

コーパスの医学英語教育への貢献： PERC Corpus プロジェクトを中心に How Corpus Linguistics Helps EMP Education, with Particular Reference to the PERC Project

演者

投野由起夫

明海大学外国語学部助教授

それでは最初に、コーパス言語学の基礎的な概念をお話しします。ここでお話しするのは、まずコーパスの定義 (definition) です。それから、コーパスにはいろいろな情報を言語的に入れ込むということを行います。1つは mark-up、もう1つは annotation といいますが、そこにはどのような情報が入るのかということをご紹介します。それから、英語のさまざまなコーパスがすでにかなりできあがっているのですけれども、どんなものがあるかをザッと紹介します。最後に、ESPとか、私はここで初めて知ったのですが、EMP というような言い方をするそうですね、その EMP Corpora はどうかというようなことを見たいと思います。

コーパスの定義と言語情報

まず definition ということですが、コーパスというのはいろいろ条件があるのです。1つ目は *New Oxford Dictionary of English* (NODE) で、いまは2版が出ていますが、その98年の definition です。“a collection of written or spoken material in machine-readable form, assembled for the purpose of studying English structures, frequencies, etc.” というので、collection であると。それから、書き言葉でもあるし、話し言葉でもあるよということです。また machine readable form というので、コンピュータが読め

るようになっているということです。そして「特定の目的のために集められている」ということです。

もう1つ、別の definition がありまして、これは私がランカスターの同じ仲間たちと一緒に書いた text なのですから、これにもいくつか似たようなことが繰り返されています。“Corpus is a collection of machine-readable authentic texts, which is a sample to be representative of a particular language or languages variety” ということです。ここでもいくつかポイントになるような要素が挙げられています。

コーパスというのは単純な寄せ集めではないのです。例えば、Google か何かで検索したような結果とは違います。コーパスの場合は、中身が何かわかっているということです。そして、収集には目的があるということです。特定の言語の編集なりバラエティ、あるいは医学英語なら医学英語というような、特定の目的をまず頭に入れて集められているということです。

続いて、コーパスにはいろいろな特徴があります。まず、言語的な特徴、corpus mark-up というのを見てみましょう。コーパスにはさまざまな text 的な情報を入力するので、すけれども、普通はコーパスはただの plain text なのです。そこにさまざまな情報を入れる入れ方があります。大きく分けると、mark-up と annotation とに分かれています。“mark-up” というのは、text に関する情報を入れるのです。

演者紹介：投野由起夫氏 (とうの・ゆきお)

東京学芸大学卒業。同大学大学院英語科教育専攻(英語教育学)、英国ランカスター大学大学院言語学科博士課程修了、Ph.D. (コーパス言語学)。東京都立航空工業高等専門学校英語科専任講師、東京学芸大学教育学部英語科教育講座専任講師等を経て、2001年より現職。

大量の言語データをコンピュータで分析して言語研究を行うコーパス言語学を専門とし、コーパス言語学に基づく第2言語習得(特に語彙習得)を応用したNHK教育テレビ「100語でスタート! 英会話」の講師としても活躍。



Features: corpus mark-up

Mark-up: a system of standard codes inserted into a document stored in electronic form to provide information about the text itself.

```
<header>
<textnum>0057</textnum>
<filedesc>
  <title>my family</title>
  <name>Hanako Yamada</name>
  <grade>10</grado>
  <date>1999-07-10</date>
</filedesc>
<textdesc>
  <medium>essay</medium>
  <domain>informative</imaginative>
  <genre>student writing</genre>
  <region>Japanese EFL</region>
</textdesc>
```

スライド1

これはだいたいヘッダセクションに入れるのですが、例えばそのtextがだれによって書かれたかとか、何年に出版されたかとか、どこのウェブから取ってきたか、そのようなtextに関する情報です。information about textということです。

もう1つは、textのstructureがあります。例えば、ここからここまでがparagraphなのだとか、ここまでがsectionだとか、ここがAbstractだとか、論文だったらここがConclusionだとか、そのようなparagraph,あるいは論文なりtextの構造がありますが、そういうものを記号で入れるというようなことをします。

もう1つは“corpus annotation”です。annotationという用語は、言語的な情報を入れるときに使います。textの言語的特徴という、例えば品詞のタグです。これは動詞だとか、名詞だとか、そのような情報です。

2つ目はラマタイゼーション(lemmatization)といいまして、これは見出し語です。例えばbring, brought, bringingとか、そのように活用するときに、全部これはbringなのだとか、そのような見出し語情報を入れたりします。

“semantic tagging”というのは、意味領域を入れたりします。それから“parsing”,これは構文解析データです。

例えば、スライド1はcorpus mark-upの例です。これは私がつくっているlearner corpusの例ですけれども、text自体は出ていないけれども、textのdescriptionとかfile descriptionの部分をヘッダに格納しているという、このようなフォーマットです。ここを検索して、コンピュータが、例えば皆さんが目的とするようなtextをコンピュータが探すときの手がかりなわけです。例えば、医学英語でこういう領域のtextだけ取ってこいとか指示できるようになっているわけです。

同じように、その論文の構造みたいなことに関する情報をmark-upしておけば、abstractだけ取って来るとか、あるいはreviewのところの表現だけ取って来るとか、そんなこともできるかもしれません。

Corpus annotation: POS tagging

POS information (British National Corpus)

```
<s n="26">
<w PNP>it <w VBZ>is <w AV0>now <w DT0>more <w CJS>than
<w AV0>ever <w VVB>clear <w CJT>that <w AT0>every <w NN1>section
<w PRF>of <w NN1>society <w VVZ>needs <w TO0>to <w VBD>be
<w AJ0>involved <w PRP>in <w VVG>responding <w PRP>to <w NN1>AIDS
<g PUN>, <w PRP>including <w AT0>the <w NN2>churches<g PUN>.
<s n="27">
<w NPD>ACET <w VBZ>is <w AT0>a <w AJ0>Christian <w NN1>initiative
<w VVN>supported <w PRP>by <w DT0>all <w NN2>denominations<g PUN>.
<s n="28">
<w PNP>it <w VBZ>is <w AT0>the <w NN2>churches <w CJT>that
<w VVB>provides <w DPS>our <w NN2>volunteers<g PUN>;
<w PRP>without <w DPS>their <w NN1>support <w PNP>we <w VM0>would
<w XX0>not <w VBD>be <w AJ0>able <w TO0>to <w VVI>provide <w AT0>a
<w NN1>service <w AV0>at all<g PUN>.
```

スライド2

Corpus annotation: lemmatization

POS information (CLAWS vertical format)

