# Bayesian Symbol-Refined Tree Substitution Grammars for Syntactic Parsing

Hiroyuki Shindo
NTT CS Labs.

Yusuke Miyao
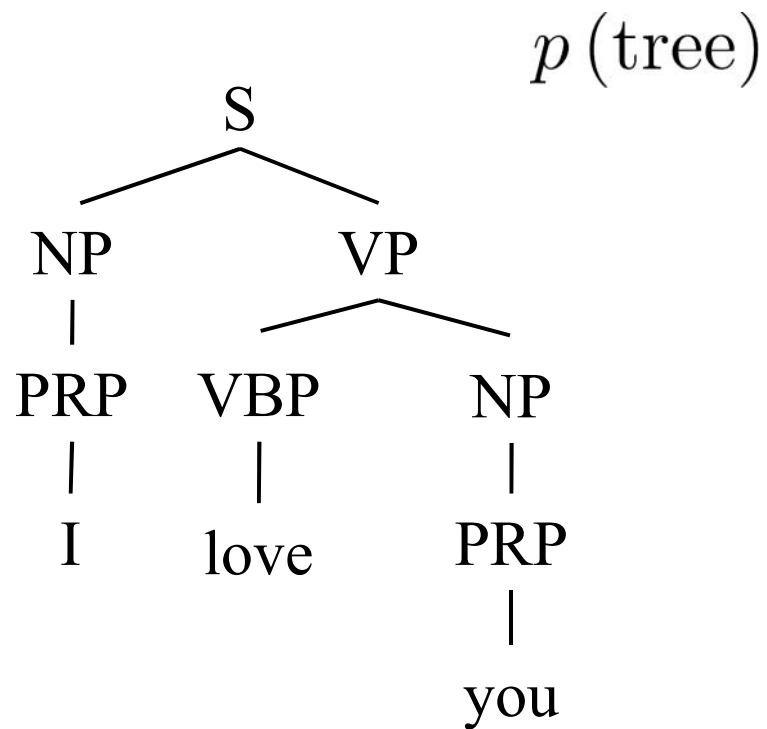NII

Akinori Fujino
NTT CS Labs.
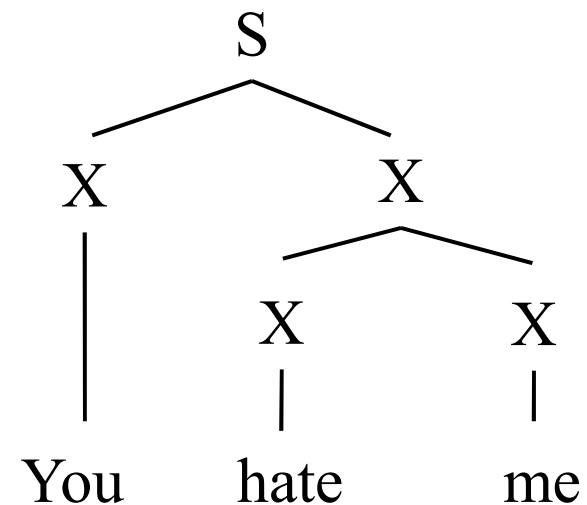
Masaaki Nagata
NTT CS Labs.

July 10
ACL 2012

# Task: Statistical Constituent Parsing

**Training**

**Testing**

$p\,(\text{tree})$

```
            S
          /   \
        NP      VP
        |      /  \
       PRP   VBP   NP
        |     |     |
        I    love  PRP
                    |
                   you
```

```
            S
          /   \
        X        X
        |      /   \
       You    X     X
              |     |
             hate   me
```
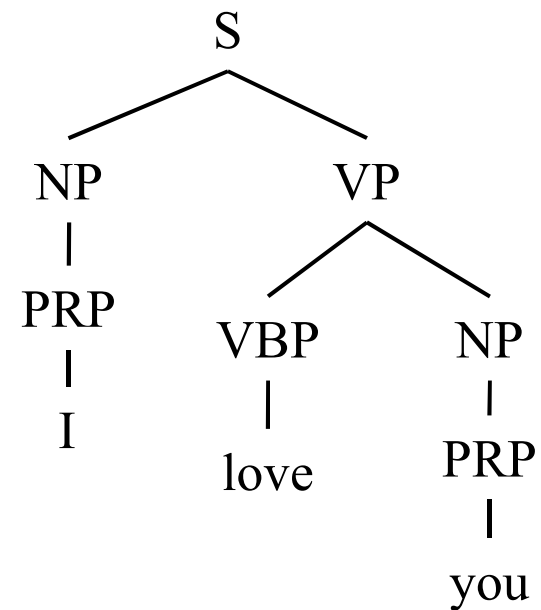
F-score:  ?

# Previous Work

- Naive CFG-based parser does NOT perform well...

