
Unsupervised and Semi-supervised learning of Nonparametric Bayesian word segmentation

Daichi Mochihashi

NTT Communication Science Laboratories, Japan

daichi@cslab.kecl.ntt.co.jp

NOE Statistical Machine Learning Seminar

2011-1-19 (Wed)

ISM

Word segmentations in NLP

- Segment string into “words”

```
% echo “やあこんにちは, 統数研はどうですか。”
```

```
| mecab -O wakati
```

```
やあ こんにちは, 統数研 どうですか。
```

```
(やあこんにちは, 統数研はどうですか。) ✕
```

- Very important problem for unsegmented languages like Japanese, Chinese, Thai, Hebrew, ...

Chinese

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(I) ヘルプ(H)

http://www.tsinghua.edu.cn/qhdwzy/board_xxqk2.jsp?board=12&bid2=1201&pagen ☆ Google

 清华大学
Tsinghua University

学校概况 | 院系设置 | 管理机构 | 科学研究 | 教师队伍 | 人才培养 | 招生就业 | 虚拟校园

学校沿革
历任校长
统计资料

校长致辞



喜迎百年华诞 再铸新的辉煌

——2011年新年献辞

● 校长 顾秉林

一世纪沧桑砥砺，一百年春华秋实。此时此刻，2010年的余晖散去，2011年的曙光在前！此时此刻，清华正在送走她的第一个百年，迈向新百年的征程！在辞旧迎新、钟声回荡之际，请允许我代表学校，向清华全体同学、教职员工、离退休人员和海内外广大校友，向长期关心支持清华发展的各界人士，致以崇高的敬意和新年的祝福！

一百年来，清华大学的发展始终与国家民族的命运休戚与共，形成了优良的精神传统和鲜明的办学特色。一代代清华人“自强不息、厚德载物”，涌现出众多学术大师、兴业英才和治国栋梁，为中国社会进步和世界文明发展作出了重要贡献。特别是近年来，在国家的大力支持下，学校致力于世界一流大学建设，积极探索中国特色的“大学之道”，各项事业不断取得新的进展，正在跻身世界一流大学的行列。

大学之道，育人为本。一年来，以“清华新百年人才培养的使命与战略”为主题的第23次教育工作讨论会顺利举行，全校师生就推动办学优势转化、培养拔尖创新人才进一步取得共识。“清华学堂人才培养计划”以及多项教育教学改革措施相继实施。招生工作大力推进多元评价、兼顾拔尖与公平，生源质量进一步提高。就业毕业生中，超过80%选择到国家重要行业和领域建功立业。外国留学生规模不断扩大，结构进一步优化，外国研究生在学规模居全国高校首位。


大学之道，学术为魂。一年来，我校积极面向国际学术研究前沿和国家重大战略需求开展高


Thai

iGoogle - 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(I) ヘルプ(H)

http://www.chula.ac.th/ Google

 **ขอเชิญชาวจุฬาร่วมนิมกอยหลังสู่การเป็นเจ้าภาพ “จามจุรีเกมส์”**
กีฬามหาวิทยาลัยแห่งประเทศไทย ครั้งที่ 38 วันที่ 15-22 มกราคม 2554


 **Chulalongkorn University**
จุฬาลงกรณ์มหาวิทยาลัย
เสาหลักของแผ่นดิน

มหาวิทยาลัย
อันดับหนึ่ง
ของประเทศไทย
สามปีติดกัน

- ข่าวทั่ว
- ค้นหาด้วยเสิร์จเอนจินจิวจุฬาร
- เมนู


- หน้าหลัก
- สมัครรับจดหมาย
- บุคลากร
- ศิษย์เก่า


- สภามหาวิทยาลัย
- ติดต่อเรา
- แผนที่เว็บไซต์





ศูนย์ประสานงานสื่อมวลชน
แข่งขันกีฬามหาวิทยาลัยแห่งประเทศไทย ครั้งที่ 38

จุฬาร เปิดศูนย์ประสานงานสื่อมวลชน (Press Center) การแข่งขันกีฬามหาวิทยาลัยฯ ครั้งที่ ๓๘ “จามจุรีเกมส์”



 **สนับสนุนการเกิดไทยสาร**
ทางมหาวิทยาลัยจามจุรีเกมส์ ครั้งที่ 38
ระหว่างวันที่ 15-22 มกราคม 2554
by Pichai Sathaporn, Jiraporn Sathaporn, Jiraporn Sathaporn

 **ศาสตราจารย์**

 **ศูนย์ประสานงานสื่อมวลชน**
แข่งขันกีฬามหาวิทยาลัยแห่งประเทศไทย ครั้งที่ 38

การรับสมัคร คณะ สำนักวิชาและสถาบัน พุ่มาเขื่อน แนะนำจุฬา สื่อสารองค์กร จุฬากับนานาชาติ การค้นคว้าวิจัย กิจกรรมและการบริการสังคม

Persian



The image shows a screenshot of the IUT website in Persian. The browser window title is "isfahan university of technology - Google 検索 - Mozilla Firefox". The address bar shows "http://www.iut.ac.ir/". The website header features the IUT logo and the name "دانشگاه صنعتی اصفهان" in Persian calligraphy. Below the header, there are navigation links for "کتابخانه", "آموزشهای الکترونیکی", "سامانه الکترونیکی دروس", "سامانه اتوماسیون اداری", "سامانه گلستان", and "پست الکترونیکی". The main content area includes a search bar, a language selector (EN فارسی), and a list of news items. The news items are:

- اطلاعیه پذیرش و ثبت نام موقت دانشجویان مقطع دکتری ورودی نیمسال دوم 1389-90
- تغییر از اساتید دانشگاه صنعتی اصفهان در ششمین جشنواره تجلیل از پژوهشگران استان اصفهان
- قلمب علمی مطالعات گودگی آب و خاک دانشگاه صنعتی اصفهان در نمایشگاه دستاوردهای پژوهش و فناوری کشور
- بیانیه دانشگاه صنعتی اصفهان به مناسبت "اسلگرد نهم دی ماه روز حملنه و بصیرت"
- سومین شماره هفته نامه الکترونیک دانشگاه صنعتی اصفهان منتشر شد

At the bottom of the news section, there are four categories: "درباره ی دانشگاه", "افتخارات", "رویدادها", and "اخبار". On the right side, there are four sections:

- [دانشجویان]**: خدمات دانشجویی امور فرهنگی مرکز مشاوره
- [اساتید و کارکنان]**: اداره رفاه اداره کارگزینی بهداشت و درمان اساتید
- [سرویس ها و خدمات]**: سرویسهای چندرسانه ای رتیلد با مسنولین
- [فناوری اطلاعات]**: مرکز فناوری اطلاعات هسته محتوای دیجیتال برنال

(Isfahan university of technology, Iran)

Word segmentations in NLP

- Crucial first step for unsegmented languages like Japanese, Chinese, Thai, Hebrew, ...
 - Very important and fundamental problem
 - Especially: **Chinese** (1,200,000,000 speakers)
 - SIGHAN word segmentation Bakeoff: 2003-2010
 - Many intensive research!

Word segmentation (2)

- Learning methods so far: *Supervised*

```
# S-ID:950117245-006 KNP:99/12/27
* 0 5D
一方 いっぽう * 接続詞 * * *
、 * 特殊 読点 * *
* 1 5D
震度 しんど * 名詞 普通名詞 * *
は は * 助詞 副助詞 * *
```

Standard “Kyoto corpus”:
38400 **hand-segmented**
newspaper sentences

- “Correct” segmentations with **huge human effort**
 - CRF, SVM, Maxent, .. classifiers to learn
- “Correct” segmentations for **speech? Twitter?**
Unknown language..?
 - |女御|更衣|あ|また|さ|ぶら|ひ|た|ま|ひける|中|に|、|...

MeCab analysis of “*The Tale of Genji*”, AD1000

Overview

- Unsupervised learning of word segmentation
 - HPYLM² prior + Blocked MCMC
 - “Words” for even unknown language!
- Semi-supervised learning of word segmentation
 - Systematic integration with CRF
 - Markov \leftrightarrow semi-Markov conversion (general)
 - MCMC-EM like algorithm for inference

Unsupervised Word Segmentation

- **Basic idea:** Maximize the probability of the segmentation \mathbf{w} of a string s :

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} p(\mathbf{w}|s) .$$

- Ex. $p(\text{he sings a song}) > p(\text{hes ingsao ng})$
 - *No dictionary, no “correct” supervised data*
 - Find the “most natural segmentation” of a string
- **Note:** Exponential number of candidates
 - A sentence of 50 characters:
 $2^{50} = 1,125,899,906,842,624$ different segmentations

Probability of a sentence: n -grams

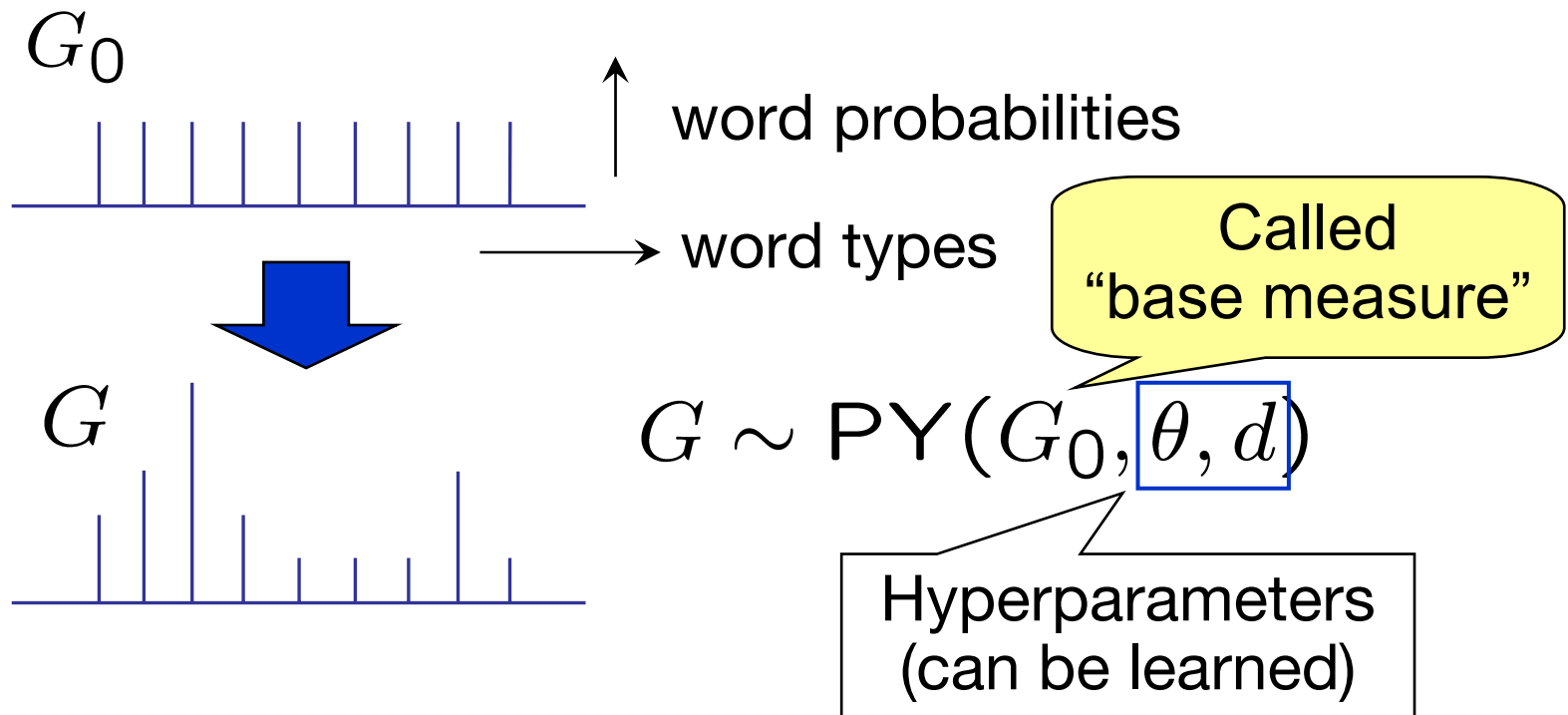
$p(\text{she likes music})$

$$= p(\text{she}|\wedge) p(\text{likes}|\text{she}) p(\text{music}|\text{likes}) p(\$|\text{music})$$

