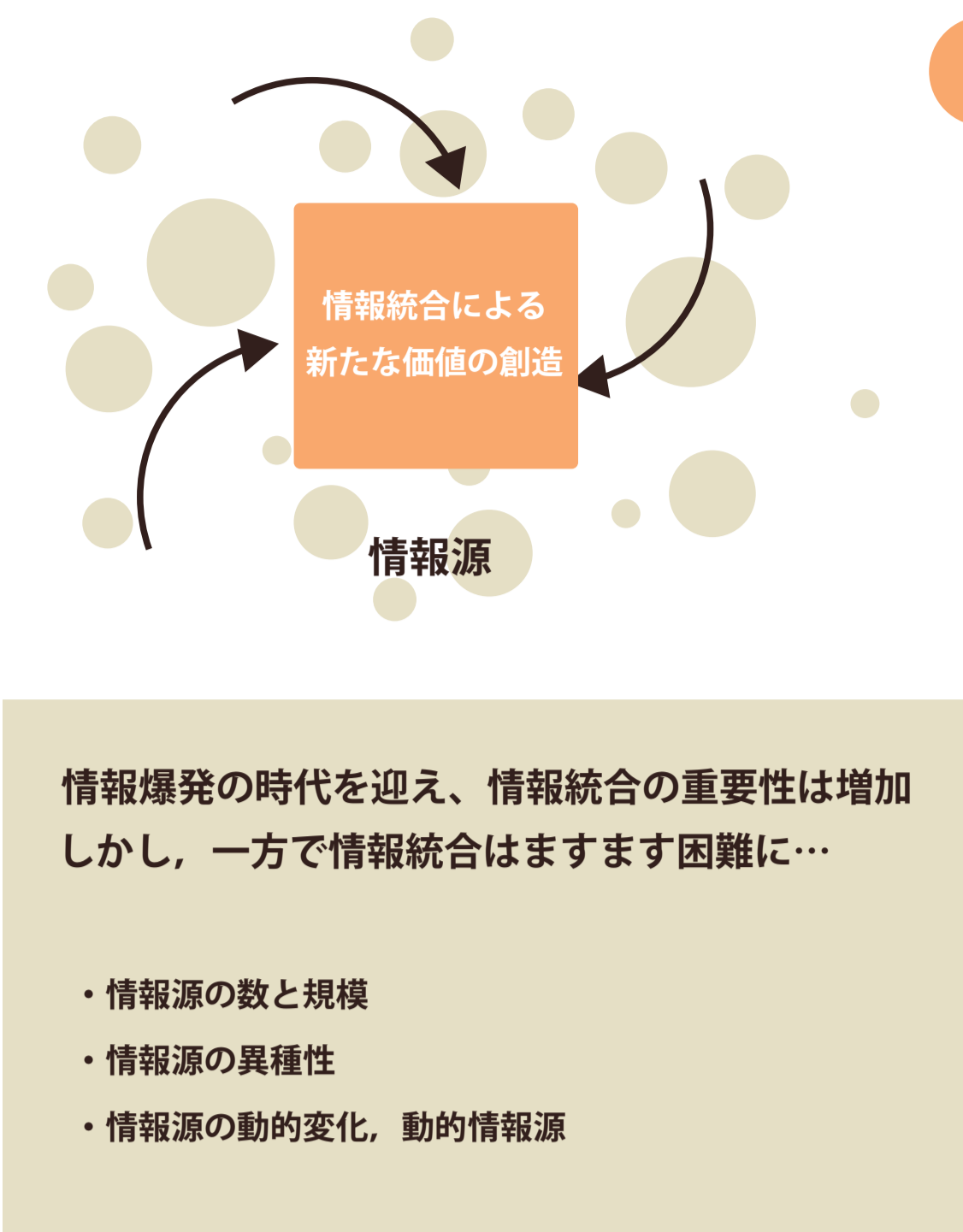


A01-03 能動的リソースマイニングに基づく異種情報統合基盤の研究

研究代表者：北川博之（筑波大学） 分担者：天笠俊之，森嶋厚行，川島英之（筑波大学）



下記の3つの視点より研究を推進

マイニングと情報統合に関する応用研究

- ソーシャルブックマークを用いた Web ランキング
- 類似検索に基づく異種 XML 統合
- 連続的モニタリングによる Web 一貫性管理

マイニングのための要素技術に関する研究

- トランザクションデータに対する外れ値検出
- ストリームに対する外れ値検出
- XML データに対する OLAP
- 時系列文書クラスタリング
- 移動体統計情報抽出

