

ブートストラッピングによる スキーマ抽出の基礎検討

大阪大学大学院 基礎工学研究科
野村 慎太郎 中根 史敬
土方 嘉徳 西田 正吾

背景(1/2)

Web上の半構造化文書

属性名

- メーカー型番 : 170957I
- タイプ : B5ノート
- CPU : Core 2 Duo
- メモリ : 1GB
- HDD : 120GB

属性値



データベースのように整形されておらず、データを計算機で扱えない



スキーマ(属性名の組)の抽出
属性名に対応する属性値の抽出

比較的多くの
研究が
行われてきた

ほとんど研究が
行われていない

が必要

背景 (2/2)

PCについてのページ

ThinkPad X60 1709GDJの価格 1~15位

順位	価格	差額	送料	在庫	店頭
1位	¥192,800	最安	1,000	有	
2位	¥192,150	+26,350		有	

価格比較サイト
順位, 価格, 差額...

インテル® Core™ 2 Duo プロセッサ

製品シリーズ:	
製品番号:	
Microsoft® Office Edition 2003 搭載モデル※	
Microsoft® Office Pers 2007 搭載モデル※1	なし
英語モデル:	1709GDE
初期導入済OS※2:	Windows® XP Professional 正規版 SP2 (1709GDEは英語版)
認識OS※3:	Microsoft® Windows® 2000 Profes Windows® XP Home Edition 正規版 Windows® XP Professional 正規版 Windows Vista™ Business 正規版
メーカー:	インテル® Core™ 2 Duo プロセッサ
メーカー動作周波数:	2GHz

メーカーのサイト
製品番号, プロセッサ...

ユーザーレビュー

デザイン	★★★★★	4.6
処理速度	★★★★★	5.0
グラフィック性能	★★★★☆	3.4
拡張性	★★★☆☆	4.2
使いやすさ	★★★★☆	4.8
携帯性	★★★★☆	4.7
バッテリー		
液晶		
満足度		

レビューサイト
デザイン, 使いやすさ...

Webサイトの目的により、
属性は大きく異なる

関連研究 (1/2)

属性値抽出手法

- 人手で作成した抽出ルールを利用[Appelt 93][Baumgartner 01]
- 人手で抽出位置にタグ付けをしておき, 機械学習によりルールを学習[Kushmerick 97][Ambite 98][山田 02]
- 少量の値を与えておき, ルールと抽出値を交互に繰り返し学習 (ブートストラッピング) [Riloff 99][Yangarber 02][Ciravegna 04][楠村 07]

多様な形式をもつWeb文書においても
ある程度の精度で抽出できる

関連研究 (2/2)

スキーマ抽出手法

- フォームに隣接するテキストを属性名として抽出
[Raghavan 01][Zhang 04]
- 表の論理的構造を特定することにより属性名を抽出
 - 人手でセルにラベル付けした事例を用いる [Tengli 04]
 - セル中のテキスト間の類似度を用いる [Chen 00]
- クラス語との共起パターン, 共起頻度から属性名を特定
[Tokunaga 05]

