

実文書を自然言語処理技術と 適切に繋ぐ技術の重要性

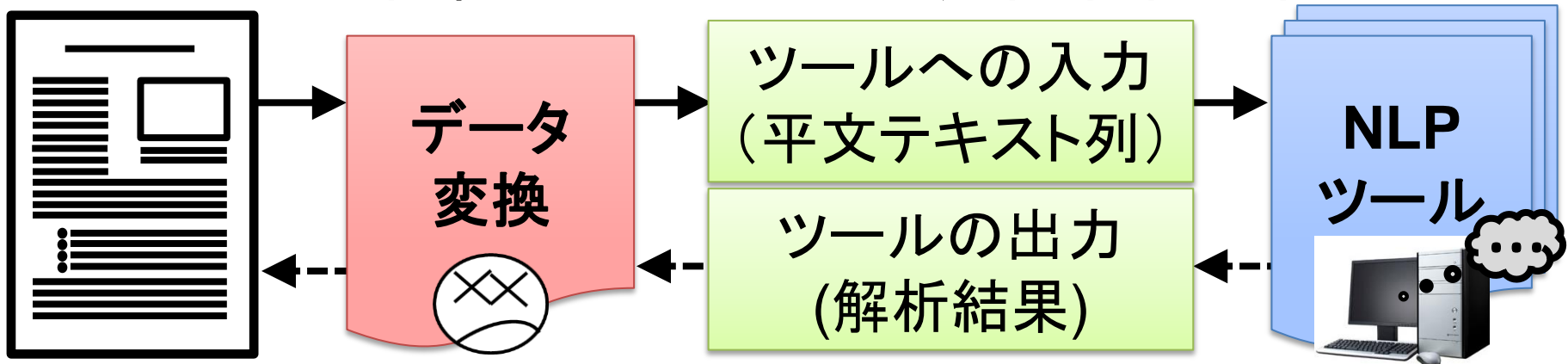
原 忠義 トピチゴラン

宮尾 祐介 相澤 彰子

(国立情報学研究所)

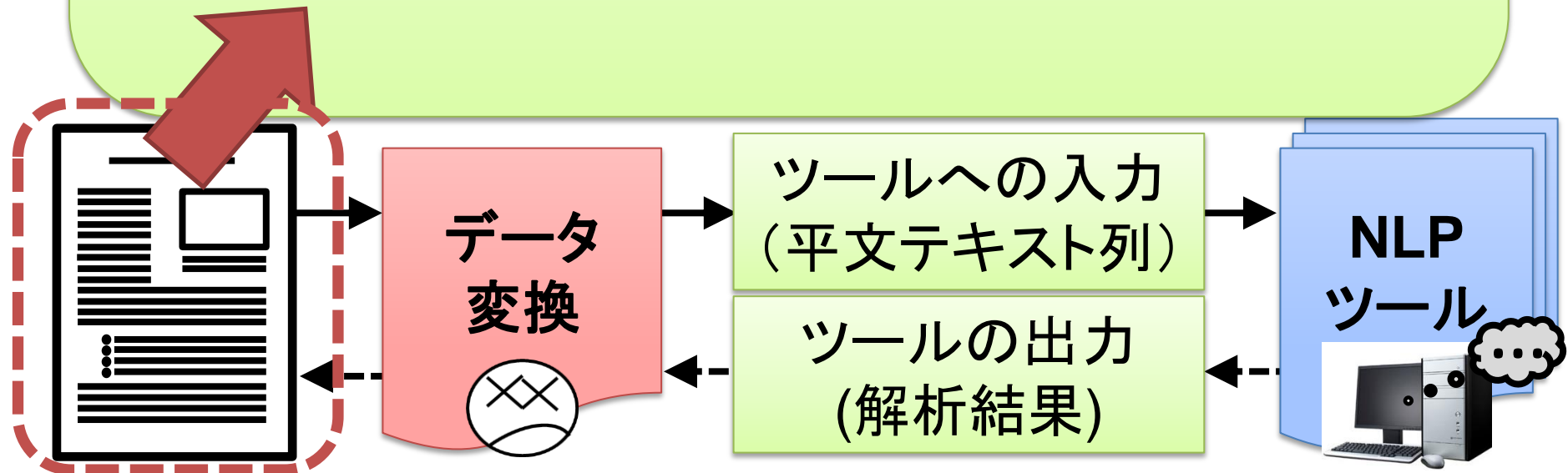
本研究の目的

- 実文書をNLPツールで解析したいが...
 - 実文書テキスト: 様々な形で構造化
 - NLPツール: 多くが平文テキストを入力として想定
- データ変換が必要 → 利用者に委ねられる
 - その都度の変換作業は負担
 - 適切に行わないとツールの効率・性能が低下



本研究の目的

New UI is shown. The UI is more useful than XYZ in [3]*, and
*Notice that



本研究の目的

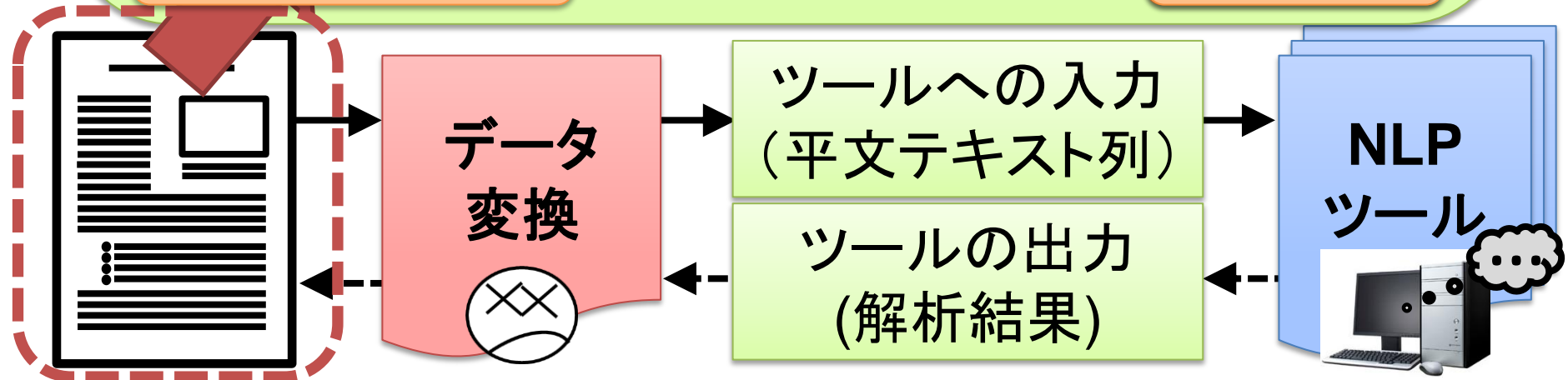
太字, 斜字体

文献参照リンク

New UI is shown. The
UI is more useful than
XYZ in [3]*, and ...
*Notice that ...

索引アンカー

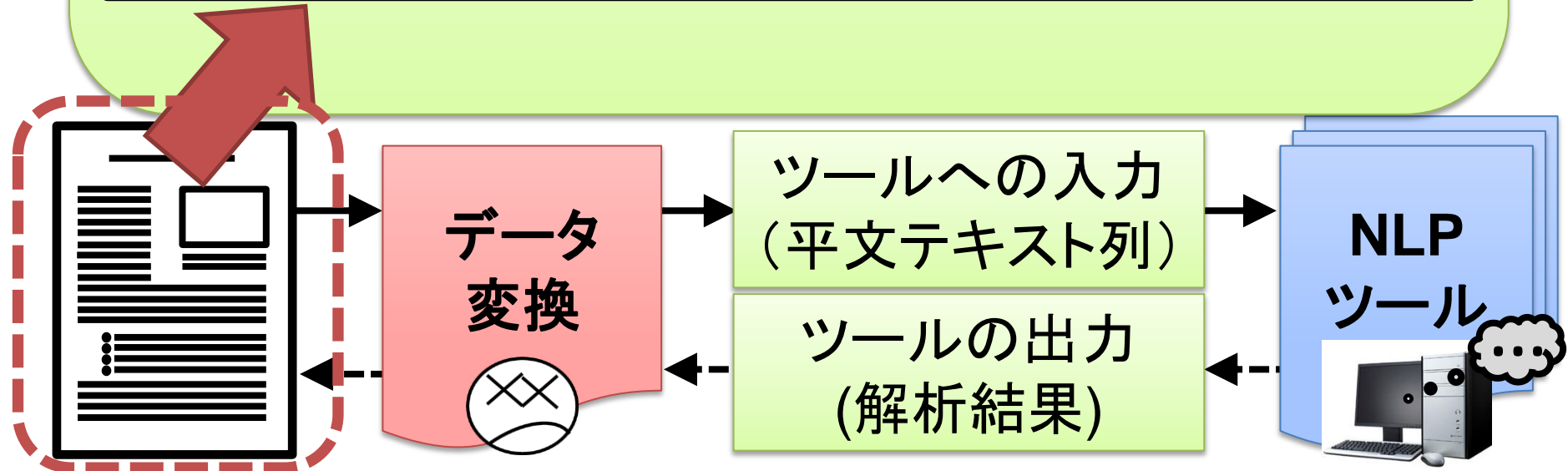
脚注



本研究の目的

タグ付テキスト

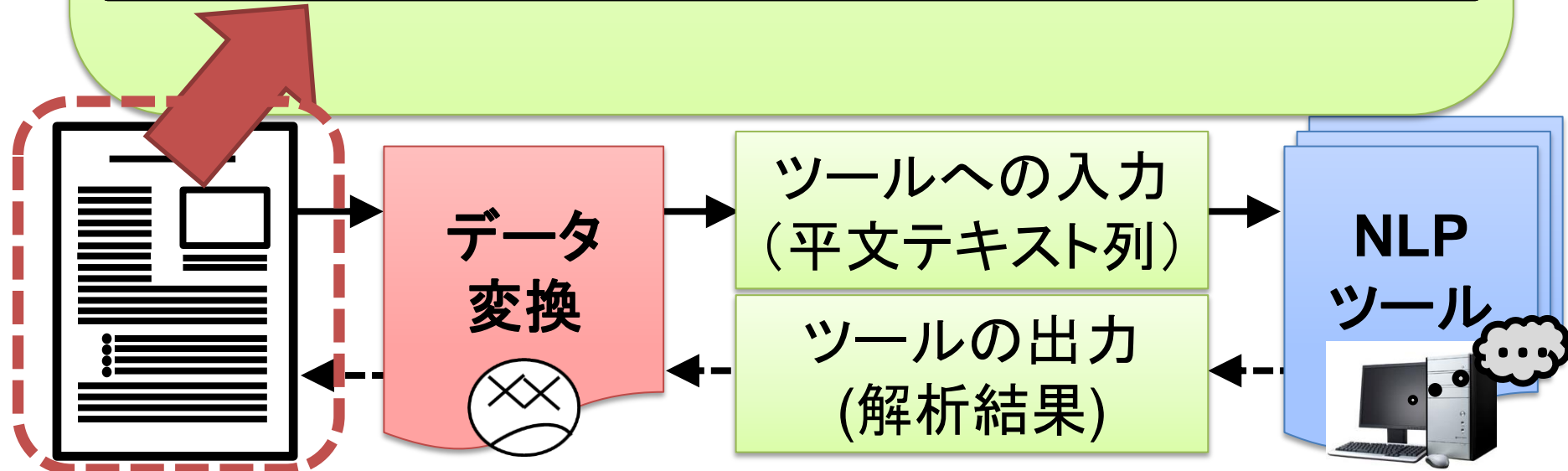
<text>New UI</text> is shown. The UI is more useful than XYZ **<index>##</index>** in **<cite>[...]</cite>** **<note>Notice that ...** **</note>**, and ...



