

医学分野のオンライン・データベースにおける
索引作業の比較評価

A Comparative Evaluation of
Indexing in Medical Online Databases

神 門 典 子, 上 田 修 一, 土 井 彩 子
Noriko Kando, Shuichi Ueda, Ayako Doi

Résumé

This study attempts to compare the indexing in medical online bibliographic databases, *EMBASE*, *JMEDICINE*, and *MEDLINE*. We conducted tests concerning “the extent to which indexing terms link related documents”, “the extent to which indexing terms discriminate among these sets within the database”, and “the extent to which terms discriminate finely among individual documents”, through the practical method suggested by White and Griffith in their paper “Quality of indexing in online data bases” (*Information Processing and Management*, 23(3): 211-224 (1987)). We also conducted qualitative comparison on indexing terms assigned by these databases, with reference to the characteristics of their vocabularies and size of each database. Based upon the comparison, some advantages and disadvantages of indexing and/or controlled vocabularies of three databases, several factors affect to the quality of indexing, and effectiveness of “devices” for indexing in the large operative databases are discussed.

I. 索引作業の評価

A. 索引作業の評価に関する研究

B. White と Griffith による索引作業の比較評価

II. 医学データベースにおける索引作業の比較評価

A. 調査目的と調査の枠組み

B. 対象データベースにおける主題表現の種類

神門典子: 慶應義塾大学大学院文学研究科博士課程図書館・情報学専攻, 東京都港区三田 2-15-45.

Noriko Kando: Graduate School of Library and Information Science, Keio University, 2-15-45, Mita, Minato-ku, Tokyo.

上田修一: 慶應義塾大学文学部図書館・情報学科, 東京都港区三田 2-15-45.

Shuichi Ueda: Professor, School of Library and Information Science, Keio University, 2-15-45, Mita, Minato-ku, Tokyo.

土井彩子: 日本航空株式会社東京空港支店第2旅客部.

Ayako Doi: Japan Air Lines, Passenger Traffic Haneda Airport.

1992年2月3日受付

C. 調査方法

III. 結果と考察

A. 索引語の定量的な比較評価

B. 索引語の質的な検討

C. 付与索引語からみた各データベースの索引作業

IV. 統制語彙との関係を踏まえた評価方法の検討

A. 定量的な評価指標の検討

B. 定量的な評価指標と統制語彙の特性

C. 統制語彙の特徴からみた索引語の質的な検討

V. 総括

I. 索引作業の評価

A. 索引作業の評価に関する研究

索引作業は、主題から文献を検索する場合、検索の成否に影響を与える主要な要因のひとつであり、それを評価することは図書館・情報学において重要な課題の一つである。

索引作業の評価に関しては、すでに多くレビューされている¹⁻⁴⁾。従来から、索引作業の評価は、主として、検索実験における検索効率、あるいは、特定のデータベースにおける索引作業の一貫性の調査^{5),6)}によってなされることが多かった。

索引作業はそれ自体が最終目標ではないので、索引作業自体を比較するよりも、情報検索というコンテキストのなかで検索効率によって評価を行なう方がより実際的である⁴⁾という観点もある。しかも、実際の検索では、索引語だけでなく、文献の標題や抄録中の語句などを用いたフリーテキスト検索を含め、さまざまなアクセスポイントを用いているので、索引語だけに限定せずに、実際に検索を行なった結果を比較する方が実情にそくしていると考ええる立場もある。

しかし、検索実験では、検索の成否や再現率および精度などで示される検索効率に、次のような要因が影響していると考えられる。すなわち、①サーチャが利用者のニーズを的確に把握しているか、②サーチャの経費や能力の差から生じる探索戦略の質、③索引作業の質や使用語彙の特性などのデータベースの質、④検索結果のスクリーニングが適切か、などさまざまな要因が考えられ⁴⁾、検索実験の結果からその要因を明らかにすることは困難である。また、再現率の算出法や検索文献の適合性の判定をどのように行なうか⁷⁾、検索効率の尺度が一

貫していないなど多くの問題点が指摘されている^{1),2)}。

また、比較的小規模の実験用システムにおける検索実験が多く行なわれている⁴⁾。一方、実際に稼働している大規模システムでは、検索集合が大きくなりすぎる出力過多 (output overload) が生じるなど、小規模システムとは異なるふるまいが想定されるため、大規模な実稼働データベースにおける評価の重要性が指摘されている⁷⁾。

さらに、統制語彙は、従来は再現率を高めるデバイスとして考えられていたが、大規模な実稼働データベースでは、精度を上げるためのデバイスとして機能するという報告もある⁸⁾。「出力過多」に対処するには、検索の精度を上げることがひとつの有効な方法である。したがって、フリーテキスト検索と統制語彙の索引語とを組み合わせで検索を行なっている現状においても、検索の精度を上げ、出力過多を防ぐ手だてとして、統制語彙は重要な役割を果たしているといえる。そこで、統制語彙による索引作業の結果だけを取り出して比較検討することは、実務面からみても意義があると考えられる。

索引作業の一貫性に関しては、一般的に、一貫性が高いと索引作業の質も高いと考えられている。しかし、索引作業の質は、「いかに内容を表現しているか」という点で判断すべきであり、質と検索効率と一貫性とは関連はあるが別のものである⁹⁾という指摘もある。また、文献の内容は、利用者の視点から、利用者にとって関心がある内容を表現するべきであるという指摘も多い^{4),10)}。

一方、同じ分野を対象とする複数のデータベースが存在し、それらを利用するには、各データベースの特徴や索引作業の特性を明らかにすることも必要である。

同じ分野を対象とする複数データベースの索引作業の比較に関して、Steven G. Watkins¹¹⁾ や Judith Lin-

genfelter ら¹²⁾は、検索の例や出力レコードの例を示して比較を行なっている。MaryEllen Siever¹³⁾は、オンライン検索に関する文献が特定の雑誌に集中しているというビブリオメトリックス研究の結果を踏まえ、それらの雑誌の掲載論文に付与された索引語を調べて集計することによって、オンライン検索に関する文献に対する索引作業を複数データベース間で比較している。

また、Howard D. White と Berver C. Griffith¹⁴⁾は、医学行動科学領域に関して *MEDLINE* を評価した一連のプロジェクトにおいて、同じ領域を対象とする複数データベースの索引作業の比較評価法を提示している。これは、主題に共通性がある文献群に対して、複数のデータベースが付与した索引語を比較する方法で、検索システムとは独立して索引作業だけを、大規模な実験システムにおいて評価することを目指した方法であり、種々のデータベースや主題にも一般化が可能である。

以上のごとく、索引作業を評価するには、さまざまな方法や観点がある。その中で、本稿では、大規模な実験システムにおいて、さまざまな主題領域に適用可能な方法で、検索システムとは独立して、統制語彙を用いた索引作業の結果を、同じ領域の複数データベース間で比較評価するという立場をとり、White と Griffith¹⁴⁾の方法を用いることとする。以下に White と Griffith による索引作業の比較評価の方法を概説する。

B. White と Griffith による索引作業の比較方法

1. 比較評価の観点と方法の概要

White と Griffith¹⁴⁾は、索引作業を、①関連のある文献を結びつける、②データベース全体の中から特定の文献群を区別する、③その文献群の中から特定の1文献を区別するという3つの次元でとらえ、それぞれに対応する「索引語の共通性 (spanning)」, 「索引語の識別性 (discrimination)」, 「1文献あたりの付与索引語数」の3つの観点を設定している。比較評価の基準として、共引用関係を用いて、主題に共通性がある文献からなる文献群を選定し、各文献群中の文献に付与した索引語を、これらの3つの観点から比較検討している。

「索引語の共通性」は、索引語が関連する文献を結びつける程度を意味し、文献群中の文献に共通の索引語が付与されている程度によって示す。具体的には、文献群中の全ての文献に付与されている索引語数と半数以上の文献に付与されている索引語数とを計算し、それらを

「索引語の共通性」の指標としている。

「索引語の識別性」は、索引語が、データベース全体から、主題に共通性がある特定の文献群を区別する程度を意味し、データベース全体ではあまり使用していない索引語ほど識別性が高くなる。「識別性指標 (discrimination index)」は以下の式によって算出している。

$$\text{識別性指標} = \frac{1}{\log_{10} A} \quad (1)$$

ただし、 A は当該データベースにおけるその索引語の付与回数 (postings)

文献群の半数以上の文献に付与されている索引語で識別性指標が 0.25 以上のものを「共用識別語 (s/d terms)」としている。これは、当該文献群中では多くの文献に付与しているが、データベース全体では使用回数が少なく、当該文献群を検索するのに有効な索引語と位置づけている。なお、この識別性指標は、Salton ら¹⁵⁾の「語の識別性値 (Term Discrimination Value)」とは異なる指標であり、大規模な実験システムで容易に算出できる。

「1文献あたりの付与索引語数」は、索引語が、主題の共通性によって結びつけられた文献群中から特定の1文献を識別する程度を示す。すなわち、各文献群に付与された異なり索引語数の1文献あたりの値が大きいほど、個々の文献の主題がさまざまな索引語によって表現されていることになり、当該文献群の中で特定の1文献を識別する程度が大きくなると位置づけている。

比較は、DIALOG で提供されている *MEDLINE* と他のデータベースとの間で行ない、*MEDLINE* と *EMBASE* は4文献群、*MEDLINE* と *BIOSIS* は4文献群、*MEDLINE* と *PsycINFO* は10文献群において比較評価をしている。

さらに、4文献群に関する質的な検討結果も報告している。すなわち、文献群ごとに、2文献以上に付与した索引語を列挙し、それらを完全一致 (exact matches)、関連 (related terms)、不一致 (unmatched terms) に分けて比較することによって、索引作業の特徴を検討し、統制語彙に追加することが望ましい候補語も提案している。

2. 方法に関する議論

この White と Griffith の方法に対して、文献群の選定法と識別性指標の算出法に関して以下のような議論がある。

a. 文献群の選定法に関する議論

文献群の選定法に関して、White と Griffith 自身は、比較すべき索引作業と独立した指標に基づいて主題の共通性を定めることが必要であるとし、文献群の規模は3～8 文献が適切であるとしている。共引用関係のほか、主題知識に基づいた専門家の判断、または、レビュー論文や同一の著者によって発表された一連の論文などによる文献群の選定も可能であるとしている¹⁴⁾。Lancaster⁹⁾は、引用関係による文献間の結びつきの妥当性を積極的に認めない人もいるので、標題中に同じ語句を含んでいるかなどの、共引用以外の、作業がもっと容易な方法を用いてもよいのではないかと述べている。

