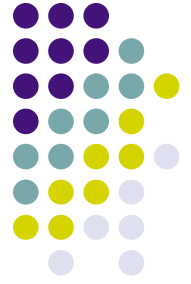


MeCab

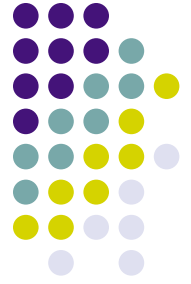
汎用日本語形態素解析エンジン

工藤 拓



アジェンダ

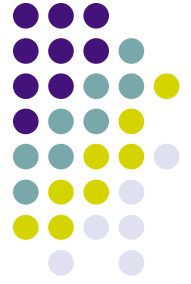
- 形態素解析の技術
 - 辞書引きのアルゴリズム、データ構造
 - 曖昧性の解消
- MeCab の開発裏話
 - 歴史
 - 設計方針
- 汎用テキスト変換ツールとしての MeCab
 - 恐ろしく汎用的!
 - 「意外な」使い方
- これから



形態素解析

- 文を単語に区切り、品詞を同定する処理
 - 全文検索 Spam フィルタリング 人工無能...
- 以下の3つの処理
 - 単語への分かち書き(tokenization)
 - 活用語処理(stemming, lemmatization)
 - 品詞同定(part-of-speech tagging)

すもも	名詞,一般,****,すもも,スモモ,スモモ
も	助詞,係助詞,****,も,モ,モ
もも	名詞,一般,****,もも,モモ,モモ
も	助詞,係助詞,****,も,モ,モ
もも	名詞,一般,****,もも,モモ,モモ
の	助詞,連体化,****,の,ノ,ノ
うち	名詞,非自立,副詞可能,***,うち,ウチ,ウチ
。	記号,句点,****,。 ,。 ,。



形態素解析の技術

- 基本的な処理: 辞書から単語を引いて、与えられた文と照合し、最も自然な単語列を求める
 - 辞書引き
 - 入力文は単語毎に区切られていない
 - どの文字列を辞書引きするか自明ではない
 - 曖昧性の解消
 - すべての可能な単語の組合せから(何らかの基準で)最適な単語列を発見する
 - 基準の定義



日本語処理のための辞書の要件

- 単語の区切りが明確でないので、先頭から何文字までが単語なのかわからない

奈良先端科学技術大学院大学情報科学研究科



- 単純な方法(hash)だと (文長)² の辞書引きが発生!

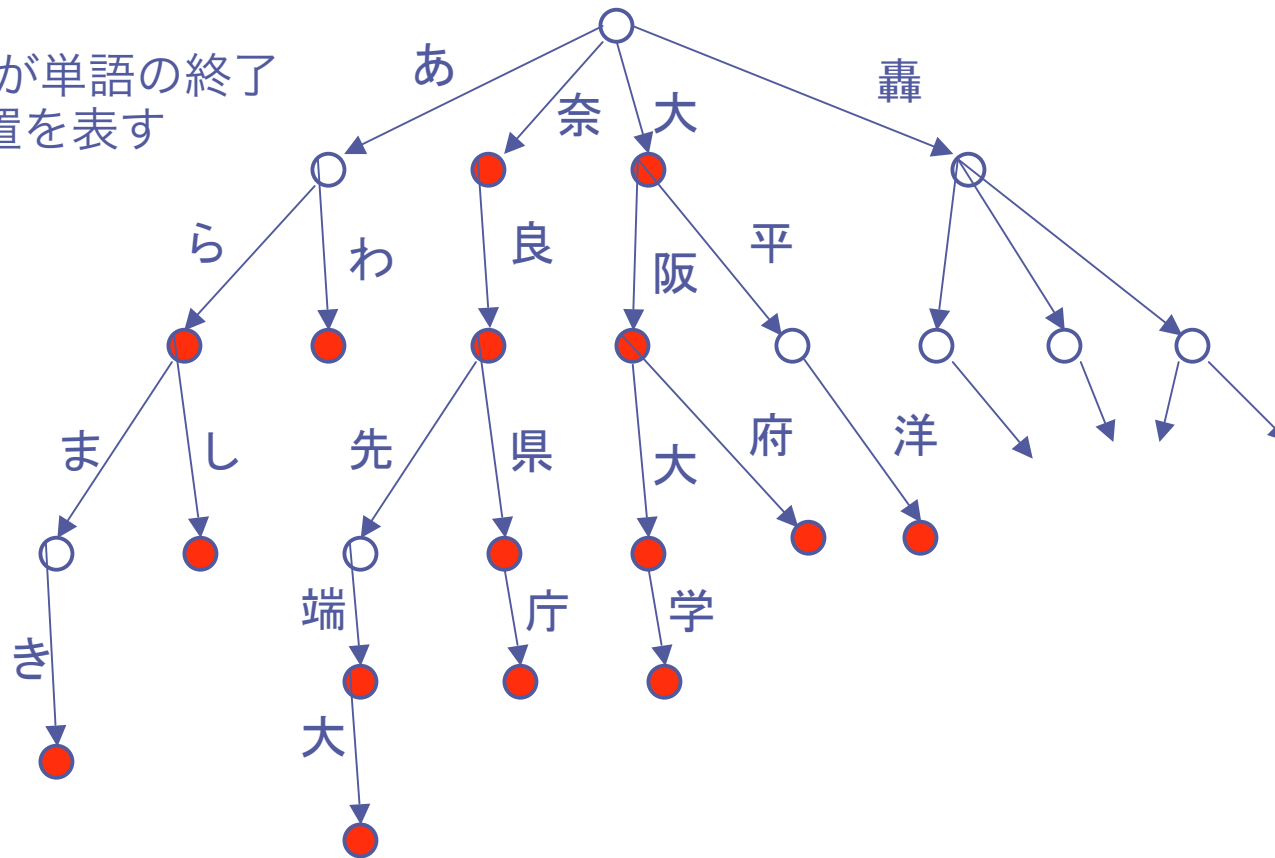
```
$str = "奈良先端科学技術大学院大学情報科学研究科";  
for (my $i = 0; $i < length($str); ++$i) {  
  for (my $j = 1; $j < length($str) - $i; ++$j) {  
    my $key = substr($str, $i, $j);  
    if (defined $dic{$key}) {  
      ...;  
    }  
  }  
}
```

*dbm, RDB は辞書として使えない



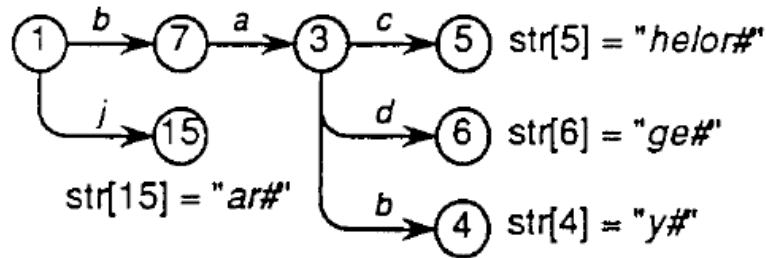
辞書検索のためのデータ構造：TRIE

- 赤丸が単語の終了位置を表す



- 対象文字列の先頭から文字を順番にたどるだけ
- 辞書引き終了のタイミングが自動的にわかる
- 入力文字列の長さに比例した時間 $O(\text{文長})$ で探索が可能
- さまざまな実装