0000001	002	-----			PUNC
0000002	001	NULL	<start>		PUNC
0000003	001	NULL	<e>		PUNC
0000003	010	FFY	You		PUNC
0000003	020	VM	must		must
0000003	030	VVI	leave		leave
0000003	040	RT	now		now
0000003	041	,			PUNC
0000003	050	RR	otherwise		otherwise
0000003	051	,			PUNC
0000003	060	FFY	you		you
0000003	070	VM	will		will
0000003	080	VBI	be		be
0000003	090	JJ	late		late
0000003	100	IF	for		for
0000003	110	APFG	your		your
0000003	120	JJ	social		social
0000003	130	NN2	studies		study
0000003	140	NN1	class		class
0000003	141	.			PUNC

スライド3

今度は言語的な annotation というのはどんなものがあるかという、品詞の tagging というのがあります。品詞タグ付与ですね。これは、いま、世界で最も大きいバランスのとれた、「均衡 corpus」といいますけれども、その1億語の corpus がイギリス英語であるのですが、その British National Corpus (BNC) というものの内容です。

スライド2を見ていただくと、1個1個の単語の前の部分に angle blanket といいますが、<w PNP>とか書いてありますね。こういうものがタグなわけです。この品詞のタグを手がかりに、一括で、例えば動詞の形を全部取って来るとか、その連鎖を取って来たりということもできるのです。

続けて、同じようなことに lemma の情報をつけることができます。例えば、スライド3では1行に1単語ベースのフォーマットになっています。You/must/leave/now/otherwise/you/will/be/late/for/your/social/studies/class.とか。そうすると、1単語ごとに品詞がついて、これが lemma になっているわけです。studies というところが study になったりしています。そうすると、見出し語検索をすれば、複数形でも単数形でも一括で取って来たり、そのようなことが全部できるわけです。

Corpus annotation: semantic tagging

■ Semantic tags (SEMTAG vertical format)

```
0000009 082 -----
0000009 090 FFIS1 I 28mf
0000009 091 VEM 'a 25 A3+
0000009 100 VVCK going T1.1.3[14.2.1
0000009 110 TO to T1.1.3[14.2.2 25
0000009 120 VVI call Q2.1 Q2.2 X3 A10+
0000009 130 PFY you 28mf
0000009 140 II about 25
0000009 150 APFGE my 28
0000010 010 EN1 best 32mf B18 Y23 39+
0000010 020 NN1 family 34c A4.1c
0000010 021 ,
0000010 030 AT the 25
0000010 040 KN2 Saitos 299
0000010 041 .
```

スライド4

Features of corpus-based approaches

- Various language statistics (frequencies and distributions)
- Machine-readable ⇒ Quick search for words/phrases
- Genre balance ⇒ Genre analysis
- Computer analysis ⇒ Find typical patterns not idiosyncrasies
- Huge amount of data ⇒ Finding patterns Statistical analysis
- POS/Syntactic info ⇒ more possibilities for noiseless data
- Cluster information ⇒ collocation pos n-grams, etc.

スライド6

Corpus annotation: Parsing

■ Suzanne Corpus (POS, lemma, parsing)

```
A01:0250j - PFHf It it [S[N]a.N]e]
A01:0250k - VVDt urged urgo [Vd.Vd]
A01:0250m - CST that that [Fn%:o.
A01:0250n - AT the the [Ns:s.
A01:0250p - NN1c city city .Ns:s]
A01:0250q - YL <ldquo>-
A01:0250r - VVDv +take take [V.V]
A01:0260a - NN12 steps stop [Np:o.Np:o]
A01:0260b - TO to to [Ti:c[V].
A01:0260c - VV0t remedy remedy .V]
A01:0260d - YR <rdquo>
A01:0260e - DD1f this this [Ns:o.
A01:0260f - NN1c problemproblem.Ns:o[Ti:c]Fn%:o]S]
A01:0260g - YF +. .O]
```

スライド5

もう少し進んだ研究をしている人たちもいます。彼らは semantic field というタグをつけているのです(スライド4)。例えば、ランカスターには USAS (UCREL Semantic Analysis System) というタグセットがあって、これは Longman の Lexicon という thesaurus の dictionary があるのですけれども、その類義語辞書のような意味領域を単語に振っているのです。このようなことをすることによって、例えば、一般の人が読んでいる health と、専門領域の health ということを仕分けるようなことが技術的にできるようになる可能性があります。

また、context word のタグの意味領域を見ることで、専門的なことを言っているのか、一般的なことを言っているのかを text で分別するような技術が、ここから出てくる可能性があります。

スライド5は parsing のデータです。“It” 以下がずっと元の文なのですが、これが lemma です。その隣の “it” 以下が品詞です。ここは全部構文解析した頭語情報が入っているのです。これは Suzanne Corpus というコーパスの一部なので、このようにセンテンスの NODE から、NP, VP のかたまりなどということが入って来たりす

ると、構造的な情報を検索したりすることもできます。

このような情報が入ってくると、いろいろなことの検索が可能になってくるわけです。皆さんはもうわかると思いますが、こういうデータを入れたものから、コーパス言語学というのはさまざまな頻度、そして分布の情報を取ってきます(スライド6)。

まず、コンピュータに入っていますから、非常に高速に単語やフレーズのサーチができます。それだけではなく、ジャンルのバランスなどをとってあれば、こういう領域ではこんな使い方、みたいなこともできます。

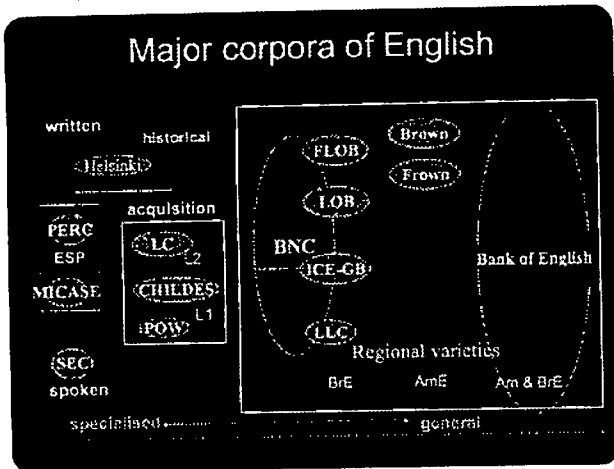
次に、コンピュータ処理しますから、人間が気づかないようなパターンを取って来たりすることができます。

それから、ものすごく量が多くなってくると、いろいろな統計処理などをして、さまざまなデータの summary の仕方ができます。

それから、単純に単語レベルで検索していたり、あるいは文字列処理しているだけだと正確に取ってこれられないのですが、それに品詞や構造、構文解析のデータがついていれば、もっとノイズの少ない、正確なサーチができるようになります。

また、1単語だけではなくて、クラスタの情報も取ってこられます。ですから、どういう単語とどういう単語が連結しているかとか、どういう品詞とどういう品詞がよくくっついてくるかとか、品詞と単語の連鎖なども取ってこることができます。

以上がコーパスのおおよその特徴なのですが、英語のコーパスというのは世界中にどのぐらいあるかということ、いまはものすごくたくさんあるのです。前は自分でリストをつくっていたのですが、いまはあまりにもたくさんいろいろなコーパスがあるので、自分でリストをつくっているのが面倒くさくなってしまったぐらい、たくさんあります。



