

条件付確率場とベイズ階層言語モデルの 統合による半教師あり形態素解析

持橋大地* 鈴木潤 藤野昭典

NTTコミュニケーション科学基礎研究所

daichi@cslab.kecl.ntt.co.jp

言語処理学会2011

2011-3-10(Fri), 豊橋技術科学大学

最初に

- 公式のPDFは完成度が充分でないため、
ご興味のある方は下の完成版をお取り下さい:
<http://chasen.org/~daiti-m/paper/nlp2011semiseg.pdf>
– または, タイトルで検索で簡単に見つかります

“不自然”言語の解析

- 従来の新聞記事コーパスでは扱えない言語データが増えている (口語, 新語, 新表現, ...)

Twitter



ichijohisato ひさっちゃん

@vi_hazuki 目がはなせないでしょ(´▽`)♡しないからさやかあたりを狙ってみる(*`ω`)えではないですよ。ただ普段から小説の言葉などしまうもので・・・まいったなあ。照

29 seconds ago



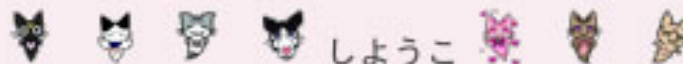
Hybrid_Soul_ 雑種魂 ~Hybrid Soul~

ごめんなさああああい(´;ω;`)!!!

45 seconds ago

Blog

いぬまるだしwwwwwまる出しwwwwww
声だして笑ってしまうwwwwwすさまじいww
物ですwwwwwおっおwwwwwゲーム、バク
ンプ等々つぼすぎるwwwww
3以降もAmazoったー！！！！HPが回復する
(´▽`)



– 全部人手で辞書登録... ? (Brain damaged!)

話し言葉の解析

CSJ話し言葉コーパスの一部

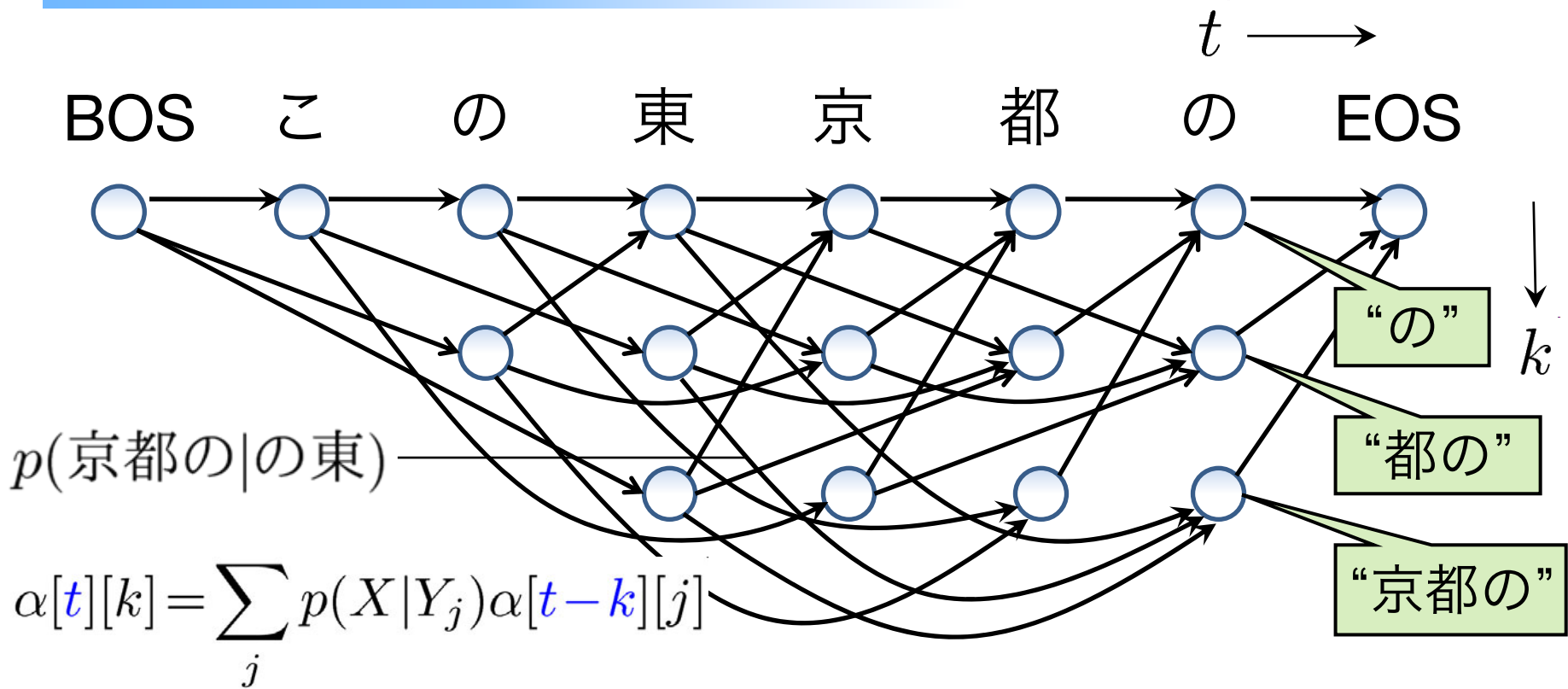
何て言うんでしょよねその当時はそう思われていたっていうことを全部ドラマにしちゃうっていうところそういうところとかが面白くて凄く見てたんですけど最初の方は凄くまともで緋色の研究とか四人の署名とかそういうのから始まってたんですけどそれはもうとにかく原作に沿っててこういうこれこれホームズってこれってというのが見たくて私はずっと見てます見てましたで凄くこれは何かNHKが放送されてる時から...

- 多数の口語表現 (“そいで”, “ってさあ”, “ちゃって”...)
 - 無数のバリエーションが存在
- 自然な話し言葉の音声認識や音声科学のために、**長期的にきわめて重要**

教師なし形態素解析 (持橋+, ACL2009)

- 生の文字列だけから, 階層ベイズで「単語」を学習
 - モデル: NPYLM (Nested Pitman-Yor LM)
 - 1 神戸では異人館街の二十棟が破損した。
 - 2 神戸では異人館街の二十棟が破損した。
 - 10 神戸では異人館街の二十棟が破損した。
 - 50 神戸では異人館街の二十棟が破損した。
 - 100 神戸では異人館街の二十棟が破損した。
 - 200 神戸では異人館街の二十棟が破損した。

NPYLM as a Semi-Markov model



- **Semi-Markov** HMM (Murphy 02, Ostendorf 96)の教師なし学習+MCMC法
- 状態遷移確率(nグラム)を超精密にスムージング

教師なし→半教師あり学習

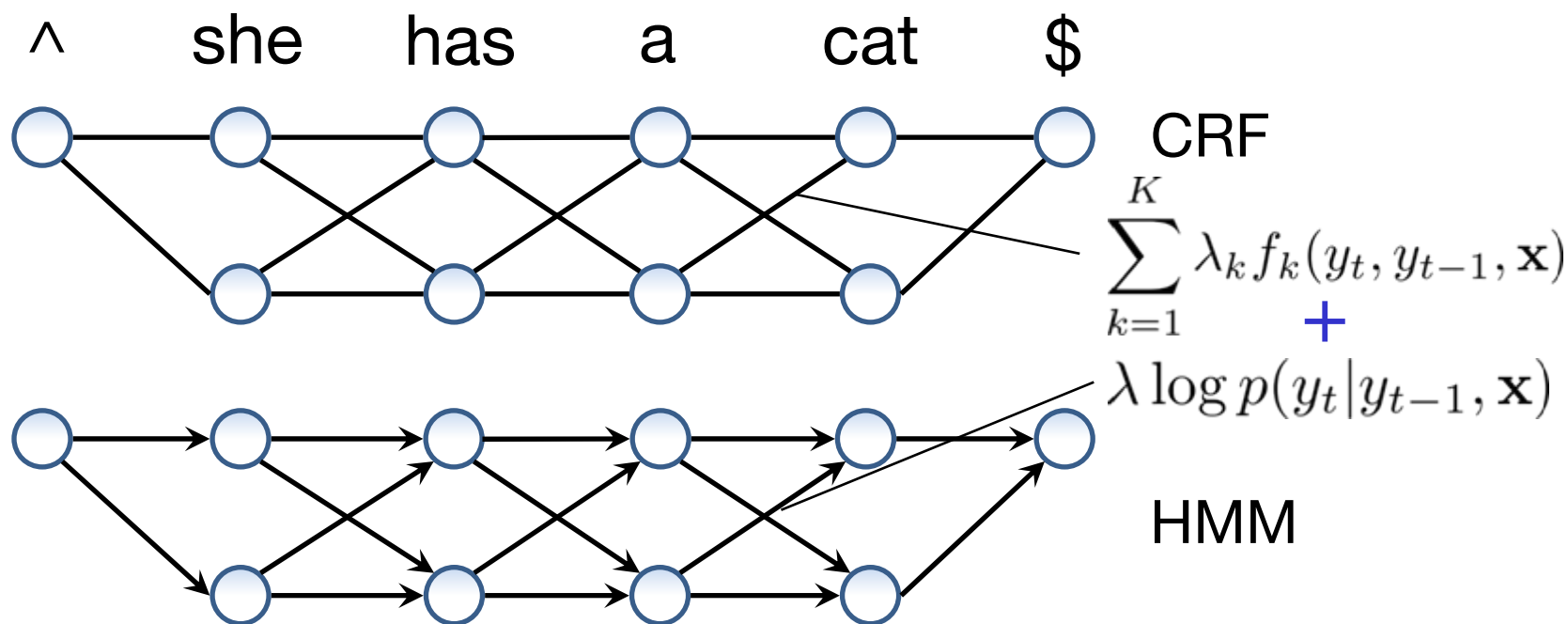
- 教師なし学習は低頻度語, 人間の基準に弱い
 - “阪急桂駅”, “首都グロズヌイ”
 - “歌う”→“歌う”, “静かな”→“静かな”
- 生成と識別モデルの結合確率モデル: JESS-CM

$$p(\mathbf{y}|\mathbf{x}; \Lambda, \Theta) \propto p_{\text{DISC}}(\mathbf{y}|\mathbf{x}; \Lambda) p_{\text{GEN}}(\mathbf{y}, \mathbf{x}; \Theta)^\lambda$$