Double-Array(ダブル配列) TRIE



{#= 1, a=1, b=2, c=3..}

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
BASE	4	0	1	-15	-1	-12	1	0	0	0	0	0	0	0	-9
CHECK	0	0	7	3	3	3	1	0	0	0	0	0	0	0	1

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
TAIL	h	e	l	o	r	#	?	?	a	r	#	g	e	#	y	#

POS=17

An efficient implementation of Trie Structures,
Aoe et al 92 より引用

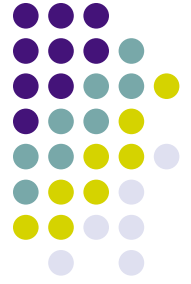
```
int n = 1
for (int i = 0; i < strlen(key); ++i) {

    int k = BASE[n] + charcode(key[i]);
    if (CHECK[k] != n) break;

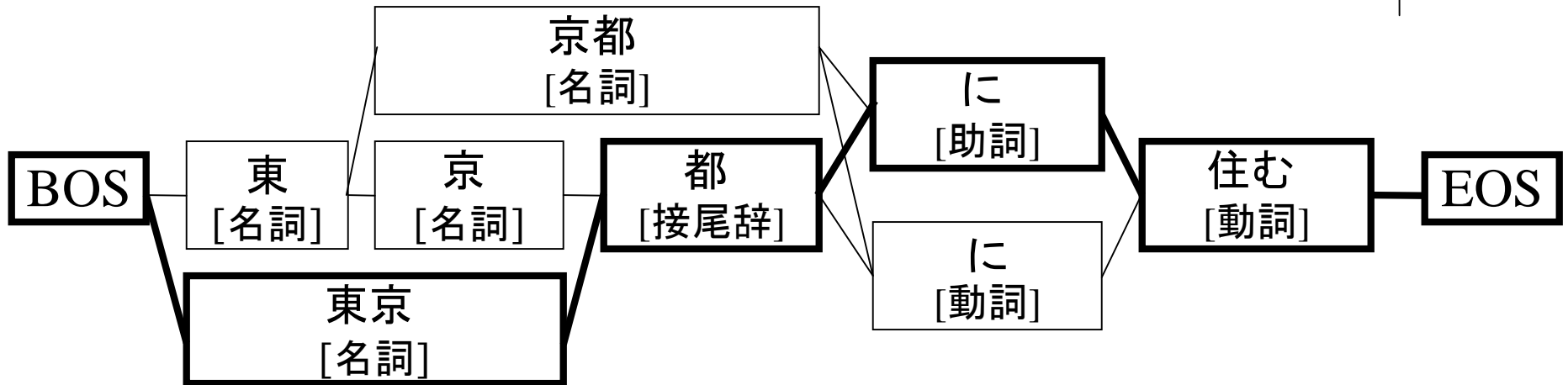
    // 見つかった!
    if (BASE[k] < 0)
        printf ("%d\n", -BASE[k]);

    n = k;
}
```

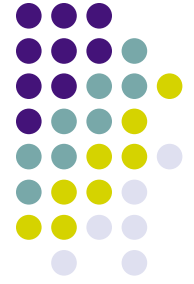
- MeCabに採用 (後に ChaSen も)
- 利点: (知る限り)最も高速
- 欠点: 辞書サイズが大きい, 構築が面倒



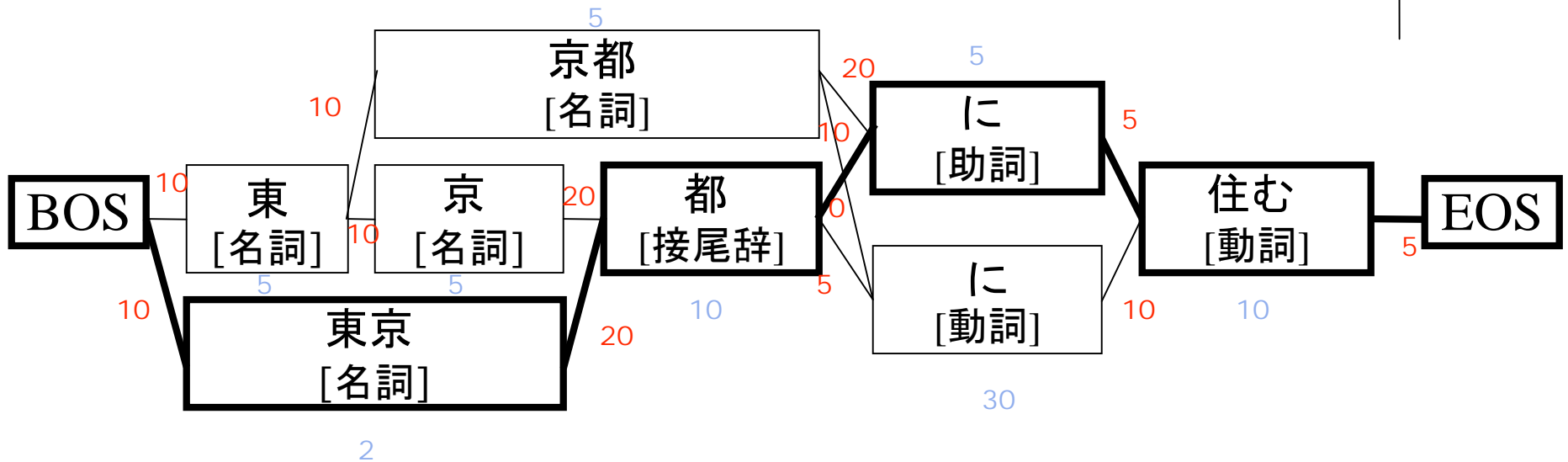
曖昧性の解消



- 規則(ヒューリスティックス)に基づく手法 (80年代)
 - 最長一致: 長い単語を優先 (KAKASI)
 - 分割数最小: 文全体の単語の数を最小にする候補
 - 文節数最小: 文全体の文節数を最小にする候補
- 多くの場合曖昧性を解決できない