本研究の目的

平文？テキスト

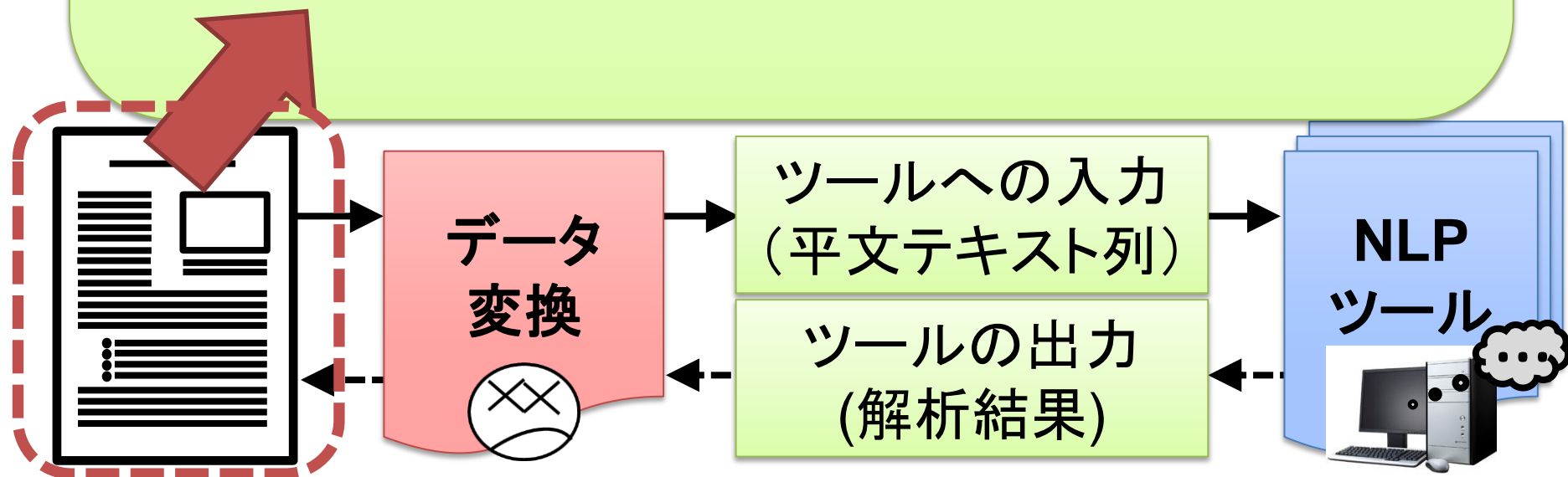
New UI is shown. The UI is
more useful than XYZ ##
in [...] Notice that ...
. , and ...



本研究の目的

平文？テキスト

New UI is shown. The UI is more useful than XYZ## in [...] Notice that, and ...



本研究の目的

平文？テキスト

New UI is shown. The UI is more useful than XYZ## in [...] Notice that ... , and ...

解析対象でない断片

文の埋め込み

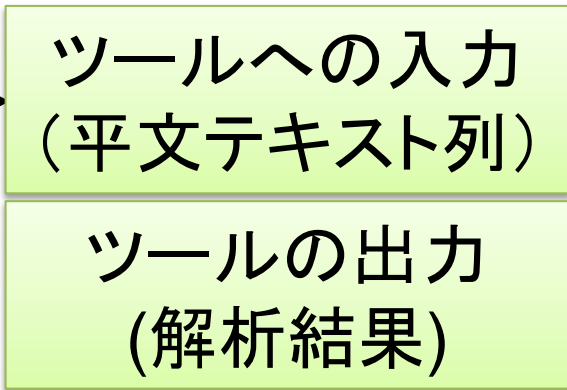
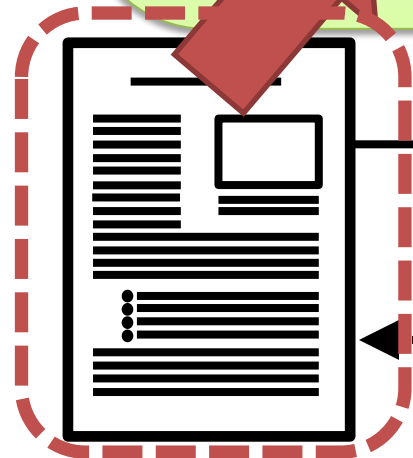
自然言語以外の構造

データ
変換

ツールへの入力
(平文テキスト列)

ツールの出力
(解析結果)

NLP
ツール



本研究の目的

平文？テキスト

New UI is shown. The UI is more useful than XYZ## in [...] Notice that ... , and ...

解析対象でない断片

文の埋め込み

自然言語以外の構造

データ
変換

ツールへの入力
(平文テキスト列)

ツールの出力
(解析結果)

NLP
ツール

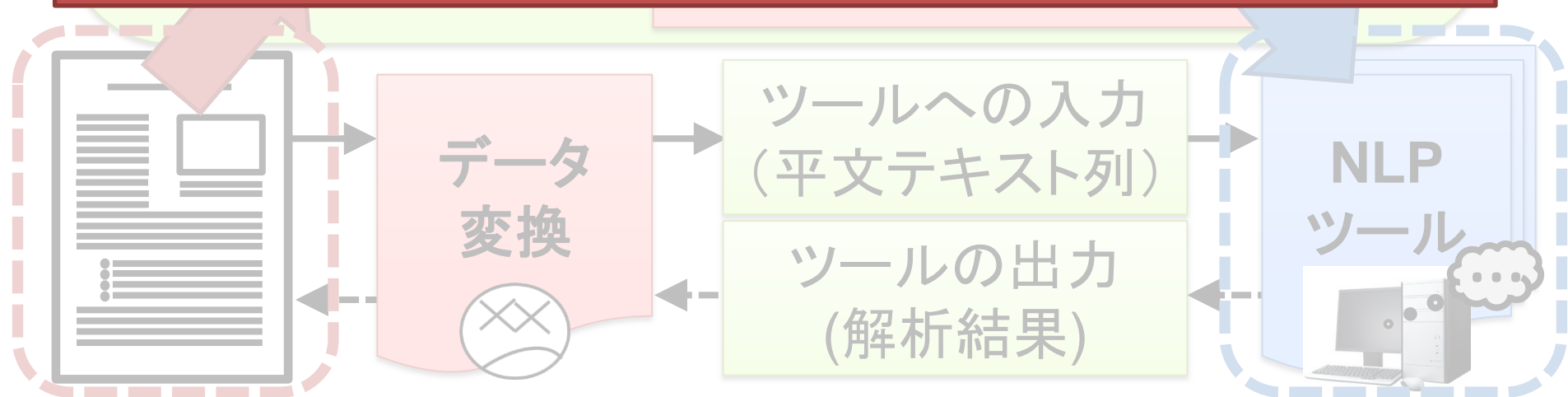


本研究の目的

平文？テキスト

New UI is shown. The UI is more useful than XYZ## in [...] Notice that ... , and ...

適切な変換の重要性を定量的に検証
→ 実文書解析の議論促進へ



本発表の概要

- 関連研究 – データ変換への注目度の低さ
- 研究方針 – XML文書上での重要性の検証
- 提案枠組 – タグ分類に基づく平文テキスト抽出
- 実験
 - 対象文書からの平文テキスト列抽出の実施
 - NLPツールの高被覆・効率的な適用実現の検証
- 議論
 - 実文書をNLP技術と適切に繋ぐ技術の重要性

対象文書⇔NLPツール入出力間の データ変換に関する関連研究

- 統一的手法に取り組んだ研究は見当たらず
- ツール側からの変換ツール提供
 - 構文解析器*^{1,2}: 品詞タグ付文 → タガーを同梱
→ タガー自体も1文毎に区切られた平文を想定
- 様々な文書・ツールでの解析を扱う統合枠組
 - UIMA*³ (その実装*^{4,5,6}) および GATE*⁷
→ NLPツールをシステムに統合する必要あり

*¹ C&C (Clark et al., 2007) *² Enju (Ninomiya et al., 2007) *³ Ferruci et al., 2006

*⁴ RASP4UIMA (Andersen et al., 2008) *⁵ U-compare (Kano et al., 2011)

*⁶ Kachako (Kano, 2012) *⁷ Cunningham et al., 2013

本研究の方針

【目的】適切なデータ変換技術の重要性を検証

【方針1】対象をXML文書に絞る

– 平文を超える構造情報はタグで表されると仮定

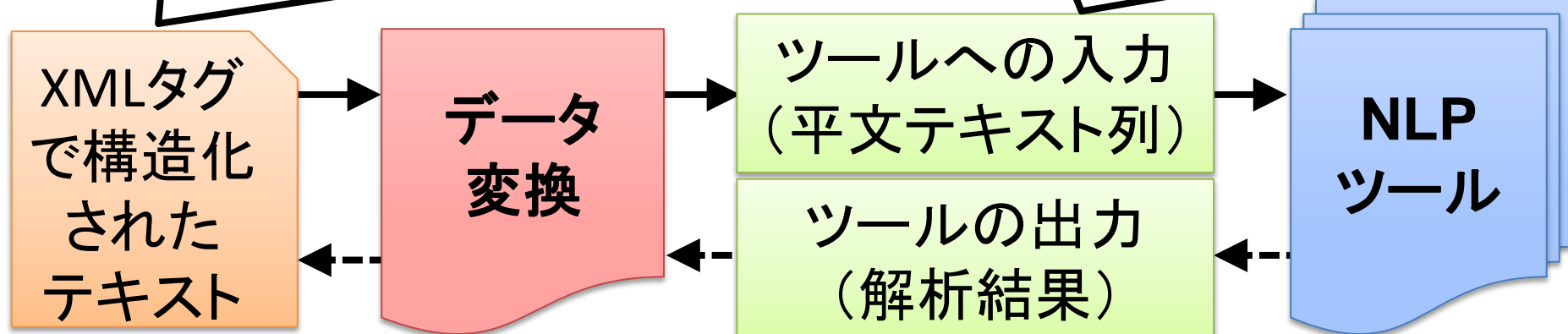
【方針2】既存のNLPツールを、それ自体に手を加えず対象文書に適用可能にする枠組を提案

→有用性の検証から、データ変換の重要性を定量的に示す

提案枠組概要

<p> In our case, we use the CTT (Concur Task Tree) <cite>[<bibref bibrefs="paterno-ctte-2001"/>]</cite>.</p>

In our case, we use the CTT (Concur Task Tree) [1].

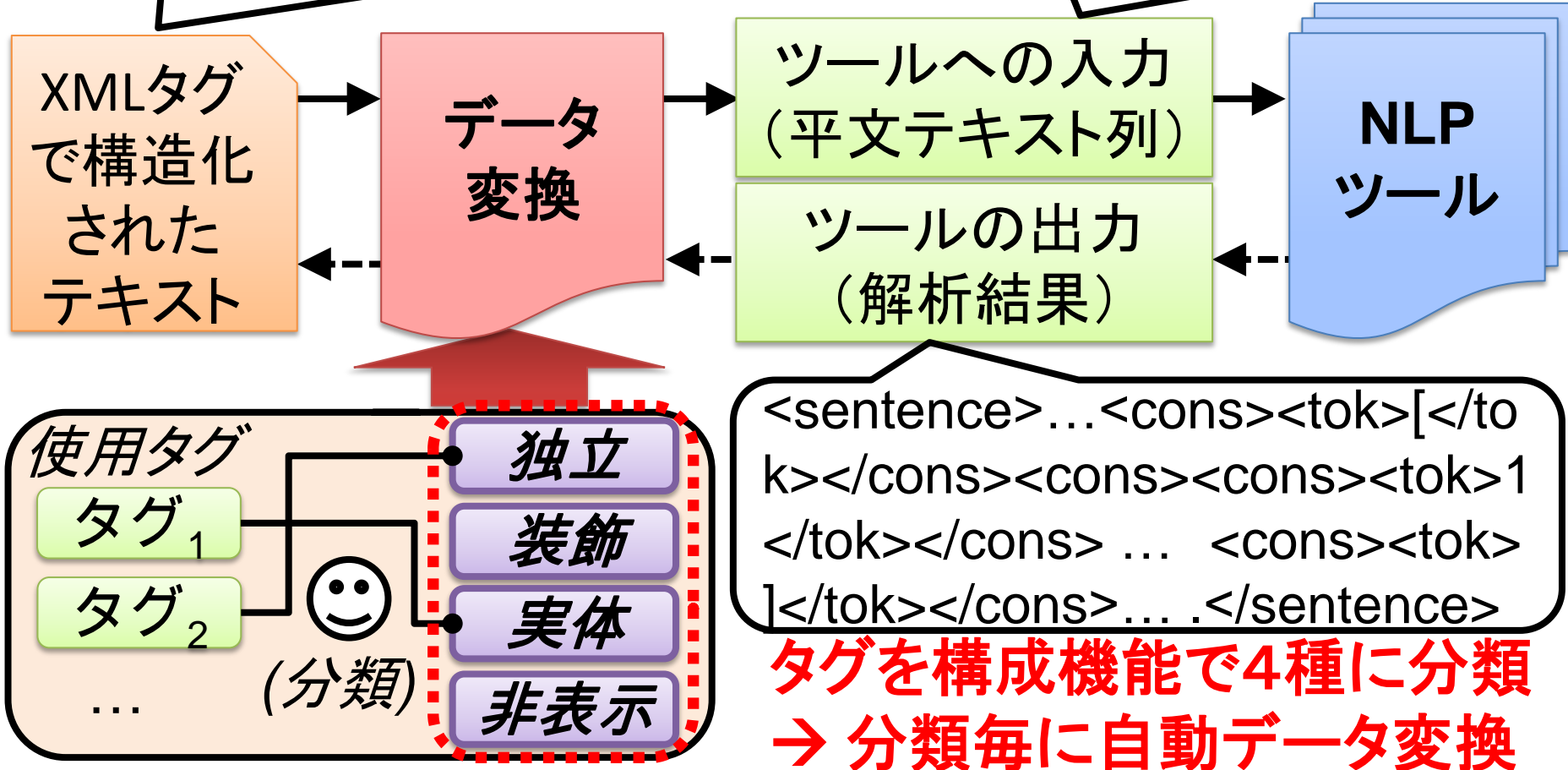


<sentence>...<cons><tok>[</tok></cons><cons><cons><tok>1</tok></cons>... <cons><tok>]</tok></cons>...</sentence>

提案枠組概要

<p> In our case, we use the CTT (Concur Task Tree) <cite>[<bibref bibrefs="paterno-ctte-2001"/>]</cite>.</p>

In our case, we use the CTT (Concur Task Tree) [1].



