

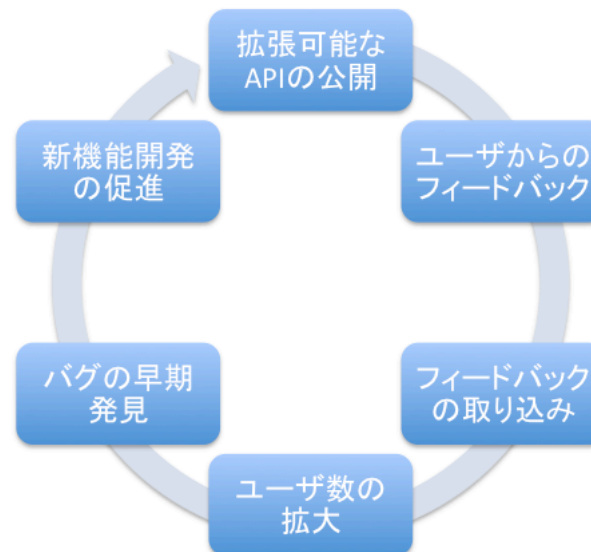


Solrにおける日本語処理の現状

株式会社 ロンウイト
www.rondhuit.com

Apache Lucene/Solr

- Apache Lucene
 - Javaで書かれたオープンソースの全文検索ライブラリ
- Apache Solr
 - Luceneをベースに開発されたオープンソースの全文検索サーバ