本研究の目的

既存のスキーマ抽出手法

- ・ 特定の形式(フォーム、テーブル)のみが対象
- ・ 人手によるラベル付けが必要

Web ページにおける記述の多様性に対応し、
網羅的にスキーマを抽出できない

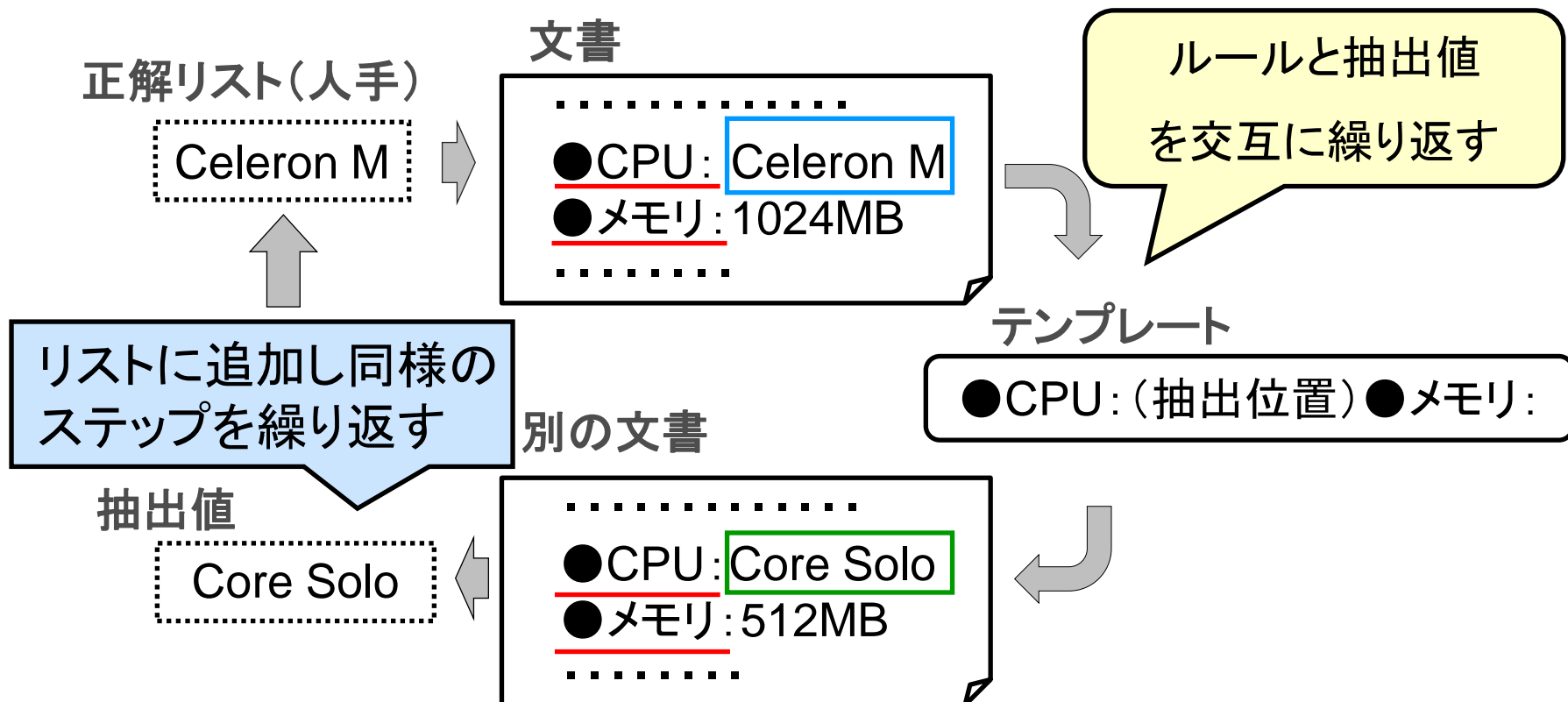


本研究

ブートストラッピングによりスキーマの自動抽出を目指す

ブートストラッピング

人手で属性値を少数与えておき、抽出値とテンプレートを繰り返し学習することで大量の属性値を得る手法



属性名抽出における問題点

属性値抽出の場合

CPU: Celeron M
メモリ: 1024MB



テンプレートに**属性名**が含まれ、多様性は低い

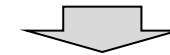


属性名抽出の場合

CPU: Celeron M
メモリ: 1024MB



テンプレートに**属性値**が含まれ、多様性が高い



適合する確率が低くなる

—:テンプレート

属性名抽出の提案方法

2つの属性名の前後にあるテキストで、
完全一致する部分をテンプレートとする

検索語の例

『CPU』 『メモリ』

```
<DIV class=content id=product>
<UL>
<LI> <B> メーカー型番 </B> 170957I
<LI> <B> タイプ </B> B5ノート
<LI> <B> CPU </B> Core 2 Duo
<LI> <B> メモリ </B> 1GB
<LI> <B> HDD </B> 120GB </UL> </DIV>
```

属性名、属性値が
書かれる構造は
同じであることが多い

テンプレート

(抽出位置):