識別モデル 生成モデル 重み入も学習

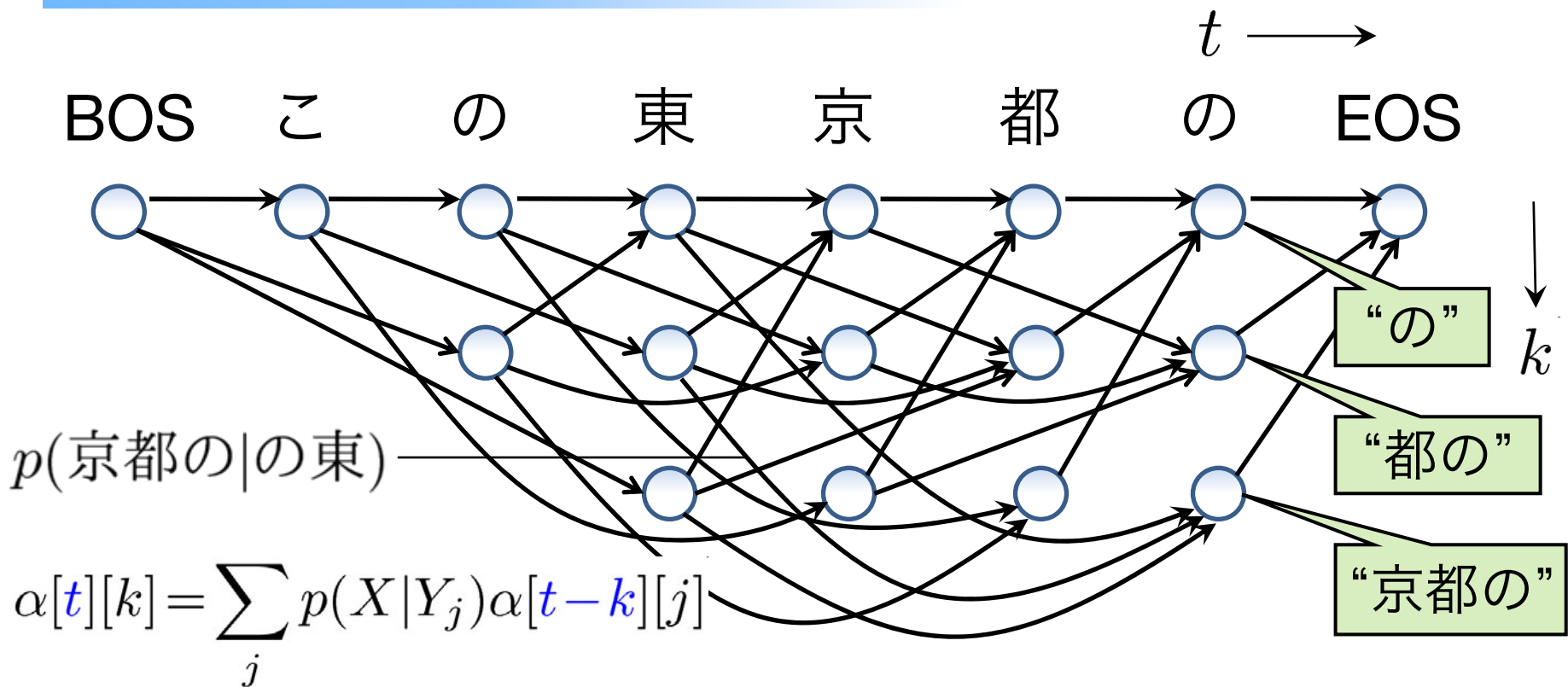
- joint probability model embedding style semi-supervised conditional model (鈴木+ ACL08/09)
- 現在世界最高性能の半教師あり学習
 - CRF/HMM, CRF/Naive Bayesなど

JESS-CM on CRF/HMM (鈴木+, ACL2008)



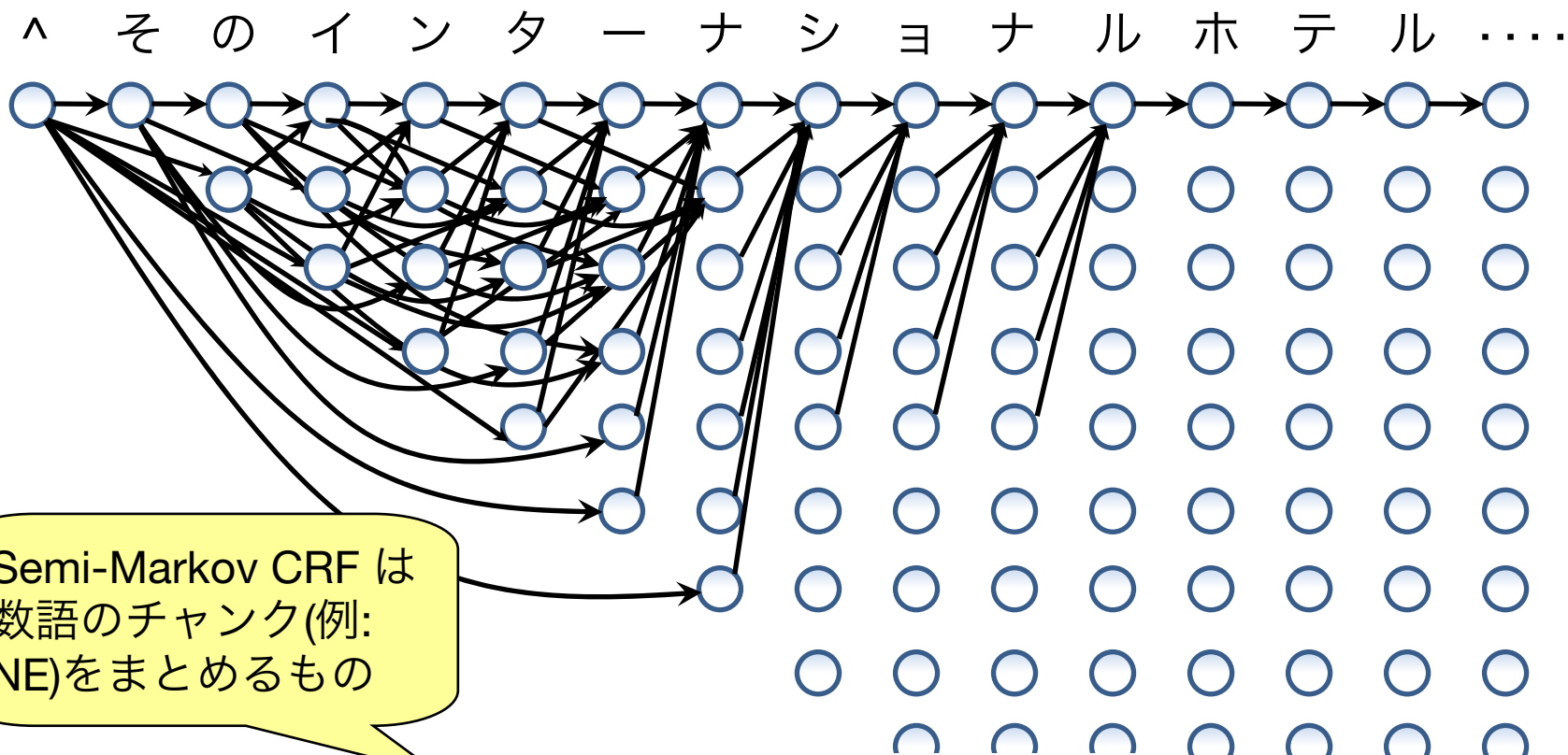
- 同じグラフィカルモデル上で, パスのコストを重みつきで足し合わせる
→ CRFとHMMを交互に学習(「教え合う」)
- NPYLMの場合は? (Semi-Markov)