よくある質問／問題

- あいまい検索はできますか？
- 形態素解析とN-gramどちらを使えばいいですか？
- ハイライトが正しく動作しません。
- ハイライトがずれます。

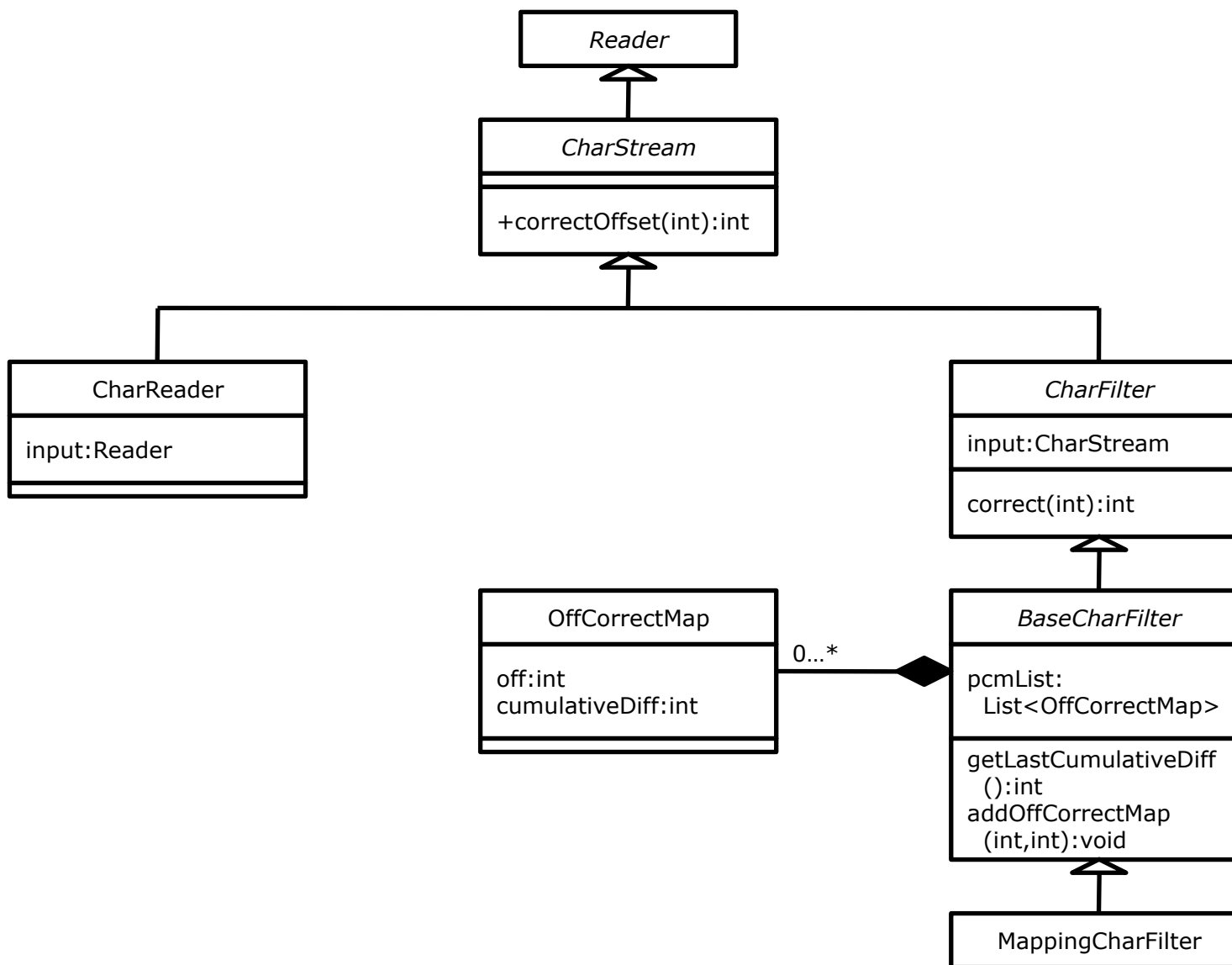
あいまい検索はできますか？

- あいまい検索とは？

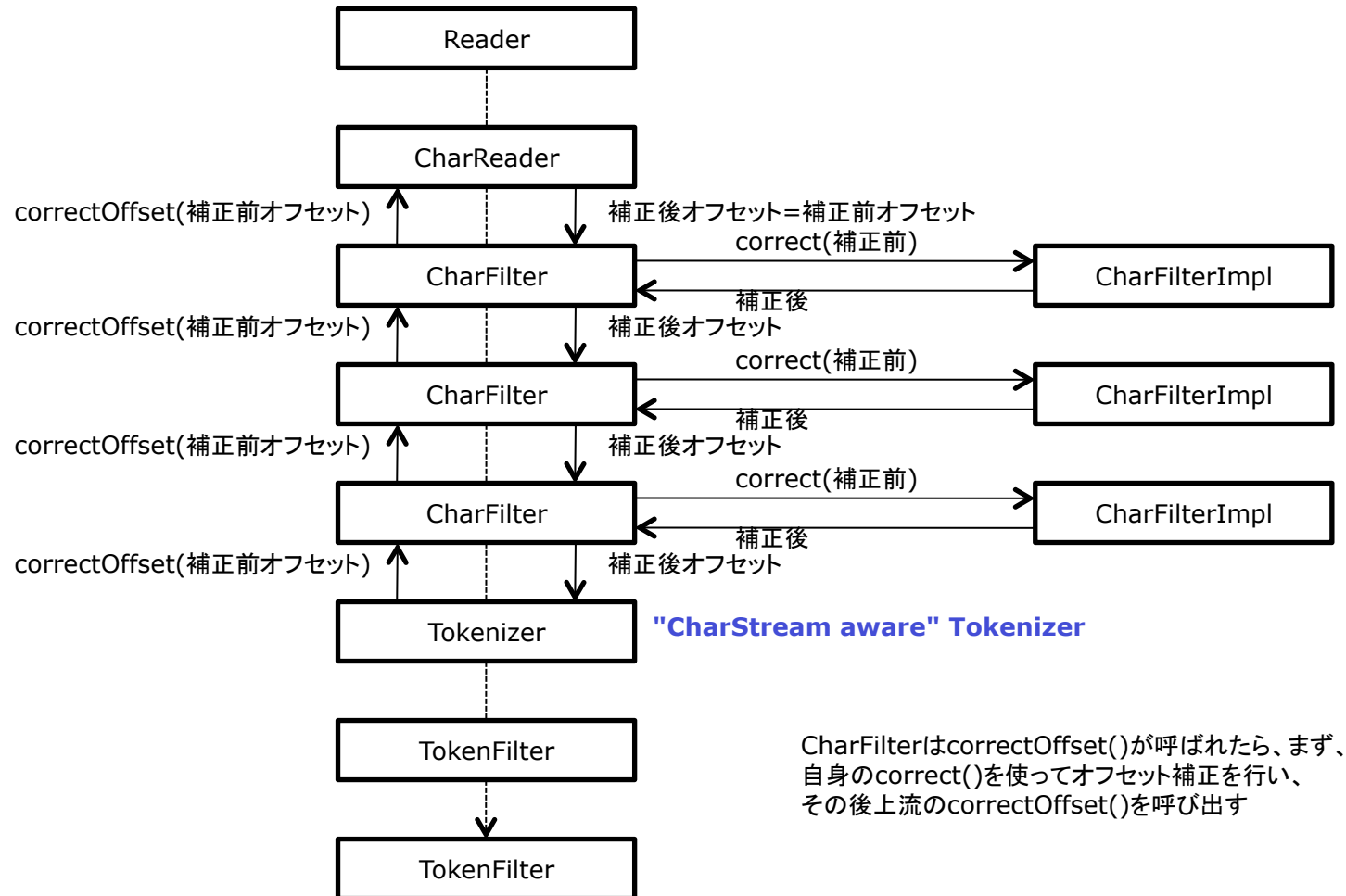
- ユーザによって「あいまい」の定義はさまざま

- 全角・半角混じり文の検索
- カタカナ語の表記揺れ
- 正式名称 ⇔ 短縮名称検索
- いすゞ ⇔ いすゞ
- 結合文字検索
- 西暦 ⇔ 和暦検索
- ワイルドカード検索 / 正規表現検索 / あいまい (Fuzzy) 検索
- Google 「もしかして・・・」
- Google サジェスト

CharFilter – クラス図



文字オフセット補正の動き



MappingCharFilter

- <charFilter/>のmapping属性で指定されたファイルに従った「文字マッピング」を実行する。