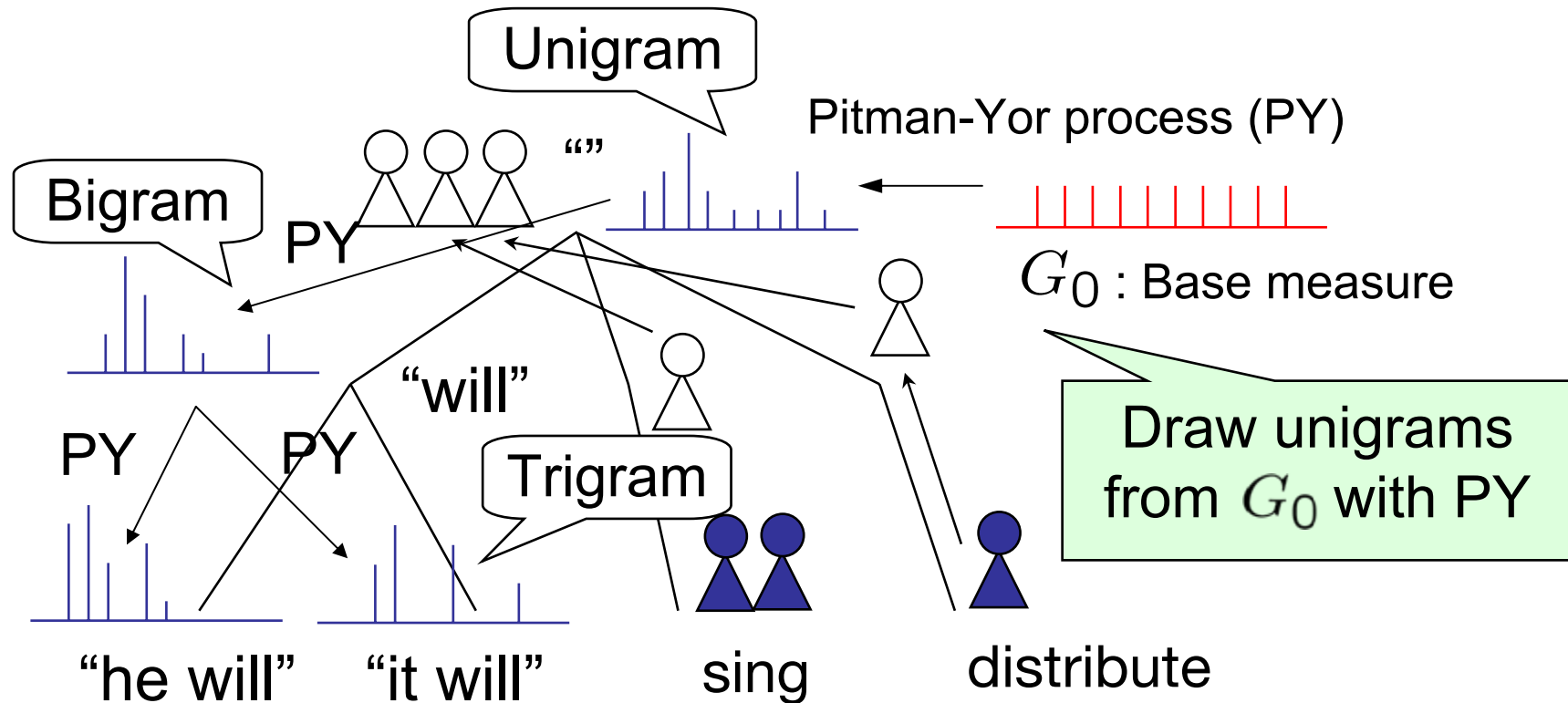
- Markov model on words: very common in NLP
 - First introduced by Shannon (1948)
 - PCFG is ok too.. but simple model will suffice
 - PCFG on Twitter yellings??
 - Probability tables are mostly 0
 - Hierarchical smoothing is necessary
 - Every substring could be a “word”
- ➡ Bayesian treatment: HPYLM (yw’s talk)

Pitman-Yor n-gram model

- The Pitman-Yor (=Poisson-Dirichlet) process:
 - Draw distribution from distribution
 - Extension of Dirichlet process (w/ frequency discount)

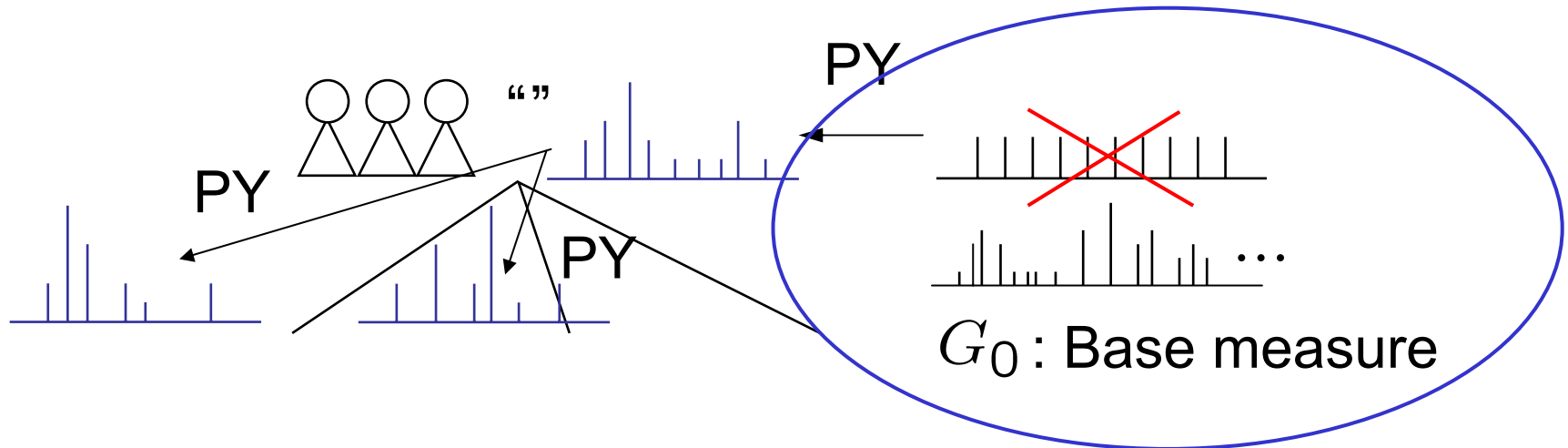


Hierarchical Pitman-Yor n-gram



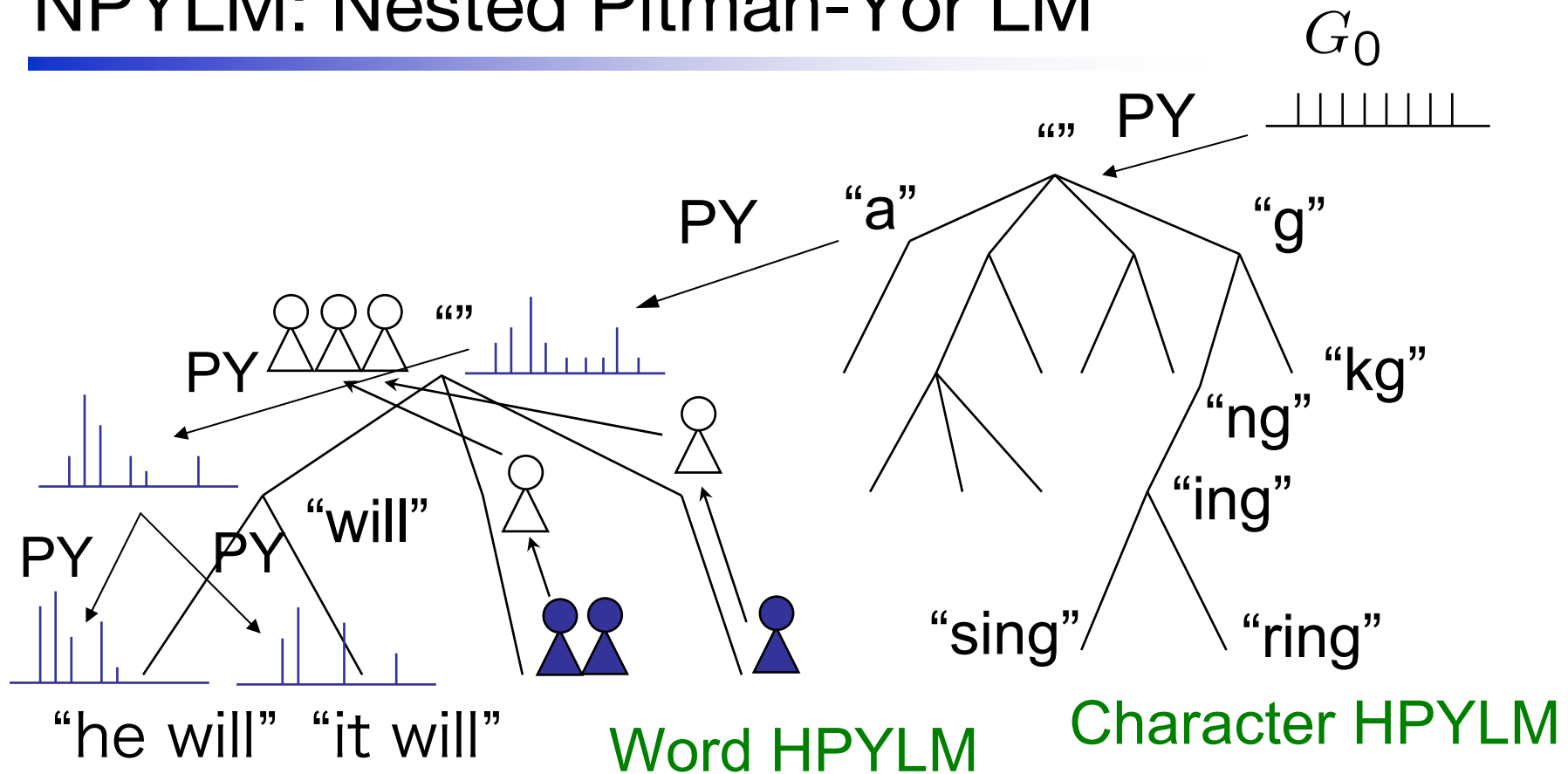
- Kneser-Ney smoothing is an approximation of hierarchical Pitman-Yor process (Teh, ACL 2006)
 - HPYLM = “Bayesian Kneser-Ney n-gram”

Problem: Word spelling



- Possible word spelling is not uniform
 - Likely: “will”, “language”, “hierarchically”, ...
 - Unlikely: “illbe”, “nguag”, “ierarchi”, ...
- Replace the base measure using character information
 - Character HPYLM!

NPYLM: Nested Pitman-Yor LM



- Character n-gram embedded in the base measure of word n-gram
 - i.e. hierarchical Markov model
 - Poisson word length correction (see the paper)

Inference and Learning

- Simply maximize the probability of strings
 - i.e. minimize the perplexity per character of LM
- X : Set of strings s_1, s_2, \dots, s_N
 Z : Set of hidden word segmentation indicators

z_1, z_2, \dots, z_N

$$p(X) = \prod_n p(s_n)$$

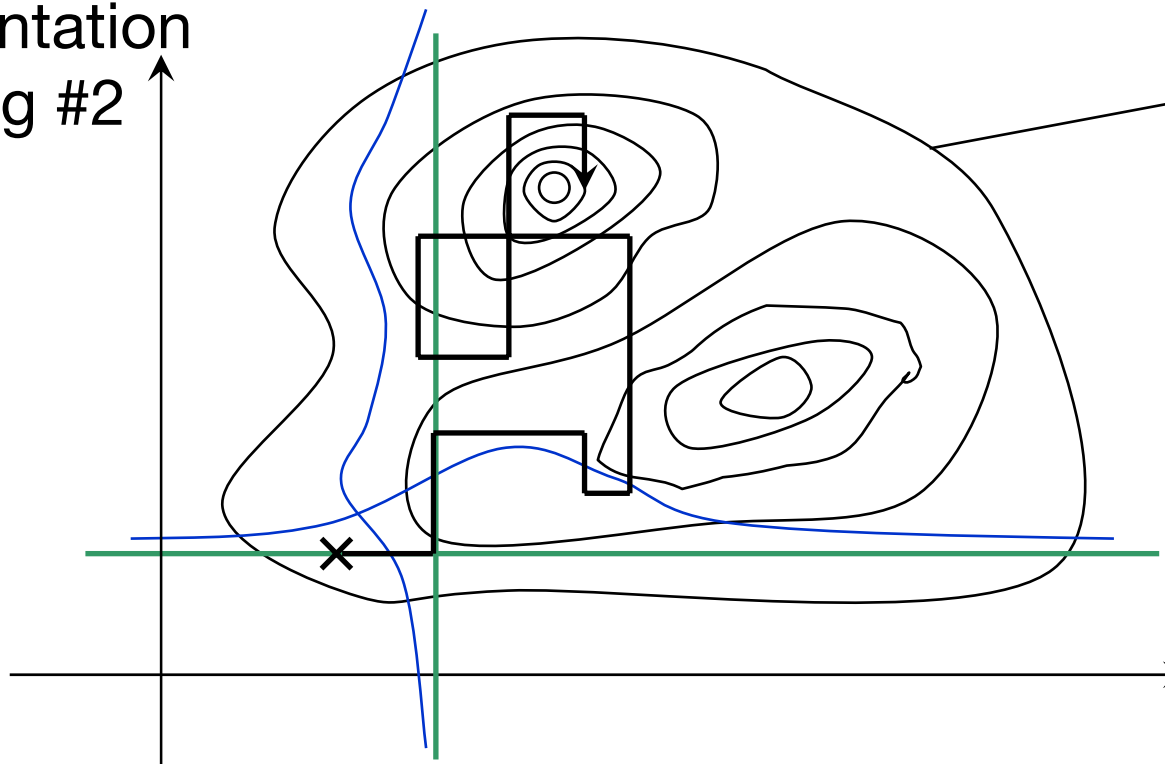
$$p(s_n) = \sum_{z_n} p(s_n, z_n)$$

Hidden word segmentation
of string s_n

- Notice: *Exponential possibilities* of segmentations!

