

違法・有害情報を検出するための 技術開発について

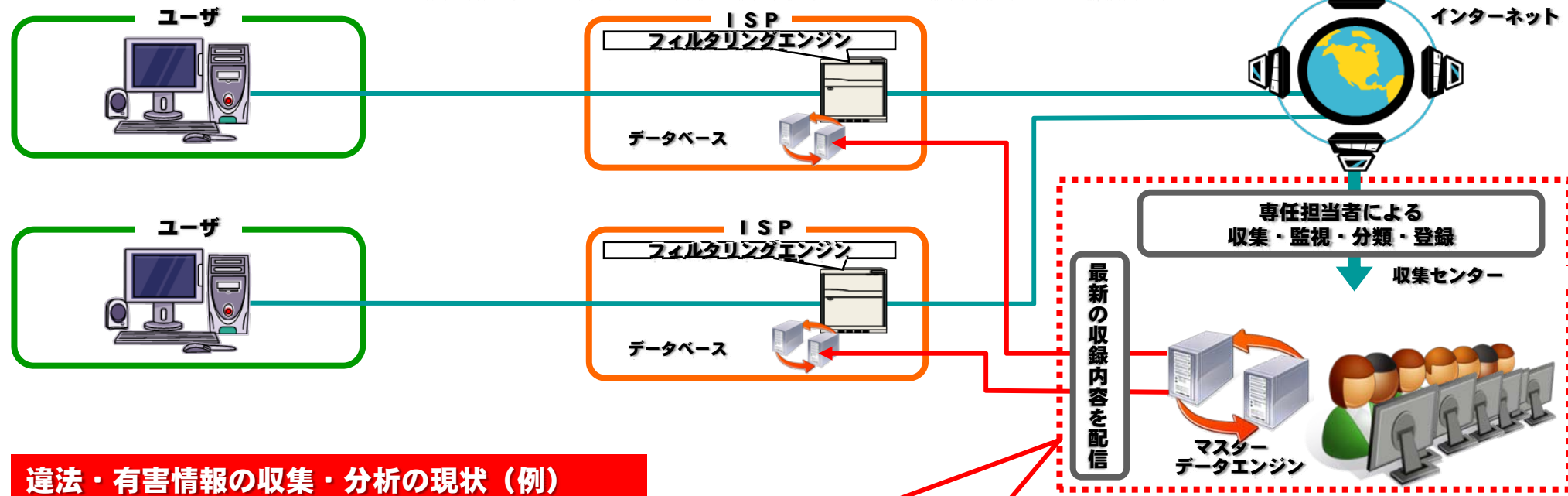
独立行政法人 情報通信研究機構

木俵 豊

2008年4月2日

違法・有害情報検出の現状と課題

インターネット



違法・有害情報の収集・分析の現状（例）



違法・有害情報の削除等への対応の迅速化・効率化を図るためには以下の課題の解決が必要。

- ✓ 違法・有害情報の検出の迅速化
- ✓ 違法・有害情報の検出負担の軽減
- ✓ ネット上の情報量の急増に対する対応の限界の克服



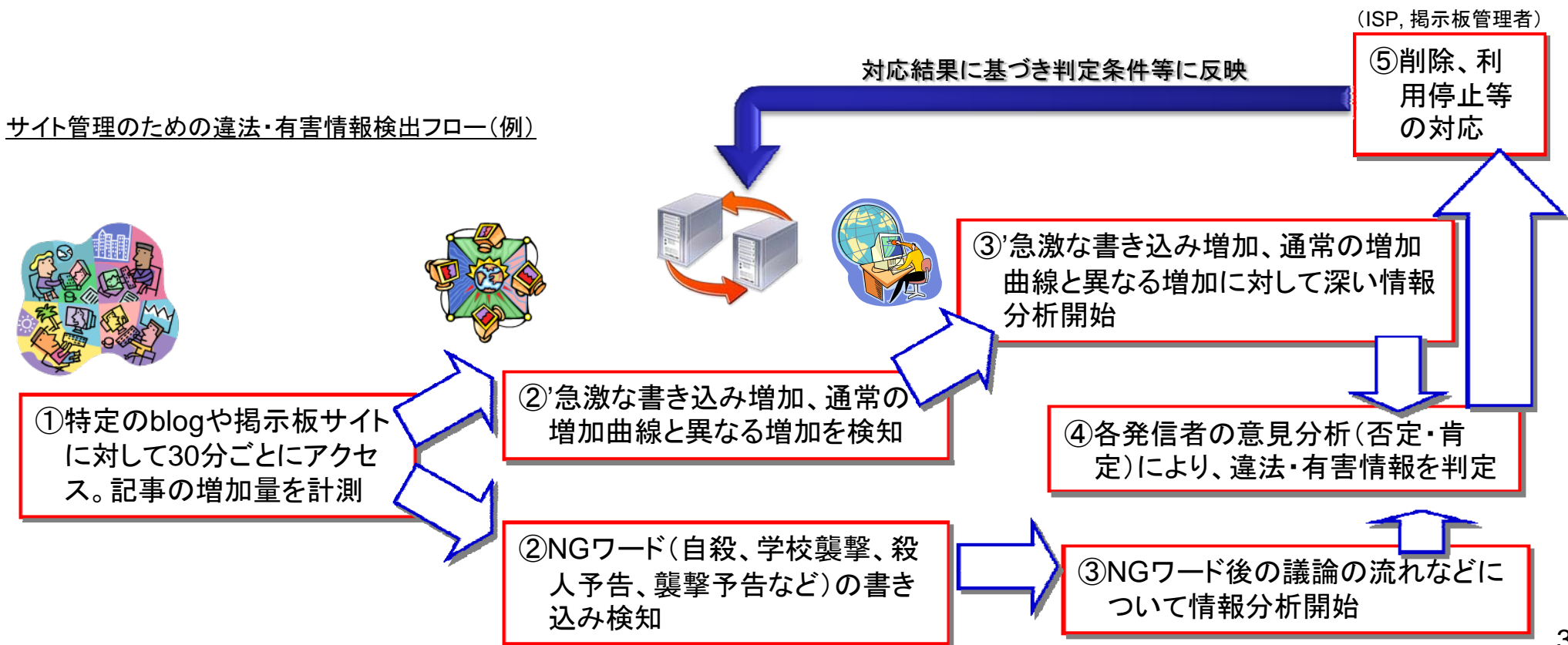
違法・有害情報対策に資する
研究開発を支援

インターネット上の違法・有害情報の検出を行うための技術開発 NICT

時々刻々新たに流通するコンテンツを監視するには、効果的なコンテンツ・チェック方式の高度化が不可欠。現在は、単語レベルで一一致したものを検出するのみで意味解析がされないため、不要なサイトまで検出する場合や違法・有害な情報が検出対象から除外される場合がある。

このため、自然言語技術を活用してWebコンテンツの構造分析を行い、違法・有害情報の検出精度を向上させるための研究開発に、インターネットサービスプロバイダー、コンテンツプロバイダー、関係団体、情報通信研究機構等が協力して取り組むことが必要。