スライド7

コーパスの分類

大きく分けると、書き言葉のコーパス (written corpus), それから spoken corpus とに分かれます (スライド7)。それから、どちらかという general な汎用のコーパスと, specialized corpus といいますけれども, 特殊領域のコーパスに分かれるのです。汎用のほうは, 大規模なものが多いです。BNCは1億語, Bank of Englishがいまは4億数千万語とされていますが, つい最近, OxfordがOxford English Corpusというのをつくってしまして, 来年ぐらいに完成する予定なのです。そのコーパスのサイズは10億語です。ものすごい規模のコーパスの構築がいま盛んになってきています。

特殊コーパスのほうは, 集めにくかったり, 技術的に spoken だと書き起こしに手間がかかったり, 小さいものが多いです。ESP関係は, このあたりの領域に位置づけられるわけです。EAPなども含めれば, いくつか公開されているものがありますけれども。

医学英語というのは, 実は私は門外漢で, 今回も発表の中にたくさん難しい単語が出てきて, 私の知らない単語がたくさんあったので, 先生方に「これは何という意味?」と聞かれると恥づかしいな, と思ったりしたのですが, Medical English Corporaというのは, 私の知るかぎりではあまりたくさん公開はされていないのです。

1つは, 歴史的な研究をしている人たちが公開しているものがあります。The Corpus of Early English Medical Writingというのがあります。これはCD-ROMでJohn Benjaminから出ているのですけれども, ご存じの方もいるかもしれませんが, 昔どのように英語を使っていたかという, 歴史的なものです。

あとは, コーパス言語学の学会に行っても, ESP関係で Medical English という発表は多いですが, ほとんどが自分でインターネット等で取ってきているものばかりです。ですから, コピーライトをクリアしているような試みはあまりないのです。私的な institute の中でキープしていたりして, だれも使えないというデータがたくさんあります。

Text Selection

Journal Citation Reports (JCR)

<http://www.isinet.com/isi/products/citation/jcr>

- It presents quantifiable statistical data that provides a systematic, objective way to determine the relative importance of journals within their subject categories.
- 5,700 journals in the Science Edition
- **Impact factor:** provides a way to evaluate or compare a journal's relative importance to others in the same field

Selection was made by choosing the top 20% of the journals with the highest impact factor in each field

JCR classification of subject fields is used

スライド8

PERC と CPE

そこで, われわれはそのような壁を1つ打破しようということで, いま, 大きなプロジェクトを立ち上げています。それがPERCという団体と, CPEというプロジェクトなのです。これをザッとご紹介したいと思います。

PERC (Professional English Research Consortium) は2002年に, 割と広く Professional English というものに関するさまざまな研究をするような association として発足しました。コーパスだけではなくて, 例えば教材作成やテスト作成, 広く Professional English の研究というものをする団体です。 (<http://www.perc21.org>)

一応, 日本でできたのですけれども, いまは会津大学の Thomas Orr 先生がヘッドになっています。世界的には, 例えばミシガンの John Swales とか, Ulla Connor とか, そのような ESP 関係の大御所も PERC のメンバーとして活動してくださる予定です。

いま, PERC の活動の中心になっているのが, この CPE (Corpus of Professional English) ということで, 科学技術系の英語に関してコーパスづくりをしようということになって, 私がコーパス言語学のほうで, Professional English の専門家ではないのですけれども, コーパスづくりという意味で参画しています。

そして, さまざまな形のリサーチを, このデータを使ってやっつけようではないかということで始まっているのです。ただ, なかなか時間がかかって大変な作業なので, まだ完成には程遠いのですが, 第1段階はほぼ終わってきたのでご紹介しています。CPEは, もともとはさまざまなコーパスの text のタイプをバランスよく集めようということを考えています。例えば, 医学なら医学なのだけれども, journal だけではなくて, この分野のさまざまなことをサンプルとして1つのコーパス・データとして持っているようなことをイメージしています。ただ, 最初は academic journals を中心に収集しています (スライド8)。

text セレクションは, Journal Citation Reports というもの

JCR Fields

- | | |
|--------------------------|---------------------------------------|
| ☒ Medicine | ☒ Electrical & Electronic Engineering |
| ☒ Biology | ☒ Computer Science |
| ☒ Food Science | ☒ Telecommunications |
| ☒ Environmental Sciences | ☒ Nuclear Science |
| ☒ Mathematics | ☒ Material Science |
| ☒ Physics | ☒ Metallurgical Engineering |
| ☒ Chemistry | ☒ Construction & Building Technology |
| ☒ Engineering | ☒ Civil Engineering |
| ☒ Earth Science | ☒ Fisheries |
| ☒ Agriculture | ☒ General Sciences |
| ☒ Oceanography | |
| ☒ Forestry | |

スライド9

を元にして、ここのインパクトファクターが非常に高い journal を選定して、サイエンスのエディションだと 5,700 誌ぐらいがこのレポートに載っているのですけれども、そのうち上位 20% の journal を選んでいます。そしてその出版元に Consortium として手紙を書いて、〇〇年から〇〇年までの全部の article (あるいは論文) を研究用に提供してほしい、というようなことをしています。

JCR の fields は、大体 22 fields (スライド 9) になります。medicine は、そのうちの 1 つなのです。科学技術関係の割と総合的なコーバスというようなデザインになっています。

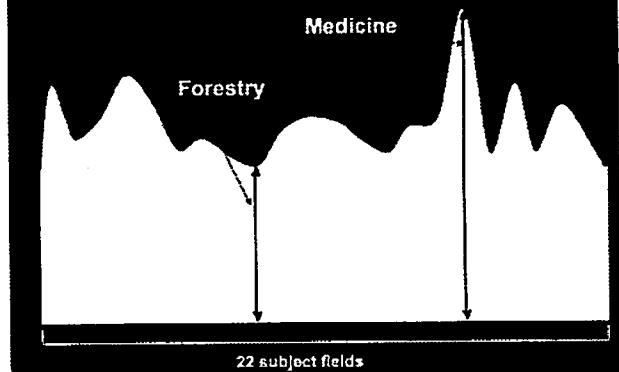
バランスに関しては、実は 22 fields をやってみると、JCR でも随分違うのです。例えば、medicine は非常に JCR に出てくる率が高いとか、あるいは forestry は低いとか、そのようなことがあります (スライド 10)。

そこで、コーバスのデザイン的には、まず academic journal text を、例えば 30 万語ぐらい必ずどの領域にも集めておく。そのあと、70 万語ぐらいは他の text type を集めて、少なくとも 100 万語規模のものを 22 fields 集めると。そうすると、2,200 万語になりますから。このようなもので、100 万語単位の比較が科学技術分野に関してできるようなものを、まずベースに持っています。そして、それに集められるだけ上乘せしておいて、使える方は全部 medicine を使いたいとか、そのようなことがあってもよいだろうというようなデザインになっています。

