

How do we choose the grammar complexity?

Grammar induction:

How many grammar **symbols** (NP, VP, etc.)?

?

She heard the noise

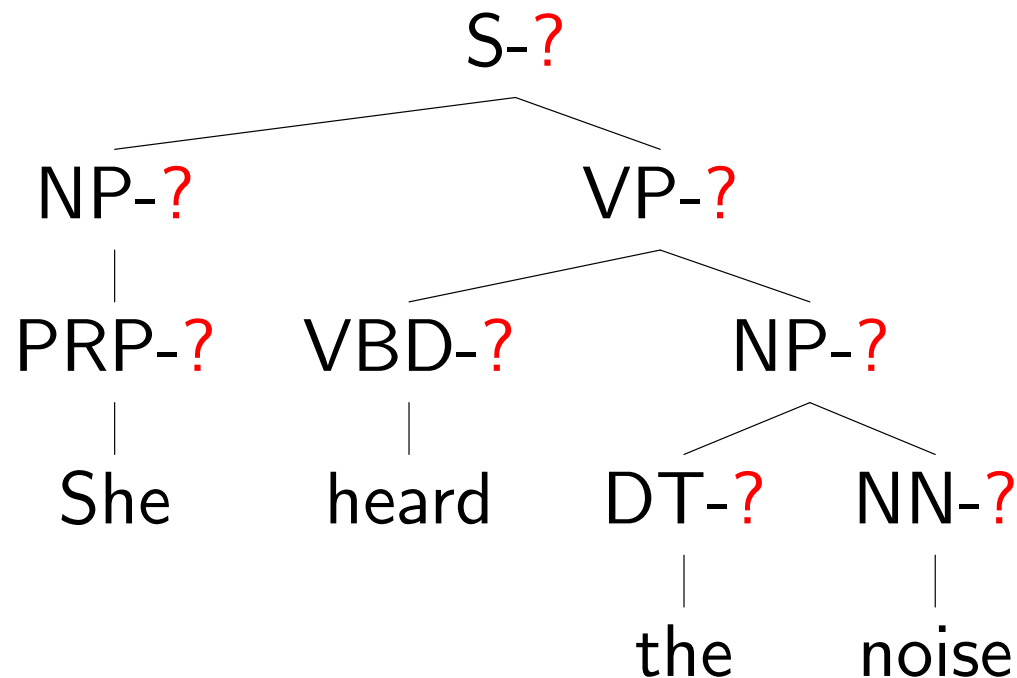
How do we choose the grammar complexity?

Grammar induction:

How many grammar symbols (NP, VP, etc.)?

Grammar refinement:

How many grammar subsymbols (NP-*loc*, NP-*subj*, etc.)?



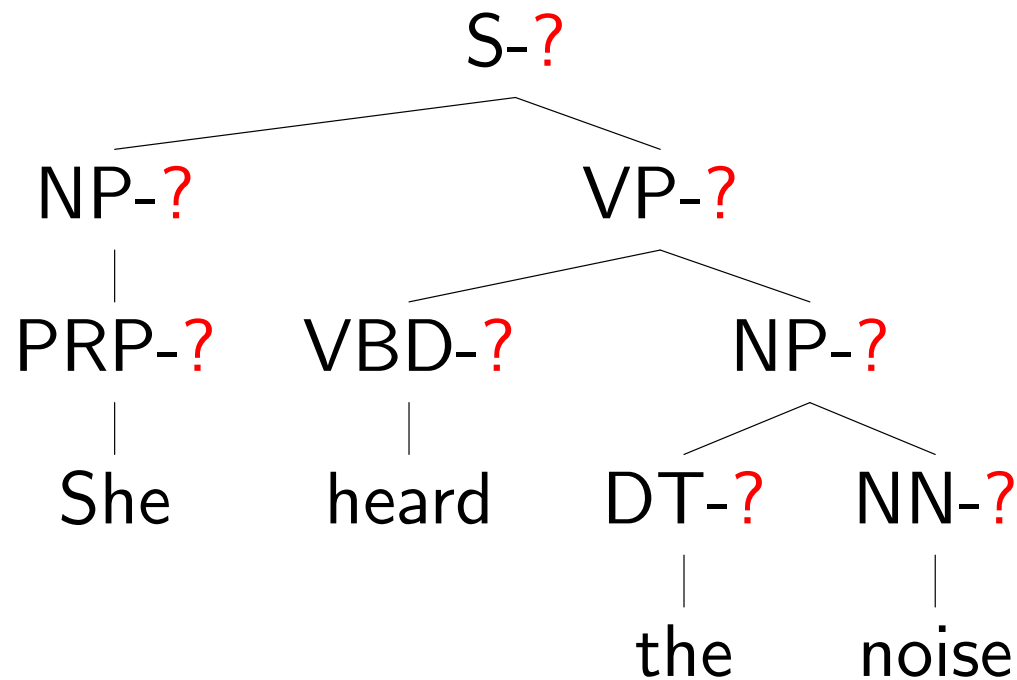
How do we choose the grammar complexity?

Grammar induction:

How many grammar symbols (NP, VP, etc.)?

Grammar refinement:

How many grammar subsymbols (NP-*loc*, NP-*subj*, etc.)?



Our solution: the HDP-PCFG allows number of (sub)symbols to adapt to data

A motivating example

True grammar:

$S \rightarrow AA \mid BB \mid CC \mid DD$

$A \rightarrow a_1 \mid a_2 \mid a_3$

$B \rightarrow b_1 \mid b_2 \mid b_3$

$C \rightarrow c_1 \mid c_2 \mid c_3$

$D \rightarrow d_1 \mid d_2 \mid d_3$

A motivating example

True grammar:

$S \rightarrow AA \mid BB \mid CC \mid DD$

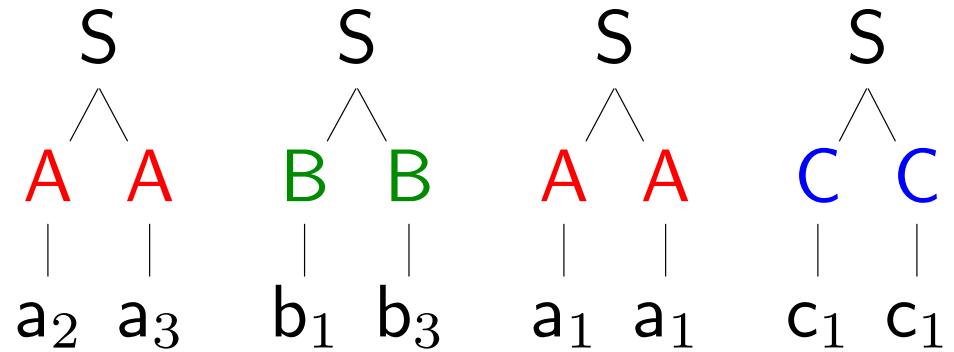
$A \rightarrow a_1 \mid a_2 \mid a_3$

$B \rightarrow b_1 \mid b_2 \mid b_3$

$C \rightarrow c_1 \mid c_2 \mid c_3$

$D \rightarrow d_1 \mid d_2 \mid d_3$

Generate examples:



A motivating example

True grammar:

$S \rightarrow AA \mid BB \mid CC \mid DD$

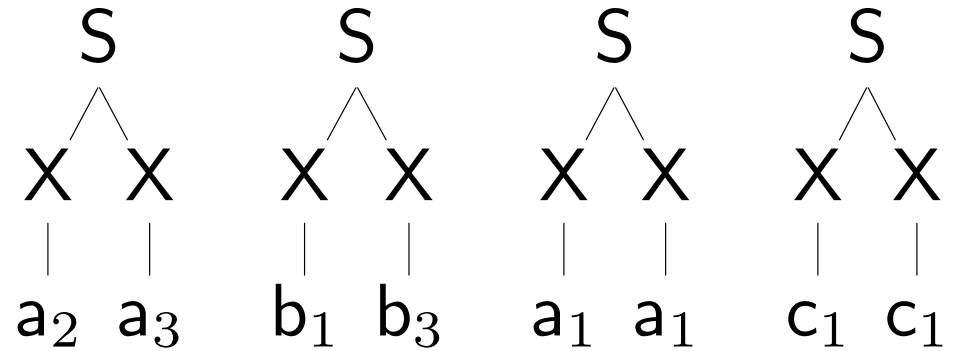
$A \rightarrow a_1 \mid a_2 \mid a_3$

$B \rightarrow b_1 \mid b_2 \mid b_3$

$C \rightarrow c_1 \mid c_2 \mid c_3$

$D \rightarrow d_1 \mid d_2 \mid d_3$

Collapse $A, B, C, D \Rightarrow X$:



A motivating example

True grammar:

$S \rightarrow AA \mid BB \mid CC \mid DD$

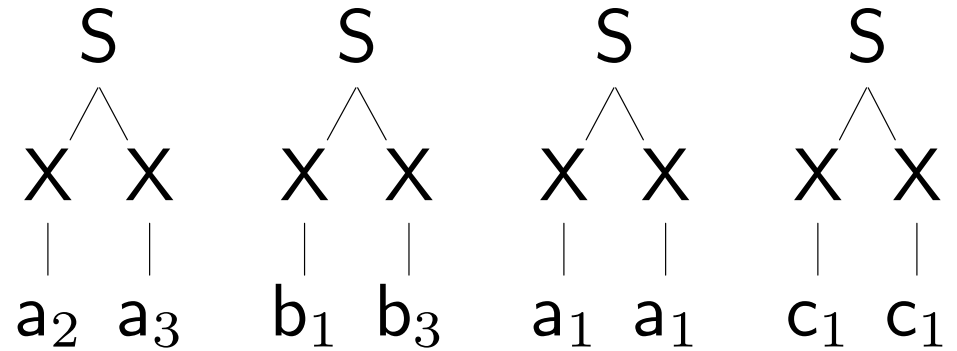
$A \rightarrow a_1 \mid a_2 \mid a_3$

$B \rightarrow b_1 \mid b_2 \mid b_3$

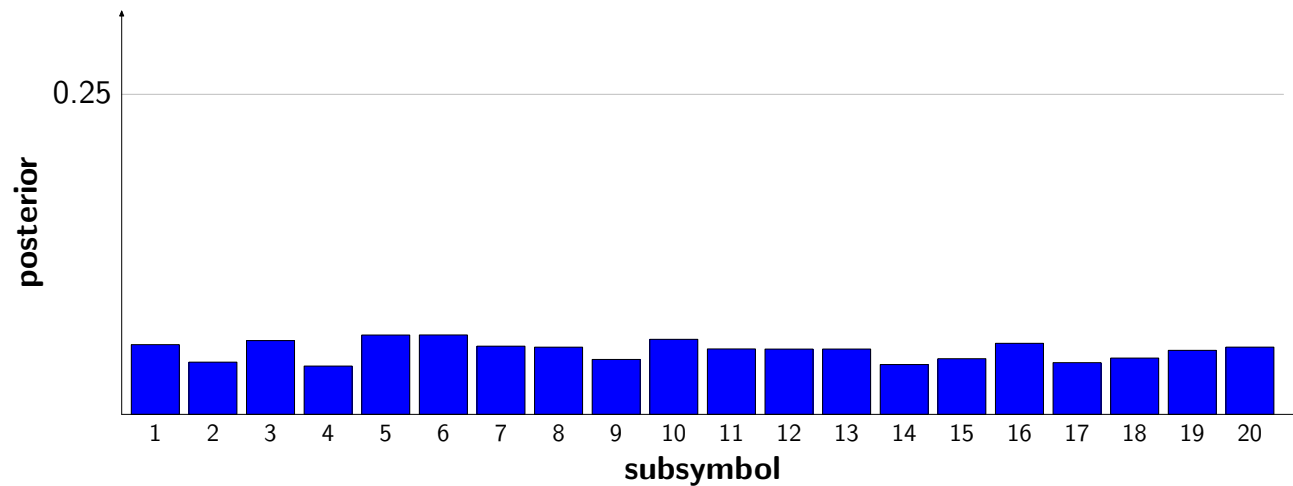
$C \rightarrow c_1 \mid c_2 \mid c_3$

$D \rightarrow d_1 \mid d_2 \mid d_3$

Collapse A,B,C,D \Rightarrow X:



Results:



standard PCFG

A motivating example

True grammar:

$S \rightarrow AA \mid BB \mid CC \mid DD$

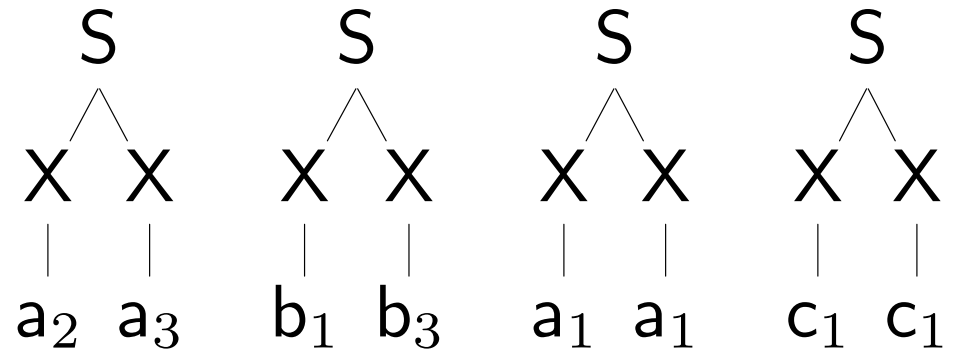
$A \rightarrow a_1 \mid a_2 \mid a_3$

$B \rightarrow b_1 \mid b_2 \mid b_3$

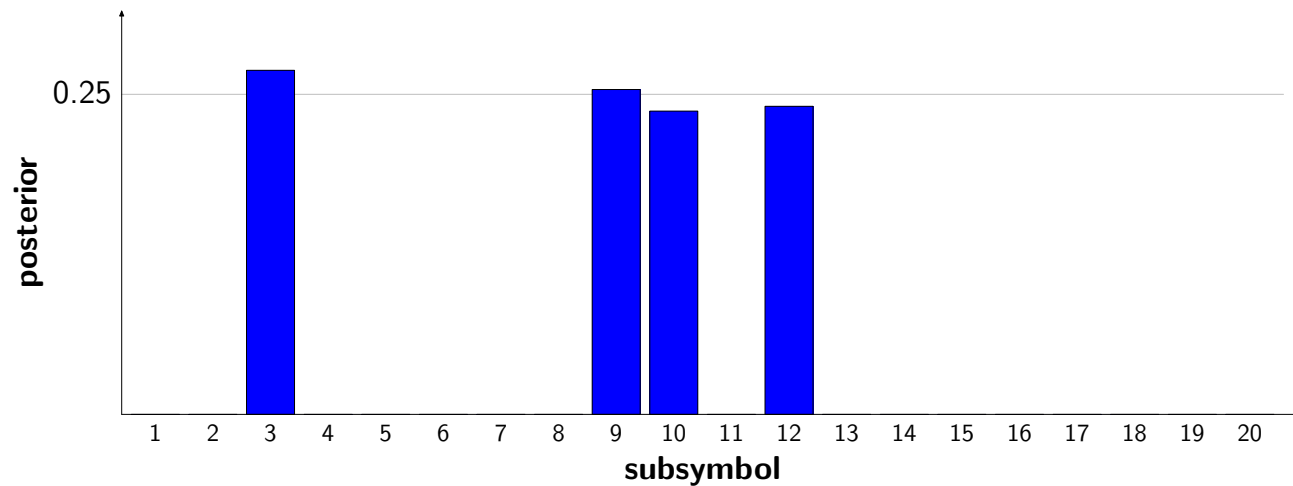
$C \rightarrow c_1 \mid c_2 \mid c_3$

$D \rightarrow d_1 \mid d_2 \mid d_3$

Collapse A,B,C,D \Rightarrow X:



Results:



HDP-PCFG

The meeting of two fields

Grammar learning 

Lexicalized

[Charniak, 1996]

[Collins, 1999]

Manual refinement

[Johnson, 1998]

[Klein, Manning, 2003]

Automatic refinement

[Matsuzaki, et al., 2005]

[Petrov, et al., 2006]

The meeting of two fields

Grammar learning

Lexicalized

[Charniak, 1996]

[Collins, 1999]

Manual refinement

[Johnson, 1998]

[Klein, Manning, 2003]

Automatic refinement

[Matsuzaki, et al., 2005]

[Petrov, et al., 2006]

Bayesian nonparametrics ○○○○...

Basic theory

[Ferguson, 1973]

[Antoniak, 1974]

[Sethuraman, 1994]

[Escobar, West, 1995]

[Neal, 2000]

More complex models

[Teh, et al., 2006]

[Beal, et al., 2002]

[Goldwater, et al., 2006]

[Sohn, Xing, 2007]

The meeting of two fields

Grammar learning

Lexicalized

[Charniak, 1996]

[Collins, 1999]

Manual refinement

[Johnson, 1998]

[Klein, Manning, 2003]

Automatic refinement

[Matsuzaki, et al., 2005]

[Petrov, et al., 2006]

Bayesian nonparametrics

Basic theory

[Ferguson, 1973]

[Antoniak, 1974]

[Sethuraman, 1994]

[Escobar, West, 1995]

[Neal, 2000]

More complex models

[Teh, et al., 2006]

[Beal, et al., 2002]

[Goldwater, et al., 2006]

[Sohn, Xing, 2007]

Nonparametric grammars

[Johnson, et al., 2006]

[Finkel, et al., 2007]

[Liang, et al., 2007]

The meeting of two fields

Grammar learning 

Lexicalized

[Charniak, 1996]

[Collins, 1999]

Bayesian nonparametrics ○○○○...

Basic theory

[Ferguson, 1973]

[Antoniak, 1974]

Ma

Our contribution

- Definition of the HDP-PCFG
- Simple and efficient variational inference algorithm
- Empirical comparison with finite models on a full-scale parsing task

Auto

Nonparametric grammars

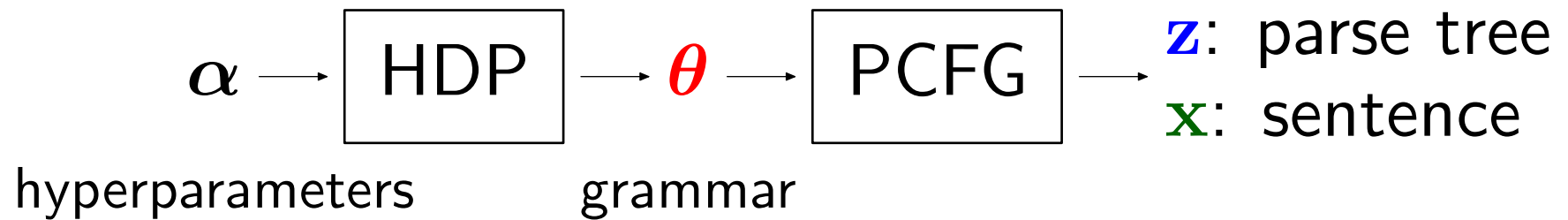
[Johnson, et al., 2006]

[Finkel, et al., 2007]

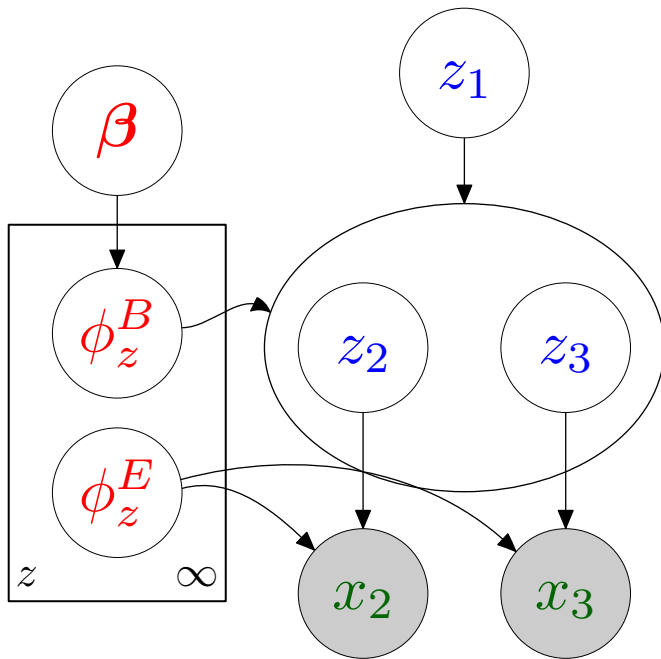
[Liang, et al., 2007]

Bayesian paradigm

Generative model:



HDP probabilistic context-free grammars



HDP-PCFG

$\beta \sim \text{GEM}(\alpha)$ [generate distribution over symbols]

For each symbol $z \in \{1, 2, \dots\}$:

$\phi_z^E \sim \text{Dirichlet}(\alpha^E)$ [generate emission probs]

$\phi_z^B \sim \text{DP}(\alpha^B, \beta\beta^T)$ [binary production probs]

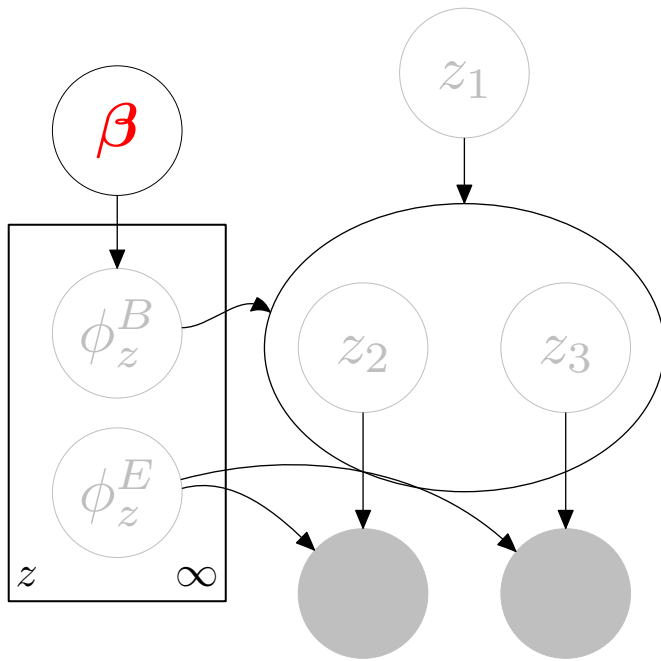
For each nonterminal node i :

$(z_{L(i)}, z_{R(i)}) \sim \text{Multinomial}(\phi_{z_i}^B)$ [child symbols]

For each preterminal node i :

$x_i \sim \text{Multinomial}(\phi_{z_i}^E)$ [terminal symbol]

HDP probabilistic context-free grammars



HDP-PCFG

$\beta \sim \text{GEM}(\alpha)$ [generate distribution over symbols]

For each symbol $z \in \{1, 2, \dots\}$:

$\phi_z^E \sim \text{Dirichlet}(\alpha^E)$ [generate emission probs]