テンプレートをWeb全体には適用せず、
現在のページにのみ適用して属性名を抽出する

スキーマ抽出の問題点

- (1)用語解説のようなページも検索に適合してしまう
- (2)多様なスキーマを抽出できない
- (3)テンプレートがページ中の様々な位置に適合する可能性がある



問題点と解決方法(1/3)

2つの属性名によるページ検索

検索語

『CPU』 『メモリ』

CPUとは

GPUとは、Central Processing Unitの略で「中央処理装置」と呼ばれます。私ですら、これだけでは何のことなのかさっぱりわかりませんでした。

人間でいうと、頭の回転スピードのことですね。CPUのスピードが速いことになり、それだけ早く処理することができるというわけです。

メモリとは

一方、メモリというのは、専門用語でいうと、RAM(Random Access Memory)をあらわしています。

机の上で、といたほうがいいのかもわかりません。

用語解説

製品情報

メーカー	Lenovo (IBM)
型番	1954-2MJ
製品シリーズ	ThinkPad T60 キー
OS	Windows XP Professional
CPU	インテル Core Duo T2300E 1.6GHz
メモリ	512MB (PC2-5300 DDR2 SDRAM) 2GB +512MBメモリ

必要でないページ

必要なページ

用語解説のようなページも検索に適合してしまう



『属性名 + キー』でページを検索

問題点と解決方法 (2/3)

メーカー	Lenovo (IBM)
型番	1954-2MJ
製品シリーズ	ThinkPad T60
OS	Windows XP Professional
CPU	インテル Core Duo T2300E 1.6
メモリ	512MB (PC2-5300 DDR2 SDRAM) 2GB

新たな
属性名

メーカー, 型番,
製品シリーズ, OS

Let's note LIGHT R7

本体サイズ: B5相当

CPU: Core 2 Duo

メモリ: 1GB

画面サイズ: 10.4インチ

画面解像度: 1024x768

HDD: 120GB

対応: 無線LAN

本体サイズ, 画面サイズ,
画面解像度, HDD・対応

異なるキーを用いると多様なスキーマを獲得できる



キーもブートストラッピングで獲得

問題点と解決方法 (3/3)

テンプレートがページ中の様々な位置に適合する可能性がある

テンプレート

(抽出位置):