最小コスト法



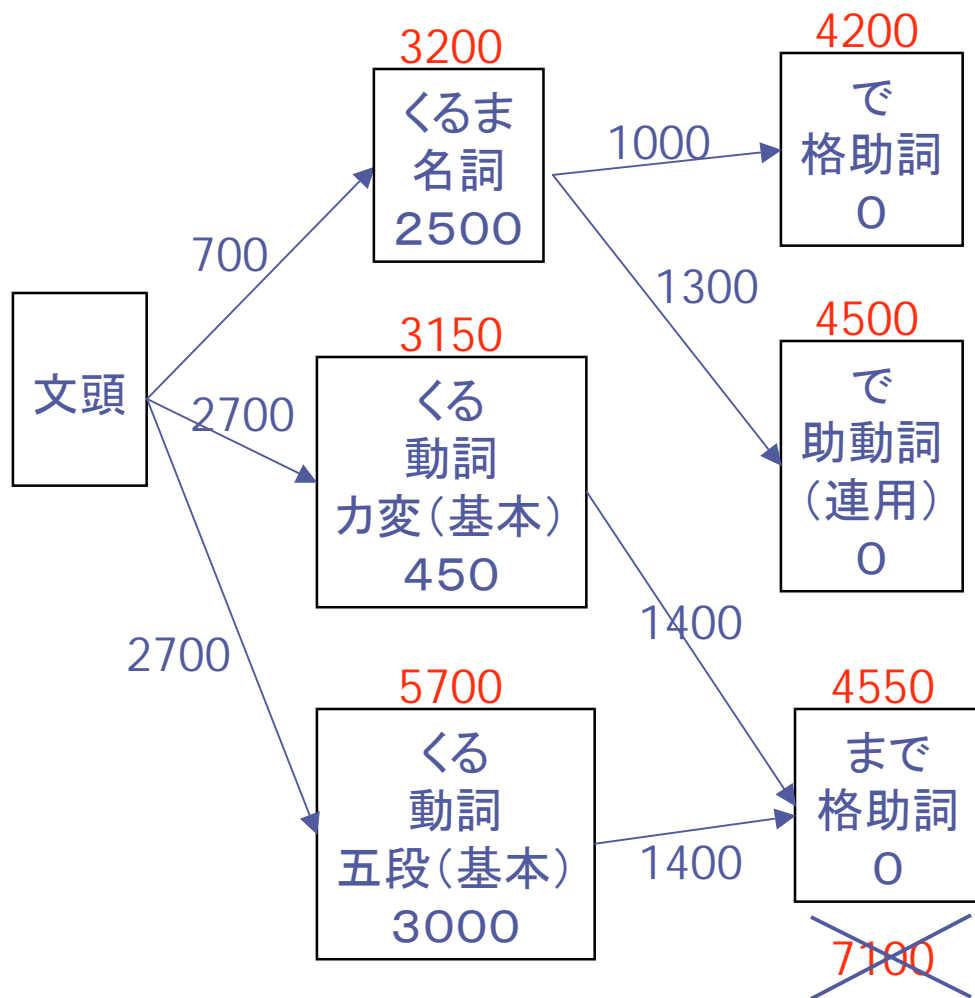
接続コスト: 二つの単語のつながりやすさ

生起コスト: 一つの単語の出現しやすさ

- 接続コストと生成コストの和が最小になる解
- コストはなんらかの方法で決定 (後述)
- Viterbi アルゴリズム (動的計画法の一種) で $O(\text{文長})$ で探索可能

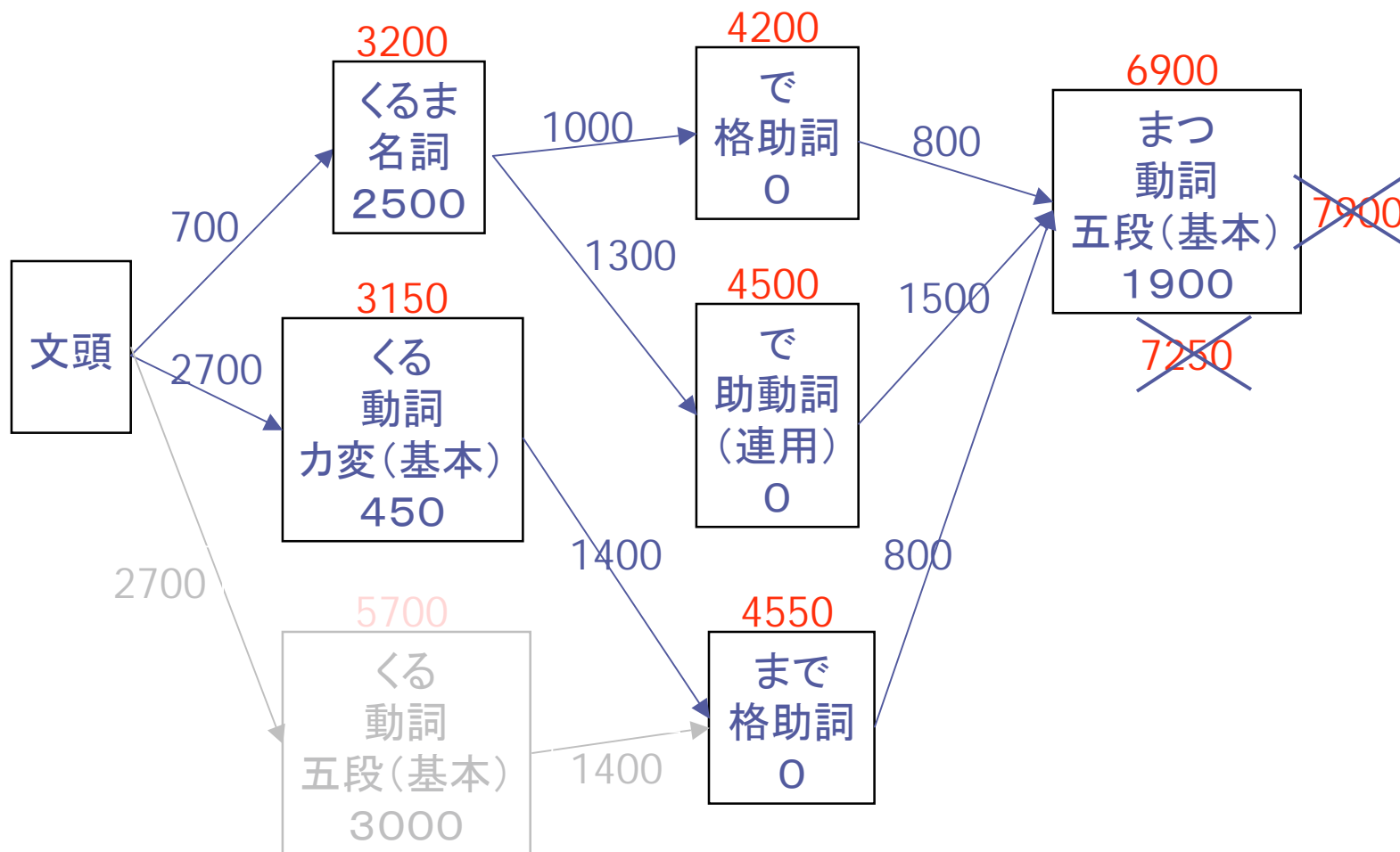


最小コスト法 (Viterbi アルゴリズム)



