

**e-Society**  
 文部科学省リーディングプロジェクト 基盤ソフトウェアの総合開発

## 研究のご紹介

東京大学 生産技術研究所 教授  
 文部科学官  
 情報処理学会副会長  
 文部科学省特定領域研究「情報爆発」領域代表  
**喜連川優**

特定領域研究  
情報爆発

### MySpace

2006年各国人口

国際連合統計部による各掲載年の7月1日現在の推計人口

ワールドビジネスサテライト 2007. 7. 25

**e-Society**  
文部科学省リーディングプロジェクト

## Socio-Sense

### =社会のセンサーとしてのWEB=

研究代表者  
 喜連川 優  
 (東京大学生産技術研究所戦略情報融合国際研究センター)

**e-Society**  
文部科学省リーディングプロジェクト

### 実社会の射影としてのウェブ ウェブは社会のセンサー

Webは実世界の転写構造を形成

目的:ウェブ情報の高度利用システムの構築(WEBの時空間解析)

**e-Society**  
文部科学省リーディングプロジェクト

ウェブ情報の空間的・時間的分析による高度利用を実現

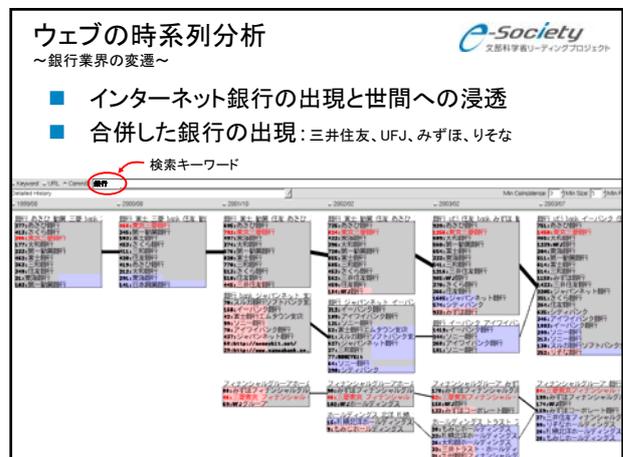
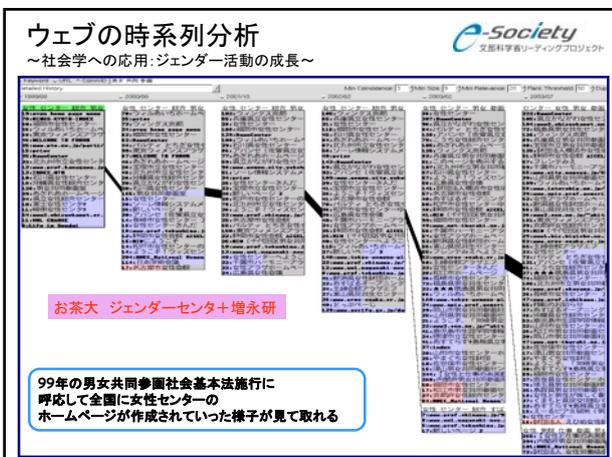
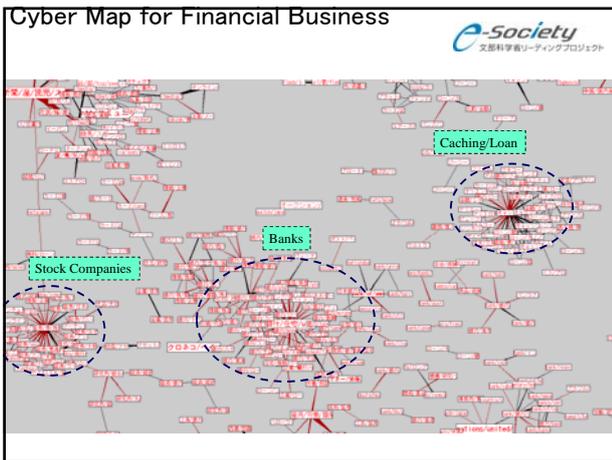
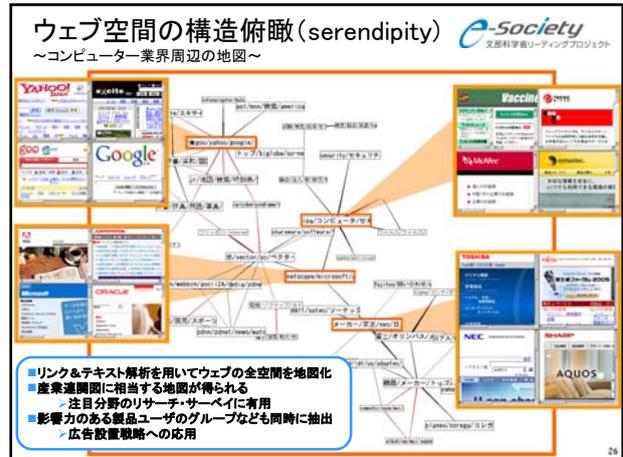
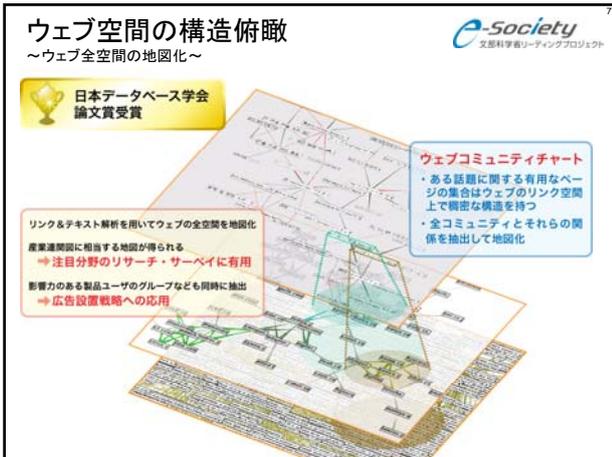
**e-Society**  
文部科学省リーディングプロジェクト