タグや記号だけから
構成され, 短い

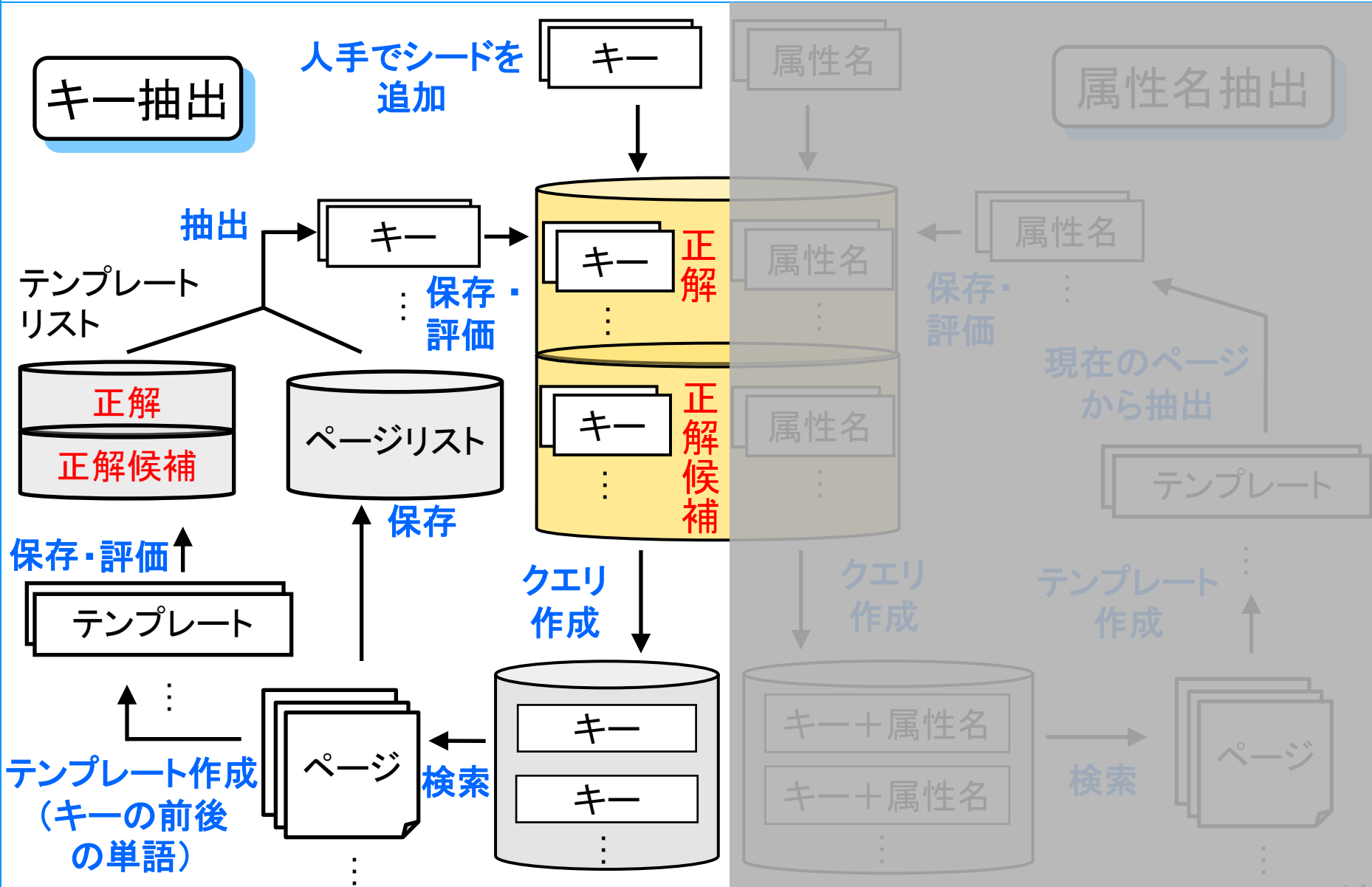
```
<DIV class=content id=product>
<UL>
<LI><B>メーカー:</B>
<LI><B>タイプ:</B>
<LI><B>CPU:</B>Core 2 Duo
<LI><B>メモリ:</B>1GB
<LI><B>HDD:</B>120GB</UL></DIV>
.
<LI><B>関連商品:</B>.....
<LI><B>お問い合わせ:</B>.....
```