複数データベースの索引作業を比較している他の研究をみると、Clara M. Chu と Isola Ajiferuke¹⁶⁾ は授業の参考文献リストを用い、Siever¹³⁾ はビブリオメトリクス研究によって、特定の数誌に文献が集中していることが示されているトピックをとりあげ、その数誌の雑誌の全掲載論文をそのトピックに関する文献群として用いている。いずれも、索引作業とは独立した基準で主題に共通性がある文献群を選定しているが、適用できるトピックが限定され、あらゆる場合に一般化できる方法ではない。

それに対し、引用関係は、さまざまなトピックに適用可能な方法である。しかも、その主題の専門家である著者が、各自の関心に基づいて引用して利用したことに基づく関係であり、利用者からみた文献の内容の類似性が何等かの形で反映されていると考えられる。さらに、共引用関係に関しては、文献の類似性を示す尺度としての妥当性が示されている¹⁷⁾。以上により、引用関係に基づいて主題に共通性のある文献群を選定することは、作業量は多いが、さまざまな主題に対して適用可能であり、索引作業とは独立した基準に基づいて主題の共通性を認定し、文献の利用者である著者からみた関係を反映していることから、妥当な方法であると考えられる。

b. 識別性指標の算出法に関する議論

Isola Ajiferuke と Clara M. Chu¹⁸⁾ は、White と Griffith の識別性指標の算出法は、データベース規模の違いを無視していると批判し、代わりに次式を提案している。

$$\text{識別性指標'} = \frac{\text{当該データベースにおけるその索引語の付与回数}}{\text{当該データベースのレコード総数}} \quad (2)$$

しかし、前述のごとく、大規模な実稼働システムでは検索集合が大きくなりすぎる「出力過多」が大きな問題であり、データベース規模の違いは重要な問題である。したがって、(2) 式のようにデータベース規模で正規化するのは出力過多を考える場合にはふさわしくない。また、出力過多かどうかの判断は、検索結果が絶対数として何件であるかということよりも、数件の単位なのか、数十件か、あるいは数万件かというように、検索結果が1 件ずつ適合性を判断できる桁数の規模であるかどうかが必要になる。このことから、(1) 式のごとく、対数を用いた指標の方が実際に即していると考えられる。

II. 医学データベースにおける索引作業の比較評価

A. 調査目的と調査の枠組み

日本で利用できる主要な医学分野のデータベースである MEDLINE、EMBASE、JMEDICINE (JICST・医中誌国内医学文献ファイル) の索引作業の特徴を明らかにすることを目的とする。

医学は、研究者はもとより臨床医にとっても、文献情報の必要性が認識され、早期から情報検索への関心が高かった領域である。一方、わが国で出版された医学文献を広く収録している JMEDICINE は、2 種類のファイルをマージしたもの¹⁹⁾で、MEDLINE や EMBASE に比べて歴史が浅く、まだ評価が定まっていない。以上により、この3 つのデータベースを比較することは意義があると思われる。

調査は、White と Griffith¹⁴⁾ の方法を応用して、実際に稼働しているオンライン・データベースにおいて、付与されている索引語を調べ、それらを比較することによって、索引作業を比較評価することとする。比較の基本となる「主題に共通性がある」文献群は、比較する索引作業とは独立した指標によって求めることが必要である。ここでは、さまざまな領域に適用可能であり、利用者から見た文献間の関係を反映していると考えられる引用関係を用いることとする。

この対象文献群中の文献に、各データベースが、付与した索引語を、「索引語の共通性」「索引語の識別性」「付与索引語数」という3 つの観点から定量的に比較し、統制語彙との関連を踏まえながら質的に検討する。

また、第1 表のごとく、対象としたデータベースはそれぞれ、特徴のある主題検索のためのアクセスポイントとなる主題表現が提供されている。したがって、各デー

第1表 調査対象となる医学文献データベースにおける主題表現

	EMBASE	JMEDICINE	MEDLINE
<u>ディスクリプタ</u>			
統制語彙	MALIMET ¹	JICST 科学技術用語シソーラス	MeSH
デバイス 副標目	[リンク ('88~)] ¹	—	副標目
重みづけ ²	○	—	○
上位語の自動付加	—	○ ³	— ⁴
<u>タグ</u>	EMTAG (item index)	—	チェックタグ専用語 ⁵
<u>分類</u>	EMCLAS (section headings)	JICST 分類コード 医中誌分類コード ⁶	—
非統制語	アイデンティファイア 製造者名 製品名	自然語キーワード (準ディスクリプタ) (医中誌の索引語 ⁷)	—
<u>統制語彙</u>			
語数	MALIMET ¹ 約60万語 ('88) (優先語25万)	JICST 科学技術用語シソーラス 48,196語 ('87) (優先語 38,407)	MeSH (優先語 15,126)
階層構造	[部分的に ('88~)] ¹	○	○
見出し語の形	倒置なし, 単数形 小文字 省略形は最小限 “”, “.” は省略	省略形あり	倒置あり, 複数形 大文字 “”, “.” あり (検索時は省略可)
<u>索引作業</u>			
索引作成者	医学専門家	情報員・校閲員	索引作成者
索引作成者の用語彙	自然語 ⁷	統制語彙	統制語彙
<u>データベース規模</u>			
1985年出版の文献の収録数	274,380 件	206,537 件	307,265 件

- 1: 調査対象文献は主に 1985~86 年に各データベースに収録されたので, EMBASE に 1988 年以降に導入された新しいデバイスはで [] 示した。
- 2: 重みづけは重要な論点と周辺的な論点の 2 段階。「重要な論点」に対して付与した索引語は, 冊子体でアクセスポイントとなり, DIALOG では主ディスクリプタ (major descriptors) として扱われ, アスタリスク (*) が付けられている。なお, MeSH では, 冊子体の Index Medicus でもデータベースでも使用できる語を major descriptors, それよりも下位の概念を表し, データベースでは使用できるが, 冊子体では使用しない語を minor descriptors と称しているが, 本稿では, DIALOG の用語に合わせ, 各文献の主要な論点を表しているディスクリプタを「主ディスクリプタ」と呼ぶこととする。また, 1991 年以降の MeSH では, 両者の区別はなくなっている。
- 3: 自動付加された索引語を含めない狭い検索をする場合は, 索引語の前に @ をつける。
- 4: 文献の主要な論点に対して minor descriptors を付与した場合, 冊子体の Index Medicus では, 自動的にその上位の major descriptors に置き換えられる。なお, MeSH における minor descriptors については注 2 参照のこと。
- 5: チェックタグとしてしか使用されず, 冊子体ではアクセスポイントにならないもの。
“ANIMAL, CASE REPORT”, “COMPERATIVE STUDY”, “FEMALE”, “HUMAN”, “IN VITRO”, “MALE”, “SUPPORT, NON-U.S. GOV'T”, “SUPPORT, U.S. GOV'T, NON-P.H.S.”, “SUPPORT, U.S. GOV'T, P.H.S.” の 10 語。
- 6: 医中誌基本ファイルから抽出したレコード (主として会議抄録など) のみ。
なお, 医中誌の索引語は『医学用語シソーラス』(医学中央雑誌刊行会編) による。
- 7: 付与された自然語の索引語を, MALIMET と照合し, 自動的に優先語に変換する。MALIMET 中の語と照合しなかった語は, 編集者が検討し, 適宜, MALIMET に追加する。
出典: 文献 19) 20) 21) 23) および「JICST 科学技術シソーラス 1987」, 青木仕 (医学図書館, 36(3): 133-44 (1989))

データベースにおける主題表現にはどのようなものがあるかを検討し、共通して比較できる範囲を確定した。

B. 対象データベースにおける主題表現の種類

1. EMBASE の主題表現

MALIMET (Excerpta Medica's Master List of Medical Terms) は、優先語約 25 万語、その同義語などを含めると総計約 60 万語からなる膨大なディスクリプタの典拠リストである。調査対象とした 3 つのデータベースの統制語彙の中で最も優先語数が多く、特に約 10 万語以上の化学物質名を含んでいるという特徴がある²⁰⁾。索引作業は、医学の専門家が自然語で索引語を付与し、それを *MALIMET* と照合して自動的に優先語に変換している。照合しなかった語は、編集者が検討し、適宜 *MALIMET* に追加する²¹⁾。

ディスクリプタは、当該文献の重要な内容を表すクラス A と周辺的な内容を表すクラス B とに分ける簡単な重みづけがなされている。クラス A のディスクリプタは、冊子体の *Excerpta Medica* においてアクセスポイントとなり、DIALOG で提供されるオンライン・データベースではアスタリスク (*) が付けられ、主ディスクリプタ (Major descriptors) として扱われている²⁰⁾。

タグ・フィールドは、EMTAG、または、事項索引 (Item Index) という。これは、文献の種類、年齢、研究の種類、性別などの重要な側面、病因病理遺伝的側面、処置、実験法、機器、実験対象、身体部位、投薬経路、地域など、医学・薬学文献で広く使用する概念を表す 200 余語からなる。これは、コンピュータ検索において広い検索を容易にするためのもので、冊子体の *Excerpta Medica* では使用しない。EMTAG の一部は *MALIMET* から自動生成される。分類は、EMCLAS といい、階層的な構成をもつ冊子体の *Excerpta Medica* で抄録が掲載されている部分 (クラス) を表す。およそ 6,500 のクラスがあり、それぞれのクラス名を表すことばで検索することができる²⁰⁾。

そのほか、非統制語としては、被験者数などディスクリプタでは表現できないが文献の内容を知る手がかりとなるアイデンティファイアや、薬の製造者名・製品名がある。冊子体 *Excerpta Medica* の索引記入は、アクセスポイントとなるクラス A のディスクリプタの下に、各文献ごとに、付与された全てのディスクリプタとアイデンティファイアが列挙され、簡単な抄録のように当該文献の内容を示す働きがある²¹⁾。

なお、1988 年 1 月から、*MALIMET* は、*MeSH* の副標目に相当する「リンク」と使用頻度の高いディスクリプタからなる *MiniMALIMET* に階層構造をもつ「EMTREE コード」を導入しているが²²⁾、今回の調査は、1988 年以前にデータベースに収録された文献を対象としているので、この EMTREE コードやリンクは対象に含まれない。

2. JMEDICINE の主題表現

日本科学技術情報センター (以降 JICST とする) は、1981 年から「JICST 国内医学文献ファイル」を提供し、1983 年分以降は、それに医中誌基本データベースから JICST ファイルに収録されていないレコードを抽出して加え、*JMEDICINE* として提供している¹⁹⁾。したがって、*JMEDICINE* は、JICST の索引基準に従って作成したレコードと「医中誌基本データベース」から抽出したレコードとからなっており、主題表現もそれぞれ別の体系に従っている。

JICST の索引基準に従って作成されたレコードには、「JICST 科学技術文献ファイル」から医学分野のレコードを抽出したものと、医学中央雑誌刊行会に索引作業を委託して作成したものがあるが、いずれも、「JICST 科学技術用語ソーラス」を用い、同一の索引基準に従って索引作業が行なわれ、上位語の自動付加がなされている¹⁹⁾。「JICST 科学技術用語ソーラス」は、医学のみでなく、科学技術領域全般を対象としたソーラスである。

