

研究テーマ(有村)

■ 北大大学院情報科学研究科, 教授

■ 専門:

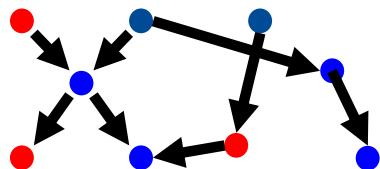
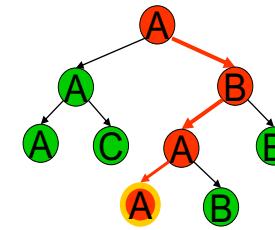
- データマイニング
- 情報検索(とくに全文テキスト索引)
- 計算学習理論(機械学習)

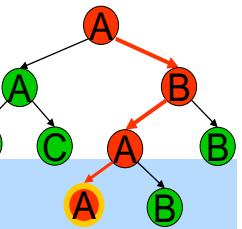
■ 興味があること

- 膨大なデータから, 人間に役立つ情報と知識をとりだすこと
- 高速なアルゴリズム(プログラム)を設計すること

■ 趣味

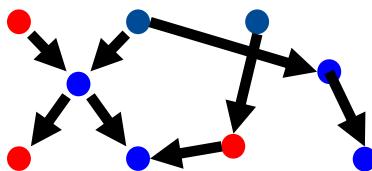
- 読書・音楽鑑賞・ハイキング・家庭菜園の手伝い





■ 情報検索(「文字列」-)

- 文字列, 木, グラフ, 集合族, . . .
- 検索・比較・発見・圧縮
- アルゴリズムとデータ構造



■ データマイニング

- 列挙アルゴリズム: 多数の解をもれなく計算する
- 構造・パターンマイニング

■ 機械学習(計算学習理論)

- 信頼されるAI(公平性・説明性)
- 離散機械学習(決定木ほかk)
- 効率良い学習アルゴリズム

やりたいこと

- 膨大なデータから、人間に役立つ情報と知識をとりだす
- 高速なアルゴリズム(プログラム)を設計する

対象

- 文字列・木・グラフ
 - 例:テキスト・遺伝情報, etc.
- 空間データ
 - 例:地理・移動・音楽・???データ

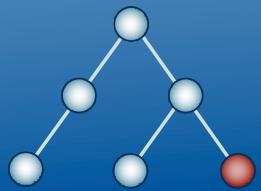


■ どちらかというと基礎

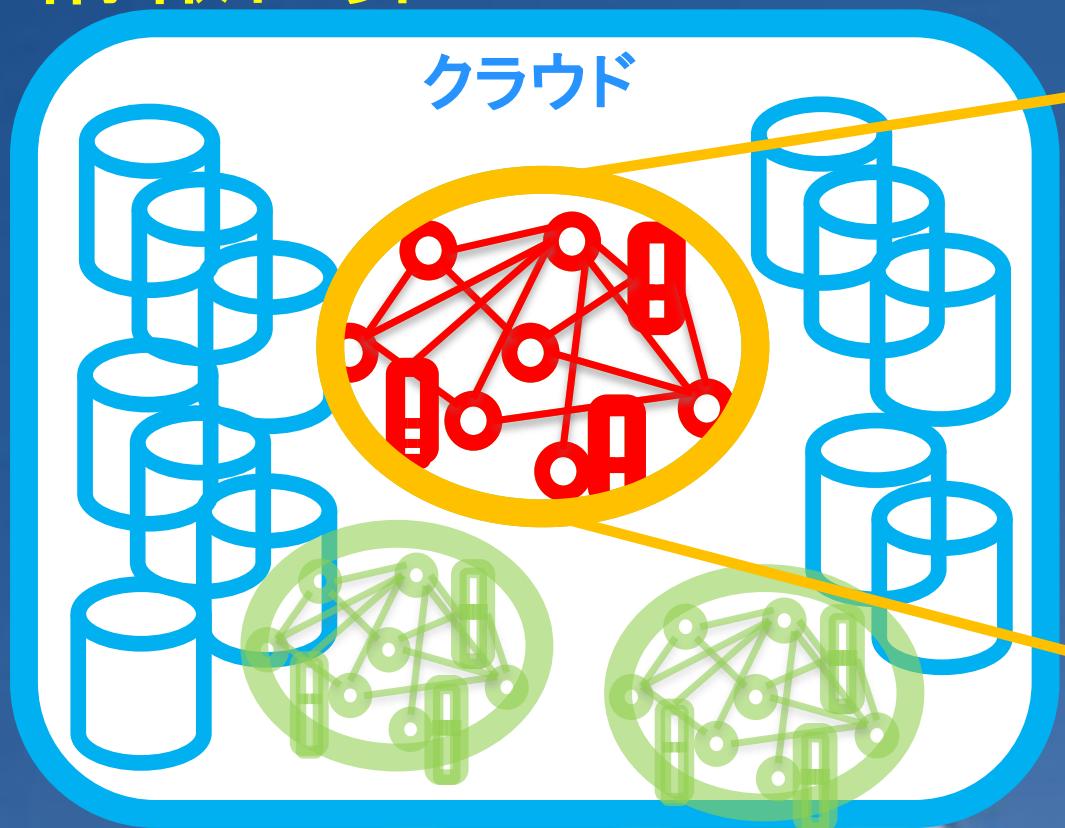
- アルゴリズムが中心(道具)
- 情報は「形」「道具」の学問
- 先生方は二つの専門を持っている
 - 道具(How) :
 - 対象(What) : Web, Networks, 生命科学,
機械学習, 検索, データ解析. . .
- 現代は, いろいろなタイプの人と一緒ににはたらくのがふつう: 基礎・応用・人間/社会

■ こんな人

- 考えるのが好き and/or 手を動かす



情報世界



- 大量のデータ
- 多数のCPU
- 高速なネットワーク
- 膨大な計算

「集中」

実世界



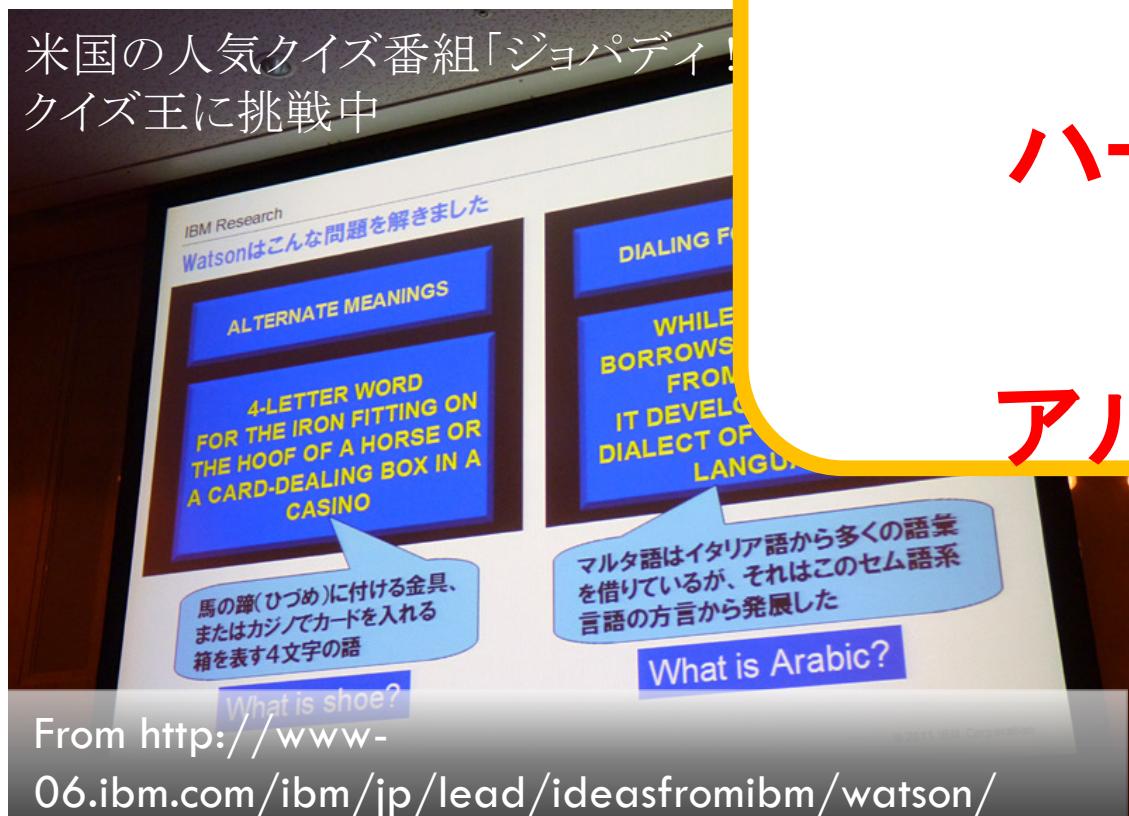
- さまざまなデバイス
- 多様な人間活動と応用
- 多様で非均一な時空間
- 不完全で複雑なデータと情報

