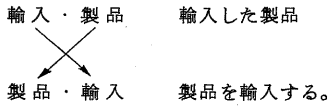
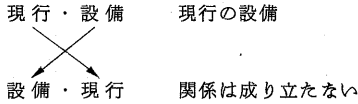


例2



十数万件の知識データが整理されて集められているが、このデータを利用することにより知識データを増加させることが出来る。但し例1, 例2のように全ての場合に逆転が成立するわけではない。

例3



逆転したデータが知識データ・ファイルに有るか否か事前に調べ、有るものは省き、残りを分析すればよい。しかし調べるのは人間であり、調査員の語彙能力、分野の知識に大きく依存する。調査員の選択が重要である。約2割程度の増加が見込める。

② 助詞、助動詞を利用し、KWICを用いて知識データを抽出する。

助詞、助動詞としては次のものを考えている。
が、を、に、へ、と、から、より、により、の、する、した、に、に対する、に、関する、……

例

- 安全を守る。
- 安全を確保する。
- 方法を拡張する。
- 機能を拡張する。

のポイントとして、工期・経済性・安全を考慮し、技術革新、環境保全、安全を守るためにも、環境を汚染から守ることを重視することが、健康と安全を確保するとともに企業目的にもかなうことと、改良する様に開発されてきた多くの項目を確保する試験機関。
UL 日常生活の安全
作業の安全
守るためには管理が不可欠であり、環境を考慮したシステムを紹介。
5にあたっては、人間・機械系全体の安全を考えて設計し、製作することが大切であり、通路の歩行、運搬の安全

り、特異な制約問題にたいする雑音の推定を拡張させる定理を定式化する。
Seidel法を拡張し、その収束性を吟味する方法を開発
幾何学的計画法を拡張し、多項式 (posynominal) 拡張しコンテナ基地を開設し、クレーンも拡張した。
それと本年航行期までに倉庫を拡張したものである。
ゼロ和二人ゲームを拡張した外そうテーブルを用いて加速試験
ゲーム1の機能を含むと同時に大幅に機能を拡張した新しいインデックス法により
品質管理のために、ラムダサンプリング法を拡張して、より強力な結果を与える。
"クセス法のシリンドライディングスの概念を拡張して、より強力な結果を与える。
これ
introl3_ [1] (' 65) の方法を拡張して求めた。
制御問題において、よく知られた結果を拡張することによって、新しく導かれた形
に必要な一般的な定理を Finch の方法を拡張することにより得ている。
より一般的な問題へ、求めた結果を拡張することが可能である。
各種因子の影響に関する研究、収束範囲を拡張する試みを行なった。
過去の歴史的データに基づいた報告データを拡張する必要がある。

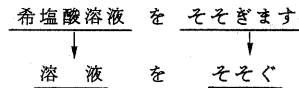
図1 助詞、助動詞を中心としたKWIC

KWICは機械的に作成することができる。このKWICの例で説明すると“を”の前の“安全”“を”の後の“拡張”という語を別々に抽出しておき、プログラムで前接語、後接語のテーブルにより機械的に知識データを抽出することができる。また、この抽出されたデータを集約することにより頻度情報も得られる。

この方法による語と語の関係の抽出は全て機械的に出来るのではなく、次の点に注意しなければならない。

- 長単位用語から基本概念語を抽出しなければならない。

例



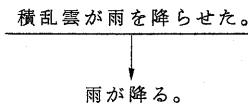
- 動詞を終止形にしなければならない。
- 助詞、助動詞等の直前、直後の語が必ずしも語と語の関係を持つものではない。

例

~方法をつぎつぎに実行した。
ごみ以外を集めたもの。

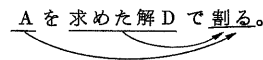
- 文中に現われた格助詞は表層的なもので語と語の関係の適切に表わしていない場合がある。