一方、「医中誌基本ファイル」から抽出したレコードには、「JICST 科学技術用語ソーラス」を用いた索引作業はなされず、ディスクリプタ・フィールドはない。「医学用語ソーラス」(医学中央雑誌刊行会発行) によって付与された索引語はすべて自然語キーワードとして扱われ、JOIS では、「医学用語ソーラス」の階層構造を利用した検索はできず、上位語の自動付加もない。分類は医中誌分類コードである。なお、この医中誌基本ファイルから抽出したレコードは、当該年の *JMEDICINE* レコード数のおよそ 6 割にあたり、その 80% は会議録である¹⁹⁾。

3. MEDLINE の主題表現

MEDLINE の索引作業は、米国国立医学図書館の訓練を受けた索引作成者が行なっている。*MeSH* (Medical Subject Headings) は、医学領域の文献を索引することを目的とするソーラスであり、15 のカテゴリーに分けられた階層構造を持つ。*MeSH* では、通常のディスクリプタを主標目 (main headings または main

descriptors) といい、これに適宜、副標目 (subheadings) を組み合わせて使用することができる。副標目は、qualifiers ともいい、主目標が当該文献中で果たしている役割を明確にしたり、主目標の意味をある側面に限定する働きを持っている。索引に使用するものは 75 語あり、それぞれ、疾患カテゴリー中の主標目に付加するもの、化学物質カテゴリー中の主標目に付加するものなど、組み合わせられる主標目のカテゴリーがあらかじめ決まっている²³⁾。

ディスクリプタは、当該文献の主要な論点を表すものと周辺の論点を表すものとに分ける、簡単な重みづけがなされている。主要な論点を表すものは、冊子体の *Index Medicus* においてアクセスポイントとなり、オンライン・データベースではアスタリスク (*) が付けられ、主ディスクリプタとして扱われている。

また、妊娠、ヒトの年齢、研究対象、性別、研究の種類、医学史の年代、研究助成など、医学文献に多く出現する概念を表す索引語は、チェックタグ (Check Tags) としてあらかじめまとめられ、索引時に当てはまる語を選ぶようになっている。このなかには、チェックタグとしてしか使用されず、冊子体の *Index Medicus* ではアクセスポイントにはならないチェックタグ専用のものがある (表 1 欄外の註参照)²³⁾。DIALOG では、タグ・フィールドに、このチェックタグ専用語はいっている。

4. 比較する索引語の範囲

以上のごとく、各データベースの主題表現を検討した結果、3 つのデータベースをなるべく共通の条件で比較するために、ディスクリプタとタグの両者を比較対象の「索引語」と捉えることとした。*JMEDICINE* にはタグがないが、*EMBASE* や *MEDLINE* でタグとなっている“HUMAN (ヒト)”, “MALE (男性)”などの語が、*JMEDICINE* のディスクリプタに含まれていることからタグも含めることとした。また、*JMEDICINE* で自動付加された上位語は除外する。*MEDLINE* では、定量的な比較においては、副標目は考慮せず、主標目のみを対象とすることとする。

C. 調査方法

1. 文献群の選定

はじめに、主題に共通性のある文献からなる、基準となる文献群を選定した。*JMEDICINE* の収録対象は日本で出版された文献であるのに対し、*EMBASE* と *MEDLINE* は日本で出版された文献の収録が少ないの

第 2 表 調査対象文献群の主題と文献群を構成する文献数

文献群	主 題	文献数
A	成長ホルモン放出因子	6
B	NK 細胞とモノクローナル抗体	6
C	ウイルス蛋白質とウイルス抗原	5
D	神経繊維	3
E	グリノサミノグリカン	3
F	硝子体の膠質細胞	5
G	硝子体のフルオロフォトメトリー	3
H	突発性門脈圧こう進症	3

で、まず、3 つのデータベースが共通して収録している文献を調べ、その中で主題に共通性のある文献群を選定することとした。

具体的には、「日本科学技術関係逐次刊行物総覧 1988」²⁴⁾ を用いて、対象とした 3 つのデータベースと *SciSearch* が共通して収録している医学分野の雑誌を調べた。これらの雑誌の全掲載論文が収録されるわけではない。そこで、*EMBASE* (DIALOG)・*JMEDICINE* (JOIS)・*MEDLINE* (DIALOG) を用いて、1985 年にこれらの雑誌に掲載された文献を検索し、3 つのデータベースが共通して収録しているおよそ 1200 件の文献を得た。なお、作業上の制約から調査対象を 1 年間に限定し、*JMEDICINE* が安定して収録を開始し、引用関係の調査も可能な年代として 1985 年を選んだ。

さらに、この 1200 件の中で主題に共通性がある文献からなる文献群を選定するために書誌結合関係を調べた。*SciSearch* (DIALOG) を用いて、これらの文献が引用している文献 25,921 件を検索し、1,293 組の書誌結合関係を認定した。そこから書誌結合が強い 92 の文献群を認定し、その中で調査に適した規模とされている 3 件以上 8 件以下の文献から成る 8 つの文献群を調査対象とした。これらは、第 2 表に示すごとく、文献群 A から H とした。

なお、当初、White と Griffith と同様に、3 つのデータベースが共通して収録している文献中の共引用関係を調べたが、共引用関係は非常に少なく、文献群を選定することができなかった。調査の基準として、主題内容に共通性がある文献から成る文献群は不可欠であり、索引作業とは独立した、さまざまな領域に一般化できる方法によって文献の主題内容の類似性を示すことが望ましい。そこで、本調査では共引用関係と同様に、引用関係、すなわち利用者である研究者からみた文献間の何等かの

医学分野のオンライン・データベースにおける索引作業の比較評価

関係を累積した尺度である書誌結合関係を用いることとした。共引用関係は年代と共に変化するのに対し、書誌結合は対象文献が出版された時点の固定した関係を示すことが大きな違いであるが、本稿では同一年に出版された文献に限定していることから、年代による変化を反映しないという書誌結合の問題点の影響は少ないと考える。

2. 分析項目と手順

文献群 A~H 中の各文献を、EMBASE, JMEDICINE, MEDLINE で検索し、各データベースが付与し

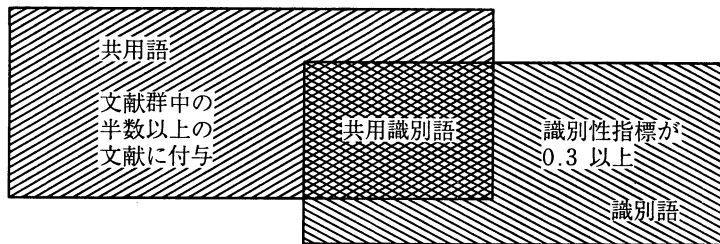
た索引語を出力した。その結果を White と Griffith¹⁴⁾の方法に従って、以下の点に関して分析した。

a. 文献群における索引語の共通性 (spanning)

データベースごとに、各文献群において、当該文献群中の全ての文献に付与している索引語を「完全共用語」、半数以上に付与している索引語を「共用語」とし、それぞれの数を計数した。

b. 索引語の識別性 (discrimination) と共用識別語

White と Griffith の調査では、データベースの対象年代がそれぞれ 5 年、6 年、7 年、20 年近くと大きく



索引語	付与文献数	共用語	識別性指標	共用識別語
<u>EMBASE</u>				
growth hormone releasing factor	6	○	0.434	○
growth hormone	5	○	0.344	○
drug efficacy	4	○	0.223	×
hypophysis	3	○	0.313	○
protirelin	3	○	0.346	○
gonadorelin	2	×	0.338	×
<u>JMEDICINE</u>				
ヒト	6	○	0.369	○
成長ホルモン	5	○	0.434	○
GRH	5	○	0.562	○
静脈内投与	4	○	0.294	×
ホルモン調節	3	○	0.699	○
血しょう中濃度	2	×	0.333	×
男性	2	×	0.333	×
TRH	2	×	0.440	×
GnRH	2	×	0.463	×
<u>MEDLINE</u>				
Peptide Fragments	6	○	0.309	○
Somatotropin-Releasing Hormone	6	○	0.431	○
Somatotropin	4	○	0.328	○
Adult	3	○	0.270	×
Insulin-Like Growth Factor I	2	×	0.572	×
Thyrotropin-Releasing Hormone	2	×	0.362	×
		↑ 文献群の半数 (3) 以上の文献に付与	↑ 共用語で 識別性指標 0.3 以上	

第 1 図 共用語・共用識別語の関係 文献群 A (文献数 6) における例

異なっていた。オンライン・データベースは、作成年代によって適宜分割して提供されることがあり、その対象年代の範囲は提供機関や時期によって異なる流動的なものである。しかも、作成年あるいは出版年によって検索を限定することは、比較的容易である。そこで、識別性指標は、White と Griffith の式 (1 式) に文献の出版年による限定を加えて、年次別に算出することとした。今回は 1985 年に出版された文献を対象としているので、次式のごとく、各データベース中の 1985 年に出版された文献への当該索引語の付与回数の対数の逆数とした。

$$\text{識別性指標}_{85} = \frac{1}{\log_{10} A_{85}} \quad (3)$$

ただし、 A_{85} は、その索引語が当該データベース中の 1985 年に出版された文献に付与された回数

また、White と Griffith¹⁴⁾ は、5 年から約 20 年の間の 40 万件から 198 万件を収録するデータベースを対象として、識別性指標が 0.25 以上、すなわち、索引語の付与回数が 1 万回以下の語を識別語とした。今回は 1 年間に限定して識別性指標を算出することから対象年数の比率を考慮して、0.3 という閾値を設定した。文献群毎に、識別性指標が 0.3 以上の共用語である「共用識別語」の数を計数した。共用語と共用識別語の関係を、第 1 図に示した。

c. 付与索引語数

文献群ごとに、各データベースが付与した索引語の異なる総数、および、その 1 文献あたりの平均を計数した。

以上の 3 つの指標を用いて、3 つのデータベースの索引作業を定量的に比較した。さらに、より詳細に検討するために、文献群毎に、各データベースが付与した索引語を質的に比較検討した。以下にその結果を示す。

III. 結果と考察

A. 索引語の定量的な比較評価

第 3 表に、文献群ごとに、文献数と、データベース別の完全共用語数、共用語数、共用識別語数、付与索引語数を示した。ただし、EMBASE において、ディスクリプタとタグに同一の語が付与されていた場合は、ディスクリプタのみで計数した。