「分散」

大量データと次世代情報技術

ワトソン君

- IBMリサーチ (2011/02/16)
- クイズ番組で人間に勝利！
- 100万冊の本を読んで回答
- 人工知能と自然言語 アルゴリズム、検索の技



クラウド計算
世界中の情報を計算！

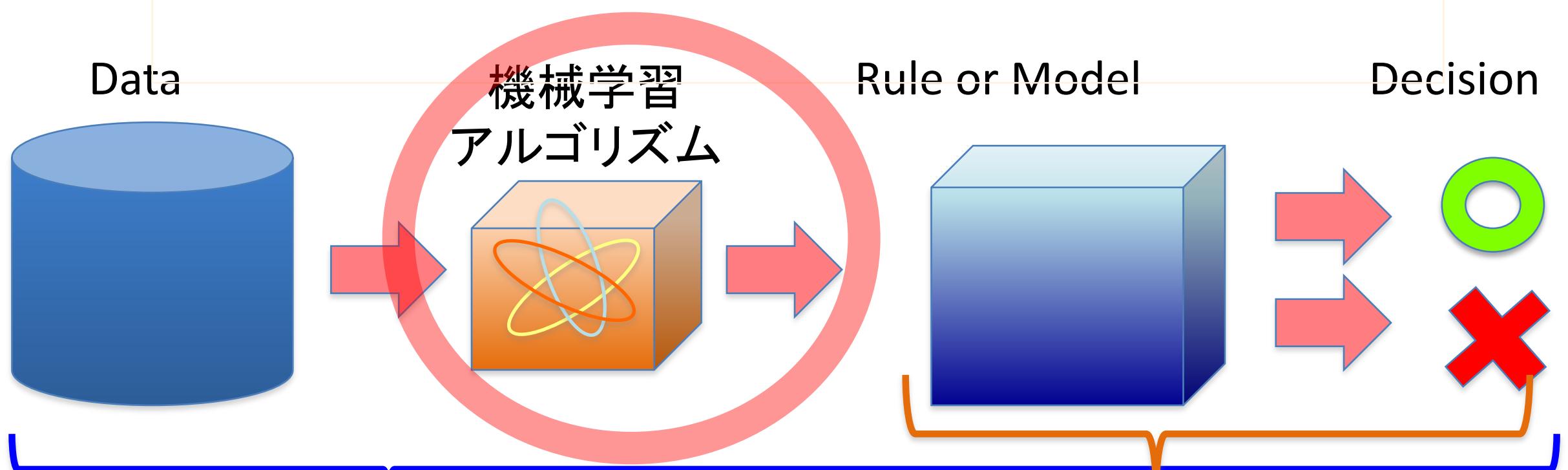
次世代IT(AI)の原動力

ハードウェア!
データ!
アルゴリズム!



現代の人工知能

データから規則を学習する**機械学習技術**
が中核技術になっている

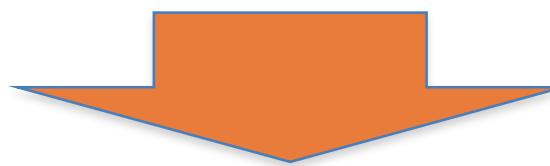


現代のAI（データ指向AI）はデータに基づく意思決定のすべての局面をつかう

伝統的なAIは規則を使った意思決定だけをつかう

次世代の発見科学技術

- * 最新の人工知能技術
- * 21世紀はDeep Learningだけで

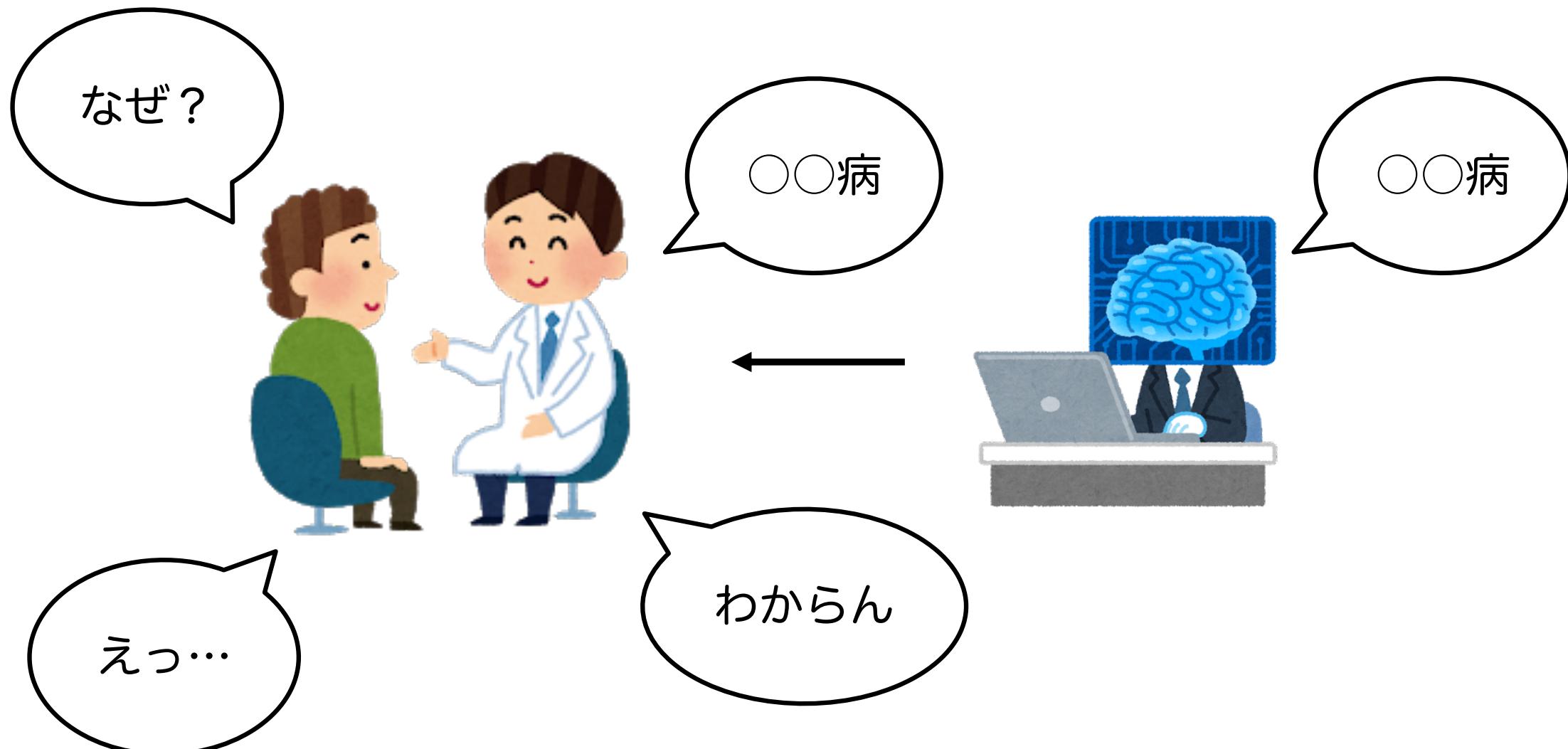


有村君、21世紀
の知識獲得は
Deep Learningで
決まりかね？

- 21世紀の社会で、コンピュータが
人間と働くために必要な
人工知能技術とは何か？

有川節夫先生
(私の先生)

信頼されるAIの基盤技術



どうやってAIの導きだす回答を
信頼させられるものにするか？

研究1：予測モデルに対する反事實説明 (Kanamori+, IJCAI2020)

Accepted at
IJCAI 2020

DACE: Distribution-Aware Counterfactual Explanation by Mixed-Integer Linear Optimization

-
- 金森 憲太朗 (北海道大学)
 - 高木 拓也 (富士通研究所 / 理研AIP)
 - 小林 健 (富士通研究所 / 理研AIP / 東京工業大学)
 - 有村 博紀 (北海道大学)

Take-home message: 機械学習と説明可能性

- 機械学習の目的の一つは、入力空間 \mathcal{X} と出力空間 \mathcal{Y} を上手に対応づける関数（**予測モデル**） $h: \mathcal{X} \rightarrow \mathcal{Y}$ を見つけること。
- **実応用上の要求:** 予測の**説明可能性**（解釈可能性）
 - ▶ 知識発見: モデルの予測結果や傾向から、新たな知見や仮説を見つけてたい。
 - ▶ 信頼性: 医療や司法などの意思決定に応用するには、予測根拠の提示が必要。