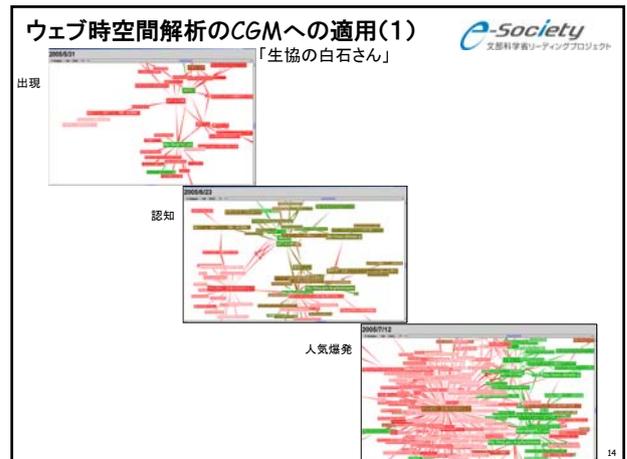
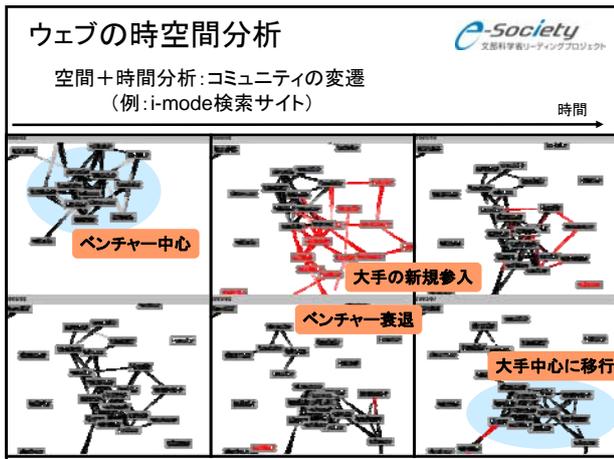
### 日本語ウェブアーカイブの構築

- 9年間にわたり100億ページ規模の日本語ウェブページを集積し、継続期間および規模において**アジア圏最大級**のウェブアーカイブを構築
- 各URLの更新頻度に応じた収集技術を開発し、1日~1年の**可変周期収集**を実現