1. 索引語の共通性

完全共用語となったディスクリプタは、1 文献群あたりの平均で、MEDLINE が 1.3 語、JMEDICINE が 1.1

語、EMBASE が 0.8 語であり、EMBASE が他よりやや少なかった。タグも含めると、EMBASE と MEDLINE の完全共用語数はほぼ同程度であり、JMEDICINE より多かった。特に、MEDLINE では、8 文献群中 7 文献群に完全共用語があり、「関連のある文献を結びつける」働きが非常に強い索引語が含まれていると考えられる。しかし、いずれの文献群も完全共用語数は少なく、文献群によってデータベースの順位が異なり、一定の傾向を述べることはできない。

共用語となったディスクリプタも同様に、1 文献群あたりの平均で、MEDLINE が 4.6 語、JMEDICINE が 4.5 とほぼ同程度に多く、それに対し EMBASE は 2.0 語と少なかった。MEDLINE で副標目をつけた形でみた場合も 3.9 語であり、EMBASE のおよそ 2 倍であった。すなわち、MEDLINE と JMEDICINE のディスクリプタは、EMBASE のディスクリプタよりも共用語数が多く、付与されたディスクリプタが「関連のある文献を結びつける」働きが強いと考えられる。

また、タグとの合計では、共用語は 1 文献群あたりの平均で EMBASE が 7.8 語、MEDLINE が 7.5 語とほぼ同数となり、いずれも JMEDICINE より多かった。すなわち、EMBASE では、タグは「関連性のある文献を結びつける」働きが強く、ディスクリプタの弱点を補っていると考えられる。

2. 索引語の識別性

共用識別語も、やはり EMBASE が少なかった。しかし、共用語中の共用識別語の割合を調べると、EMBASE が 81.3%、JMEDICINE が 80.6%、MEDLINE が 67.6% であった。すなわち、EMBASE のディスクリプタは、共用語が少なく、「関連性のある文献を結びつける」働きは弱い、共用語の 1 語 1 語が「データベース全体の中から特定の文献群を区別する」働きは他のデータベースより強いといえる。逆に、MEDLINE のディスクリプタは、「関連性のある文献を結びつける」働きは強いが、共用語 1 語ずつが「特定の文献群を区別する」働きが弱い。

他方、EMBASE と MEDLINE では、文献の主要な論点を表す語を主ディスクリプタとして扱う、簡単な重みづけを行なっている。この主ディスクリプタの共用語は、8 文献群の合計で、MEDLINE が 9 語、EMBASE が 11 語であった。MEDLINE ではその全てが共用識別語であったのに対し、EMBASE ではそのうち共用識別語は 10 語であった。すなわち、MEDLINE では、主デ

医学分野のオンライン・データベースにおける索引作業の比較評価

第3表 各データベースの索引作業の定量的比較

文献群 <文献数>	データ ベース名	完全共用語			共用識別語			付与索引語数 (異なり)			付与索引語 数 (延べ)
		D (D+T)	D (D+T)	D (D+T)	D (D+T)	D (D+T)	D (D+T)	D*	D (D+T)	1 文献あた り平均 D (D+T)	1 文献あた り平均 D (D+T)
A <6>	EMBASE	1 (2)	5 (9)	4 (4)	30	42 (63)	7.0(10.5)	9.8(17.0)			
	JMEDICINE	1 (1)	5 (5)	4 (4)	—	33 (33)	5.5(5.5)	9.0(9.0)			
	MEDLINE	2 (3)	4 (8)	3 (3)	11	35 (43)	5.8(7.2)	8.8(13.1)			
B <6>	EMBASE	0 (2)	2 (7)	1 (1)	15	20 (38)	3.3(6.3)	4.8(11.6)			
	JMEDICINE	1 (1)	5 (5)	4 (1)	—	45 (45)	7.5(7.5)	11.3(11.3)			
	MEDLINE	0 (1)	5 (8)	1 (1)	16	29 (35)	4.8(5.8)	7.3(10.1)			
C <5>	EMBASE	0 (0)	0 (7)	0 (0)	16	31 (52)	6.2(10.4)	7.6(16.0)			
	JMEDICINE	0 (0)	2 (2)	2 (2)	—	39 (39)	7.8(7.8)	9.4(9.4)			
	MEDLINE	0 (1)	2 (4)	2 (2)	17	56 (60)	11.2(12.0)	11.2(13.4)			
D <3>	EMBASE	0 (0)	0 (3)	0 (0)	12	16 (32)	5.3(10.6)	5.3(11.6)			
	JMEDICINE	0 (0)	3 (3)	3 (3)	—	32 (32)	10.6(10.6)	11.7(11.7)			
	MEDLINE	0 (0)	3 (5)	2 (2)	7	32 (37)	10.6(12.3)	12.0(14.3)			
E <3>	EMBASE	0 (2)	1 (6)	1 (1)	8	10 (17)	3.3(5.7)	5.3(8.3)			
	JMEDICINE	2 (2)	5 (5)	5 (5)	—	20 (20)	6.7(6.7)	8.3(8.3)			
	MEDLINE	3 (5)	7 (10)	5 (5)	8	19 (23)	6.3(7.7)	9.7(12.6)			
F <5>	EMBASE	0 (0)	1 (7)	1 (1)	21	27 (49)	5.4(9.8)	7.0(14.6)			
	JMEDICINE	0 (0)	5 (5)	4 (4)	—	43 (43)	8.6(8.6)	11.2(11.2)			
	MEDLINE	0 (3)	7 (10)	4 (4)	13	49 (53)	9.8(10.6)	13.2(16.0)			
G <3>	EMBASE	3 (8)	4 (11)	4 (4)	12	13 (24)	4.3(8.0)	6.7(14.3)			
	JMEDICINE	2 (2)	6 (6)	3 (3)	—	21 (21)	7.0(7.0)	9.7(9.7)			
	MEDLINE	5 (6)	5 (8)	5 (5)	8	25 (30)	8.3(10.0)	11.0(14.6)			
H <3>	EMBASE	2 (5)	3 (12)	2 (2)	9	16 (30)	5.3(10.0)	7.0(15.6)			
	JMEDICINE	3 (3)	5 (5)	4 (4)	—	25 (25)	8.3(8.3)	11.0(11.0)			
	MEDLINE	0 (2)	4 (7)	3 (3)	7	18 (24)	6.0(8.0)	7.3(11.0)			
平均	EMBASE	0.8(2.4)	2.0(7.8)	1.6(1.6)	15.4	21.9(38.1)	5.1(9.0)	6.7(13.9)			
	JMEDICINE	1.1(1.1)	4.5(4.5)	3.6(3.6)	—	32.3(32.3)	7.6(7.6)	10.1(10.1)			
	MEDLINE	1.3(2.6)	4.6(7.5)	3.1(3.1)	10.9	32.9(38.1)	7.7(9.0)	10.0(13.0)			

注 D: ディスクリプタ

D+T: ディスクリプタとタグの合計

D*: 主ディスクリプタ

ィスクリプタが付与索引語に占める割合は小さいが、識別力が強く、「データベース全体から特定の文献群を区別する」働きが強い。

また、タグは、いずれも識別性指標が 0.3 未満であり、共用識別語はなかった。タグは、「関連のある文献を結びつける」働きは強いが、「データベース全体の中から特定の文献群を区別する」働きは弱いといえる。

3. 付与索引語数

各データベースが付与した文献群ごとの異なりディス

クリプタ数は、1 文献あたりの平均で *MEDLINE* が 7.7 語、*JMEDICINE* が 7.6 語と多く、*EMBASE* は 5.1 語と少なかった。*MEDLINE* は、副標目をつけた形で計数すると 9.6 語とさらに多かった。また、タグも含めると *MEDLINE* と *EMBASE* は 1 文献あたりの平均がともに 9.0 語であったが、副標目を考慮すると、*MEDLINE* が 10.9 語と最も多かった。*MEDLINE* の索引作業は「文献群の中から特定の 1 文献を区別する」働きが強く、副標目はさらにその働きを強めている。

B. 索引語の質的な検討

以上、索引語を、共通性、識別性、付与索引語数という観点から比較した。しかし、調査対象数が少なく、しかも、統制語彙の特性やデータベース規模の差異が、これらの指標に影響を与えていると考えられる。そこで、統制語彙の特性を考慮しながら、各データベースが付与した索引語を詳しく比較することとした。共用識別語数が多かった文献群 A, E, G をとりあげた。

第 4 表から 6 表は、文献群 A・E・G における各データベースが付与した索引語の比較と文献群を構成する文献の書誌事項を示している。それぞれ、文献群中の 2 件以上の文献に付与したディスクリプタとタグを全て列挙した。2 件未満の文献に付与した索引語でも比較に必要な場合は () に入れて示した。JMEDICINE の索引語は日本語なので、「JICST 科学技術用語シソーラス日英対訳リスト」²⁵⁾ で示された英訳の索引語を用いて比較した。この表を分析することにより、各データベースの索引作成作業や索引言語の特性をみることができる。

各データベースが付与したこれらの索引語を比較し、それぞれの関係を White と Griffith¹⁴⁾ と同様に「一致」「関連」「不一致」の 3 段階に分類した。分類の基準は以下のように定めた。すなわち、「一致」とは、2 つのデータベースで付与した索引語が全く同一の場合である。ただし、大文字と小文字、単数形と複数形、倒置形と倒置していない語句とは、同一と見なした。「関連」とは、互いに意味的に関係があることを示す。具体的には、比較する何れかのデータベースの統制語彙中で USE 参照で示された同義語もしくは準同義語関係が中心であり、これらの関係は実線で示した。省略形と完全な綴りとの関係もこれに含めた。そのほか、シソーラスで上位・下位関係が示されているものや、シソーラスでは関係が示されていないが意味が近いと判断できるものも「関連」とし、点線で示した。「一致」「関連」以外の関係を「不一致」とした。

1. 文献群 A における索引語の比較 (第 4 表)

文献群 A は、A1 から A6 の 6 文献からなっている。第 4 表のごとく、いずれのデータベースでも、ヒト (“human”, “ヒト”, “HUMAN”), 成長ホルモン放出因子 (“growth hormone releasing factor”, “GRH”, “SOMATOTROPIN-RELEASING HORMONE”), 成長ホルモン (“growth hormone”, “成長ホルモン”, “SOMATOTROPIN”) を表す索引語を多くの文献に付与していた。「一致」および「関連」の関係をみると、デ

ータベースごとに、概念を表すラベルとしてのディスクリプタは異なる場合もあるが、だいたい同様の概念を索引していた。しかし、それぞれの索引語を付与した文献数や「不一致」になった索引語を検討することによって各データベースの特徴がいくつか示唆される。

a. 索引語を付与した文献数の差

1) 「ヒト」を表す索引語

MEDLINE と JMEDICINE では A 群中の全文献に付与しているのに対し、EMBASE では文献 A2 には付与していなかった。この A2 は、ヒト成長ホルモン放出因子を仔牛に投与した実験に関する文献であり、実験対象はヒトではなく仔牛である。3 つのデータベースともこの文献 A2 に、仔牛やウシを表す索引語を付与しているが、それとともにヒト成長ホルモン放出因子を用いていることから JMEDICINE と MEDLINE では “HUMAN” (“ヒト”) を付与したと考えられる。それに対し、EMBASE では “human” は、Item Index の「実験対象・被験者 (研究対象となった生物)」というグループに属しており、研究対象としてヒトが用いられた文献にのみ付与していると考えられる。

