 (ネンダンセイハシクウカジョウ ノ センカリョウイキ ノ オンパデンパン) ←
 ☐ I010 Green関数と基準モード展開法とを用いて浅海の音波伝播を調べた。理想流体の境界が上面の圧力零、下面は粘弹性固体で、流体の音速は深さに依存するとして解析した。バルト海の実験海域のパラメータにより数値計算し、送受波器間水平距離と伝搬損失との関係をプロットした。結果は実験と定性的に一致し、海底の組成が音波伝播に重要なことがわかった。
 ☐ J010 000313 * 音波伝播④オンパデンパン
 ☐ J020 000840 * 水中音響④スイチュオング・ウ
 ☐ J030 035336 * 海底④カイティ
 ☐ J040 000999 静弾性④ネンダンセイ
 ☐ J050 000302 音速④オンソク
 ☐ J060 021666 水深④スイシン【カサ】
 ☐ J070 010455 依存性④イゾンセイ
 ☐ J080 007946 バルト海④バルトカイ
 ☐ K001 015216 01 波動伝播④ハドウデンパン
 ☐ K002 046038 02 伝播④デンパン
 ☐ K003 000292 01 音響学④オンキョウガク
 ☐ K004 002882 02 物理学④ブツリガク
 ☐ K005 004281 03 科学④カガク【science】
 ☐ K006 013895 01 位置④イチ
 ☐ K007 000427 01 機械的性質④キカイテキセイツ
 ☐ K008 036399 01 弾性波速度④ダンセイヒソクド
 ☐ K009 025003 02 波の速度④ナミノソクド
 ☐ K010 000951 03 速度④ソクド【スピード】
 ☐ K011 014185 01 深さ④フカサ
 ☐ K012 014135 02 長さ④ナガサ
 ☐ K013 042740 03 狹向学量④キカガクリョウ
 ☐ K014 031902 01 北東大西洋④ホクトウタイセイヨウ
 ☐ K015 031801 02 北大西洋④キクタイセイヨウ
 ☐ K016 031900 03 大西洋④クイセイヨウ
 ☐ K017 020834 04 海洋名④カイヨウメイ
 ☐ L010 浅海④ンカイ
 ☐ M010 粘弹性半空間④ネンダンセイハシクウカン
 ☐ M020 浅海④ンカイ
 ☐ M030 領域④リョウイキ
 ☐ M040 音波伝播④オンパデンパン
 } reading
 } free human indexing
 } automatic indexing from Japanese title

000041 81/11/05

- ☒ A010 890390021④ G20740810421810506④1039715☒
☒ B010 1WC03020T④621.785☒
☒ B510 1G03031M03035☒
☒ C010 a103 0040RU 800770829R143AAJ SUN80 Izv Vyssh Uchebn Zaved Chern Meta
11④Izv Vyssh Uchebn Zaved Chern Metal☒
☒ C020 0368-0797 I V U M A 1④0363-0797* ☒
☒ C030 ④④8④84-88☒
☒ E010 ОЛЬШАНСКИЙ В МОЛЬШАНСКИЙ В ☒ original title
☒ E020 ГРИНБЕРГ В ЯФГРИНБЕРГ В Я ☒ translated title
☒ F010 Определение оптимальной калорийности топлива при нагреве.☒ automatically transl.
☒ G010 加熱の際の燃料の最高発热量の決定☒
☒ H010 (カネツノサイネンリョウノサイトキハツネンリョウノケッテイ)☒ title (reading)
☒ I010 金属の加熱条件、熱交換及び相対的熱損失が知られている時に、消費される燃料の価格を最小にできる加熱炉内で
の燃料の最高発热量を決めるアルゴリズムを提案した。最高条件で金属を加熱するがについては、純粋なガスを
利用する際に消費燃料の価格を最低にできることを明らかにした☒
☒ J010 004549④ 金属材料④キンゾクザイリョウ☒
☒ J020 000399④ 加熱④カネツ【ヒーティング】☒
☒ J030 003507④* 燃料④ネンリョウ☒
☒ J040 010692④* 热量④ネンリョウ☒
☒ J050 035795④ 热処理条件④ネンショリジョウケン☒
☒ J060 001131④ 伝熱④デンネツ☒
☒ J070 003501④ 热効率④ネツコウリツ☒
☒ J080 043197④ エネルギーコスト④エネルギーコスト☒
☒ J090 000652④* 最適化④サイトキガ☒
☒ J100 000180④ アルゴリズム④アルゴリズム☒
☒ J110 003493④ 気体燃料④キタイネンリョウ☒
☒ J120 011320④ 純度④シュンド☒
☒ K001 042893④01④ 条件④ジョウケン☒
☒ K002 000614④01④ 効率④コウリツ [efficiency]☒
☒ K003 002197④01④ 原価④ゲンカ☒
☒ K004 046014④01④ 改変④カイヘン☒
☒ K005 003507④01④ 燃料④ネンリョウ☒
☒ K006 046017④01④ 度④ド☒
☒ M010 加熱④カネツ☒
☒ M020 燃料④ネンリョウ☒
☒ M030 最適④サイトキ☒
☒ M040 発热量④ハツネンリョウ☒
☒ M050 決定④ケッテイ☒
- reading ↗
- human indexing }
up-word pasting } by thesaurus
automatic indexing from Japanese }
title