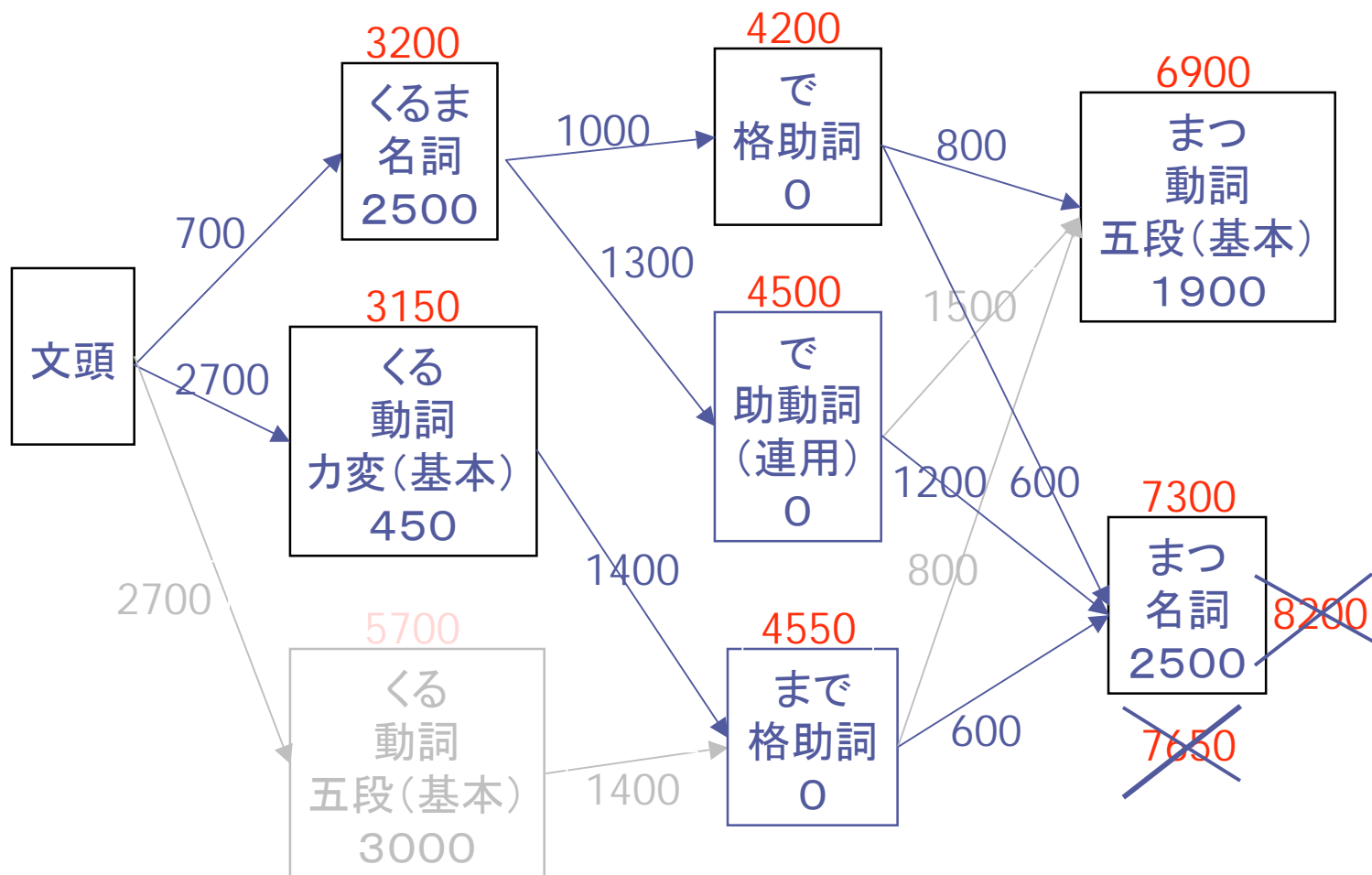


最小コスト法 (Viterbi アルゴリズム)