累積ページ履歴数

ページ更新回数



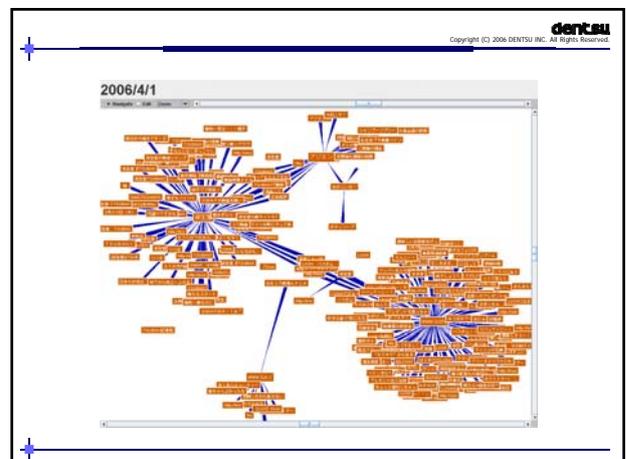
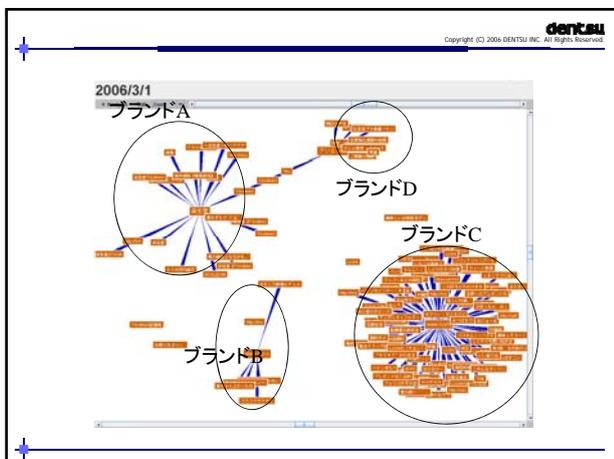
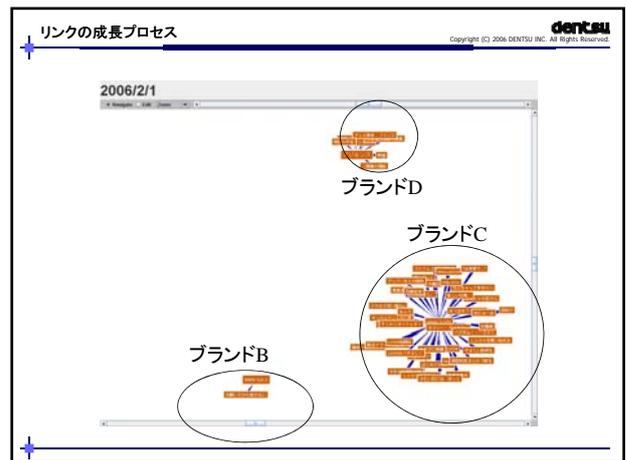