- mapping.txtのサンプル

```
# Syntax:
# "source" => "target"
# "source".length() > 0 (source cannot be empty.)
# "target".length() >= 0 (target can be empty.)

# example:
# "À" => "A"
# "\u00C0" => "A"
# "\u00C0" => "\u0041"
# "ß" => "ss"
# "\t" => " "
# "\n" => ""

# À => A
"\u00C0" => "A"

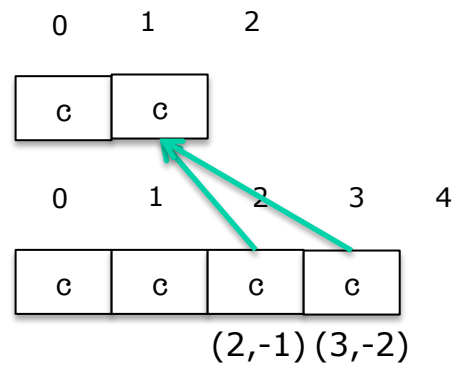
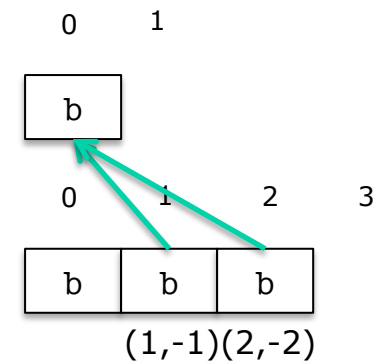
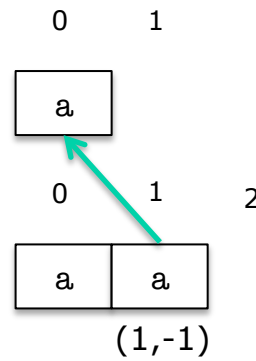
# Æ => AE
"\u00C6" => "AE"

# Ç => C
"\u00C7" => "C"
```

MappingCharFilter(オフセット補正)