2) 「成長ホルモン」を表す索引語

文献 A1, A2, A4, A6 には 3 つのデータベースが共通して付与し、A5 には EMBASE だけが付与している。A3 は成長ホルモン放出因子の欠乏によって引き起こされる下垂体性小人症の治療法に関する文献であるが、JMEDICINE では “成長ホルモン” を付与しているのに対し、EMBASE では “human growth hormone” と成長ホルモンの欠乏を意味する “growth hormone deficiency” を付与していた。このように、EMBASE のディスクリプは他のデータベースより特定のであるが、意味が似ている用語を多く含み、主題に関連がある文献に対してこれぞ特定の別索引語を付与する場合がある。これが索引語の共通性の低さの一因と考えられる。なお、MEDLINE と JMEDICINE では、この A3 の文献に下垂体性小人症を意味するディスクリプタを付与していた。

3) 「男性」を表す索引語と「成人」を表す索引語

MEDLINE では男性と女性の両方を扱っている文献も “MALE” を付与しているのに対し、JMEDICINE では男性のみを扱っている文献のみに “男性” を付与している。「成人」 (“ADULT”) も同様であった。EMBASE では、Item Index の語彙中に “adult”, “male”, “female” 等の語が含まれているが、ここでは使用してい

第4表 文献A群における索引語の比較 (文献数6)

一致	EMBASE	JMEDICINE	MEDLINE
	○<T>Human [1 3456]	◎ヒト [123456]	○<T>HUMAN [123456]
	◎growth hormone [12 456]	◎成長ホルモン [1234 6]	
		男性 [1 4]	○<T>MALE [1 456]
		(成人 [1])	○ADULT [1 4 6]
関	(growth hormone deficiency [3])		◎SOMATOTROPIN [12 4 6]
	(human growth hormone [3])		◎—BLOOD [1 4 6]
	◎growth hormone releasing factor [123456]	◎GRH [123 56]	—SERELOGY [2 4]
			◎SOMATOTROPIN-RELEASING HORMONE [123456]
			—ADMINISTRATION & DOSAGE [2 4]
			◎—PHARMACOLOGY [12 4 6]
	◎protirelin [23 6]	—TRH [2 6]	—THYROTROPIN-RELEASING HORMONE [2 6]
			—PHARMACOLOGY [2 6]
	○<T>Intravenous drug administration [12 56]	○静脈内投与 [12 4 6]	—(INJECTIONS, INTRAVENOUS (1))
	gonadorelin [3 6]	GnRH [4 6]	(GONADORELIN—PHARMACOLOGY [1])
	somatomedine [3]	(ソマトメジン [1])	(SOMATOMEDINS—BLOOD [1])
	(somatomedine c [1])		INSULIN-LIKE GROWTH FACTOR I [1 3]
			—BLOOD [1 3]
一致	○<T>Endocrine gland [123456]	◎ホルモン調節 [234]	◎PEPTIDE FRAGMENTS [123456]
	○<T>Priority journal [12345]	血しょう中濃度 [4 6]	—ADMINISTRATION & DOSAGE [2 4]
	○drug efficacy [2 456]		◎—PHARMACOLOGY [12 4 6]
	○hypophysis [1 4 6]		○<T>SUPPORT, NON-U.S. GOV'T [1 3456]
	<T>Benign neoplasm [3 6]		○<T>FEMALE [23 56]
	<T>Human experiment [1 4]		<T>ANIMAL [2 5]
	<T>Immunological procedures [2 5]		<T>COMPERATIVE STUDY [2 6]
	<T>Normal humans [1 4]		<T>SUPPORT, U.S. GOV'T, P.H.S. [1 4]
	<T>Radioisotope [2 5]		
	<T>Therapy [3 6]		
A1:	Takano, K.; Honda, N.; Shizume, K.; Hizuka, N.; Ling, N. C. Plasma growth hormone and somatomedin-C responses to continuous growth hormone-releasing factor infusion in normal adult men. Endocrinol. Jpn., 32(2): 287-93 (1985)		
A2:	Hodate, K.; Johke, T.; Ohashi, S. Growth hormone, thyrotropin and prolactin responses to simultaneous administration of human growth hormone-releasing factor and thyrotropin releasing hormone in the bovine. Endocrinol. Jpn., 32(3): 375-83 (1985)		
A3:	Takano, K.; Hizuka, N.; Shizume, K. Effects of hGRF treatment of a patient with hGRF deficiency. Endocrinol. Jpn., 32(4): 511-6 (1985)		
A4:	Hotta, M.; Shibasaki, T.; Masuda, A.; Imaki, T.; Wakabayashi, I.; Demura, H.; Ling, N.; Shizume, K. The inter- and intra-subject variabilities of plasma GH response to human growth hormone-releasing hormone (1-44) NH2 in men. Endocrinol. Jpn., 32(5): 673-80 (1985)		
A5:	Minami, S.; Wakabayashi, I.; Tonegawa, Y.; Sugihara, H.; Akira, S. Production of antisera to growth hormone-releasing factor: usefulness in radioimmunoassay and passive immunization. Endocrinol. Jpn., 32(6): 907-16 (1985)		
A6:	Hanew, K.; Sato, S.; Sasaki, A.; Shimizu, O.; Murakami, O.; Yoshinaga, K. Comparative study on the responses of plasma GH to synthetic GH-releasing factor and other stimulatory and inhibitory agents in patients with acromegaly. Tohoku. J. Exp. Med., 145(2): 161-6 (1985)		

<T> はタグ、それ以外はディスタクリプタ。◎ は共用識別語。○ は共用語。--- 一致・use 参照。--- 上位下位関係、関連語等。[] 内の数字は当該索引語が付与された文献の番号、番号の下線は主ディスタクリプタを示す。
MEDLINE における --- は副標目を示す。

ない。非統制語であるアイデンティファイアのフィールドには、“4 healthy male” (A1), “girl of 12” (A3), “13 healthy men” (A4), “10 patients” (A6) が付与されていた。これらのアイデンティファイアは文献の内容を推し量るにはふさわしいが、多様な語を使用しており、網羅的な検索は困難である。

4) 「成長ホルモン放出因子」を表す索引語

これは、この文献群を結びつけている中心的な概念である。*EMBASE* と *MEDLINE* では、これを表す索引語を A 群中の全ての文献に付与しているのに対し、*JMEDICINE* では A4 に付与していなかった。A4 は、標題にも human growth hormone-releasing hormone という語句を含み、成長ホルモン放出因子に関して述べていると考えられる。そのほか、*JMEDICINE* は、“ソマトメジン”、“男性”、“成人”などの索引語でも、他のデータベースより付与した文献数が少なかった。このような付与文献数の差異は、それぞれのデータベースにおける索引作成方針や索引語の適用規則の影響もあると考えられる。

5) 薬剤の投与方法に関する索引語

静脈内投与に関連する索引語では、*MEDLINE* の統制語彙である *MeSH* には“INJECTIONS, INTRAVENOUS (静脈注射)”しかなく、他のデータベースと比べて付与した文献数が少なかった。*MeSH* の主標目には、治療法などの操作を表す語は比較的少ないが、副標目がそれを補う働きをしている。例えば、この文献群中では、薬物の投与方法を示す“—ADMINISTRATION”, 血中濃度を示す“—BLOOD”, 薬物を治療目的で使用することを示す“—THERAPEUTIC USE”などの副標目が多く用いられていた。

6) 「ソマトメジン」とその下位概念を表す索引語

「JICST シソーラス」では“ソマトメジン”というディスクリプタのみであるが、*MeSH* では“SOMATOMEDINS”の下位語に“INSULIN-LIKE GROWTH FACTOR I”があり、*MALIMET* ではこれに相当するものとして“somatomedine c”がある。このソマトメジンに関して、*JMEDICINE* は付与した文献数が他より少なく、特定の概念を表す索引語もなかった。

b. 「不一致」の索引語

EMBASE では、身体部位や広い概念を表す語が多かった。これらの語は、「一致」や「関連」関係になった索引語の上位の概念や関連のある概念を表しているものが多く、タグも含めた索引語の付与数は多かったが、必

ずしも幅広く概念を索引した網羅的な索引作業がなされているわけではない。“priority journal”というタグは文献の種類を示し、他のデータベースとは異なっている。

JMEDICINE では、“ホルモン調節”と“血しょう中濃度”が不一致であった。このような操作や状態を表す語は、化学物質や身体部位、疾患名などと比べて、概念の規定のしかたが難しく、各データベースごとに扱い方が異なる場合が多いと考えられる。たとえば、*MEDLINE* では、このような概念を表すのに副標目を活用する場合が多く、*JMEDICINE* で“ホルモン調節”や“薬物療法”を付与した文献には、*MEDLINE* では個々のホルモンを表す主標目に“—ADMINISTRATION”や“—THERAPEUTIC USE”という副標目を付加し、*JMEDICINE* で“血しょう中濃度”を付与した文献には、*MEDLINE* では成長ホルモンの血中濃度を示す“SOMATOTROPIN—BLOOD”を付与している。

MEDLINE では、“PEPTIDE FRAGMENTS”が他とは異なる索引語である。“GRH”, “GnRH”, “TRH”, “成長ホルモン”などはみなペプチドホルモンであり、*MeSH* ではこのディスクリプタによってこの文献群の結びつきがより強くなっている。また、研究助成機関を示すタグや研究の種類を示す“COMPARATIVE STUDY”などのタグは、他のデータベースとは異なる特徴的なものである。

2. 文献群 E における索引語の比較 (第 5 表)

第 5 表のごとく、3 つのデータベースは、共通して、グリコサミノグリカン (“glycosaminoglycan”, “GLYCOSAMINOGLYCANS”, 上位の概念である“ムコ多糖類”), ウサギ (“rabbit”, “RABITS”), 子宮 (“uterus”, “UTERUS”) に関連する索引語を多くの文献に付与していた。*MEDLINE* と *JMEDICINE* では、“HYALURONIC ACID” (“ヒルアロン酸”), “CHONDROITIN SULFATES” (“コンドロイチン硫酸”) も一致した。

1) 「グリコサミノグリカン」を表す索引語

MEDLINE と *EMBASE* では“GLYCOSAMINOGLYCANS” (“glycosaminoglycan”) を付与しているが、*JMEDICINE* ではシソーラスにグリコサミノグリカンに相当する語がなく、準ディスクリプタとして、自然語キーワードでグリコサミノグリカンを付与している。*JMEDICINE* が付与したディスクリプタである“ムコ多糖類” (“MUCOPOLYSACCHARIDES”) は、*MeSH* では“GLYCOSAMINOGLYCANS”の上位語

