

学術情報流通を実現する技術(1)

--要素技術(検索、DB等の基盤技術に
特化した話を中心に)--

Code4Lib JAPAN コアメンバー

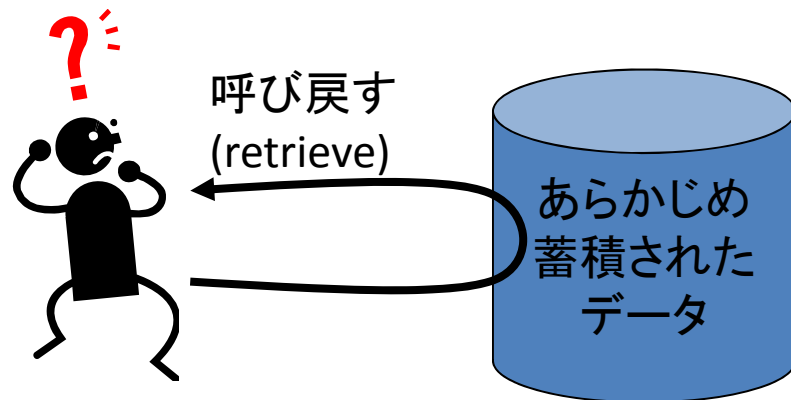
国立教育政策研究所 教育研究情報センター

研究員 江草由佳

twitter:@yegusa

情報検索とは

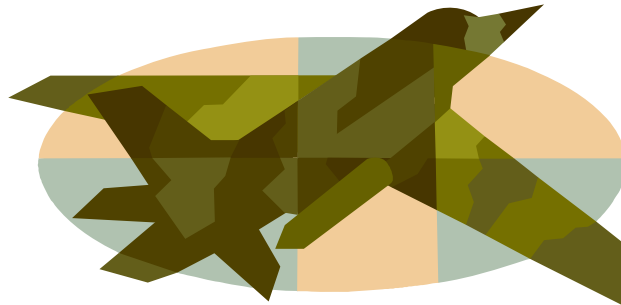
- 情報検索
 - IR: information (storage and) retrieval
 - 情報(information) を呼び戻すこと(retrieval)
 - 元は*i*nformation storage and *r*etrieval 情報の蓄積と検索
 - 1950年にムーアーズ(Calvin N. Mooers)が初めて定義
 - 1960年代に広く使われるようになる
 - (search: これも「検索」と訳すが。。。)



retriever(レトリバー):
獲物をくわえて戻って
くるように訓練された猟犬

データベースの起源

- 1950年代
- 米国国防総省が戦力に関する**情報を保管、集中管理**するためコンピュータを使ったライブラリーを開発
- **データの基地**(data base)から由来



データベースの定義(1)

- 著作権法二条十の三

- 論文、数値、図形その他の情報の集合物であって、それらの情報を電子計算機を用いて検索することができるように体系的に構成したもの

- 日本工業規格(JIS)

- 適用業務分野で使用するデータの集まりであって、データの特性とそれに対応する実態の間の関係とを記述した概念的な構造によって編成されたもの(X0017)

- 特定の規則に従って電子的な形式で、一か所に蓄積されたデータの集合であって、コンピュータでアクセス可能なもの(X0807)