観察: 4種のテキスト構成機能

タグ領域の表示スタイルを変更

装飾

表示されるテキスト以外の記述

非表示

`<text>New UI</text>` is shown. The UI is more useful than XYZ `<indexmark>...</indexmark>` in `<cite>[...]</cite>` `<note>Notice that ... </note>` and ...

実体

内部が非自然言語の原理で構成されている表示オブジェクト

独立

統語的に独立した領域として表示されるテキスト

NLPツール適用の直感

タグは
無視できそう

装飾

解析対象には
ならなさそう

非表示

`<text>New UI</text>` is shown. The UI is more useful than XYZ `<indexmark>...</indexmark>` in `<cite>[...]</cite>` `<note>Notice that ... </note>` and ...

実体

内部は解析
できなさそう

独立

別個に解析した
方が良さそう

データ変換戦略 (1/3): NLPツール入力への変換

タグを除去

タグとタグ領域を除去

装飾

非表示

`<text>New UI</text>` is shown. The UI is more useful than XYZ `<indexmark>...</indexmark>` in `<cite>[...]</cite>` `<note>Notice that ... </note>` and ...

実体

独立

タグ領域を代替の平文字列と置換

タグ領域を分離

データ変換戦略 (1/3): NLPツール入力への変換

オフセット情報
を保持

タグとタグ領域を除去

非表示

text

New UI is shown. The UI is more useful
than XYZ `<indexmark>...</indexmark>` in
`<cite>[...]</cite>` `<note>Notice that ... </note>` and ...