a => aa
b => bbb
cc => cccc

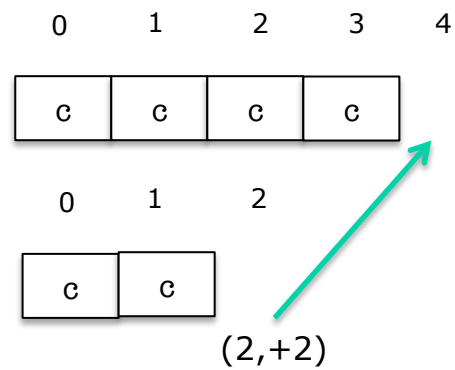
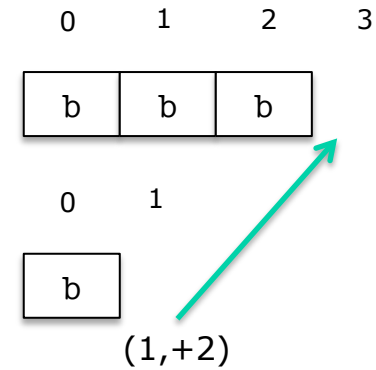
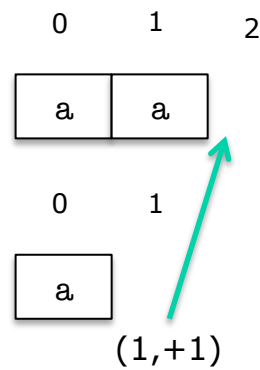
文字数が増える場合



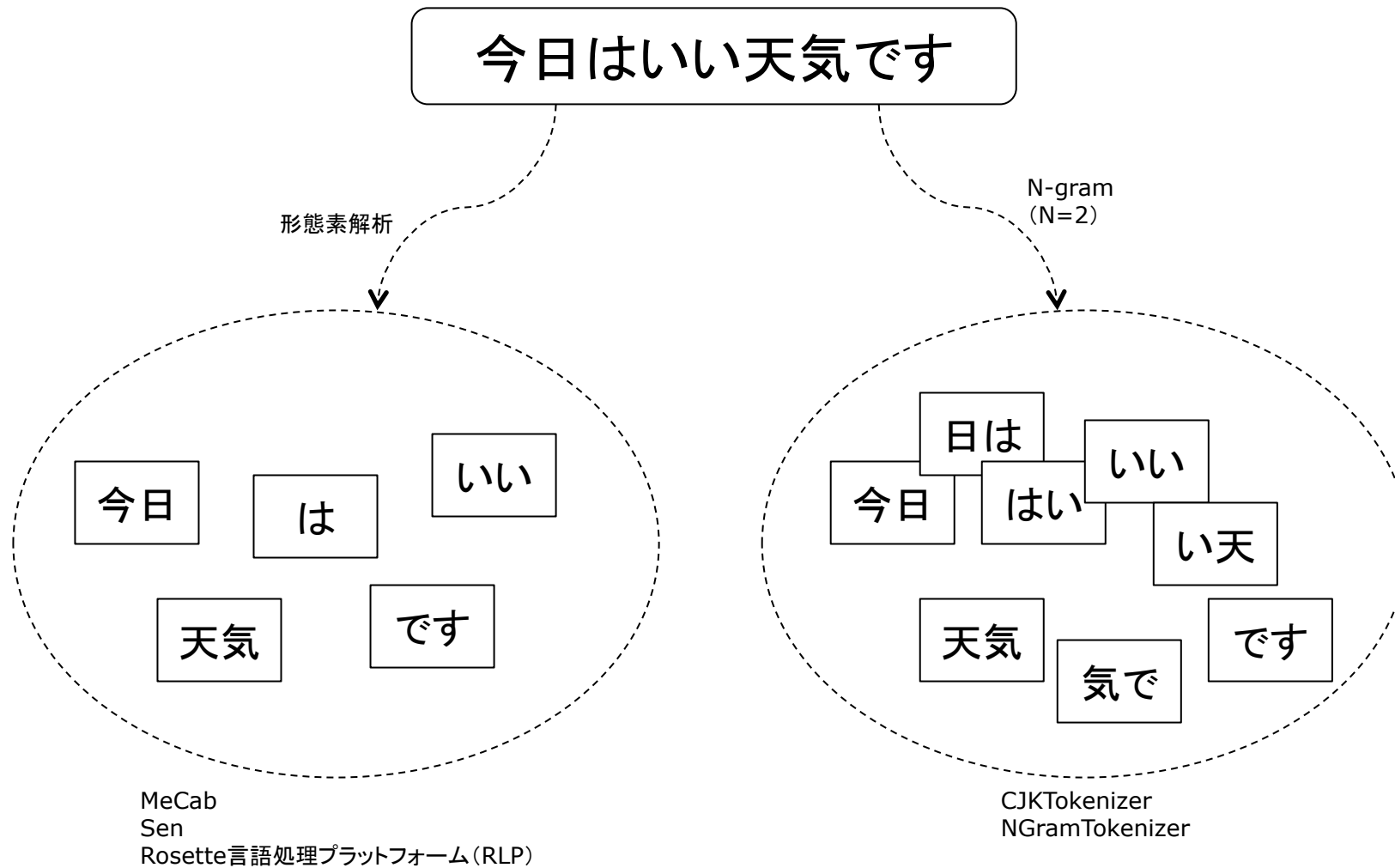
MappingCharFilter(オフセット補正)

```
aa => a  
bbb => b  
cccc => cc
```