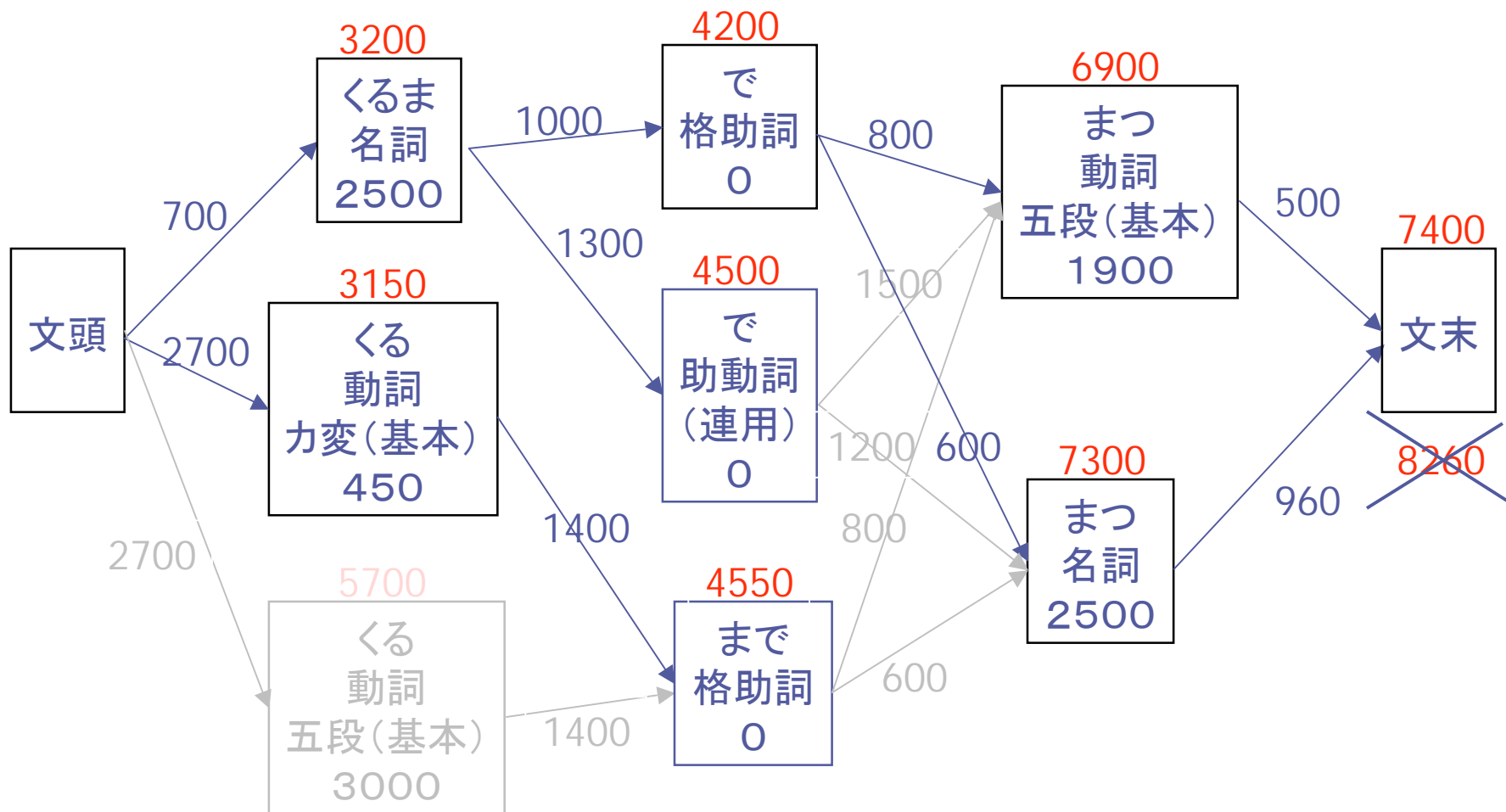


最小コスト法 (Viterbi アルゴリズム)



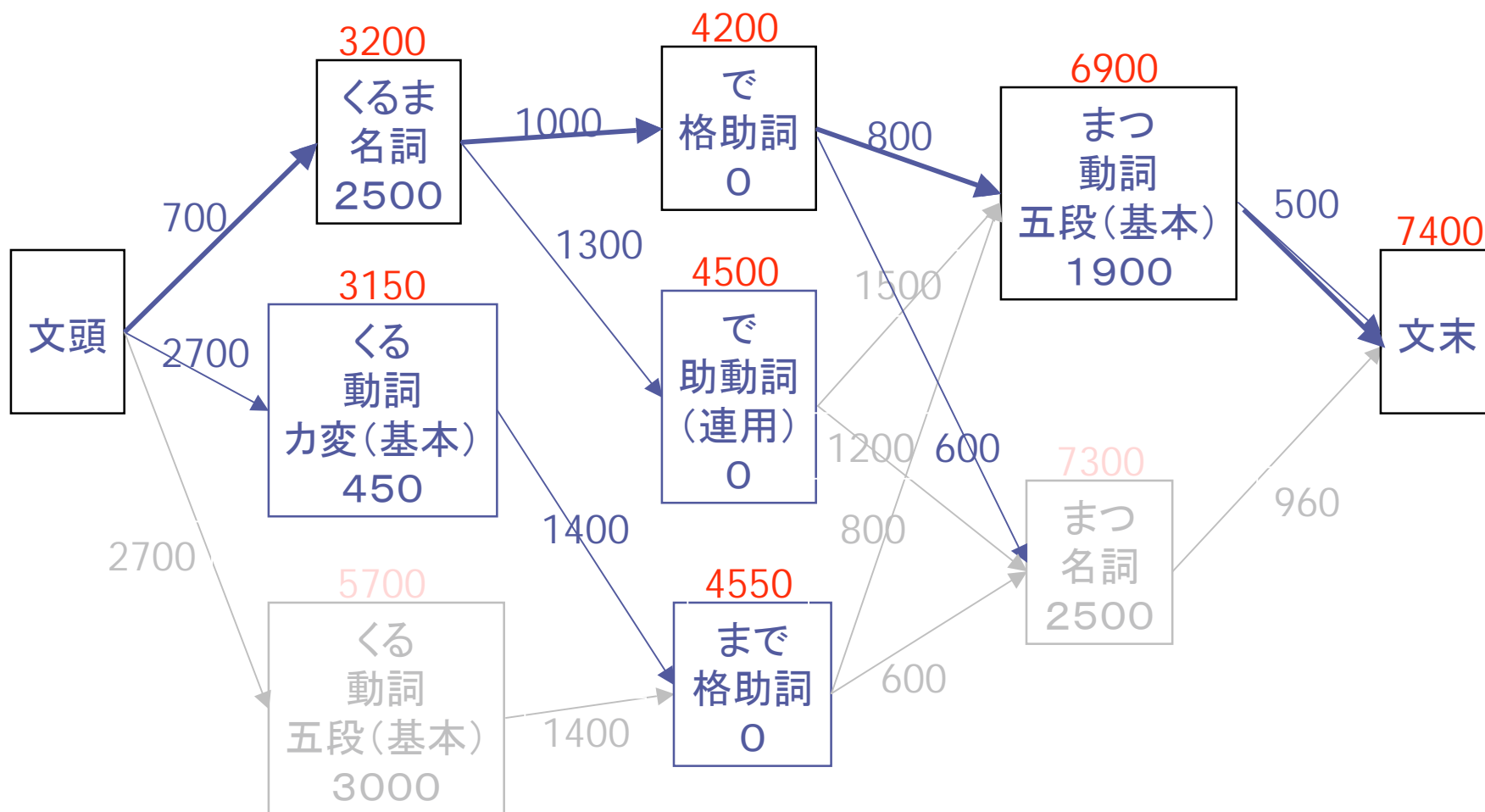


最小コスト法 (Viterbi アルゴリズム)