実体

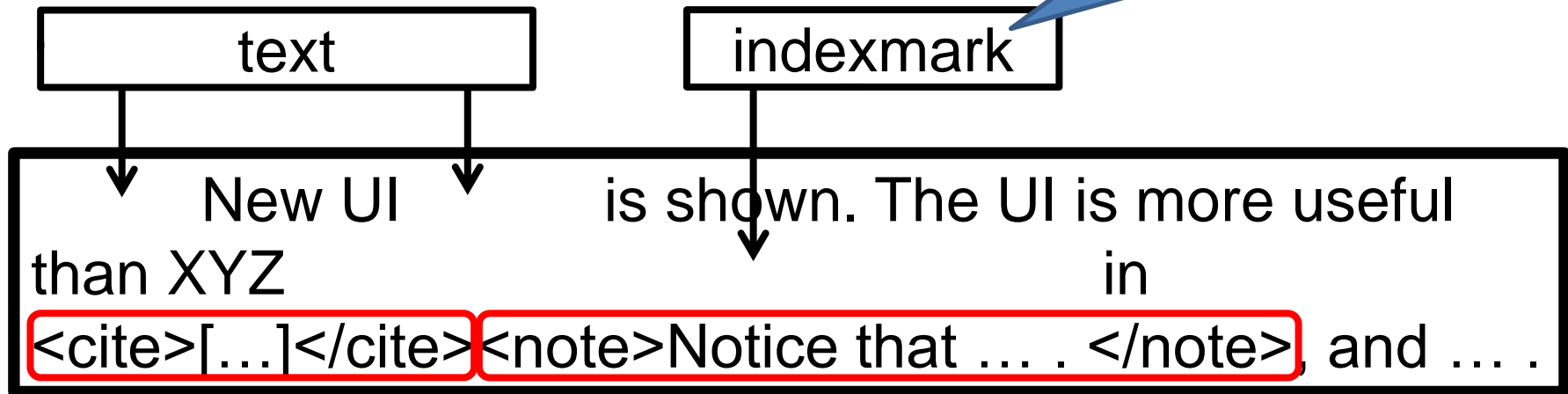
独立

タグ領域を代替の平文字列と置換

タグ領域を分離

データ変換戦略 (1/3): NLPツール入力への変換

オフセット情報
を保持



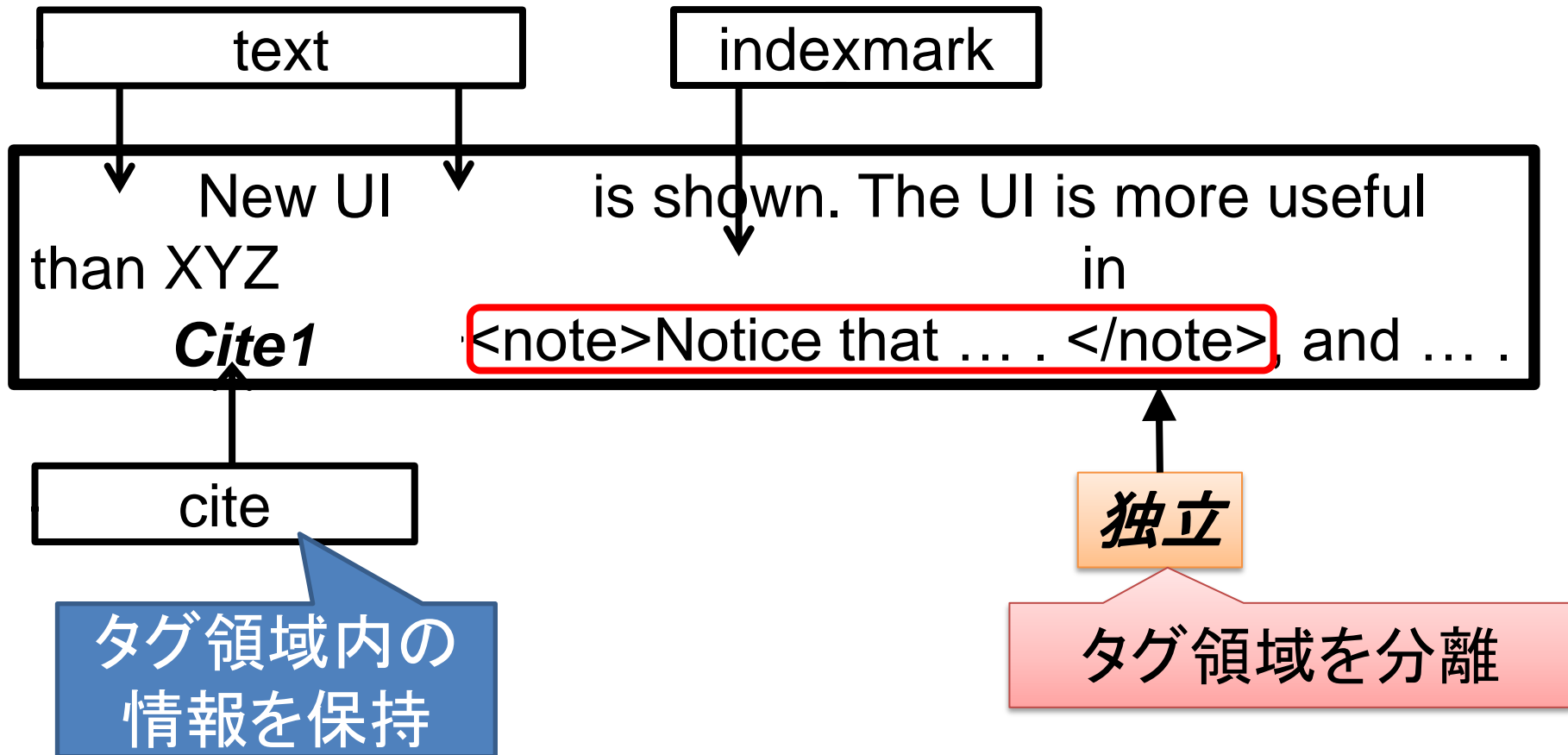
実体

タグ領域を代替の平文字列と置換

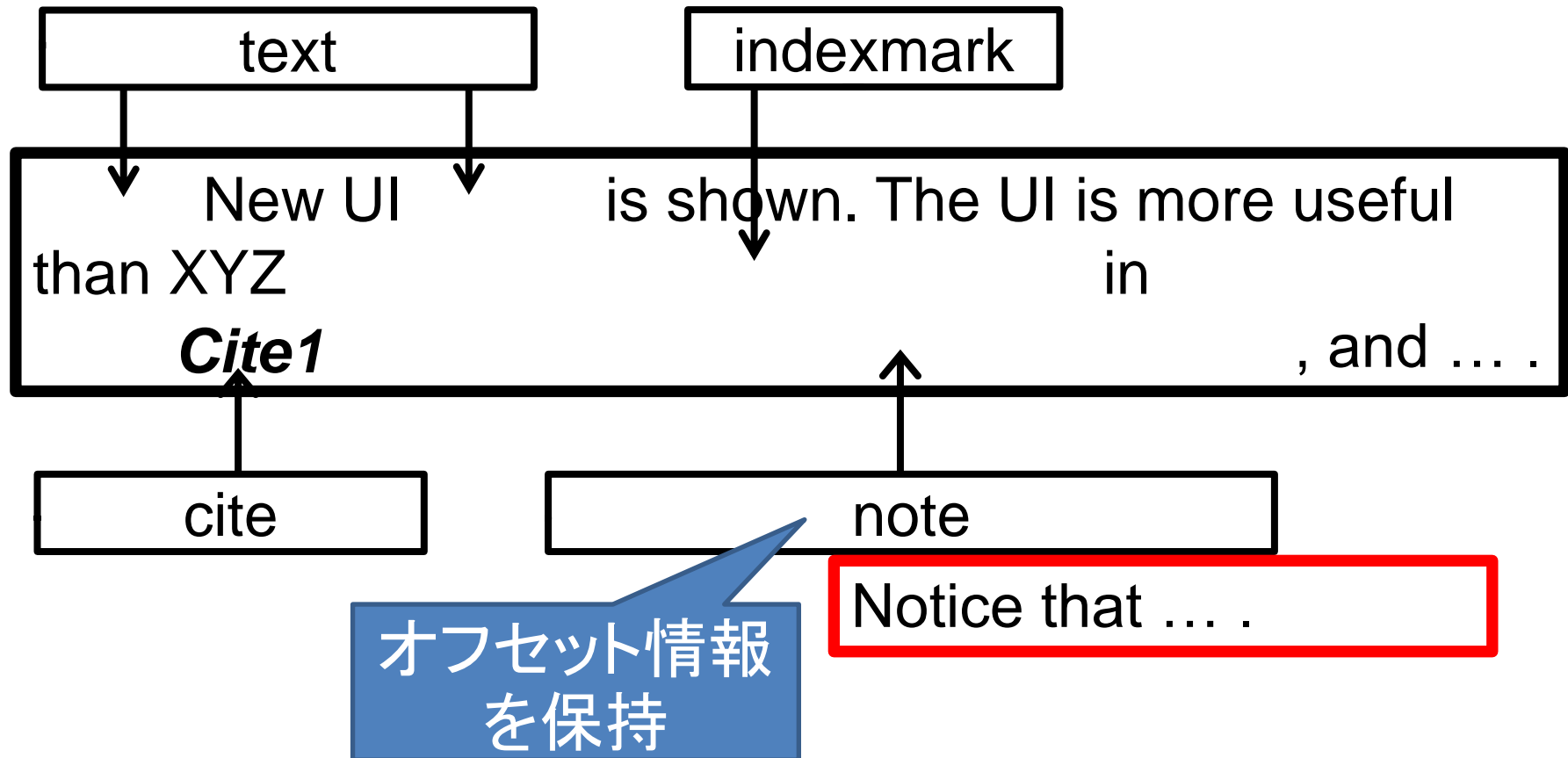
独立

タグ領域を分離

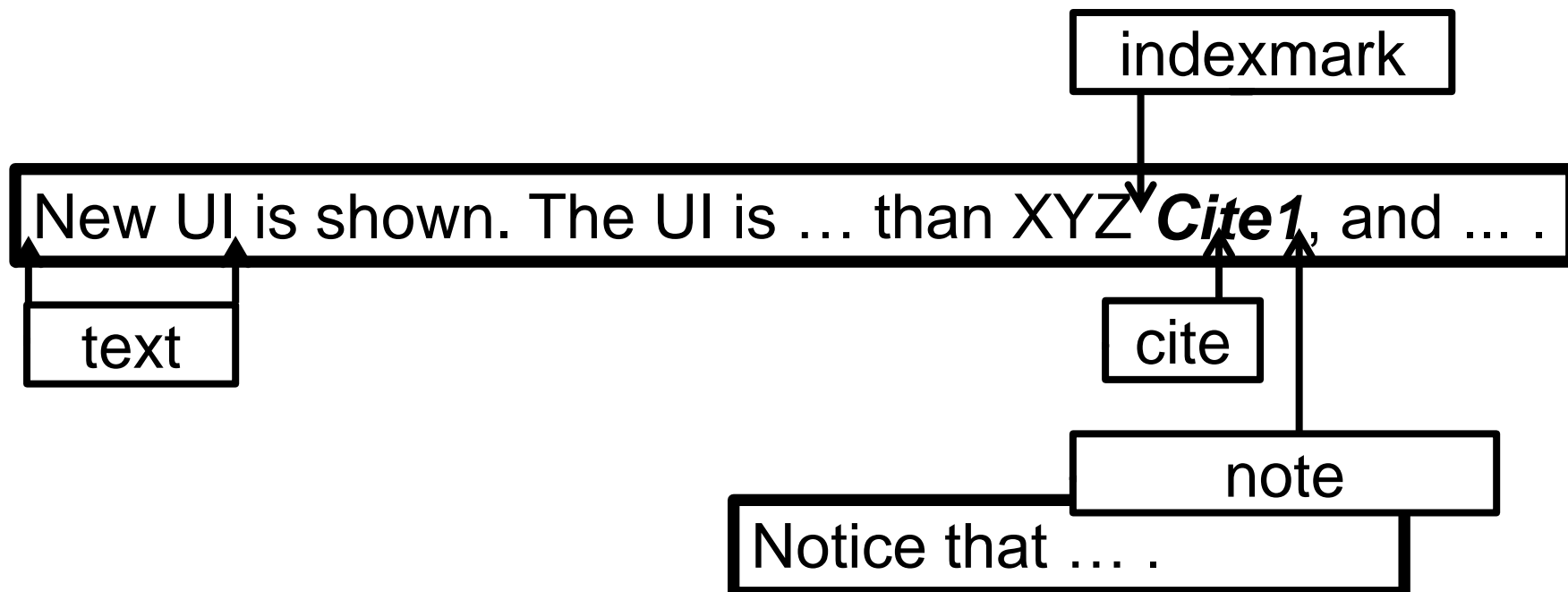
データ変換戦略 (1/3): NLPツール入力への変換



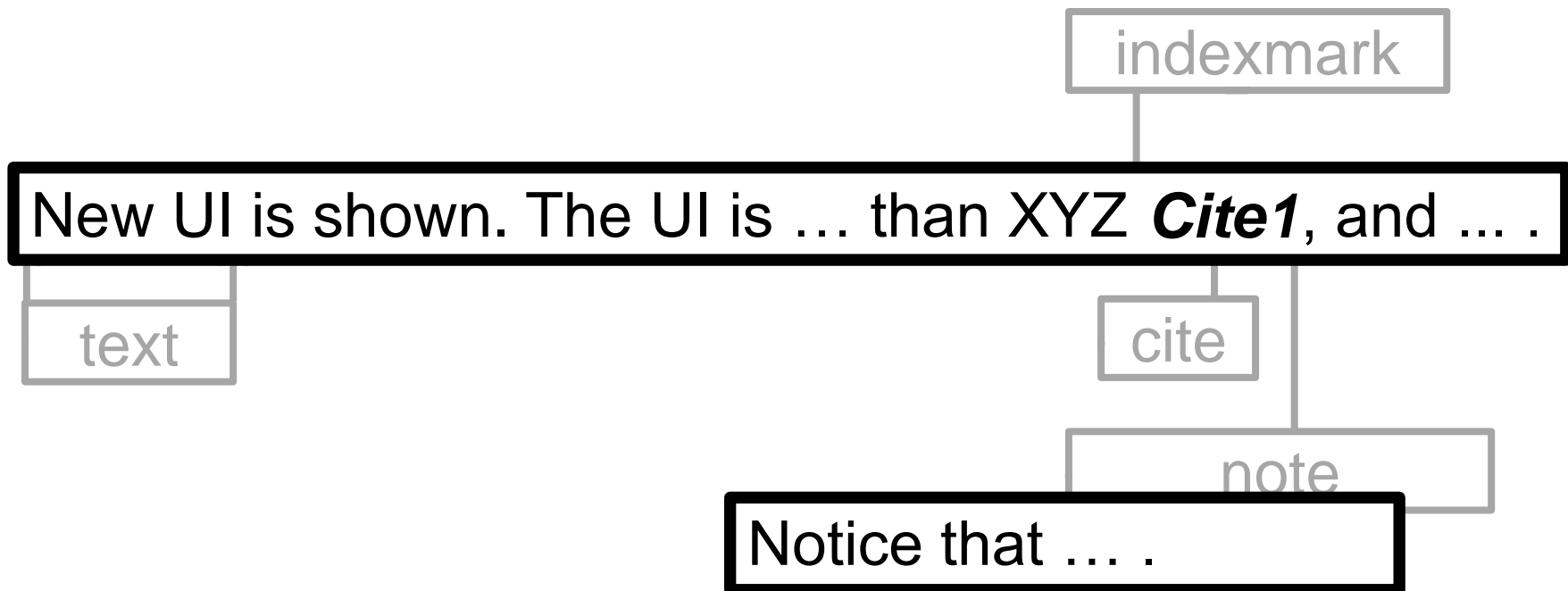
データ変換戦略 (1/3): NLPツール入力への変換



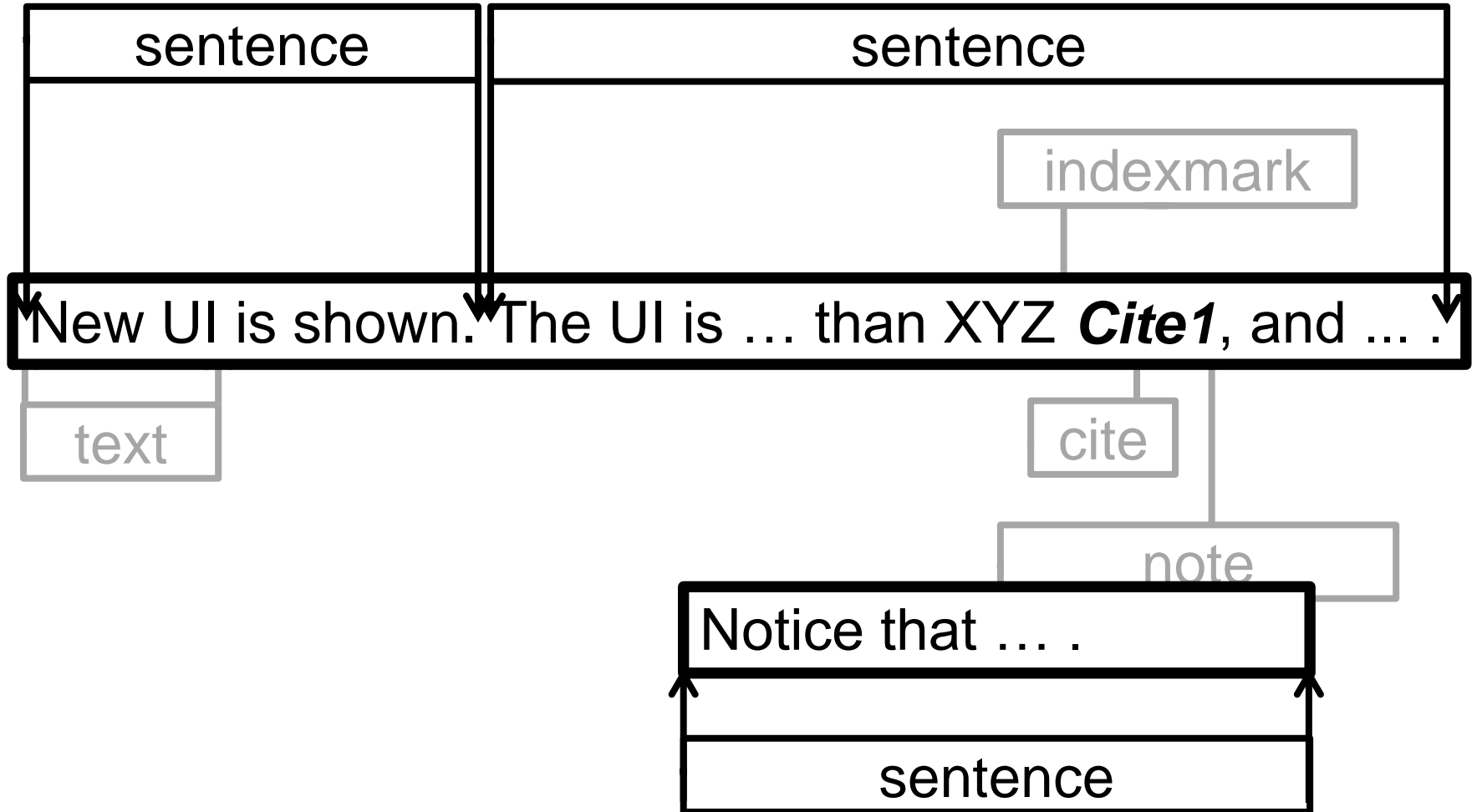
データ変換戦略 (1/3): NLPツール入力への変換



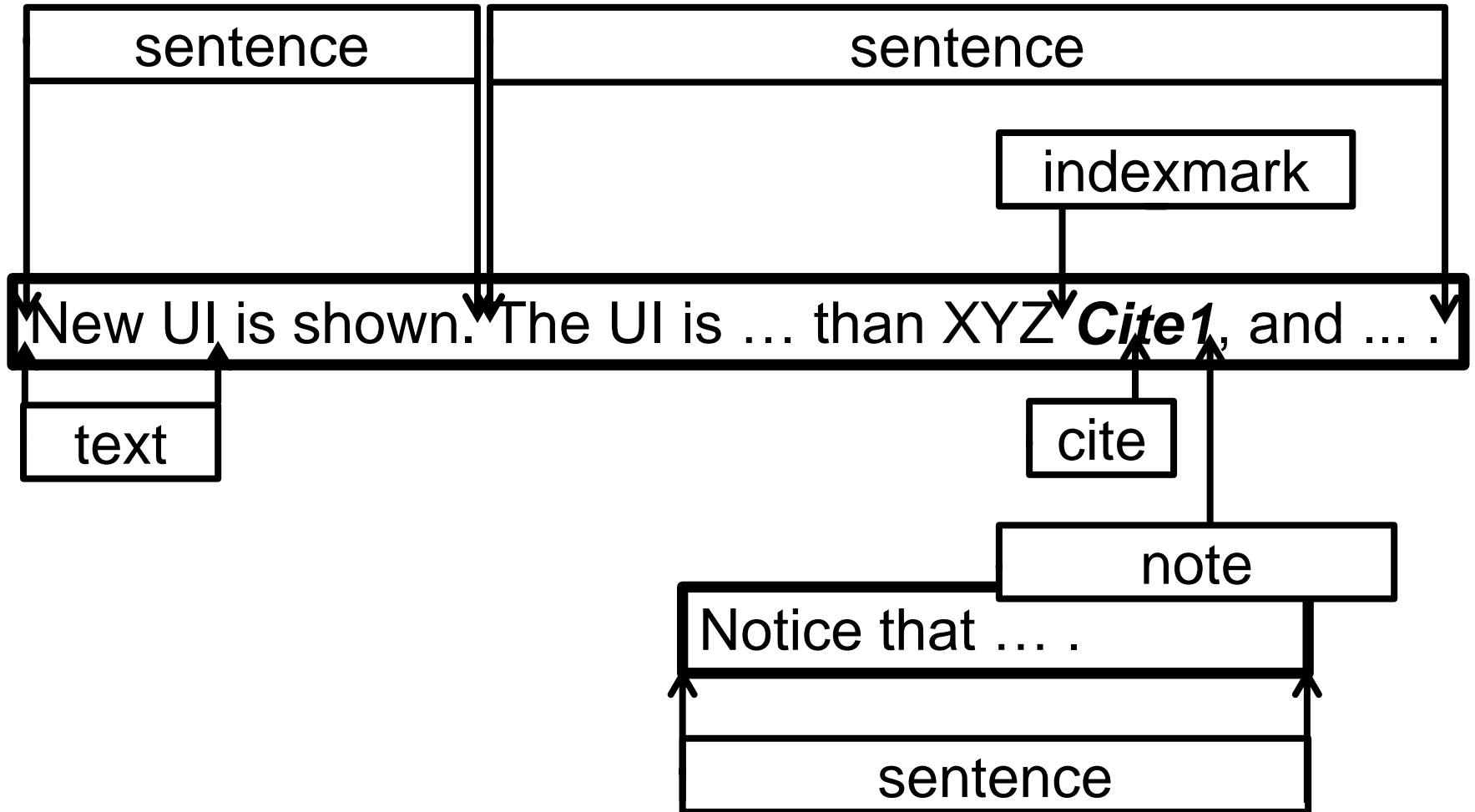
データ変換戦略 (1/3): NLPツール入力への変換



データ変換戦略 (2/3): NLPツールの適用 (文区切り器)