Blocked Gibbs Sampling

Segmentation
of String #2



Probability
Contours of
 $p(X,Z)$

Segmentation
of String #1

- Sample word segmentation block-wise for each sentence (string)
 - High correlations within a sentence

Blocked Gibbs Sampling (2)

- Iteratively improve word segmentations: $\text{words}(s)$ of s

0. For $s = s_1 \cdots s_N$ do
 parse_trivial(s, Θ).

Whole string is
a single “word”

1. For $j = 1..M$ do

 For $s = \text{randperm}(s_1 \cdots s_N)$ do

 Remove $\text{words}(s)$ from NPYLM Θ

 Sample $\text{words}(s) \sim p(w|s, \Theta)$

 Add $\text{words}(s)$ to NPYLM Θ

 done

 Sample all hyperparameters of Θ

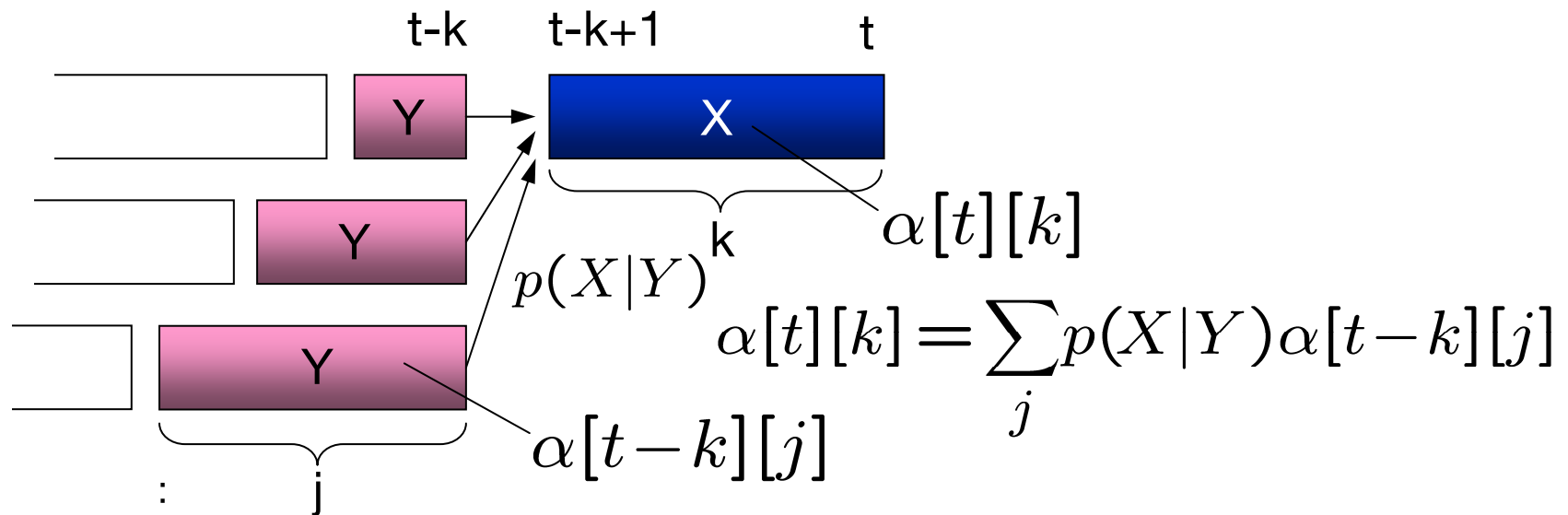
done

Gibbs Sampling and word segmentation

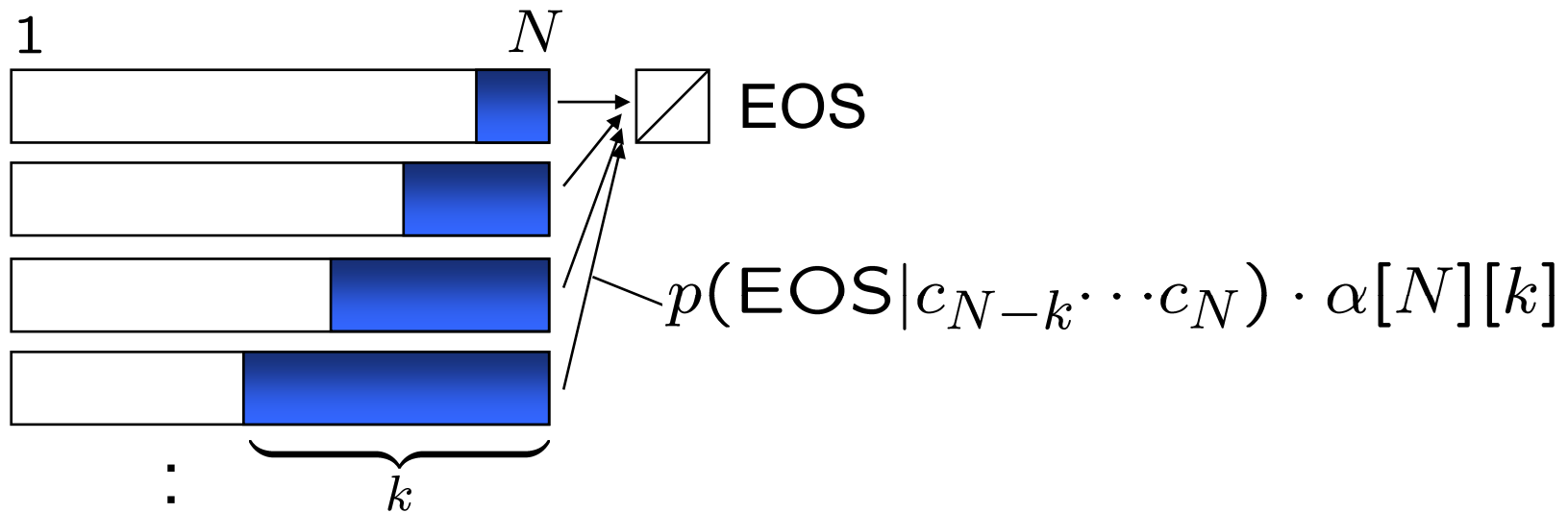
- 1 神戸では異人館 街の 二十棟 が破損した。
 - 2 神戸 では 異人館 街の 二十棟 が破損した。
 - 10 神戸 では 異人館 街の 二十棟 が破損した。
 - 50 神戸 では異人 館 街 の 二十棟 が破損した。
 - 100 神戸 では 異 人館 街 の 二十棟 が破損した。
 - 200 神戸 では 異人館 街 の 二十棟 が破損した。
- Iteratively resample word segmentations and update language models accordingly.

Sampling through Dynamic Programming

- Forward filtering, Backward sampling (Scott 2002)
- $\alpha[t][k]$: inside probability of substring $c_1 c_2 \dots c_t$ with the last k characters constituting a word
 - Recursively marginalize segments before the last k

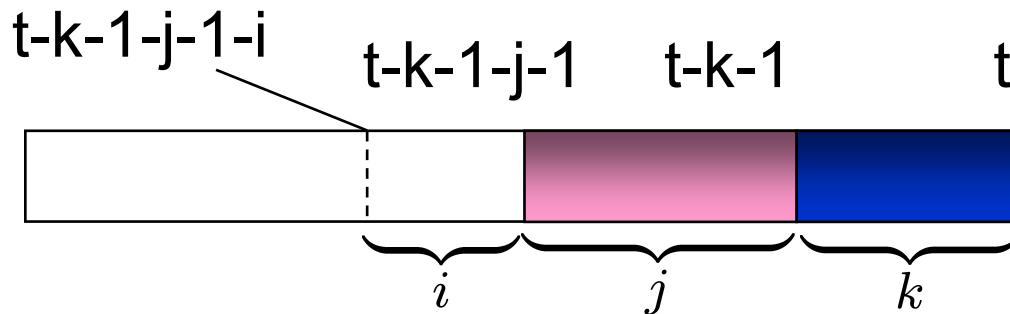


Sampling through Dynamic Programming (2)



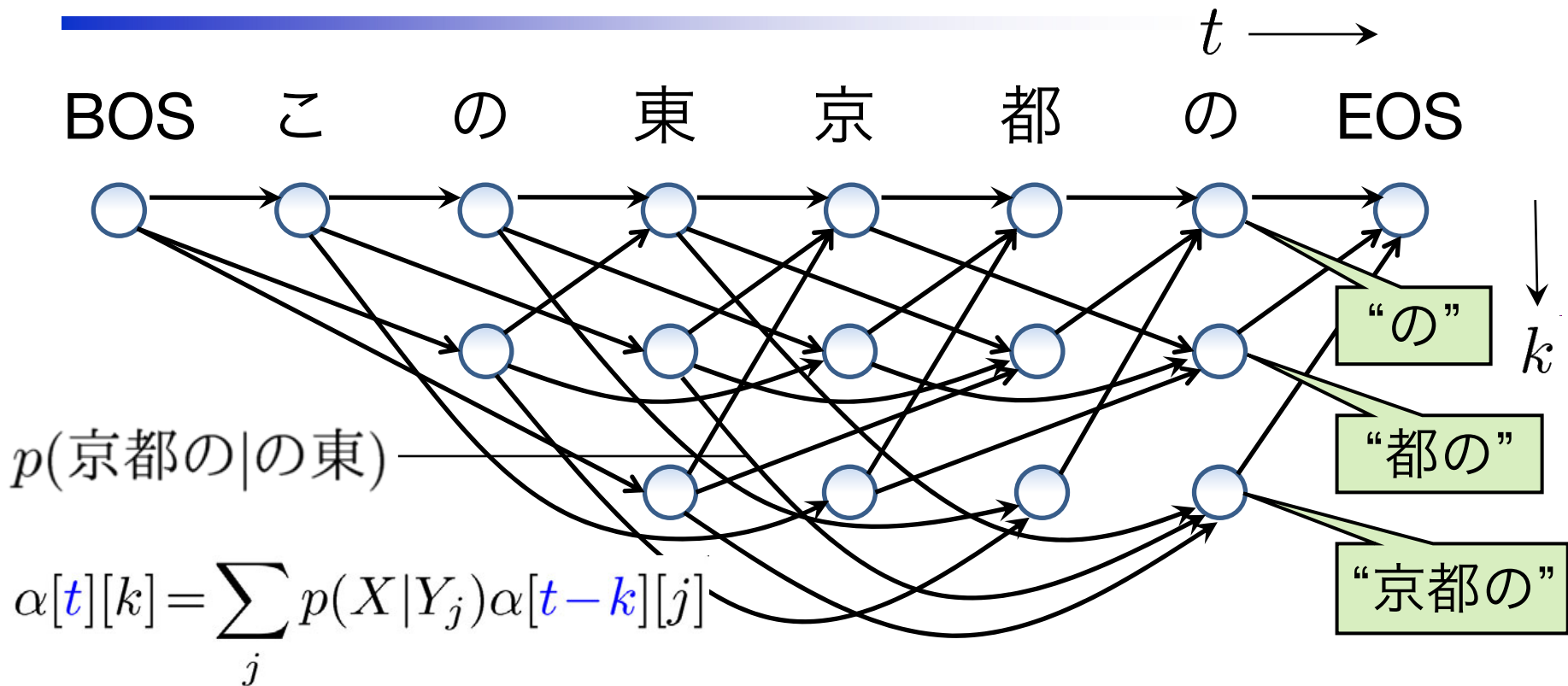
- $\alpha[N][k]$ = probability of the entire string $c_1 \dots c_N$ with the last k characters constituting a word
 - Sample k with probability to end with EOS
- Now the final word is $c_{N-k} \dots c_N$: use $\alpha[N-k-1][k']$ to determine the previous word, and repeat

The Case of Trigrams



- In case of trigrams: use $\alpha[t][k][j]$ as an inside probability
 - $\alpha[t][k][j] =$ probability of substring with the final k chars and the further j chars before it being words
 - Recurse using $\alpha[t-k-1][j][i]$ ($i = 0 \dots L$)
- >Trigrams? Practically not so necessary, but use Particle MCMC (Doucet+ 2009 to appear) if you wish