文字数が減る場合



形態素解析とN-gram



形態素解析とN-gram

	形態素解析	N-gram
長所	<ul style="list-style-type: none">• 自然言語の単語に近い• N-gramよりもインデックスサイズを小さくできる• 単語の品詞が抽出できる	<ul style="list-style-type: none">• 実装が容易• 辞書不要• 新語への対応が容易
短所	<ul style="list-style-type: none">• 新語への対応が困難• 辞書のメンテナンスが必要• N-gramよりも検索漏れが生じやすい	<ul style="list-style-type: none">• 形態素解析よりもインデックスサイズが大きくなりがち• Nより短い文字が検索漏れとなる
用途	<ul style="list-style-type: none">• NLP• 固有表現抽出	<ul style="list-style-type: none">• 特許検索• ブログ検索

単語の境界がはっきりしない例



製造部門長谷川

ハイライトのトラブル

- ハイライトとは？

自然言語処理

検索

ウェブ全体から検索 日本語のページを検索

[自然言語処理 - Wikipedia](#) ☆

検索語の強調表示

ハイライトスニペット

[編集] 基礎技術. **自然言語処理**の基礎技術にはさまざまなものがある。**自然言語処理**はその性格上、扱う言語によって大きく処理の異なる部分がある。現在のところ、日本語を処理する基礎技術としては以下のものが主に研究されている。 ...

[基礎技術 - 処理内容とその限界](#) - [具体的な課題](#) - [統計的自然言語処理](#)

ja.wikipedia.org/wiki/自然言語処理 - [キャッシュ](#) - [類似ページ](#)

- トラブルその1

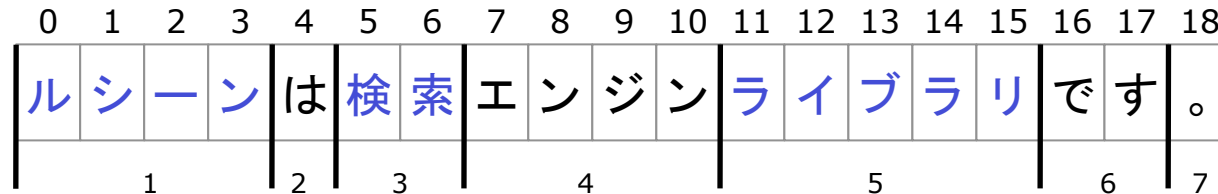
「N-gramを使うとハイライトがおかしくなります」

- トラブルその2

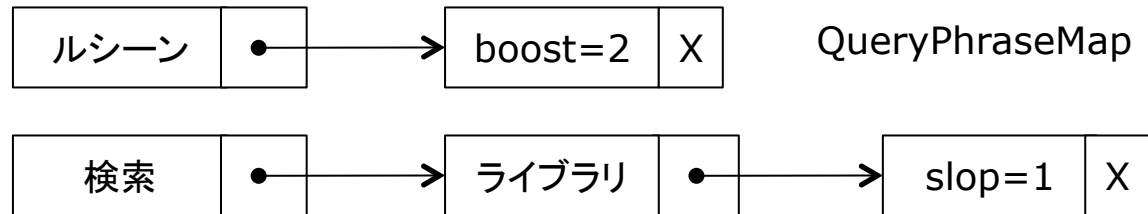
「FVHを使うとハイライトがずれます」

FastVectorHighlighter

q=ルシーン^2 "検索 ライブラリ"~1



```
public class QueryPhraseMap {
    boolean terminal;
    int slop;
    float boost;
    Map<String, QueryPhraseMap> subMap;
}
```



FieldTermStack

ルシーン (0,4,1)
 検索 (5,7,3)
 ライブラリ (11,16,5)