NPYLM as a Semi-Markov model



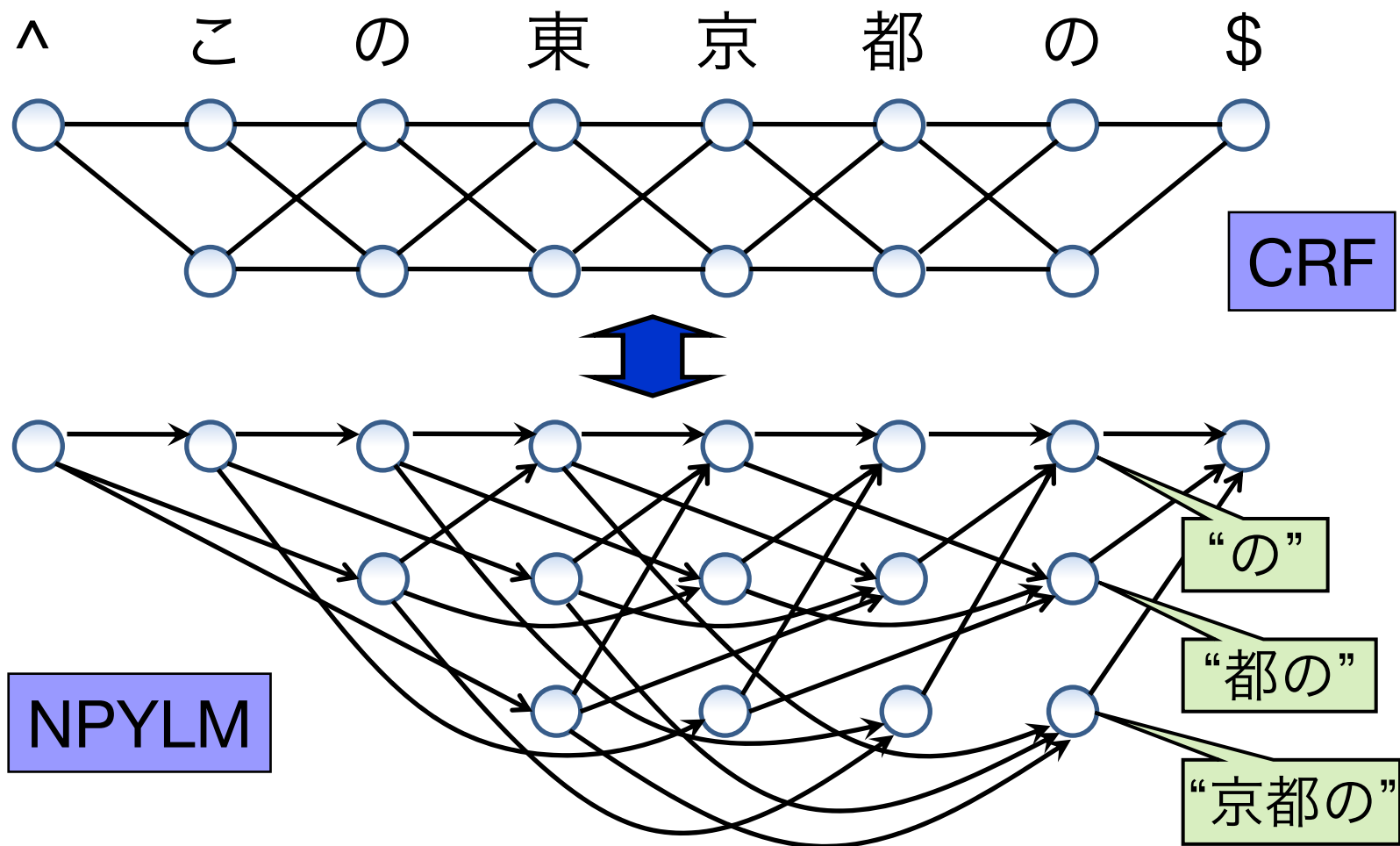
- Semi-Markov HMM (Murphy 02, Ostendorf 96)の教師なし学習+MCMC法
- 状態遷移確率(nグラム)を超精密にスムージング

Semi-Markov CRF (NIPS 2004)を適用?



- **膨大なメモリ** (1GB→20GB)
- (教師あり学習の)精度: 高々 **95%**
 - 単語のみ、文字レベルの情報なし

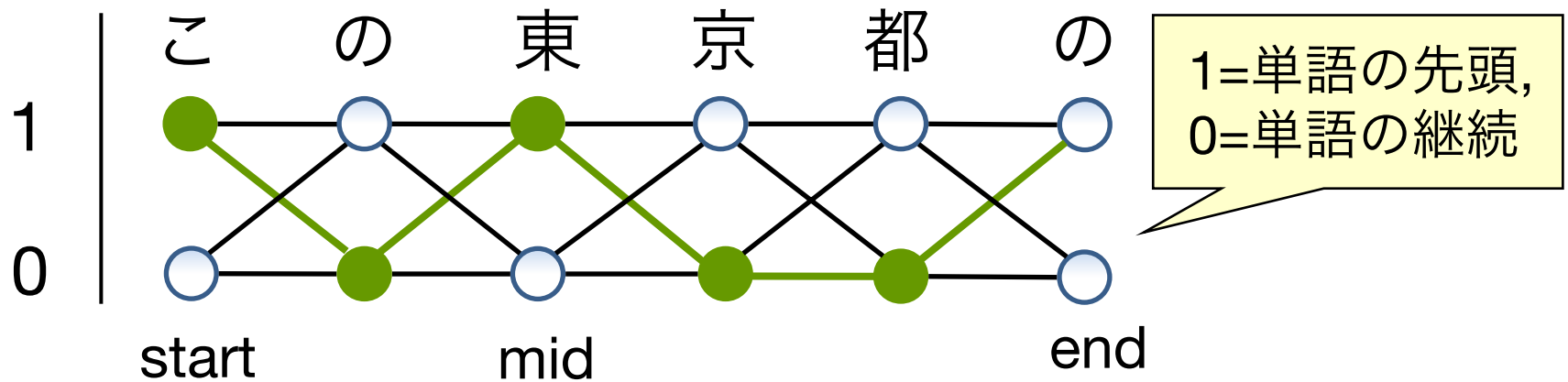
Markov CRF \leftrightarrow Semi-Markov LMの学習



- どうやって2つの違うモデルを統合するか？

CRF→NPYLM

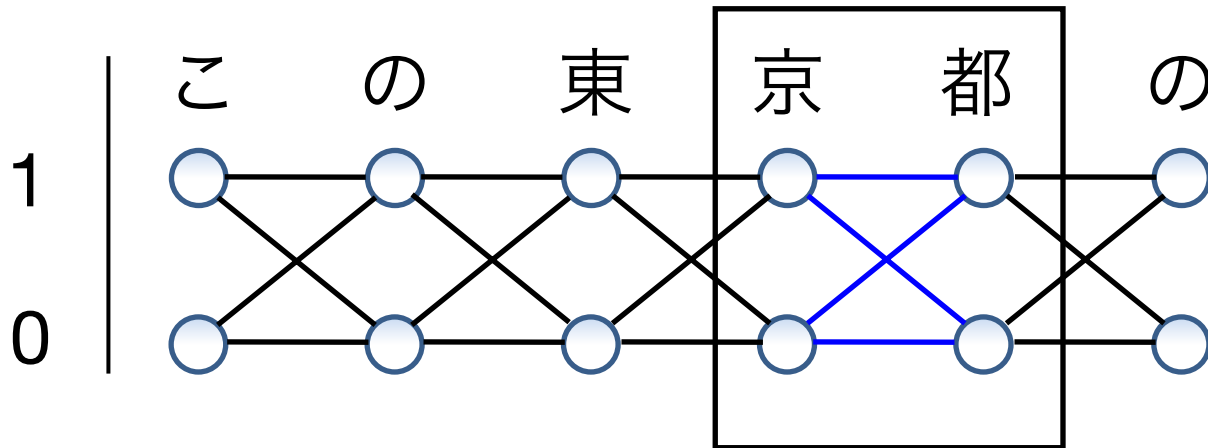
- Andrew+(EMNLP 2006)で既出、易しい
 - CRF→semi-Markov CRFの変換
 - $p(\text{“この”} \rightarrow \text{“東京都”})$



- 上のパスに沿って素性の重みを足し合わせる
- $\gamma(\text{start, mid, end})$
 $:= \gamma(\text{start, mid}) + \gamma(\text{mid, end})$

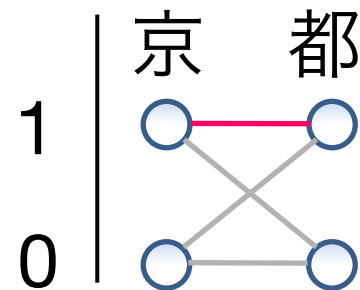
NPYLM→CRF (1)

- 難しい!!