NPYLM as a Semi-Markov model



- Unsupervised learning of Semi-Markov HMM (Ostendorf 96, Murphy 02)
- State transition = word transition with an intensive smoothing w/ NPYLM + MCMC

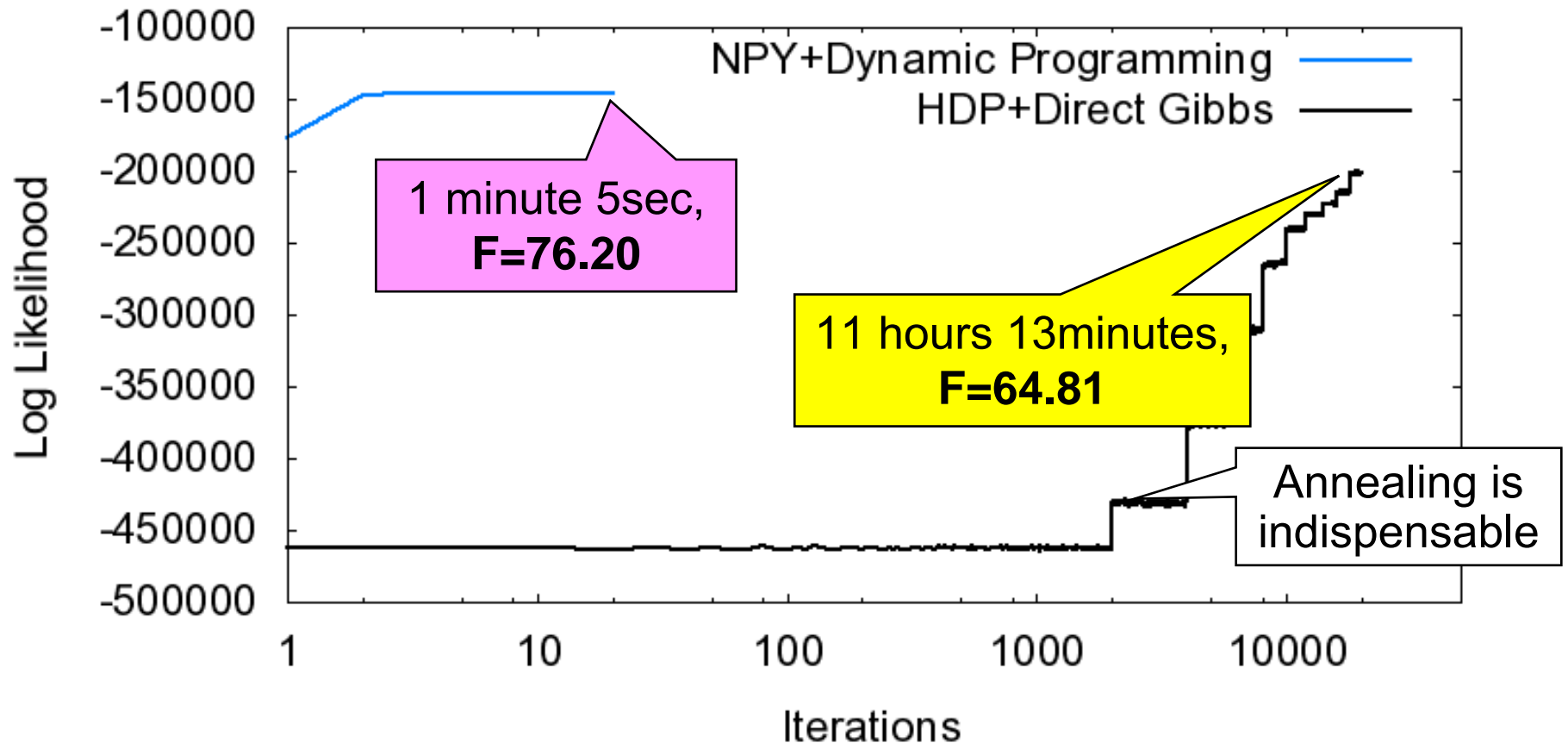
English Phonetic Transcripts

- Comparison with HDP bigram (w/o character model) in Goldwater+ (ACL 2006)
- CHILDES English phonetic transcripts
 - Recover “WAtsDIs” → ”WAts DIs” (What’s this)
 - Johnson+(2009), Liang(2009) use the same data

<i>Model</i>	P	R	F	LP	LR	LF
NPY(3)	74.8	75.2	75.0	47.8	59.7	53.1
NPY(2)	74.8	76.7	75.7	57.3	56.6	57.0
HDP(2)	75.2	69.6	72.3	63.5	55.2	59.1

- **Very small** data: 9,790 sentences, 9.8 chars/sentence

Convergence & Computational time



- NPYLM is very efficient & accurate! (600x faster here)

Chinese and Japanese

Perplexity per character

Model	MSR	CITYU	Kyoto
NPY(2)	0.802 (51.9)	0.824 (126.5)	0.621 (23.1)
NPY(3)	0.807 (48.8)	0.817 (128.3)	0.666 (20.6)
NPY(+)	0.804 (38.8)	0.823 (126.0)	0.682 (19.1)
ZK08	0.667 (—)	0.692 (—)	—

- MSR&CITYU: SIGHAN Bakeoff 2005, Chinese
- Kyoto: Kyoto Corpus, Japanese
- ZK08: Best result in Zhao&Kit (IJCNLP 2008)

Note: Japanese subjective quality is much higher (proper nouns combined, suffixes segmented, etc..)

Arabic

- Arabic Gigawords 40,000 sentences (AFP news)

الفلستيني بسبب تظاهرة لانصار حركة المقاومة الاسلامية حماس
و اذا تحقق ذلك فان كيسلو فسكيه قد حاز ثلاثه جري في ابرز ثلاثة

صحية
+ قائد
الا يقل

Google translate:

“Filstinebsbptazahrplansarhrkpalmquaompalaslami
phamas.”

وقالت دانيل تومسون التي كتبت السيناريو. وقد استغرق اعداد خمسة اعوام. "تاريخي

↓ NPYLM

الفلستيني بسبب تظاهرة لانصار حركة المقاومة الاسلامية حماس
و اذا تحقق ذلك فان كيسلو فسكي يكون قد حاز ثلاثه جري في ابرز ثلاثة

صحية
سطينية
مالا يقل

Google translate:

“Palestinian supporters of the event because of the
Islamic Resistance Movement, Hamas.”

وقد استغرق اعداد ه خمسة اعوام . و قال ت دانيل تومسون التي " تاريخي

English (“Alice in Wonderland”)

first, she dreamed of little Alice herself, and once again the tiny hands were clasped upon her knee, and the bright eager eyes were looking up into hers -- she could hear the very tones of her voice, and see that queer little toss of her head to keep back the wandering hair that would always get into her eyes -- and still as she listened, or seemed to listen, the whole place around her became alive the strange creatures of her little sister's dream. The long grass rustled at her feet as the white rabbit hurried by -- the frightened mouse splashed his way through the neighbouring pool -- she could hear the rattle of the tea cups as the March hare and his friends shared their never-ending meal, and the shrill voice of the queen...



first, she dream ed of little Alice herself , and once again the tiny hand s were clasped upon her knee , and the bright eager eyes were looking up into hers -- she could hear the very tone s of her voice , and see that queer little toss of her head to keep back the wandering hair that would always get into hereyes -- and still as she listened , or seemed to listen , the whole place a round her became alive the strange creatures of her little sister 's dream. the long grass rustled at her feet as the whiterabbit hurried by -- the frightened mouse splashed his way through the neighbour ing pool -- she could hear the rattle of the tea cups as the march hare and his friends shared their never -endingme a l , and the ...

Conclusion of the first part

- Completely unsupervised word segmentation of **arbitrary language** strings
 - Combining word and character information via hierarchical Bayes
 - **Very efficient** using forward-backward+MCMC
- Directly optimizes Kneser-Ney language model
 - N-gram construction **without any “word” information**
 - Sentence probability calculation **with all possible word segmentations marginalized out**
 - Easily obtained from dynamic programming

Problems of Unsupervised segmentation?

- Optimize n -gram language models
 - Must be optimized for different tasks
 - For machine translation: Nguyen+ (COLING 2010)
 - For speech recognition : Neubig+ (Interspeech 2010)
- Not always fit for human standards
 - Ex. Inflection, proper nouns, human preference
 - “咲か”-“咲き”-“咲く”, “盧前大統領”, “その”
 - Remedy:
 1. Make the generative model more complex
 2. *Semi-supervised learning*
 - human standards are usually closed and general

Semi-supervised learning : JESS-CM

- JESS-CM (Suzuki+ ACL-HLT 2008):

better than Druck+
(ICML2010)

“joint probability model embedding style semi-supervised conditional model”

- *highest performance semi-supervised learning on POS tagging, NE chunking, dependency parsing*

- Model:

Discriminative model

Generative model

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

- Product model: discriminative and generative
- However no naïve product
 - “model weight” λ is included in Λ
 - Θ is influenced by Λ through learning (i.e. recursive)

JESS-CM : inference

Discriminative model

Generative model

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

- If the discriminative model is log-linear (like CRF):

$$p_{\text{DISC}}(\mathbf{y}|\mathbf{x}) \propto \exp\left(\sum_{k=1}^K \lambda_k f_k(\mathbf{y}, \mathbf{x})\right)$$

- Then the model is again loglinear:

$$p(\mathbf{y}|\mathbf{x}) \propto \exp\left(\lambda \log p_{\text{GEN}}(\mathbf{y}, \mathbf{x}) + \sum_{k=1}^K \lambda_k f_k(\mathbf{y}, \mathbf{x})\right)$$

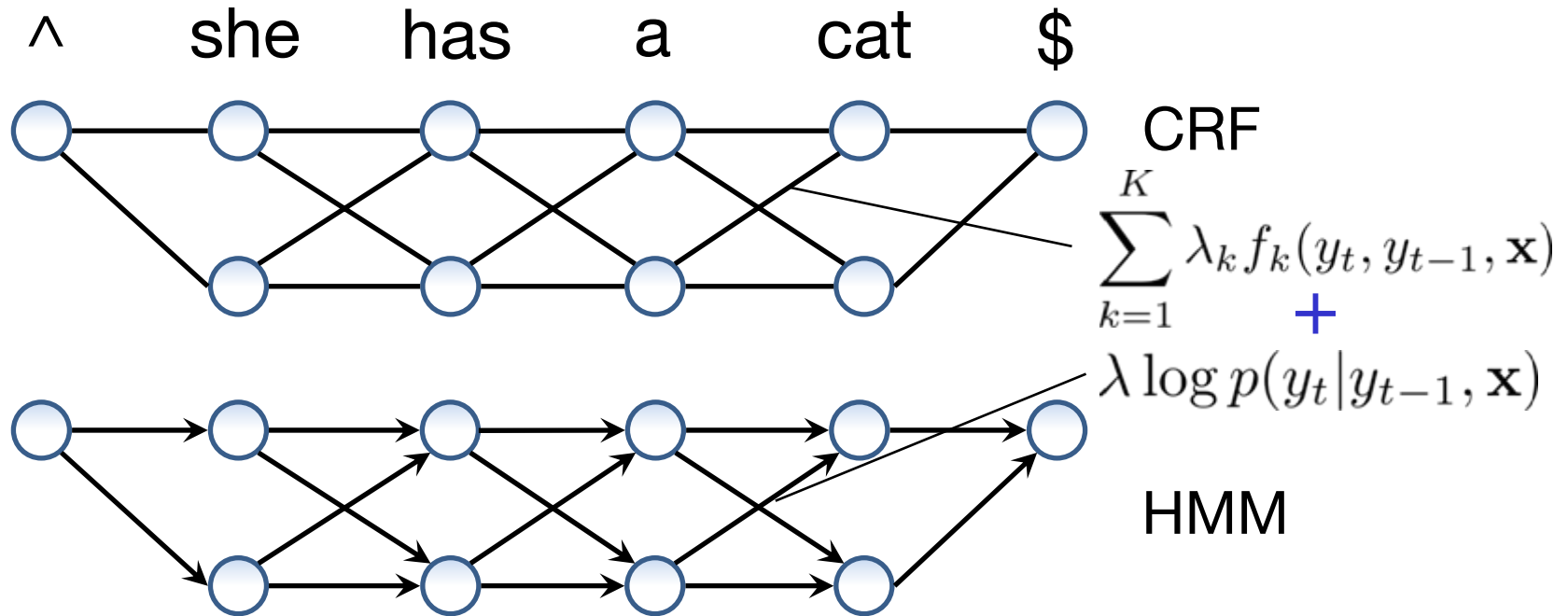
Additional feature!