  - Coarse symbol annotations

  - Strong independence assumption

Previous approaches:

a) CFG with automatic symbol refinement
   [Matsuzaki et al. 05, Petrov et al. 06]

b) Tree substitution grammars (TSG)
   [Cohn et al. 09, Post et al. 09]

# a) CFG with automatic symbol refinement

- Idea:  split symbols into subcategories
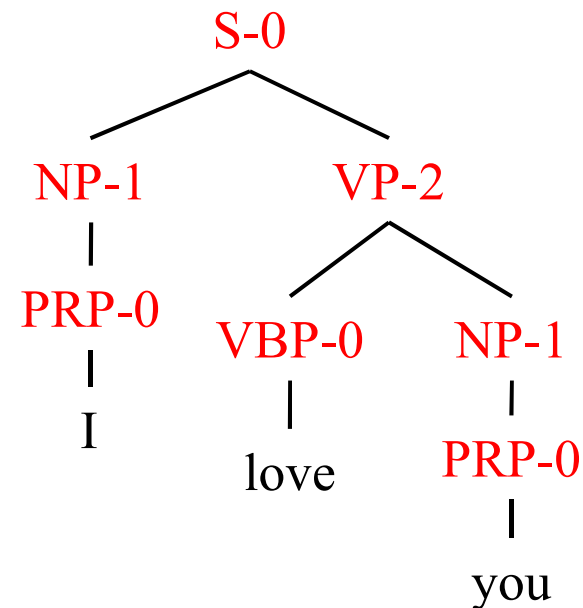  (based on the likelihood)

- Inference:  EM algorithm

# a) CFG with automatic symbol refinement

- Idea: split symbols into subcategories
  (based on the likelihood)
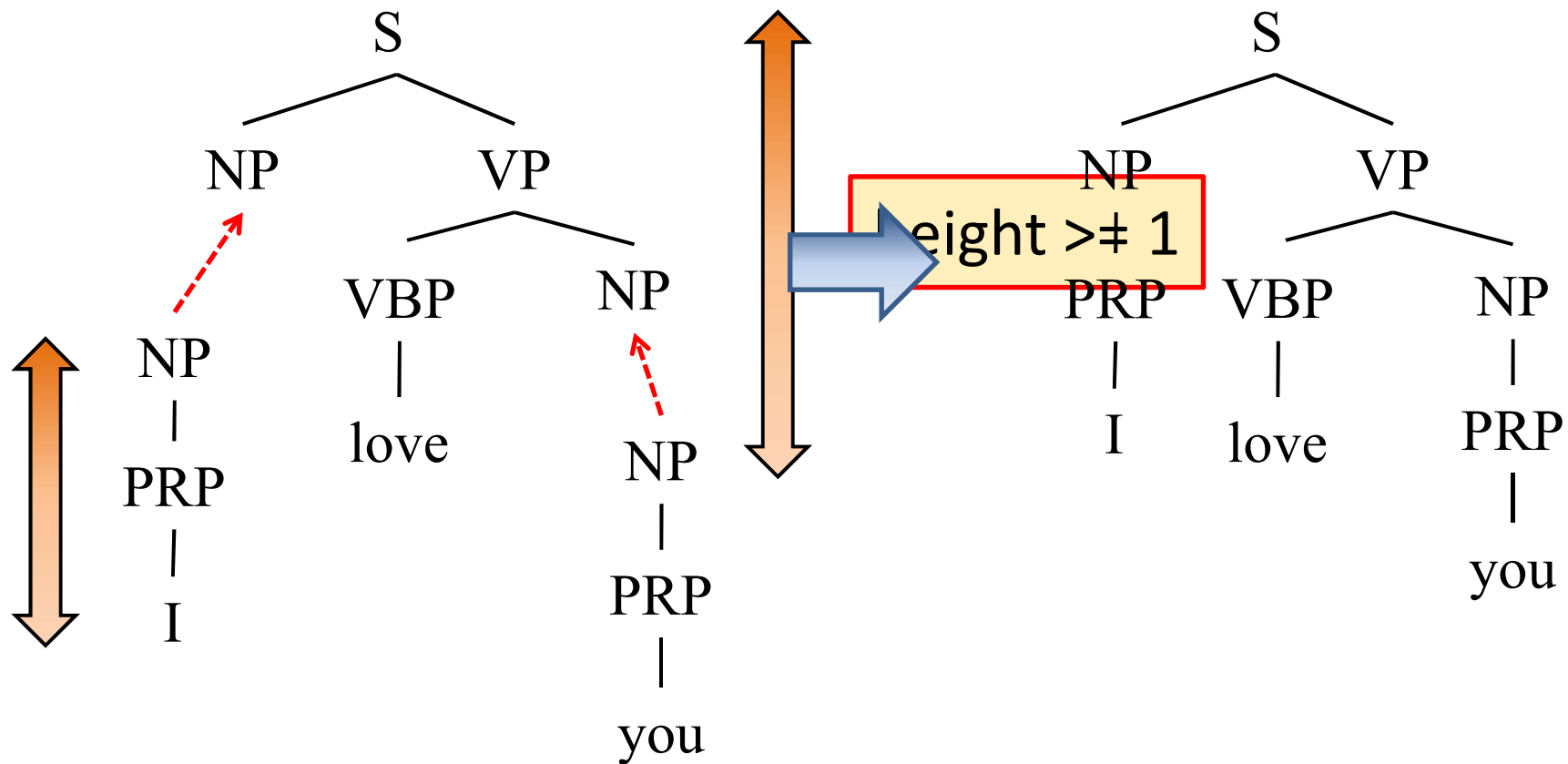
- Inference: EM algorithm