データ変換戦略 (3/3): 元のタグの復帰



データ変換戦略 (3/3): 元のタグの復帰

```
<sentence><text>New UI</text> is shown.</sentence>  
<sentence>The UI is more useful than XYZ<indexmark  
>...</indexmark> in <cite>[...]</cite><note><sentence>  
Notice that ... . </sentence></note>, and ... .</sentenc  
e>
```

タグ分類と変換戦略まとめ

分類	分類基準	変換戦略
独立	周囲のテキストから統語的に独立した領域	(A)を(B)から分離→(A), (B)にツール適用後, (A)を(B)に復帰
装飾	領域内の表示スタイルのみ変更	(A')のみをテキスト(B)から除去→ツール適用後に(A)を(B)に復帰
実体	テキスト構成要素として扱われる最小のオブジェクト単位を表す	(A)を(C)で置換→ツール適用後に(C)を(A)に復元
非表示	表示方法の設定や追加情報を記述する	(A)を(B)から除去→ツール(A)を(B)に復帰

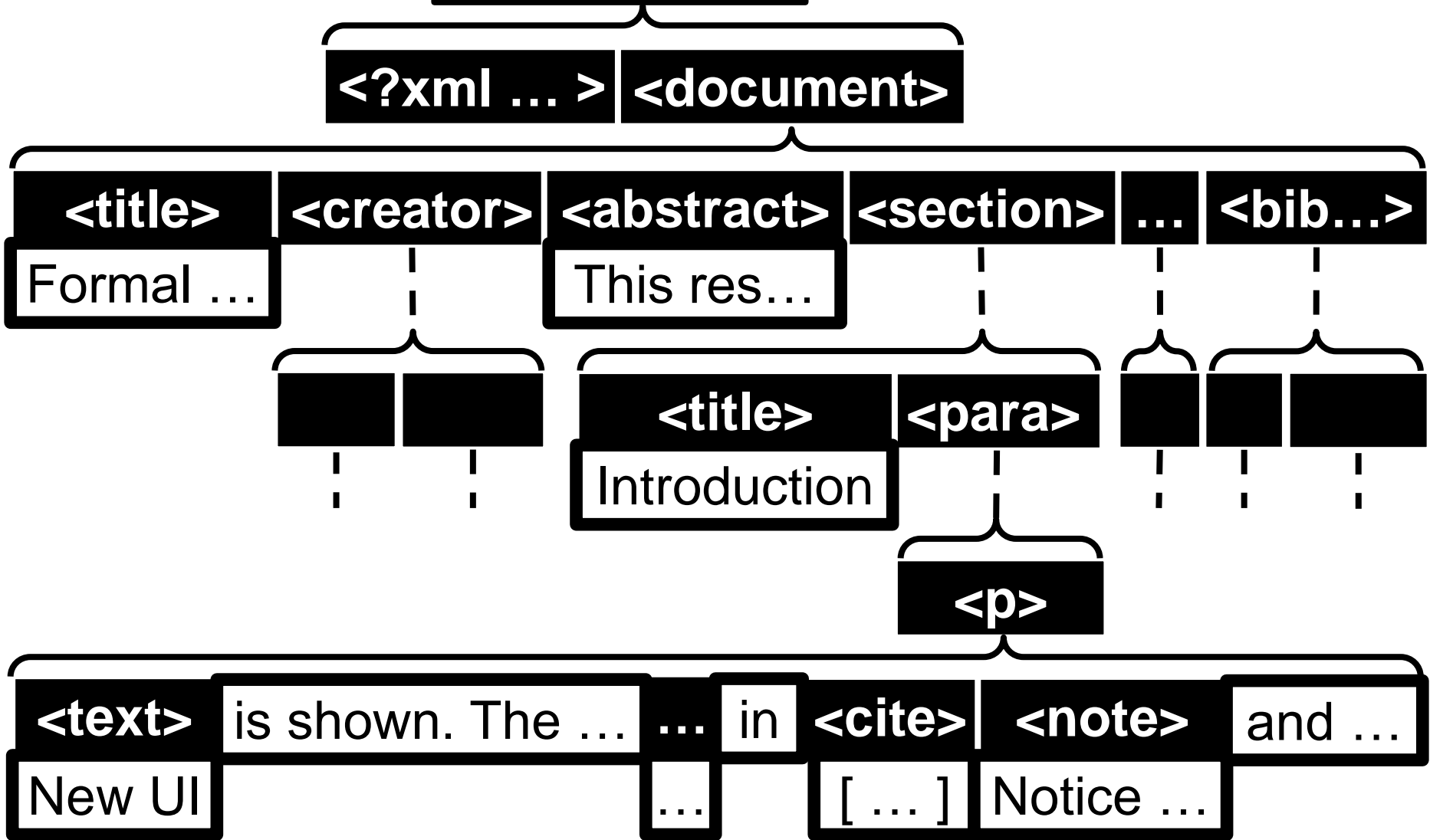
(A): タグとタグ領域 (A'): タグ (B) 元テキスト (C) 代替文字列

タグ分類の効率性

```
<?xml ...>
<document ...>
  <title>Formal approaches ... </title>
  <creator> ... </creator>
  <abstract><p>This research ... </p></abstract>
  <section><title>Introduction</title>
    <para><p><text>New UI</text> is shown. The
UI is more useful than XYZ<indexmark> ...
</indexmark> in <cite>[ ... ]</cite><note>Notice
that ... </note> and ... .</p></para>
  </section>
  <section> ... </section>
  <bibliography> ... </bibliography>
</document>
```

XML文書を直接観察して
タグ分類するのは非効率的

観察: XML文書の多層構造



タグ分類の手続き

- 最上層のタグから, 既に分類済みのタグ領域を展開

XML文書

`<?xml ... >` `<document>`

`<title>`

`<creator>`

`<abstract>`

`<section>`

...

`<bib...>`

Formal ...

This res...

`<title>`

`<para>`

Introduction

タグ分類の手続き

- 最上層のタグから, 既に分類済みのタグ領域を展開

XML文書

`<?xml ... >` `<document>`

`<title>`

Formal ...

`<creator>`

`<abstract>`

This res...

`<section>`

...

`<bib...>`

`<title>`

Introduction

`<para>`

タグ分類の手続き

- 最上層のタグから, 既に分類済みのタグ領域を展開

独立: `<document>`



非表示: `<?xml ... >`

XML文書

`<?xml ... >` `<document>`

`<title>`

`<creator>`

`<abstract>`

`<section>`

...

`<bib...>`

Formal ...

This res...

`<title>`

`<para>`

Introduction

タグ分類の手続き

- 最上層のタグから, 既に分類済みのタグ領域を展開
独立: `<document>`

XML文書

非表示: `<?xml ... >`

`<?xml ... >` `<document>`

`<title>`

`<creator>`

`<abstract>`

`<section>`

...

`<bib...>`

Formal ...

This res...

`<title>`

`<para>`

Introduction

タグ分類の手続き

- 最上層のタグから, 既に分類済みのタグ領域を展開

XML文書



独立: `<document>` `<section>`

`<title>` `<abstract>`

非表示: `<?xml ... >` `<creator>`

`<bibliography>`

`<?xml ... >` `<document>`

`<title>`

`<creator>`

`<abstract>`

`<section>`

...

`<bib...>`

Formal ...

This res...