さまざまな出版社にコンタクトをとって、そして非常に大きな反響がありました。例えば、いま大御所だったら Elsevier とか、IEEE とか、そのようなところからもコピーライトの許諾を得ています。ただ、彼らはテキストファイルで提供してはくれないのです。勝手にウェブページにアクセスして、取っていいよということは言うのですけれども。

われわれがずっとこの 2 年ぐらいやっているのは、PDF ファイルを e-journal のサイトから取って、そしてそれを全部テキスト処理でダウンロードしたものをコンバートして、

Balance of each field



スライド10

CPE-Medicine

- | | |
|--|---|
| ☒ ALLERGY | ☒ MEDICINE, OPHTHALMOLOGY |
| ☒ MEDICINE, DENTISTRY, ORAL SURGERY & MEDICINE | ☒ MEDICINE, OTORHINOLARYNGOLOGY |
| ☒ MEDICINE, EMERGENCY MEDICINE | ☒ MEDICINE, PARASITOLOGY |
| ☒ MEDICINE, ENDOCRINOLOGY & METABOLISM | ☒ MEDICINE, PERIPHERAL VASCULAR DISEASE |
| ☒ MEDICINE, GASTROENTEROLOGY & HEPATOLOGY | ☒ MEDICINE, PSYCHIATRY |
| ☒ MEDICINE, GENERAL & INTERNAL | ☒ MEDICINE, RADIOLOGY, NUCLEAR MEDICINE & MEDICAL IMAGING |
| ☒ MEDICINE, GERIATRICS & GERONTOLOGY | ☒ MEDICINE, RESEARCH & EXPERIMENTAL |
| ☒ MEDICINE, HEALTH CARE SCIENCES & SERVICES | ☒ MEDICINE, SUBSTANCE ABUSE |
| ☒ MEDICINE, IMMUNOLOGY | ☒ MEDICINE, TOXICOLOGY |
| ☒ MEDICINE, MEDICAL INFORMATICS | ☒ MEDICINE, VETERINARY SCIENCES |
| ☒ MEDICINE, NEUROSCIENCES | ☒ NEUROSCIENCES |
| ☒ MEDICINE, OBSTETRICS & GYNECOLOGY | ☒ OTORHINOLARYNGOLOGY |
| ☒ MEDICINE, ONCOLOGY | ☒ PHARMACOLOGY & PHARMACY |
| | ☒ VETERINARY SCIENCES |

Tokens (running words) in text 3,155,412

スライド11

そしてクリーニングするという作業をしています。この部分の作業がものすごく大変で、これで何本も論文を書いているぐらいのものなのですけれども、PDF からテキストへの自動変換というのは非常に難しいのです。簡単だ、ソフトもあるよ、と言う人がよくいるのですけれども、そんな生易しいものではなくて、さまざまなゴミが出るのです。そのようなものをクリーニングしながら、皆さんがある程度使えるようなフォーマットにまでするのが一苦勞という作業を、ずっとしてきました。

今回、CPE の first release というのがそろそろ出せるような状態になってきました。これは academic journal だけです。ですから、先ほどのバランスはとれていません。ただ、すばらしいことに、すべてコピーライトは許諾を得ています。そして、さまざまな出版社から、いまのところ 3,700 編の学術論文を取れています。いまの規模ですと、1,700 万語ぐらい、これは科学技術英語全体ですけれども、それが 22 fields に分かれています (スライド 11)。

medicine はどんな内容になっているかという、GSER のブレイクダウンなのですけれども、スライド 9 のような分野のものが取れていて、journal はまたこの下に何種類か

Term	Count	Percentage	Other	Frequency
PATIENTS	8715	0.22	0.07	20,000,000
CELLS	8391	0.20	12.92	0.00
URINE	2,029	0.05	0.06	0.00
DISEASE	2,272	0.05	1.27	0.00
CLINICAL	7,203	0.18	0.1	4,979.37
TREATMENT	4,211	0.11	0.11	4,100.11
CELL	5,844	0.15	10.24	0.00
INFECTION	1,791	0.04	0.11	3,070.97
MORE	2,211	0.05	1.29	3,794.46
WOMEN	1,974	0.05	0.4	9,703.40
DRUGS	1,323	0.03	0.0	3,427.16
DOSE	1,836	0.04	1.44	3,493.71
BLOOD	2,144	0.05	1.57	3,301.07
TRAIN	1,299	0.03	0.0	3,311.03
RISK	2,032	0.05	1.72	3,289.20
CANCER	1,257	0.03	0.0	3,234.74
COCAINE	831	0.02	0.0	3,104.30
PATIENT	1,830	0.04	0.23	3,091.00
VACCINE	873	0.02	0.0	2,962.73
WEEK	20,029	0.05	0.4	2,693.44
THESE	1,023	0.02	0.0	2,725.64
ORAL	1,838	0.04	0.1	2,674.30
AMINO	1,197	0.03	0.0	2,674.30
HUMAN	2,543	0.06	0.12	2,674.30
SMALL	1,710	0.04	0.0	2,614.46
RECEPTOR	1,072	0.02	0.01	2,472.00
INFECTED	1,213	0.03	0.0	2,465.00
A	1,520	0.04	1.06	2,436.14
LDL	969	0.02	0.0	2,294.43
WHO	1,882	0.04	0.11	2,278.81
WAYS	1,072	0.02	0.0	2,200.01

スライド 12

入っているのだと思います。

現在、われわれがつかっている全体のコーパスの中で、medicineのセクションは300万語ぐらいです。ですから、それを大きいと思うか、小さいと思うかは皆さん次第だと思いますけれども、少なくともこのぐらいの規模のものがいま集まっています。

CPEのavailabilityということですが、このfirst releaseというのは、まず最初はPERC researchersという人たちが使って初期の研究をします。public releaseというのは、web-basedなインターフェースで、将来、来年の初めぐらいにオープンにしたいと思って、いま作業中です。小学館コーパスネットワーク (<http://www.corpora.jp>) という、小学館が企業としてこの科学技術関係の部分のサポートしてくれているので、そのようなサイトをオープンする予定です。

もう皆さんも知っているかもしれませんが、BNCとか、COBUILDのデータとか、大きなコーパス・データが、この小学館のサイトではweb-basedで利用できるようになっています。ですから、web-basedな、コーパスのポータルみたいな感じに使っていこうということです。

コーパスとEMP

では、残りの20分ぐらい、コーパスとEMPという話をしたいと思います。先ほどのようなデータができてきて、われわれがそれにアクセスできるようになると、どのような分析が可能で、そしてそれを元にすると、例えばどんな可能性があるのだろうかというようなことを、ちょっと考えてみたいと思います。

実は私の分野というのは、コーパス言語学といってもバリエーションにコーパス処理の言語学的な部分ではないのです。私はどちらかというと、言語教育とのリンクを考えるようなことが専門なので、医学英語教育という意味では、私のやっているようなことを医学英語に適用するようなことを

blood: grammatical relations

Triplets:

- <subject, flow, blood>:
1 "blood" is the subject of "flow"
- <object, suck, blood>:
1 "blood" is the object of "suck"
- <modifies, transfusion, blood>:
1 "blood" modifies "transfusion"
- <a_modifier, blood, cold>:
1 "cold" is an adjective modifier of "blood"

スライド 13

やっごらんになったらよいか、と思うことがいくつかあります。そういう研究分野では、実はコーパスというのは、皆さんもおわかりだと思いますが、直接的に利用する場合と、間接的に利用する場合があるのです。そういうことについて、ちょっとお話ししたいと思います。

それから後半のほうでは、available corporaということ、すでにつくられているコーパスを使うというようなことと、あとはDIY (do it yourself) corporaですね、自分でコーパスをつくるような技術が、最近発達してきているのです。これはたぶん、インターネット等に興味のある方ならば、必ず役に立つ技術だと思いますから、ぜひメモでもしてってください。

まず、間接的にはどんなことが可能かということの例をいくつかお見せします。まずはEMPのmaterials developmentということ、たぶん語彙リストみたいなものをつくっている方がいますね。コーパスをうまく使えば、非常に有効に語彙リストがつかれます。それから、そこから発展してlexicon database, 単なるword listではなくてlexiconとしてそういうもののデータをつくっていくということも可能になります。それから、reference toolsだったら辞書や文法書、それから教科書、テキストですね、それから言語テスト。先ほど、言語テストの話をしてきたみたいですが、

一般の英語教育のほうでされているようなことをちょっと見て、それを元に、医学英語だったらどうかというようなことを考えてみたいと思います。

コーパスを使った語彙リストにはいろいろな技術があるので、いまは非常に目的特化したtextを集めます。そして、そのtextでどんな単語が使われているかを自動抽出するような技術があります。それはreference corpusと違って、大量の一般的な英語の文章と比較することによって、その専門分野で特によく使われている特徴語を抽出するような、そういう手法があるのです。それぞれのコーパスから語彙リストをつくって、医学英語のコーパスから

Behavior of "blood" (1)

		BNC	CPE-Med
blood + V	run	91	0
	flow	68	0
	come	51	0
	take	38	taken (1)
	pour	34	0

• In PERC-Med Corpus, the word "blood" is seldom followed directly by a verb.

スライド 14

頻度表をつくる。そして、そのreferenceになる大型コーパスから、無色の、割とニュートラルなリストと比べてみるわけです。そういうことによって、medical Englishに特徴的な語彙を取り出すということです。このようなことを、いまはコーパスではソフトウェアで簡単に行うことが、技術的には可能になっているのです。

これには、いろいろな統計の手法を使って抽出する、例えば代表的なものはlog likelihoodとか、そういう相対頻度のずれの度合いの大きいものを抽出するのを使います。

そうしますと、例えば先ほどのmedicineのセクションをCPE(科学技術英語)全体と比べてみたりすると、medicineの300万語のセクションのキーワードは、スライド12のような感じで出てくるのです。medicineのtext properに出てくる特徴語が、特に特徴の強いものから順番に並んだりします。なかには専門用語、専門用語との中間、機能語的なものももう少し下へ行くと出てきますけれども、医学英語の特徴をこういう形で機械的に抽出してみても、そしてわれわれが気がつかなかったようなポイントをコンピュータ的に取り出すようなことが可能になります。

これだけで、単にword listを分野ごとに細かくつくる、というような発想になりますね。それだと単語単体のことになってしまうので、どのようにほかの単語と一緒に使うかというようなことがわかりません。

そこで、もう少し細かく文法的な関係、これをgrammatical relationsというのですけれども、そのようなものをしっかり見たり、それからcollocation分析みたいなことをすると、もっと単語の使い方が詳しくわかるようになります。

例えば、"blood"という単語を例にとると、一般の英英辞典であるMacmillanの辞書では、こんなことが書いてありました。1つめは普通の「血液」です。2つめは「血統」とか、町、地域、グループとしての、仲間意識としての「血」みたいなことでしょうか。そして3つめがviolenceやdeathを意味するという一般の英語で使われる"blood"の定義があります。

Behavior of "blood" (2)

☞ blood + noun: most common usage

- | | |
|------------------|---------------------|
| - blood flow | - blood transfusion |
| - blood pressure | - blood aqueous |
| - blood samples | - blood gas |
| - blood vessels | - blood ocular |
| - blood cells | - blood retinal |
| - blood donation | - blood loss |
| - blood glucose | - blood levels |

スライド 15

そして、一般英語のコーパスを見てみますと、"blood"というのは割と"blood"単体で出てきて、そして動詞と結合するようなものがたくさん類文で出てきます。ですから、一般の英語で"blood"と使う使い方というのは、動詞とのcollocationなどがとても多いのです。

ところが、これをもっと詳細に分析してみるわけです。例えば"blood"というのがどのような関係をしているかというのを、tripletという形で取り出すわけです(スライド13)。「bloodはflowのsubject(主語)になる位置関係にある」というのを<subject / flow / blood>というtripletとして取り出します。このようなことを、機械的にやるわけです。そうすると、"blood"という単語の使われる場面や状況を、grammatical relationsの束で表そうとするわけです。こういうことをたくさんやると、"blood"という単語がどのように使われるかということをもっとtextを分析した結果としてデータベース的に取り出すことができるのです。こういうものを取り出すと、かなりメリハリがわかってきます。動詞との連結の例を見るだけでもかなり違うのです。

BNCではこういう動詞との結合がたくさん出てくるのですが、おもしろいことに、medicineのコーパスではこういう例はほとんど出てこないのです(スライド14)。つまり、"blood"単体で結びつく動詞との連鎖は、medicalなtextではほとんど出てこないのです。逆にどういものが出てくるかというと、"blood"というのはほとんどがcompoundで出てくるのです。

スライド15は名詞との連鎖で、大量に先ほどのmedicineのテキストから抽出したもののなのですけれども、この"blood flow"のflowは動詞ではなく、名詞のflowなのです。このようなmedicalなテキストでは、"blood + noun"という形での結合が圧倒的に多いのです。単純に動詞とbloodが結びついている例というのは、ほとんどありませんでした。

これを見ると、われわれが知っている"blood"などという単語でも、医学では全く違う使い方をするということがわかります。こういう部分というのは、案外われわれにとっ

Sketch Engine

- Adam Kilgarriff (Univ. of Brighton, Lexicom)
- His Sketch Engine makes it possible to integrate grammatical information extracted from corpora into a web-based lexicon database.
- Easy to use, very fast
- <http://sketchengine.co.uk>