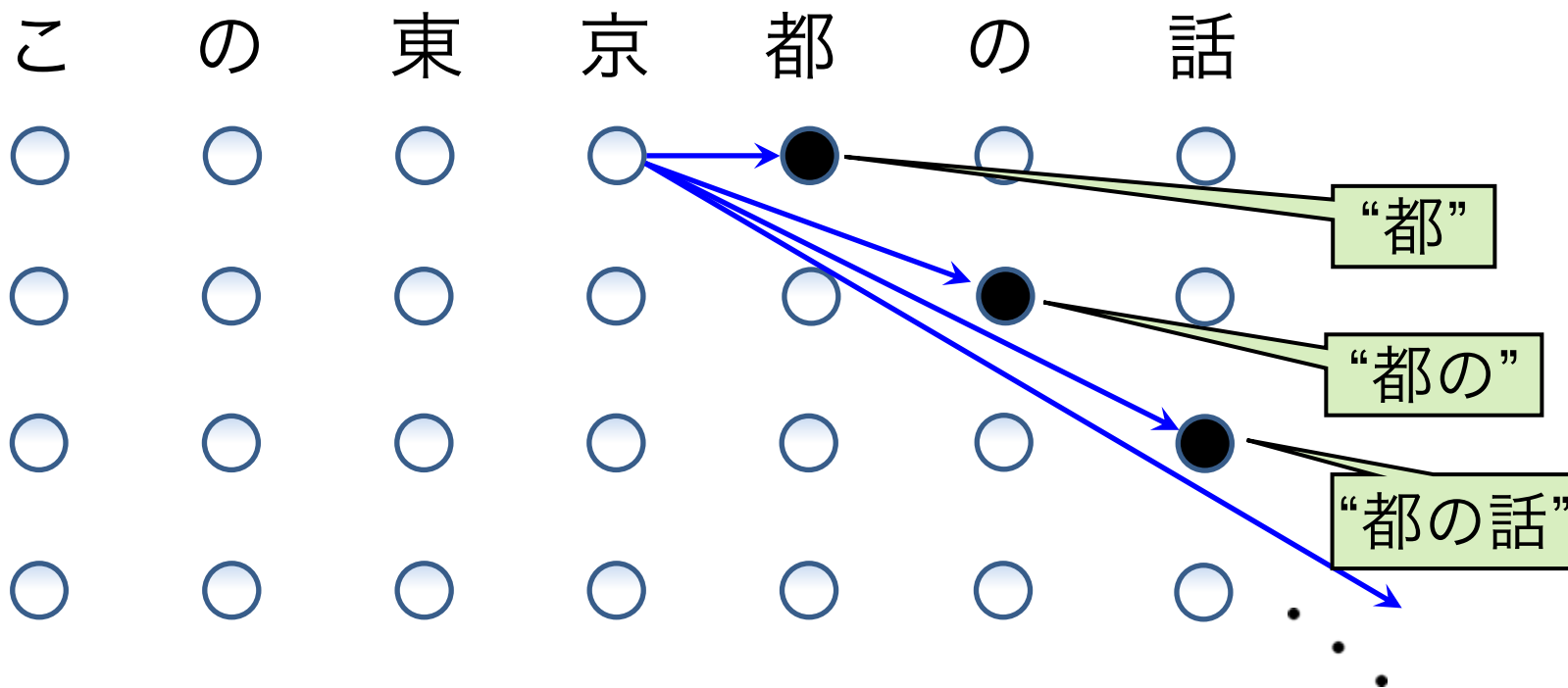
- 4通りのパス: $0 \rightarrow 0$, $0 \rightarrow 1$, $1 \rightarrow 0$, $1 \rightarrow 1$
- モデルがもしMarkov(=HMM)だった時の重みは、文 x が与えられれば、Semi-Markovの確率を複雑に足し合わせることで**計算可能!**

NPYLM→CRF (2)

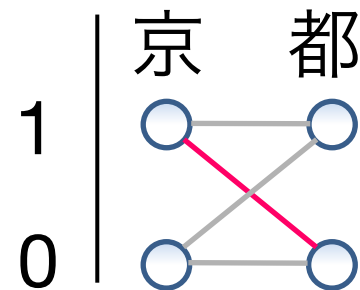


- Case 1→1 :

1→1 = “京→都”, “京→都の”, “京→都の話”, ...

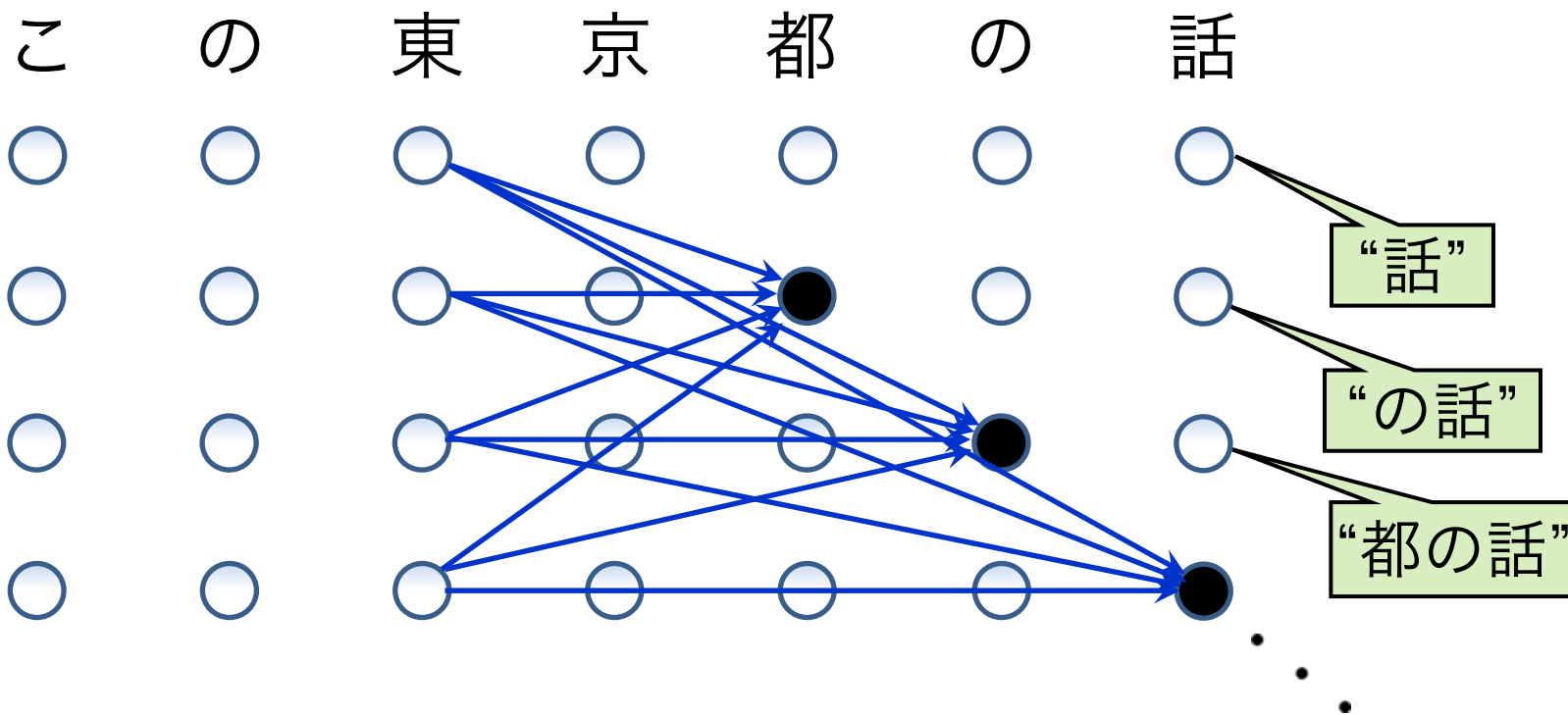


NPYLM→CRF (3)



- Case 1→0 :

1→0 = “東→京都”, “の東→京都”, “この東→京都”,
“東→京都の”, “の東→京都の”, “この東→京都の”,
“東→京都の話”, “の東→京都の話”, ……



NPY→CRF: Code example

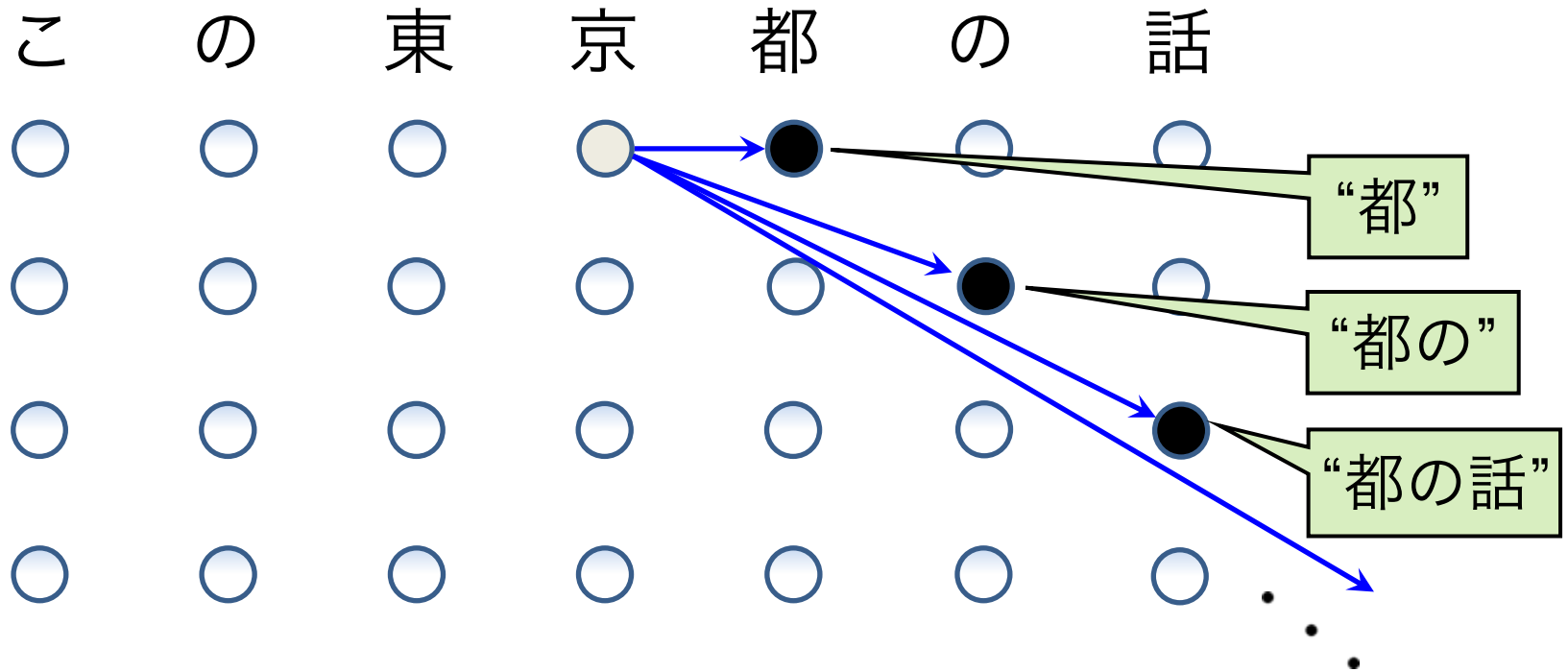
- “0→0”のポテンシャルを計算するC++コード

```
double
sentence::ccz (int t, HPYLM *lm)
{
    wstring w, h;
    int i, j, k, L = src.size();
    double z = 0;

    for (k = 0; k < MAX_LENGTH - 2; k++) {
        if (!(t + 1 + k < L)) break;
        for (j = 2 + k; j < index[t + 1 + k]; j++) {
            w = src.substr(t + 1 + k - j, j + 1);
            if (t + k - j < 0) { /* (t + 1 + k - j) - 1 */
                h = EOS;
                z += lm->ngram_probability (w, h);
            } else {
                for (i = 0; i < index[t + k - j]; i++) {
                    h = src.substr(t + k - j - i, i + 1);
                    z += lm->ngram_probability (w, h);
                }
            }
        }
    }
    return z;
}
```