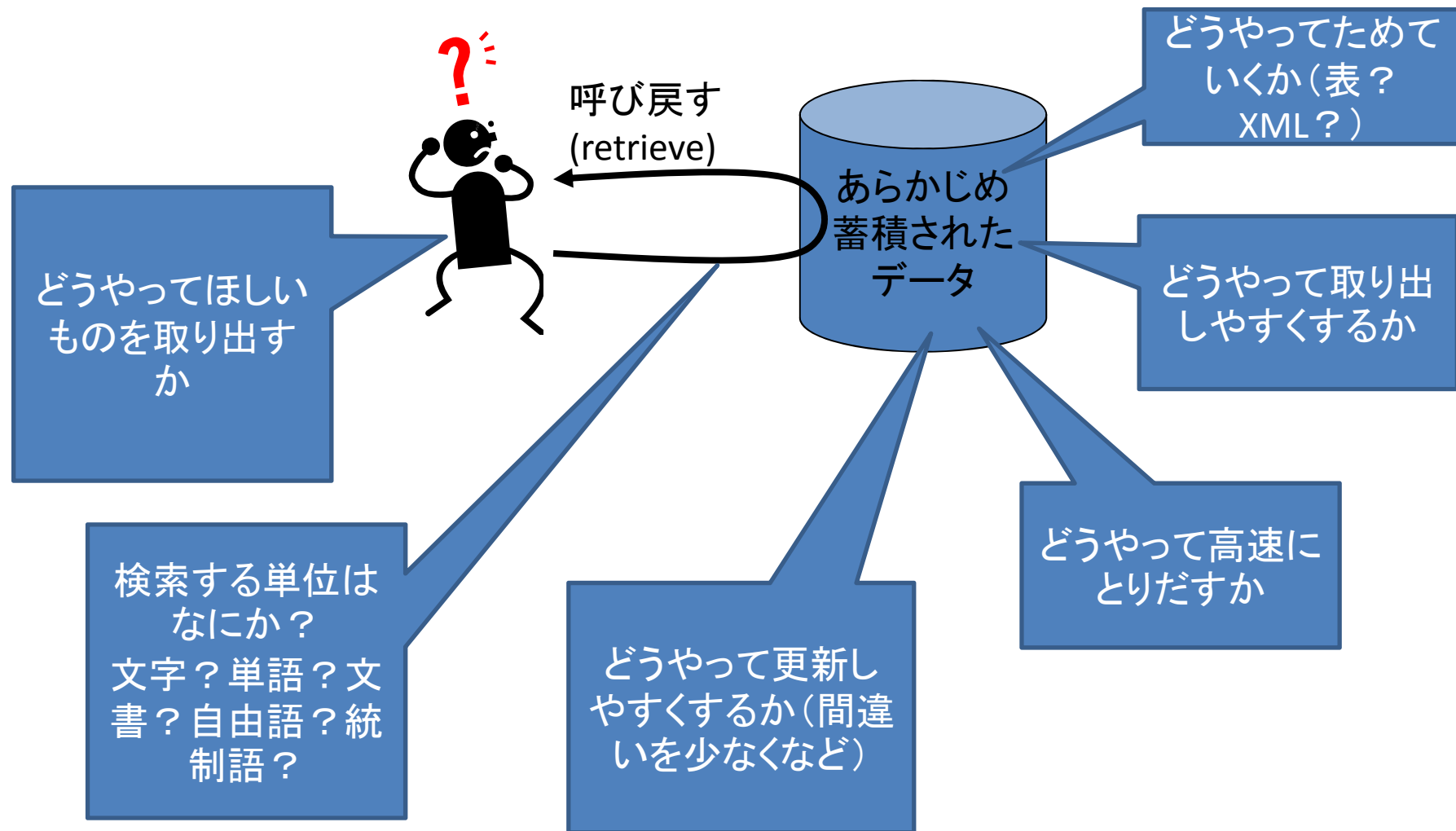
データベースの定義(2)

--日本のデータベースの特徴--

- データベースとは”コンピュータを用いて検索できる”ことが重要である。情報が電子メディアに蓄積され、コンピュータ、携帯情報端末(PDA)、地上波テレビ端末などを使用して検索できる状態になっている。
- データや情報がコンピュータ処理できるように体系的に整理され、統合化・構造化されて蓄積・保存されており、必要な情報だけを部分的に取り出せる。
- 蓄積情報の検索や更新が容易に行えるよう、効率化を図ったものである

一方、ヨーロッパにおけるデータベースの定義では、コンピュータを使用するかしないか、電子的であるかどうかについては特に限定していない

検索にまつわる様々な観点(一部)