- Inference = maximize the objective

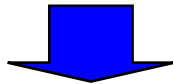
$$\log p(D|\Lambda, \Theta) = \log p(\mathbf{Y}_s|\mathbf{X}_s; \Lambda, \Theta) + \log p(\mathbf{X}_u; \Lambda, \Theta)$$

- Fix Θ and optimize Λ on $\mathbf{Y}_s, \mathbf{X}_s$
- Fix Λ and optimize Θ on \mathbf{X}_u

JESS-CM on CRF-HMM (Suzuki+ ACL08)

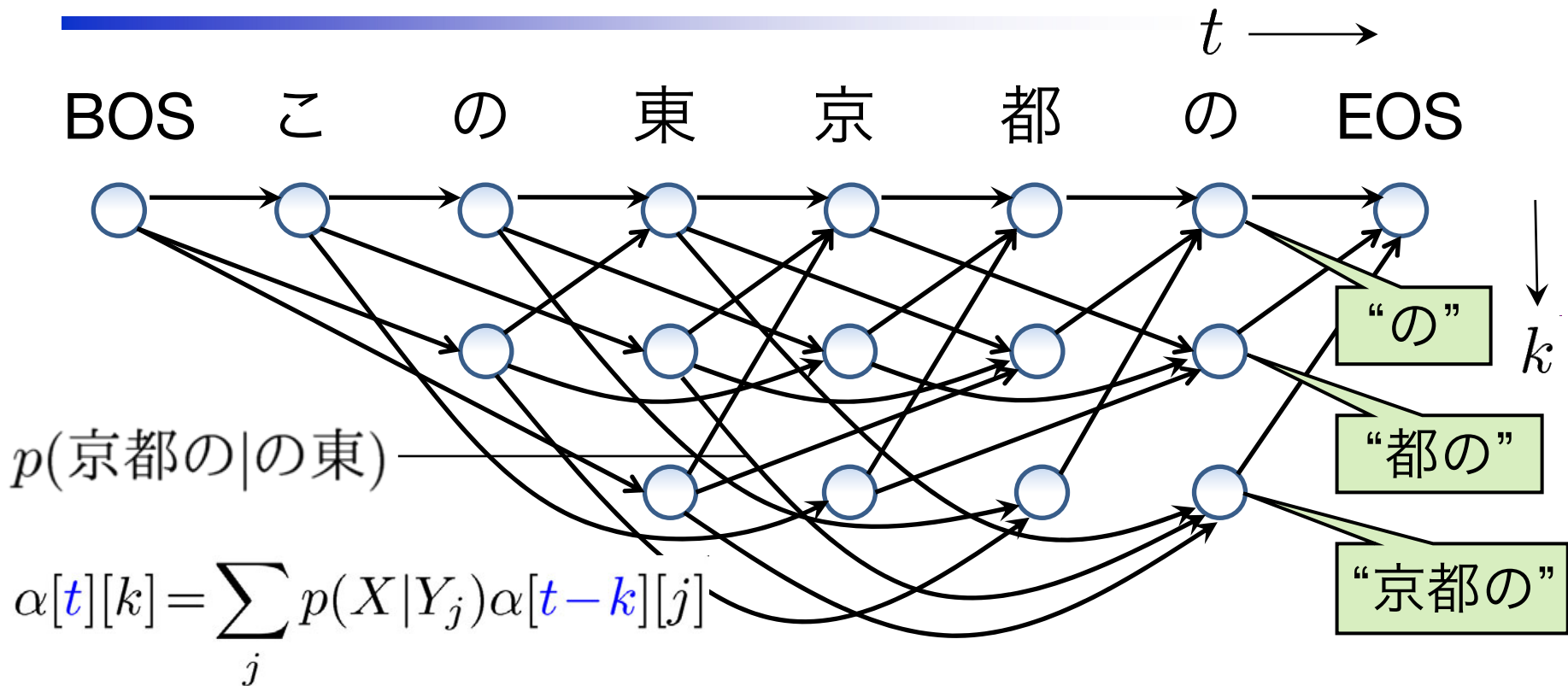


- Sum the corresponding weights of the path *on the same graphical model*



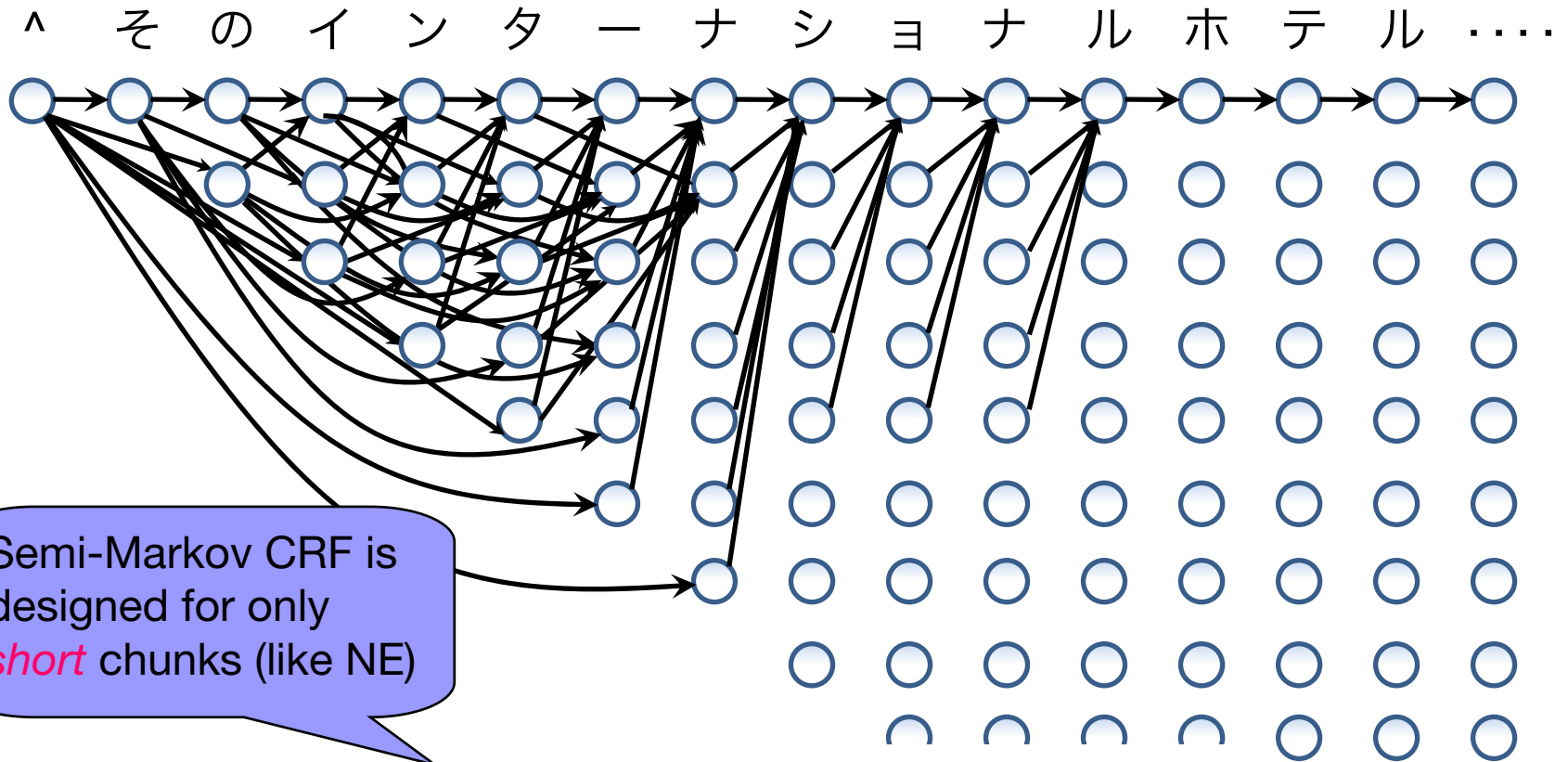
Interleave CRF and HMM optimization

NPYLM as a Semi-Markov model



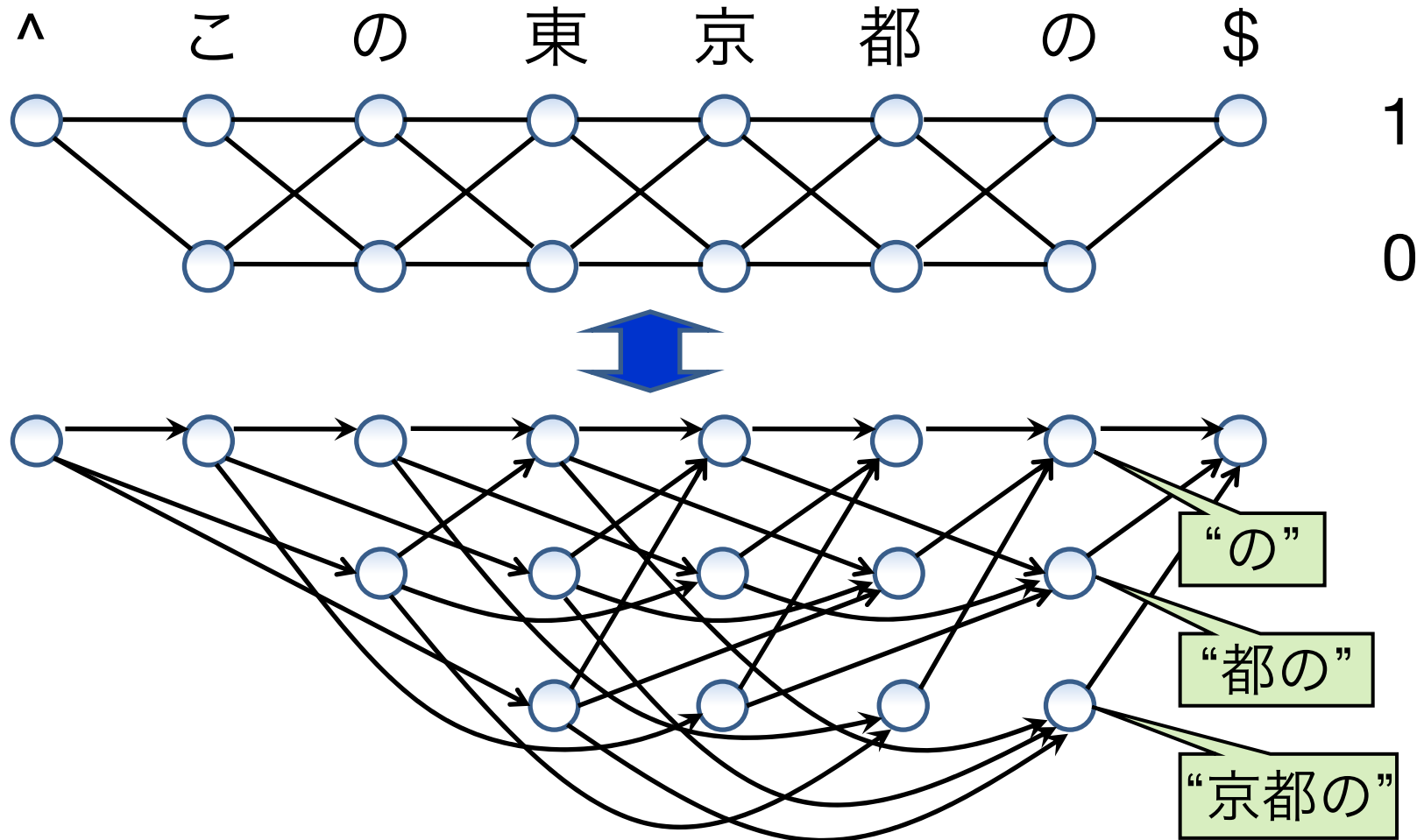
- NPYLM is not Markov, but semi-Markov
- State transition = word transition with an intensive smoothing w/ NPYLM + MCMC
 - CRF combination?

Semi-Markov CRF (NIPS 2004)?