$\phi_z^B \sim \text{DP}(\alpha^B, \beta\beta^T)$ [binary production probs]

For each nonterminal node i :

$(z_{L(i)}, z_{R(i)}) \sim \text{Multinomial}(\phi_{z_i}^B)$ [child symbols]

For each preterminal node i :

$x_i \sim \text{Multinomial}(\phi_{z_i}^E)$ [terminal symbol]

HDP-PCFG: prior over symbols

$$\beta \sim \text{GEM}(\alpha)$$

$\alpha = 1$



GEM



HDP-PCFG: prior over symbols

$$\beta \sim \text{GEM}(\alpha)$$

$\alpha = 1$



GEM



HDP-PCFG: prior over symbols

$$\beta \sim \text{GEM}(\alpha)$$

$\alpha = 1$



GEM



HDP-PCFG: prior over symbols

$$\beta \sim \text{GEM}(\alpha)$$

$\alpha = 1$



GEM



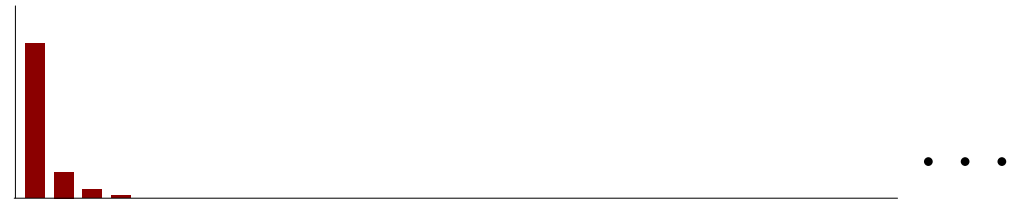
HDP-PCFG: prior over symbols

$$\beta \sim \text{GEM}(\alpha)$$

$\alpha = 0.5$



GEM



$\alpha = 1$

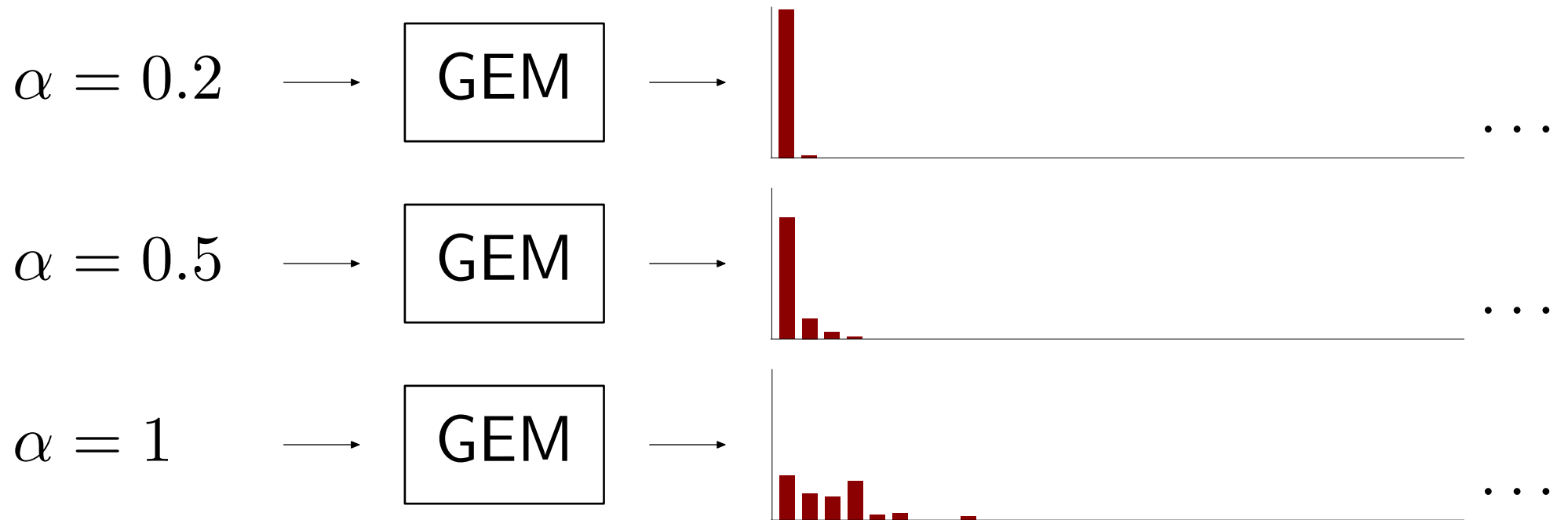


GEM



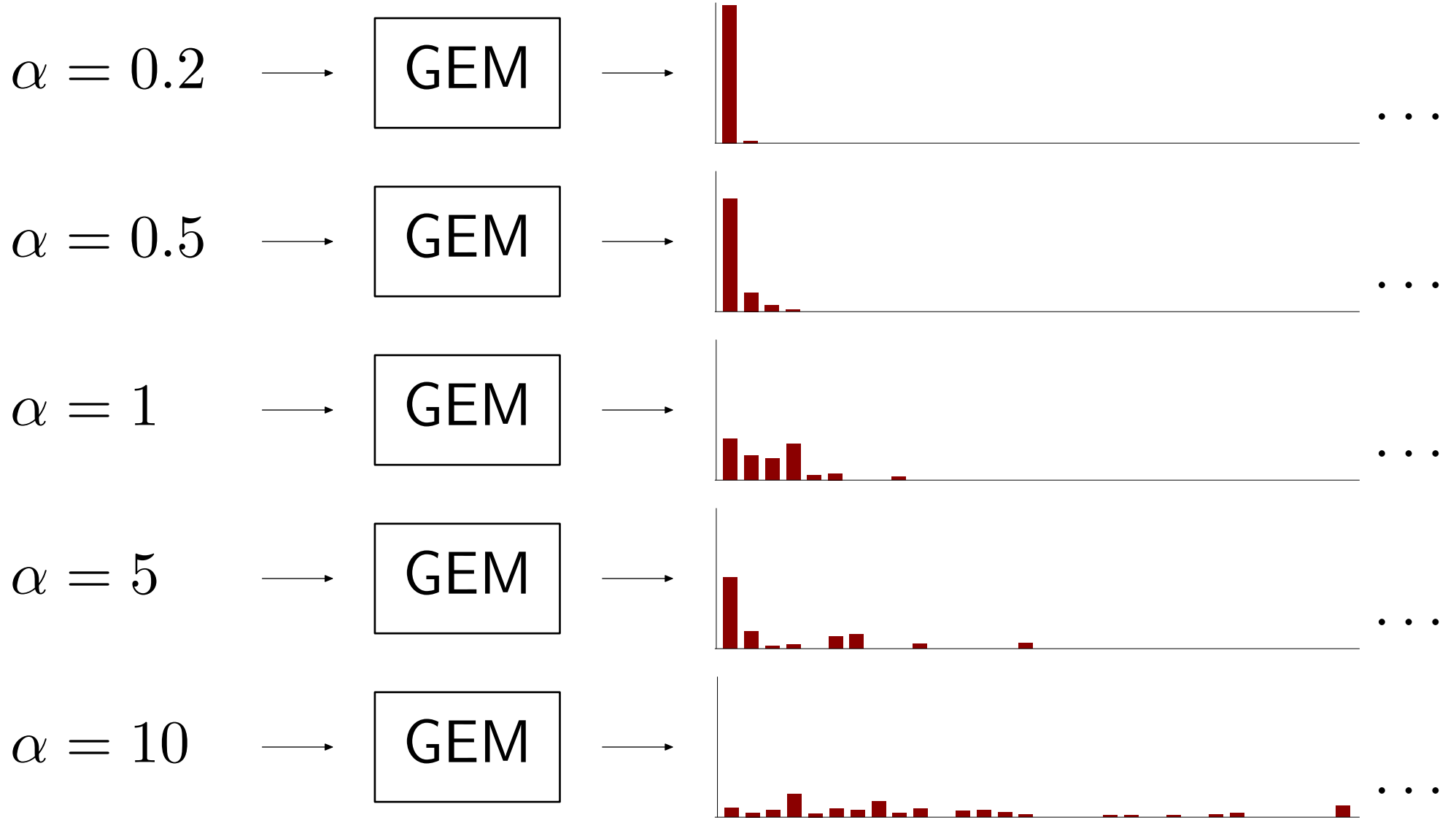
HDP-PCFG: prior over symbols

$$\beta \sim \text{GEM}(\alpha)$$

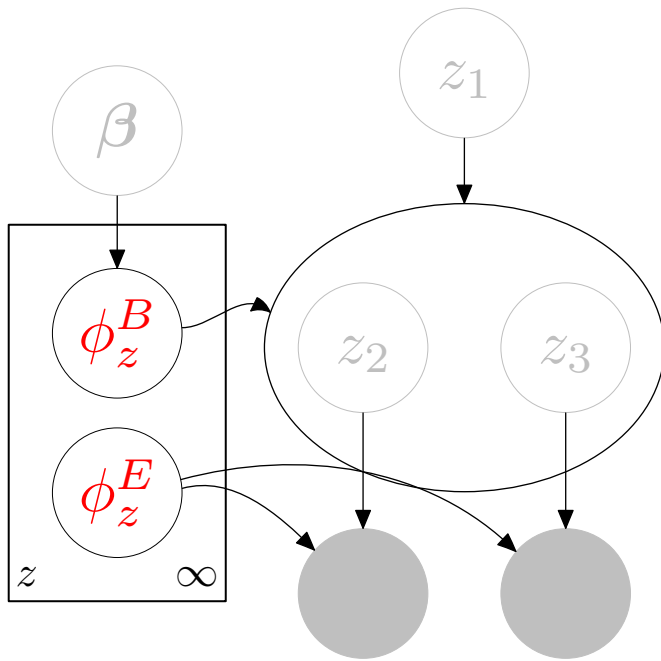


HDP-PCFG: prior over symbols

$$\beta \sim \text{GEM}(\alpha)$$



HDP probabilistic context-free grammars



HDP-PCFG

$\beta \sim \text{GEM}(\alpha)$ [generate distribution over symbols]

For each symbol $z \in \{1, 2, \dots\}$:

$\phi_z^E \sim \text{Dirichlet}(\alpha^E)$ [generate emission probs]

$\phi_z^B \sim \text{DP}(\alpha^B, \beta\beta^T)$ [binary production probs]

For each nonterminal node i :

$(z_{L(i)}, z_{R(i)}) \sim \text{Multinomial}(\phi_{z_i}^B)$ [child symbols]

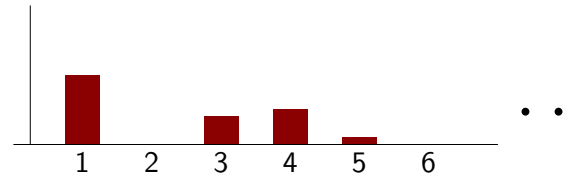
For each preterminal node i :

$x_i \sim \text{Multinomial}(\phi_{z_i}^E)$ [terminal symbol]

HDP-PCFG: binary productions

Distribution over symbols (top-level):