例



格助詞を変えなければならない場合がある。
使役、受身(自発、尊敬)は注意しなければならない。

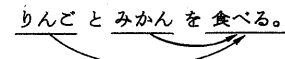
- 語の係り受け関係が少し複雑な場合



“Aを求めた”ではなく“Aを割る。”を抽出しなければならない。

- 並列関係

例



直前直後の関係だけでは“りんごを食べる”は抽出できない。

- 機械的に処理できなかったデータをリストし、手作業で分析しなければならない。

幾つかの問題点、例外はあるがこれらは発生する割合が少ないので、手作業による修正を行うとしている。この方法によって“を”を中心としたKWICを作り分析中である。対象としたデータは日本科学技術情報センターの抄録文であり、データ量は約79万件(KWICの行数)である。

このデータの中から重複をまとめ、約20万件程度の知識データが得られる予定である。“を”以外の助詞、助動詞へも拡大する予定である。

79万語の“を”を中心としたKWICから前接語、後接語を抽出する方法を考えなければならない。

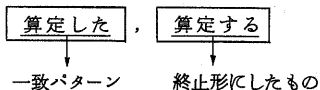
前接語の抽出方法

- (i) KWICの前接語が含まれている部分の文字列から漢字、片仮名、記号()の部分文字列を抽出する。
- (ii) (i)で抽出した文字列を最後の文字から分類し、同一のものは集計し、まとめ、頻度情報を付ける。
- (iii) (ii)で集約したDataのうち頻度2以上のものをリストする。(5万行程度になる)
- (iv) (iii)でリストした中から基礎的概念語を手作業で抽出する。(約2万語程度になる予定)

後接語の抽出方法

- (i) KWICの後接語が含まれている部分の文字列から8文字を抽出する。
- (ii) 最初の文字から漢字コード順に分類し、同一のものはまとめ1件にする。約8万行程度になる。
- (iii) 集約したDataからKWICのパターンと一致させる文字列と終止形に変形する文字列を抽出する。

例



前接語、後接語をテーブルに納めプログラムで語と語の関係を抽出する。抽出できないものは2つのテーブルを順次増やすことにより処理することができる。また処理できないものはリストし手作業によってテーブルを改良しなければならない。

前接後、後接語のための作業テーブルを次に示す。

両理論	3
改良理論	2
乱子乱流理論	2
計量理論	2
運動量理論	2
混合距離理論	2
行列理論	14
反応行列理論	2
核分裂理論	2
配向緩和理論	2
確率論	7
化学量論	16
記録	17
メンテナンス記録	2

表2 前接語作業表の一部

与えるか。	7
与えるかという	3
与えるかどうか	5
与えるかについて	17
与えるかもしれない	4
与えるかを確立	3
与えるかを検討	7
与えるかを考察	6
与えるかを調べ	13
与えるが、5 -	3
与えるが、この	11
与えるが、これ	11

表3 後接語作業表の一部

この方法は機械的方法が多いため試行錯誤が行える。費用が安く、処理速度が早いという特徴がある。全体的処理図は図2に示す。

③ 3文字漢字列の分割による知識データの獲得

4文字漢字列の分割と同様に3文字漢字列を分析することにより知識データを得る。3文字漢字列を分析すると次のようになる。

- (1) 3文字で1つの語となり分割が不可能なもの。

例 不可決、不思議

- (2) $A(B_1 + B_2)$ のように分解され、 AB_1, AB_2 が意味のある語となるもの

例 輸・出入、 国・内外

- (3) (A1 + A2) のように分解され A1B, A2B が意味のある語となるもの。

例 入出・力 上下・院

- (4) SA, AS A が主要な語で、それ等に接頭語(S), 接尾語(S)が付いたもの。

例 未・解決, 近代・化

- (5) AB A, Bともに主要な語でそれが結合したもの。

例 金・鉱山, 障害・児

- (6) A・B・C A, B, Cが主要な語で、それらが並列的に結合したもの。

例 年・月・日, 市・町・村, 英・独・仏

大きくわけると上記6つのようになる。これに各語が持つ品詞成分を付け加えることにより、さらに細分類することができる。この考え方は東京女子大学水谷静夫の分析方法を用いている。このような方法を基礎にして延 169,163 件、種類 33,155 件のデータを分析中である。機械的に収集した 3 文字漢字列であるため、この中から有効なデータを抽出することができる。