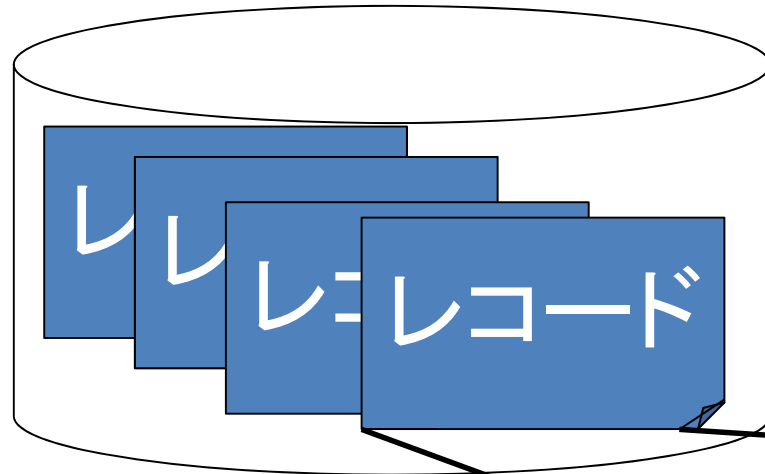


どうやってためていくか？表？XML？

-- データベースの種類 --

- リレーショナルデータベース (RDB)
 - 表としてデータを扱う
- オブジェクトデータベース
 - オブジェクトとしてデータを扱う
- XMLデータベース
 - XMLとしてデータを扱う
- 全文データベース
 - いろいろある、なんでもRDBにすればよいというものではない
 - 書誌データ、全文データはRDBには向いていない