スライド 16

て、というか普通の英語を勉強してきた者にとっては難しい点というか、あまり気がつかない点ではないかと思います。

実は、これを大量にデータベースとして皆が使えるような研究をしている人がいるのです。Brighton大学のAdam Kilgarriffという人ですけれども、この人はLongmanの辞書、Macmillanの辞書、そしていまはOxfordなどを手伝っている自然言語処理(natural language processing)の研究者で、lexical databaseなどに非常に詳しい人です。この人はSketch Engine(<http://sketchengine.co.uk>)という先ほどのgrammatical relationsみたいなものをウェブ上で検索できるようなシステムをつくっているのです(スライド16)。

これは非常にパワフルなもので、例えば“medicine”という単語を検索すると、これはBNCなので、ESPのテキストではないのですが、それでも、例えば“medicine”というのは、“object of”ですから、こういう動詞の目的語になって現れる。“study medicine”, “practice medicine”, “prescribe medicine”とか、そのような目的語の関係。こちらはこういう形容詞がくっつくか。それから、どんなcompoundをつくるかなどということが、データベースの形でガーッと出るので。

そうすると、この“medicine”という単語の使われ方がチャートになっていて、そしてこういうところのstudyの42番というところをクリックすると、具体的に例文に飛ぶわけです(スライド17)。そうすると、どのような使われ方をするか一度にわかるわけです。これをESPのコーパスでやったら、とてもおもしろいのではないのでしょうか。そうすれば、一般のgeneral Englishではなくて、ESPのtext, EMPのtextだったらどうなるかというようなことがわかれば、その差がかなりあると思うのです。そういう差を研究するようなことも、必要だと思います。

いま、コーパスを使った辞書というのは国内外でたくさん出ているのです。ですから、たぶん先生方も先ほどのようなことを元にして研究すれば、もっと使いやすいEMP用の辞書というのをしてくれるはずなのです。いまのEMP用の

スライド 17

辞書は、私はよくわかりませんが、専門用語の辞書は単にterminologyの辞書ですよね。そういうものが圧倒的に多いと思うのです。ですから、本当の使い方がよくわからないものが多いと思うのです。

例えば“infection”という単語は、一般の英英辞典ではan ear infectionという用例が1つ、それからあとはmild/slight/severeなどの程度を表す形容詞と一緒に使われているとか、この程度の情報量です。これを医学生が聞いても、とてもではないけれども、彼らは“infection”の医学における全体像はわかりません。

ところが、先ほどのようなデータベースを駆使して研究した結果、何かしらそういうものを元にしたテキストをつくれば、例えばですけども、“infection”の2つの意味で、特にこちらのほうの意味では、どのような名詞とcompoundをつくるかとか、adjectiveだったらどんな種類のもの結びつくかなどということを整理事たりできるかもしれません(スライド18)。これは、ちょっと不正確なものもあるかもしれませんが。ですから、これは1つの例だと思って見てください。こちらでは、動詞との結びつきです。そういうものも、“infection”にくっつく動詞などがあって、これがもしも訳語と簡単な例文つきなどで出れば、相当活用できるresourceになるのではないのでしょうか。

コーパスを用いた辞書だけではなく、いまは文法書もどんどん出始めてきているのです。先生方はどうお考えになるかわかりませんが、医学英語には医学英語なりの文法があるのでしょうか。そういうことは、ESPなどでもいろいろな議論があると思います。もしも一般の文法とは少し違う文法の特徴があるとしたら、そのようなものをコーパスで調べて、コーパスに基づくような医学英語の文法書があってもよいかもしれません。maybe ENP grammarということですね。

私は、コーパスを使ったテキストをたくさん出しています。NHK関係もそうですし、ほかにもいろいろなところから出ています。こういうものは、先ほどのiPodとか、iTunesとか、ああいうところにも私の教材が載っているの

EMP dictionary entries

in·fec·tion [n.]

1 [countable] a disease that affects a particular part of your body and is caused by bacteria or a virus:

[noun + infection] *HIV; pylorus; HIV-1; virus; tract; chest; throat; ear*
[adj. + infection] (1) *viral; bacterial; urinary; respiratory; fungal; candidal; trichomonal; chronic; genital*
(2) *secondary; heavy; opportunistic; acute; recurrent*

2 [uncountable] when someone is infected by a disease:

[verb + infection] *prevent; acquire; diagnose; eradicate; transmit; cause; treat; combat; associate; develop; catch; eliminate; spread; suffer; avoid*

スライド 18

ですけれども、そういうコーパスに基づくテキストという発想があれば、ひょっとしたらコーパスに基づくEMPのテキストもあってよいかもしれません。私のような中高生モード、一般モードでなくても、もっと医学生向きの何かがつくれるかもしれません。

そして、Language Test Developmentということですが、実はいまは世界的な規模でいくと、CambridgeのUCLES, USAではETS, この辺が大量にコーパス・データを集めています(スライド19)。5年ぐらい前には、そういうことはしていませんでしたが、最近では彼らは非常にそういうことについての重要性を感じ始めていて、テストを受けた人たちのexam responsesとかエッセーなどを大量にアーカイブしているのです。そして、それを元に、いろいろな形のことをしています。例えば、自分たちのつくったテストの開発、そしてvalidationのためにそのようなものを参考にしたリ、それからエッセーを自動gradingするということですね。これはETSなどがいま盛んにレポートを発表していますけれども、こんなこともやっています。

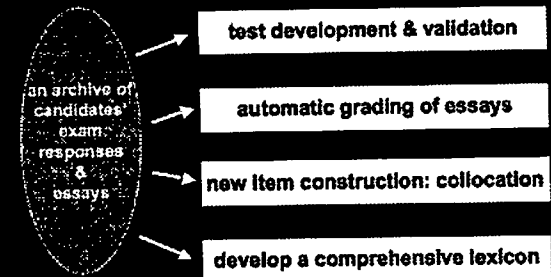
それから、新しいitem constructionということのを collocation でやるようなことを、UCLESでやっています。つまり、collocationの知識で測ったほうが、単語単体の知識よりも弁別性があるというようなことを、彼らはレポートで言っているのです。

それから、comprehensive lexiconということ、このようなものを元にして、どのような単語の知識のデータベースができれば、それを元にテストが開発できるか、そのようなものを総合的につくろうというような発想も、テスト分野の人たちの間には出てきています。

いまの繰り返しになってしまいますが、general Englishとmedical関係のEnglishとの間の違いも、当然あるでしょうね。そして、それを習っている医学生が、あるいはお医者さんになる人たちが習っている過程で、どんなつまずきや間違いをするか。きょうも学習者コーパス(learner corpus)の発表がいくつかありましたけれども、そのような発想で、こういう3つのデータを組み合わせて研究するよう

Language test development

■ UCLES (Cambridge) & ETS (USA)



スライド 19

な可能性が出てきます。つまり、目標となるESPのターゲットの様子をしっかりと研究するだけではなくて、そこまでに行くまでにどのようなプロセスを経ていて、どこがつかずいているかというようなことを、learner corpusで調べるわけです。これは、spokenだったり、writtenだったりしますけれども。こんなことを組み合わせていくと、EMPのさまざまな materials developmentに非常にプラスになるような素材をたくさん集めることができます。