- 予測結果の説明（のようなもの）を抽出する研究が注目されている。
 - ▶ 例) 特徴量重要度: 各特徴量が予測にどのくらい貢献したかを示す指標。
 - ▶ いまのところ説明可能性（解釈可能性）にちゃんとした定義は存在しない。

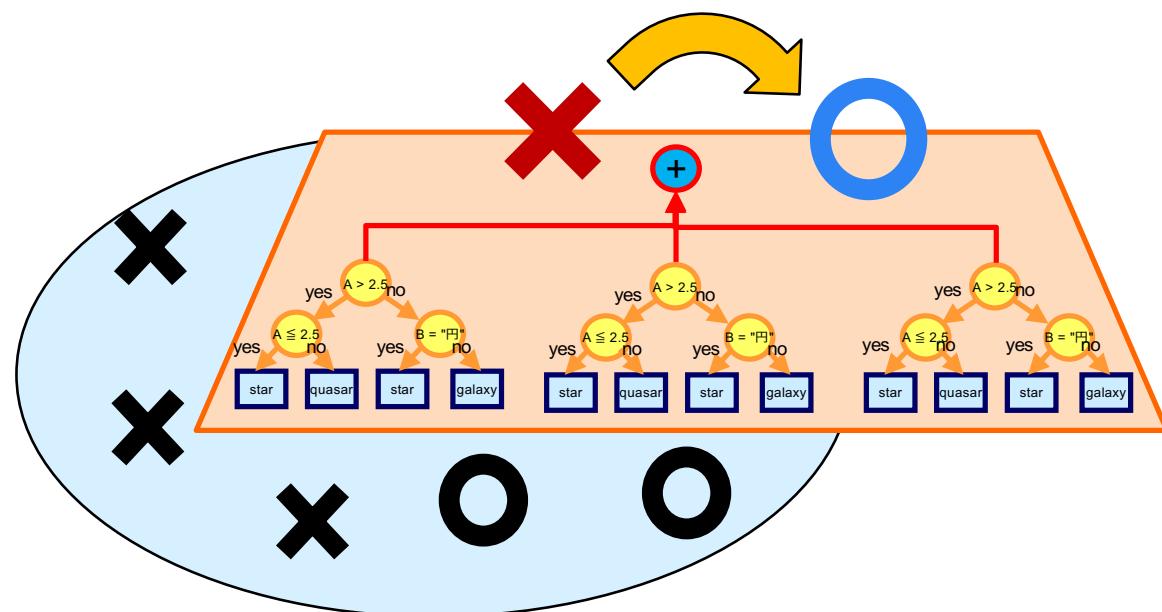
研究: 人に優しい人工知能技術

K. Kanamori, T. Takagi, K. Kobayashi, H. Arimura, IJCAI-PRICAI 2020 (to appear)

機械学習の解釈性

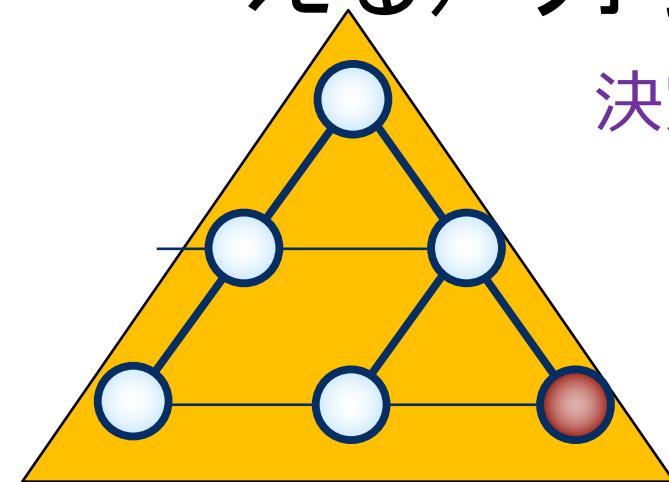
- 解釈性・説明性・公平性
- 反事実説明 (Kanamori+, IJCAI2020)

「予測を×から○へ改善するためには、状況をどう変るべきか？」



離散最適化に基づく学習

- 解釈しやすい離散的ルールを学習する
 - トップK・計数(数を数える)・列挙・乱択
- 決定木を列挙!

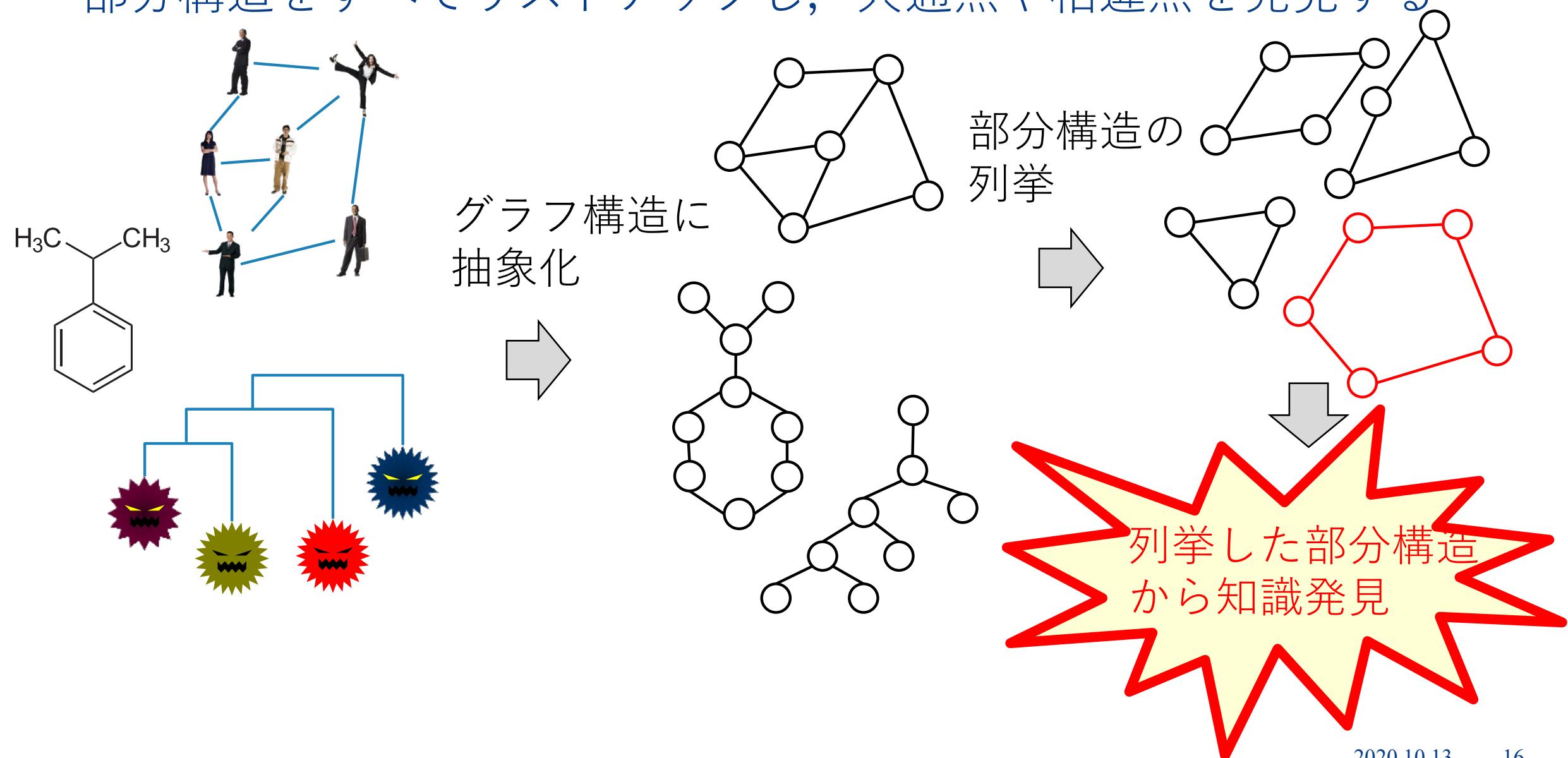


- 制約
- サイズ
 - 葉の例数
 - 精度

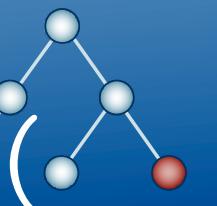
研究: グラフアルゴリズム & 列挙アルゴリズムによる知識発見

どんな構造が含まれているか, すべて挙げよ!