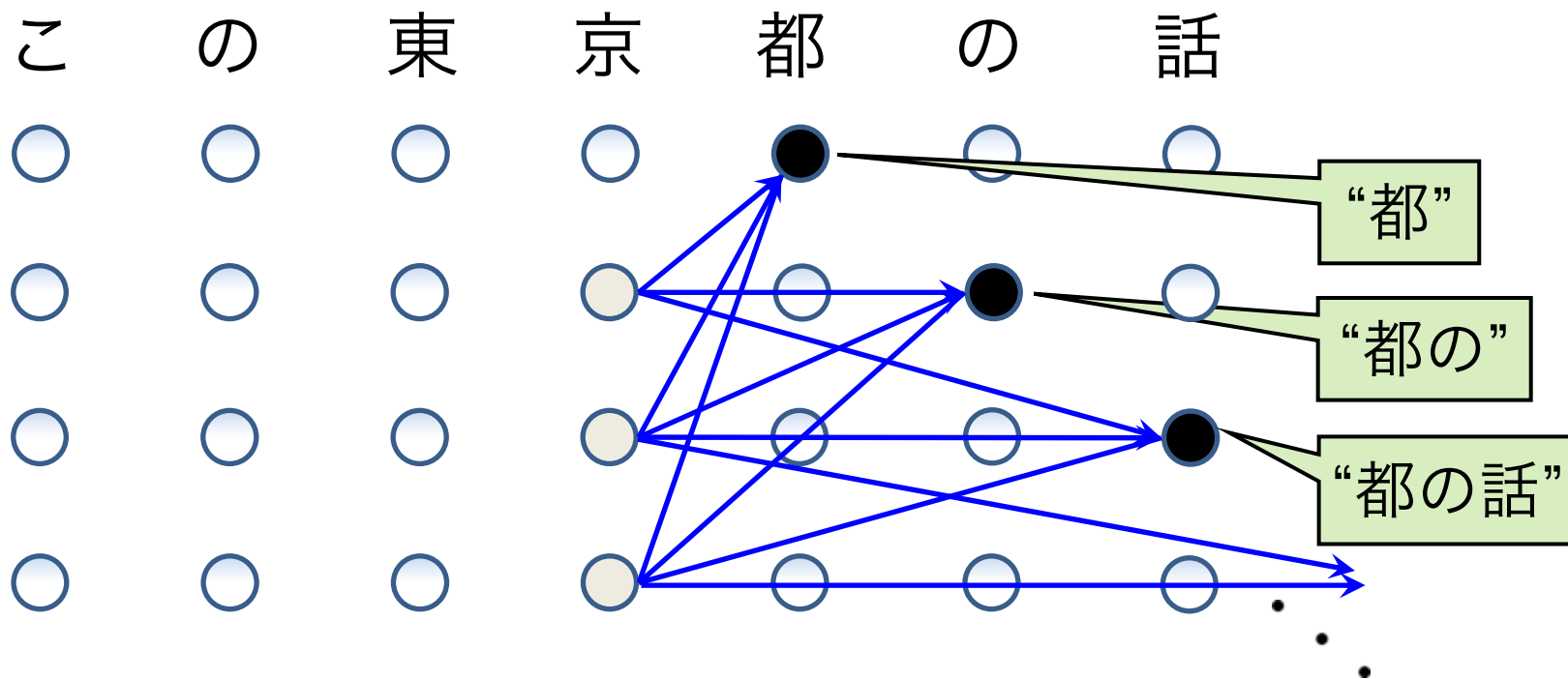

What are we doing? (1)

- グラフィカルモデル上で、 $1 \rightarrow 1$ に対応する確率の和をとる：



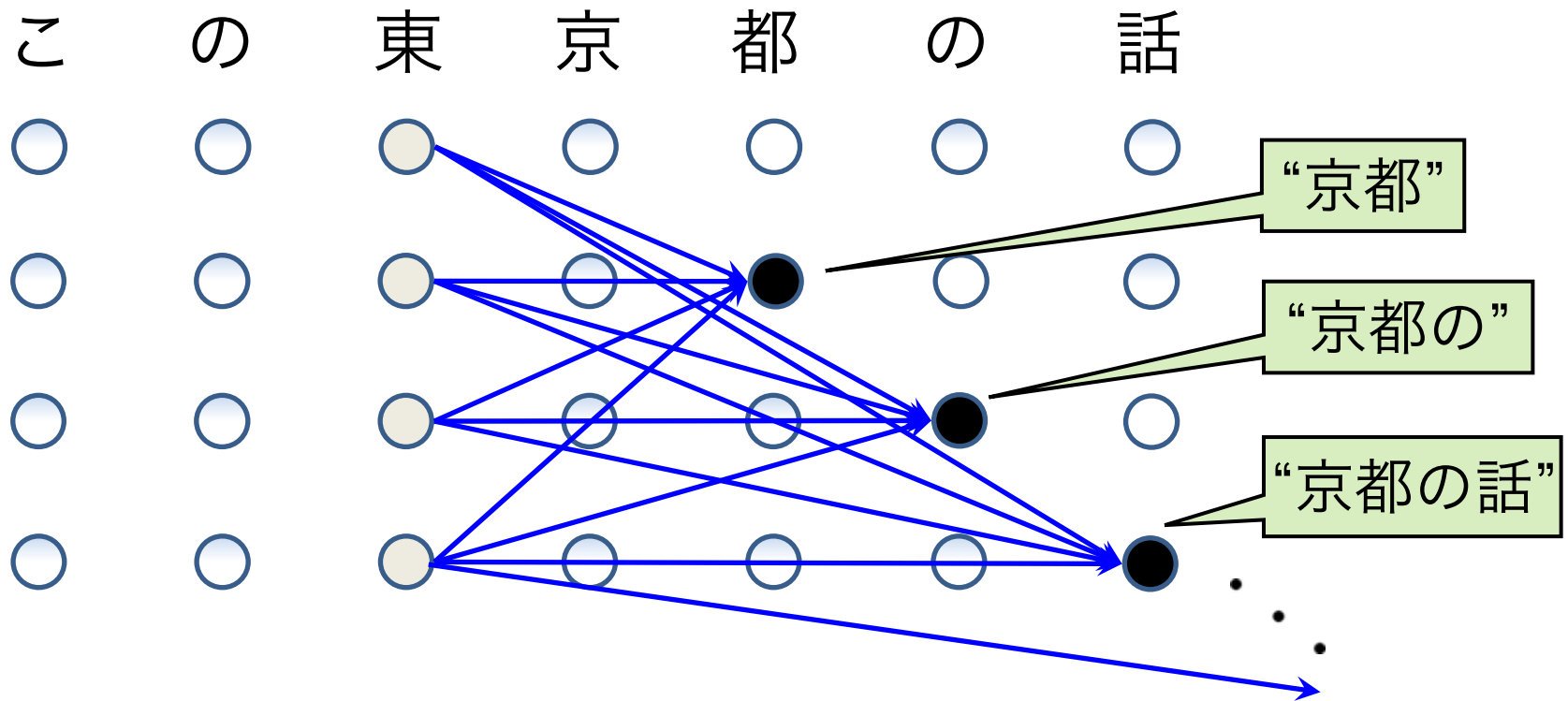
What are we doing? (1)

- グラフィカルモデル上で、 $0 \rightarrow 1$ に対応する確率の和をとる：



What are we doing? (1)