Refined PCFG rules:

S-0    →    NP-1  VP-2    $p_0$

NP-1   →    PRP-0    $p_1$

PRP-0  →    "I"    $p_2$

VP-2   →    VBP-0  NP-1    $p_3$

# b) Tree Substitution Grammars (TSG)

- Idea:  Allow arbitrarily large tree fragments

- Inference:  MCMC sampling

# Motivation

Two approaches are complementary

[Cohn et al. 09, Bansal & Klein 10]

- Clustering context

- Learning structure

Proposal:

Symbol-Refined TSG (SR-TSG)

# SR-TSG



TSG

```
        S
      /   \
    NP      VP
    |      /  \
    NP   VBP   NP
    |    |     |
   PRP  love   NP
    |          |
    I         PRP
               |
              you
```

SR-TSG

```
        S-0
      /     \
   NP-0      VP-2
    |       /    \
   NP-0   VBP-0  NP-1
    |      |      |
  PRP-0  love   NP-1
    |            |
    I          PRP-0
                |
               you
```

# SR-TSG

- Latent variables of SR-TSG:

| tree fragments | + | symbol subcategories |
|---|---|---|

## Search space is huge !

Challenge:

Fully automatic learning

⟷ TSG + refinement
[Bansal & Klein 10]

S-0
NP-1    VP-2
NP-1    VBP-0    NP-1
PRP-0    love    NP-1
I    PRP-0
you

| Probabilistic model |
|---|

| Inference |
|---|

# Probabilistic Model of SR-TSG

# Overview of Probabilistic Model

What we need:

1. Probability distribution over refined tree fragments

Pitman-Yor process

[Pitman & Yor 96]

2. Framework for back-off smoothing

Large $\rightarrow$ small tree fragments

3-level hierarchy

# 3-Level Hierarchy

complex ← simple

Pitman-Yor process

Sparse!

1. SR-TSG

2. SR-CFG

3. RU-CFG



S-0
├── NP-1
└── VP-0
    ├── VBP-3
    └── NP-2

S-0
├── NP-1
└── VP-0

S
├── NP-1
└── VP-0

# Pitman-Yor Process (PYP)

- A prior for non-parametric Bayesian model

- Useful for modeling data with power-law distribution

- Closely related to Chinese Restaurant Process (CRP)

# Pitman-Yor Process (PYP)



1. SR-TSG

PYP

2. SR-CFG

A
```
        S-0
       /    \
    NP-1    VP-0
           /    \
        VBP-3   NP-2
```

A1
```
      S-0
     /    \
  NP-1    VP-0
```

A2
```
       VP-0
      /    \
   VBP-3   NP-2
```

Eq. $p(A) = \alpha \cdot count(A) + \beta \cdot P_0(A)$   Base distr.

"Rich get richer" effect

Chinese Restaurant Process

$\propto p'(A1) \times p'(A2)$

smoothing

$$\text{PYP}\left(d, \theta, P_0\right)$$

A
S-0

NP-1    VP-0

VBP-3    NP-2

(A)

1

$\times P_0\left(A\right)$

# Chinese Restaurant Process

$$\mathrm{PYP}\,(d, \theta, P_0)$$

A

$$1 - d \qquad \theta + d$$

S-0

NP-1     VP-0

VBP-3     NP-2

# Chinese Restaurant Process

$$\text{PYP}\,(d, \theta, P_0)$$



$$2 - d \qquad \theta + d$$

$$\times P_0\,(B)$$

- Prob. of choosing a table $\propto \begin{cases} \text{count}\,(k) - d & \text{occupied} \\ \theta + \#\,\text{tables} \cdot d & \text{new table} \end{cases}$

- "Rich get richer" effect

# Chinese Restaurant Process

$$\text{PYP}(d, \theta, P_0)$$

S-0

NP-1    VP-0

NP-1    VBP-3    NP-2

| A | B | |
|---|---|---|
| $2 - d$ | $1 - d$ | $\theta + 2d$ |

PRP-0

I

- Prob. of choosing a table $\propto \begin{cases} \text{count}(k) - d & \text{occupied} \\ \theta + \# \text{ tables} \cdot d & \text{new table} \end{cases}$

- "Rich get richer" effect

# Chinese Restaurant Process

$$PYP\,(d, \theta, P_0)$$



A      B

$$3 - d \qquad 1 - d \qquad \theta + 2d$$

S-0

NP-1    VP-0

NP-1     VBP-3    NP-2

PRP-0

I

- Prob. of choosing a table $\propto \begin{cases} \text{count}\,(k) - d & \text{occupied} \\ \theta + \#\,\text{tables} \cdot d & \text{new table} \end{cases}$