どうやって取り出しやすくするか？ --レコードと検索フィールド--



検索フィールド名

検索フィールド値

検索フィールド

論題：

Reading—速読・多読
について考える

著者名：

清水由理子

請求記号：

P343-5C2-14

掲載誌名：

獨協大学外国語教育研究14

発行年月：

1995.12

掲載ページ：

p.273～282

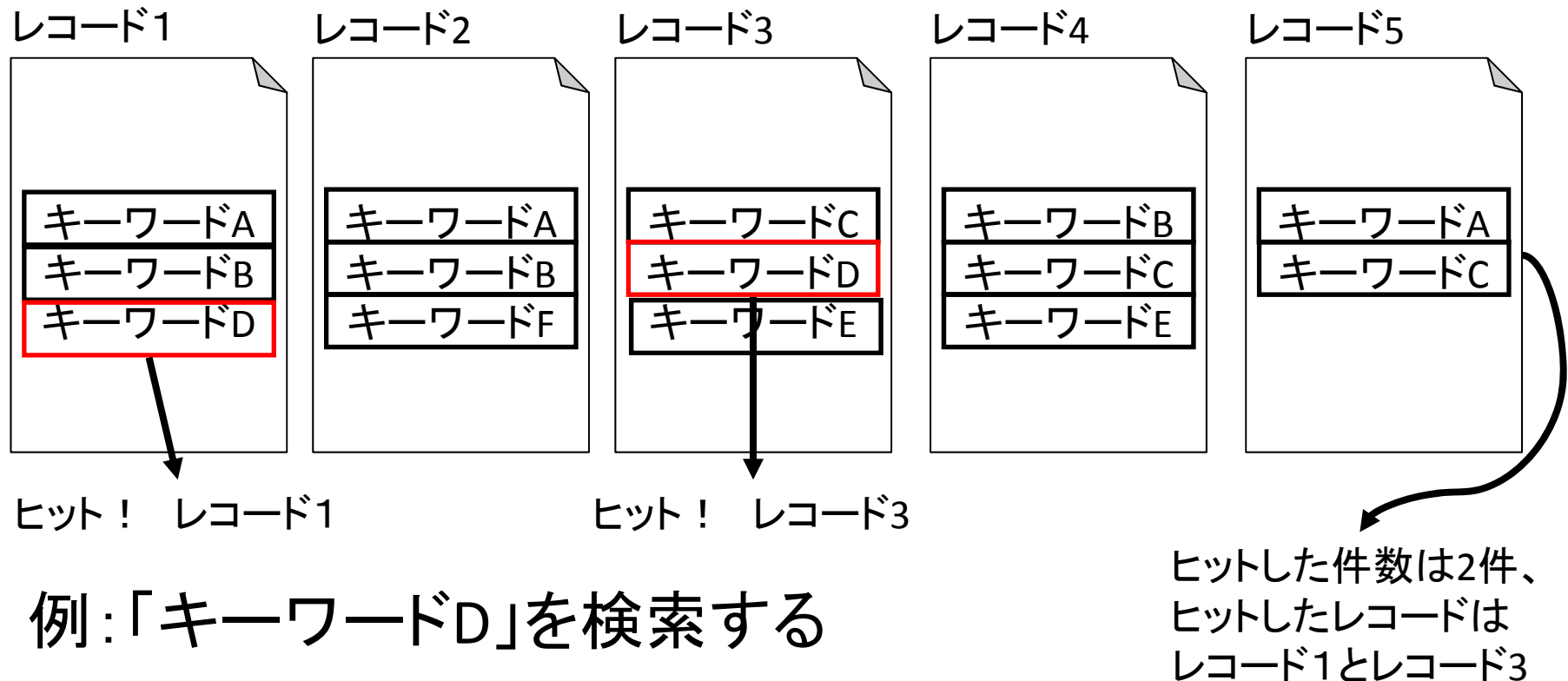
登録日：

19970930

どうやって高速にとりだすか(1)

--なんにも仕掛けがないと。。。--

- レコードを最初から最後まで順番に検索
- レコードが多くなると時間がかかってしまう



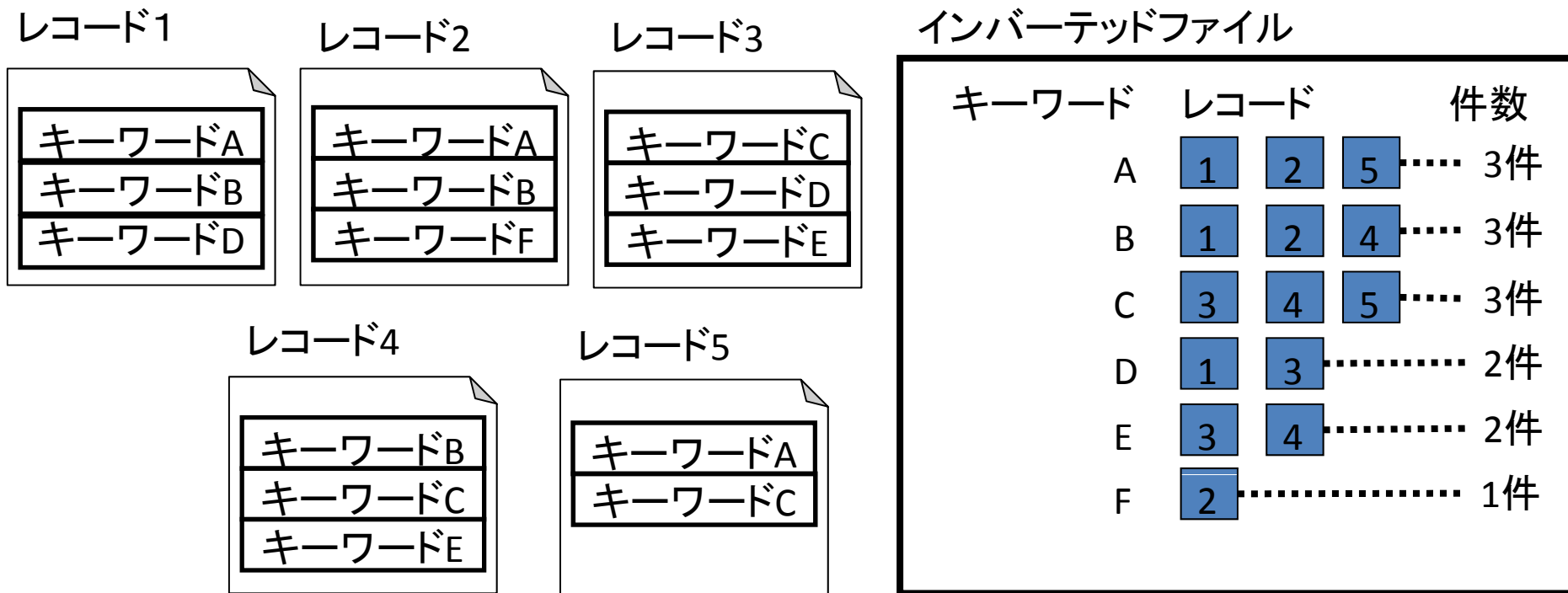
どうやって高速にとりだすか(2)