部分構造をすべてリストアップし, 共通点や相違点を発見する

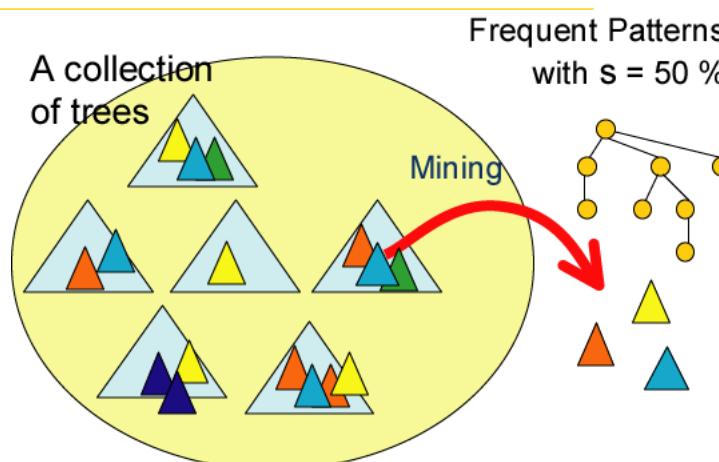


研究テーマ: 半構造マイニング「最右拡張技法」(



高速半構造データマイニングエンジン

- Discovering all frequent sub-structures from a collection of labeled trees
- Extendible to most statistical functions



- **FREQT:** Efficient ordered tree mining engine (SIAM DM'02)

Applications

Japanese Text Mining
(Morinaga, Arimura et al.)

Rightmost Expansion & Ordered Tree Enumeration Trees (SIAM DM'02, IEEE ICDM'02, DS'03)

FREQT: DEWS'02優秀論文賞 (H14年6月)

Unot: DEWS'04優秀論文賞 (H16年6月)

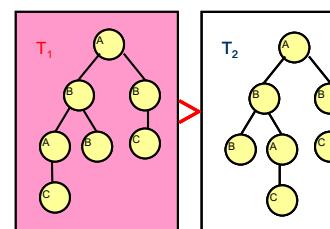


「最右拡張技法」に基づく高速な木パターン発見エンジン

- 半構造データの特徴的部分構造の発見
- その後のグラフマイニングの基礎技術に
- 理論と実装: FREQT, StreamT, Unot
- SIAM DM'02, PKDD'02, IEEE ICDM'02, DS'03
- 公開&応用

trees (Left-Biased Tree)

(Lexicographically largest trees over depth-label sequence encoding)



(0A,1B,2A,3C,2B,1B,2C) (0A,1B,2B,2A,3C,1B,2C)

Difficulties in un-ordered tree mining

Exponentially many isomorphic patterns

■ Unique represen-

■ Enumeration with

duplicates in a un-

■ Efficient comput-

occurrences

Google Scholar
引用数(2019.4)
[主要文献] 596件

Awarded

MaxMotif: Awarded DEWS2004 (H17.6)

StreamT: '03AI学会大会優秀賞 (H15.6)

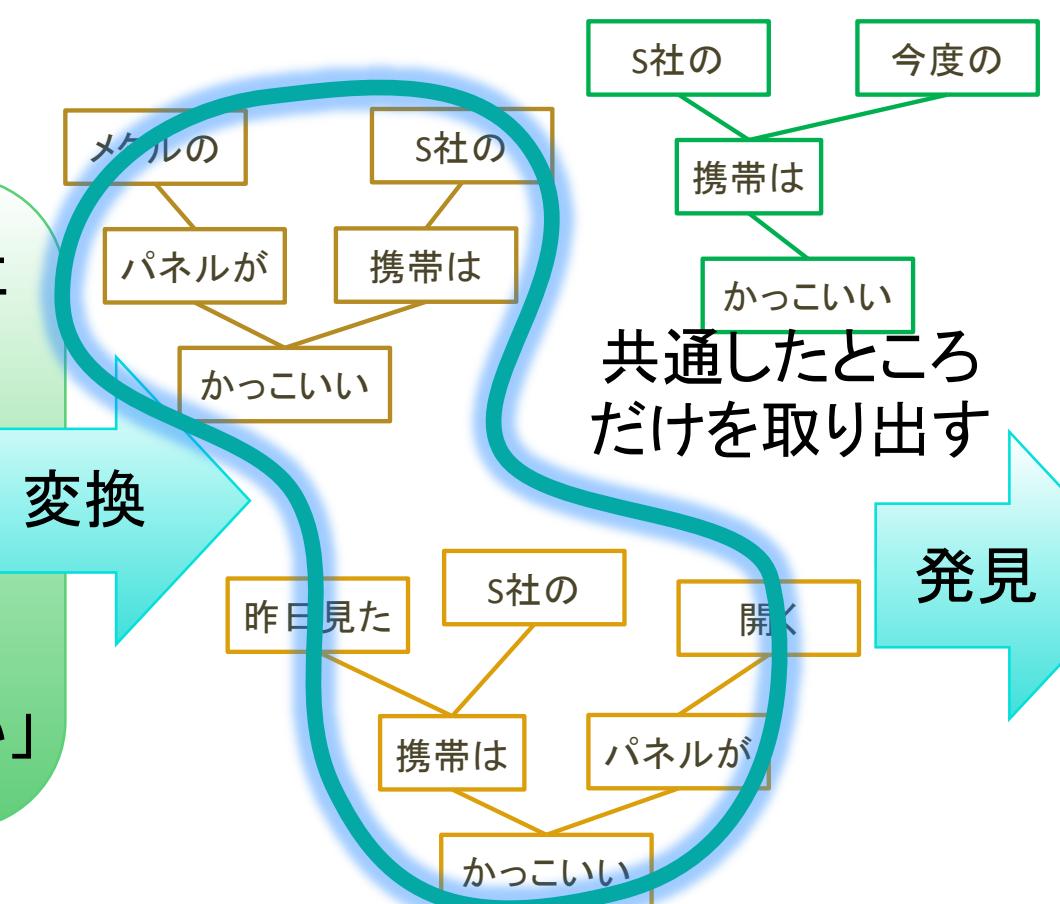
コンピュータの学習の例

テキストマイニング：

- * まず大量のテキスト（ブログでの評判や自由記述アンケート）を解析し、文の構造を「木」としてとりだす。
- * 次に、共通する構造を発見（マイニング）すると、文書に隠れた共通の意見が取り出せる。

入力のいろいろな文章

「メタルのパネルが、S社の携帯はかっこいい」
 「今度のS社の携帯はパネルがかっこいい」
 「昨日見たS社の携帯は開くパネルがかっこいい」



この仕組みは、企業でテキストデータの解析に使われた。

人工知能技術の鍵！ = 機械学習技術

NIPS 2017(機械学習研究最大の学会)

- Thirty-first Conference on Neural Information Processing Systems, Dec. 4-9, Long Beach Convention Center
- Attracting over 8,000 registered attendees, up 2,000 from last year