- "Rich get richer" effect

# Chinese Restaurant Process

# Summary of Probabilistic Model

- Probability distribution over SR-TSG fragments

  - Pitman-Yor process as a prior

- 3-level hierarchy for back-off smoothing

  - SR-TSG $\leftarrow$ SR-CFG $\leftarrow$ RU-CFG

- Parsing accuracy(f-score):  91.1% on English PTB

# Inference

# Overview of Inference

**What we want:**

# Overview of Inference

- MAP estimation: $\underset{Z}{\text{argmax}}\ p\left(\mathbf{Z}\,|\,\mathbf{T}\right)$

- MCMC sampling for inference

| tree fragments | | symbol subcategories |
|:---:|:---:|:---:|
| tree fragments | **+** | symbol subcategories |

**Stepwise training:**

1.      Fix                     Train

2.      Train              Almost fixed

# Overview of MCMC Sampling

S
NP          VP
|          /    \
PRP      VBP     NP
|         |       |
I        love    PRP
                  |
                 you

# Overview of MCMC Sampling

# Overview of MCMC Sampling

S-2

NP-0         VP-1

PRP-1   VBP-3   NP-2

I      love    PRP-1

you

# Overview of MCMC Sampling

S-0
NP-2            VP-2
PRP-0   VBP-1       NP-1
I        love       PRP-0
                     you

# Overview of MCMC Sampling

S-2
NP-0          VP-0
PRP-1   VBP-3   NP-1
I       love    PRP-1
                you

# Overview of MCMC Sampling

```
                         S-0
                     /        \
                  NP-0        VP-2
                   |          /    \
                 PRP-0    VBP-0    NP-1
                   |        |       |
                   I       love   PRP-0
                                    |
                                   you
```

# Overview of MCMC Sampling

# Overview of MCMC Sampling

# Overview of MCMC Sampling

S-0

NP-0            VP-2

NP-0                              VP-2

PRP-0                    VBP-0            NP-1

I                       love            PRP-0

you

# Overview of MCMC Sampling



S-0

NP-0          VP-2

VBP-0     NP-1

NP-0          love      PRP-0

PRP-0                    you

I

# Overview of MCMC Sampling

# Overview of MCMC Sampling

# Inference of Symbol Refinement

**Problem:**

 - Gibbs sampler is inefficient

   Update only one variable at a time

# Inference of Symbol Refinement

Proposal:

- 3 types of blocked samplers to find better solution

For each MCMC iteration...

a) Run <u>sentence</u> sampler

b) Run <u>subtree</u> sampler

c) Run <u>restaurant</u> sampler

# a) Sentence Sampler

- Metropolis-Hastings (MH) algorithm [Johnson et al. 07]

For each sentence...

1. Run inside algorithm
2. Sample a derivation
3. Accept or reject the sample

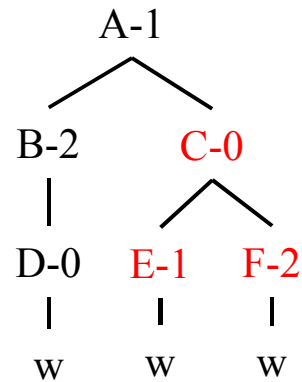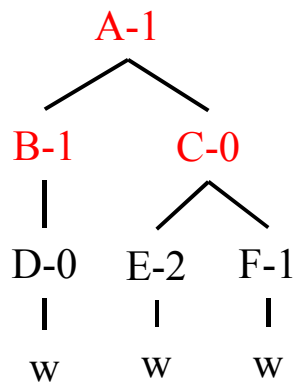# a) Sentence Sampler