④ 慣用語による知識データの収集

日本語の中には慣用的に用いられる句が多い。これらを構成している一語、一語を取りあげて、その意味を考えても全体の句の意味には結び付かない場合が多い。意味の深い部分に立ち入れば、又、語の連想から各語の意味を結合すれば全体が合成できるものもある。しかし、これを機械的に行うには、かなり無理な面がある。

例 反吐が出る → 嫌悪感を感じる。
鼻が高い → 自慢する。

慣用語をおきかえ可能な語に変えることを考える。このような方法を取れば機械翻訳等に便利である。

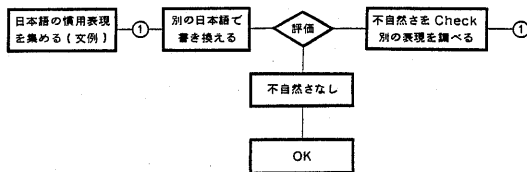


図3 慣用語の表現変え

慣用語表現を機械的に収集するにはどのようにすればよいかという問題がある。そこで語に定義された数個の意味だけを与え(あまり深い意味はもちいない。), これを基準にして個々の語の意味から合成されるものと慣用語表現の持つ意味を照合すれば判断することができる。又、慣用語表現は語と語の結び付きが強いものであるから語と語の共起と考えられる。このことを利用し機械的に抽出することができる。

慣用語表現は既に多くの人々が慣用語の研究として出版物にまとめている。ここでは出版物を利用し研究の材料とする。

慣用語表現を分析するために、次の二つのファイルを作成する。

(1) 慣用語表現ファイル

このファイルは慣用語表現と、その属性を入力する。

(2) 慣用語表現例文ファイル

この例文ファイルは例文と慣用語表現に対応する NO を持っている。

これら2つのファイルの項目は次のようなものを予定している。

• 慣用語表現ファイルの項目

1. 登録 NO (SEQ. NO)
2. 慣用語表現
3. 読み
4. 分かち書き
5. 自立語の見出し
6. 文、句の区別
7. 品詞、又は品詞相当のもの
8. 活用
9. 語尾
10. 置換可能語
11. 置換にあつての文字列操作1(前)
12. 置換にあつての文字列操作2(後)
13. 訳語
14. 共起語
15. その他区分

• 慣用語表現例文ファイルの項目

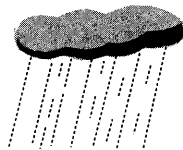
1. 登録 NO
2. 例文
3. 慣用語表現ファイルの対応 NO
4. その他

置換可能語は慣用表現をできるだけ簡単な表現に書き換えるためのものである。

慣用表現を置換可能な表現に変えた場合、追加したり、削除する文字が有るか否か、有ればその手順は常に一定か、例文で検証する。不自然であれば置換可能な語に変える。

このほか、自然言語の表現の中にはある特定の現象や事象については特定の言葉を用いるということがある。語には共起性の強いものがある。

例えば



「雨が降る。」を
考えてみると
「雨」のかわりに
「水、水滴、H₂O、
雨水」に変える表現は適切でない、また「降る」のかわりに「落ちる、落下する」と変えても良い表現ではない。

「雨が降る。」は一種の慣用表現とみなすことができる。

計算機に言語を処理させるにはこのような語と語の共起についても多量のデータを集めておかなければならない。

以上4つの新しい知識データの獲得方法について述べた。この4つの中で②のKWICを用いて知識データを得る方法が有効な手段と考えられる。

4. 収集した知識データの分析と利用方法

日本科学技術情報センターの抄録文を分析し、表1の知識データを得た。知識データの収集過程における各種資料から判断し、4文字漢字列の延80%に当る漢字列を収集しているとみなすことができる。(抄録文に関して)この知識データは応用研究に充分対応することができる。