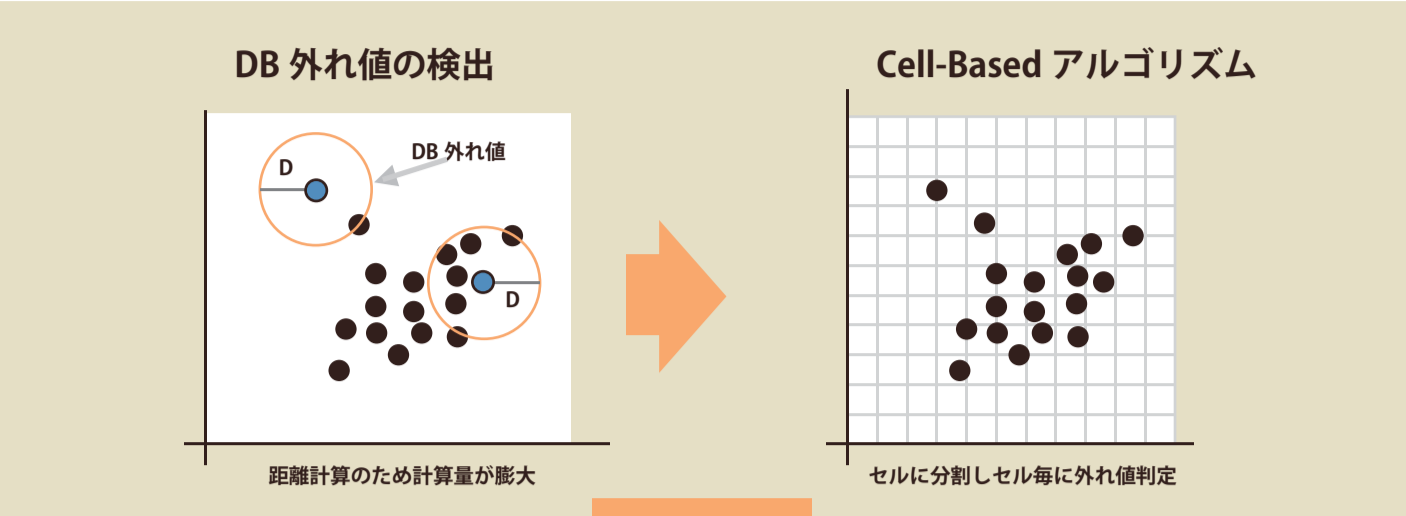
情報統合基盤システム

- 情報源の動的選択機能を有する情報統合基盤システム
- 多次元ストリーム用高性能索引機構

ストリームに対する外れ値検出

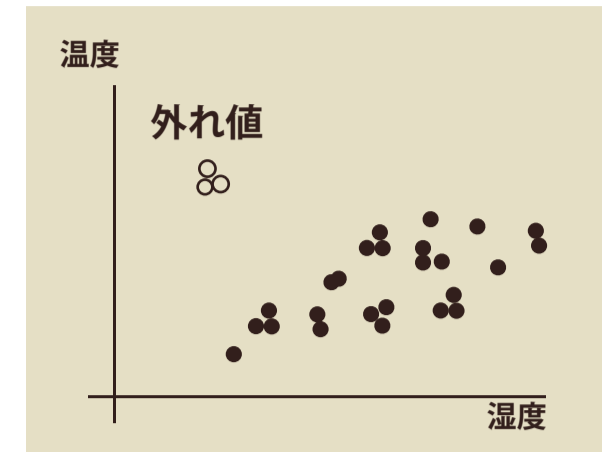
DEWS'08 優秀論文賞
DEXA'08 採録

- 目的：データストリームの特徴を活かして効率的に外れ値の検出を行う
- 手法：DB 外れ値を検出する Cell-Based アルゴリズム (E.M.Knorr) をもとした差分処理を行う手法