- **Enormous** memory (1GB→20GB)
- (Supervised) precision: at most **95%**
 - Only words, no character-level information

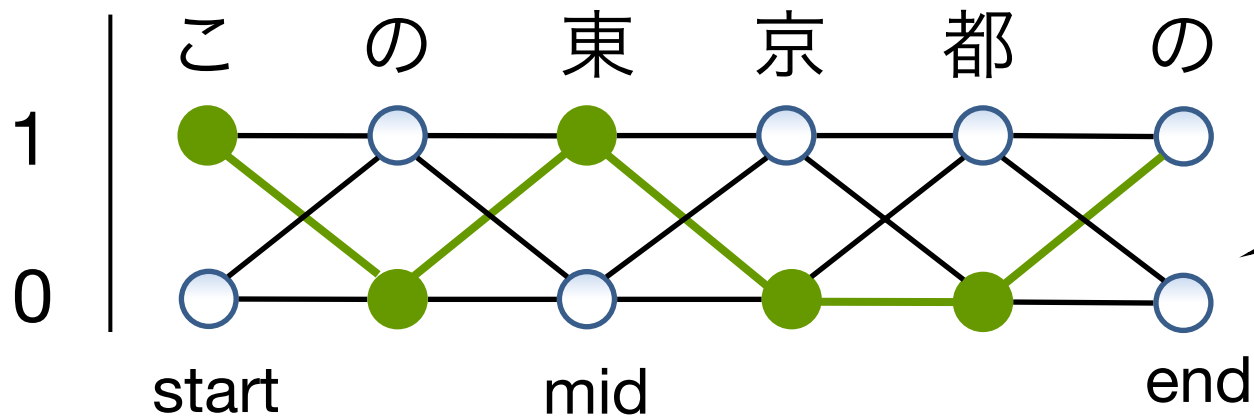
Learning Markov CRF \leftrightarrow Semi-Markov LM



- *How to combine two different models?*
 - CRF \rightarrow NPYLM, NPYLM \rightarrow CRF

CRF→NPYLM

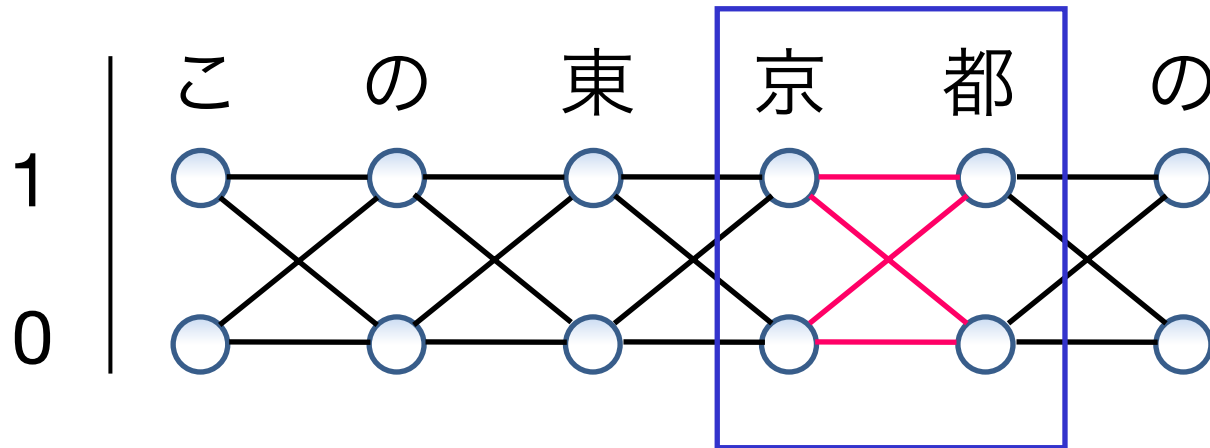
- Easy, proposed in (Andrew+ EMNLP 2006)
 - CRF→semi-Markov CRF
 - $p(\text{“この”} \rightarrow \text{“東京都”})$



- Summing up the weight of features along the path
- $\gamma(\text{start}, \text{mid}, \text{end})$
 $:= \gamma(\text{start}, \text{mid}) + \gamma(\text{mid}, \text{end})$

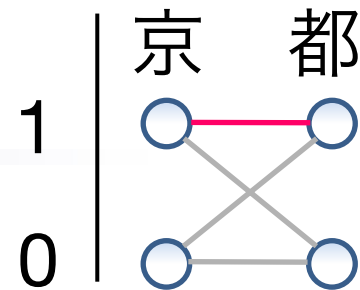
NPYLM \rightarrow CRF (1)

- Nontrivial !



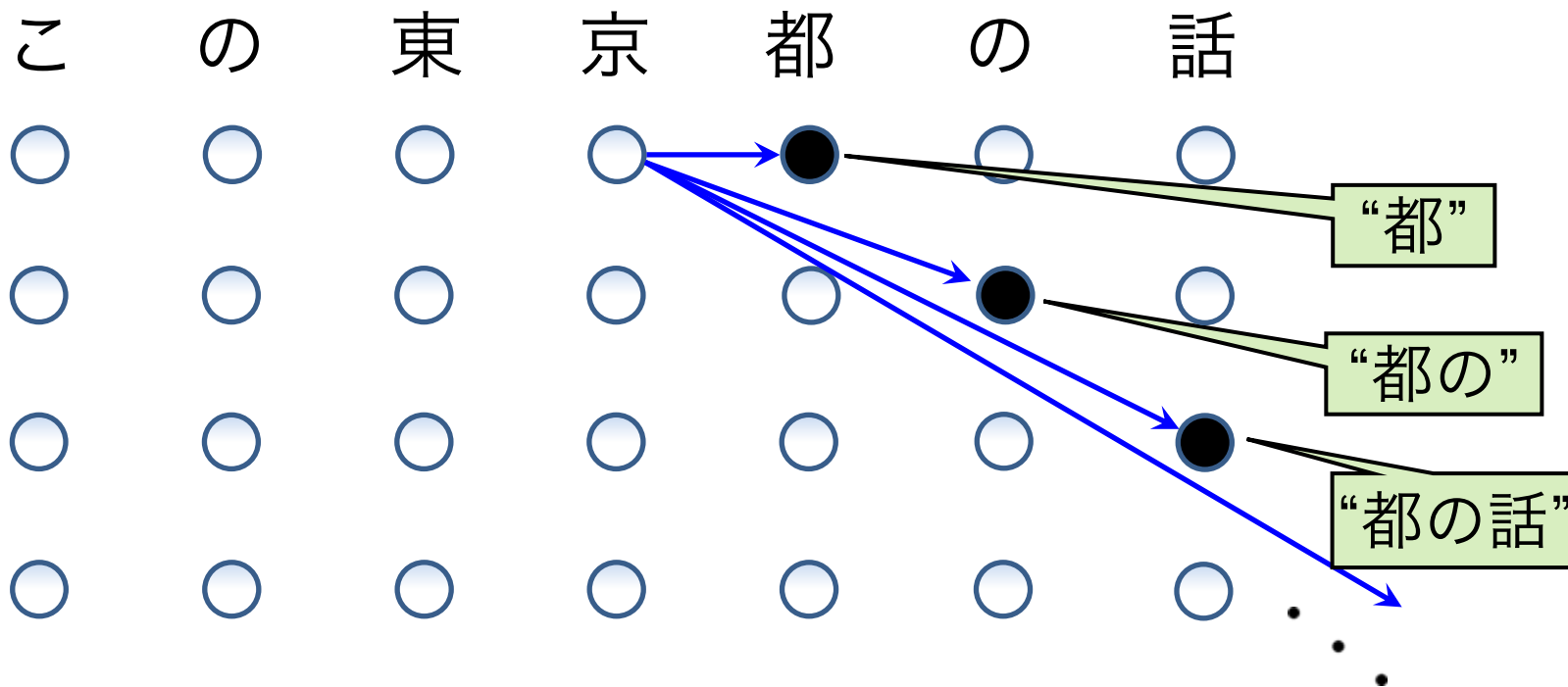
- Four patterns: $0 \rightarrow 0$, $0 \rightarrow 1$, $1 \rightarrow 0$, $1 \rightarrow 1$
- Pretend if the model is Markov HMM (not semi-Markov)
- When sentence \mathbf{x} is given, we **can** compute the corresponding potentials by intricately summing up NPYLM probabilities!

NPYLM→CRF (2)

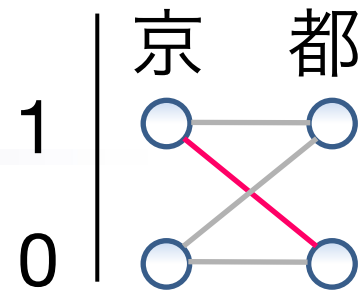


- Case 1→1 :

1→1 = “京→都”, “京→都の”, “京→都の話”, ...

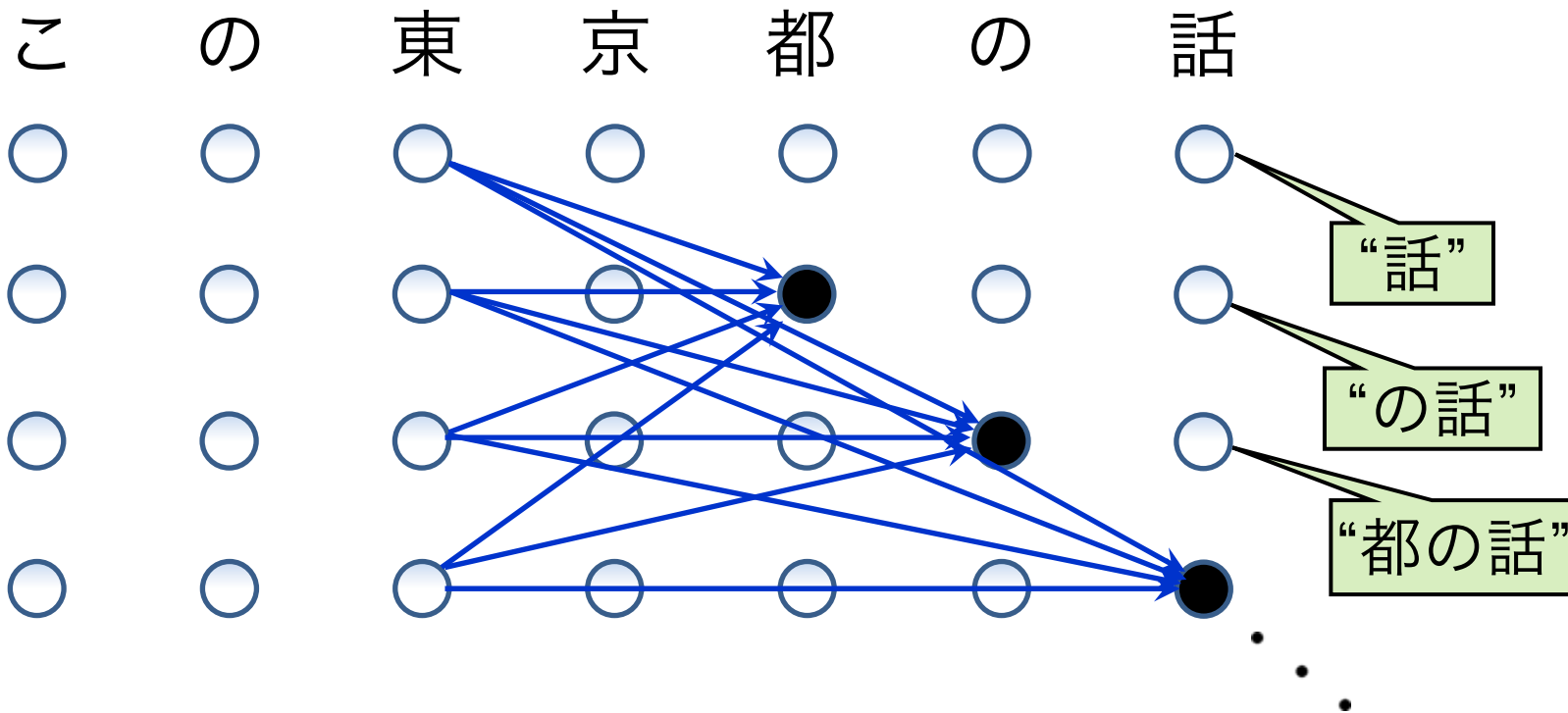


NPYLM→CRF (3)