① 4文字漢字列の分析内容

4文字漢字列として使用頻度の高いものを示すと表4のようになる。また、前接漢字列の種類は約4,700種類、後接漢字列の種類は約3,200種類であり、その使用頻度の高いものをあげると表5、表6のようになる。この収集した知識データを見易くするために1つの語についての接続状態を調べてみると表7のような結果を得た。これにより各語の接続の関係が明らかになる。

② 概念の結合関係の分析

概念の結合関係、概念と動詞の結合関係について分析することはこの研究の一つの課題である。前接語、後接語を幾つかの概念(時間、場所、空間、道具、…等)に分類し、これらの結合関係の強弱を見つけ出すことである。

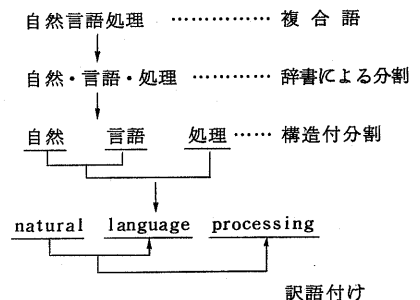
このためには次のような方法を考えている。前接語、後接語に類語新辞典(角川)のコードを割り付ける。同一コードの組を集めることにより、語のグループができる。これによりカテゴリー化ができる。類語新辞典の分類は10進分類を基本とし、1,100のカテゴリーに分類されている。それゆえ、語と語の関係は $1,000 \times 1,000 = 100$ 万のカテゴリー空間のどこかに入る。約28,000件のデータでは2%程度の部分空間に入る。17万件のデータでは約10数%程度の部分空間を満す。(重複があるため。)

これにより概念の結合関係がわかるし、今後新しく出現する語と語の結合を予想することもできる。

③ 複合語の構造付分割と訳語への利用

複合語の翻訳にあたっては、翻訳の基礎となる複合語の分割が重要である。ただ単純に分割するのではなく構造を持った状態に分割することが必要である。

例1



この構造を持った漢字列の分割のためには次の5つの方法が考えられる。

- 1) 最小語数に分割する。
- 2) 分割された各語の係り受けは非交差である。
- 3) 語と語の知識データを使う。
- 4) 語の持つ特殊性を利用する。
(以上、以前、等、中、…)
- 5) 複合語の分割パターンの頻度を利用する。
等が考えられる。

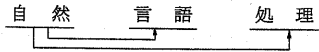
例1の場合には知識データとして次の2つの知識データが有ることを仮定している。

知識データ：言語(ヲ)処理(スル)

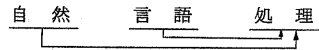
自然(ナ)言語

もし、これらの知識データがなければ次のような色々な構造が考えられる。

例2



例3



複合語の翻訳にあたっては正しい構造に分割することが第一であるが、次に分割された各語がどのような訳語を取るかということが重要である。(この考え方の前提は分割された語が1対1の訳語を取るという場合である。)多量の複合語に訳語を付けるということについては研究の第一歩を踏み出した状況である。

④ 格構造の部分構造を抽出する。

格構造の分析のために知識データを利用することを考える。日本語の機械処理では格構造を結合価文法で表わしたものが多く用いられている。これは文の構造を把握するために重要な役割をはたしている。しかし、現段階では用語に対する文型パターンが少なすぎる。これを充実するためには、知識データをこれら文型パターンにあてはめ、文型パターンを増やさなければならない。

例 吸う：N[hum]がN[Con]をV
{hum→人間, Con→具象物}

しかし、“脱脂綿が水を吸った。”という文はこの文型にはあてはまらない。我々は実際の抄録文から使われているサンプルを集め格構造の部分構造を抽出し、より強固な結合価文法の構築をめざしている。