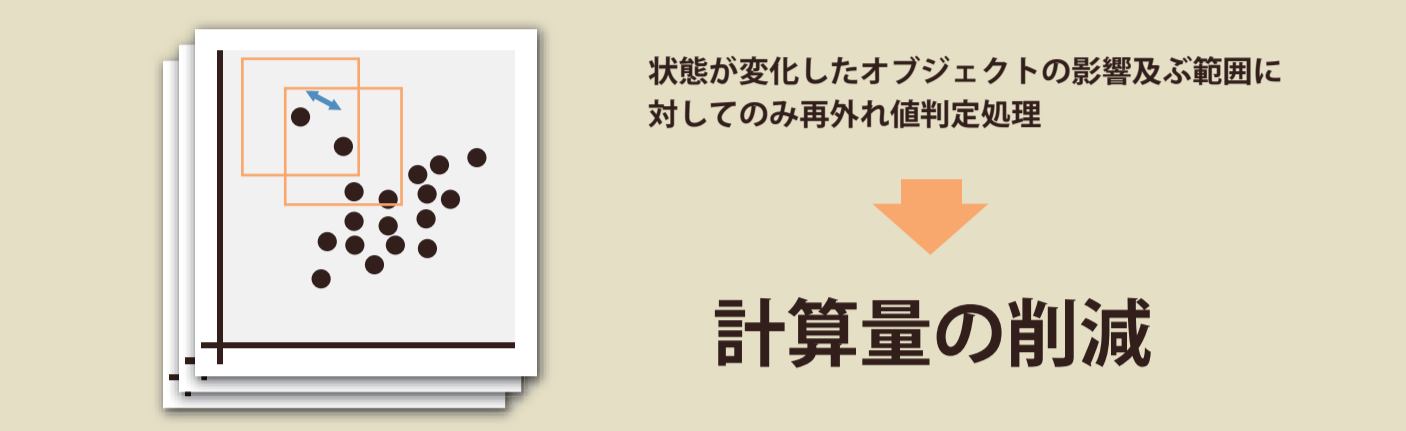


例：センサの観測値

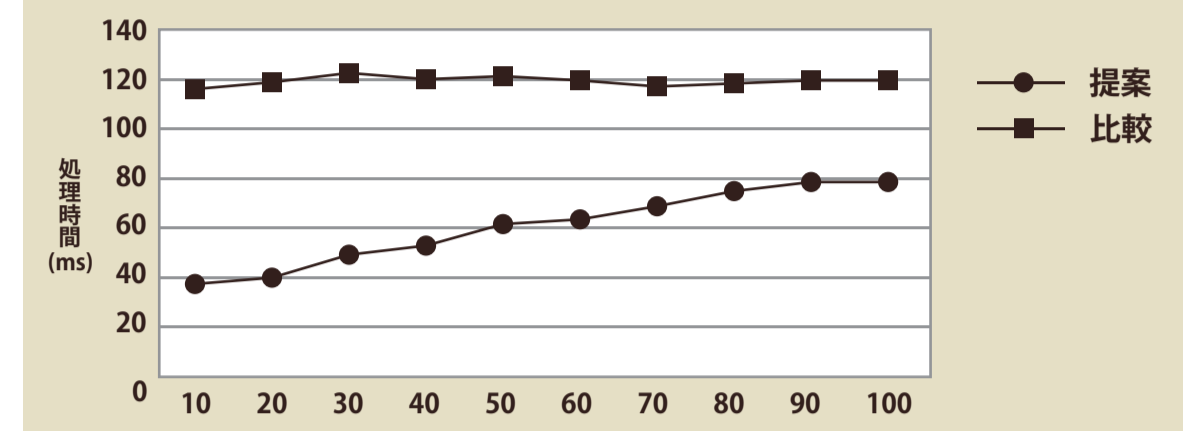
- データストリーム特徴
- データの分布が直前時刻のものと類似する場合が多い



提案手法：データストリームに対する外れ値検出アルゴリズム



実験結果 移動体の割合の変化



外れ値トランザクション検出

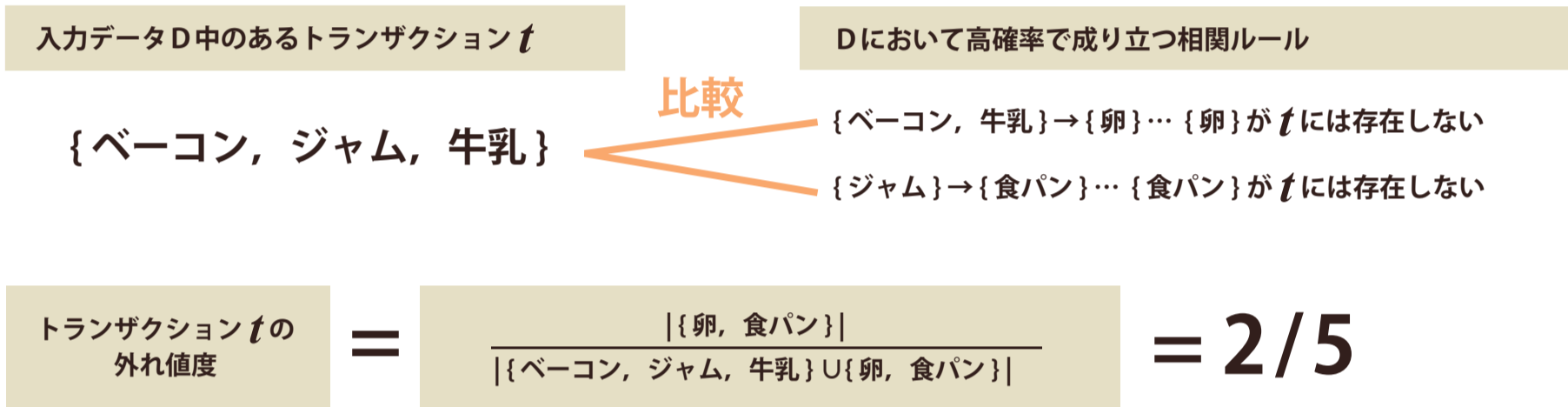
WAIM'08 最優秀学生論文賞

トランザクションデータの規則性から大きく逸脱したトランザクションを検出する枠組みの提案