この部分だけから抽出したい

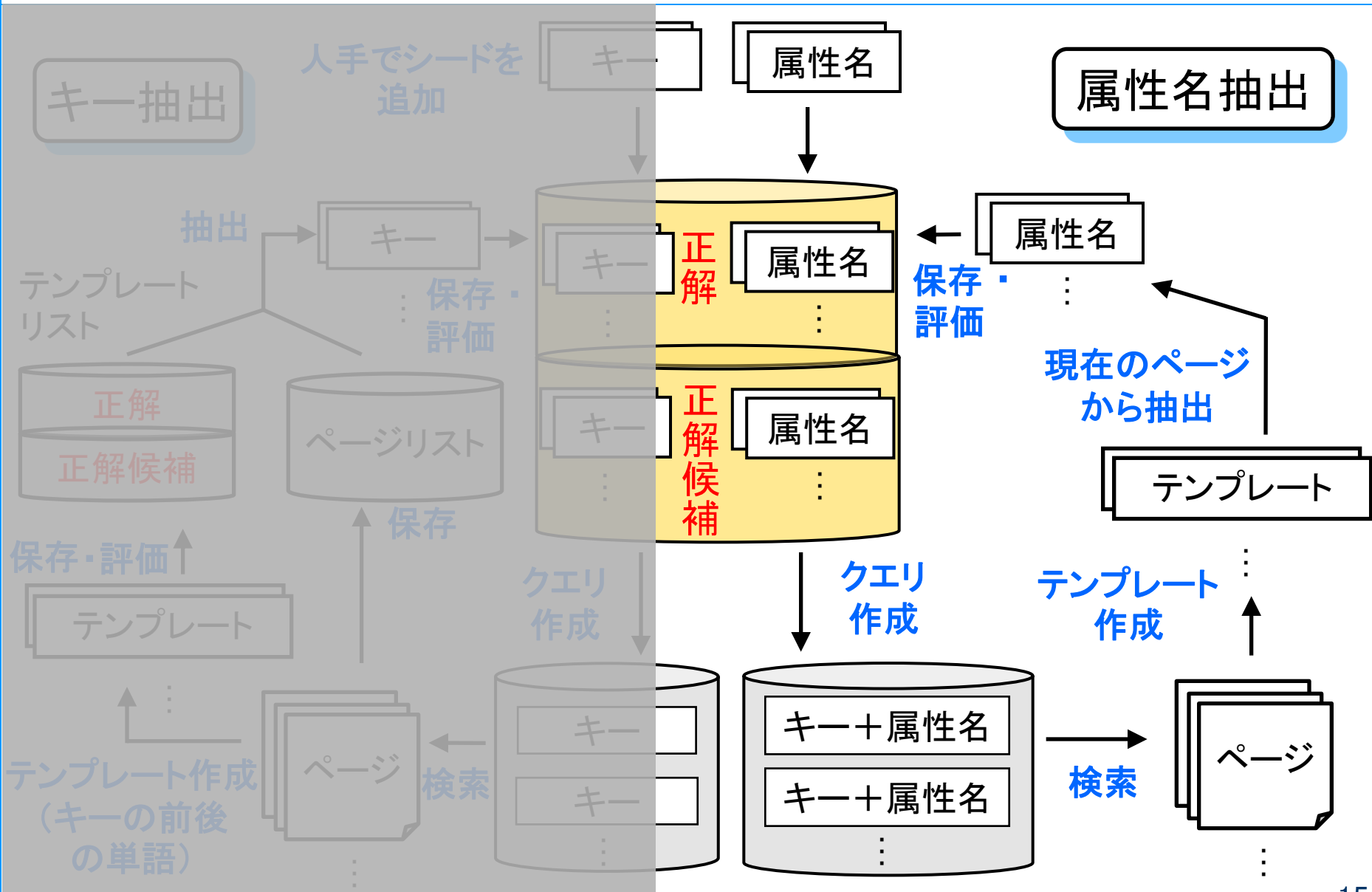
これらにも適合

→ 文書構造を参照し, 抽出を行う範囲を限定

キー抽出処理の流れ



属性名抽出処理の流れ



キー抽出実験

目的 精度・抽出数が最も良くなるキー抽出パラメータの調査

- ・ テンプレート単語数: テンプレートに用いる単語の数
- ・ キー更新数: 各サイクルで正解とするキーの数
- ・ テンプレート更新数: 各サイクルで正解とするテンプレートの数

パラメータ	各パラメータの値 (中央値)					
テンプレート単語数	1	2	3	4	5	6
キー更新数	1	2	5	10	20	50
テンプレート更新数	10	20	50	100	200	500