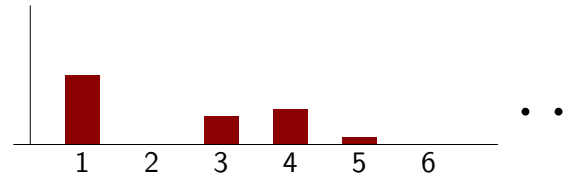
$$\beta \sim \text{GEM}(\alpha)$$



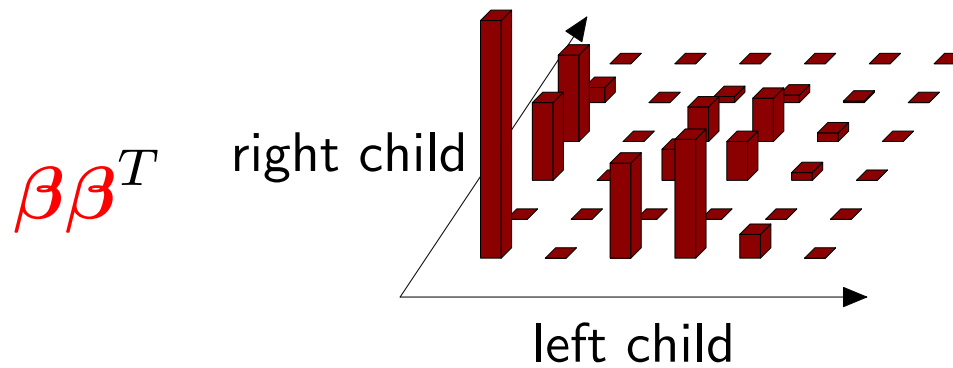
HDP-PCFG: binary productions

Distribution over symbols (top-level):

$$\beta \sim \text{GEM}(\alpha)$$



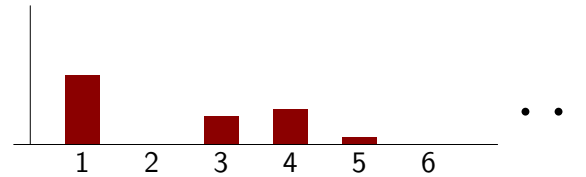
Mean distribution over child symbols:



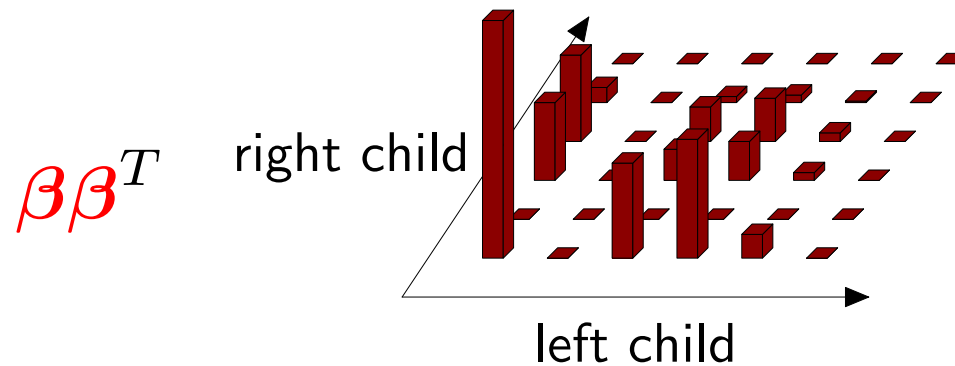
HDP-PCFG: binary productions

Distribution over symbols (top-level):

$$\beta \sim \text{GEM}(\alpha)$$

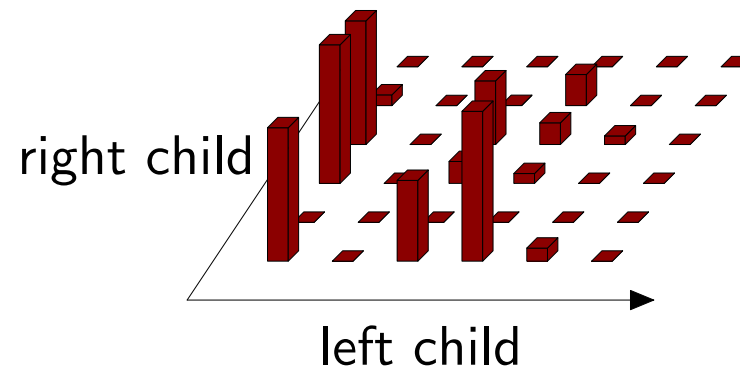


Mean distribution over child symbols:



Distribution over child symbols (per-state):

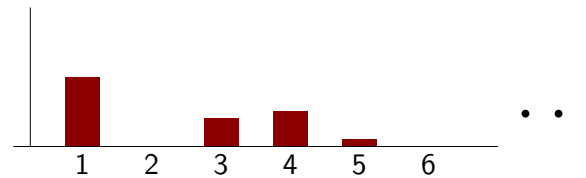
$$\phi_z^B \sim \text{DP}(\alpha^B, \beta\beta^T)$$



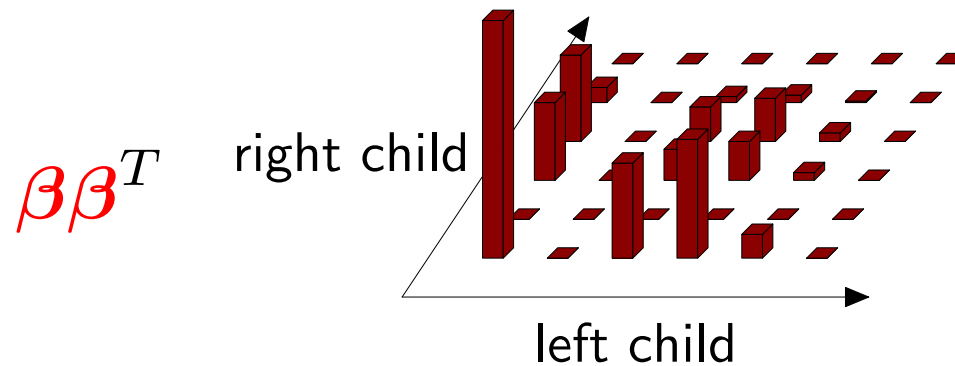
HDP-PCFG: binary productions

Distribution over symbols (top-level):

$$\beta \sim \text{GEM}(\alpha)$$

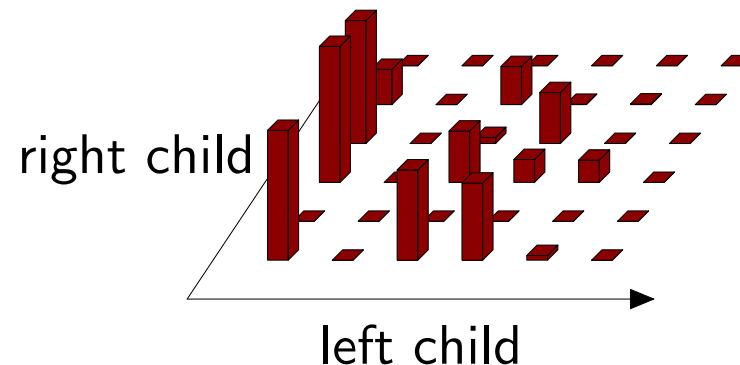


Mean distribution over child symbols:



Distribution over child symbols (per-state):

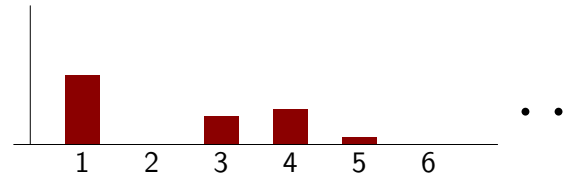
$$\phi_z^B \sim \text{DP}(\alpha^B, \beta\beta^T)$$



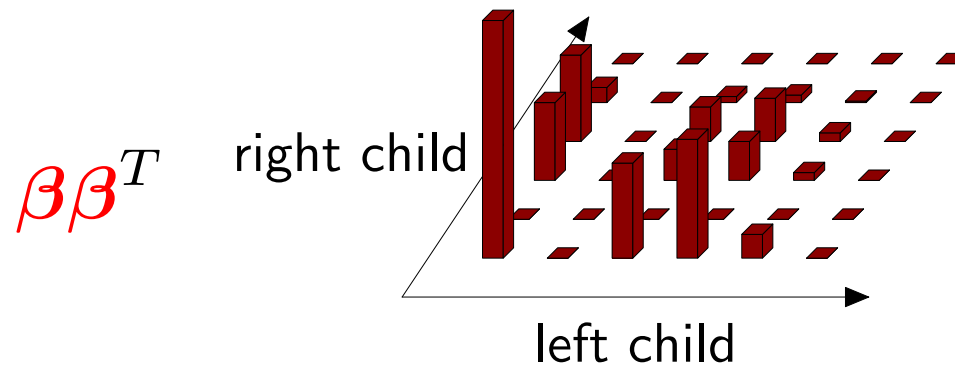
HDP-PCFG: binary productions

Distribution over symbols (top-level):

$$\beta \sim \text{GEM}(\alpha)$$

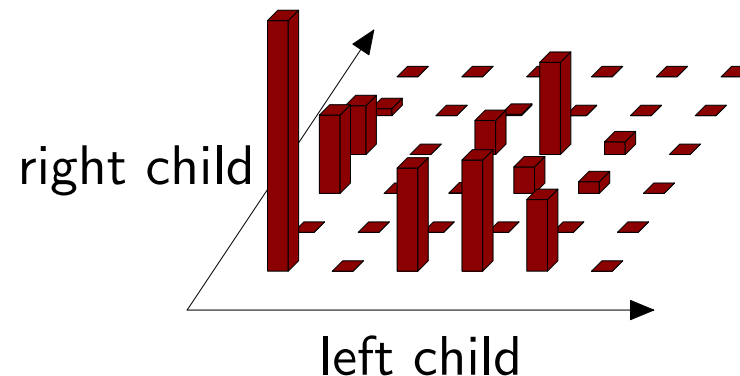


Mean distribution over child symbols:



Distribution over child symbols (per-state):

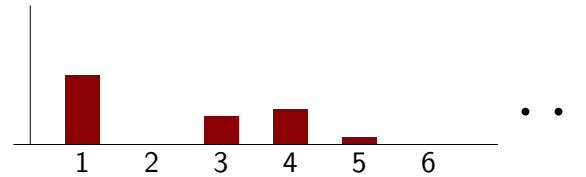
$$\phi_z^B \sim \text{DP}(\alpha^B, \beta\beta^T)$$



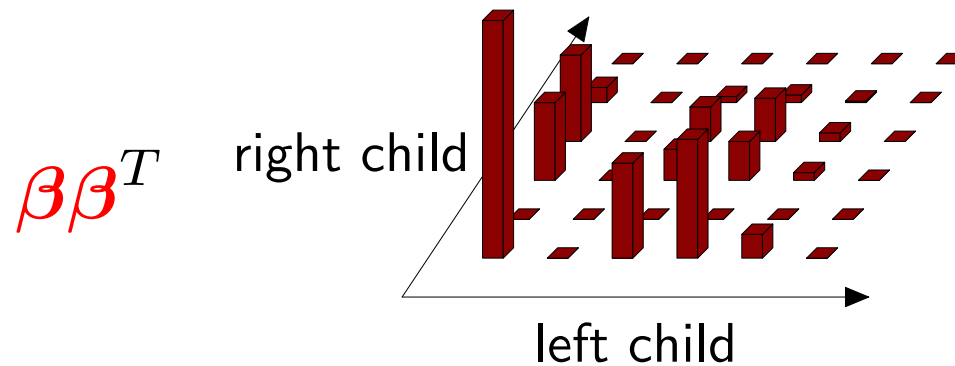
HDP-PCFG: binary productions

Distribution over symbols (top-level):

$$\beta \sim \text{GEM}(\alpha)$$

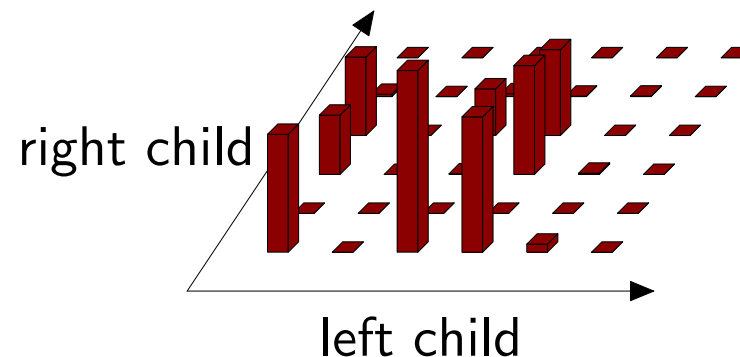


Mean distribution over child symbols:

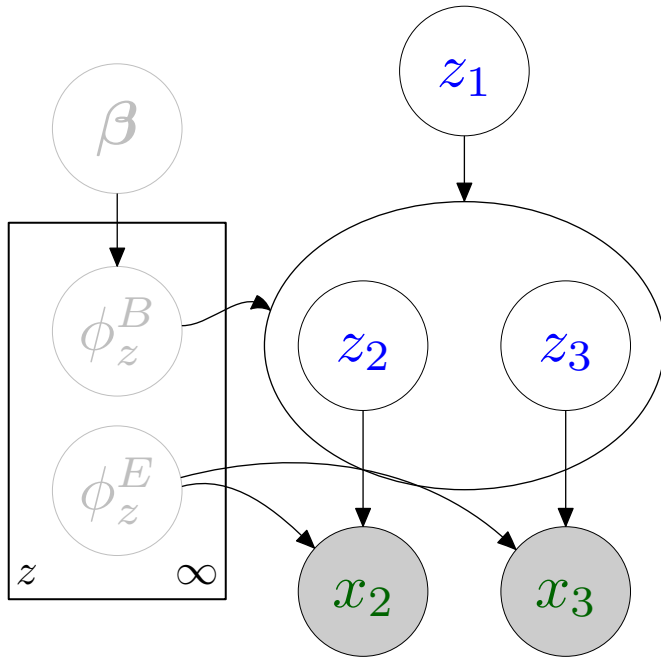


Distribution over child symbols (per-state):

$$\phi_z^B \sim \text{DP}(\alpha^B, \beta\beta^T)$$



HDP probabilistic context-free grammars



HDP-PCFG

$\beta \sim \text{GEM}(\alpha)$ [generate distribution over symbols]

For each symbol $z \in \{1, 2, \dots\}$:

$\phi_z^E \sim \text{Dirichlet}(\alpha^E)$ [generate emission probs]

$\phi_z^B \sim \text{DP}(\alpha^B, \beta\beta^T)$ [binary production probs]

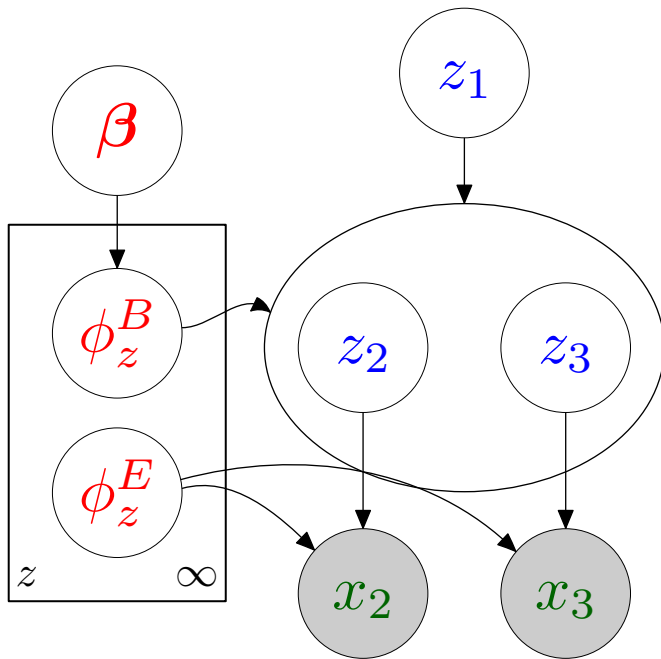
For each nonterminal node i :

$(z_{L(i)}, z_{R(i)}) \sim \text{Multinomial}(\phi_{z_i}^B)$ [child symbols]

For each preterminal node i :

$x_i \sim \text{Multinomial}(\phi_{z_i}^E)$ [terminal symbol]

HDP probabilistic context-free grammars



HDP-PCFG

$\beta \sim \text{GEM}(\alpha)$ [generate distribution over symbols]

For each symbol $z \in \{1, 2, \dots\}$:

$\phi_z^E \sim \text{Dirichlet}(\alpha^E)$ [generate emission probs]

$\phi_z^B \sim \text{DP}(\alpha^B, \beta\beta^T)$ [binary production probs]

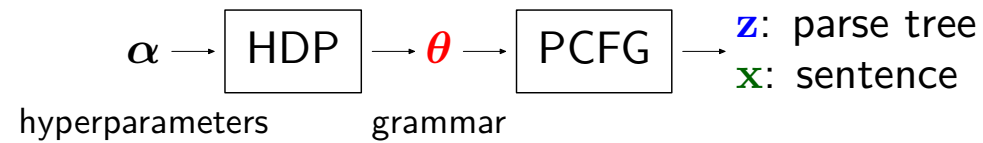
For each nonterminal node i :

$(z_{L(i)}, z_{R(i)}) \sim \text{Multinomial}(\phi_{z_i}^B)$ [child symbols]

For each preterminal node i :

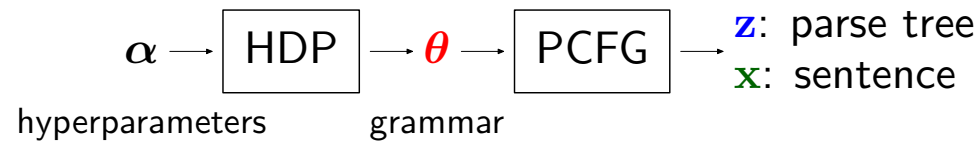
$x_i \sim \text{Multinomial}(\phi_{z_i}^E)$ [terminal symbol]

Variational Bayesian inference



Goal: compute posterior $p(\boldsymbol{\theta}, \mathbf{z} \mid \mathbf{x})$

Variational Bayesian inference

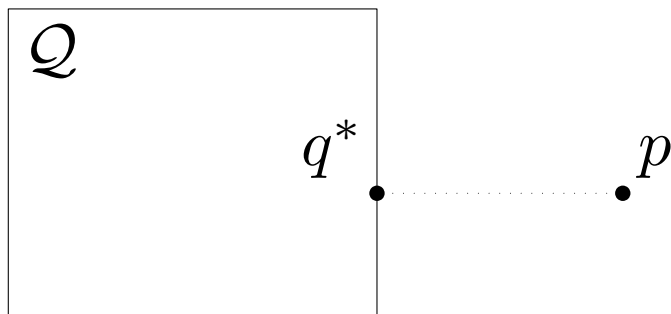


Goal: compute posterior $p(\theta, \mathbf{z} \mid \mathbf{x})$

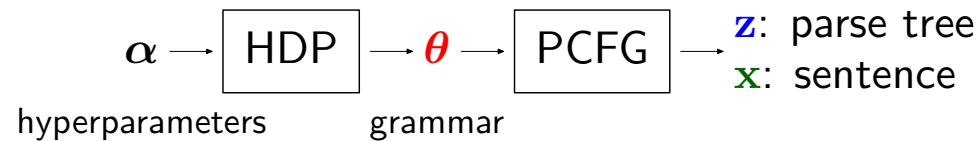
Variational inference:

approximate posterior with
best from a set of tractable
distributions \mathcal{Q} :

$$q^* = \operatorname{argmin}_{q \in \mathcal{Q}} \operatorname{KL}(q \parallel p)$$



Variational Bayesian inference

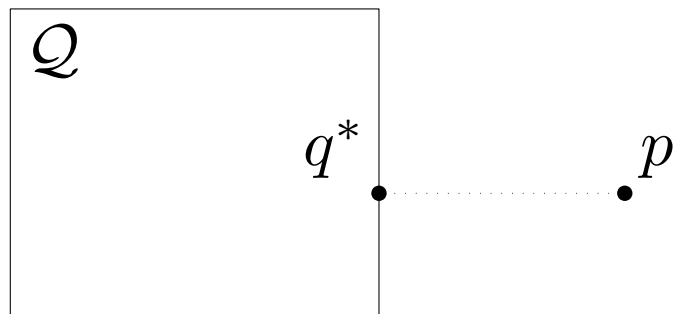


Goal: compute posterior $p(\theta, \mathbf{z} \mid \mathbf{x})$

Variational inference:

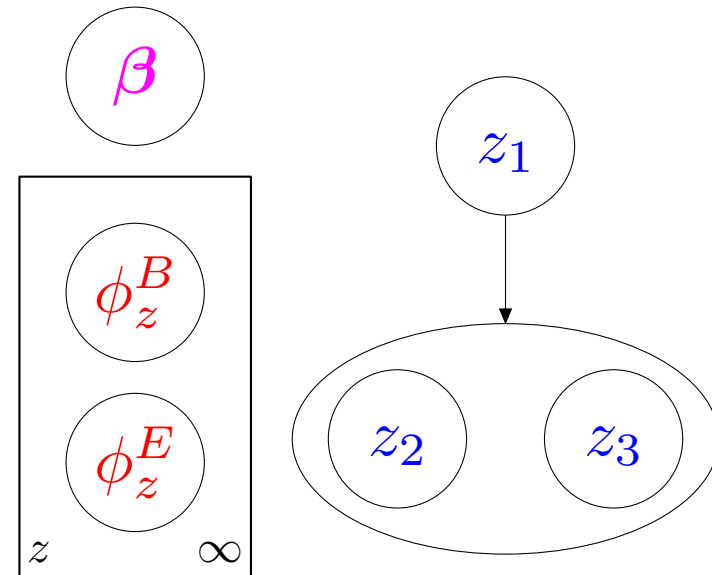
approximate posterior with best from a set of tractable distributions Q :

$$q^* = \operatorname{argmin}_{q \in Q} \text{KL}(q \parallel p)$$



Mean-field approximation:

$$Q = \left\{ q : q = q(\mathbf{z})q(\beta)q(\phi) \right\}$$



Coordinate-wise descent algorithm

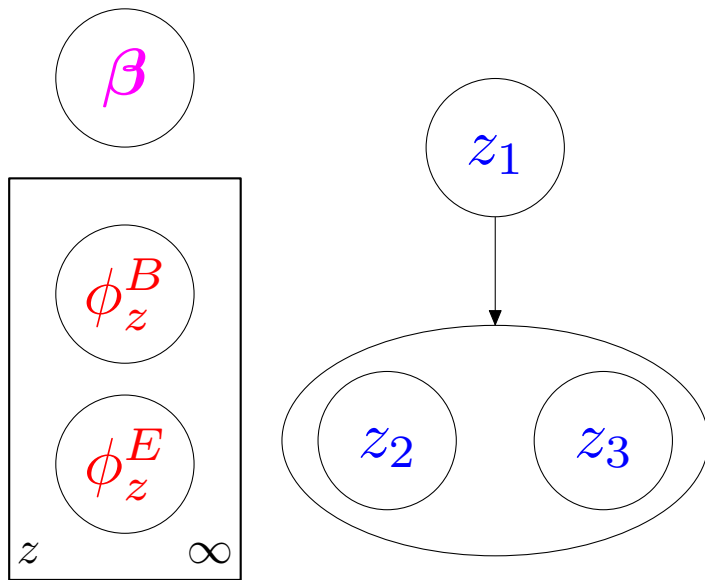
Goal: $\operatorname{argmin}_{q \in \mathcal{Q}} \operatorname{KL}(q || p)$