Registration



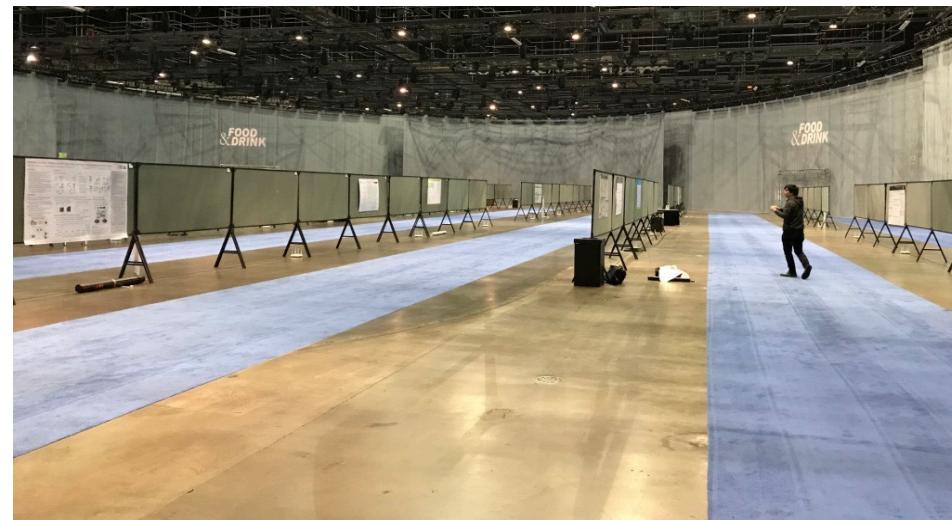
exhibition booth



session



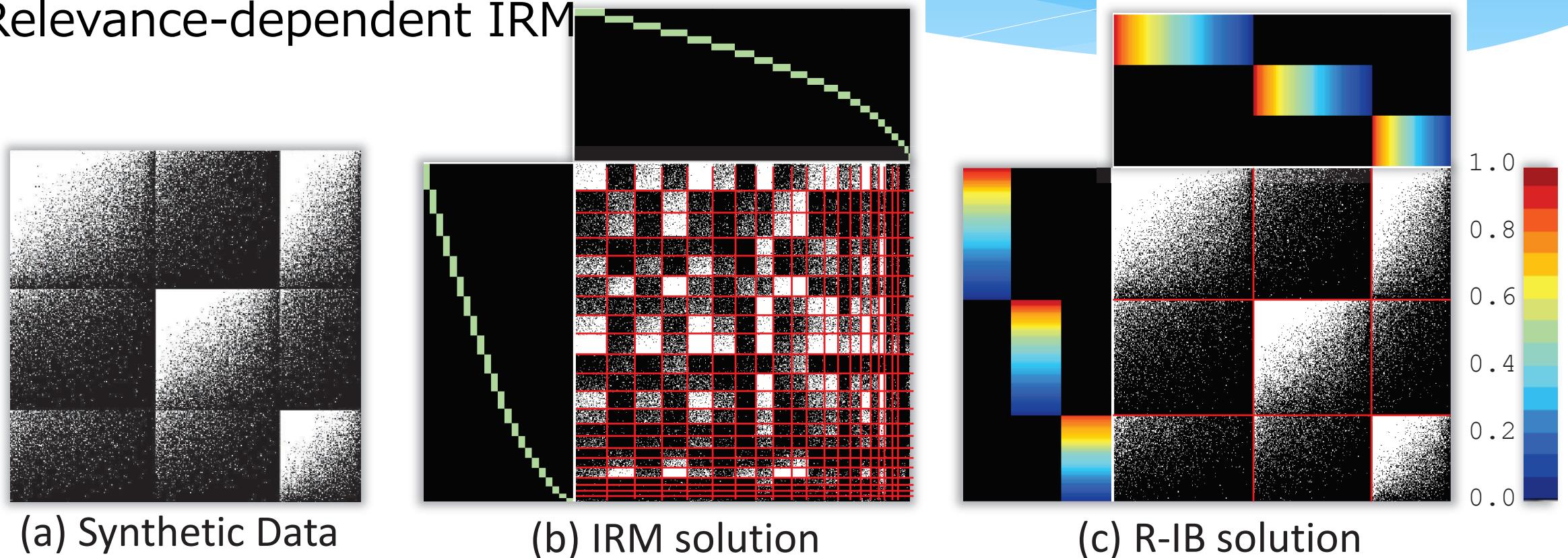
Poster site



研究: ベイズ学習・トピックモデル

(Ohama, Iida, kida, Arimura, IJCAI 2017, NIPS2017)

Relevance-dependent IRM



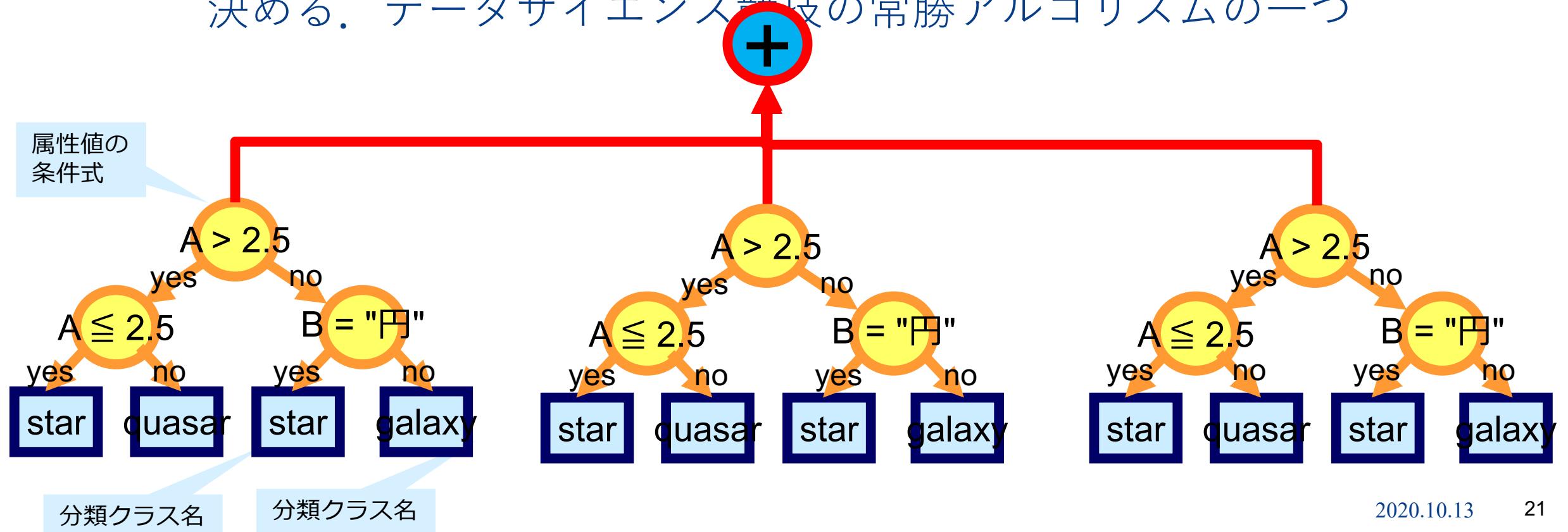
- 従来手法(IRM)は、個々の人や商品の「活動度」に影響されて本来存在しない、小さなクラスに分かれている。
- 提案手法(R-IB)は、個々の活動度に影響されずに、大きな(正しい)クラスタ構造を見つけている

Figure 3: Synthetic example: (a) 500×500 relational data; (b) IRM solution; (c) R-IB solution. In (b) and (c), the left and top matrices indicate z_1 and z_2^\top in a 1-of-K representation, respectively. Colored areas on matrices z_1 and z_2^\top indicate relevance parameters for corresponding objects. For intuitive understanding, a relevance parameter $\theta \in [0, +\infty]$ is transformed into a probability as $1 - \exp(-\log(2) \times \theta)$, so that the probability is 0.5 when $\theta = 1$.

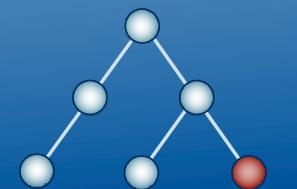
研究: 機械学習の研究

* 決定木(decision tree)と木アンサンブル(勾配ブースティング, ランダムフォレスト)による自動分類

- ▼ 決定木: 対象を, その属性の値に従って分類する. 基本は「すごろく」と同じ. 条件式の値 (YES/NO) にしたがって, 根から葉までたどる. ゴールが分類クラス.
- ▼ 木アンサンブル: 多数の決定木の予測を, 多数決して全体の予測を決める. データサイエンス競技の常勝アルゴリズムの一つ