これにより、違法・有害情報の検出を迅速に行うとともに、検出の負担を軽減し、もってネット上の情報の適正化を推進することが可能。



違法・有害情報の検出の現状と将来像

サイトから違法・有害
情報を出さない

プロバイダによる
フィルタリング

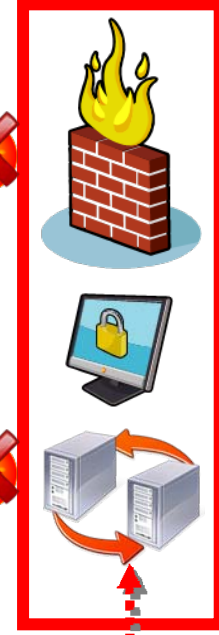
サイト運用者



インターネット



フィルタリング



ネット情報収集

違法・有害情報を
主に人手により検出

違法・有害情報の
フィルタリング情
報に反映

情報分析技術を活用し
検出の効率化を支援
・迅速な検出
・検出負担の軽減

当初はこちらで
の利用をター
ゲットとして研究

違法・有害情報の
所在情報を提供
又は自らサイト内
をチェック



自然言語技術を活用するメリット

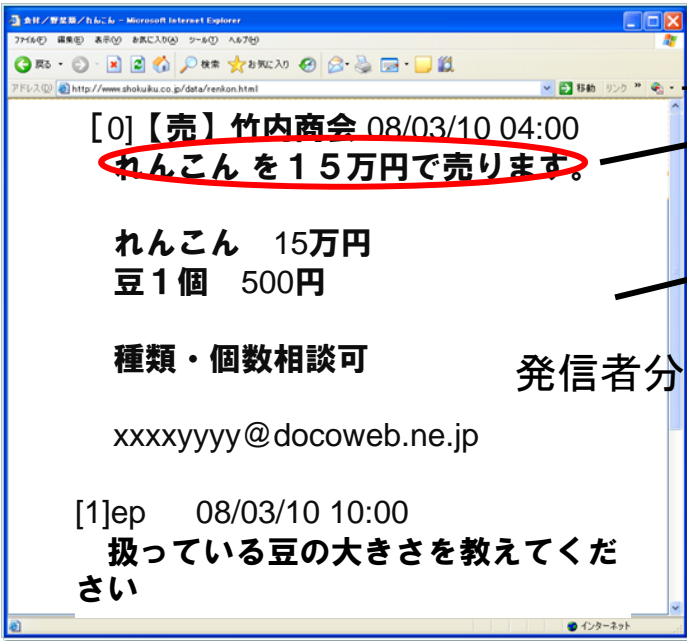
課題: 単語「れんこん」だけでは、「食品」と「拳銃の隠語」とを区別ができない。

自然言語技術を活用して、「誰が何をどうする」を抽出することで、違法・有害ページの候補を迅速に高精度で判定が可能。

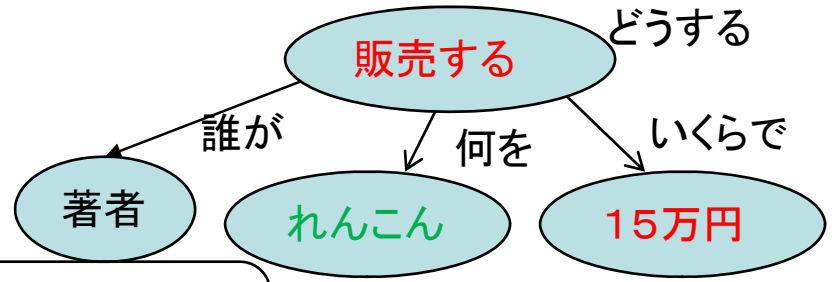


- 言語解析: 著者が**事実を陳述**
「れんこんは**ビタミンが豊富である**」
- 発信者分析 / ページ外観分析
発信者: **株式会社 日本食品薬化**
掲載場所: **企業サイト**
連絡先: **住所 + 電話番号あり**

発信者が企業で、有害表現パターンを含まない事実・意見のページであることを認識



言語解析



発信者: 竹内商会: **匿名(実世界の企業と認定できない)**
掲載場所: 掲示板
連絡先: 電子メールのみ

発信者分析 / ページ外観分析

知識データベースと照合することで

- **匿名(実世界での連絡先がない)の人が高額で販売していること**
- **れんこん**は隠語で拳銃を意味することから違法である可能性があることを認識



自然言語技術を活用するメリット

人手で作成・チェックした知識(有害表現パターン、有害語)を元に、自然言語処理技術を用いて大規模Webデータから知識を自動拡張することで、知識データベースの半自動構築が可能

□知識データベース： 有害表現パターンのデータベースの例
「誰が」「何を」「いくらで」「どうする(販売する)」のパターンと有害度の関係を記述

誰が	何を	いくらで	どうする	有害度
匿名 (住所なし)	(拳銃) (覚せい剤) (口座)	高額[1万以上]	販売する	違法の可能性が 極めて高い
匿名	*	高額[1万以上]	販売する	18才以下制限
匿名	*	*	販売する	15才以下制限
企業(住所あり)	(書籍)	*	販売する	12才以下制限
...