- Update a sentence at a time

# b) Subtree Sampler
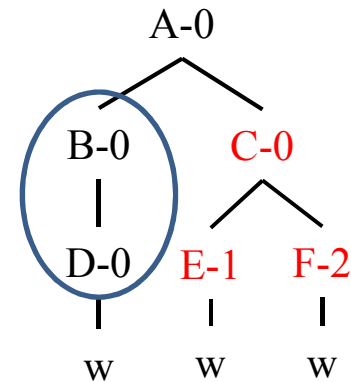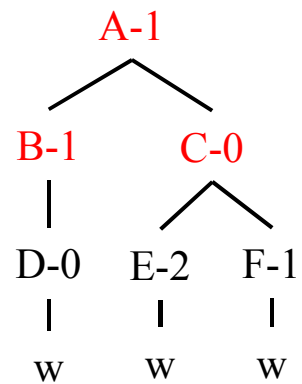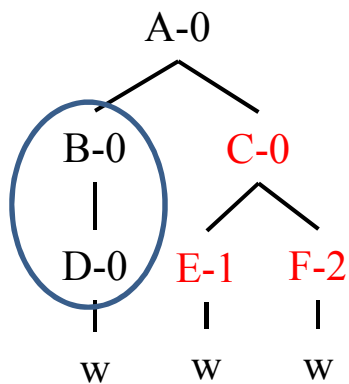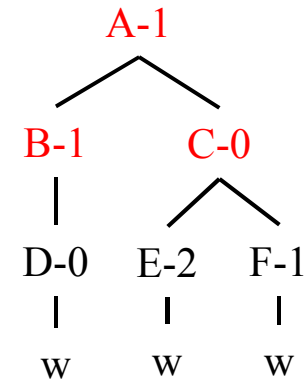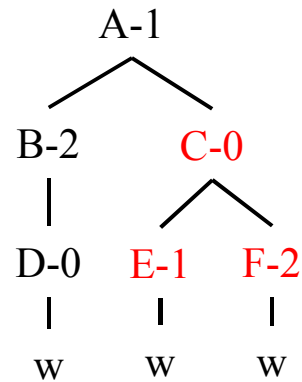
- Update the same type of subtree simultaneously

# b) Subtree Sampler

- Update the same type of subtree simultaneously

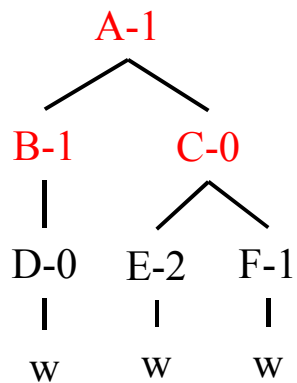# b) Subtree Sampler

- Update the same type of subtree simultaneously
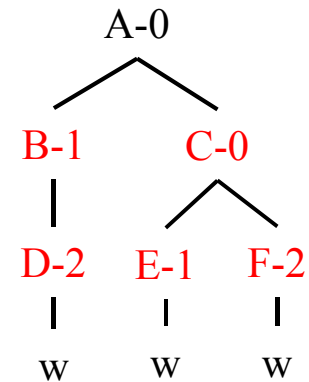
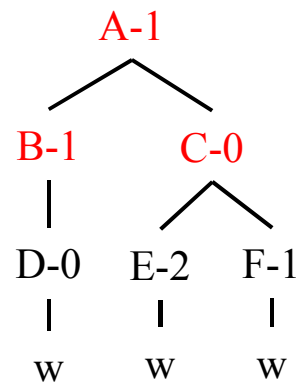# b) Subtree Sampler

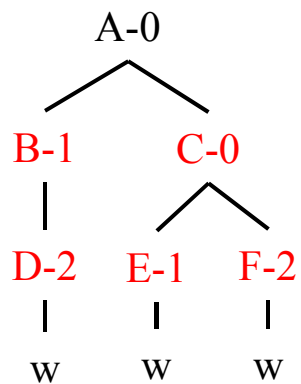- Update the same type of subtree simultaneously

# b) Subtree Sampler

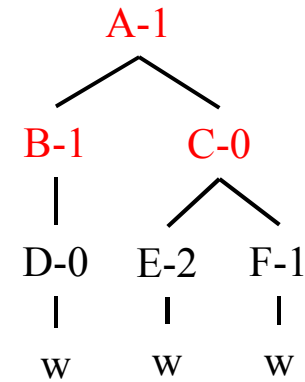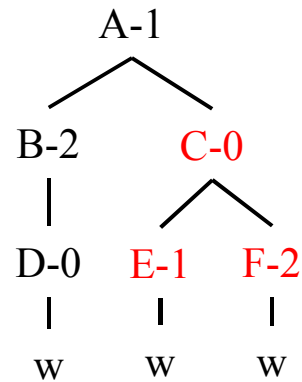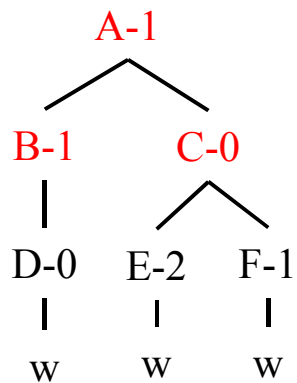- Update the same type of subtree simultaneously