第5表 文献群Eにおける索引語の比較 (文献数3)

	EMBASE	JMEDICINE	MEDLINE
一致 関連 不一致	◎glycosaminoglycan [12]		◎GLYCOSAMINOGLYCANS [123] ◎—ANALYSIS [12] ◎—METABOLISM [1 3]
	○〈T〉Rabbits and heres [1 3] (rabbit [1])	◎ウサギ [1 3] ◎コンドロイチン硫酸 [123] ◎ヒアルロン酸 [23] ◎子宮 [23]	○RABBITS [1 3] +CHONDROITIN SULFATES [123] ◎—ANALYSIS [12] ◎HYALURONIC ACID [23] + (UTERUS—ANALYSIS [2])
	○〈T〉Female genital system [23]	◎ムコ多糖類 [123] (高速液体クロマトグラフィー [1])	○CHROMATOGRAPHY, HIGH PRESSURE LIQUID [1 3]
	○〈T〉Animal tissue, cells or cell components [123] ○〈T〉Chemical procedures [123] ○〈T〉Nonhuman [1 3]		◎DERMATAN SULFATES [123] ◎—ANALYSIS [12] ◎MYOMETRIUM [23] ○〈T〉ANIMAL [123] ○〈T〉FEMALE [23] ○〈T〉SUPPORT, NON-U.S. GOV'T [123]

- E1: Munakata, H.; Isemura, M.; Aikawa, J.; Kodama, C.; Yosizawa, Z. Glycosaminoglycans from renal brush border membranes of rabbits. *Tohoku. J. Exp. Med.*, 145(4): 353-8 (1985)
- E2: Munakata, H.; Isemura, M.; Kodama, C.; Yosizawa, Z. Glycosaminoglycans of porcine uteri. *Tohoku. J. Exp. Med.*, 147(1): 73-5 (1985)
- E3: Munakata, H.; Isemura, M.; Aikawa, J.; Kodama, C.; Yosizawa, Z. Changes of glycosaminoglycan composition of uterine myometrium of rabbit induced by female sex steroids. *Tohoku. J. Exp. Med.*, 147(1): 77-81 (1985)

〈T〉はタグ、それ以外はディスクリプタ。◎は共用識別語。○は共用語。—— 一致・use 参照。--- 上位下位関係、関連語等。[] 内の数字は当該索引語が付与された文献の番号、番号の下線は主ディスクリプタを示す。

MEDLINE における —— は副標目を示す。

になっている。

2) 「子宮」を表す索引語

E3 は、論文の題名から子宮筋層 (myometrium) に関して述べていると思われる。それに対し、MEDLINE では“MYOMETRIUM”を主要な論点を表す主ディスクリプタとして付与しているのに対し、JMEDICINE ではシソーラス中に子宮筋層に相当する語がなく、その上位概念にあたる“子宮”を付与している。EMBASE では、E3 の文献に子宮や子宮筋層に関係するディスクリプタを付与していなく、タグとしてさらに上位概念に相当する女性の生殖器系をしめす“female genital system”を付与しているだけであった。

3) データベースごとの傾向

EMBASE は、この文献群においては、全体の付与索引語数が少なく、他のデータベースと比べて“glycosaminoglycan”や“uterus”などの索引語を付与した文献数も少なく、索引語の共通性が低い。また、より特定の別の索引語を付与したために共通性が低下しているわけでもない。さらに、“rabbit”とタグの“rabbits and hares”, “uterus”とタグの“female genital system”のように、ディスクリプタと意味が近い、あるいはその上位概念を表しているタグとを同じ文献に重複して付与している場合が多い。したがって、付与された索引語が表現している概念は付与された索引語数ほど多くないと考えられる。したがって、この文献群では、

EMBASE の索引作業の特定性と網羅性はやや低い。しかし、このような場合にも、タグは、索引もれを防ぎ、索引語の共通性、すなわち「関連のある文献を結びつける」働きを高める機能を果たしている。

JMEDICINE は、この文献群において、“ムコ多糖類”や“子宮”の例が示すように、シソーラスに特定の語がないために、結果として索引作業の特定性が低くなっている。また、自然語の準ディスクリプタとして、“グリコサミノグリカン”、“高圧液体クロマトグラフィ”のような特定の語を付与している場合もあり、検索の一助となる。

MEDLINE は、この文献群においては、他のデータベースより付与索引語数が多いことから、表現している概念の種類も多いと考えられる。特定の索引語が多く付与され、索引語の共通性も高かった。

3. 文献群 G における索引語の比較 (第 6 表)

第 6 表のごとく、ヒト (“human”, “HUMAN”) と硝

子体 (“vitreous body”, “VITREOUS BODY”) は 3 つのデータベースで一致して付与している。さらに、EMBASE と MEDLINE では、角膜の検査に用いる薬剤であるフルオロセイン (“fluorescein”, “FLUORESCINS”) も一致している。

1) 「フルオロフォトメリー」を表す索引語

EMBASE では、“fluorophotometry” という特定の索引語を付与しているが、MEDLINE では fluorophotometry の上位概念を表す “FLUOROMETRY” と “PHOTOGRAPHY” を付与し、この 2 語の組合せで fluorophotometry という概念を表している²³⁾。

2) 「診断」を表す索引語

EMBASE ではタグに “diagnosis”, JMEDICINE では“診断”を付与した文献に対し、MEDLINE では、検査に使用した薬剤である “FLUORESCIENS” にその薬剤を診断のために使用したことを示す “—DIAGNOSTIC USE” という副標目を付加している。MEDLINE

第 6 表 文献群 G における索引語の比較 (文献数 3)

	EMBASE	JMEDICINE	MEDLINE
一 致	◎vitreous body [123]	◎硝子体 [1 3]	◎VITREOUS BODY [123] ◎—METABOLISM [123]
	○<T>Human [123]	◎ヒト [23]	○<T>HUMAN [123]
	◎fluorescein [1 3]		◎FLLUORESCINS [123] ◎—DIAGNOSTIC USE [123]
関 連	◎fluorophotometry [123]	◎蛍光 X 線分析 [123]	◎FLUOROMETRY [123] ◎PHOTOGRAPHY [123]
不 一 致	◎blood retina barrier [123]	○眼疾患 [123]	◎PERMEABILITY [123]
	○<T>Blood and hemopoietic system [123]	○生体機能検査 [12]	○<T>FEMALE [23]
	○<T>Diagnosis [123]	○予後 [23]	○<T>MALE [23]
	○<T>Priority journal [123]		
	○<T>Visual system [123]		
	○<T>Clinical article [23]		
	○<T>Peripheral vascular system [12]		

G1: Maurice, D. M. Theory and methodology of vitreous fluorophotometry. Jpn. J. Ophthalmol., 29(2): 119-30 (1985)

G2: Krupin, T.; Waltman, S. R. Fluorophotometry in juvenile-onset diabetes: long-term follow-up. Jpn. J. Ophthalmol., 29(2): 139-45 (1985)

G3: Miyake, K. Vitreous fluorophotometry in aphakic or pseudophakic eyes with persistent cystoid macular edema. Jpn. J. Ophthalmol., 29(2): 146-52 (1985)

<T> はタグ、それ以外はディスクリプタ。◎ は共用識別語。○ は共用語。—— 一致・use 参照。--- 上位下位関係、関連語等。[] 内の数字は当該索引語が付与された文献の番号、番号の下線は主ディスクリプタを示す。

MEDLINE における —— は副標目を示す。

では、副標目を用いることにより、索引語間の関係が明確であり、検索時の誤結合を防ぐと同時に、より特定の概念を表現していると考えられる。

3) 「不一致語」とデータベースごとの傾向

EMBASE は、この文献群では、付与索引語は他のデータベースより少ないが、完全共用語と共用識別語が多く「関連のある文献群を結びつけ」、「データベース全体から特定の文献群を区別する」力は充分である。しかも、「fluorophotometry」や「blood retina barrier」など特定の索引語を付与しており、文献の主題を充分に表現していると考えられる。

それに対し、*JMEDICINE* では、「ヒト」、「硝子体」、「診断」などの索引語では他データベースより付与文献数が少なく、そのため完全共用語数が少ない。また、これらに関連する他の索引語を代わりに付与しているわけでもない。同様に、女性と男性の両方を扱っている G2 と G3 には「女性」も「男性」も付与せず、ヒトと動物との両方を扱っている G1 には「ヒト」を付与していない。さらに、「fluorescien」などの物質名や「blood retina barrier」などにあたる索引語を付与していない。その一方で、「予後」や「診断」、またこの文献群ではそれぞれ 1 文献ずつのみであり、表中には示していないが、「理論」や「方法論」などの一般的で大きな概念を示す索引語を付与していた。

MEDLINE は、この文献群では、完全共用語が多く、付与索引語数も多いが、その中には、「FLUOROMETRY」と「PHOTOGRAPHY」のように索引作業の特定性を高めるために付与数が多くなっているものも含まれており、単純に、付与数が多いから網羅的な索引作業であるとはいえない。その一方で、副標目により、より特定の概念を表現することが可能になっている。

なお、*MEDLINE* の統制語彙である *MeSH* は、今回の調査では *EMBASE* だけにあった「BLOOD RETINA BARRIER」を 1987 年から、より特定の索引語がなかった「FLUOROPHOTOMETRY」を 1990 年から新ディスクリプタとして採用している。

C. 付与索引語からみた各データベースの索引作業

以上の分析から明らかになった、各データベースの索引作業の特徴をまとめると、以下のようになる。

1. *EMBASE*

EMBASE は、ディスクリプタのみでは、他のデータベースよりも、共用語数、共用識別語数、一文献あたり

の付与数ともに少なかったが、共用語に対する共用識別語の比率は高かった。したがって、ディスクリプタは「関連のある文献を結びつける」働きは弱い、共用語となったディスクリプタ 1 語 1 語の識別性は比較的高いといえる。タグも含めた索引語全体では、付与索引語数、共用語数ともに多かった。タグは、「関連性のある文献を結びつける」働きが強く、ディスクリプタの弱点を補う働きがあるが、「データベース全体から特定の文献群を区別する」働きは非常に弱い。

質的に検討すると、*EMBASE* のディスクリプタは、自然語に近く、極めて特定性の高いものが多かった。したがって、文献群 G の「fluorophotometry」の例のように付与数が他より少なくても特定の語で十分に内容を表現している場合もあった。一方、文献群 A の「growth hormone」と「human growth hormone」の例にみられたように、他のデータベースでは同一の索引語を付与している文献に対しても、より特定の、あるいは意味が非常に近い別のディスクリプタを付与しているために共用語数が少なくなっている例もあった。このような場合は、ディスクリプタの共通性が低くても、より特定の索引語を用いて適切に、より特定の文献の内容が表現されていると考えられる。