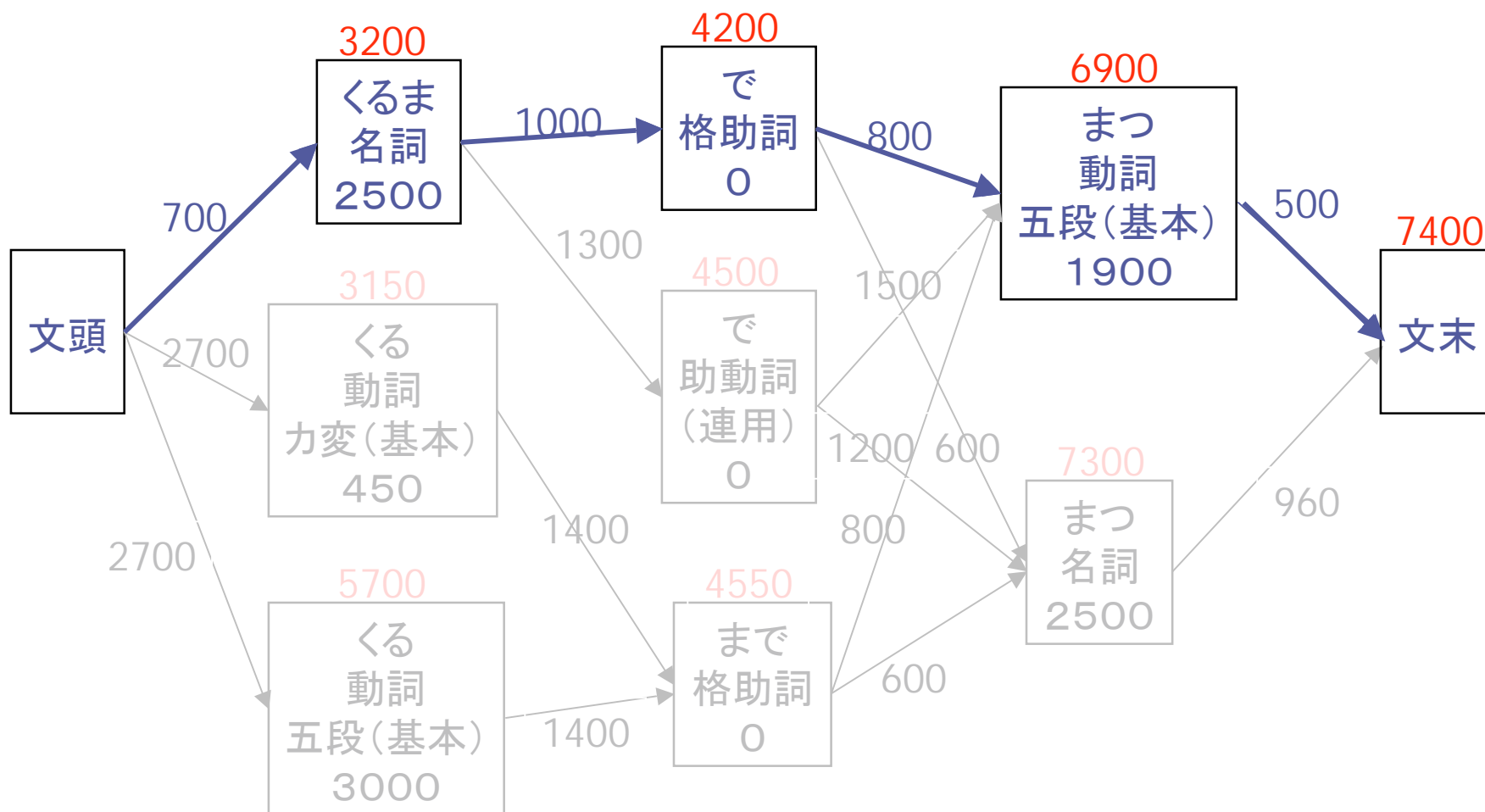


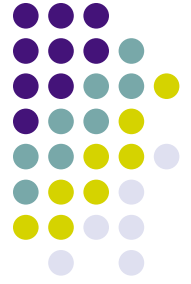
最小コスト法 (Viterbi アルゴリズム)





最小コスト法 (Viterbi アルゴリズム)





コストの決定方法

- 人手でガンバル (90年代はじめ)
 - 試行錯誤の連続, かなり大変
 - 客観的評価が難しい
- 統計処理
 - 大量の生テキストから推定
 - 超低コスト
 - 質に問題がある (全文検索目的だったら可能かも)
 - 正解データを人手で作ってデータから推定
 - 現代の形態素解析器の主流
 - 低コスト
 - これから...
 - 大量の生テキスト + 少量の正解データ + 統計処理

正解データ作成ツール (VisualMorphs)



VisualMorphs -p property.vm -a analyzer.vm -c C:\WINDOWS\デスクトップ\test.txt

ファイル PartOfSpeech Inflection

全文解析 ▲ ▲ 5 じゃあ京都行くまでに通行止めとかないのかなあ

部分解析 × ● 単語に区切られていない

切り出し ▼ ▼

見出し語 品詞 活用 最大コスト

基本形 全文コスト 10015.0

読み 全文解析幅 5000.0

発音 部分解析幅 10000.0

破棄

に 助詞 格助詞 一般 0

区 区切ら 動詞 自立 2950

れ 動詞 接尾 0

て 助詞 接続助詞 6

い 動詞 非自立 0

ない 助動詞 0

BOS/EOS 0

単語 名詞 一般 3929

に 助詞 副詞化 0

区 名詞 一般 3556

切ら 動詞 自立 1980

て 動詞 非自立 1445

い 動詞 自立 1344

単 接頭詞 名詞接続 2315

語 名詞 接尾 一般 2044

区 名詞 接尾 助数詞 2097

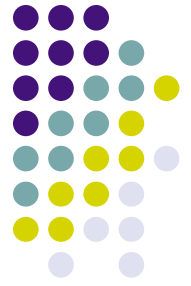
単 名詞 一般 3755

語 名詞 接尾 一般 2045

語 名詞

672 672 658 658 425 425 162 162 392 392 270 270 21 21 330 330

1572 828 2243 537 746 162 392 1430 279 672 828 3440



Conditional Random Fields

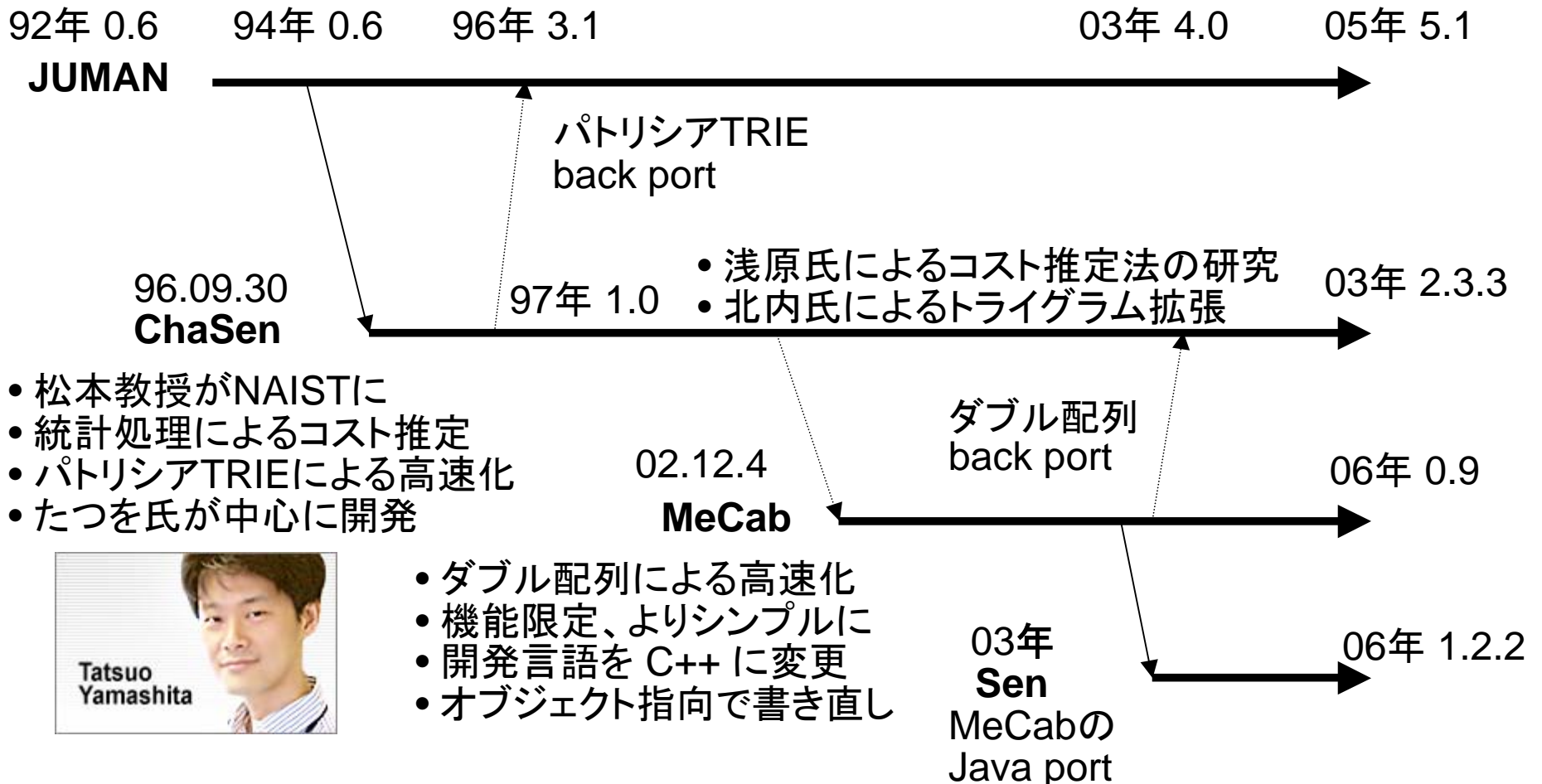
- MeCab 0.90 から採用された統計的コスト推定アルゴリズム
- コスト推定モジュールを同封
- ChaSenが採用している手法に比べて高性能
 - 1/3 程度の正解データで同程度の性能
- CRFの基本的な考え方
 - 「**正解のコスト < 残りの解のコスト**」
が満たされるようにコスト値を探索
 - ラティス中のコストをカナヅチで調整するイメージ
 - 数値最適化問題に帰着