- 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

- C++ code for summing probabilities to compute

“0→0” potential:

```
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)

- *Mathematically*, this computation is a **marginalization**

- By definition,

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

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

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

- Then we can marginalize:

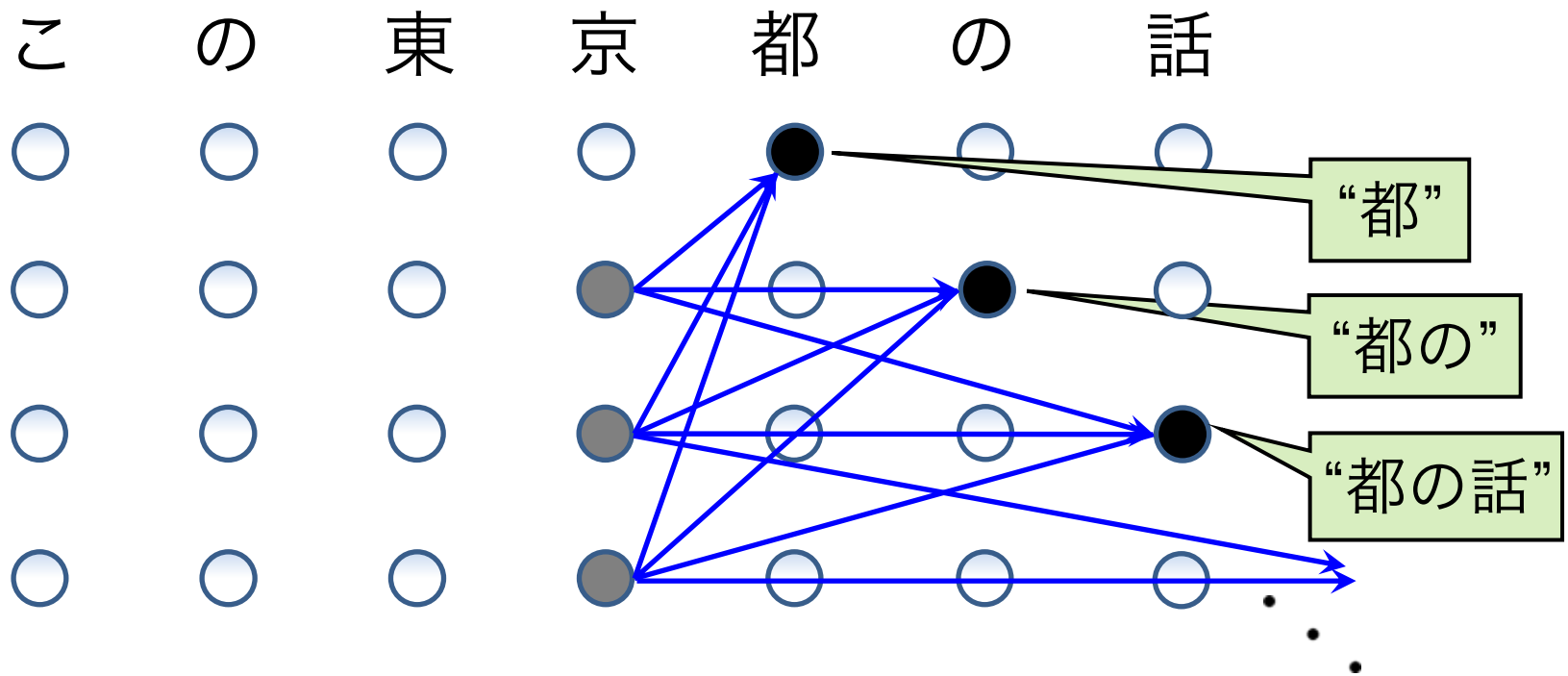
$$p(z_t = 1, z_{t+1} = 1) = \sum_k p(\underline{z_t = 1}, \underline{z_{t+1} = 1}, \dots, \underline{z_k = 1})$$

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

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

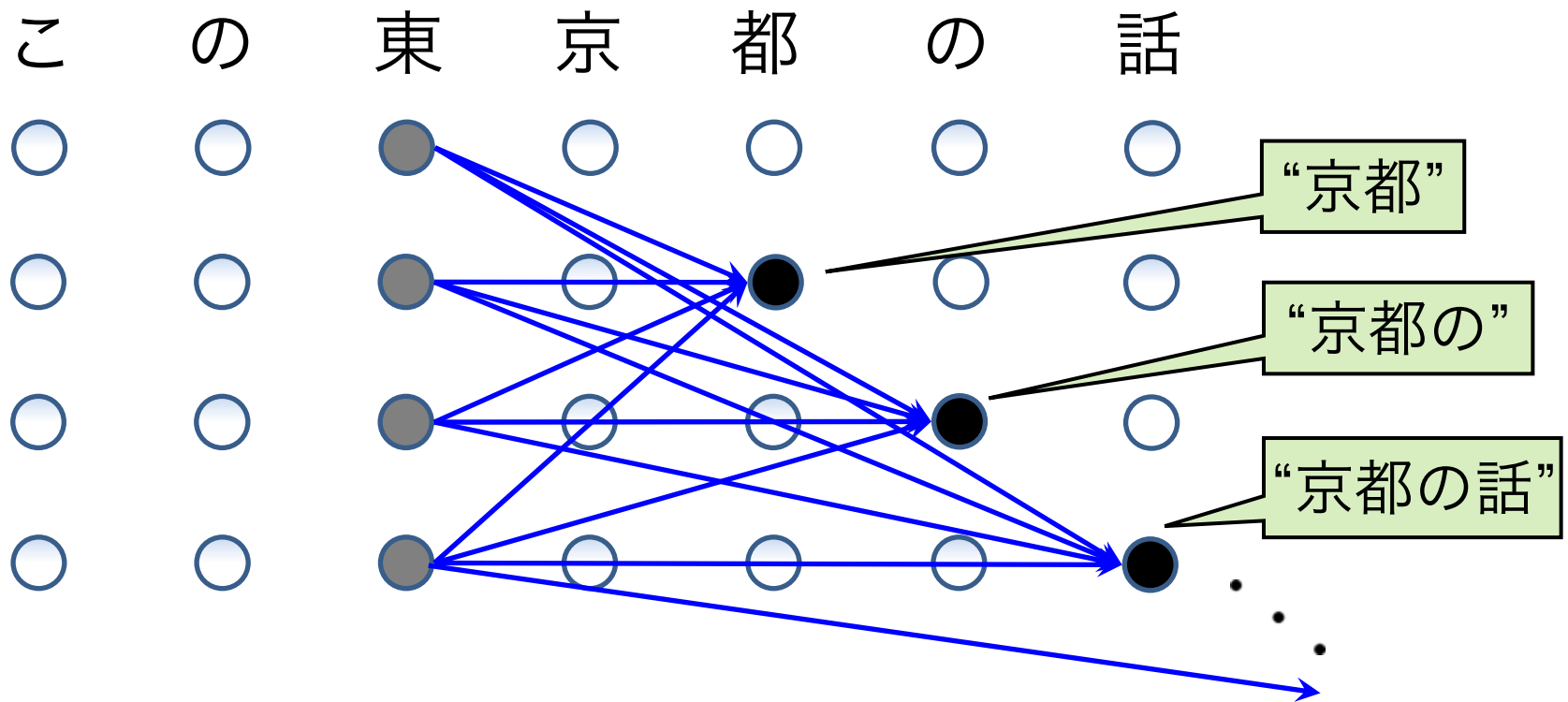
What are we doing? (2)

- Graphically, summing the collection of paths for the case of $0 \rightarrow 1$:



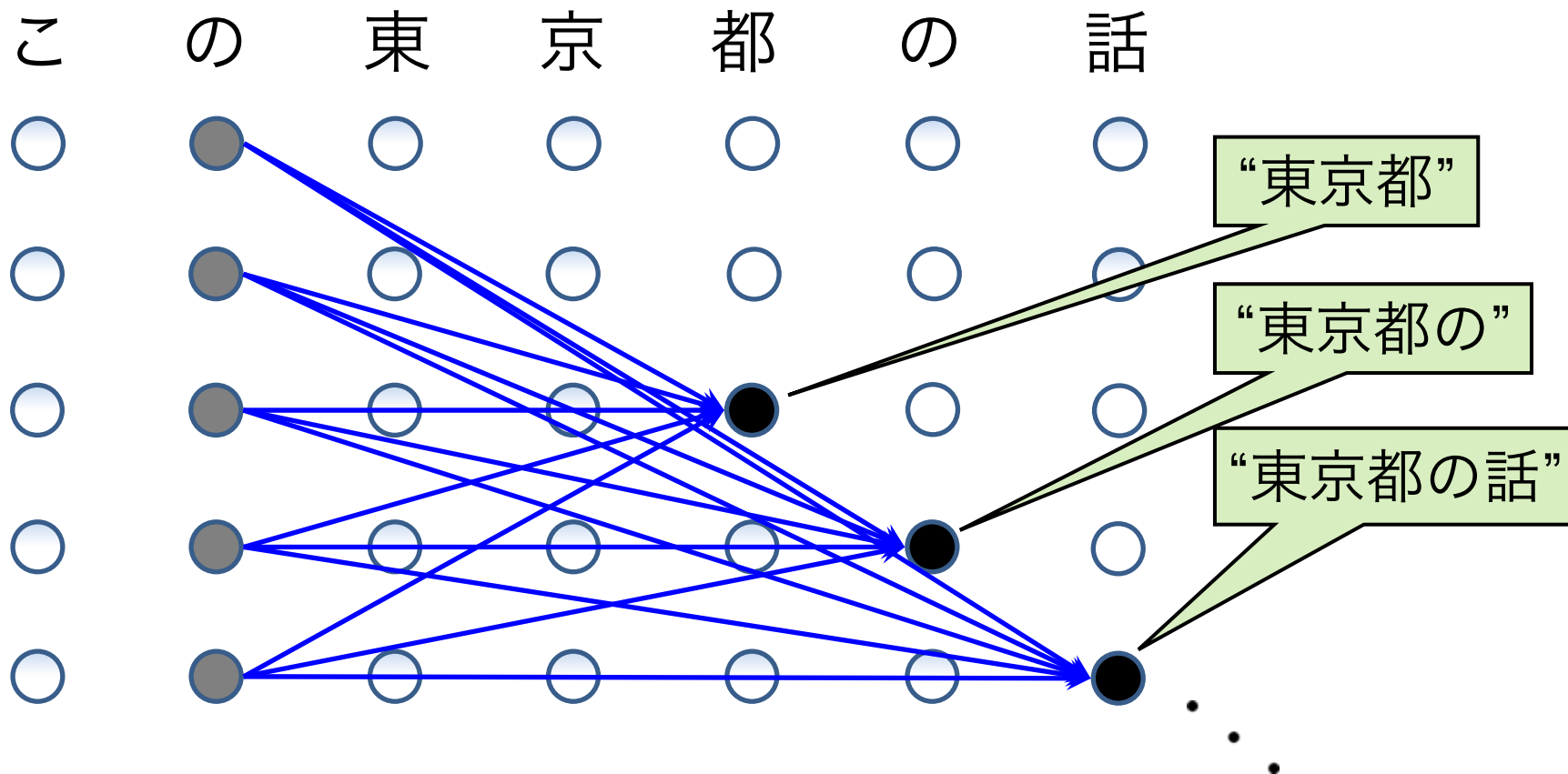
What are we doing? (2)

- Graphically, summing the collection of paths for the case of $1 \rightarrow 0$:



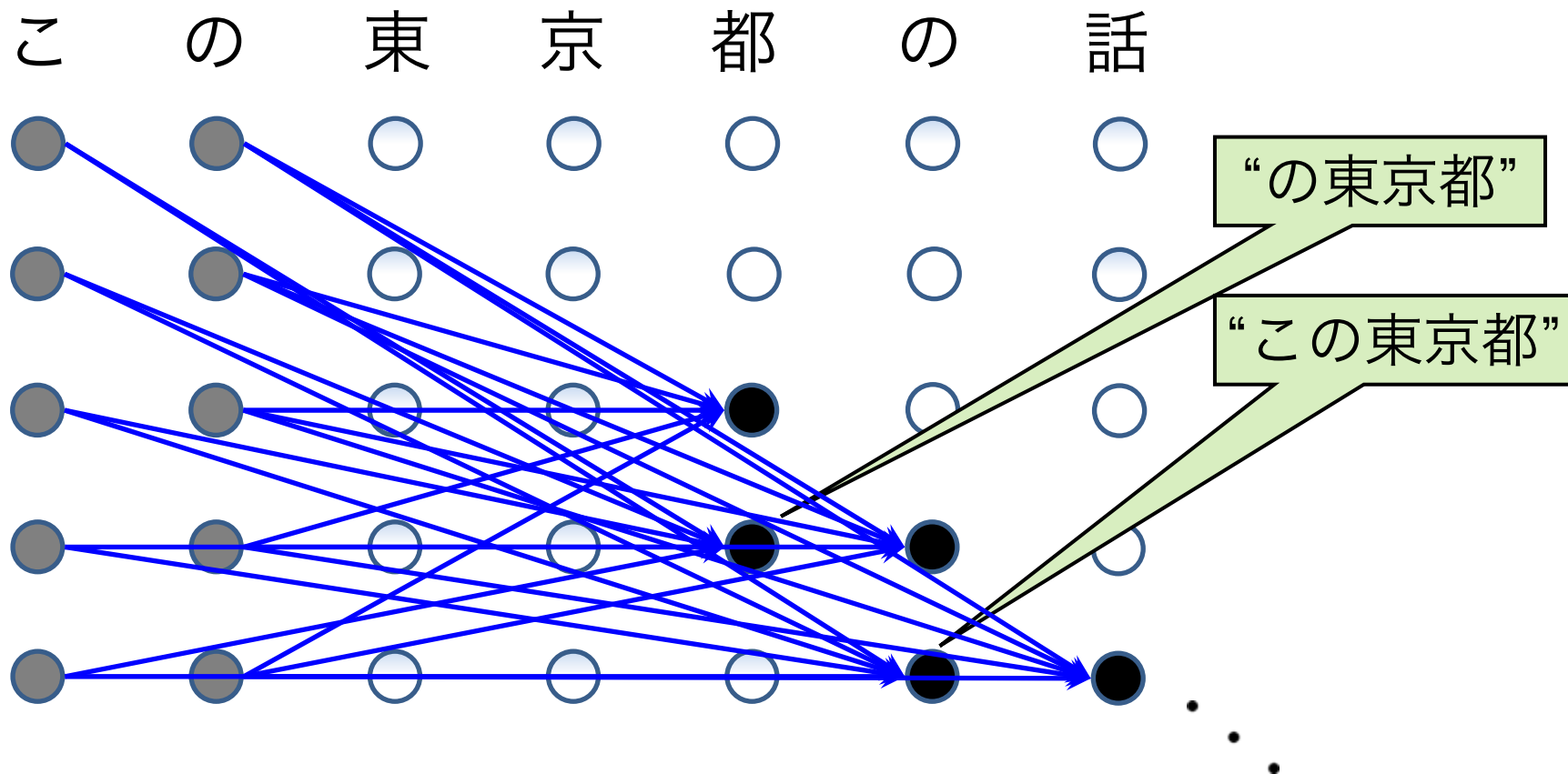
What are we doing? (2)

- Graphically, summing the collection of paths for the case of $0 \rightarrow 0$:



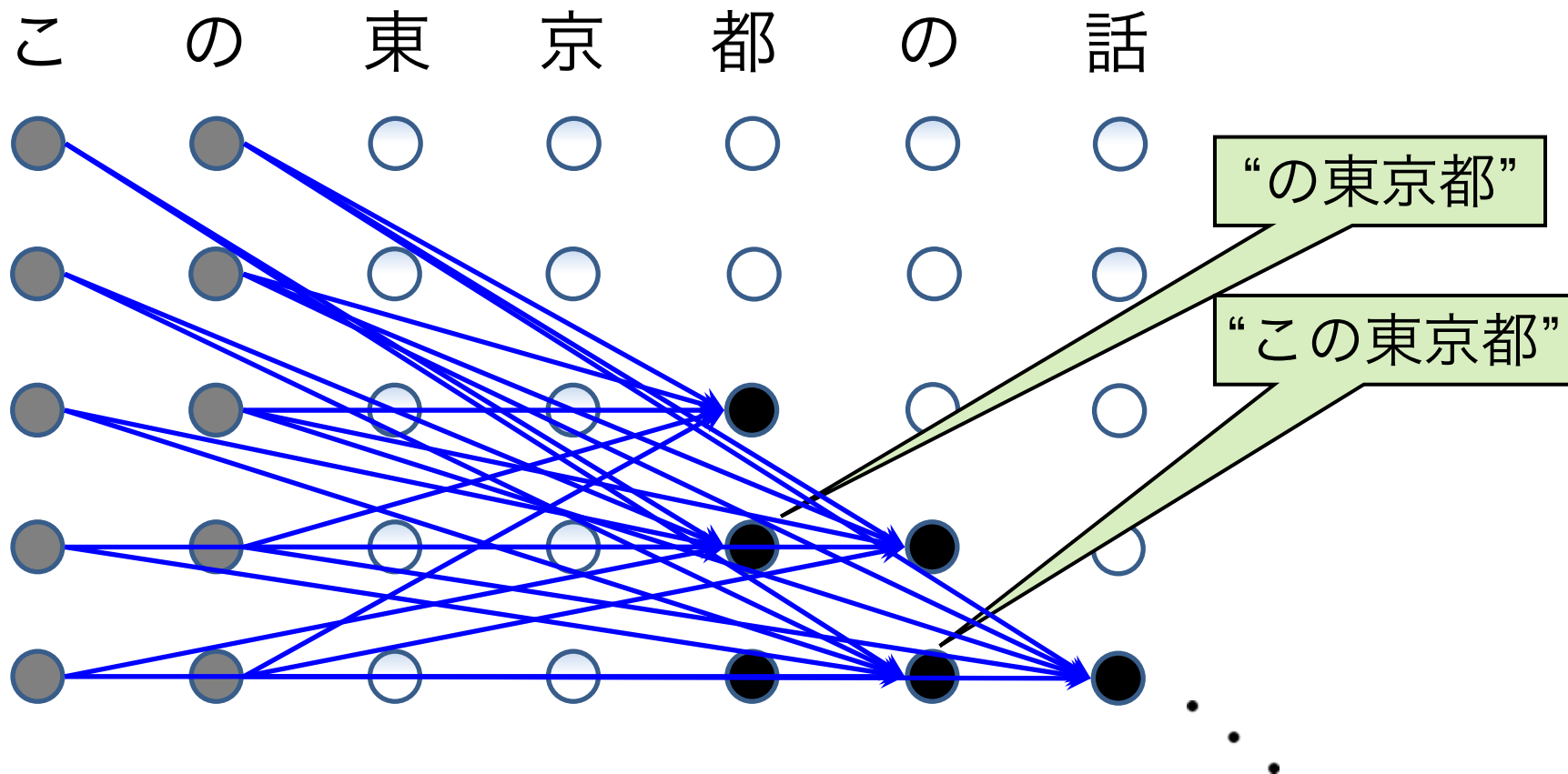
What are we doing? (2)

- Graphically, summing the collection of paths for the case of $0 \rightarrow 0$:



What are we doing? (2)

- Graphically, summing the collection of paths for the case of $0 \rightarrow 0$:

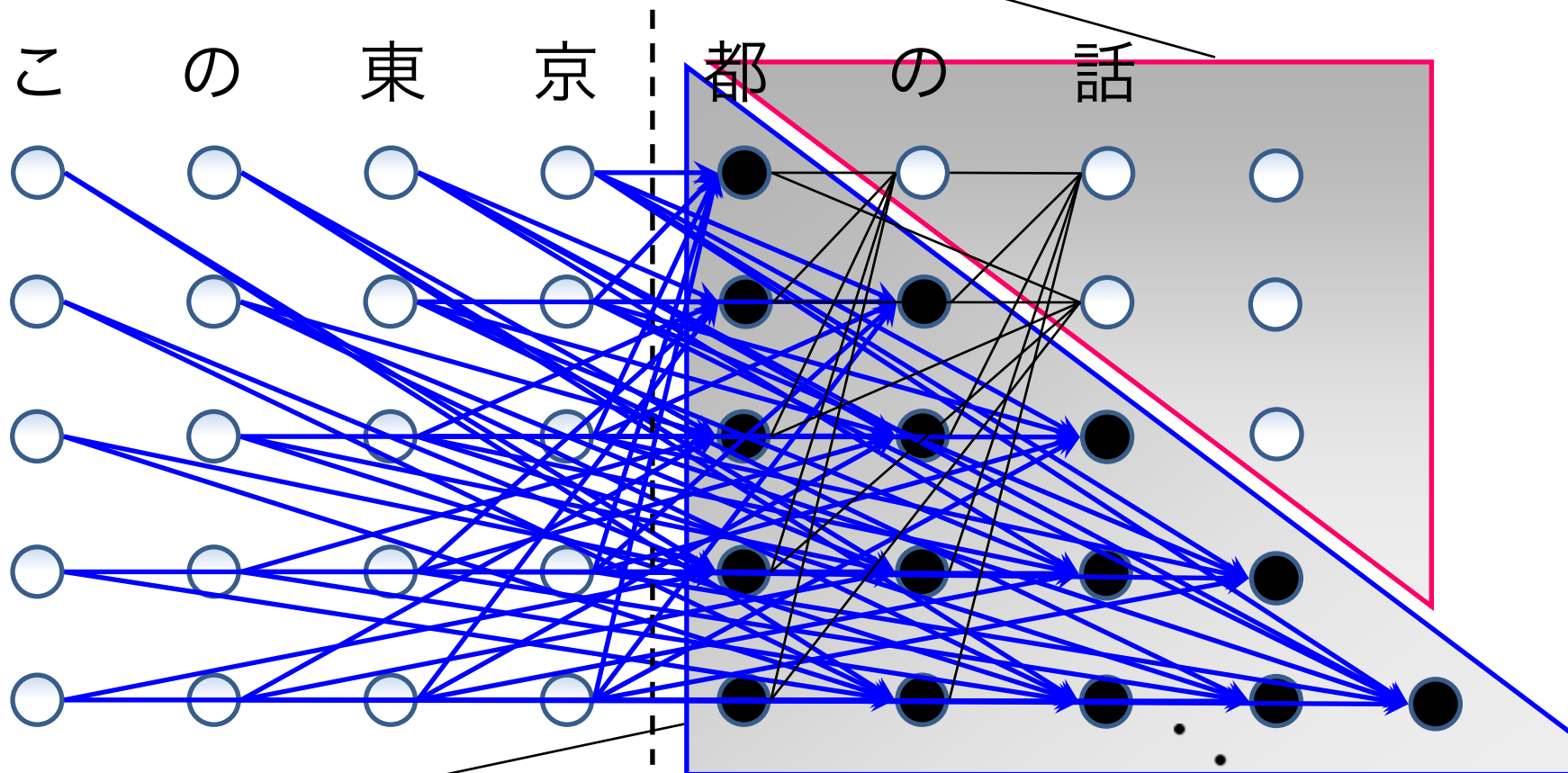


What are we doing? (2)

- Graphically, divide the paths at the section into four bins:

Section

Irrelevant Area



Relevant area, crossing the section

Experiments (still ongoing)

- Sina Microblog (新浪微博) 

Tremendous!

 - Chinese equivalent of Twitter, 94,800,000 users
 - Full of non-standard Chinese sentences
- Japanese blog (Shokotan blog)
 - Famous for being full of jargons, new proper nouns, ..
- CSJ (Corpus of Spontaneous Japanese)
 - by National Institute of Japanese language and linguistics, next to ISM
- SIGHAN word segmentation Bakeoff 2005
 - Public dataset, intensively studied, newspaper

“Shokotan Blog” segmentations

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 w w w w w w お茶 w w w w w w イト カワユス w w w w w w イト カワユス w w w w w w (^ω^)(^ω^)(^ω^) 深夜までお疲れさ マミタス(°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^)

- Supervised : Kyoto corpus 37,400 sentences
- Unsupervised: Shokotan blog 40,000 sentences