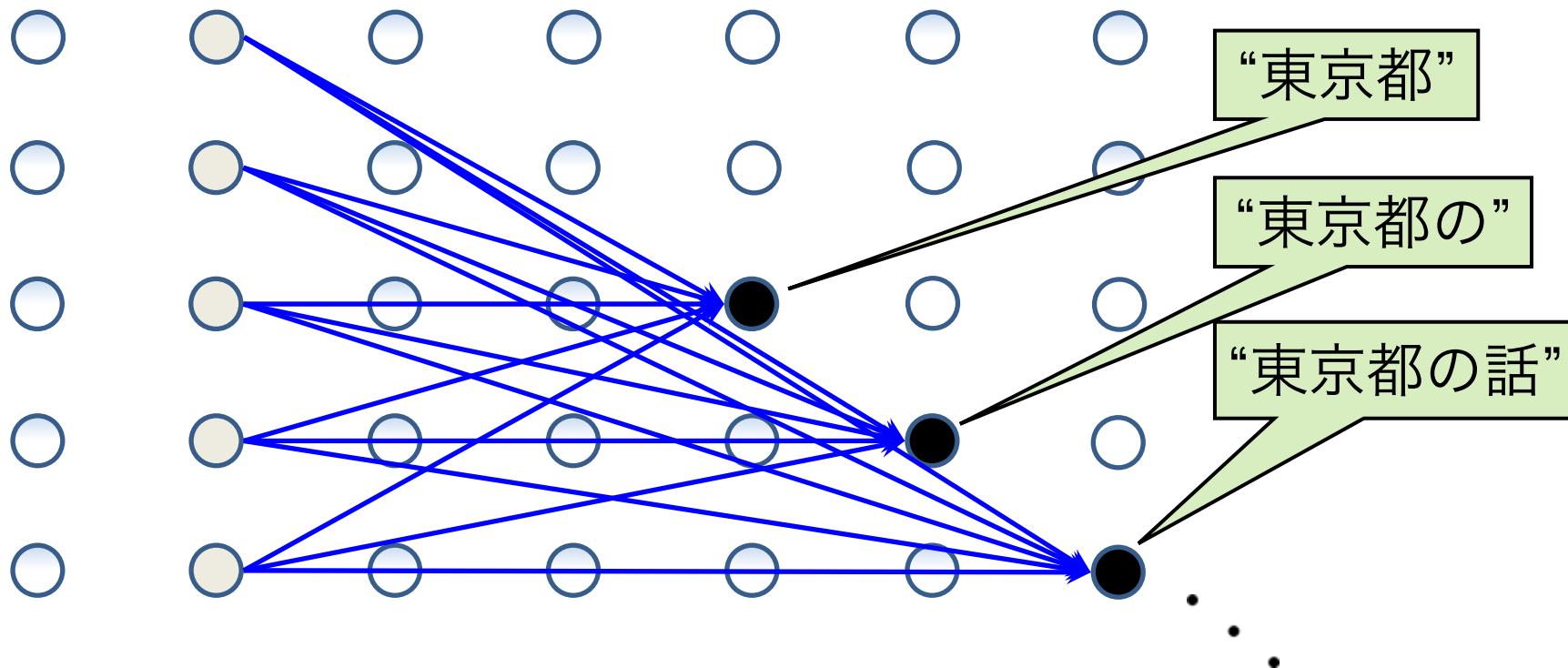
- グラフィカルモデル上で、 $1 \rightarrow 0$ に対応する確率の和をとる：



What are we doing? (1)

- グラフィカルモデル上で、 $0 \rightarrow 0$ に対応する確率の和をとる：

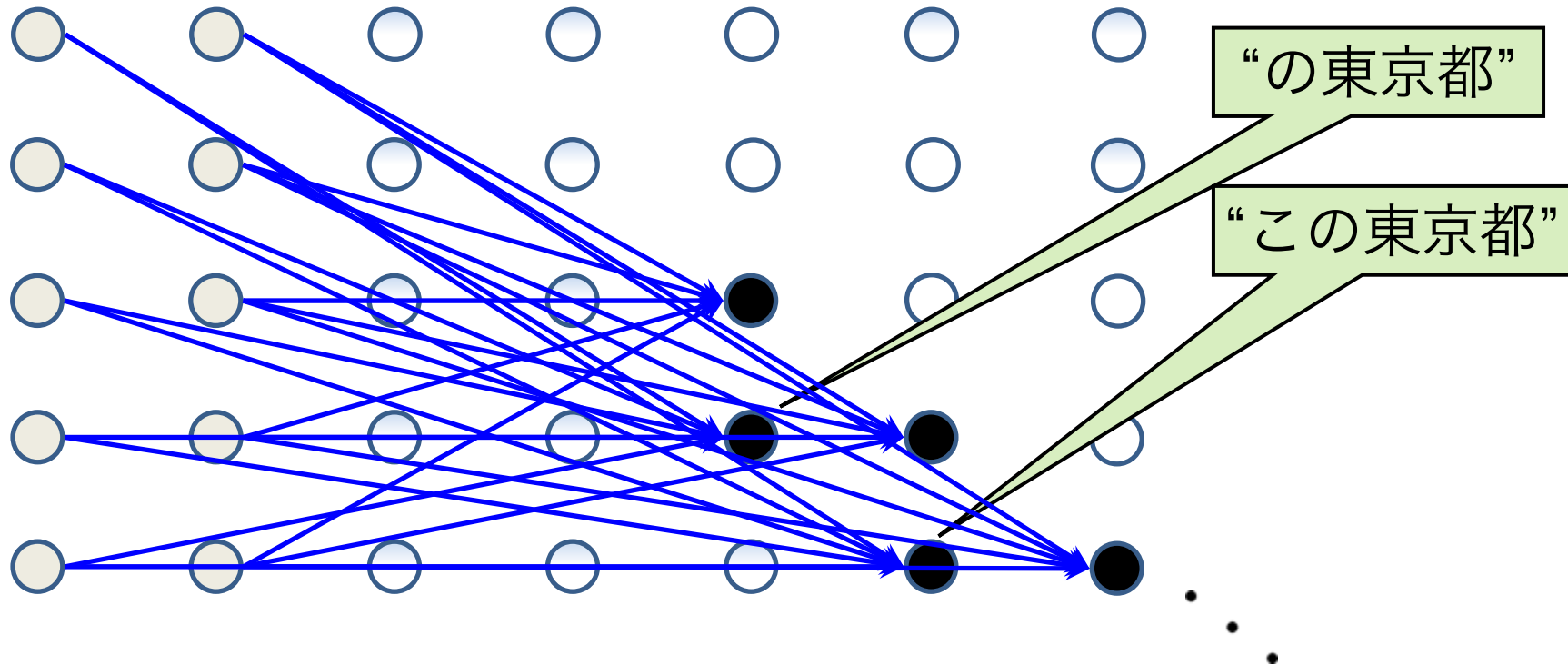
こ の 東 京 都 の 話



What are we doing? (1)

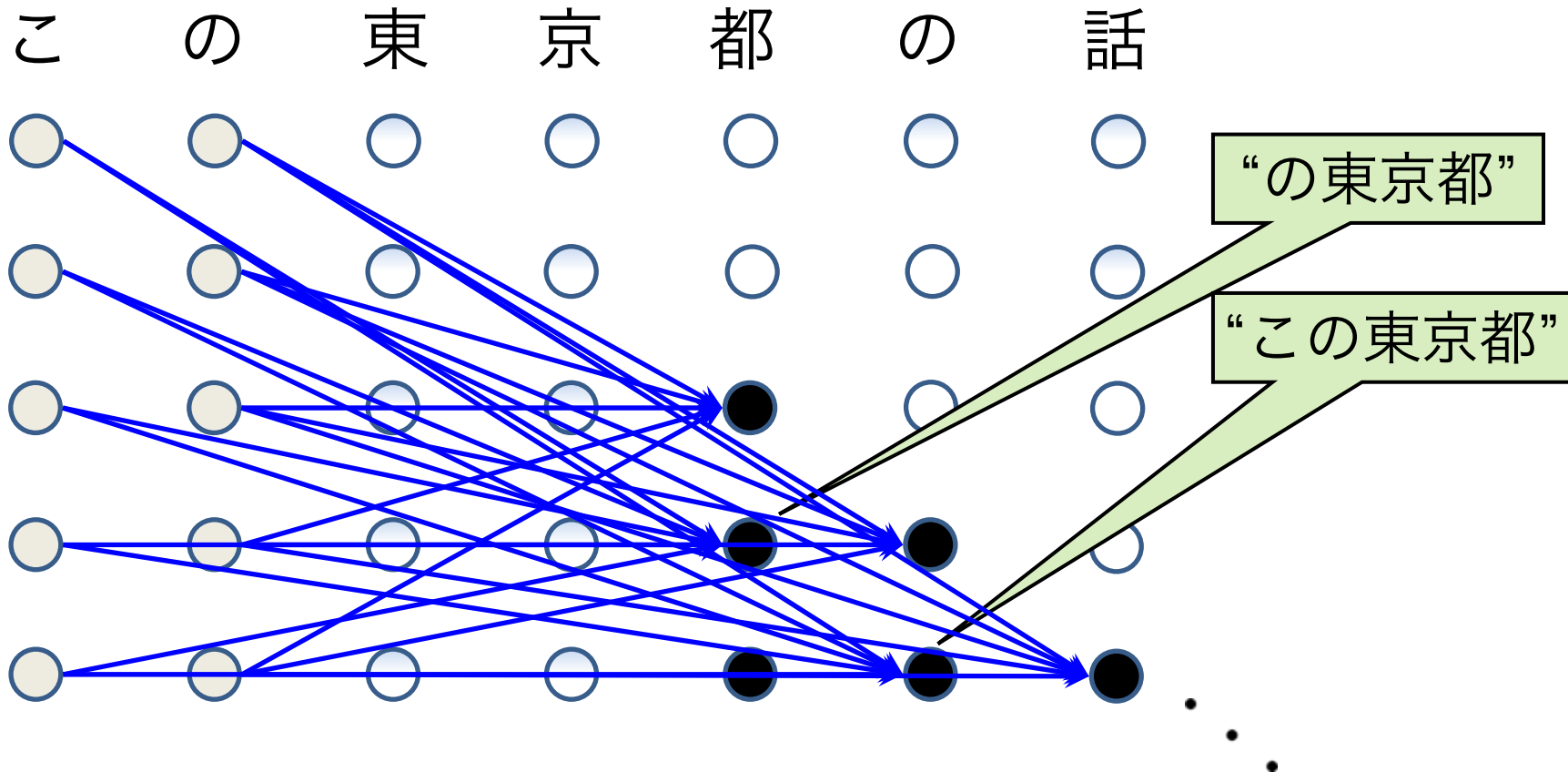
- グラフィカルモデル上で、 $0 \rightarrow 0$ に対応する確率の和をとる：