Web as Corpus

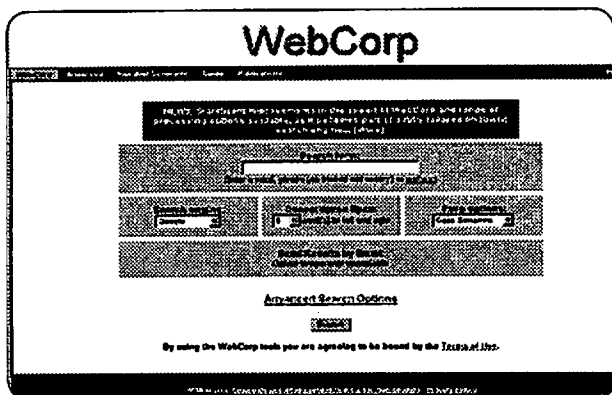
すでにあるコーパスを使うというのは、先ほどのPERCのコーパスなどを使えばよいのですが、そうではなくて、自分でつくりたいという方がいると思うのです。実は、こういう技術が最近ものすごく発達してきているので、ちょっとそのことをご紹介します。

実は、web as corpora,あるいはweb as corpusというような動きが、コーパス言語学者の中にあるのです。2006年度のAdam Kilgarriffによれば、Googleでインデックスされているduplicate-free(データベースの形になっていない、つまりstaticな形でインターネットに載っている)なtextは、英語の場合は10 thousand billion wordsです。どのぐらいか、日本語だと何兆と言ってよいかわからなくなってしまうぐらい多い量ですけれども、このぐらいがestimateであるので、本当はどうかかわからないぐらいたくさんの量があります。

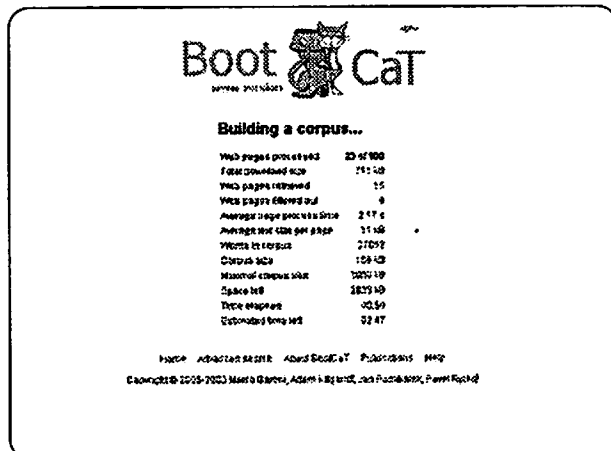
これを何とかして利用しなければ、というようなムーブメントがだんだん強くなってきて、実はGoogleに勤めているFranz Ochさんという人が、250 billion wordのトレーニング・コーパスというのをつくったと発表しました。すごい量ですね。ですので、こういうことが実際に研究としては可能になってきています。

そこで、いくつか既存の技術で皆さんが検索できるようなcorpusの使い方があるかというようなことを、ちょっと説明します。

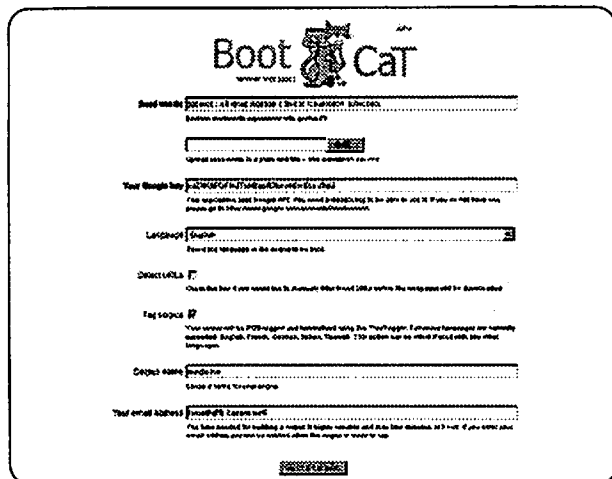
1つはWebCorp(<http://webcorp.org.uk>)という、これは



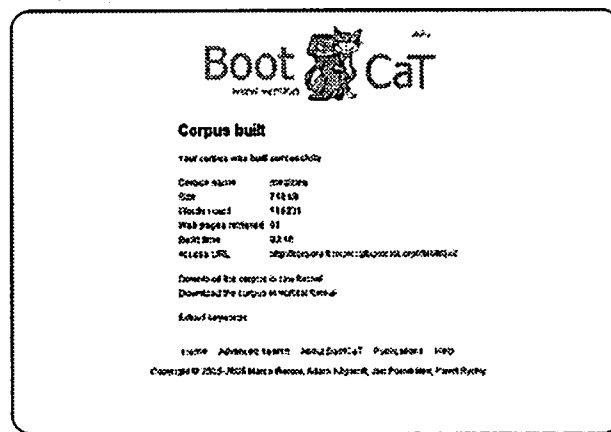
スライド 20



スライド 22



スライド 21



スライド 23

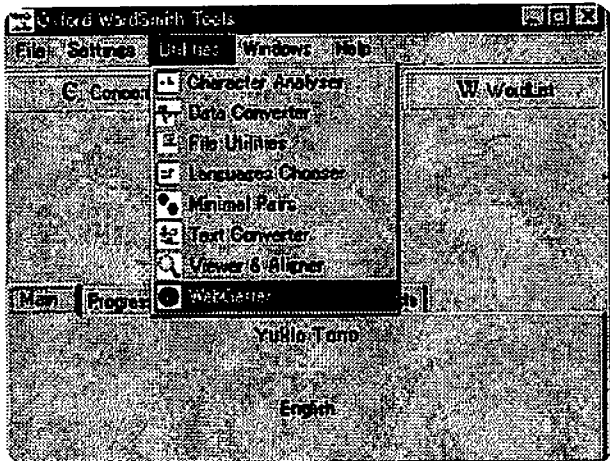
Antoinette Renouf という人が中心でやっています (スライド 20)。Liverpool 大学から最近 Central England 大学へ移りましたが、これはウェブ上の Google などの resource を元にキーワードを検索して、それを表示するようなサービスです。それも、ウェブページをコーパス的に使おうというような発想のインターフェースなのです。

WebCorp を立ち上げて search term というところに自分でキーワードを入れるわけです。これを組み合わせたりすることもできますし、自分なりにパターンを書いたりすれば、かなり医学的な text だけを絞って取ってくるようなことができます。サーチエンジンを選んだり、いろいろなことをやって、submit を押すのです。ただ、ウェブをクロールして集めてくるので、ちょっと時間がかかります。ですから、WebCorp の 1 つの問題は、遅いということなのです。ただし、チョンとやっておいて、少しどこかへ行ってから戻ってきたりすると、このようになっているわけです。中には「何だこりゃ」と思うようなものも出てくるかもしれませんが。そういう場合には、もともとの URL へ行って、そこが信頼できるかどうかということをチェックしてくればよいわけです。しかし、少なくともこういう形でたくさんのフレーズを一度に集めたりするようなことは、イ