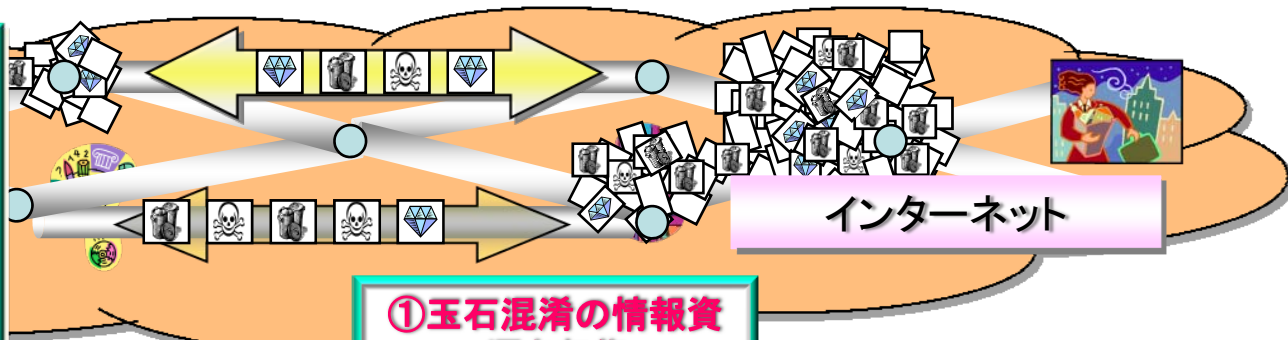
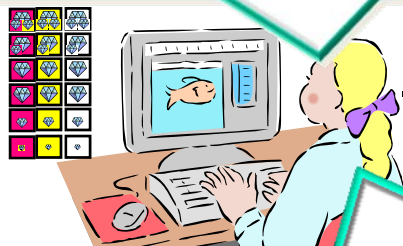
□知識データベース： 有害語の同意語辞書の例

意味	同意語(隠語)				
拳銃	ちゃか	れんこん	レーニン	北京ダック
覚せい剤	シャブ	S	スピード	クリスタル
...					