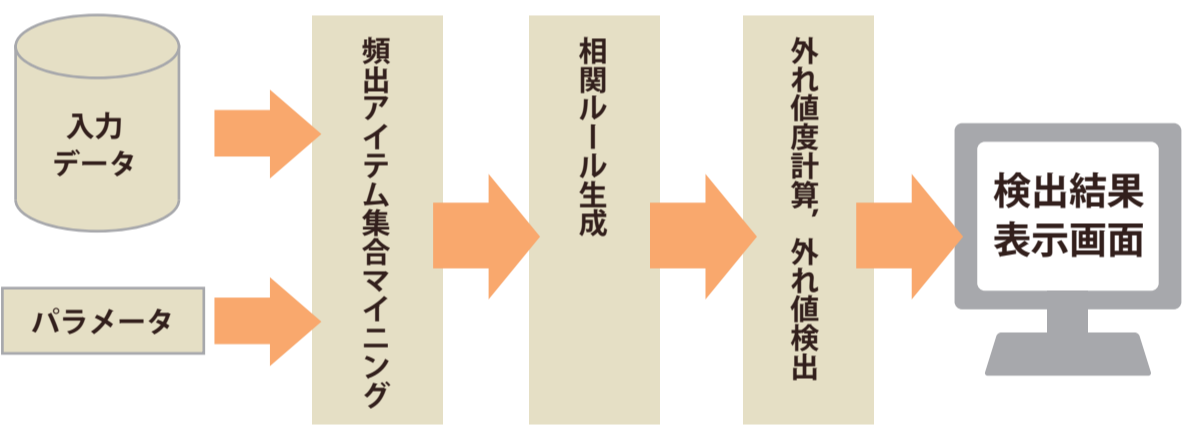
- 高確度の相関ルールに基づく外れ値度の導入
- 検出アルゴリズムの提案

外れ値度

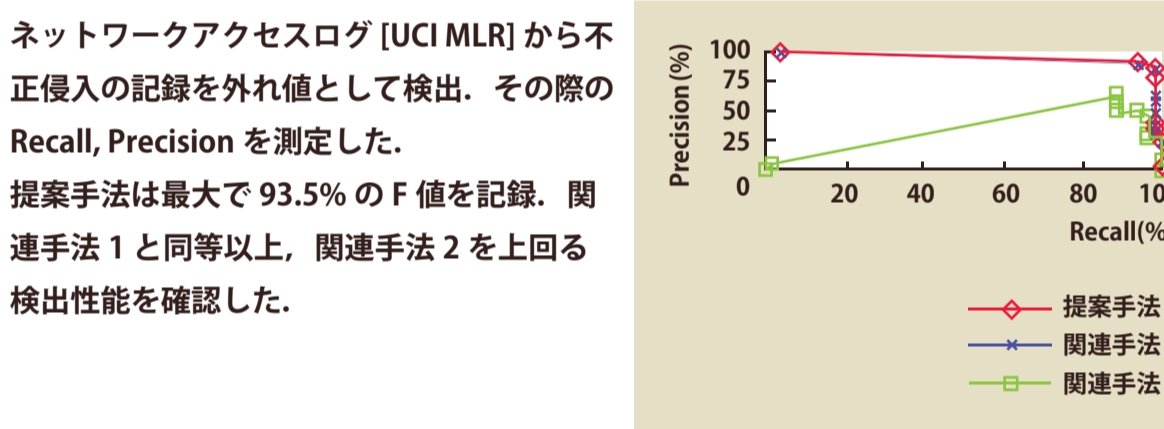
トランザクションの希少性 = アイテム集合と共起すべきアイテムが存在しない



検出アルゴリズム



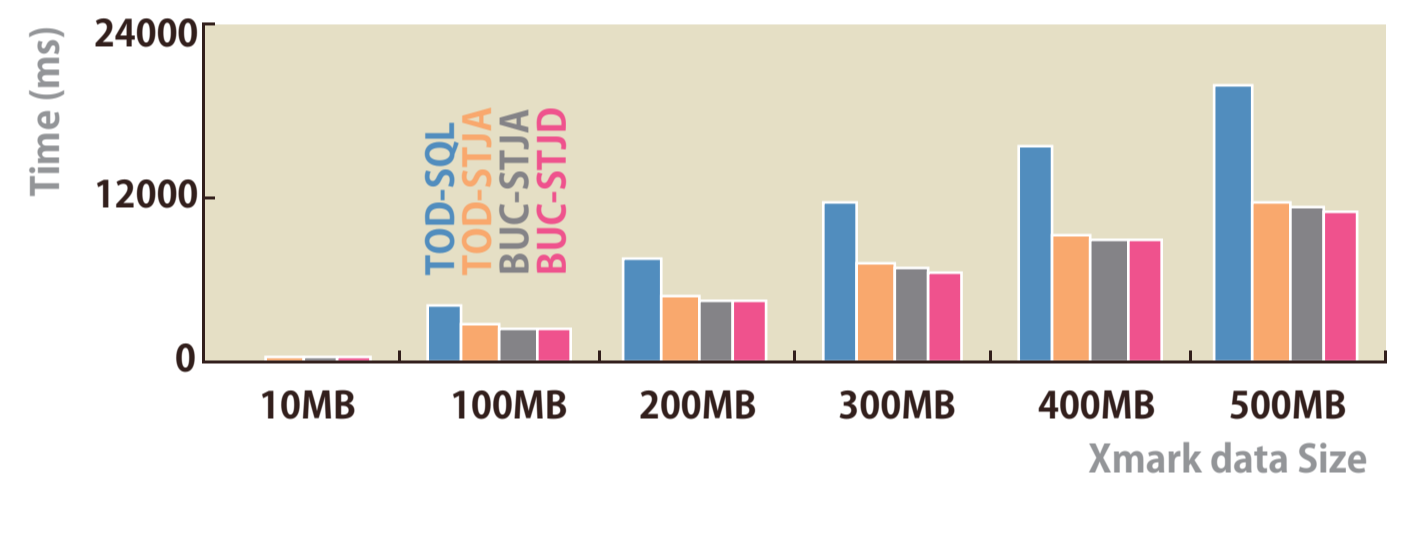
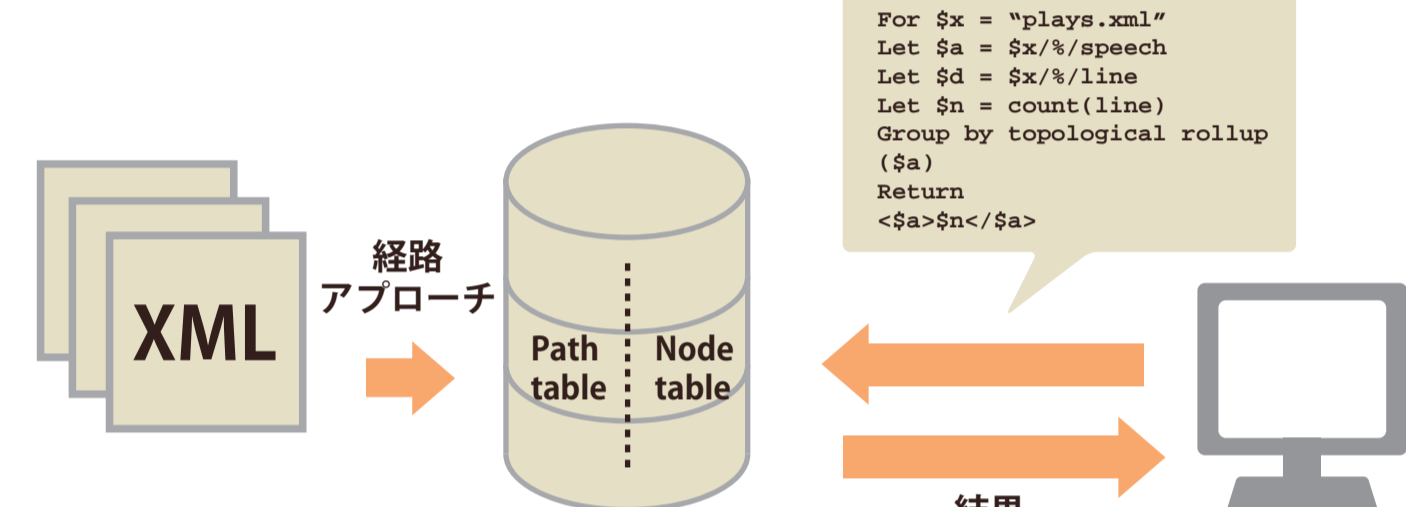
検出精度