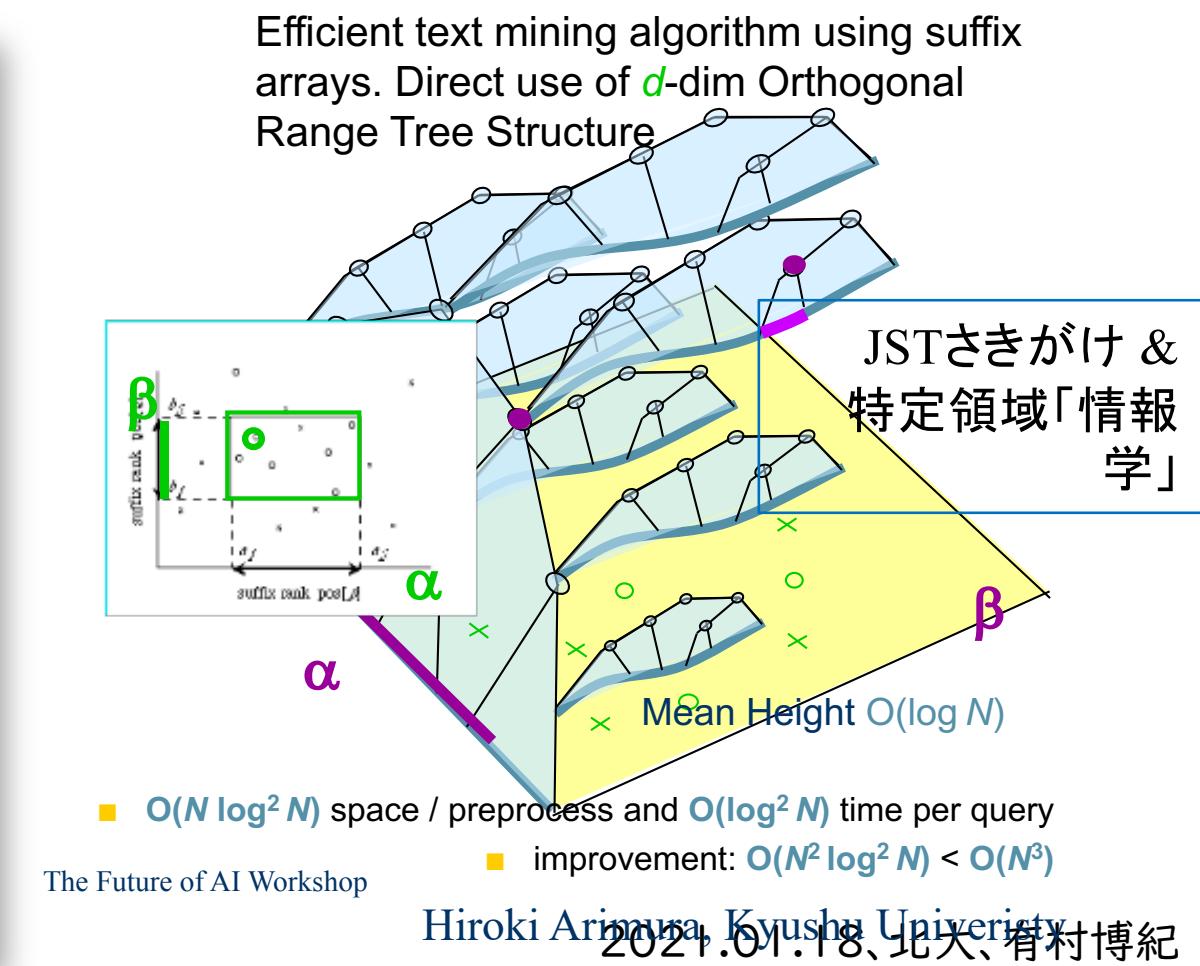
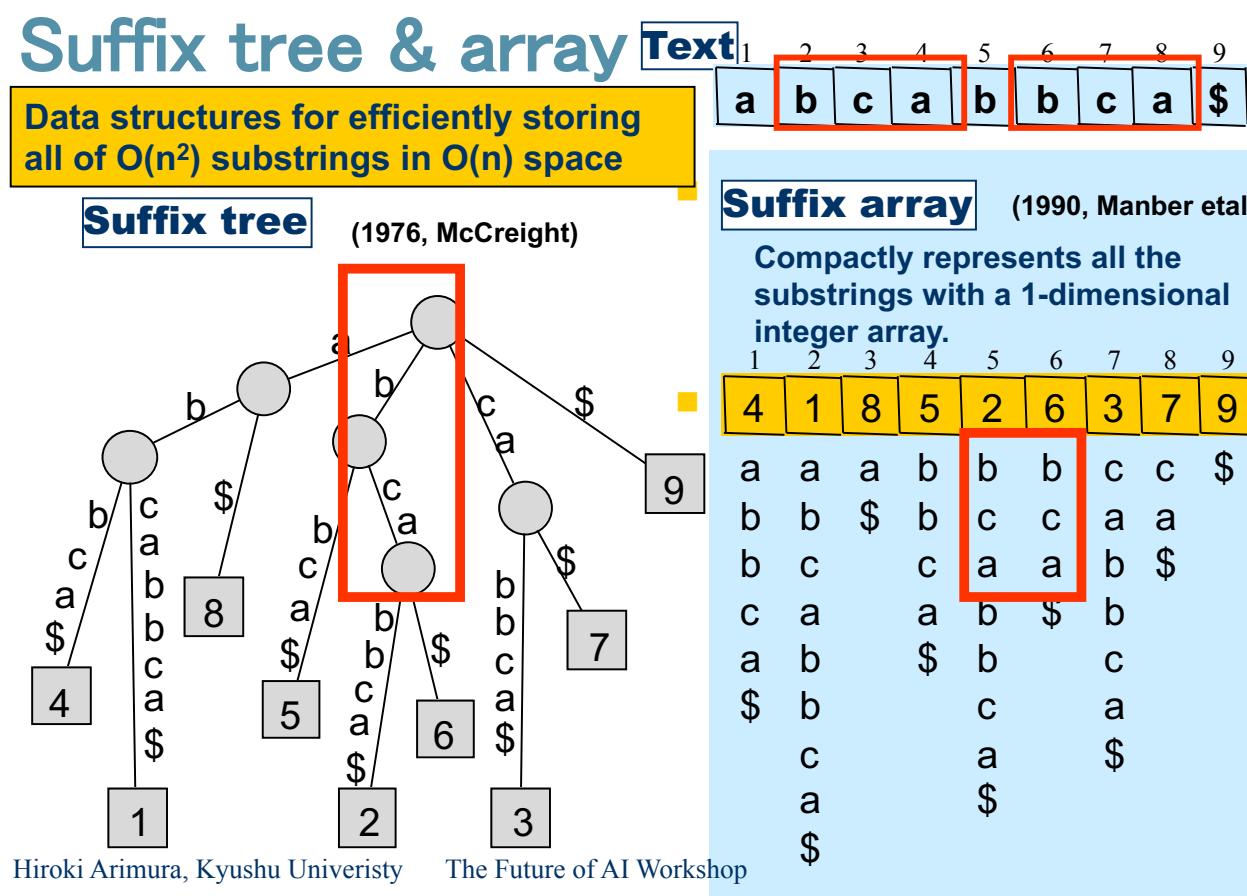


研究テーマ: ビッグデータの検索



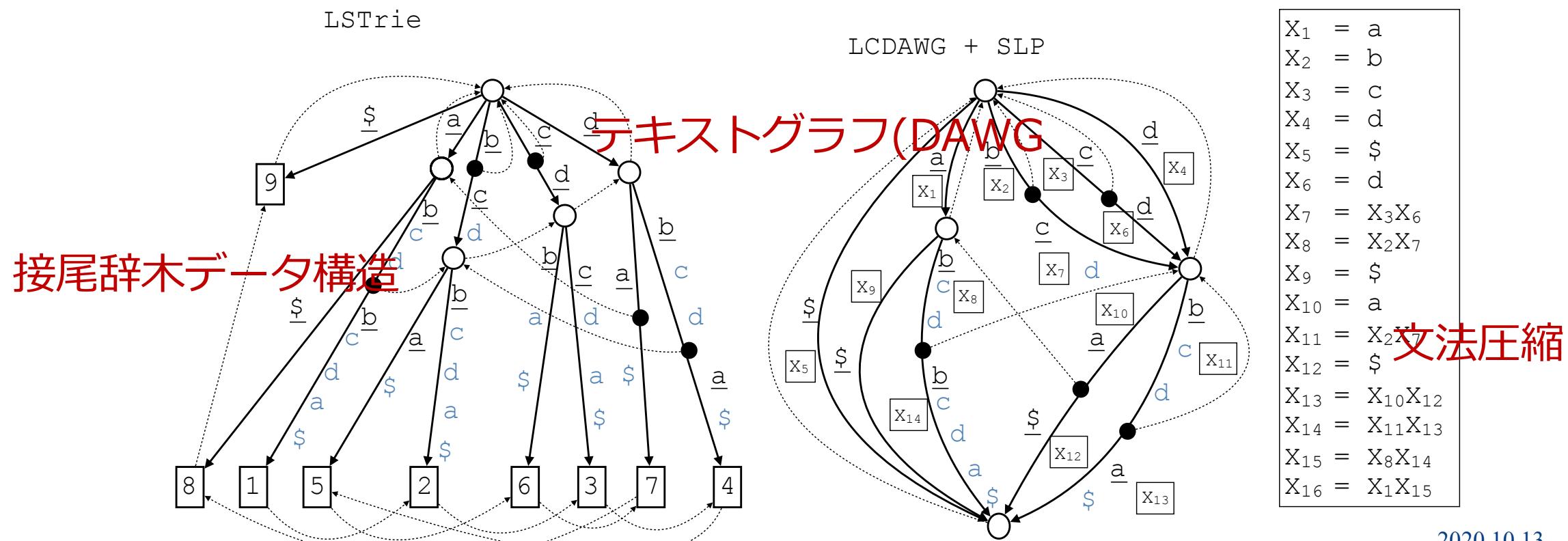
- 世界初の線形LCP配列計算アルゴリズム (CPM2001)
- 研究分野ベスト10引用論文として表彰: 2010年に分野トップ会議(CPM)での「接尾辞木データ構造40周年記念」集会で表彰
- CACM紙の「接尾辞木データ構造40周年記念」特集で言及された: 「数十年以上の未解決問題を肯定的に解決し, 接尾辞配列と接尾辞木の相互変換のMissing Linkをつないだ」。