$$q = q(\mathbf{z})q(\boldsymbol{\phi})q(\boldsymbol{\beta})$$

\mathbf{z} = parse tree

$\boldsymbol{\phi}$ = rule probabilities

$\boldsymbol{\beta}$ = inventory of symbols



Coordinate-wise descent algorithm

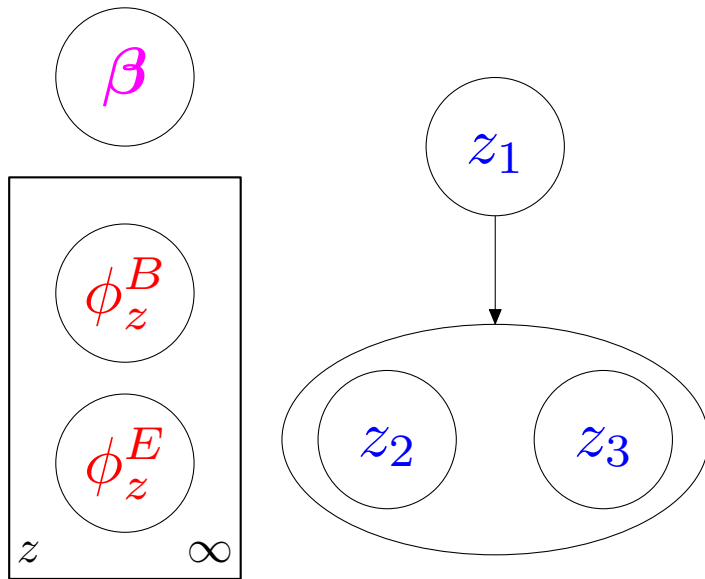
Goal: $\operatorname{argmin}_{q \in \mathcal{Q}} \operatorname{KL}(q || p)$

$$q = q(\mathbf{z})q(\phi)q(\beta)$$

\mathbf{z} = parse tree

ϕ = rule probabilities

β = inventory of symbols



Iterate:

- Optimize $q(\mathbf{z})$ (E-step):

- Optimize $q(\phi)$ (M-step):

- Optimize $q(\beta)$ (no equivalent in EM):

Coordinate-wise descent algorithm

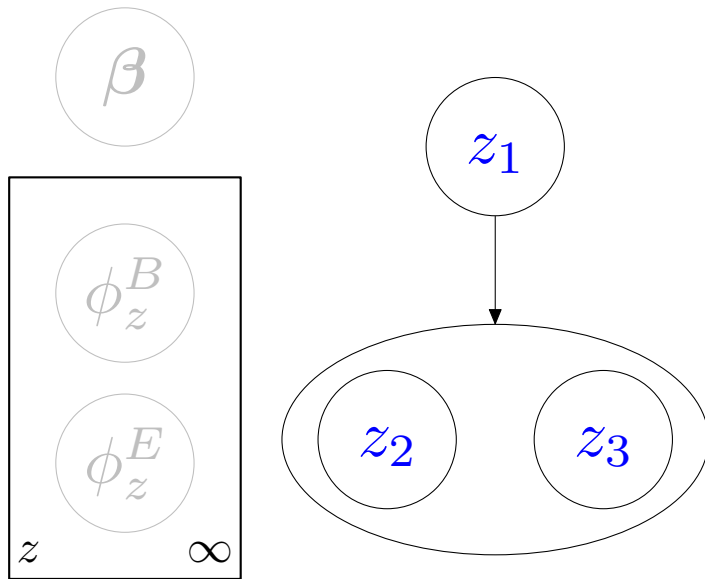
Goal: $\operatorname{argmin}_{q \in \mathcal{Q}} \operatorname{KL}(q || p)$

$$q = q(\mathbf{z})q(\phi)q(\beta)$$

\mathbf{z} = parse tree

ϕ = rule probabilities

β = inventory of symbols



Iterate:

- Optimize $q(\mathbf{z})$ (E-step):
 - Inside-outside with rule weights $W(r)$
 - Gather expected rule counts $C(r)$
- Optimize $q(\phi)$ (M-step):
- Optimize $q(\beta)$ (no equivalent in EM):

Coordinate-wise descent algorithm

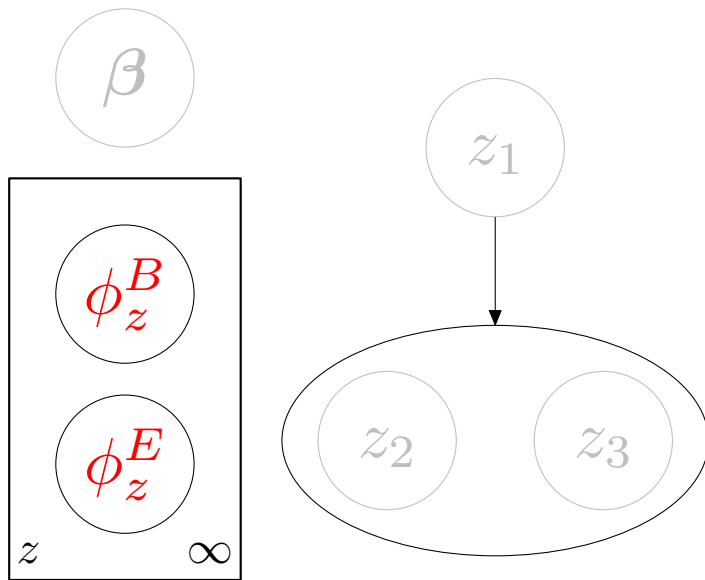
Goal: $\operatorname{argmin}_{q \in \mathcal{Q}} \operatorname{KL}(q || p)$

$$q = q(\mathbf{z})q(\phi)q(\beta)$$

\mathbf{z} = parse tree

ϕ = rule probabilities

β = inventory of symbols



Iterate:

- Optimize $q(\mathbf{z})$ (E-step):
 - Inside-outside with rule weights $W(r)$
 - Gather expected rule counts $C(r)$
- Optimize $q(\phi)$ (M-step):
 - Update Dirichlet posteriors (expected rule counts + pseudocounts)
 - Compute rule weights $W(r)$
- Optimize $q(\beta)$ (no equivalent in EM):

Coordinate-wise descent algorithm

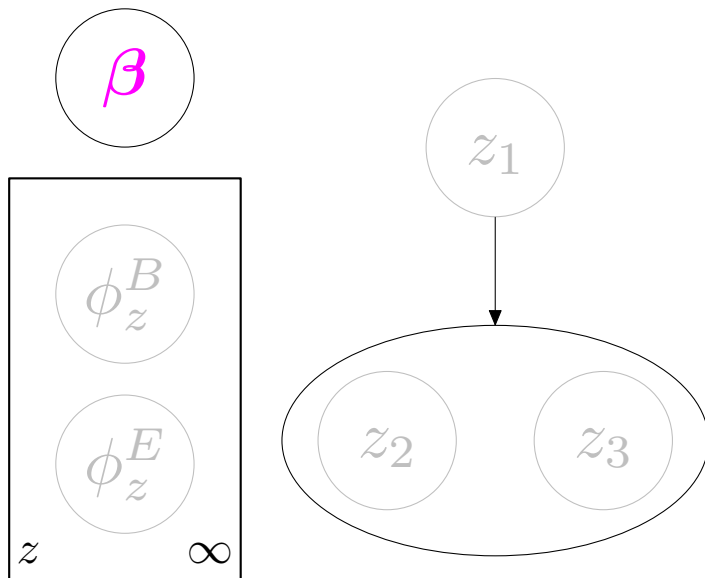
Goal: $\operatorname{argmin}_{q \in \mathcal{Q}} \operatorname{KL}(q || p)$

$$q = q(\mathbf{z})q(\phi)q(\beta)$$

\mathbf{z} = parse tree

ϕ = rule probabilities

β = inventory of symbols



Iterate:

- Optimize $q(\mathbf{z})$ (E-step):
 - Inside-outside with rule weights $W(r)$
 - Gather expected rule counts $C(r)$
- Optimize $q(\phi)$ (M-step):
 - Update Dirichlet posteriors (expected rule counts + pseudocounts)
 - Compute rule weights $W(r)$
- Optimize $q(\beta)$ (no equivalent in EM):
 - Truncate at level K (set the maximum number of symbols)
 - Use projected gradient to adapt number of symbols

Rule weights

- Weight $W(r)$ of rule r similar to probability $p(r)$
- $W(r)$ unnormalized \Rightarrow extra degree of freedom

Rule weights

- Weight $W(r)$ of rule r similar to probability $p(r)$
- $W(r)$ unnormalized \Rightarrow extra degree of freedom

EM (maximum likelihood):

$$W(r) = \frac{C(r)}{\sum_{r'} C(r')}$$

Rule weights

- Weight $W(r)$ of rule r similar to probability $p(r)$
- $W(r)$ unnormalized \Rightarrow extra degree of freedom

EM (maximum likelihood):

$$W(r) = \frac{C(r)}{\sum_{r'} C(r')}$$

EM (maximum a posteriori):

$$W(r) = \frac{\text{prior}(r) - 1 + C(r)}{\sum_{r'} \text{prior}(r') - 1 + C(r')}$$

Rule weights

- Weight $W(r)$ of rule r similar to probability $p(r)$
- $W(r)$ unnormalized \Rightarrow extra degree of freedom

EM (maximum likelihood):

$$W(r) = \frac{C(r)}{\sum_{r'} C(r')}$$

EM (maximum a posteriori):

$$W(r) = \frac{\text{prior}(r) - 1 + C(r)}{\sum_{r'} \text{prior}(r') - 1 + C(r')}$$

Mean-field (with DP prior):

$$W(r) = \frac{\exp \Psi(\text{prior}(r) + C(r))}{\exp \Psi(\sum_{r'} \text{prior}(r') + C(r'))}$$

Rule weights

- Weight $W(r)$ of rule r similar to probability $p(r)$
- $W(r)$ unnormalized \Rightarrow extra degree of freedom

EM (maximum likelihood):