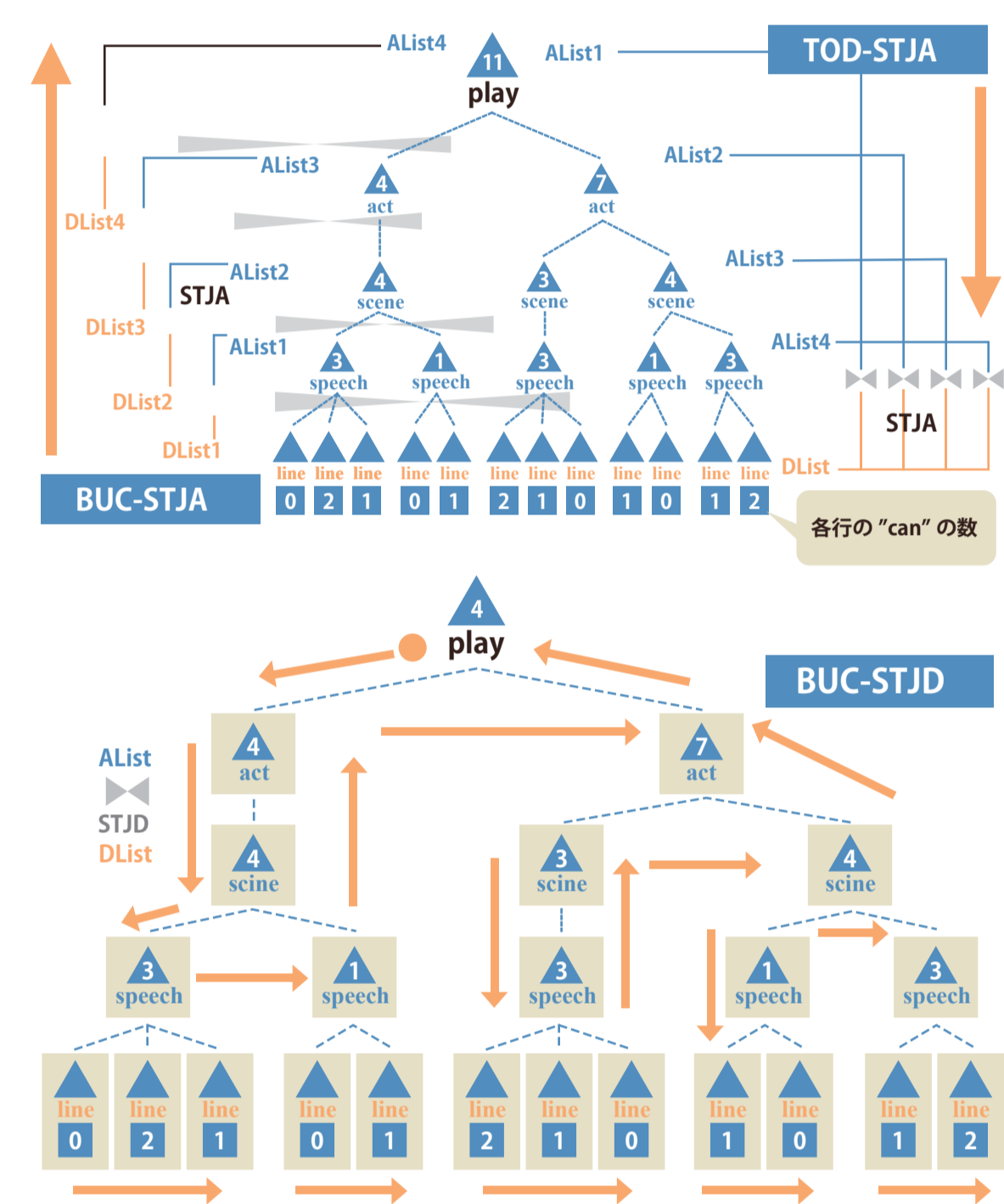
XMLデータの多次元分析

- XML はデータ流通とストレージに広く使われているが、解析処理はまだ行われていない
- XML-OLAP 技術の構築に挑む
- TOPOLOGICAL ROLLUP を Stack Tree Join を用いたアプローチで効率的に実現
- Top Down (TOD-STJA)
- Bottom-up (BUC-STJA)
- Bottom-up (BUC-STJD)