Google Scholar
引用数(2019.4)
[主要文献2] 535件



研究: 文字列アルゴリズムとデータ構造

- * 大規模テキスト・系列データを効率よく圧縮・検索・解析するためのアルゴリズムとデータ構造を研究する
- * 文字列・正規表現・木構造・NoSQLの検索アルゴリズム
- * 接尾辞木と接尾辞配列, BWT, テキストグラフ(DAWG)（二分探索木や, ハッシュ表, 領域検索, RMQなどのテキストへの拡張）





Optimally Computing Compressed Indexing Arrays Based on the Compact Directed Acyclic Word Graph

Hiroki Arimura¹

Shunsuke Inenaga²

Yasuaki Kobayashi¹

Yuto Nakashima²

Mizuki Sue¹

1) Graduate School of IST,
Hokkaido University, Japan

2) Department of Informatics,
Kyushu University, Japan

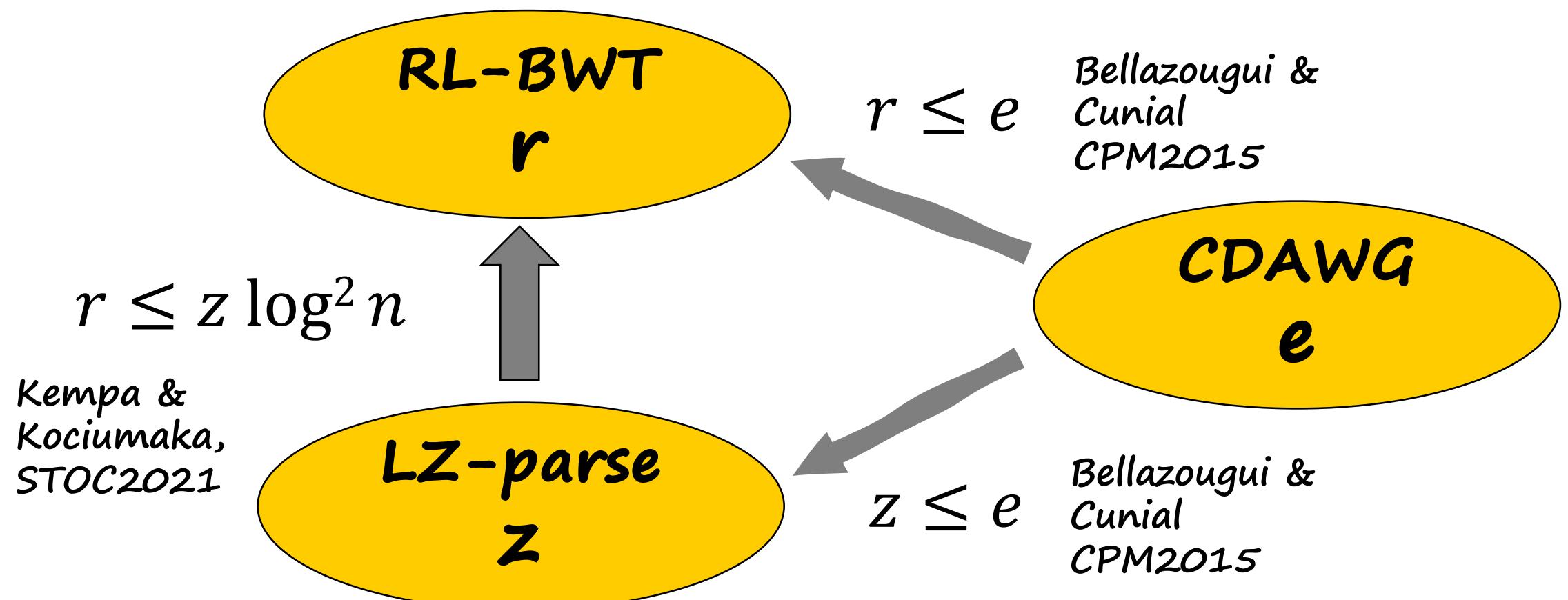
For details and proofs, visit the arXiv site of this talk:

- Manuscript at <https://arxiv.org/abs/2308.02269> ;
- Slide pdf at “Code section” or Github <https://ikndeva.github.io>)

This work is partly supported by MEXT Grant-in-Aid for Basic Research A, 2000-2004, Japan



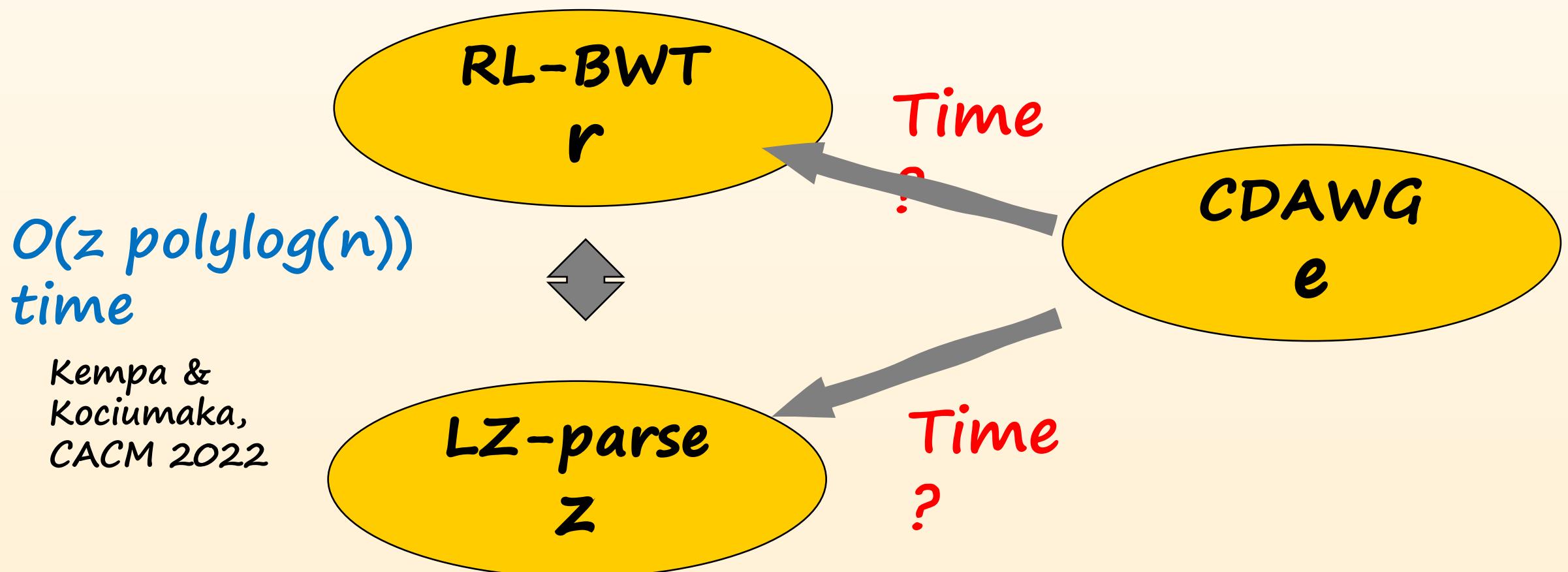
- Consider the relationship between their sizes
 - has been studied so far.





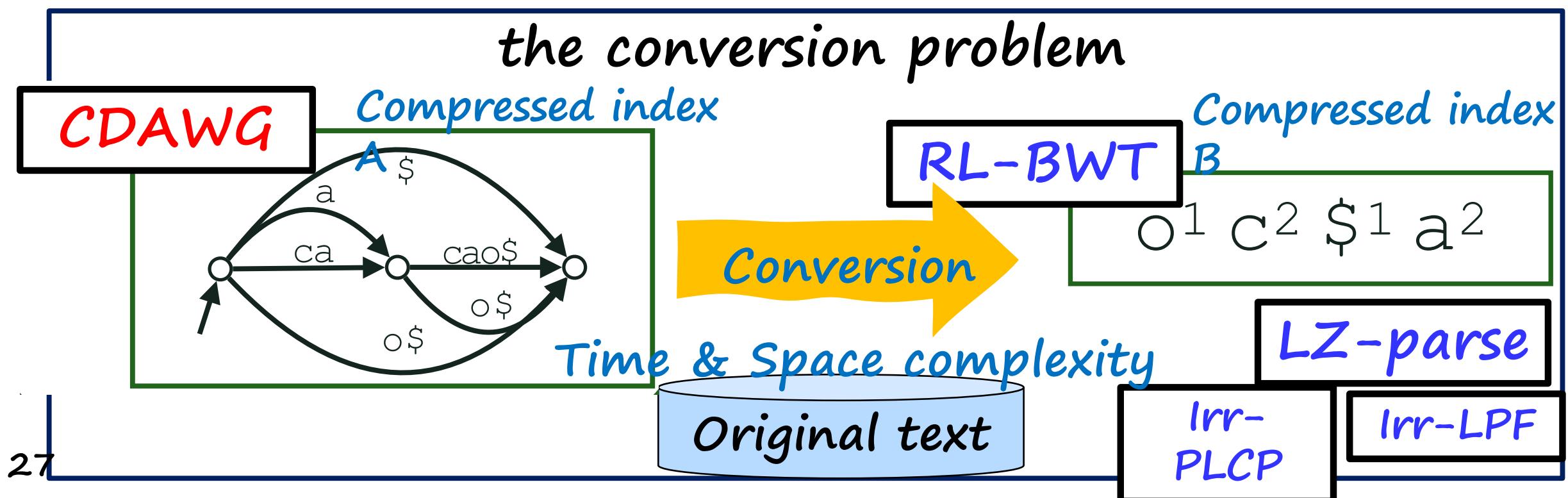
Backgrounds: Previous work

- The time and space complexities of conversions
 - have not been studied very much