一方、*MALIMET* は、*EMTREE* コードが導入される以前は、上位・下位などの索引語間の関係があらかじめ規定されていないディスクリプタの典拠リストであり、シソーラスではなかった。したがって、このような特定のディスクリプタだけで検索を行なう場合、網羅的に検索するには、下位概念や類似の概念を表すディスクリプタ全ての論理和をとる必要があり、困難である。逆に、タグとディスクリプタ、あるいはディスクリプタ同士で、非常に意味に近い語や上位下位関係にある語を重複して付与している例が多くみられた。その場合、たとえ付与索引語数が多くても、索引語が表現している概念の種類は索引語数より少なく、文献中の概念を幅広く抽出した網羅的な索引作業がなされているとはいえない。

これに対し、1988 年から使用頻度の高いディスクリプタからなる *MiniMALIMET* に階層構造を持った *EMTREE* コードが導入されているが、上位・下位などの索引語間の関係を明示することになり、網羅的な検索や不必要な索引語の重複を避ける上で有効であると思われる。また、タグも、網羅的な検索に有効である。

さらに、文献群 E のように、共用語数、付与索引語

数、同一索引語を付与した文献数共に他のデータベースより少なく、しかも特定の語も付与していない例もみられたが、その理由は明かではない。

2. JMEDICINE

JMEDICINE は、ディスクリプタだけで比較すると、共用語数、共用識別語数、付与索引語数ともに多く、バランスが取れた索引作業に思われる。タグを含めて比較すると、共用語数や付与索引語数は他より少ないが、共用識別語はやはり 3 つのデータベースの中で最も多かった。

質的に検討すると、共用語数の多さは、ひとつには語彙の特性性が低いためであり、識別性の高さは、データベース規模が小さいことと、各索引語を付与した文献数が他より少ないことが影響していると考えられる。すなわち、JMEDICINE では統制語彙は、科学技術全般を対象とし、医学領域で使用するディスクリプタは比較的少なく、用語の特定性も低い。また、そのため、“ソマトメジン”、“子宮筋腫”、“ムコ多糖類”などのように、他のデータベースではより特定の別の索引語をそれぞれ付与している文献に対しても、より上位の同一の索引語を付与している場合があった。このように語彙の特定性の低さが索引語の共通性の高さに影響を及ぼしていると考えられる。

また、JMEDICINE の 1985 年出版の文献収録数は、第 3 表に示したごとく、EMBASE と MEDLINE より少ない。しかも「JICST シソーラス」を用いた索引作業がなされるのはその 4 割程度であり、MEDLINE や EMBASE のおよそ 1/3 と考えられる。今回用いた識別性指標は、索引語の使用回数から算出しているので、データベースの規模、すなわち当該シソーラスを用いて索引語を付与する 1985 年出版の文献収録数が小さければ、索引語 1 語 1 語の使用回数が少なく、そのために各語の識別性指標が、規模の大きい他のデータベースよりも高くなっていると考えられる。

さらに、JMEDICINE では各索引語を付与した文献の数が、他のデータベースより少ない場合が多くみられた。“男性”、“成人”、“ソマトメジン”（以上第 4 表参照），“硝子体”、“ヒト”（以上第 6 表参照）などがその例である。これは、索引作成方針や索引語の適用範囲の違いによって生じていると考えられる。これによって、検索もれが生じることも考えられるが、その一方で、個々の索引語の使用回数が少なくなり、索引語の識別性指標は高くなると考えられる。

JMEDICINE では、索引語の上位語の自動付加を行ない、より広い網羅的な検索については考慮しているが、重みづけは副標目など、より特定の検索するためのデバイスはない。今後、データベースの年間収録レコード数が拡大する、あるいは、医中誌基本ファイルからの抽出分にも同一のシソーラスから索引語を付与する場合には対象レコードが増加し、索引語の識別性がかなり低下すると予想される。それに対処するには、統制語彙中の医学文献の索引に使用する用語を増やすと共に、重みづけや副標目などのデバイスを導入して、より柔軟な検索を可能にすることが必要であろう。

3. MEDLINE

MEDLINE では、ディスクリプタのみでも、タグを加えた場合でも、共用語数と 1 文献あたりの付与索引語数が多かった。さらに、副標目を付加した形で考えると、付与索引語数は非常に多くなり、「関連のある文献を結びつける」と共に、「文献群の中から特定の文献を区別する」働きが非常に強いといえる。しかし、共用識別語数は、EMBASE より多いが、JMEDICINE より少なく、共用語中の共用識別語の割合は最も低く、索引語 1 語 1 語では識別力、すなわち「データベース全体の中から特定の文献群を区別する」働きは弱い。

しかし、質的に検討すると、以下のような特徴がみられた。①各索引語を付与した文献数が他のデータベースよりやや多い場合が多く、索引語の適用規則が広く、網羅的に付与されている、②索引語の特定性は、EMBASE よりやや低い場合があった、③索引作業の特定性を高めるために、1 つの概念を複数の索引語で表現している場合がみられた、④新しい索引語の取入れが、EMBASE と比べてやや遅い。すなわち、網羅的に索引語を付与し、しかも特定性を高めるために索引語を多く付与している場合もあるために、全般的に、索引語の使用回数が多く、そのことが共用語数の多さにも影響をし、識別性指標を低くしていると考えられる。しかし、1 文献あたりの付与索引語数が多く、副標目も用いられているため、索引語 1 語 1 語の識別性は低くても、索引語を組み合わせた場合の識別性は相対的に高くなっていると考えられる。

索引語の重みづけや副標目の利用、チェックタグは、いずれも、統制語彙の規模を拡大させずに、より特定の検索を可能にし、個々の状況や目的に応じて柔軟な検索を可能にするのに有効に機能していると思われる。

IV. 統制語彙との関係を踏まえた 評価方法の検討

A. 定量的な評価指標の検討

上述のごとく、質的に検討すると、索引語の共通性、識別性、付与索引語数という定量的な評価の観点からは、さまざまなものの影響を受けていることが示唆された。

「索引語の共通性」に関しては、共用語数は、①特定の語がないために上位語で代用する、②特定の語があれば1語で表現できる概念を表すために複数の索引語を付与する、③上位語または類似語を重複して付与している、④索引方針や索引語の適用範囲が広いなどの影響によって増加すると考えられる。

したがって、共用語数が少なく、「関連のある文献を結びつける」働きが弱いと評価された場合でも、必ずしも索引作業の質が低いわけではない。質的に検討し、それが索引もれのためであるか、あるいは語彙の特性の影響であるかなどの要因を明らかにする必要がある。

「索引語の識別性」は、統制語彙の規模が大きい、あるいは、データベースの規模が小さい場合には、一般に高くなる。それだけでなく、索引もれがあったり、索引方針や索引語の適用範囲などの影響により同一文献群に対する索引語の付与数が他のデータベースより少ない場合にも、それぞれの索引語の使用回数が少なくなり、識別性指標が高くなると考えられる。したがって、識別性に関しても、質的に検討することが必要である。

「一文献あたりの索引語の付与数」は、従来は、文献中の索引すべき概念を漏れなく、かつ、幅広く索引しているかどうかを示す「索引作業の網羅性 (exhaustivity)」を示す指標として捉えられてきた。しかし、「索引語の共通性」の②の例と同様に、索引作業の特性性を高めるために付与数が増える場合や③のように非常に近い意味を表す索引語を重複して付与する場合なども考えられる。したがって、「索引語の付与数」は、「索引作業の網羅性」とは一致しない指標である。

また、この「一文献あたりの索引語の付与数」が多いと、一般に、各索引語の使用回数が増加するために、「索引語の識別性」は低下し、逆に「索引語の共通性」を高める可能性が高い。

B. 定量的な評価指標と統制語彙の特性

これらの3つの指標は、統制語彙の規模や特性とも密接な関連がある。すなわち、統制語彙の規模が大きい場

合には、一般に特定の用語を多く含んでいると考えられる。そのために、前述のごとく、「索引語の識別性」は高くなると考えられるが、「共用語数」や「一文献あたりの付与索引語数」が他のデータベースより少なくなる。一方、統制語彙の規模が小さい場合には、「一文献あたりの付与索引語数」と「共用語数」が多くなり、各索引語の使用回数が増えるために「索引語の識別性」が低くなる。

また、シソーラスのようにあらかじめ上位・下位などの索引語間の関係が明示された統制語彙と比べ、関係が明示されていない語彙では意味的に近い索引語を重複して付与することがあり、付与数が多くなりやすい。

さらに、統制語彙の規模は、データベースの規模、すなわち当該語彙が対象とする収録レコード数との関係で捉えるべき問題である。すなわち、統制語彙の規模が小さくても対象とする収録レコード数が少なければ十分な「索引語の識別性」が得られ、逆に統制語彙の規模が大きくても対象レコード数も多い場合には「索引語の識別性」が低くなると考えられるからである。

したがって、統制語彙の規模に対して、データベースの規模が大きい場合には、「一文献あたりの付与索引語数」と「共用語数」が多くなり、「索引語の識別性」が低くなる。このようなデータベースでは、「索引語の識別性」を高めるためには、①統制語彙の用語数を増やす、②副標目を導入する、③索引語に重みづけを行なう、④研究の種類など主題以外の特性を表す語彙を増やすことなどが考えられる。しかし、語彙の規模が大きくなるほど、その維持管理や索引語の付与に一貫性を持たせることが次第に困難になる。特に、統制語彙内で十分に語間の意味的な関係が規定されていない場合は、包括的な検索が困難になり、それを助ける手だてが必要である。副標目を導入した場合は、統制語彙の規模を増加させずに、識別性を高めることが可能である。しかも事前に、主標目と副標目とを結合させるため、検索時に概念の誤結合によるノイズを防ぐ効果もある。ただし、副標目を導入しても、表現できる主となる概念の種類は一定なので語彙を拡大することも必要である。

また、データベースの利用者の状況や利用目的によって、適切だと考えられる検索集合の規模は異なり、「出力過多」だと判断されるレベルも異なる。したがって、利用者の状況や目的に応じた柔軟な検索を可能にすることが必要であり、副標目・重みづけ・索引語の階層化・タグなどの主題以外の文献の特徴を示す手段・リンク・

ロールなどのさまざまなデバイスを導入することが有用である。どのようなデバイスを導入するかは、領域やデータベースの利用者の特性を踏まえ、索引作業の負荷と検索上の効果との関係において検討すべきである。

一方、統制語彙の規模に対してデータベースの規模が小さい場合には、「索引語の識別性」は高くなるが、「索引語の共通性」が低くなる。したがって、このような場合には、索引作業の特定性を維持しながら、網羅的な検索も容易にするために、①統制語彙に階層構造をもたせ、それを用いて下位概念を含めた包括的な検索を可能にする、②索引語の上位概念を表すようなタグを導入することなどが考えられる。