System Overview



TOPOLOGICAL ROLLUP により "can" をカウント

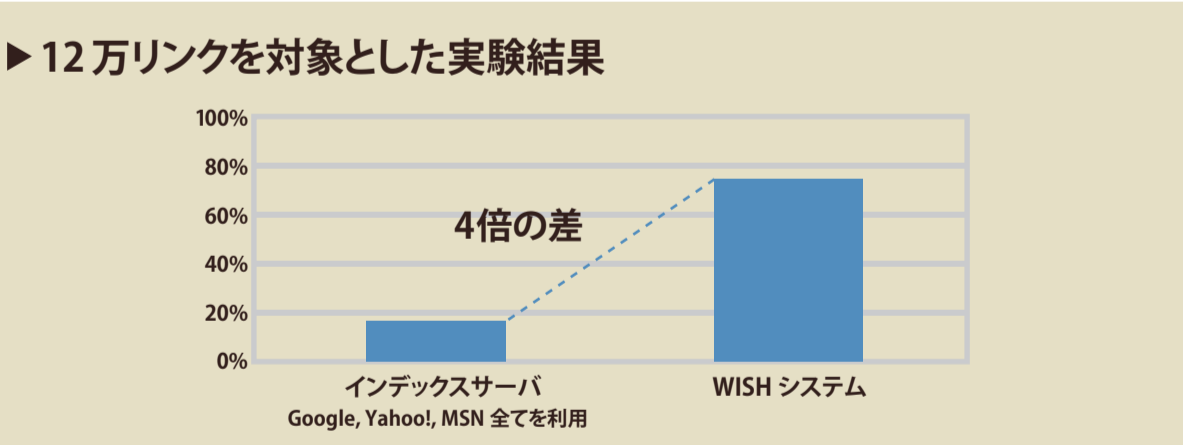
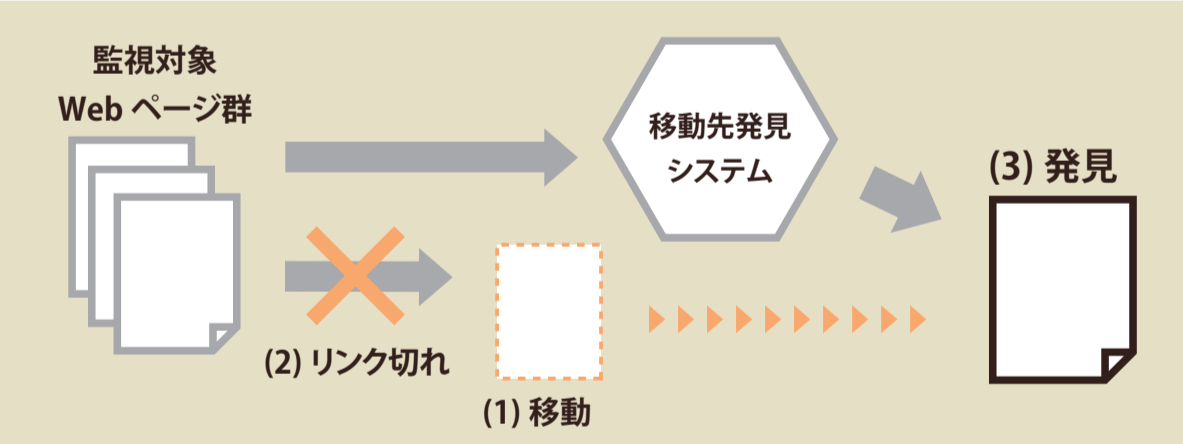