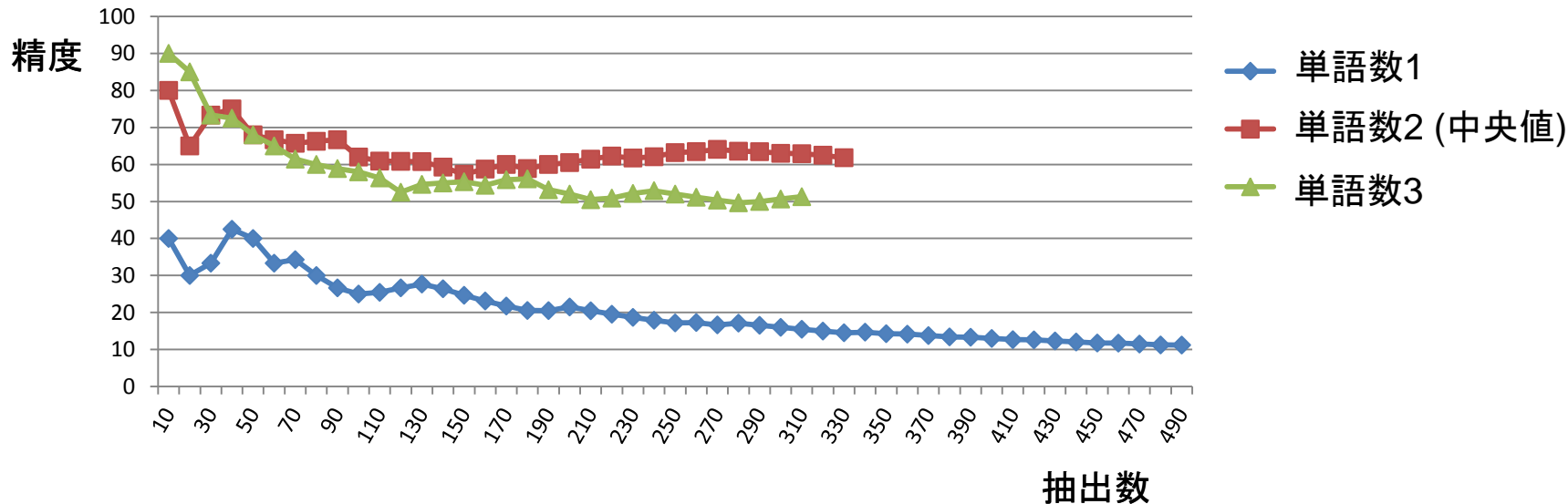
1つのパラメータを変化させて実験を行う
他2つのパラメータは中央値をとる

5万ページで
実験終了



結果(1)

テンプレート単語数を変化させた場合の精度の推移



- 抽出数が30までは、単語数=3が良い
- 5万ページ取得時には、単語数=2が良い

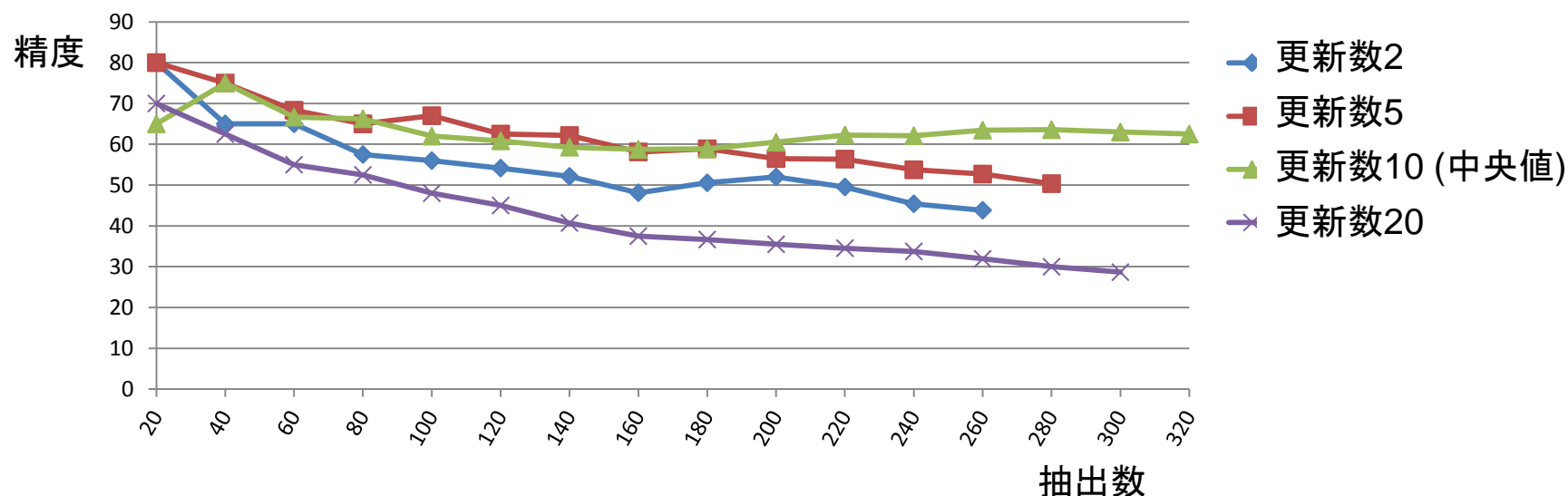


パラメータ: テンプレート単語数
サイクル: 初期3、以降2

が良い

結果(2)