Shokotan blog: words

- Excerpt of random words with frequency ≥ 2

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

Sina microblog (新浪微博)

One of the most popular sites in China/HK/Taiwan

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

學校附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(○=^·ェ·)

。┌┐ 喵~o(=∩ω∩=)m

梦。。。混乱的梦。。。清晰的梦。。。。。

- Supervised : MSR 87000 sentences (Mandarin)
- Unsupervised : Sina API, 98700 sentences

Corpus of Spontaneous Japanese (CSJ)

- Experimental setting
 - Supervised: Kyoto corpus+CSJ 10% (21000 sents)
 - Unsupervised: CSJ 90% (190000 sents)
 - Training: CRF (only supervised) / NPYLM (sup+unsup)
- Results w.r.t. CSJ hand segmentations
 - P/R/F=0.928/0.893/0.91→0.931/0.899/0.915
 - Generally better, but occasionally gets worse:

CRF> ええだから何て言うんでしょうか立っていると **お腹に** 力..

NPY> ええだから何て言うんでしょうか立っていると **お腹に** 力..

CRF> 神奈川 寄りとかあっちの方 **行っちゃい**ますよね 値段..

NPY> 神奈川 寄りとかあっちの方 **行っちゃい**ますよね 値段..

CRF> **力強い** 方向性あるタッチ

NPY> **力強い** 方向性あるタッチ

SIGHAN Bakeoff 2005

- Publicly available dataset of Chinese word segmentation, MSR portion
 - Mostly newspaper, *not so much suited for our task but standard*
- Data: MSR supervised (87000 sentences) + Chinese Gigaword (200000 sentences)
- Results:
 - Baseline: 97.4% F-measure with features of (Xu+ 2009)
 - MSR+Gigaword: 97.5%
 - MSR+Gigaword+dict: 97.5%

World best baseline! (closed)

97.3% on
DPLVM

Difficult to beat.. Much more data to cover test data
(5 days on Xeon W5590 to compute)

Conclusion

- Semi-supervised learning of **CRF-NPYLM**
 - Product models, each depends each
- Convert **Semi-Markov** \leftrightarrow **Markov** models
 - General observation, not only to NPYLM
- Good for blogs and Twitters
 - Maintaining accuracy on supervised data
 - Need: huge unsupervised data to improve on standard datasets
 - Parallelization or other implementation techniques

Future work

- Predicting POS (part-of-speech) simultaneously
 - “Class-based” n-gram models + Multiclass CRF
 - Class-based language models will also improve on unsupervised data
 - Ex. Predict “ギザ カワユス” → “ギザ スゴス”
- However, in the current NLP...
 - No standard theory for building class-based n-grams.. only heuristic (like Google Mozc)
 - Dealing with “hierarchy” of part-of-speech
 - Ex. Noun → Proper noun → Name of a person
 - Theory of Markov model on tree-structured classes? (imagine hLDA)

Recognize “ギザ” as adjective