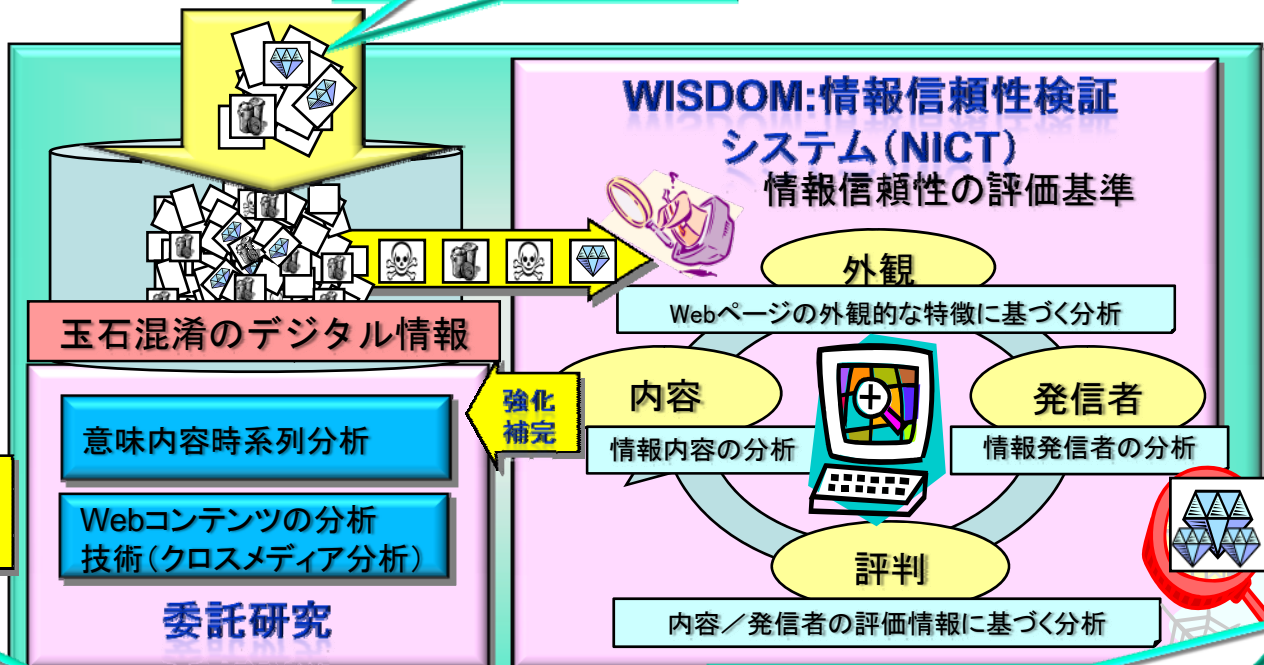
情報分析技術の研究開発

【利用例】アガリスクって本当に健康に良いのかしら？

★WISDOMを使って検索エンジンの検索リストをみると、リスト上位は販売会社ばかりよね。公的機関の内容を見る限りでは健康によいのかどうか疑問のようね。特に最近のWebページの評判は、販売サイトを覗くと90%以上否定的な意見だし、3つめの内容には、意見の矛盾があるみたい、4つめの効果を特に表現しているページの写真は、あるWebページの使い回しみたいだ。これらは信用できないなあ。



