

## 情報検索システムへのキーワード自動抽出機能の適用方法について

## 2N-4

石黒正典 雉吉秀嗣 田中一敏  
NTT 情報通信研究所

## 1.はじめに

文献情報、新聞記事情報等を提供する情報検索サービスにおいては、原情報である“原文”からキーワードを抽出し、検索用のインデックスを作成する際に、

1. 原文を読みキーワードを抽出する、
  2. 原文そのものをデータベースに登録する、
  3. 登録した原文に対応するキーワード群を入力してインデックス作成をシステムに指示する、
- という過程が全て人手を介して行なわれており、多大な工数を要していた。我々は、人手に頼っていた作業を可能な限り自動化し、システム全体の効率化を実現した。即ち、キーワード自動抽出用プログラムをシステムに組み込み、キーワード抽出からインデックス作成までを一貫してシステム側で自動化するようにした。本稿では、この自動化の方式について述べる。

## 2. 情報検索サービスにおけるインデックス

文献情報、新聞記事情報データベースを用いてキーワード検索を行なう情報検索サービスでは、インデックスは“IVL (Inverted List)”と呼ばれる構造を有する(図1)。IVLの「タグ部」の“索引語”には原文より抽出されたキーワードが登録され、「サブ部」にはRID(原文を含むレコードのID)が登録される。

文献タイトル	著者	発表年	-----	キーワード <sup>1)</sup>				原文(長大データ) <sup>2)</sup>
				キーワード1	キーワード2	-----	キーワードn	
情報システム	A氏	1989	-----	データベース	コンピュータ	-----	情報処理	原文A
知能と知識	B氏	1988	-----	人工知能	AI	-----	知識処理	原文B
；	；	；	；	；	；	；	；	；
OAとデータ処理	Z氏	1983	-----	OA	帳票処理	-----	データベース	原文Z

\*1: 「キーワード」; 原文に対してキーワード抽出を行なった結果として得られたキーワードが格納されるカラムであり配列構造である。この配列の個々の値がIVLのタグ部の索引語となる。

\*2: 「原文」; キーワード抽出対象となる原文であり、長大データとして管理される。

図2. サービス用データベースのテーブル構成例

RID	キーワード抽出結果リスト	分析不能語リスト	同義語リスト
#003	コンピュータ、情報処理、人工知能、データベース、知識処理、OA	ハイオフロード、ニューロコンピュータ、ハイラックプロトコル、ハイゲーションモジュール、ハイパーテキスト	コンピュータ、情報処理、データベース、データ管理、ディレクトリ
；	；	；	；

図3. 抽出結果格納DBのテーブル構成例

- ・キーワード抽出結果リスト：キーワード抽出用辞書から抽出されたキーワード（名詞）の一覧が格納される。
- ・分析不能語リスト：原文から名詞として検出されたが、キーワード抽出用辞書に登録されておらずキーワードとは扱われなかつた名詞の一覧が格納される。
- ・同義語リスト：キーワード抽出用辞書から抽出されたキーワードに対して、同義語辞書から当該キーワードに対する同義語として抽出された名詞の一覧が格納される。

#### 4. 実現方式

システムの全体構成並びにレコードの追加～キーワード抽出～インデックス作成までの処理の流れを図4に示す。

キーワード抽出用辞書や同義語辞書においては、学術分野の追加、技術の発展等に伴って新語が逐次追加されなければならない。このため、分析不能語として扱われた名詞のうち必要な名詞（用語）をキーワード抽出用辞書／同義語辞書に適宜追加するユーティリティを用意した。これにより各辞書を最新の状態に維持でき、キーワードも最新のものを抽出させることができた。

#### 5. 評価・今後の予定

本方式により、従来入手により行ってきたキーワード抽出からインデックス作成までの作業をシス

テム側にて自動的に行うことを可能とし、作業の効率化を実現することができた。原文の入力に当たっては今後印刷漢字O C Rをシステムに組み込むことによりさらに人手の軽減を図る。