キー更新数を変化させた場合の精度の推移



- サイクルの前半では、更新数=5が最も精度が良い
- サイクルの後半では、更新数=10が最も精度が良い

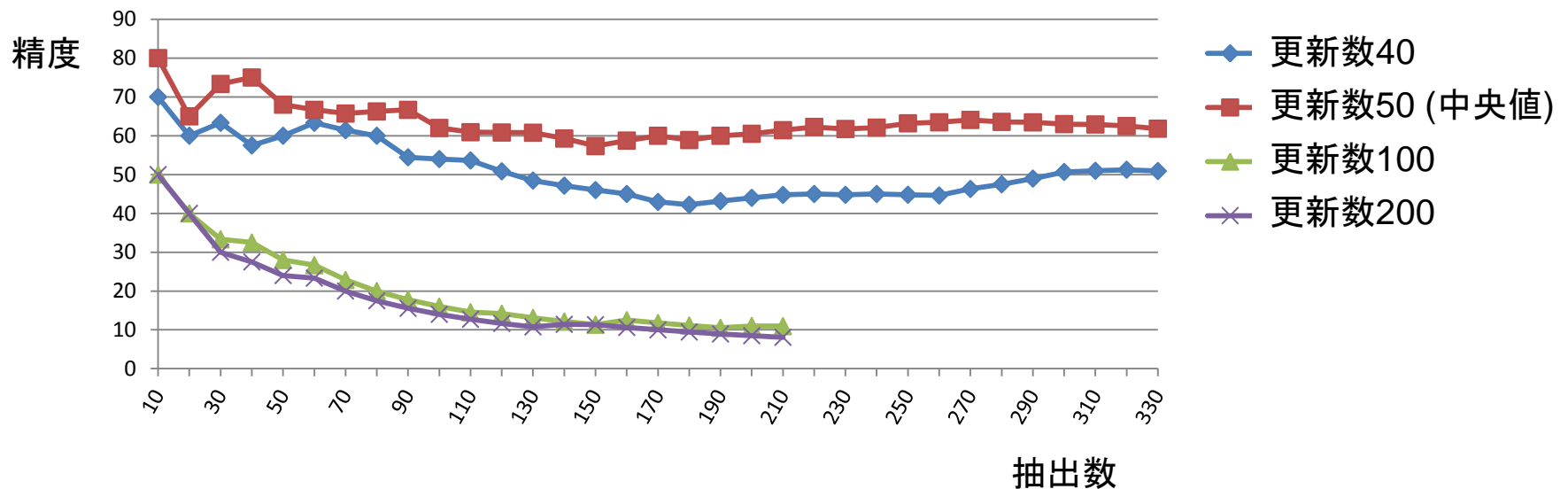


パラメータ: キー更新数
サイクル: 前半5、後半10

が良い

結果(3)

テンプレート更新数を変化させた場合の精度の推移



- 精度は更新数が 50>40>100>200 の順に良い



パラメータ: テンプレート更新数
サイクル: 終始50

が良い

まとめ

- ブートストラッピングを用いたスキーマ抽出手法の提案
 - キー+2つの属性名で検索し、前後テキストの完全一致部分からテンプレートを作成
 - 属性名とキーをそれぞれブートストラッピングで獲得
 - 文書構造を参照し、抽出を行う範囲を限定
- キー抽出パラメータと精度・抽出数の関係を5万ページで調査結果として、パラメータ固定の場合、(2、10、50)が最良
- 今後の予定：
途中再開のプログラムを追加し、実験の開始