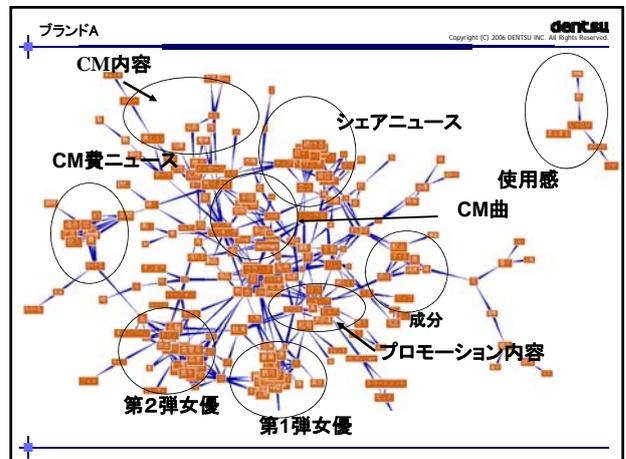
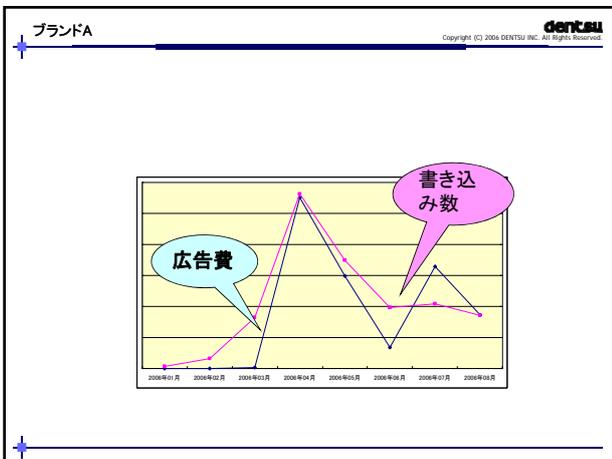
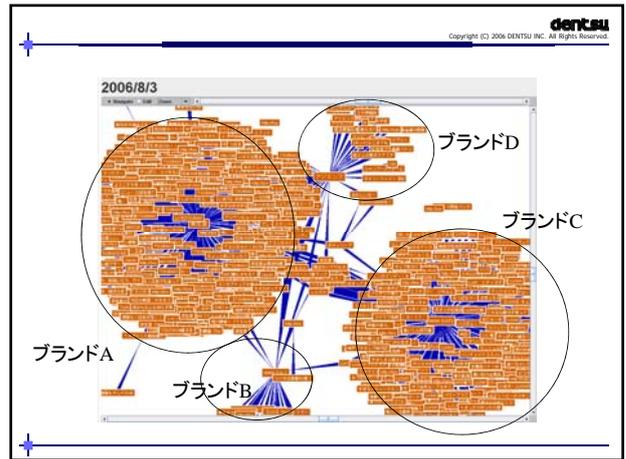
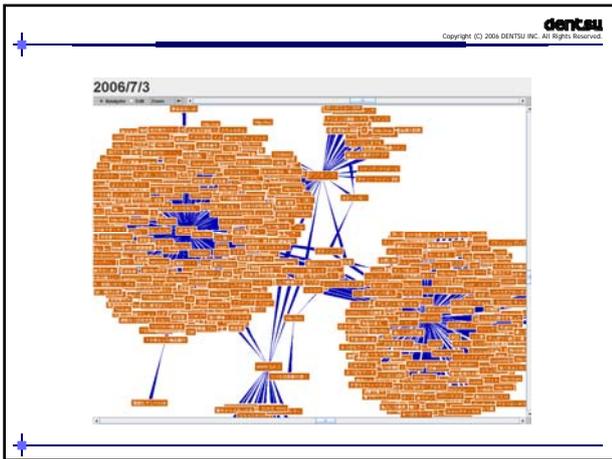
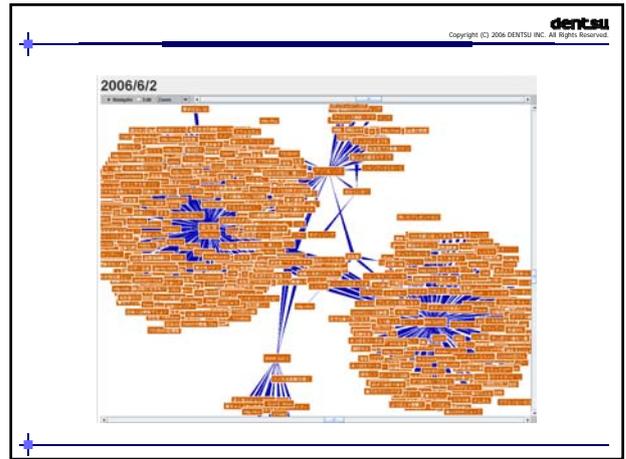
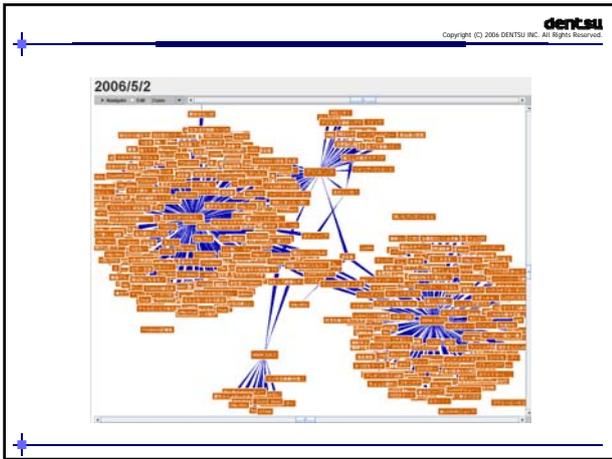
### Definition of Novelty Measure