Webコンテンツ一貫性管理

IEEE ICDE'08 採録

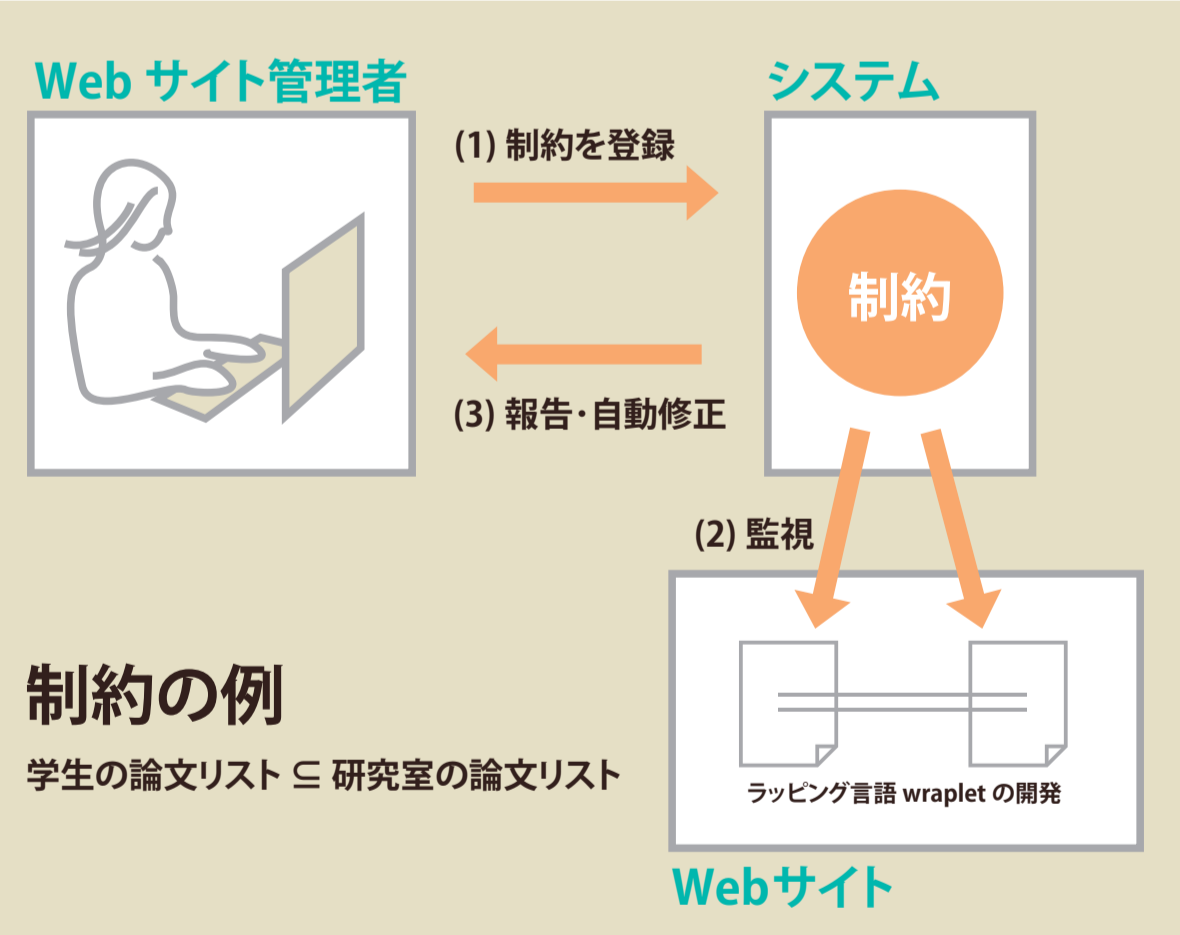
Web コンテンツ一貫性維持のためのページ移動先探索に関する研究

- Web ページの移動により生じるリンク切れの問題に着目
- ロボットにより Web ページ群を監視し、リンク切れを発見したときに Web ページの移動先を探索