# b) Subtree Sampler

- Update the same type of subtree simultaneously

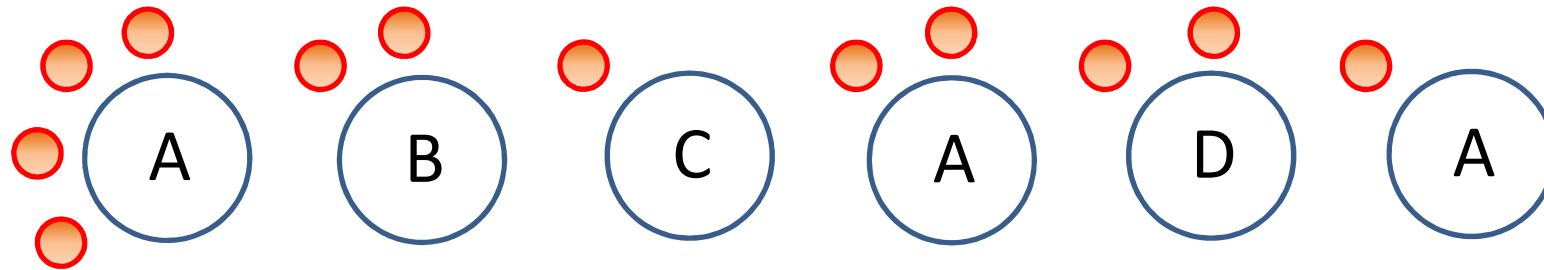# c) Restaurant Sampler

- Update # of tables

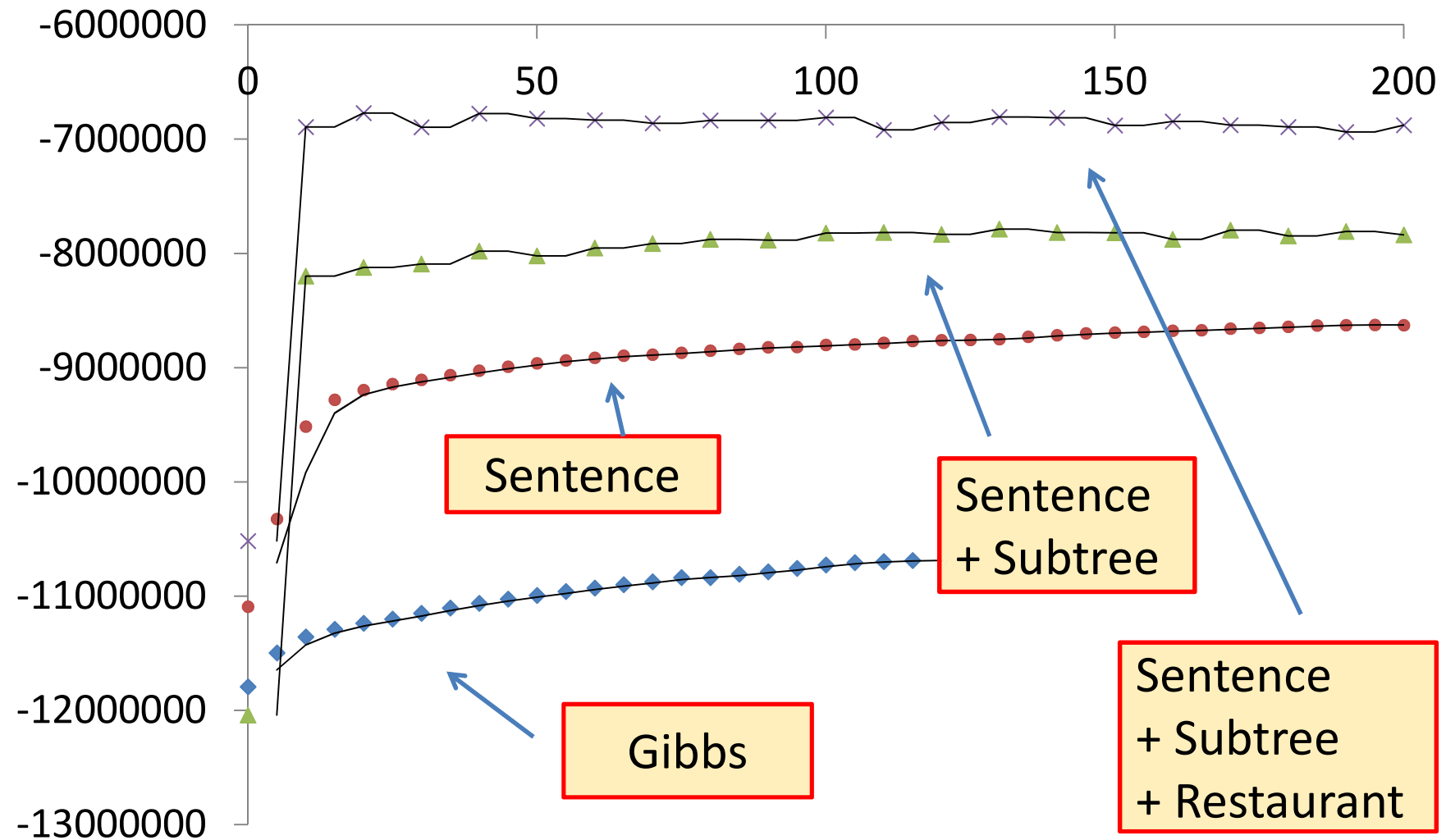# c) Restaurant Sampler

- Update # of tables



For each table label …

- Keep count       count = 6

- Remove all

- Add one by one   Helpful for optimizing # of tables

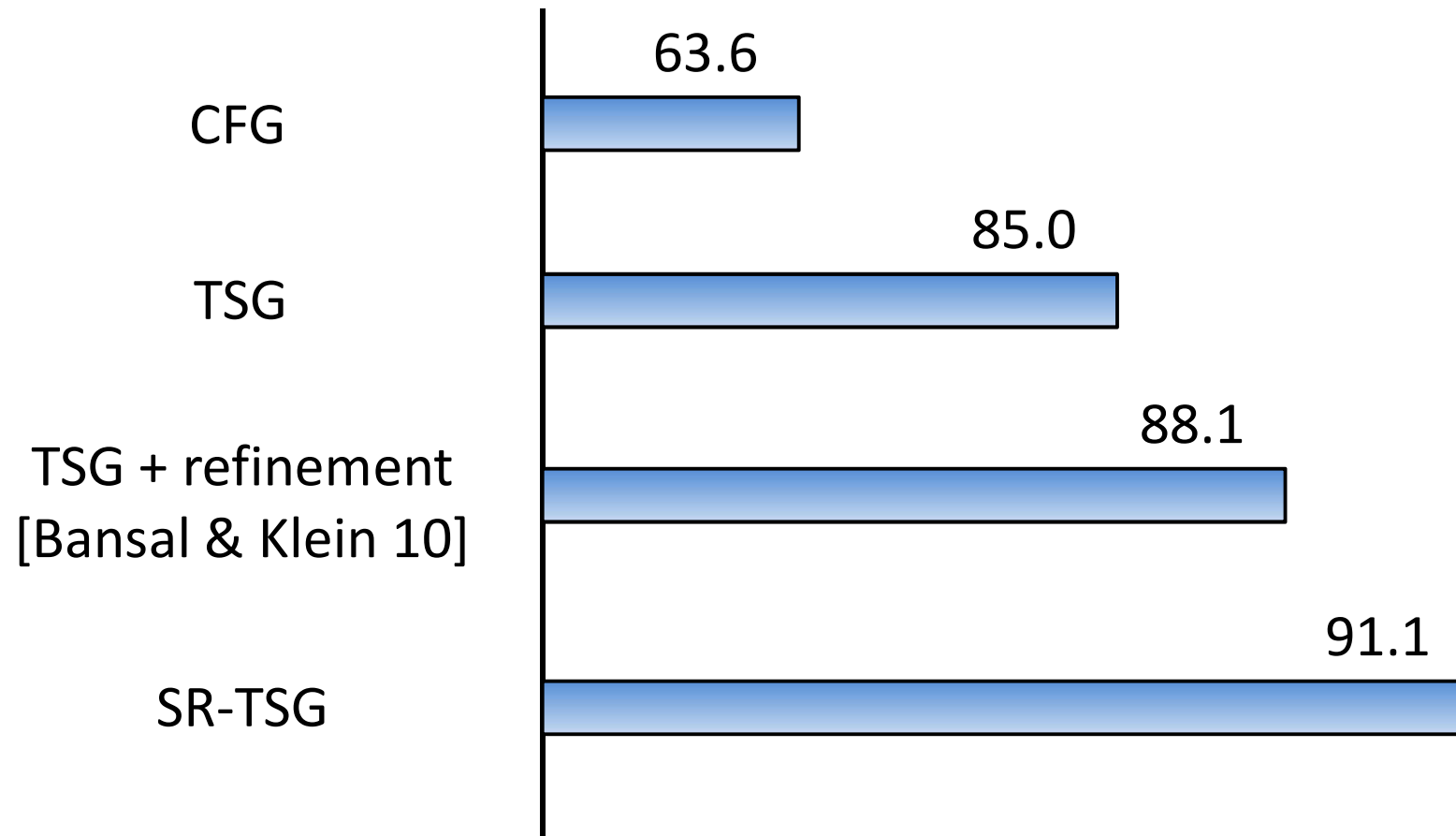# Effect of MCMC Samplers

# Summary of Inference

- MCMC sampling