⑤ 2文字漢字列の分析用基礎資料となる。

2文字漢字で構成されている語を調べると、延べ語数は非常に多いが、種類はそれほどでもない。このため2文字漢字で構成されている語の共起関係を調べることはKWICを使うにしても膨大になり大変である。そこで“語と語の関係”で得られた関係が文中に有るか否かを調べる。

これら文を抽出し知識データの有効性の判定資料として使うことができる。

分析の手順を示すと次のようになる。

- 文中に知識データによる共起関係があると予想できるもの。
 - 共起関係が認められる文である。
 - 共起関係が認められない文である。
- 文中の共起関係を知識データでは抽出できなかったもの。

このような利用方法が考えられる。

⑥ 既に考えてきた応用分野

そのほか、これまでに知識データを使つての応用分野としては次のようなものを考えている。

文字認識、音声認識の後処理に利用し精度向上をはかる。仮名漢字変換システムの同音異義語の選択、機械翻訳の多義語の選択等に役立つ。これらへの応用のためにも着実な知識データが準備されなければならない。

5. おわりに

知識データの収集、整備は日本語処理システムが一層発展するために必要である。約100万件程度のデータが集まれば、新しい研究段階になるであろう。

この知識データについては解決しなければならない多くの問題があるが少し先が見通せる状態になってきた。

知識データを利用することは高品質で、より高度の利用技術が提供されることになるが、これを実現するためには着実なデータ分析の作業が必要である。

この研究の一部は文部省科研費課題番号(60302090)(代表者 吉田 将)によって行われた。

参考文献

- (1) 田中康仁, 水谷静夫, 吉田 将 語と語の関係について 情報処理学会 自然言語処理 41-4 1984-1
- (2) 田中康仁, 吉田 将 自然言語の分析による知識データの収集 情報処理学会「自然言語処理技術」シンポジウム 1984-11
- (8) 田中康仁, 水谷静夫, 吉田 将 語と語の関係による

知識データの収集 — 自動カテゴリー化について日本科学技術情報センター第21回情報科学技術研究会論文集 1985. 3

(4) 田中康仁, 吉田 将 自然言語における知識データについて情報処理学会第31回(昭和60年度後期)全国大会 3H-9

(5) 村木 , 慣用句・機能動詞結合・自由な語結合 日本語学第4巻第1号 1985. 1 明治書院

(6) 宮地 裕 “慣用句の意味と用法” 明治書院 1982. 10

(7) 倉持, 阪田 “必携慣用句辞典” 三省堂 1985. 1

(8) 水谷静夫, 石綿敏雄 他, 「文法と意味 I」 朝倉日本語新講座3 朝倉書店

(9) 大野 普 「類語新辞典」 角川書店

(10) 田中康仁, 吉田 将 慣用表現について 情報処理学会第32回(昭和61年前期)全国大会 1S-3 1986. 3

(11) 田中康仁, 水谷静夫, 吉田 将 語と語の関係による知識データの収集について 日本科学技術情報センター 第22回情報科学技術研究会論文集 1986. 3