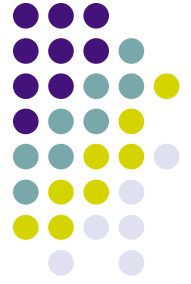


オープンソース形態素解析器

- 松本教授による prolog プロトタイプ
- 妙木,黒橋氏による C 実装
- コスト決定に2ヵ月

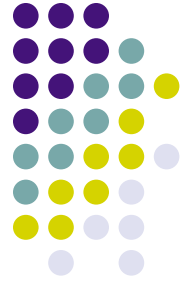
辞書の再編成





MeCab の設計方針

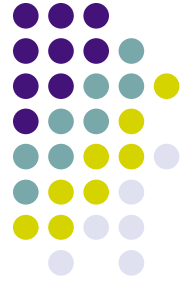
- 辞書とシステムの完全分離
 - 自然言語の複雑さはシステムではなく辞書/コストとして外部定義
 - システムは日本語を知らない
 - `grep "名詞" *.cpp` としても何も出てこない :-)
 - システムは「ひらがな」「カタカナ」の区別すら知らない (文字種の情報もすべて外部定義)
 - 他の言語も辞書さえあれば解析可能
- 解析速度を犠牲にしない
 - 事前に計算できることはすべてやっておく
 - 辞書やコスト値はすべてバイナリデータ
 - ディスクの使い方は富豪的



MeCab の設計方針

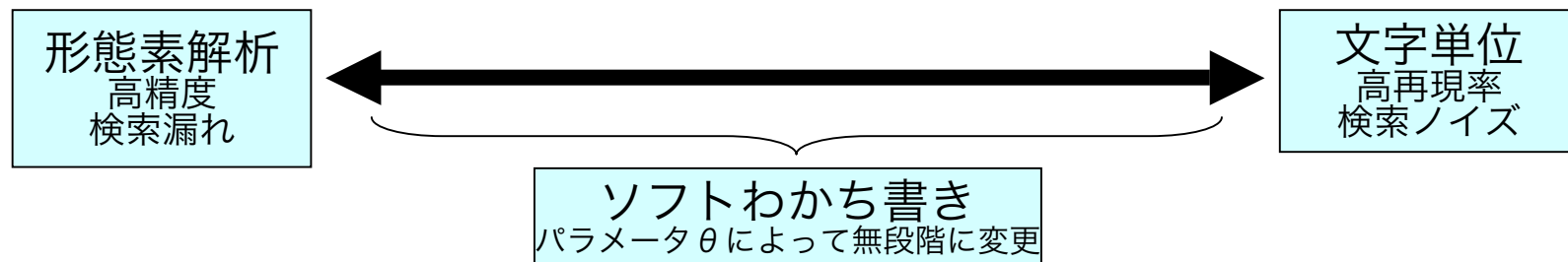
- 機能の選別
 - 前処理/後処理でできることはやらない
 - ChaSen の機能過多の反省
 - 文字コード変換, 改行処理, 連結品詞, 注釈, ChaSenサーバ
 - かわりに API を充実
 - C/C++, Perl, Java, Python, Ruby, C# ...
 - 解析器にしかできない機能を提供
 - N-best 解, 制約つき解析, ソフト分かち書き(後述)

```
use MeCab;
my $str = "すもももももものうち";
my $mecab = new MeCab::Tagger("");
for (my $n = $mecab->parseToNode($str);
     $n; $n = $n->{next}) {
    printf "%s\t%s\n", $n->{surface}, $n->{feature};
}
```



ソフトわかち書き = (形態素解析 + 文字単位解析)/2

- 全文検索におけるインデックスの単位
 - 形態素解析: 高精度, 検索漏れ
 - 文字単位/n-gram: 高再現率, 検索ノイズ
- 2つの立場を **融合, 単一化** できないか?
 - 応用によって2つの立場を無段階に選択する
 - いいとこどり