FieldPhraseList

ルシーン [(0,4)] boost=2
 検索 ライブラリ [(5,7),(11,16)] boost=1

FieldFragList

ルシーン [(0,4)]
 検索 ライブラリ [(5,7),(11,16)]
 totalBoost=3

Solrサブスクリプション・パッケージ



ロンウイットは、OSSのSolrを企業／団体様が安心してお使いいただけるよう、年間購読料形式の「サブスクリプション・パッケージ」としてご提供しています。お客様のアプリケーションの目的に応じたエディションをご選択いただけます。

	BASIC	LANGUAGE	ENTERPRISE
Solr 1.4	✓	✓	✓
バグフィックス	✓	✓	✓
各種プラグイン	✓	✓	✓
サポートサービス	✓	✓	✓
言語判別		✓	✓
類義語辞書		✓	✓
固有表現抽出			✓

Solrサブスクリプション・パッケージ



● 各種プラグイン

※提供予定の機能を含みます。また機能は順次追加されます。

- 全角／半角／新旧漢字混じり、カタカナ語表記揺れ吸収検索
- Java 6のNormalizerを使用した正規分解／互換分解／正規合成／互換合成処理による検索
- 踊り字検索
- ギャル文字検索
- 顔文字検索
- 1-gram最適化検索
- 専門用語抽出
- カタカナ⇔ひらがな⇔ローマ字検索
- 日本語文章境界を認識するハイライター
- 日本語サジェスチョン
- Heritrixアーカイブファイルローダー

Solrサブスクリプション・パッケージ



- 言語判別

- 世界およびインターネット上で使われている主要な下記25言語をサポート:

アラビア語、ベンガル語、デンマーク語、ドイツ語、ギリシャ語、英語、スペイン語、ペルシャ語、フィンランド語、フランス語、ヒンディー語、ハンガリー語、アイスランド語、イタリア語、日本語、ジャワ語、韓国語、オランダ語、ノルウェー語、ポルトガル語、ロシア語、スウェーデン語、タイ語、トルコ語、中国語

- 類義語辞書

- SynonymFilterで利用可能な分野別類義語辞書:

分野	およそのエントリ数
コンピューター	8,000
政治・経済	14,000
医療	3,000
科学	2,500
教育	3,000
芸能・スポーツ	2,000

Solrサブスクリプション・パッケージ



- 固有表現抽出 ※BASIS Technology社の固有表現抽出製品REXを用いています。
 - 人名／地名／国名／組織名などを文章から抽出
 - ベイズ理論に基づき、抽出すべき語の「文脈パターン」を事前学習

Language: Japanese
Script: Japanese (alias for Han + Hiragana + Katakana)
MIME Type: binary
Encoding: UTF-16
Length: 149

Named Entities

Named Entity (# instances)	
GPE	4
LOCATION	1
ORGANIZATION	1
PERSON	1
TEMPORAL:DATE	5

北京オリンピックは、2008年8月8日から8月24日までの期間、中華人民共和国の首都北京で開催。アジアで夏季のオリンピックが開催されるのは、1988年のソウルオリンピック以来20年ぶり。SPEEDOの水着で臨んだ北島選手は、日體大在学中の2004年に続き男子平泳ぎで二個の金メダルを獲得した。

#	Type	Phrase
1	GPE	北京
2	TEMPORAL:DATE	2008年8月8日
3	TEMPORAL:DATE	8月24日
4	GPE	中華人民共和国
5	GPE	北京
6	LOCATION	アジア
7	TEMPORAL:DATE	1988年
8	GPE	ソウル
9	TEMPORAL:DATE	20年
10	PERSON	北島
11	ORGANIZATION	日體大
12	TEMPORAL:DATE	2004年

ロンウイットのSolr技術トレーニング



- 特徴

- Lucene/Solrコミッター監修
- 演習付き
- PDFテキスト(持ち帰りOK)
- 受講後1ヶ月のアフターサポート付き

- トレーニングコース

※2010年6月より順次ご提供できる予定です。

- Solr 1.4 基礎
- Solr 1.4 応用
- Solr 1.4 運用(?)
- Solr 1.4 プラグイン開発
- Solr 1.4 DIH開発(?)

ご静聴ありがとうございました！