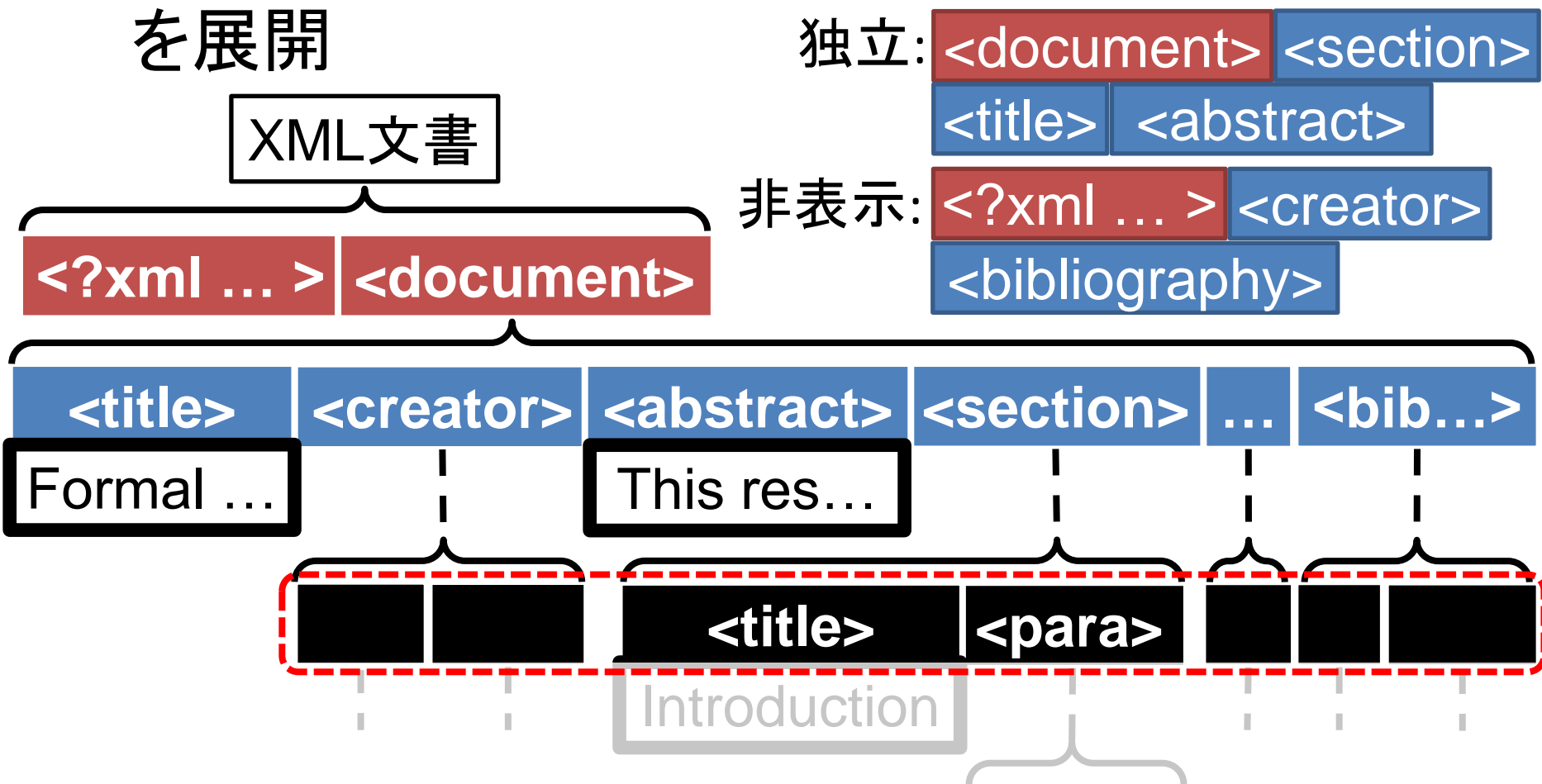
`<title>`

`<para>`

Introduction

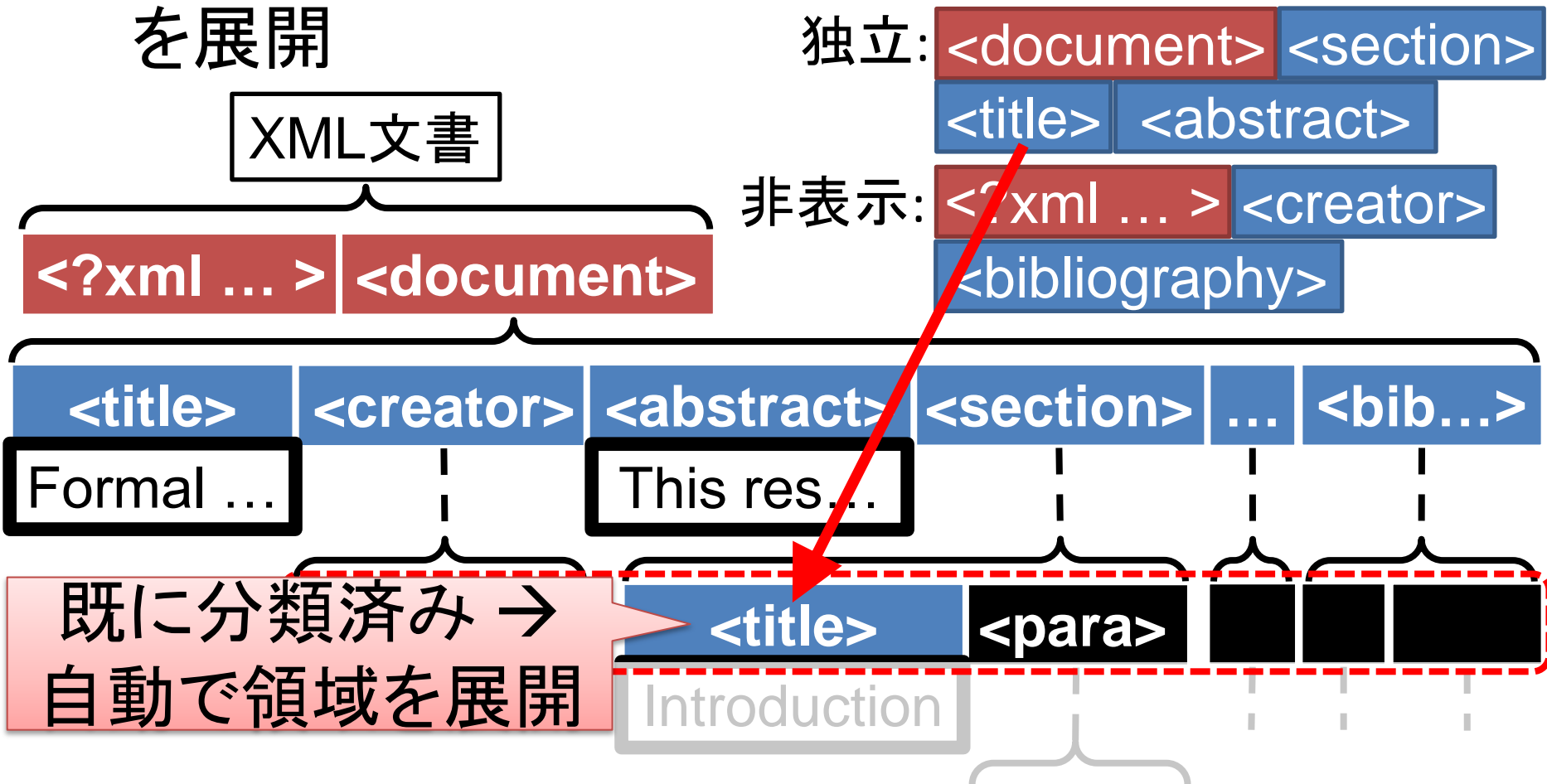
タグ分類の手続き

- 最上層のタグから, 既に分類済みのタグ領域を展開



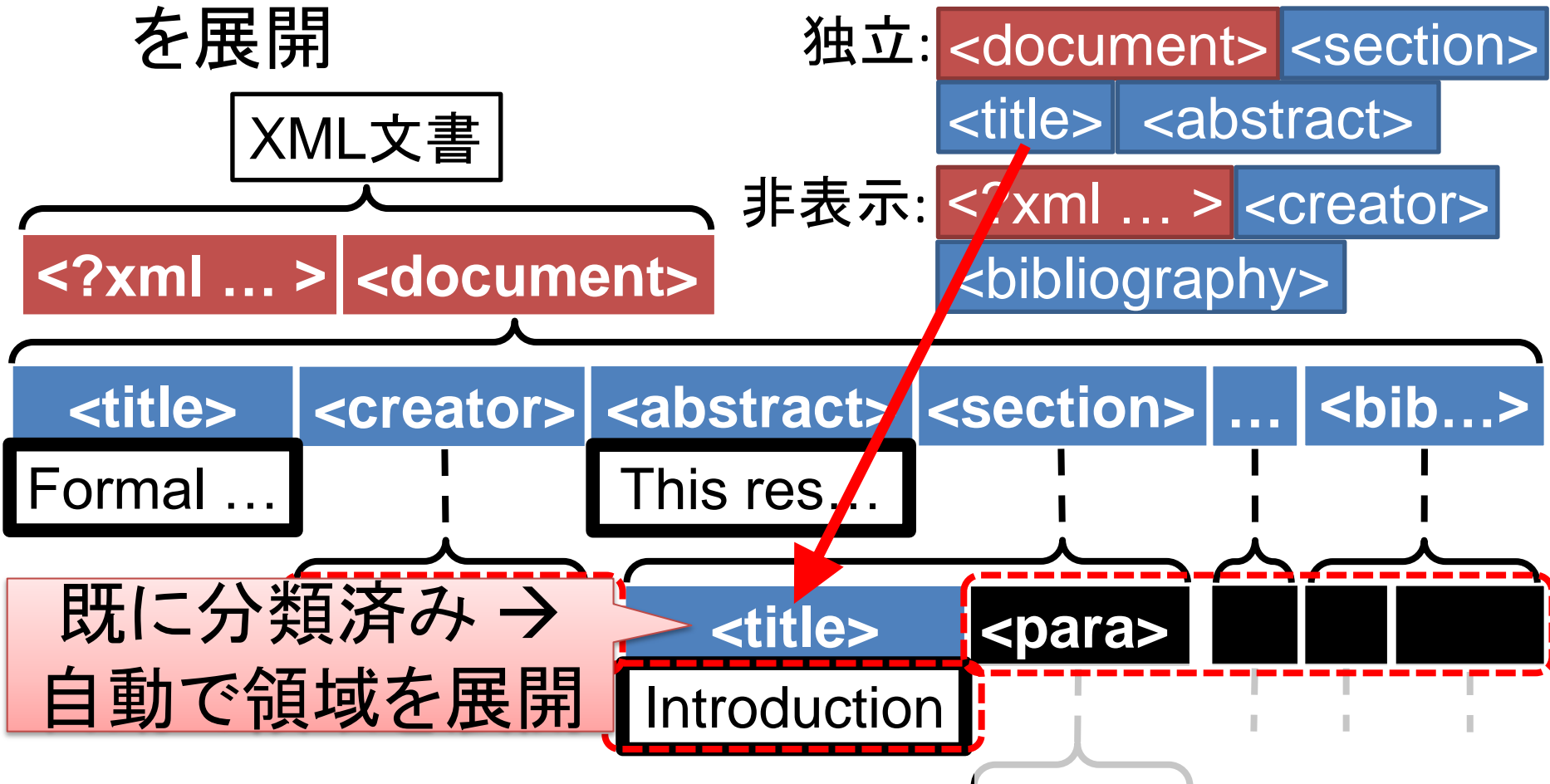
タグ分類の手続き

- 最上層のタグから, 既に分類済みのタグ領域を展開



タグ分類の手続き

- 最上層のタグから, 既に分類済みのタグ領域を展開



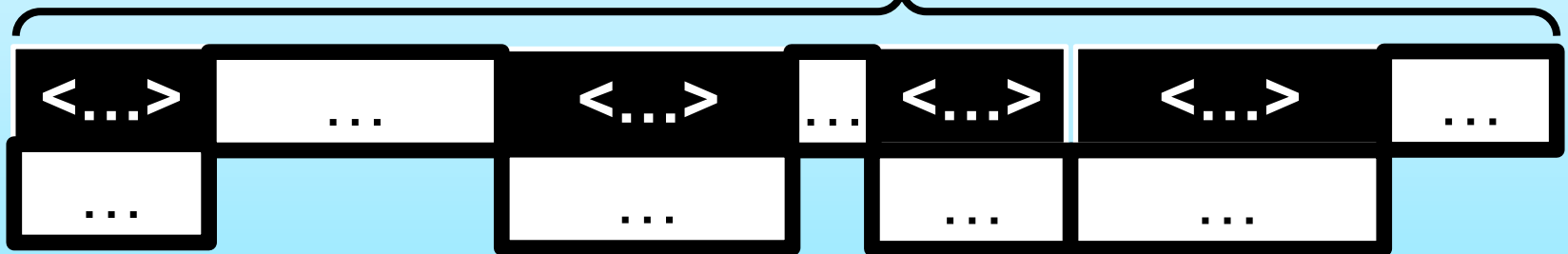
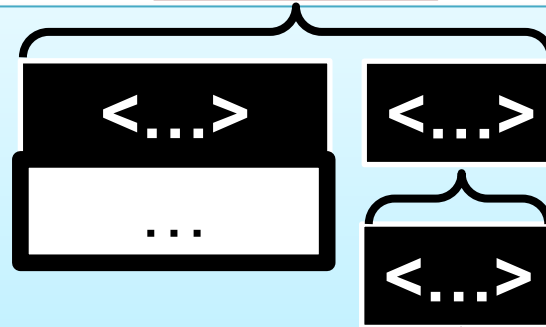
本手続きでもたらされる効率性

- 非表示／実体タグの領域は展開しない
→ 最小限の分類作業に

実体: **<Math>**

展開しない

<Math> (→ 1単語と見なされる)



変換ツールの実装

未分類タグ

Math
equationgroup
quote
text
enumerate
cite
break
ref
tabular
note
graphics
emph

分類操作

独立 >>
装飾 >>
実体 >>
非表示 >>
<< 未分類
追加属性
(1つのみ)

独立タグ

document
title
abstract
section
subsection

実体タグ

装飾タグ

p
para

非表示タグ

bibliography
creator
appendix

このタグ分類を保存してXML文書を変換

登場文脈

[10009_2006_10.xml-format-tag-removed]:

... blems of `<emph>specifying</emph>` contracts, `<emph>monitoring</emph>` their execution for performance, `<note class="footnote" mark="1"><emph>Performance</emph> in contract lingo refers to <emph>compliance</emph> with the <emph>promises</emph> (contractual commitments) stipulated in a contract; nonperformance is also termed <emph>breach of contract</emph>. </note> <emph>analyzing</emph> their ramifications for planning, pricing and other purposes prior to and du ...`

[10009_2006_10.xml-format-tag-removed]:

... `<emph><text>operational</text> semantics</emph>` is ideally suited to alleviating the above problems. `<note class="footnote" mark="2">`Our language is rendered in ordinary linear syntax. but we do not intend to limit the scope

未分類タグ

Math
cite
note
emph
graphics

変換を実行

分類操作

独立 >>

装飾 >>

実体 >>

非表示 >>

独立タグ

title
section

装飾タグ

p
para

実体タグ

非表示タグ

creator
bibliography

登場文脈

未分類タグ

Math
cite
note
emph
graphics

変換を実行

分類操作

独立 >>

装飾 >>

実体 >>

非表示 >>

独立タグ

title
section

装飾タグ

p
para

実体タグ

非表示タグ

creator
bibliography

登場文脈

[10009_2006_10.xml]: ... is ideally suited to alleviating the above problems. **<note class="footnote" mark="2">**Our language is rendered in ordinary linear syntax, but we do not intend to limit the scope of the term “language” to specify linear sequences of characters only, but to include graphical objects and the like. **</note>** Note that contracts are not put to a single use as programs are, ...

未分類タグ

Math
cite
emph
graphics
text

変換を実行

登場文脈

分類操作

独立 >>

装飾 >>

実体 >>

非表示 >>

独立タグ

title
section
note

装飾タグ

p
para

実体タグ

非表示タグ

creator
bibliography

未分類タグ

分類操作

独立 >>

装飾 >>

実体 >>

非表示 >>

独立タグ

title
section
note

装飾タグ

para
emph
text

実体タグ

Math
cite
graphics

非表示タグ

creator
bibliography

変換を実行

登場文脈

実験

- 複数の文書群に対して提案枠組を適用
 - 平文テキストを獲得
 - タグ分類の効率性を検証
- 上記平文テキストに NLP ツールを適用
 - 平文テキスト列を適切に得られることがいかに重要か、単純なアプローチと比較検証

実験設定 (1/3): 対象文書群 (5種)

文書群	ジャンル	言語	スタイル	使用 文書数
PubMed Central (PMC)* ¹	科学論文	英	XML	1,000
arXiv.org* ²	科学論文	英	XHTML	300
Wikipedia* ³ エントリ	ウェブ	英	HTML†	300
言語処理学会論文誌 「自然言語処理」* ⁴	科学論文	日	XHTML††	384(全)
	科学論文	英	XHTML††	68(全)

(† XML ファイルから生成されたもの)

(†† LaTeXML*⁵ を用い XML から変換したものが公開中)

*¹ <http://www.ncbi.nlm.nih.gov/pmc/tools/ftp/> *² <http://arxiv.org/>

*³ <http://www.wikipedia.org/> *⁴ <http://nlp20.nii.ac.jp/resources/>

*⁵ <http://dlmf.nist.gov/LaTeXML/>

実験設定(2/3): NLPツール

- 英語文書：構文解析器(2種)
 - Enju パーザ^{*1} (+ GeniaSS^{*2}): 深い統語／意味解析
 - (検索空間のメモリアーオーバーフロー＝構文解析失敗)
 - Stanford パーザ^{*3}: 句構造／依存関係解析
 - ≤50単語の文を解析(メモリ不足による強制終了対策)
 - (50単語以下 = 成功 / 50単語超 = 自動的に解析失敗)
- 日本語文書：形態素解析器(2種)
 - JUMAN^{*4} / MeCab^{*5}