- Stepwise training for the induction of latent variables

- 3 types of blocked samplers
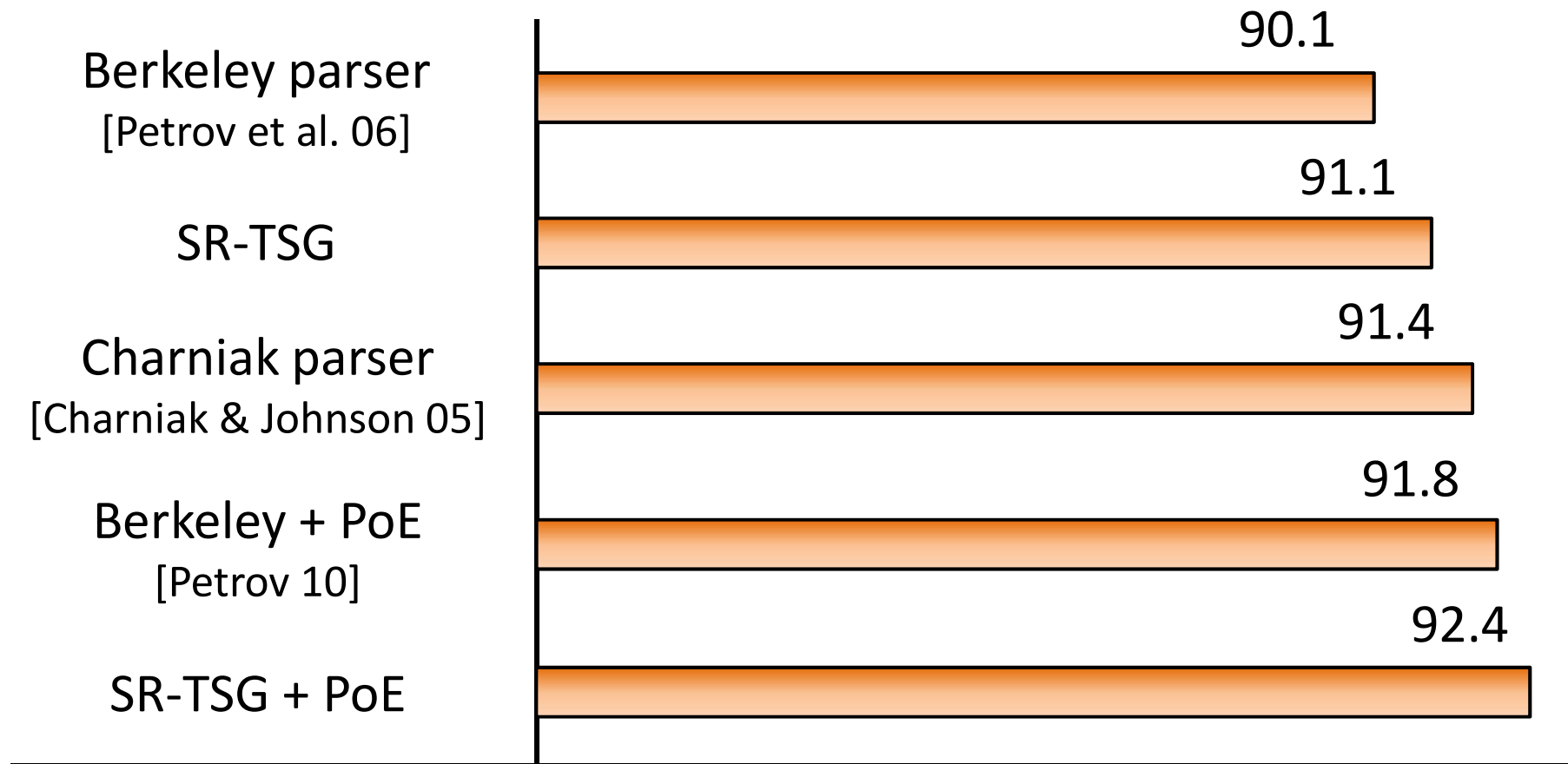
    - Effective for finding better solution

# Experiment

Parsing Accuracy

CFG — 63.6
TSG — 85.0
TSG + refinement [Bansal & Klein 10] — 88.1
SR-TSG — 91.1

# Final Result



Berkeley parser
[Petrov et al. 06]
**90.1**

SR-TSG
**91.1**

Charniak parser
[Charniak & Johnson 05]
**91.4**

Berkeley + PoE
[Petrov 10]
**91.8**

SR-TSG + PoE
**92.4**

- System combination [Zhang et al. 09]