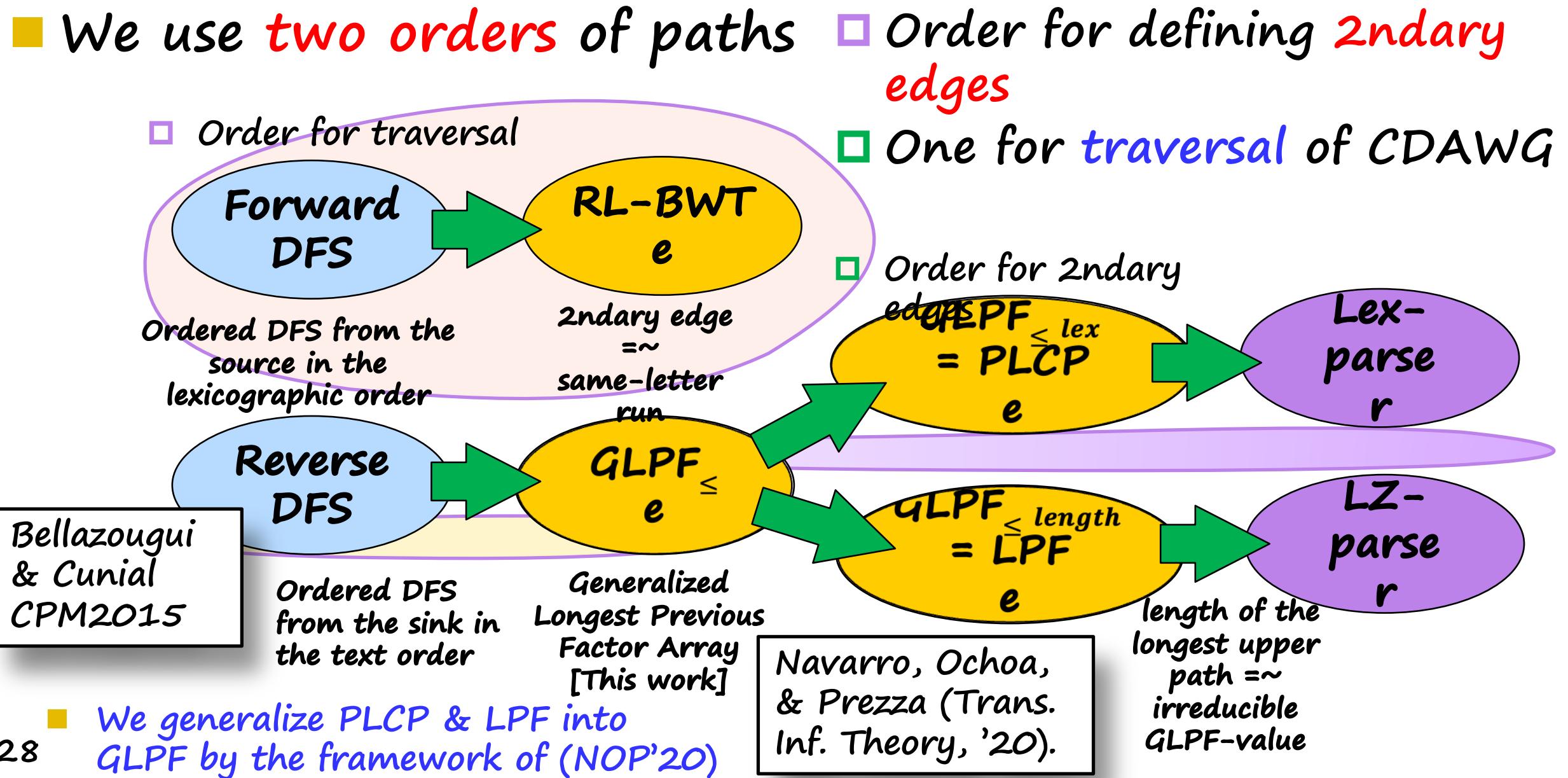


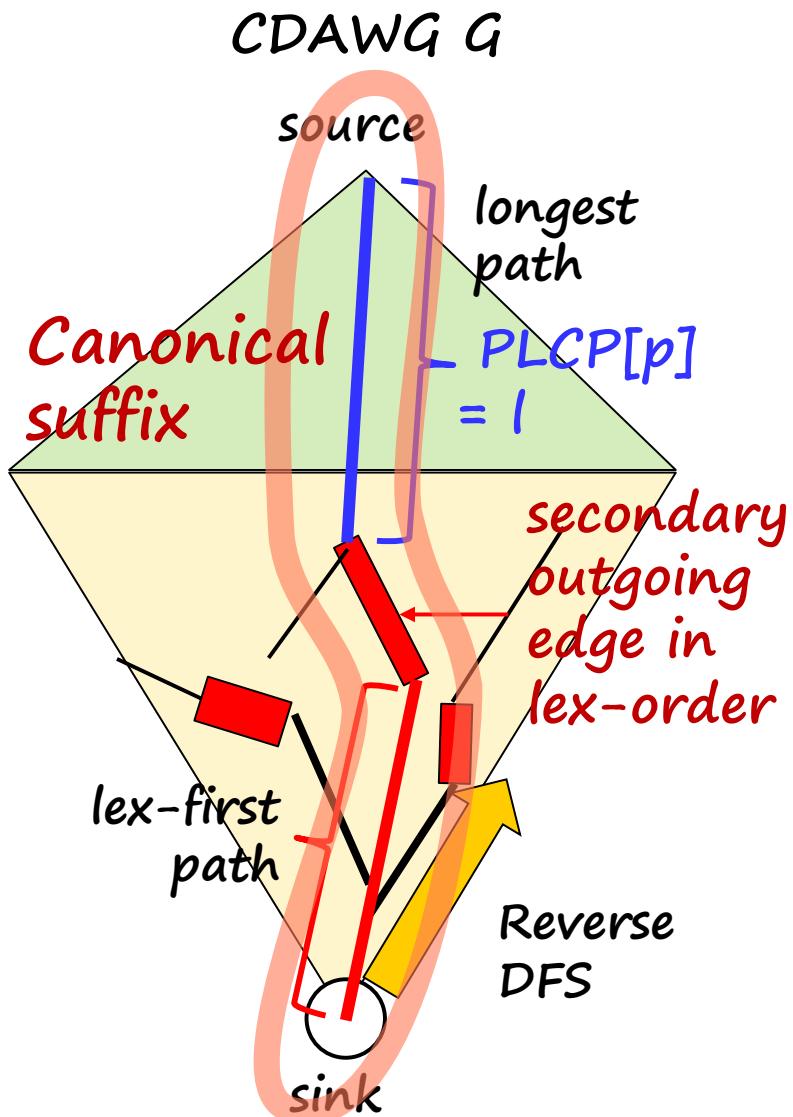
We devise efficient algorithms that solves the conversion problem from the CDAWG for a text T into various compressed indexes for T in linear time and space in the combined input/output sizes





Our approach





■ **Observation A1:** $O(e)$ secondary outgoing edge of $CDAWG(T)$ under the length-order determines the irreducible value $PLCP[p] = l$ by the length l of the longest path from the source to the corresponding branching node

■ **Observation A2:** $O(e)$ secondary outgoing edges can be enumerated in the text order of its “canonical suffix” by the reverse DFS from the sink.

We can extend the above result from PLCP to PLPF by employing the definition of 2ndary outgoing edges in length-order

おまちしています

情報知識ネットワーク研究室



有村・
{arim}@ist.hokudai.ac.jp
ex. 7678, 7679