*1 Enju (Ninomiya et al., 2007) *2 <http://www.nactem.ac.uk/y-matsu/geniass/>

*3 de Marneffe et al., 2006 *4 <http://nlp.ist.i.kyoto-u.ac.jp/index.php?JUMAN>

*5 <http://mecab.sourceforge.jp>

実験設定 (3/3): 比較手法と評価尺度

- 比較手法

- 単純除去: タグを単純除去
- 実/非のみ: 実体・非表示タグ → 提案枠組で処理
装飾・独立タグ → 単純除去
- 提案枠組: 提案枠組で全タグを処理

- 評価尺度

- 検出される文の数
- 解析時間
- 解析失敗文の数(%)

【訂正】表3

「独/非(独立＋非表示)」は
「実/非(実体＋非表示)」の
誤りです (7/1 正誤表公開)

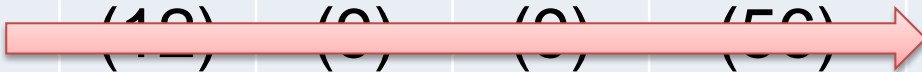
各種文書群に対して分類されたタグと 得られた平文テキスト列

文書 (数)	全タグ数 (異なり数)	分類されたタグの数 (異なり数)					獲得 列数
		独立	装飾	実体	非表示	合計	
PMC (1,000)	1,357 k (421)	32k (12)	62k (9)	48k (9)	34k (56)	177k (85)	26 k
arXiv (300)	1,969 k (210)	6k (15*)	47k (12*)	60k (8*)	8k (17*)	121k (52*)	4 k
Wiki. (300)	224 k (60)	3k (12*)	11 k (8*)	1k (28*)	11k (67*)	28k (115*)	2 k
JNL-E (68)	142 k (57)	8k (25*)	12 k (16*)	6k (9)	23k (19)	29k (69*)	6 k
JNL-J (384)	699 k (58)	50k (23*)	56 k (18*)	32k (10*)	14k (21*)	153k (72*)	38 k

各種文書群に対して分類されたタグと 得られた平文テキスト列

文書 (数)	全タグ数 (異なり数)	分類されたタグの数 (異なり数)					獲得 列数
		独立	装飾	実体	非表示	合計	
PMC (1,000)	1,357 k (421)	32k (12)	32k (9)	10k (9)	34k (50)	177k (85)	26 k
arXiv (300)	1,969 k (210)	6k (15*)	47k (12*)	60k (8*)	8k (17*)	121k (52*)	4 k
Wiki. (300)	224 k (60)	3k (12*)	11 k (8*)	1k (28*)	11k (67*)	28k (115*)	2 k
JNL-E (68)	142 k (57)	8k (25*)	12 k (16*)	6k (9)	23k (19)	29k (69*)	6 k
JNL-J (384)	699 k (58)	50k (23*)	56 k (18*)	32k (10*)	14k (21*)	153k (72*)	38 k

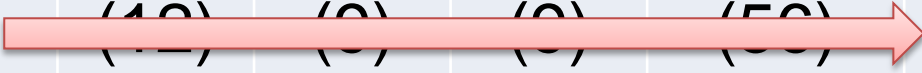
1/5 以下のみ分類



各種文書群に対して分類されたタグと 得られた平文テキスト列

文書 (数)	全タグ数 (異なり数)	分類されたタグの数 (異なり数)					獲得 列数
		独立	装飾	実体	非表示	合計	
PMC (1,000)	1,357 k (421)	32k (12)	32k (9)	10k (9)	34k (50)	177k (85)	26 k
arXiv (300)	1,969 k (210)	6k (15*)	47k (12*)	60k (8*)	8k (17*)	121k (52*)	4 k
Wiki. (300)	224 k (60)	3k (12*)	11 k (9)	11k (10)	11k (67*)	28k (115*)	2 k
JI (6)						29k (69*)	6 k
JI (384)						153k (72*)	38 k

1/5 以下のみ分類



* XHTML, HTML:
タグ名のみでは情報が抽象的
→ 属性・属性値とのペアで1つのタグ名

各種文書群に対して分類されたタグと 得られた平文テキスト列

文書 (数)	全タグ数 (異なり数)	分類されたタグの数 (異なり数)					獲得 列数
		独立	装飾	実体	非表示	合計	
PMC (1,000)	1,357 k (421)	52k (12)	52k (9)	10k (9)	34k (50)	177k (85)	26 k
arXiv (300)	1,969 k (210)	3k (15*)	17k (12*)	33k (8*)	3k (17*)	121k (52*)	4 k
Wiki. (300)	224 k (60)	3k (15*)	11k (12*)	1k (9)	11k (17*)	28k (15*)	2 k
JNL-E (68)	142 k (57)	3k (25*)	12k (16*)	3k (9)	23k (19)	29k (69*)	6 k
JNL-J (384)	699 k (58)	53k (23*)	53k (18*)	32k (10*)	11k (21*)	153k (72*)	38 k

1/5 以下のみ分類

出現数の 20% 以下のみ注目

各種文書群に対して分類されたタグと 得られた平文テキスト列

文書 (数)	全タグ数 (異なり数)	分類されたタグの数 (異なり数)					獲得 列数
		独立	装飾	実体	非表示	合計	
PMC (1,000)	1,357 k (421)	32k (12)	32k (9)	10k (9)	34k (50)	177k (85)	26 k
arXiv (300)	1,969 k (210)	3k (15*)	17k (12*)	33k (8*)	3k (17*)	121k (52*)	4 k
Wiki. (300)	224 k (60)	3k (15*)	11k (10*)	1k (1)	11k (10*)	28k (15*)	2 k
JNL-E (68)	142 k (57)	3k (25*)	12k (10*)	3k (9)	23k (10)	29k (60*)	6 k
JNL-J (384)	33k (50)	25k (25)	10k (10)	10k (10)	21k (21)	33k (12*)	38 k

1/5 以下のみ分類

出現数の 20% 以下のみ注目

実体・非表示タグで囲まれた領域の内部が
それ以上展開されなかったため

各種文書群に対して分類されたタグと 得られた平文テキスト列

文書 (数)	全タグ数 (異なり数)	分類されたタグの数 (異なり数)					獲得 列数
		独立	装飾	実体	非表示	合計	
PMC (1,000)	1,357 k (421)	32k (12)	62k (9)	48k (9)	34k (56)	177k (85)	26 k
arXiv (200)	1,969 k (210)	6k (15*)	47k (10*)	60k (9*)	8k (17*)	121k (50*)	4 k
(100)	(57)	(25)	(10)	(9)	(19)	(69*)	2 k
(100)	(57)	(25)	(10)	(9)	(19)	(69*)	6 k
JNL-J (384)	699 k (58)	50k (23*)	56 k (18*)	32k (10*)	14k (21*)	153k (72*)	38 k

文書をランダムに選択し観察 →

- ・平文テキストがNLPツールに直接入力可能
- ・文書全体を被覆可能な「適切な」文から構成されていることを確認

タグの扱い方が Enju パーザ の 性能に与える影響

文書(数)	手法	検出文数	解析時間(秒)	失敗文数(%)
PMC (1,000)	単純除去	159,327	209,783	4,721 (2.96)
	実/非のみ	112,285	135,752	810 (0.72)
	提案枠組	126,215	132,250	699 (0.55)
arXiv (300)	単純除去	74,762	108,831	2,047 (2.74)
	実/非のみ	41,265	89,200	411 (1.00)
	提案枠組	43,208	87,952	348 (0.81)
Wiki. (300)	単純除去	10,561	14,704	1,161 (10.99)
	実/非のみ	5,026	6,743	67 (1.33)
	提案枠組	6,893	6,058	61 (0.88)
JNL-E (68)	単純除去	23,196	24,881	271 (1.17)
	実/非のみ	15,606	21,304	183 (1.17)
	提案枠組	17,929	18,683	50 (0.28)

タグの扱い方が **Standordパーザ** の 性能に与える影響

文書(数)	手法	検出文数	解析時間(秒)	失敗文数(%)
PMC (1,000)	単純除去	170,999	58,865	18,621(10.89)
	実/非のみ	126,176	50,741	11,881 (9.42)
	提案枠組	139,805	63,295	11,338 (8.11)
arXiv (300)	単純除去	75,672	27,970	10,590(13.99)
	実/非のみ	48,666	24,630	5,457(11.21)
	提案枠組	50,504	26,360	5,345(10.58)
Wiki. (300)	単純除去	14,883	3,114	1,651(11.09)
	実/非のみ	6,173	2,248	282 (4.57)
	提案枠組	8,049	2,451	258 (3.21)
JNL-E (68)	単純除去	24,942	9,069	1,577 (6.32)
	実/非のみ	17,572	7,865	1,058 (6.02)
	提案枠組	19,925	10,154	892 (4.48)

タグの扱い方が Enju パーザ の 性能に与える影響

文書(数)	手法	検出文数	解析時間(秒)	失敗文数(%)
PMC	単純除去	159,327	209,783	4,721 (2.96)
	実/非のみ	135,752	135,752	810 (0.72)
	提案枠組	132,250	132,250	699 (0.55)
(300)	単純除去	108,831	108,831	2,047 (2.74)
	実/非のみ	89,200	89,200	411 (1.00)
	提案枠組	43,208	87,952	348 (0.81)
Wiki.	単純除去	10,561	14,704	1,161 (10.99)
	実/非のみ	5,026	6,745	67 (1.33)
	提案枠組	6,893	6,058	61 (0.88)
JNL-E	単純除去	23,196	24,881	271 (1.17)
	実/非のみ	15,606	21,306	183 (1.17)
	提案枠組	17,929	18,683	50 (0.28)

提案枠組: 1%未満
解析失敗 1~10% 減
= 文書被覆率の大幅改善

タグの扱い方が Enju パーザ の 性能に与える影響

文書(数)	手法	検出文数	解析時間(秒)	失敗文数(%)
PMC (300)	単純除去	159,327	209,783	4,721 (2.96)
	実/非のみ	135,752	135,752	810 (0.72)
	提案枠組	132,250	132,250	699 (0.55)
	実/非のみ	108,831	108,831	2,047 (2.74)
	提案枠組	89,200	89,200	411 (1.00)
JNL-E (68)	単純除去	23,196	24,881	271 (1.17)
	実/非のみ	15,606	21,306	183 (1.17)
	提案枠組	17,929	18,683	50 (0.28)
	実/非のみ	14,704	14,704	1,161 (10.99)
	提案枠組	6,743	6,743	67 (1.33)
	提案枠組	6,058	6,058	61 (0.88)

提案枠組: 1%未満
解析失敗 1~10% 減
= 文書被覆率の大幅改善

提案枠組:
解析時間 25~59% 減
= 大幅短縮を達成

タグの扱い方が Standord パーザ の 性能に与える影響

文書(数)	手法	検出文数	解析時間(秒)	失敗文数(%)
PMC (300)	単純除去	170,999	58,865	18,621(10.89)
	提案枠組	50,741	11,881	(9.42)
	実/非のみ	63,295	11,338	(8.11)
	提案枠組 or 実/非のみ	27,970	10,590	(13.99)
	単純除去	40,000	24,630	5,457(11.21)
	提案枠組	50,504	26,360	5,345(10.58)
JNL-E (68)	単純除去	24,942	9,069	1,577 (6.32)
	実/非のみ	17,572	7,865	1,058 (6.02)
	提案枠組	19,925	10,154	892 (4.48)
	提案枠組 or 実/非のみ			

提案枠組:
解析失敗 2~8% 減
= 文書被覆率の大幅改善

提案枠組 or 実/非のみ:
解析時間 12~28% 減
= 大幅短縮を達成

タグを単純に除去すると、構文解析器は以下により混乱

- 文の埋め込み挿入（独立タグ）
- 自然言語ではない要素で構成された表現（実体タグ）
- 文の表示に直接関係しないシーケンス（非表示タグ）

PMC (300)	単純除去	170,999	58,865	18,621 (10.89)
	実/非のみ	40,000	50,741	11,881 (9.42)
	提案枠組	50,504	63,295	11,338 (8.11)
	実/非のみ	30,000	27,970	10,590 (13.99)
	提案枠組	40,000	24,630	5,457 (11.21)
	実/非のみ	50,504	26,360	5,345 (10.58)
	提案枠組	30,000	3,114	1,651 (11.09)
	実/非のみ	30,000	2,240	282 (4.57)
	提案枠組	30,000	2,451	258 (3.21)
	JNL-E (68)	単純除去	24,942	9,069
実/非のみ	17,572	7,865	1,058 (6.02)	
提案枠組	19,925	10,154	892 (4.48)	

提案枠組:

解析失敗 2~8% 減

= 文書被覆率の大幅改善

提案枠組 or 実/非のみ:

解析時間 12~28% 減

= 大幅短縮を達成

タグを単純に除去すると、構文解析器は以下により混乱

- 文の埋め込み挿入（独立タグ）
- 自然言語ではない要素で構成された表現（実体タグ）
- 文の表示に直接関係しないシーケンス（非表示タグ）

PMC (1,000)	単純除去	170,999		58,865	18,621 (10.89)
	実/非のみ	126,176	↓	50,741	11,881 (9.42)
	提案枠組	139,805		63,295	11,338 (8.11)
arXiv (300)	単純除去	75,672		27,970	10,590 (13.99)
	実/非のみ	48,666	↓	24,630	5,457 (11.21)
	提案枠組	50,504		26,360	5,345 (10.58)
Wiki. (300)	単純除去	14,883		3,114	1,651 (11.09)
	実/非のみ	6,173	↓	2,240	282 (4.57)
	提案枠組	8,049		2,451	258 (3.21)
JNL-E (68)	単純除去	24,942		9,069	1,577 (6.32)
	実/非のみ	17,572	↓	7,865	1,058 (6.02)
	提案枠組	19,925		10,154	892 (4.48)

タグを単純に除去すると、構文解

非自然言語箇所への排除

- 文の埋め込み挿入（独立タグ）
- 自然言語ではない要素で構成された表現（実体タグ）
- 文の表示に直接関係しないシーケンス（非表示タグ）

PMC (1,000)	単純除去	170,999	58,865	18,621 (10.89)
	実/非のみ	126,176	50,741	11,881 (9.42)
	提案枠組	139,805	63,295	11,338 (8.11)
arXiv (300)	単純除去	75,672	27,970	10,590 (13.99)
	実/非のみ	48,666	24,630	5,457 (11.21)
	提案枠組	50,504	26,360	5,345 (10.58)
Wiki. (300)	単純除去	14,883	3,114	1,651 (11.09)
	実/非のみ	6,173	2,248	282 (4.57)
	提案枠組	8,049	2,451	258 (3.21)
JNL-E (68)	単純除去	24,942	9,069	1,577 (6.32)
	実/非のみ	17,572	7,865	1,058 (6.02)
	提案枠組	19,925	10,154	892 (4.48)

タグの扱い方が Enju パーザの性能に与える影響

解析時間 1~12%減

解析失敗文 0.2~0.9%減

文書	処理方法	件数	解析時間 (秒)	解析失敗文数
PMC (1,000)	単純除去	112,285	209,789	810 (0.72)
	実/非のみ	126,215	135,752	699 (0.55)
	提案枠組	74,762	132,250	2,047 (2.74)
arXiv (300)	単純除去	41,265	89,200	411 (1.00)
	実/非のみ	43,208	87,952	348 (0.81)
	提案枠組	10,561	14,704	1,161 (10.99)
Wiki. (300)	単純除去	5,026	6,743	67 (1.33)
	実/非のみ	6,893	6,058	61 (0.88)
	提案枠組	23,196	24,881	271 (1.17)
JNL-E (68)	単純除去	15,606	21,304	183 (1.17)
	実/非のみ	17,929	18,683	50 (0.28)
	提案枠組			

タグの扱い方が Enju パーザ の 性能に与える影響

解析時間 1~12%減

解析失敗文 0.2~0.9%減

文書	処理方法	件数	件数	件数	割合
PMC (1,000)	実/非のみ	112,285	135,752	810	(0.72)
	提案枠組	126,215	132,250	699	(0.55)
	単純除去	74,762	108,831	2,047	(2.74)
arXiv (300)	実/非のみ	41,265	89,200	411	(1.00)
	提案枠組	43,208	87,952	348	(0.81)
	単純除去	10,561	14,704	1,161	(10.99)
Wiki. (300)	実/非のみ	5,026	6,743	67	(1.33)
	提案枠組	6,893	6,058	61	(0.88)
	単純除去	23,196	24,881	271	(1.17)
	実/非のみ	15,606	21,304	183	(1.17)
	提案枠組	17,929	18,683	50	(0.28)

検出文数
5~37%増

独立タグ: 文区切りで検出できない文境界を示し得る
 → 文の数が増える = 文長の減少

解析時間 1~12%減

解析失敗文 0.2~0.9%減

文数	単純除去	実/非のみ	提案枠組	単純除去	実/非のみ	提案枠組	単純除去	実/非のみ	提案枠組
PMC (1,000)		112,285	126,215		135,752	132,250		810 (0.72)	699 (0.55)
	単純除去	74,762		108,831		2,047 (2.74)			
	実/非のみ	41,265		89,200		411 (1.00)			
arXiv (300)	単純除去	74,762		108,831		2,047 (2.74)			
	実/非のみ	41,265		89,200		411 (1.00)			
	提案枠組	43,208		87,952		348 (0.81)			
Wiki. (300)	単純除去	10,561		14,704		1,161 (10.99)			
	実/非のみ	5,026		6,743		67 (1.33)			
	提案枠組	6,893		6,058		61 (0.88)			
検出文数 5~37%増	単純除去	23,196		24,881		271 (1.17)			
	実/非のみ	15,606		21,304		183 (1.17)			
	提案枠組	17,929		18,683		50 (0.28)			

独立タグ: 文区切りで検出できない文境界を示し得る
 → 文の数が増える = 文長の減少

検索空間不足を抑制

解析時間 1~12%減

解析失敗文 0.2~0.9%減

文数	単純除去	実/非のみ	提案枠組	単純除去	実/非のみ	提案枠組	単純除去	実/非のみ	提案枠組
PMC (1,000)		112,285	126,215		135,752	132,250		810 (0.72)	699 (0.55)
	単純除去	74,762		108,831		2,047 (2.74)			
	実/非のみ	41,265		89,200		411 (1.00)			
arXiv (300)	単純除去	74,762		108,831		2,047 (2.74)			
	実/非のみ	41,265		89,200		411 (1.00)			
	提案枠組	43,208		87,952		348 (0.81)			
Wiki. (300)	単純除去	10,561		14,704		1,161 (10.99)			
	実/非のみ	5,026		6,743		67 (1.33)			
	提案枠組	6,893		6,058		61 (0.88)			
PMO (1,000)	単純除去	23,196		24,881		271 (1.17)			
	実/非のみ	15,606		21,304		183 (1.17)			
	提案枠組	17,929		18,683		50 (0.28)			

検出文数
5~37%増

タグの扱い方が Standard パーザ の性能に与える影響

解析時間 7~29%増

解析失敗文 0.6~1.5%減

文書	処理方法	数	失敗文数	失敗率 (%)
PMC (1,000)	実/非のみ	126,176	50,741	40.21 (9.42)
	提案枠組	139,805	63,295	45.27 (8.11)
	単純除去	75,672	27,970	36.96 (13.99)
arXiv (300)	実/非のみ	48,666	24,630	50.43 (11.21)
	提案枠組	50,504	26,360	52.20 (10.58)
	単純除去	14,883	3,114	20.92 (11.09)
Wiki. (300)	実/非のみ	6,173	2,248	36.42 (4.57)
	提案枠組	8,049	2,451	30.46 (3.21)
	単純除去	24,942	9,069	36.36 (6.32)
JNL-E (68)	実/非のみ	17,572	7,865	44.81 (6.02)
	提案枠組	19,925	10,154	50.96 (4.48)
	単純除去	24,942	9,069	36.36 (6.32)

独立タグ: 文区切りで検出できない文境界を示し得る
 → 文の数が増える = 文長の減少

解析時間 7~29%増

解析失敗文 0.6~1.5%減

文種	処理方法	数	数	数	数
PMC (1,000)	単純除去	126,176	50,741	11,881	(9.42)
	実/非のみ	139,805	63,295	11,338	(8.11)
	提案枠組	75,672	27,970	10,590	(13.99)
arXiv (300)	単純除去	48,666	24,630	5,457	(11.21)
	実/非のみ	50,504	26,360	5,345	(10.58)
	提案枠組	14,883	3,114	1,651	(11.09)
Wiki. (300)	単純除去	6,173	2,248	282	(4.57)
	実/非のみ	8,049	2,451	258	(3.21)
	提案枠組	24,942	9,069	1,577	(6.32)
PMO (1,000)	単純除去	17,572	7,865	1,058	(6.02)
	実/非のみ	19,925	10,154	892	(4.48)
	提案枠組				

検出文数
4~30%増

独立タグ: 文区切りで検出できない文境界を示し得る
 → 文の数が増える = 文長の減少

50単語以下の(解析対象となる)文の増加

解析時間 7~29%増

解析失敗文 0.6~1.5%減

文数	処理方法	検出文数	解析時間	解析失敗文
PMC (1,000)	単純除去	126,176	50,741	11,881 (9.42)
	実/非のみ	139,805	63,295	11,338 (8.11)
	提案枠組	75,672	27,970	10,590 (13.99)
arXiv (300)	単純除去	48,666	24,630	5,457 (11.21)
	実/非のみ	50,504	26,360	5,345 (10.58)
	提案枠組	14,883	3,114	1,651 (11.09)
Wiki. (300)	単純除去	6,173	2,248	282 (4.57)
	実/非のみ	8,049	2,451	258 (3.21)
	提案枠組	24,942	9,069	1,577 (6.32)
Wiki. (300)	単純除去	17,572	7,865	1,058 (6.02)
	実/非のみ	19,925	10,154	892 (4.48)
	提案枠組			

検出文数
4~30%増

タグの扱い方が JUMAN / MeCab の 性能に与える影響

JUMAN

文書(数)	手法	検出文数	解析時間(秒)	失敗文数(%)
JNL-J (384)	単純除去	96,668	122	10 (0.01)
	実/非のみ	76,277	86	8 (0.01)
	提案枠組	114,250	59	2 (0.00)

MeCab

文書(数)	手法	検出文数	解析時間(秒)	失敗文数(%)
JNL-J (384)	単純除去	97,312	7	10 (0.01)
	実/非のみ	78,461	6	8 (0.01)
	提案枠組	116,424	6	2 (0.00)

タグの扱い方が JUMAN / MeCab の性能に与える影響

JUMAN

文書(数)	手法	検出文数	解析時間(秒)	失敗文数(%)
JNL-J (384)	単純除去	96,668	122	10 (0.01)
	実/非のみ	76,277	86	8 (0.01)
	提案枠組	114,850	59	2 (0.00)

提案枠組: 解析時間51%減
= 大幅短縮を達成

提案枠組: 解析失敗 0.0%
= 文書被覆率ほぼ 100%

文書(数)	手法	検出文数	解析時間(秒)	失敗文数(%)
JNL-J (384)	単純除去	97,312	7	10 (0.01)
	実/非のみ	78,461	6	8 (0.01)
	提案枠組	116,424	6	2 (0.00)

対象文書全体の高被覆かつ効率的な処理が意味するもの

- 浅い分析(単語数のカウント等)で良いタスク
→ タグの単純除去でも事足りる
 - 文の挿入は単語数に影響せず
 - 非自然言語要素は大量のテキストで打ち消せる
 - 詳細かつ正確な分析要求(談話分析・翻訳・文法抽出など)にはより厳密な解析が必要
 - ごく狭い箇所で発せられた情報も慎重に吟味
 - その分析は、本文とは関係のない要素の混入を確実に排除したテキストに対して行う必要がある
- = 近年 NLP が確立してきたものそのもの

対象文書全体の高被覆かつ効率的

な処理が実現できる

自然言語現象の正確な解析に集中するために、
データセットは平文テキスト列の形に整備
→ 新たなタスクの発見・取り組み・成果報告

– 非自然言語要素は大量のテキストで打ち消せる

- 詳細かつ正確な分析要求（談話分析・翻訳・文法抽出など）にはより厳密な解析が必要
 - ごく狭い箇所で発せられた情報も慎重に吟味
 - その分析は、本文とは関係のない要素の混入を確実に排除したテキストに対して行う必要がある
- = 近年 NLP が確立してきたものそのもの

実文書を自然言語処理技術と適切に繋ぐ技術の重要性

- 多くの有用な成果は、本来目指す最終目的「実世界の文書の解析」へ還元させてこそ
- 本研究の実験結果が示唆するもの：
「適切な枠組で支援すれば、従来のNLPツールには、実文書テキストを『本来期待される性能で』解析する能力が既に備わっている」
- 「実文書とNLP技術を適切に繋ぐ」：
NLPをあらゆる局面で利用する際、その本来の効力をフル活用するための重要なタスク

本研究のまとめ

- XMLタグ付テキストとNLPツール入出力との間の変換枠組を提案
 - タグを4種のテキスト構成機能に分類
 - 分類に基づき平文テキスト列へ自動的に変換
 - 対象文書中の20%以下のタグを分類することで平文テキスト列を獲得成功
 - タグの単純除去に比べ、NLPツールの大幅な解析性能改善(高被覆・効率化)を確認
- 実文書をNLPと適切に繋ぐ技術の重要性示唆

今後の展望

- 提案枠組ツールの公開(近日)
 - 多種多様な実文書へのNLPツール適用の議論の提供・共有
- より柔軟・繊細な構造解析の検討
 - (例) 実体タグ内と周囲のテキストの繋ぎ方
- より多様な文書形式への適用可能性の検討
 - OCRで得られた文書・スライドデータ等との接続

ありがとうございました