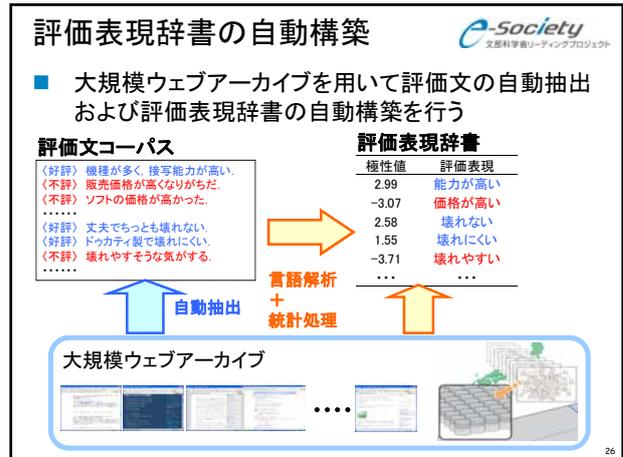
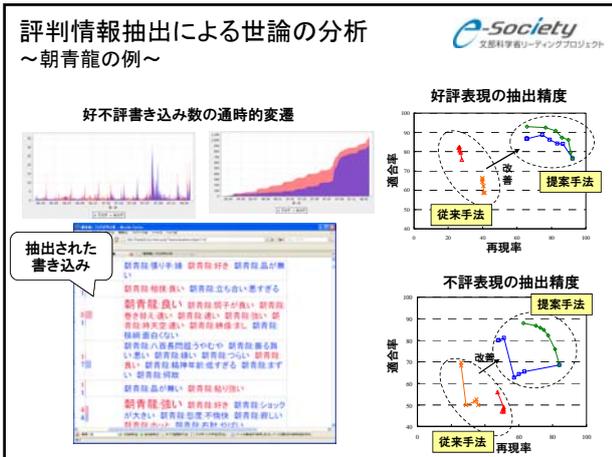
$$N(p) = (1 - \delta) \frac{\sum_{(q,p) \in I(p)} n(q,p)}{|I(p)|}$$

$$n(q,p) = \begin{cases} 1 & q \in L_2(t_k) \\ 0 & q \in O(t_k) \setminus L_2(t_k) \\ N(q) & q \in U(t_k) \end{cases} \quad (1)$$

- $\delta$ : damping factor  
– probability that there were links to  $p$  before  $t-1$







### WEBからの未知語(カタカナ用言辞書)の獲得

五段ラ行	一段	五段カ行	形容詞	形容動詞
ググる	モエる	トキメク	チャチい	サイコーだ
バモる	トロケる	ホザク	イナたい	オトコマエだ
ハラシマる	イケてる	ドツク	ヤヴァい	フルーティだ
ヒキコモる	シャガれる	シバク	カバエイ	エイジレスだ
バザる	コジャれる	ムカツク	スゲい	メタメタだ
ストロベリる	ハツチャケる	イタダク	マンドクさい	フニャフニャだ
ガンがる	デケる	ワメク	ヘコい	レビッシュだ