C. 統制語彙の特徴からみた索引語の質的な検討

各データベースが付与した索引語の質的な比較では、化学物質名・身体部位・疾患名は「一致」または「関連」が多かったのに対し、操作や状態を示す語は「不一致」が多かった。それぞれの統制語彙中でのこれらの概念を表す語の有無や位置づけが異なっていることが影響していると考えられる。

研究対象は、専門的な知識がなくても把握しやすい要素であり、「一致」しやすいと考えられるが、実際は「不一致」が多かった。*MEDLINE* では少しでもその対象を扱っていれば漏れなく付与しているのに対し、*JME-DICINE* ではその対象のみを単独で扱っているときのみに付与し、*EMBASE* ではタグの語彙中に該当する用語があるのに、あまり使用せず、アイデンティファイアとして具体的な年齢や被験者数を合わせて示すことが多いという違いがみられた。これには索引方針が影響していると考えられる。

研究の種類、文献の種類、助成金の種類などを表す索引語が付与されているかどうかは、統制語彙中にそのような用語が含まれているかどうかにか依存している。近年医学では、医学雑誌の審査²⁶⁾や論文の読み方²⁷⁾などにおいて論文の質を表すものとして、研究デザインを厳しく問う傾向が強まっている。これは、統制実験か、実験群の割当が無作為か、二重盲検法かなど、統計的に妥当な結果が得られるように配慮されているかを問うものである。研究の種類を表す索引語は、この研究デザインを示すものであり、実際の医学文献の検索においては重要な手がかりを提供するものと考えられる。

タグは、研究の種類など主題以外の側面から文献の特性を表現し、文献の利用目的に応じた柔軟な検索を可能

にする。また、「索引語の共通性」を高め、網羅的な検索を容易にする。*EMBASE* の一部のタグのように、ディスクリプタの上位概念を表している場合は、索引する概念の種類が多くなっているわけではない。

副標目は、語彙数を増やさずに索引語の特定性を高めることを可能にする。同一の主標目に対して複数の副標目の組合せがあることから、付与索引語数が増加する傾向があり、文献群の中から特定の文献を区別する働きも強まる。検索において索引語の誤結合を防ぐ働きもある。

以上のごとく、今回は、「索引語の共通性」、「索引語の識別性」、「一文献あたりの索引語の付与数」という観点から索引作業の比較評価を行なった。これらの観点を示す指標はさまざまな要因が影響し、互いに関連しており、質的に検討することも必要であった。しかし、実稼働データベースにおいて容易に比較評価ができる実用的な方法であり、これによって、各データベースの索引作業の特性を明らかにし、索引作業および統制語彙の問題点や改善方針を指摘することができると思われる。

また、統制語彙規模とデータベース規模の関係は、これらの指標に大きな影響がある。この両者の適切な関係に関しては、実務的には重要性が認識されているが、ほとんど研究がなされていない。

さらに、シソーラスはそれぞれ独自の体系を持ち、複数のシソーラス中の用語をシソーラスで規定された語間の関係やスコープノートだけによって対応づけることは困難がある。今回、主題に共通性のある文献群を対象に付与索引語を質的に検討した方法は、異なる体系の中に位置づけられている用語を互いに対応づける実際的な方法の一つであり、今後、UMLS (Unified Medical Language System) などのように複数のシソーラスの対応づけを考えていく場合にも、何等かの示唆を与えようとする。

V. 総 括

White と Griffith の方法に準拠して、日本で利用できる主要な医学文献データベースの索引作業の比較評価を行なった。その結果、以下のことが示唆された。すなわち、

1) *EMBASE* のディスクリプタは特定性が高いため、識別性は高いが、共通性が低い。網羅的な検索のためには、タグ・分類・1988 年から導入された EMTREE コードが有用であり、柔軟な検索のためにはタグ・重み

づけ・リンクも有用である。膨大な語彙の管理が今後の課題である。

2) *JMEDICINE* の索引作業は、共通性・識別性とも比較的高かったが、シソーラス規模が小さいために特定性が低く、デバイスが少ないために柔軟性が欠けている。JICST 作成分と医中誌抽出分とが 1 つの索引体系になることが望ましい。その場合、1 つの統制語彙が対象とするレコード数が拡大するため、索引語の識別性を維持するには、統制語彙を拡大して特定の用語を増やすと共に、種々のデバイスを導入することが一層必要になる。

3) *MEDLINE* は、付与索引語数が多く、デバイスの種類も多いが、語彙の規模が小さいので、データベース規模との関係において充分であるかを常に検討する必要がある。

4) 大規模な実稼働データベースでは、統制語彙が充分な大きさであると同時に、種々のデバイスが使用できることが望ましい。

5) 統制語彙規模とデータベース規模との適切な関係を、語彙の維持管理や索引付与の作業負荷をも考慮して、研究する必要がある。

なお、調査の実施、および、資料の利用に関し、東京慈恵会医科大学医学情報センターの方々に多くのご配慮をいただいた。ここに記して謝意を表します。

- 1) Bourne, Chales P. Evaluation of indexing system. *Annual Review of Information Science and Technology*. Vol. 1, p. 171-190 (1966)
- 2) Rees, Alan. Evaluation of information systems and services. *Annual Review of Information Science and Technology*. Vol. 2, p. 63-86 (1967)
- 3) Lancaster, F.W.; Elliker, Calvin; Connell, Harkness Tschera. Subject Analysis. *Annual Review of Information Science and Technology*. Vol. 24, p. 35-84 (1989)
- 4) Lancaster, F.W. Indexing and Astracting in Theory and Practice. London, the Library Association, 1991, 328 p.
- 5) Leonard, Lawrece E. Inter-indexer consistency studies 1945-1975: a review of the literature and summary of study results. University of Illinois Graduate School of Library Science, 1977, 51 p. (Occasional Papers, No. 131)
- 6) Sievert, MaryEllen; Andrews, Mark J. Indexing consistency in Information Science Abstracts. *Journal of the American Society for Information Science*. Vol. 42, No. 1, p. 1-6

(1991)

- 7) Blair, David C.; Maron, M.E. An Evaluation of Retrieval effectiveness for a full-text document retrieval system. *Communication of the ACM*. Vol. 28, No. 3, p. 289-299 (1985)
- 8) Boyce, Bert R.; Mclain, John P. Entry point depth and online search using a controlled vocabulary. *Journal of the American Society for Information Science*. Vol. 40, No. 4, p. 273-276 (1989)
- 9) Rolling, L. Indexing consistency, quality and efficiency. *Information Processing and Management*. Vol. 17, No. 1, p. 69-76 (1981)
- 10) Blair, D.C. Language and Representation in Information Retrieval. Amsterdam, Elsevier Science Pub., 1990, 335 p.
- 11) Watkins, Steven G. The IRL Life Sciences Collection and BIOSIS: a comparison of online access to the literature of biology. *Database*. Vol. 4, No. 3, p. 39-59 (1981)
- 12) Lingenfelter, Judith; Gratch, Bonnie; Chan, Betty. A comparison of five physical education indexing/abstracting services. *RQ*. Vol. 21, No. 1, p. 53-60 (1981)
- 13) Siever, MaryEllen; Verbeck, Alison. The indexing of tae literature of online searching: a comparison of ERIC and LISA. *Online Review*. Vol. 11, p. 95-104 (1987)
- 14) White, Howard D.; Griffith, Belver C. Quality of indexing in online data bases. *Information Processing and Management*. Vol. 23, No. 3, p. 211-224 (1987)
- 15) Salton, G.; McGill M. J. Introduction to Modern Information Retrieval. New York, McGraw-Hill, 1983, 448 p.
- 16) Chu, Clara M.; Ajiferuke, Isola. Quality of indexing in library and information science databases. *Online Review*. Vol. 13, No. 1, p. 11-35 (1989)
- 17) 慶應義塾大学文学部図書館・情報学科. 文献間の類似性を測定する尺度としての共引用の妥当性についての評価:「情報学」関連文献を事例として. 平成2年度慶應義塾大学学事振興基金による研究(共同研究)「主題分野のマッピングの手法に関する研究」報告書. 東京, 慶應義塾大学文学部図書館・情報学科, 1991, 77 p.
- 18) Ajiferuke, Isola; Chu, Clara M. Quality of indexing in online databases: an alternative measure for a term discriminating index. *Information Processing and Management*. Vol. 24, No. 5, p. 599-601 (1988)
- 19) JOIS 活用の手引きⅡ; データベース基礎編; (2) JICST・医中誌国内医学文献ファイル. 東京, 日本

- 科学技術情報センター, 1987, 45 p.
- 20) EMBASE Guide to EMTREE and indexing systems, Vol. 1-2. Amsterdam, Excerpta Medica/EMBASE Publishing Group, 1989.
 - 21) Stern, Barrie T. Excerpta Medica. Tokyo, Marzen, 1978, 23 p. (Marzen Database Session 専門セミナー)
 - 22) リンクは MeSH の副標目と同様に, ディスクリプタと組み合わせて使い, そのディスクリプタの表す意味をより特定のにするものである. Drug Links と Medical Links の 2 種類があり, 前者は MALIMET 中の薬学関係の語と組み合わせるもので, adverse drug reaction (副作用), drug administration (投薬) などの 15 語, 後者は医学関係の語と組み合わせるもので, diagnosis (診断法), etiology (病因) などの 12 語からなる.
EMTREE コードは, 15 のファセットに分けられ, 最大 12 段階までの階層構造を持つ. MALIMET 中の 7952 語に付与され, さらにそれらの語の下位語として EMTREE コードに関連付けられた語もあわせると約 24,000 語あり, これらを MiniMALIMET という. MiniMALIMET は, 語数では MALIMET の優先語の約 1 割に過ぎないが, 1974 年から 1987 年までに 60 回以上使用された優先語を中心に選定されたものであり, 使用頻度では約 85% を占める.
 - この EMTREE コードにより, MiniMALIMET 中の語は上位下位などの語間の関係が規定されている.
 - 23) National Library of Medicine. Medical Subject Headings: Annotated Alphabetic List; 1990. Bethesda, National Library of Medicine, 1989, 922 p.
 - 24) 国立国会図書館. 日本科学技術関係逐次刊行物総覧 1988. 東京, 国立国会図書館, 1988
 - 25) JISCT 科学技術用語 シソーラス 英日対訳リスト 1987. 東京, 日本科学技術情報センター, 1987, 816 p.
JISCT 科学技術用語 シソーラス 日英対訳リスト 1987. 東京, 日本科学技術情報センター, 1987, 733 p.
 - 26) International Committee of Medical Journal Editors. Uniform requirements for manuscripts submitted to biomedical journals. Annals of Internal Medicine. Vol. 108, p. 258-265 (1988)
 - 27) Sackett, David L.; Haynes, R. Brian; Tugwell, Peter. Clinical Epidemiology: a Basic Science for Clinical Medicine. Boston, Little Brown, 1985, 370 p.