(12) 小西友七編 英語基本動詞辞典 研究社出版 昭和55年9月

(13) 本間 茂, 山階正樹, 小橋史彦 連語解析を用いたべた書きかな漢字変換日本語文書処理 21-2 情報処理学会研究報告 85-JDP-21 1985.5

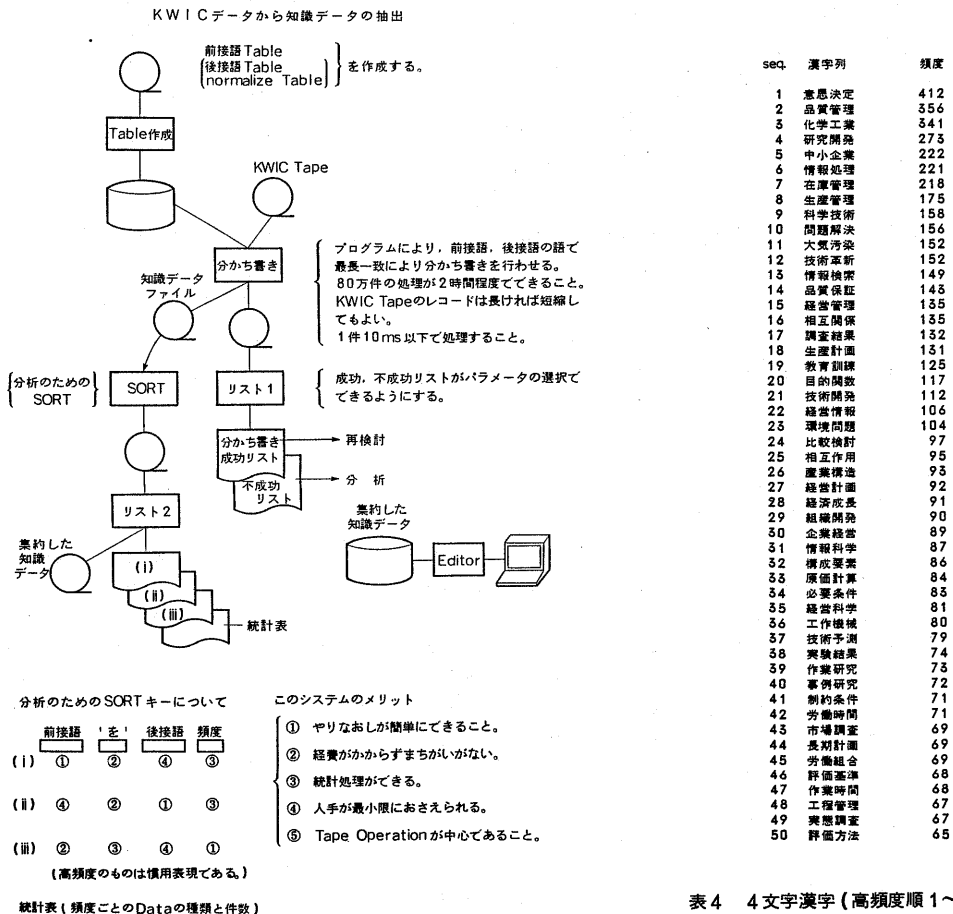


表4 4文字漢字(高頻度順1~50位)

図2 語と語の関係抽出プロセス

seq.	前接漢字列	種類	頻度
1	問題	290	1152
2	方法	282	1139
3	管理	266	2270
4	計画	262	1479
5	関係	229	776
6	作業	220	597
7	情報	212	873
8	技術	202	822
9	調査	193	684
10	時間	190	739
11	方式	189	380
12	条件	173	779
13	処理	167	661
14	活動	166	644
15	分析	164	658
16	状況	159	461
17	機能	157	501
18	構造	149	293
19	業務	148	389
20	過程	139	371
21	部門	138	515
22	能力	134	309
23	企業	133	691
24	会社	133	465
25	内容	133	387
26	調査	131	404
27	工場	130	392
28	状態	127	253
29	研究	125	67
30	組織	124	328
31	計算	117	365
32	設備	117	315
33	可能	117	305
34	教育	116	496
35	政策	114	358
36	効果	111	307
37	基準	108	392
38	分野	108	338
39	費用	106	237
40	機関	102	348
41	産業	102	259
42	機械	101	301
43	評価	99	329
44	対策	98	310
45	工業	95	712
46	期間	95	180
47	決定	92	680
48	要素	92	292
49	生産	92	228
50	特性	92	156