動詞 5段 ラ行  
はらしま・る  
ハラシマる

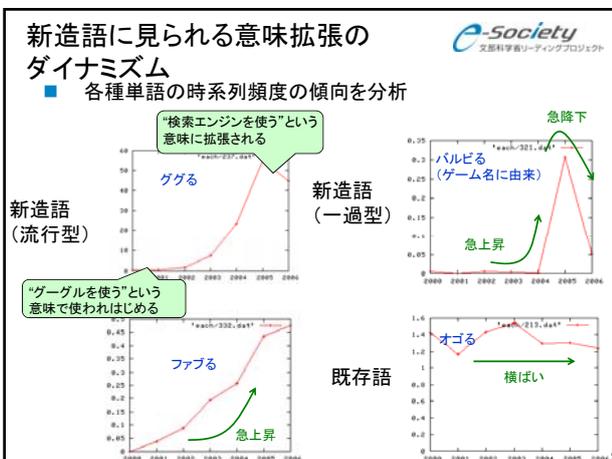
動詞 5段 ラ行  
ばもる  
バモる

動詞 5段 ラ行  
もふ・る  
モふる

### 「はらしまる」について

原稿  
縞

・ 業界でしか使わなかった言葉が発信されるようになりつつある。



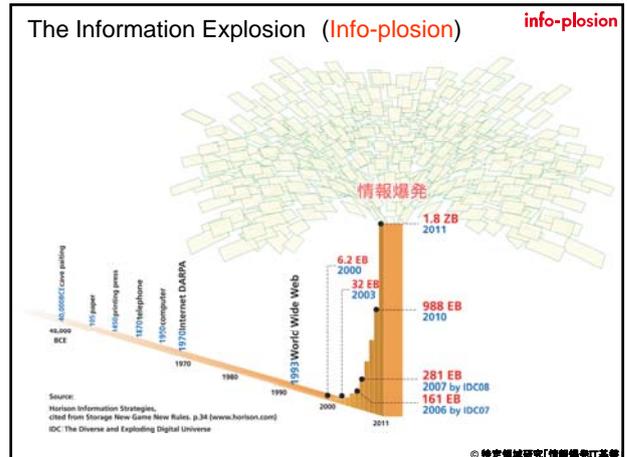
New IT Infrastructure for the Information Explosion Era  
Research Project funded by AMET Grant-in-Aid for Scientific Research in Priority Areas (FY 2005-2010)

# info-plosion

情報爆発時代に向けた新しいIT基盤技術の研究

平成17年発足 文部科学省特定領域研究

喜連川優  
東京大学 生産技術研究所



## ディープNLPオープンサーチエンジン基盤 TSUBAKI (黒橋、新里)

(皆でいじれるコンポーザブル型爆発サーチエンジン)

- 日本語ウェブ文書5,000万件を検索対象とした開放型検索エンジン基盤
  - 高度ウェブ処理用標準フォーマットによりウェブ文書を管理
  - 構造的言語処理によるインデックス
  - 無制限に利用可能なAPI
  - 透明性・再現性のある検索結果
- 計算機環境\*
  - 計算ノード32台 (計64CPU)
    - CPU: 3.60GHz × 2
    - メモリ: 2GB
    - 内蔵ディスク: 600GB (計20TB)
    - ファイルサーバー2台 (計12.5TB)

\*今年度末に、64CPU、100TBストレージ、メモリ4GBを追加・増強予定

## ディープNLPオープンサーチエンジン基盤 TSUBAKI

info-plosion 情報爆発

インドの経済発展の障害

日印経済

インド

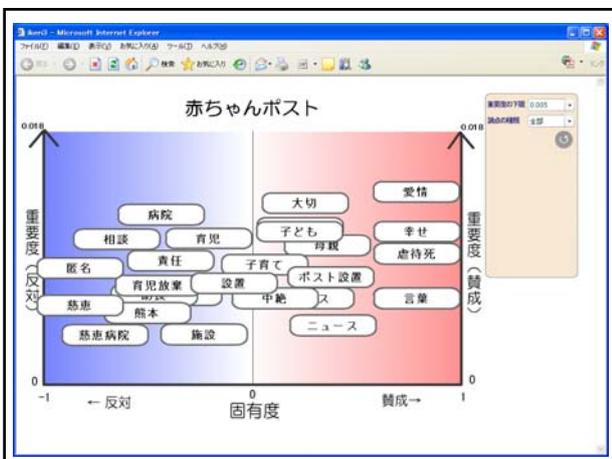
障害者の教育

人類の発展

経済発展の妨げ

障害

投資最大の障害



赤ちゃんポスト

反対 (Opposition) | 賛成 (Support)