①玉石混淆の情報資源を収集

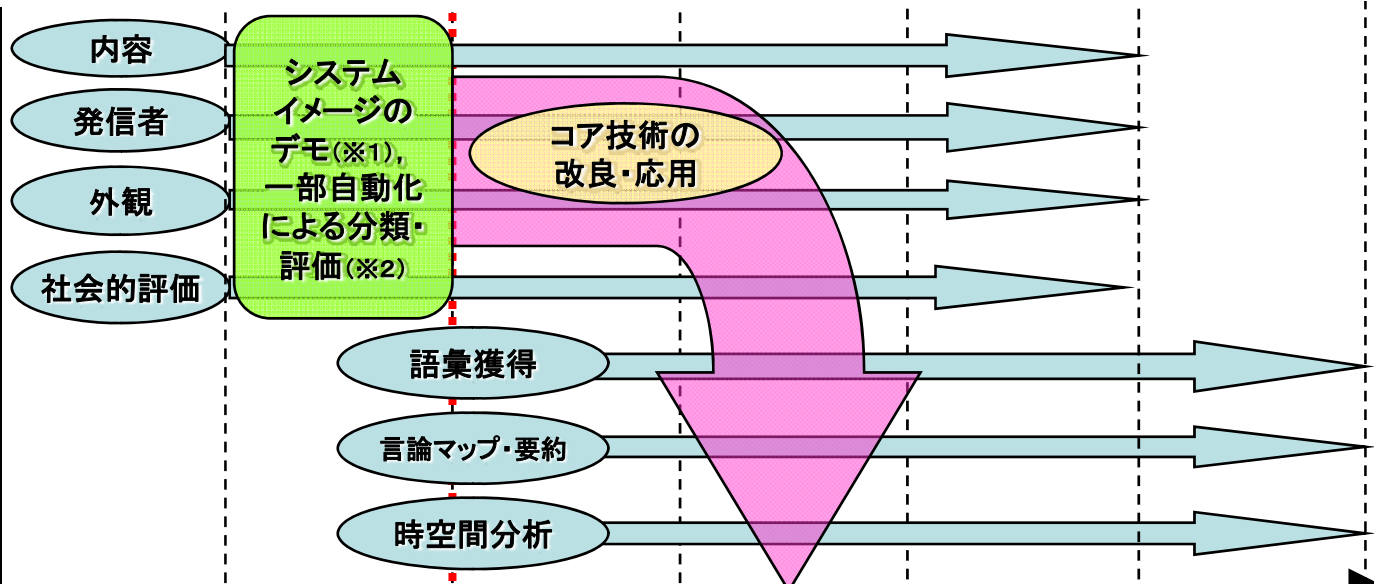


③分析結果を提示、ユーザにインターネット上の価値のある知識情報を発見させる。

②玉石混淆の情報資源を分析

情報分析研究の現状と今後

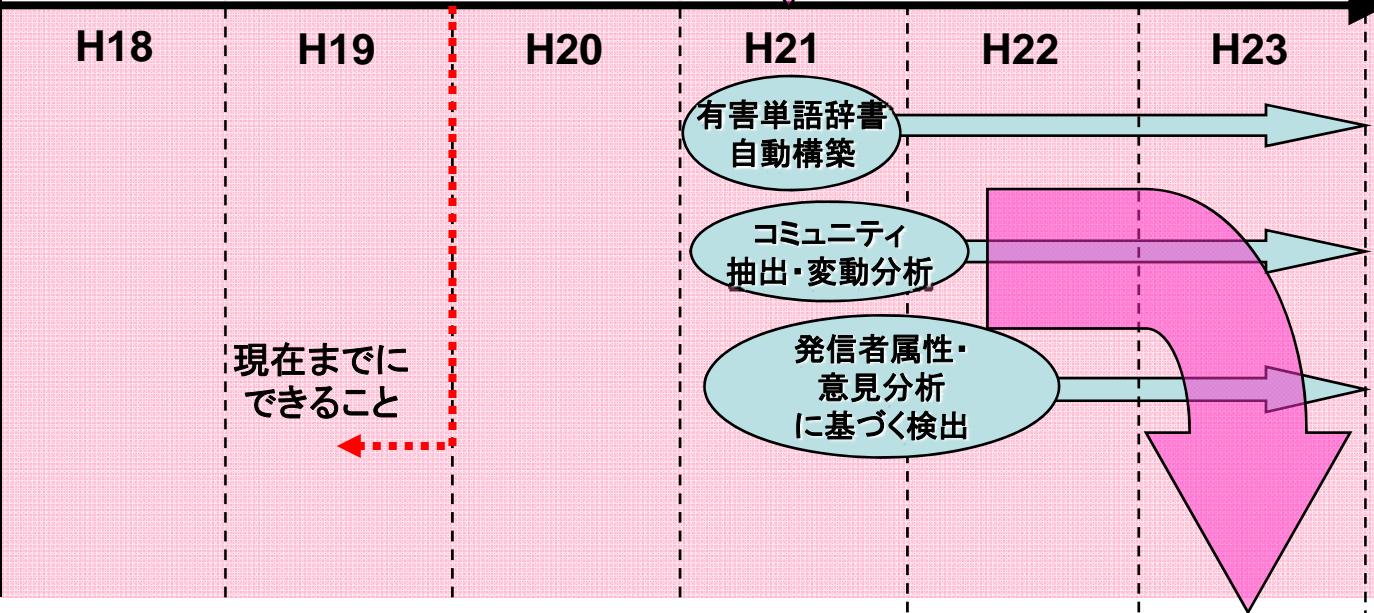
情報の信頼性評価



日本語頁
に対する
タイプ別
信頼性推定
(評判重視
型, 内容
重視型等)
(※3)

ネットワーク上
の文字、音声、
映像情報につ
いて、偽りの情
報、信頼性の
低い情報等を
分析する技術
を確立し、信頼
できる情報を
提供

違法・有害情報検出



成果
多様な有害
情報の生起を
迅速・高精度
に把握

違法・有害情
報の早期検出
や検出負担の
軽減

研究成果は順次提供

※1: 日本語2千頁規模, 基盤データ(人手)を使用
 ※2: 日本語1億頁規模, 発信者・社会的評価の一部自動分析を実装
 ※3: 日本語5億頁規模, 内容・発信者・外観・社会的評価の各観点での自動分析を実装