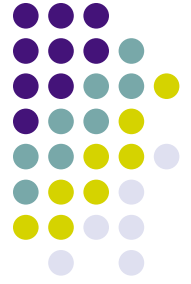


ソフト分かち書き: 動作例



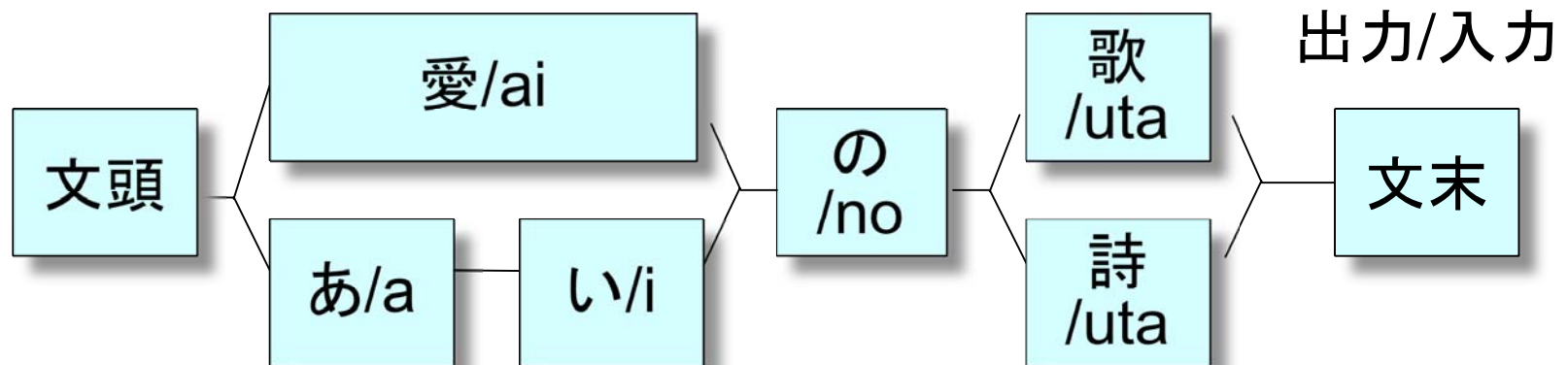
位置/単語	θ				
	10	1.0	0.5	0.1	0.01
0-2 京				■	■
0-4 京都		■	■	■	■
0-6 京都大			■	■	
0-8 京都大学	■	■	■	■	
2-4 都				■	■
4-6 大				■	■
4-8 大学		■	■	■	■
6-8 学			■	■	■

入力 「京都大学」



汎用テキスト処理ツール

- MeCab は日本語形態素解析器だけではない
 - 汎用的に作っています!
- テキスト → テキストの汎用変換ツール
 - 仮名漢字変換 (*mecab-skkserv, AJAX IME*)
 - ローマ字→ひらがな
 - 文字コード変換 (ちと強引)
 - 適切に辞書/コスト値を作れば実現可能!





MeCab の辞書

1. dic.csv (辞書定義)

```
の,166,166,8487,助詞,格助詞,一般,*,*,の,ノ,  
京都,1306,1306,1849,名詞,固有名詞,地域,一般,*,*,京都,キョウト,キョート  
桜,1304,1304,7265,名詞,固有名詞,人名,名,*,*,桜,サクラ,サク  
....
```

- 単語, 左文脈id, 右文脈id, 単語生起コスト, 素性列(CSV)
- 素性は任意の情報(品詞,活用,読み等)をCSVで記述

2. matrix.def (接続コスト定義)

1306	166	-2559
1304	1303	401
166	1304	608

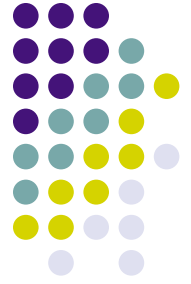


左文脈id 右文脈id 単語接続コスト

単語接続コスト
単語生起コスト

3. char.def (文字の定義)
4. unk.def (未知語処理の定義)
5. dicrc (出力フォーマット等)

AutoLink



自動的にリンクが張られる機能です
MeCabで実現できます。

1. dic.csv (辞書定義)

```
リンク,0,0,-500,http://foo.com/  
MeCab,0,0,-200,http://mecab.sourceforge.jp/  
Gree,0,0,-100,www.gree.jp  
....
```

- 接続は一状態
- 単語の長さに対し指数的に小さくなるコスト
- 素性にリンク先 URL

2. matrix.def (接続コスト定義)

```
0 0 0
```

接続は使わないので一状態
コスト 0

3. char.def (文字の定義)

4. unk.def (未知語処理の定義)

→ デフォルト 1文字 1未知語

5. dicrc

```
node-format-autolink = <a  
href="%H">%M<a>  
unk-format-autolink = %M
```

%M: 単語 (入力)

%H: 素性 (出力)

T9風 予測入力



1/あ	2/か	3/さ
4/た	5/な	6/は
7/ま	8/や	9/ら
	0/わ	

入力: 1681 → おはよう, 241 → ください
語呂合わせ: 1192 → 哀楽, 794 → 森田

1. dic.csv (辞書定義)

```
1,10,10,0,オ  
2,11,11,0,カ  
2,12,12,0,ガ  
....
```

- 単語(入力): 数字
- 文脈id: すべてのカタカナ文字に対応

2. matrix.def (接続コスト定義)

```
10 9 2505  
10 10 396  
10 11 606  
10 12 964  
...
```

- wikipedia を mecab で解析
- 単語の読みと頻度を取得
- カタカナのつながりやすさをコスト化
- 「日本語らしさ」

3. char.def (文字の定義)

4. unk.def (未知語処理の定義)

→ デフォルト 1文字 1未知語

5. dicrc

```
node-format-katakana = %H  
unk-format-katakana = %M
```

%M: 単語 (入力)

%H: 素性 (出力)



子音入力

dmdkry → だめだこりゃ
tnkyhu → てんきよほう
kdutk → くどうたく

1. dic.csv (辞書定義)

```
a,6,6,0,ア  
k,15,15,0,カ  
k,17,17,0,キ  
py,137,137,0,ピャ
```

- 単語(入力): 母音無しローマ字
- 文脈id: すべてのカタカナ文字に対応

2. matrix.def (接続コスト定義)

```
6 15 796  
6 17 675  
6 137 3121  
6 138 3121  
...
```

- wikipedia を mecab で解析
- 単語の読みと頻度を取得
- カタカナのつながりやすさをコスト化
- 「日本語らしさ」

3. char.def (文字の定義)

4. unk.def (未知語処理の定義)

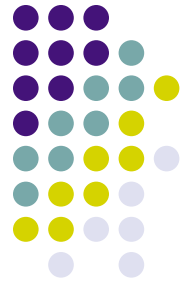
→ デフォルト 1文字 1未知語

5. dicrc

```
node-format-katakana = %H  
unk-format-katakana = %M
```

%M: 単語 (入力)

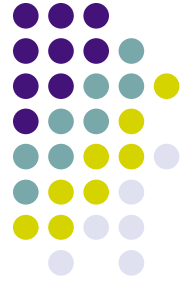
%H: 素性 (出力)



MeCab の素性フィールドの利用

- 辞書の素性は CSV なら何でも可能
- 単語にさまざまな付加情報を付与
 - 意味情報
 - 英語の訳
 - URL
 - スпамスコア
- スпамフィルタの例
 - 通常のスпамフィルタ
 - MeCab で解析 → 単語の抽出
 - 単語をキーにスпамスコア辞書をルックアップ
 - 辞書引きが2回! 辞書の付加情報として持っておけばMeCabだけでスпамスコアリングが可能

の,166,166,8487,助詞,格助詞,一般,*,*,*,の
桜,1304,1304,7265,名詞,固有名詞,人名,名,*,*,桜,サクラ
....



まとめ

- MeCabの技術
 - 辞書引き
 - 通常の hash は使えない
 - TRIE
 - 曖昧性の解消
 - 最小コスト法
 - 統計処理による正解データからの推定
- 設計方針
- 汎用性
 - テキスト変換ツール
 - 意外な使い方