seq.	前接漢字列	種類	頻度
51	環境	89	341
52	確信	88	328
53	報告	88	306
54	開発	86	819
55	結果	86	532
56	制度	85	290
57	設計	85	240
58	制御	85	221
59	資料	82	254
60	利用	82	220
61	製品	81	195
62	構成	79	154
63	関数	77	431
64	段階	77	235
65	速度	77	150
66	項目	77	137
67	目標	76	225
68	市場	76	223
69	検査	76	197
70	状態	75	125
71	商品	73	149
72	巡回	72	163
73	分布	71	319
74	機構	71	135
75	測定	70	173
76	形態	69	143
77	経済	67	298
78	体制	67	186
79	手段	67	159
80	形式	67	98
81	変化	66	186
82	手順	65	140
83	訓練	64	256
84	工程	64	242
85	機器	64	220
86	政策	64	205
87	行動	64	174
88	要因	63	132
89	部品	62	150
90	目的	61	143
91	予測	60	258
92	実験	60	110
93	効率	59	217
94	水準	59	213
95	部分	58	78
96	原価	57	150
97	領域	57	118
98	手段	57	93
99	単位	56	84
100	社会	55	240

seq.	後接漢字列	種類	頻度
1	生産	189	1231
2	企業	140	845
3	作業	156	848
4	技術	147	939
5	経営	142	1062
6	情報	140	1236
7	経済	135	698
8	管理	127	560
9	各種	122	186
10	労働	114	618
11	教育	114	524
12	自動	106	366
13	基本	105	396
14	販売	105	357
15	社会	101	455
16	研究	100	731
17	製造	99	315
18	計画	99	269
19	設計	94	287
20	開発	93	264
21	製品	92	279
22	処理	88	364
23	産業	87	290
24	一般	85	190
25	組織	84	397
26	主要	83	169
27	環境	82	579
28	安全	82	331
29	品質	81	762
30	機械	80	389
31	市場	79	282
32	製造	76	400
33	工場	73	177
34	利用	70	310
35	輸送	69	201
36	会計	68	342
37	投資	68	215
38	関連	68	152
39	使用	67	207
40	工業	66	249
41	特許	65	281
42	計算	65	255
43	標準	65	199
44	國際	61	271
45	制御	61	241
46	基礎	59	184
47	評価	58	205
48	専門	58	204
49	時間	58	200
50	設備	58	141

表5 前接漢字列種類頻度の高いもの順(1~100位)

表6 後接漢字列種類頻度の高いもの順(1~50位)

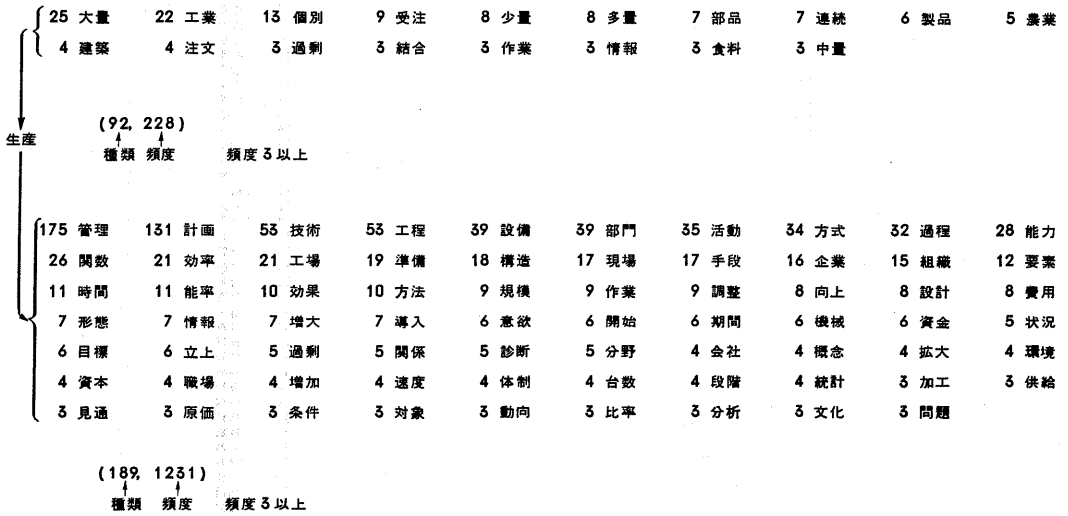


表7 語と語の接続関係(生産)