明示的な制約を利用した分散管理 Web コンテンツの一貫性維持

- 既存の Web コンテンツに対して後付けでコンテンツ一貫性管理を実現可能
- ラッピング・制約指定支援技術の開発



SBMを利用したランキング

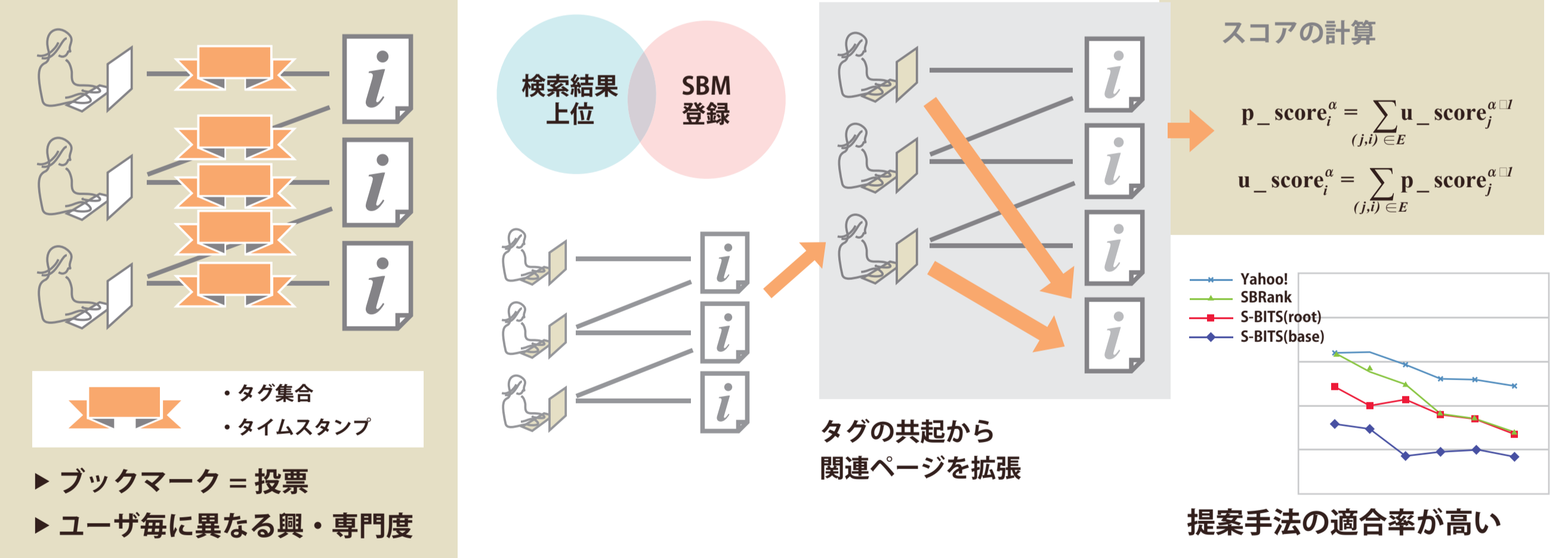
WAIM'08 採録
DASFAA'09 採録

ソーシャルブックマーク (SBM) ユーザの興味・評価、専門度 (Hub 度) を反映させたランキング

- 良いページは多くの良いユーザに参照される
- 良いユーザは多くの良いページを参照する

ソーシャルブックマーク

提案手法：S-BITS

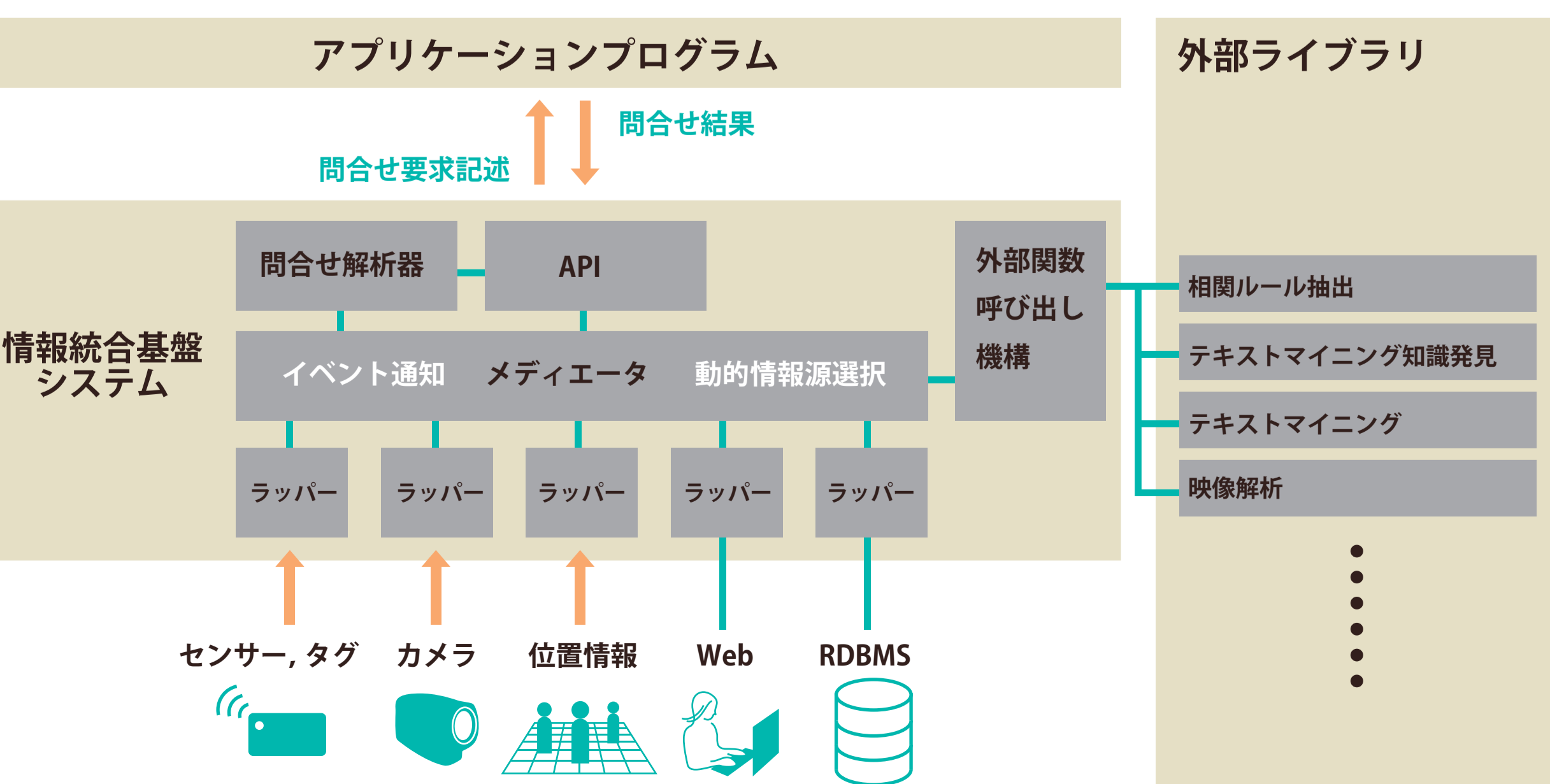


能動的情報統合基盤システム

Information Systems 条件付採録

ストリーム等を含めた情報統合基盤：StreamSpinner

- データ到着やタイマーに連動し、イベント駆動で能動的に各種統合処理を実行
- 外部関数呼び出し機構やアプリケーション記述のための Java API による拡張性
- 膨大な情報源の中から、利用者の興味に応じて接続対象を動的に選択可能



更新に最適化した空間索引

IEEE Trans. on Knowledge & Data Eng. 採録

Rsb-tree: 準バルクロードを用いる R-tree

- 移動体のように頻りに値が変わるデータへの索引機構の提案
- 準バルクロード：更新をメモリでバッファしてディスク転送コストを削減
- バッファサイズ、ページサイズ、更新/問合せ比率を変えて実験し、性能向上を確認
- I/O を要する演算を遅延 - split, remove, overflow treatment, and underflow treatment