ンターネット上でもできます。こういうインターフェイスのほうが、Google よりもごみがなくて、比較的よいかもしれません。

次に、まだ開発中なのですが、すごいものが出てきました。BootCaT というシステムです (スライド 21)。これは Bootstrapping Corpora and Terms というのですが、これは先ほどのような分野ごとの特徴語みたいなもの、「このキーワードを 10 個ぐらい入れれば、絶対に関連する text がばっちり取れる」と思うようなキーワードのリストをつくりますね。そういうものを seedwords といいます。その seedwords を元に、tuples という、いくつかペアをつくるのです。例えば、3 つずつぐらいのペアを、10 個のうち最初の 3 つ、1, 2, 3。次には 2, 3, 4。次には 3, 4, 5 というように、組み合わせた単語のリストをどんどん Google にほうり込みます。そうすると、その組み合わせで Google が取ってきたものの中の集合体みたいなものができますね。それをコーパスとして考えようという発想なのです。そして、Google のヒットを、これは API でヒットしたヒット結果を元に、自動でダウンロードするようなことをするわけです。

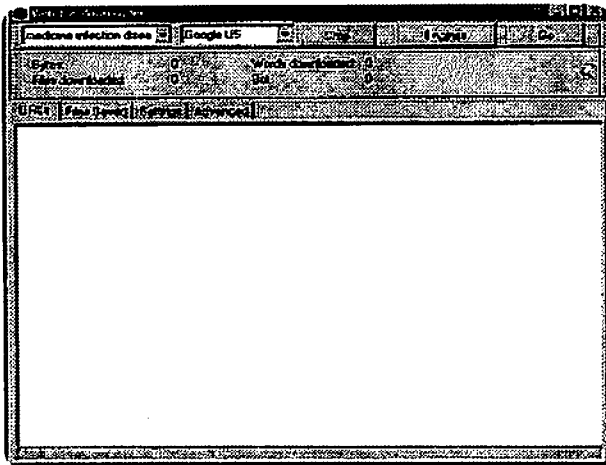
この BootCaT は現在開発途中なのですが、彼らは何とすべてをウェブインターフェースでやろうと考えています。



スライド 24



スライド 26



スライド 25

ですから、われわれ利用したい人がその seedword を入れれば、クロールしてきて、コンピュータでサイトから集めてきたウェブページを保存して、クリーニングして、そしてすべてコーパスとしてそれをフォーマットして見せてあげるところまで、全部やるのです。すごい技術でしょう。このようなことを、いま WACCI (web as corpus cool initiative) という団体が中心でやっているのです。

どのようになっているのか、お見せしますと、私が medicine のセクションから取ってきた "patient" "cell" "drug" "disease" "clinical" "treatment" "infection" などというキーワードトップの単語を入れたとします。これには Google の API を使う、Google に登録したキーが必要ですが、英語のデータを取ってきてくれ、それを "medicine" という名前にするよ、私のメールアドレスはこれですと書いて、submit するわけです。

そうすると、コーパスをつくっていますということで、ウェブページをプロセスしたのがこれだけで、というようなページが出て、自動的にこの数字がどんどんふえていくわけです (スライド 22)。そして、corpus built (コーパスができました) というメッセージが表示されれば完了です (スライド 23)。実は先ほどありましたけれども、いまの a バ

ージョンはテスト版で出ているので 3 メガバイトまでしかつくるパワーがないのです。けれども、この範囲でこの URL にアクセスすると、何と、単に text をダウンロードできるだけではなくて、すぐに検索できるインターフェイスに飛んでいけるのです。つくってくれたものを、すぐにコーパスの検索ソフトとして使用できます。そうすれば、"infection" というようなものの動き方などが、全部コーパス上でさまざまに分析することができます。

それだけではなくて、collocation の処理等もできます。例えば「動詞 + cancer という形を取ってこい」というと、動詞のところが高ライトされて、cancer とどんな動詞がくつつくかとか、そのようなこともわかったりします。

このようなインターフェイスは、実は Sketch Engine とも連動しています。将来的には、何と、先ほど私が紹介した grammatical relations みたいなものも瞬間的につくって、そしてデータベースとして同時に使えるような技術をいま開発中なのです。大量のデータにこういうことができるようになるには、まだちょっと時間がかかると思いますが、少なくとも先生方が教室内で非常に特定の文脈の医学トピックのようなものを seedword にして、この text を世界中から取ってきて、ミニ・コーパスを使って授業をしたりということは、すでに可能な時代になってきているということです。

似たようなことは、実はいちばん新しい WordSmith のバージョン 4 という、Windows のソフトでもできるのです。WebGetter という機能を使うのですが、WebGetter というのを選びますと (スライド 24)、こんな画面になって (スライド 25)、そこに先ほどの "medicine" とか "infection" というキーワードを入れるわけです。そして「Google で取ってきて」というようにやると、どんどん世界中のこういうページがあるところをクロールしていきます。そうすると、先ほどの WordSmith ですぐに検索できます。ただしこれは、先ほどのように文法のタグ等はつきません。先ほどの Boot-Cat のほうは、文法のタグまでつけるので、それがすごい

ですけれども。こんな感じで手軽に、このWordSmithなどは1万数千円のソフトですけれども、こういうものを使うだけでも、すぐにかなりなことはできます。

こういうことを元にすれば、直接医学生がソフトを使って、そしてコーパスにアクセスして、そのコーパスでいろいろ調べるようなことが現実味を帯びてきます。インターネットさえあれば、自分でキーワードを入れて、瞬時にコーパスをつくったりすることもできるかもしれません。また、そういうものを元にしてエッセーを書いたり、いろいろなことをする指導をうまく組み合わせることもできます。こういうものをdata-driven learningというのですけれども、そのようなことを実際にやっている人たちもふえてきています。

ということで、皆さんはコーパスをどう思われたでしょうか。EMPに非常に役に立つ部分があると思うのです。ま

だあまりたくさん行われてはいないかもしれませんがけれども。手間暇をどのぐらいかけるかということでは、ある程度テクニックさえ習えば、1人でもできます。それほど大変なことではないです。大変だと思ったら、われわれのようなコーパス言語学者と組んでやればよいわけです。そして、international committee, 先ほどのBootCatみたいなものは、世界中の人たちが「おれもやりたい」といって参加して、オープンソースでやっているのです。あの辺の技術は、全部無料です。ですから、そのようなことを皆さんと一緒にやれば、resourceを共有したりできますね。Medical Englishというのは、たぶん世界中共通なところが何かあるのでしょうか。そのようなことをやってみるのもよいのではないのでしょうか。どうもありがとうございました。

(2006年7月16日、ウエルシティ金沢にて収録)