こ の 東 京 都 の 話



What are we doing? (1)

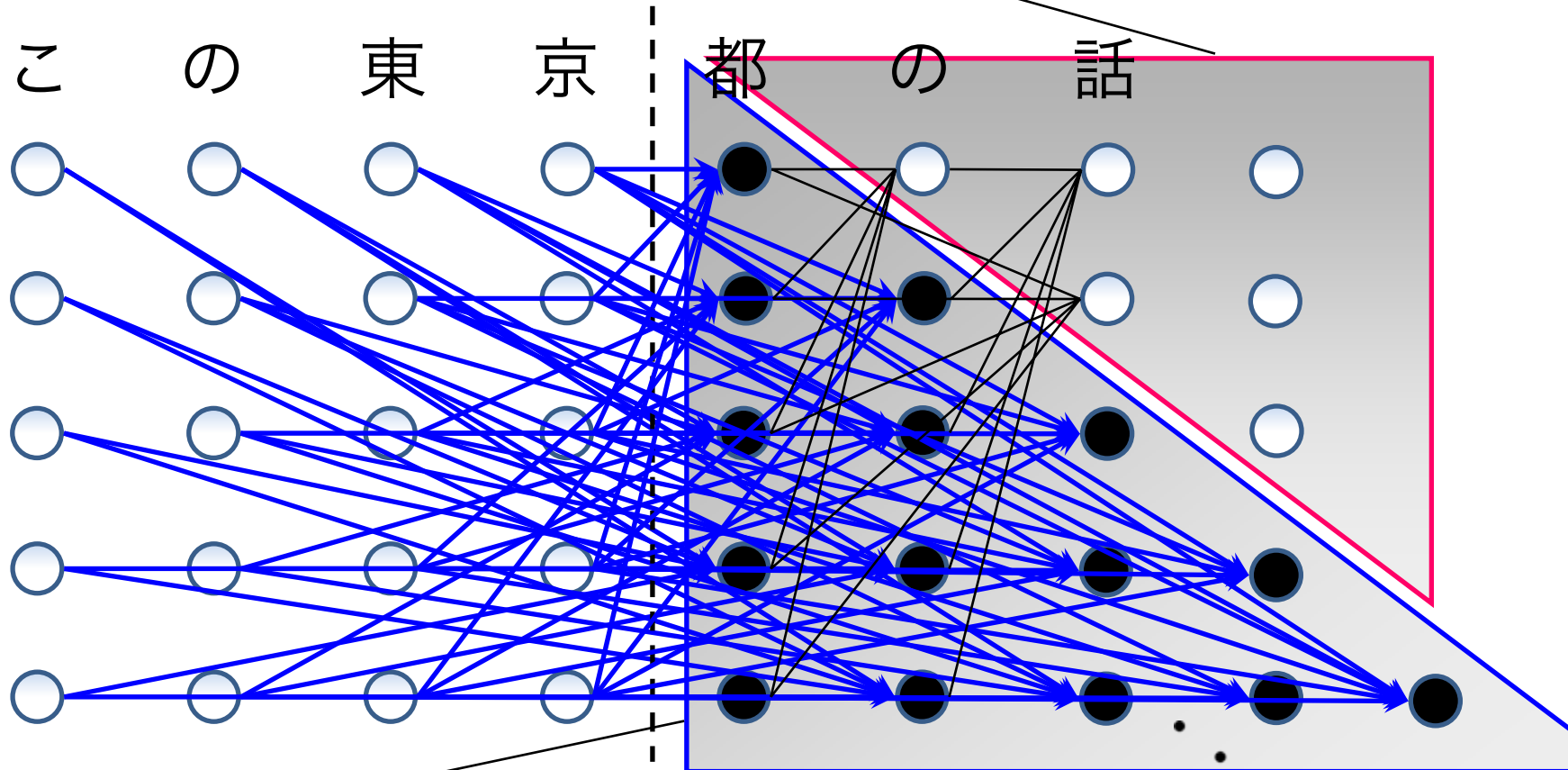
- グラフィカルモデル上で、 $0 \rightarrow 0$ に対応する確率の和をとる：



What are we doing? (1)

- DAG上で, 切断面を横断するパスを4種類に分類:

切断面 和に無関係なノード集合



和に関する終端ノード集合

What are we doing? (2)

$$p(x) = \sum_Y p(x, y)$$

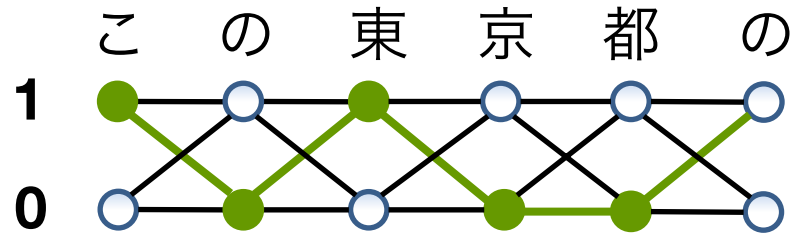
- 数学的には、この計算は確率の周辺化と等価

– 定義に従って、

$$p(c_t^{u-1} | c_s^{t-1})$$

$$\propto \gamma(s, t, u)$$

$$= p(\underline{z_s = 1}, z_{s+1} = 0, \dots, z_t = 1, z_{t+1} = 0, \dots, \underline{z_u = 1})$$



– 周辺化により、求める同時確率が得られる

$$p(z_t = 1, z_{t+1} = 1) = \sum_k p(\underline{z_t = 1}, \underline{z_{t+1} = 1}, \dots, \underline{z_k = 1})$$

$$p(z_t = 1, z_{t+1} = 0) = \sum_l \sum_k p(\underline{z_t = 1}, z_{t+1} = 0, \dots, \underline{z_k = 1}, \dots, \underline{z_l = 1})$$

$$p(z_t = 0, z_{t+1} = 0) = \sum_j \sum_l \sum_k p(\underline{z_{t-1} = 1}, z_t = 0, z_{t+1} = 0, \dots, \underline{z_l = 1}, \dots, \underline{z_j = 1})$$

Experiments (still ongoing)

- 新浪微博 (Sina Microblog)
 - 中国語圏のTwitter, 95000000人のユーザ
- 「しょこたんブログ」
 - “しょこたん語”で有名な、崩れた日本語blog
- CSJ日本語話し言葉コーパス (省略)
- SIGHAN Bakeoff 2005
 - 公開データセット、新聞記事

Tremendous!

「しょこたんブログ」の分割

中三のとき後楽園遊園地にタイムレンジャーショーを見に行きまくってたころのことそうだおね、セル画は修正に全て塗り直しとかあるだろうけどデジタルなら一発でレイヤー直せるから... 戸田さんの生歌声が最高に美しくてチャーム状態になりました。そしてザ・ピーナッツ役の堀内敬子さん瀬戸カトリーヌさんが息ピッタリに渡辺プロのそうそうたる名曲を歌いあげ、最高のハーモニーでとにかくすばらしかったです。生歌であの美しさ...。四つとも全部この川柳wwwwwwお茶wwwwwwイトカワユスwwwwwwイトカワユスwwwwww(^ω^)(^ω^)(^ω^)深夜までお疲れさまミタス(°ω°)ギャル曾根たん！最近よく一緒になると楽屋に遊びにきてくれるのでいろいろおしゃべりしてタノシス！(^ω^)今日もいろいろ話したおねイプサの化粧水がケア楽チンだし肌にギザあう！これギガント肌調子よくなりました(^ω^)

- 教師あり: 京大コーパス 37,400文
教師なし: しょこたんブログ 40,000文

しょこたんブログの「単語」

● 頻度2以上の単語をランダムに抽出

あるるるる	2	スワロフスキー	3	早いし	2
ますえ	2	わたる	11	信じろ	6
そびれちゃった	2	コマ送り	3	似てる	26
メリクリスマース	3	おおっお	7	居る	10
シクシク	3	にじむ	4	よる	85
チーム	45	簿	12	LaQua	7
ロック	11	ギギ	2	ただただ	7
キムタク	12	呼んで	29	ストロベリメロディ	21
うなあ	2	席	31	スター———トウハツツツ	2
したろう	3	100	55	ひろがって	3
去った	4	グラビア	85	しろま	3
死兆星	4	田尻	3	カワユスピンク	2
スッキリ	6	より焼き	2	海馬	3
ドバァア	2	ヒヤダルコ	3	除外	3
開催	47	永久	34	けえ	6
おく	17	ヤマト	2	なんとゆう	2

新浪微博 (Sina microblog)

中国／香港／台湾での
超有名サイト(Twitter)

今天一大早就被电话吵醒了，折磨死我了，昨天太晚睡了，早上被这电话搞的晕忽忽！

头疼，发热。。貌似感冒了，晚上睡觉不能裸睡了。要穿睡衣了。咿~？半个钟前发的围脖咋不见了咧~~只是感慨了一下今天的归途特顺嘛~~~(ノ~~~~ノ)b

下雨了，不知道广州那边有没有下雨，明天的同学聚会我去不了了，[伤心]大哭

學校付近一隻很可愛的狗狗，做了點特效[心][心][心]我們學校學生超愛牠的！！！！[哈哈]