- Self-training [Huang et al. 10]

are better than SR-TSG.

# Conclusion

**Approach:**

SR-TSG = TSG + symbol refinement
- Fully automatic learning

**Probabilistic Model:**

Pitman-Yor process + 3-level hierarchy

**Inference:**

3 types of MCMC samplers for efficient training
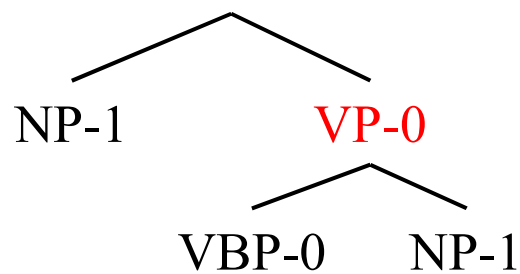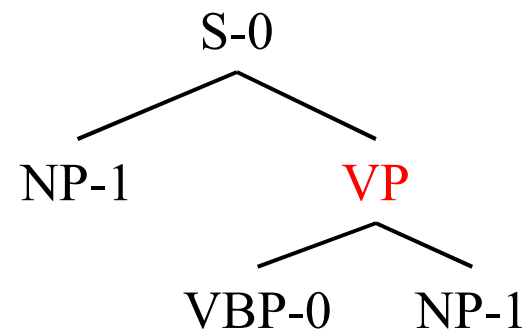
**Result:**

State-of-the-art

Thank you.

# Internal Subcategory Marginalization

- Encourage the model to find large tree fragments