子育ての代表的な意見

「子育て」に関する賛成意見 (右) と反対意見 (左)

## Webアーカイブからの辞書、 知識ベース自動作成

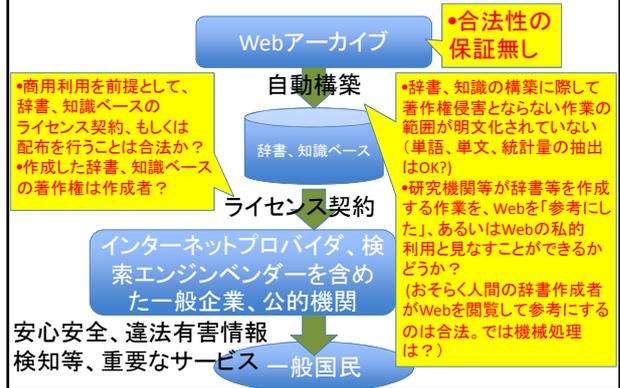
## 社会における辞書、知識ベース

- 機械学習等技術の進歩により、Webアーカイブをもとに大規模で有用な辞書や知識ベースを自動構築することは現在十分可能
- 辞書/知識ベースとその社会的意義の例
  - 違法有害情報検知で有用な辞書
    - アダルトビデオ女優の名前のリスト
    - 猥褻な日本語表現、犯罪に関係の深い日本語表現
  - 安心安全の担保で有用な知識ベース
    - 食品と有害物質の関係(ぎょうざと残留農薬、アジサイとシアン化合物による中毒)
    - 製品とその欠陥(DHAサプリメントと副作用の出血、ガスコンロとその欠陥)等々
- これらの辞書、知識ベースを、Web上で広く利用可能とすることで、**安心安全の確保、違法有害情報フリーなWeb**などを国民が享受することができる
- これら、Webから自動的に作られた辞書、知識ベースは現在合法的か？

## 辞書作成技術の例

- 「中毒」がトラブルであることを認識するには。。
  - Webアーカイブ上で「中毒」の周辺に現れる単語を収集
    - ...中毒で死亡する...
    - ...中毒で入院する...
    - ...中毒で体調を崩す...
  - そのパターンを統計的手法により自動的に学習
    - 学習された規則: 「Xで死亡する」、「Xで入院する」等のパターンに現れる単語Xはトラブルの可能性が高い
    - この規則により、
      - 癌、肝臓がん、毒、ガス中毒、交通事故、人身事故、。。。
  - おそらく、この種の表現は数万語単位で存在し、人手での収集は現実的ではない
  - また、新規に問題になるトラブルも継続的に収集する必要がある⇒自動化は必須

## 現状の法的問題点



## 今後のニーズ、技術の進化を 踏まえた法整備の必要性

- 少なくとも、現状、以下のようなニーズ、技術の可能性がある
  - 辞書、知識ベース作成によるイノベーション支援
    - 例: ISP細胞開発において有効利用可能な物質、遺伝子をWeb上に公開されたデータから自動的に列挙して知識ベースとし、体系的に実験
      - 現在はおそらく研究者がWebページを手手で閲覧して決定⇒膨大な手間
      - あらかじめ知識ベースをWebから作成しておくことで研究は加速
  - 機械翻訳の品質向上のためのWebデータ利用
    - Webをもとに膨大な対訳文を抽出、あるいは合成し、機械翻訳システムの統計的パラメータの計算で利用
    - 対訳データ抽出や合成が著作権侵害とならないか？
    - こうした利用が可能であるならば、日本文化、知的財産の発信において有用
  - 重要、かつ大量の情報を素人に分かりやすく説明するため、Web情報を一部自動改変、あるいは要約
    - おそらく現行法では非常にグレー、もしくは非合法
    - 問題は一種のデジタルデパイドであり、社会の効率、安全性を高める上で重要
- 他にも様々な技術のニーズとシーズが存在

## ニュース映像 Tivolution (井出 名大)