-- インデックス を用意する --

- 高速にデータにアクセスするために必要
- インデックス方式
 - インバーテッド(転置)ファイル
 - もっともよく使われる方式
 - Suffix Array
 - …いろいろある

インバーテッド・ファイル(1)

- 検索キーごとにレコードを集計したファイルを作成する方式
- インバーテッド・ファイル(転置ファイル、倒置ファイル)ともいう



インバーテッド・ファイル(2)

- 例:「キーワードD」を検索する
- 欠点:あらかじめインバーテッドファイルを作らなければならない、ファイルの容量増

レコード1

キーワードA
キーワードB
キーワードD

レコード2

キーワードA
キーワードB
キーワードF

レコード3

キーワードC
キーワードD
キーワードE

レコード4

キーワードB
キーワードC
キーワードE

レコード5

キーワードA
キーワードC

インバーテッドファイル

キーワード	レコード	件数
A	1 2 5	3件
B	1 2 4	3件
C	3 4 5	3件
D	1 3	2件
E	3 4	2件
F	2	1件

インバーテッドファイル(3)

--どうやってキーワードを取り出す？--

例: 「<http://nyti.ms/po6J1Z> うーむ、Aaron Swartzが昨年9月のJSTORからの大量ダウンロード容疑で逮捕・起訴とは、絶句。
#librahack」

- <http://nyti.ms/po6J1Z>
- うーむ
- Aaron
- Swa
- が
- 昨年
- 9
- 月
- の
- JSTOR
- から
- の
- 大量
- 逮捕
- ・
- 起訴
- とは
- 絶句
- #librahack

意味のある語単位で切る＝形態素解析
(日本語の文章の意味(かかり受けや品詞など)を解析して
意味のある語で区切るようにする)

インバーテッドファイル(3)

--どうやってキーワードを取り出す？--

例: 「<http://nyti.ms/po6J1Z> うーむ、Aaron Swartzが昨年9月のJSTORからの大量ダウンロード容疑で逮捕・起訴とは、絶句。
#librahack」

- <http://nyti.ms/po6J1Z>
- うー
- ーむ
- Aaron
- Swartz
- が昨
- 年9
- 9月
- 月の
- JSTOR
- から
- らの
- の大
- ウン
- ンロ
- ロー
- ード
- ド容
- 容疑
- 捕・
- ・起
- 起訴
- 訴と
- とは
- は絶
- 絶句

N-gram: 強引に、2文字ずつなどで強制的に切る方法

インバーテッドファイル(4)

--どうやって検索する?--

- 例
 - 「<http://nyti.ms/po6J1Z> うーむ、Aaron Swartzが昨年9月のJSTORからの大量ダウンロード容疑で逮捕・起訴とは、絶句。#librahack と」
- Q: [Aaron] -> Hit!
- Q: [9月] -> Hit?
- Q: [大量ダウンロード] -> Hit?
- Q: [#librahack] -> Hit???
- Q: [アーロン・シュワルツ] -> No Hit?

インデックスを使った
検索(仕組み/人)
のところで、
それぞれ工夫必要

まとめ

- 情報検索～データベースとは？
- 検索にまつわる観点
- レコードの蓄積からインデックスの構築、検索まで
- インバーテッドファイルを例に
- 参考文献：
 - Modern Information Retrieval: The Concepts and Technology behind Search (2nd Edition) (ACM Press Books), ISBN: 978-0321416919, 2011.02, p.944