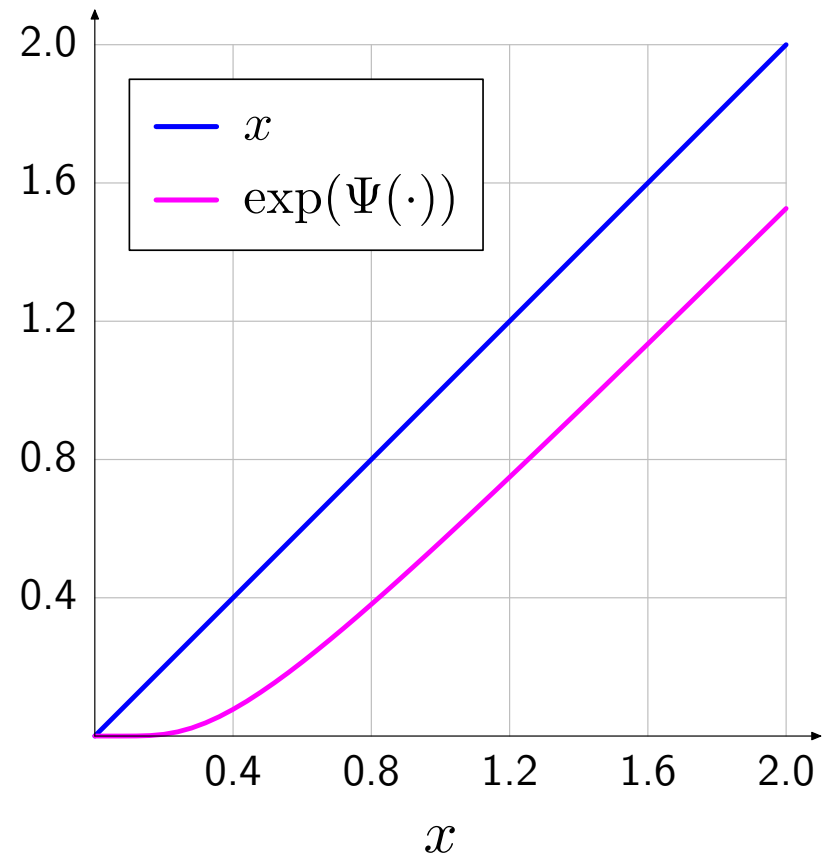
$$W(r) = \frac{C(r)}{\sum_{r'} C(r')}$$

EM (maximum a posteriori):

$$W(r) = \frac{\text{prior}(r) - 1 + C(r)}{\sum_{r'} \text{prior}(r') - 1 + C(r')}$$

Mean-field (with DP prior):

$$W(r) = \frac{\exp \Psi(\text{prior}(r) + C(r))}{\exp \Psi(\sum_{r'} \text{prior}(r') + C(r'))}$$



Rule weights

- Weight $W(r)$ of rule r similar to probability $p(r)$
- $W(r)$ unnormalized \Rightarrow extra degree of freedom

EM (maximum likelihood):

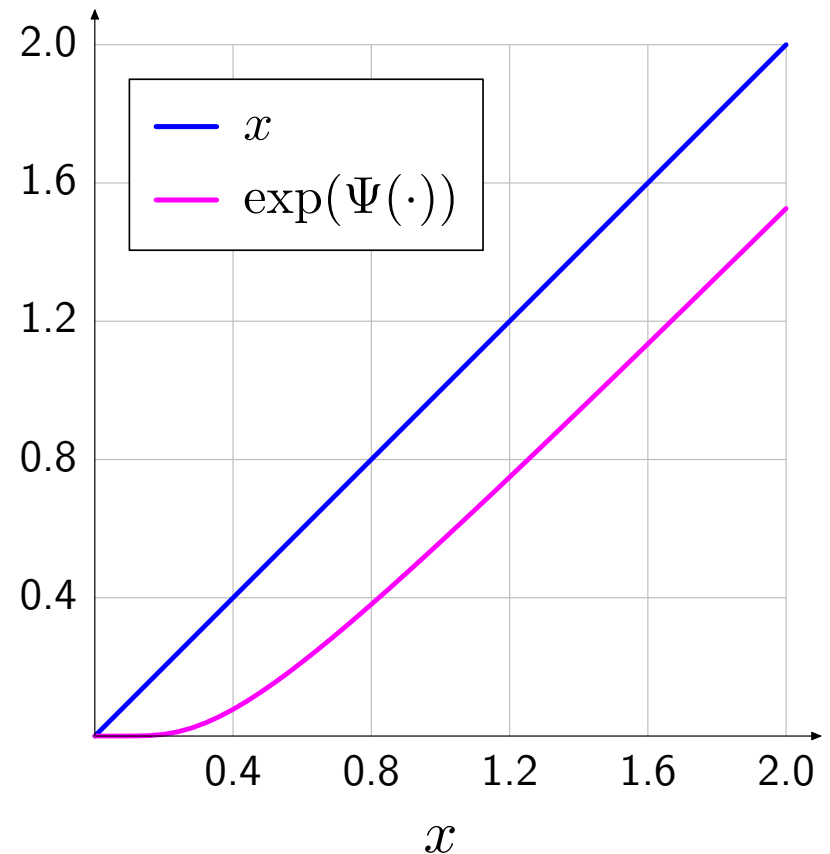
$$W(r) = \frac{C(r)}{\sum_{r'} C(r')}$$

EM (maximum a posteriori):

$$W(r) = \frac{\text{prior}(r) - 1 + C(r)}{\sum_{r'} \text{prior}(r') - 1 + C(r')}$$

Mean-field (with DP prior):

$$W(r) = \frac{\exp \Psi(\text{prior}(r) + C(r))}{\exp \Psi(\sum_{r'} \text{prior}(r') + C(r'))}$$
$$\approx \frac{C(r) - 0.5}{\sum_{r'} C(r') - 0.5}$$



Rule weights

- Weight $W(r)$ of rule r similar to probability $p(r)$
- $W(r)$ unnormalized \Rightarrow extra degree of freedom

EM (maximum likelihood):

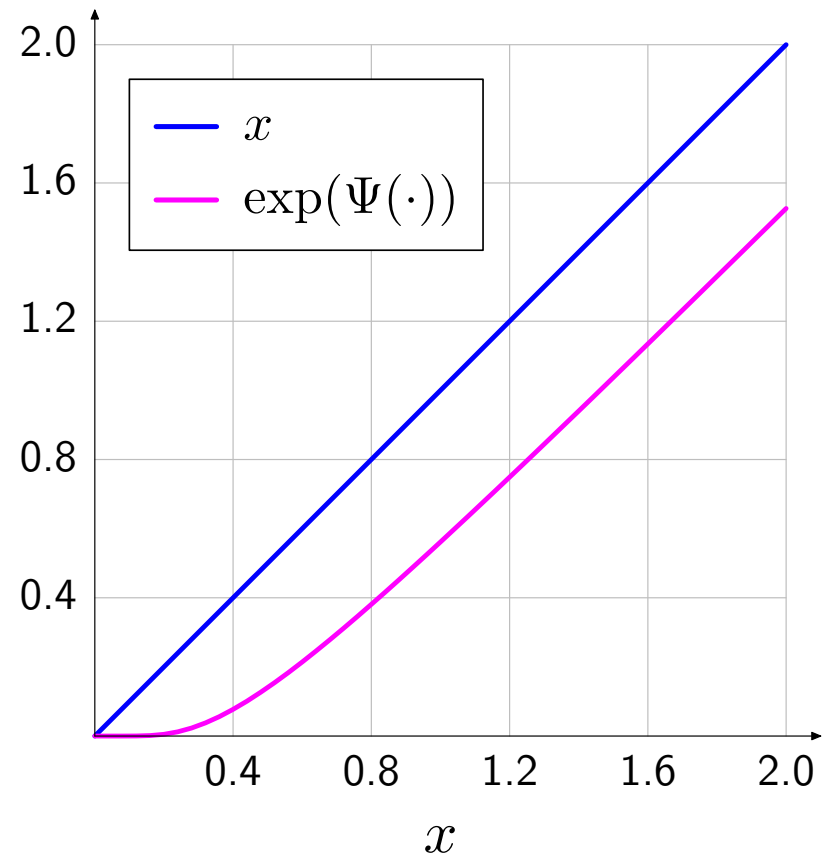
$$W(r) = \frac{C(r)}{\sum_{r'} C(r')}$$

EM (maximum a posteriori):

$$W(r) = \frac{\text{prior}(r) - 1 + C(r)}{\sum_{r'} \text{prior}(r') - 1 + C(r')}$$

Mean-field (with DP prior):

$$W(r) = \frac{\exp \Psi(\text{prior}(r) + C(r))}{\exp \Psi(\sum_{r'} \text{prior}(r') + C(r'))}$$
$$\approx \frac{C(r) - 0.5}{\sum_{r'} C(r') - 0.5}$$

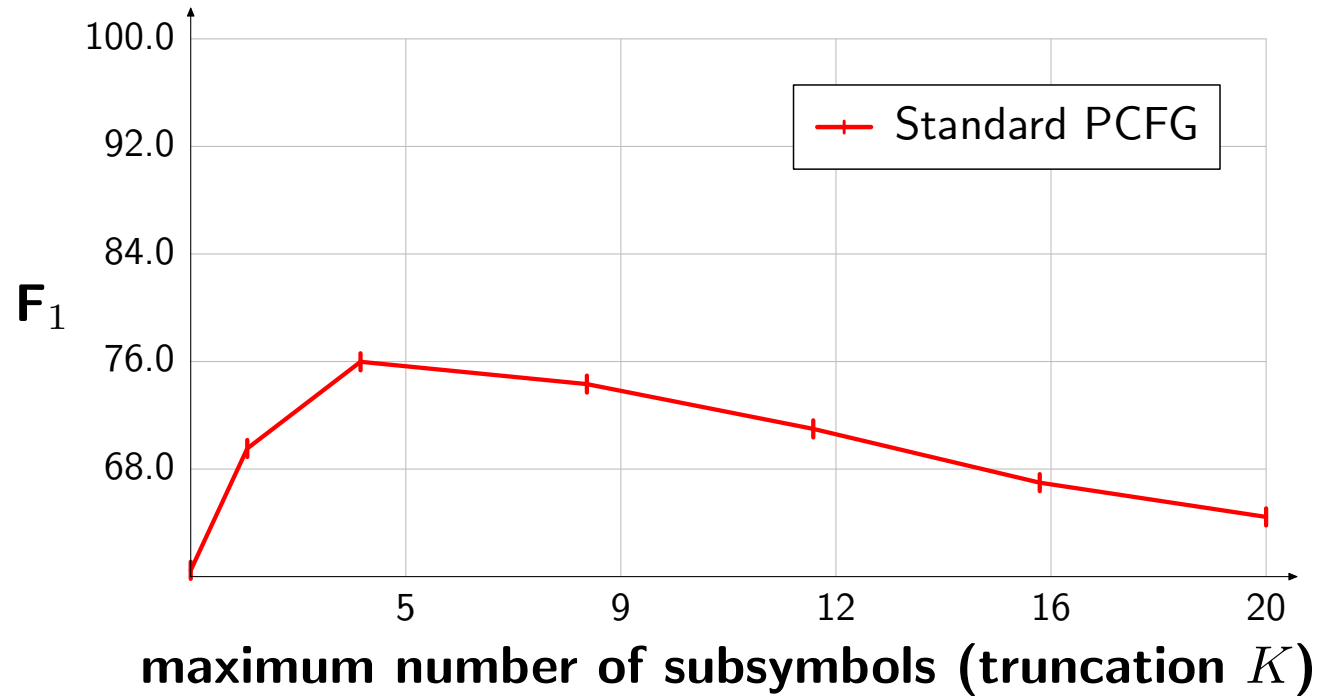


Subtract 0.5 \Rightarrow small counts hurt more than large counts
 \Rightarrow rich gets richer \Rightarrow controls number of symbols

Parsing the WSJ Penn Treebank

Setup: grammar refinement (split symbols into subsymbols)