明儿我要把中山陵搞定~~~~玛丽隔壁的~~~(ノ_ノ)

好饿啊....走！妈妈带你出去吃饭去~.....(((((((ヾ(o=^·ェ·)

o┐┌ 喵~o(=∩ω∩=)m

梦。。。混乱的梦。。。清晰的梦。。。。。

- 教師あり: MSR 87000 文 (普通話)
教師なし: Sina API, 98700 文

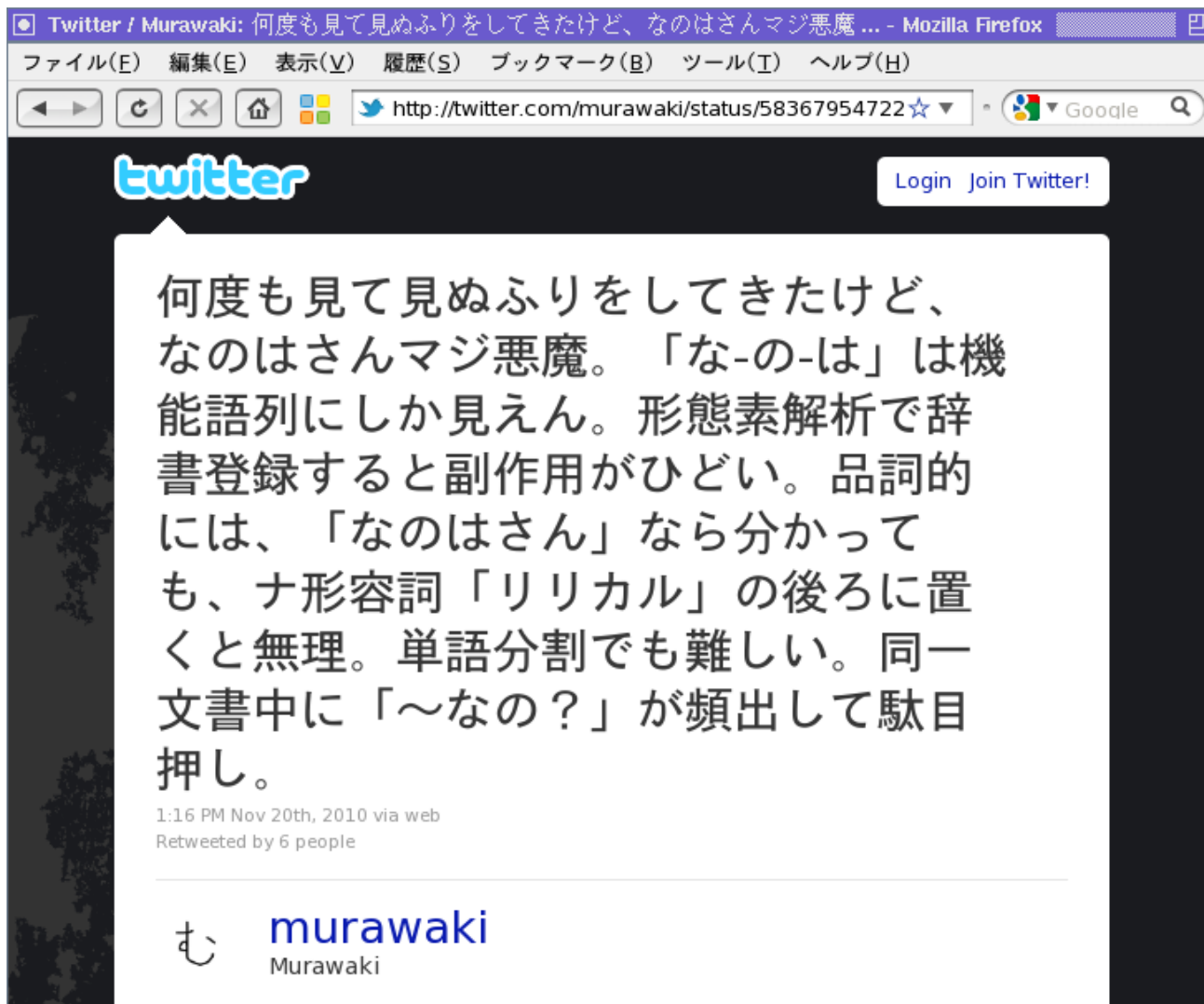
SIGHAN Bakeoff 2005

- 中国語単語分割の標準データセット
 - 新聞記事のため, 半教師あり学習の目的とは
やや異なるが...
- データ: MSR Asia 87k+ Chinese Gigaword 200k

Model	CRF	NPYCRF	+辞書
Token F値	97.4	97.5	97.5
OOV 再現率	83.5	84.1	82.1
IV 再現率	98.5	98.6	98.8

- 改善しているが, 大きな改善にはもっとデータ量が必要 (現在2004年の約半年分)
- ベースラインの97.4%は教師ありでは世界最高

解析が非常に困難な例 [注: 以下お祭りです]



Twitter / Murawaki: 何度も見て見ぬふりをしてきたけど、なのさんはマジ悪魔... - Mozilla Firefox

ファイル(E) 編集(E) 表示(V) 履歴(S) ブックマーク(B) ツール(T) ヘルプ(H)

http://twitter.com/murawaki/status/58367954722

twitter Login Join Twitter!

何度も見て見ぬふりをしてきたけど、なのさんはマジ悪魔。「な-の-は」は機能語列にしか見えん。形態素解析で辞書登録すると副作用がひどい。品詞的には、「なのさんは」なら分かって、ナ形容詞「リリカル」の後ろに置くと無理。単語分割でも難しい。同一文書中に「～なの？」が頻出して駄目押し。

1:16 PM Nov 20th, 2010 via web
Retweeted by 6 people

む murawaki
Murawaki

実際にやってみました

- データ: 2ch 魔法少女なのは総合スレ874-883
 - 10スレッド分(10000レス), 26474文
 - NPYCRF/K=12, 教師データ: 京大コーパス
- 結果: 難しい!

2004 年秋 に 放送 された 「魔法少女 リリカルなのは」と、
2005 年秋 に 放送 された 「魔法少女 リリカルなのは A's」、
2007 年春 に 放送 された 「魔法少女 リリカルなのは StrikerS」
について 語り ましょう。

なのはは基本的に個人戦だからほとんど描写されないが
なのはさんは仕事が恋人なワーカホリックです

- 不可能なのか? → No!

教師なし学習の結果

- 京大コーパス生文+なのはスレ26474文の教師なし学習

2004 年秋に放送された 「魔法少女リリカルなのは」 と、
2005 年秋に放送された 「魔法少女リリカルなのは A's」、
2007 年春に放送された 「魔法少女リリカルなのは StrikerS」
について語りましょう。

脚本 漫画 全て 都築作のなのは 出す気だったりして
なのは SS で雄犬がずっと犬形態なのはなぜ？
なんでなのはさんって俺の嫁なの？

防衛政策に不可欠なのは国民的合意であり、可能な限り...

- 一部完全ではないが、重要な部分は本来、完全に自動的に学習できる！

山形弁の解析

<http://www.nhk.or.jp/namara03-blog/>

- NHK 「今夜はなまらナイト」 の視聴者ブログ
をクロール, 重複を除いて6,000文

教師なし:

次回の放送、あっかもすんねすないがもすんねんだべけど、
いずれは全国放送に...(笑)
おばんかだっす。おらほの家では、生卵やとろろを入れて
食べます。ほんて、んまえぞ〜。
えっちゃんさ〜ん、是非やまがださきての〜。

半教師あり:

次回の放送、あっかもすんねすないがもすんねんだべけど、
いずれは全国放送に...(笑)
おばんかだっす。おらほの家では、生卵やとろろを入れて
食べます。ほんて、んまえぞ〜。
えっちゃんさ〜ん、是非やまがださきての〜。

何が問題?

- 半教師あり学習の一般的なモデル(JESS-CM):

識別モデル

生成モデル

$$p(\mathbf{y}|\mathbf{x}; \Lambda, \Theta) \propto p_{\text{DISC}}(\mathbf{y}|\mathbf{x}; \Lambda) p_{\text{GEN}}(\mathbf{y}, \mathbf{x}; \Theta)^\lambda$$

- 2つのモデルの補間重みが常に固定!

- 言語モデルの高精度な補間と同様な問題

- Jelinek-Mercer スムージング = 線形補間

$$p(w|v) = \lambda_1 p(w) + (1 - \lambda_1) p(w|v) \quad \leftarrow \text{重みは常に一定}$$

- 階層Bayes (Dirichlet) スムージング (MacKay 1994)

$$p(w|v) = \frac{f(v)}{f(v) + \alpha_0} \hat{p}(w|v) + \frac{\alpha_0}{f(v) + \alpha_0} \bar{\alpha}_w$$

HMM的な仕組みの導入が必要

重みが動的に変化!

まとめ

- CRFと教師なし形態素解析を組み合わせた、半教師あり形態素解析
 - JESS-CM法の枠組で「教え合う」
 - Semi-Markov(単語) \leftrightarrow Markov(文字) の情報相互変換
- 人手の基準を守りつつ、未知語を完全に自動的に認識
- モデル補間重みの文脈依存化が今後の課題 [半教師あり学習の一般的な問題]