また、I V Lのようなインデックス構成を有するシステムの場合は、従来ロック期間が長くなること等によりデータベース更新は夜間のバッチ処理にて行われるのが一般的であったが、我々はロック期間を極小化させオンラインサービス中ににおけるデータベースの更新及びI V Lの更新を実現した<sup>3), 4)</sup>。これにより、キーワード抽出からインデックス作成までをオンラインサービス中に実施することも可能となる。これについては、今後検索処理への影響（応答時間等）の面から検証する予定である。

#### 〔参考文献〕

- 1) 村田、寺中他：マルチメディア・データベースに向けたデータモデルについて、データベースシステム研究会(43-1) 1984.9.17
- 2) 岸本、長浜他：マルチメディアDBMSにおける長大データ管理について、データベースシステム研究会(52-1) 1986.3.10
- 3) 坪井、川下：情報検索システムに於けるインバーテッドファイルのオンライン更新実現手法、情報処理学会第33回全国大会(1H-2)
- 4) 石黒、青木：情報検索システムにおける検索結果の矛盾回避手法、情報処理学会第33回全国大会(1H-3)

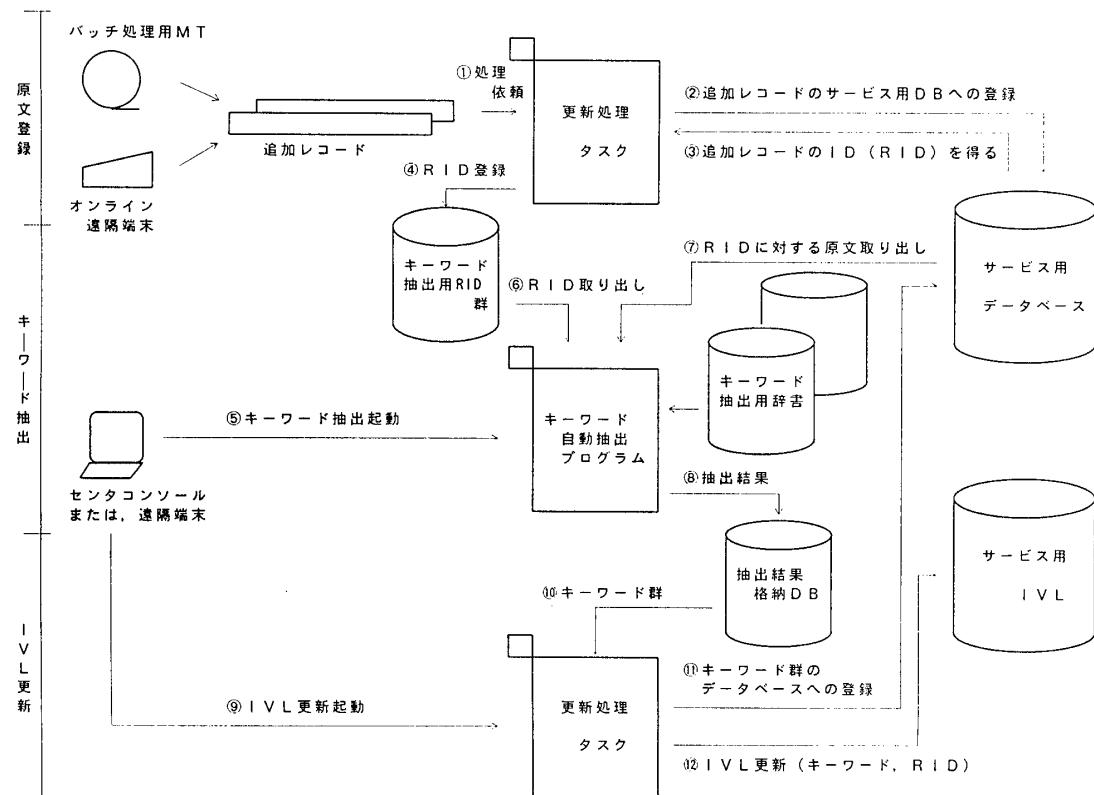


図4. システム構成／処理の流れ