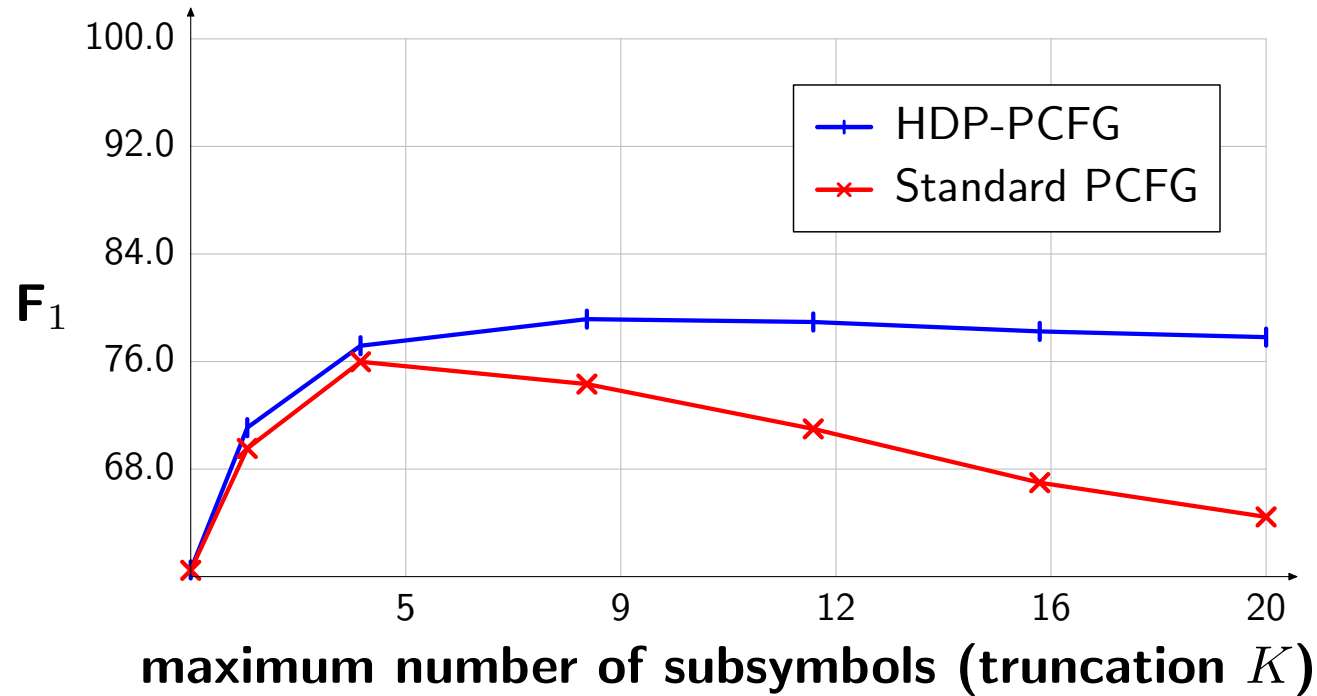
Training on one section:



Parsing the WSJ Penn Treebank

Setup: grammar refinement (split symbols into subsymbols)

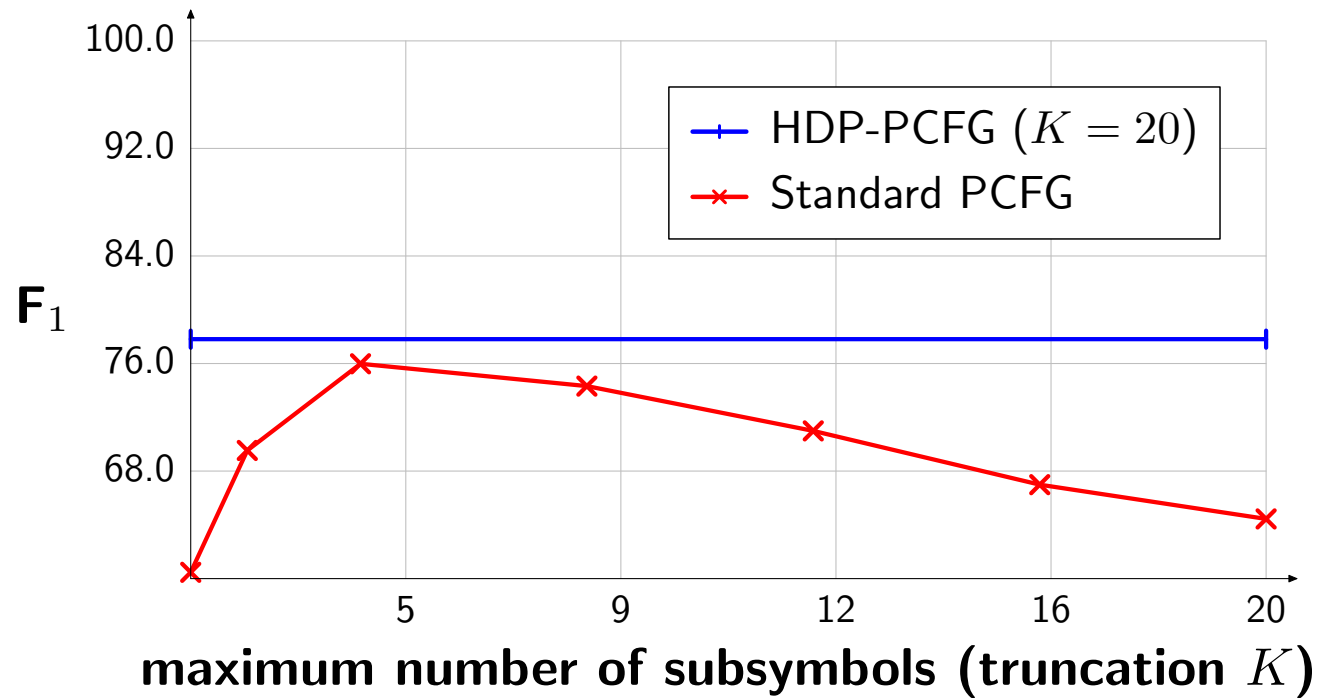
Training on one section:



Parsing the WSJ Penn Treebank

Setup: grammar refinement (split symbols into subsymbols)

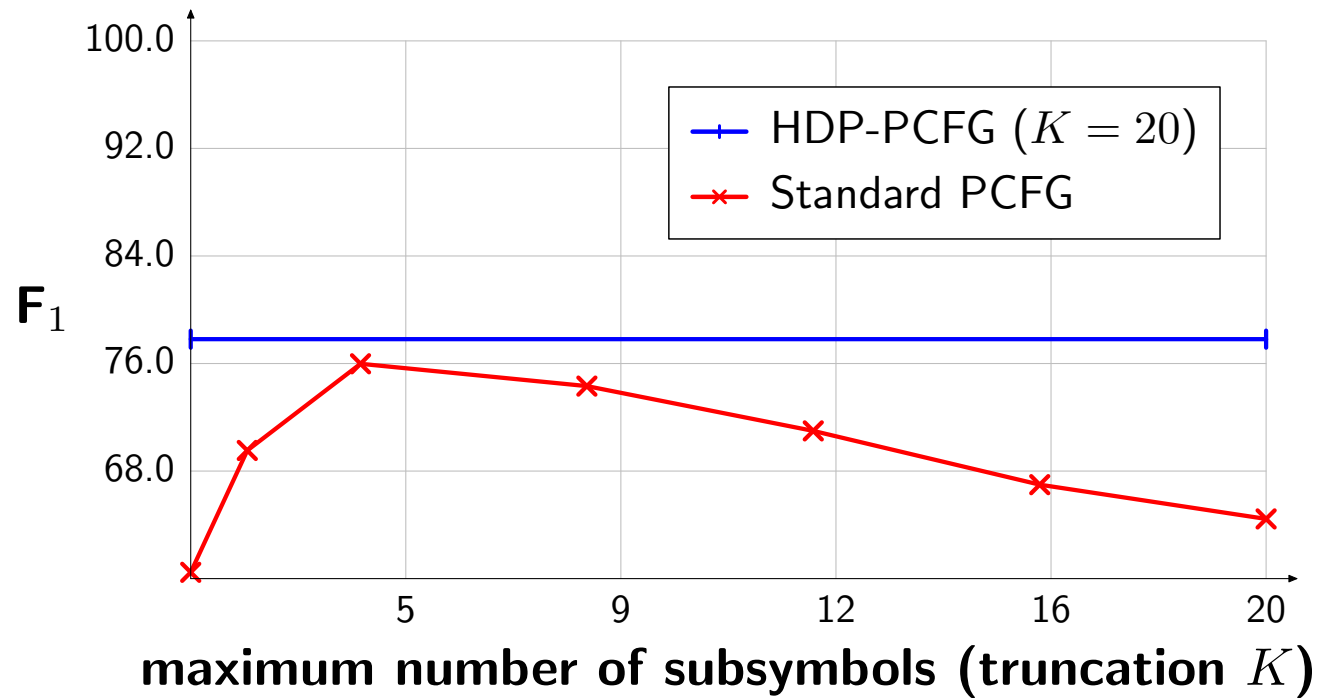
Training on one section:



Parsing the WSJ Penn Treebank

Setup: grammar refinement (split symbols into subsymbols)

Training on one section:



Training on 20 sections: ($K = 16$)

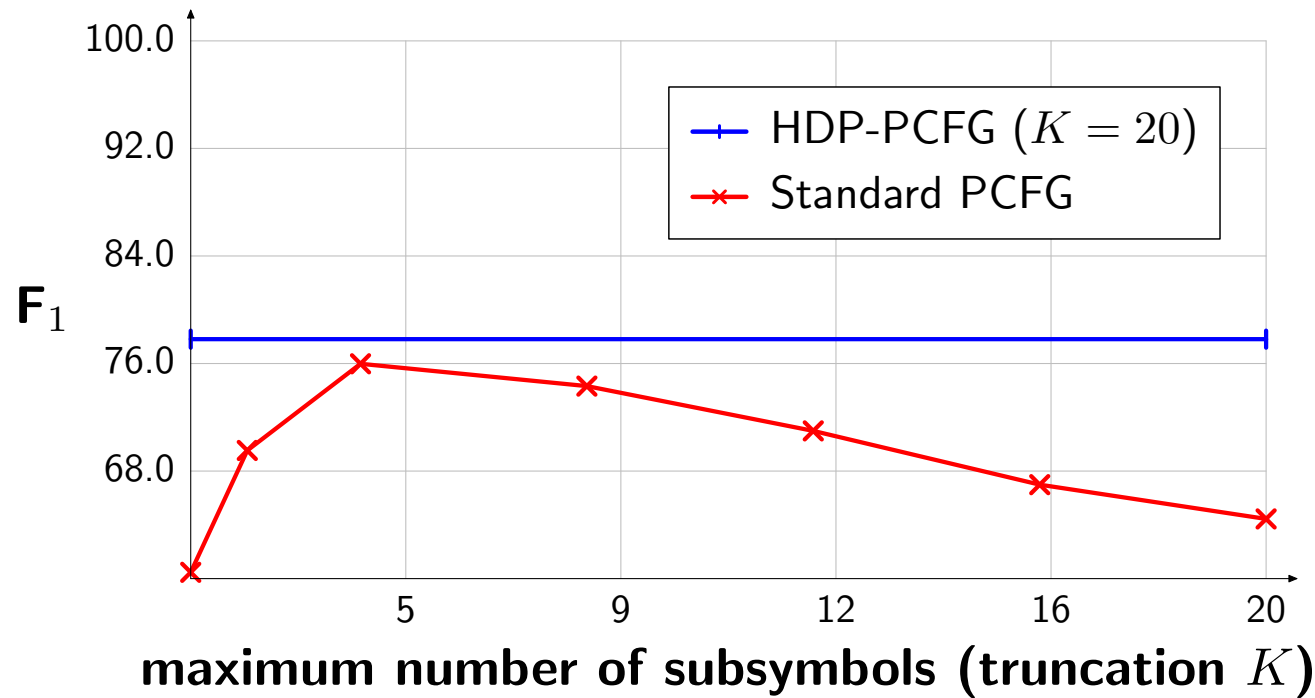
Standard PCFG: 86.23

HDP-PCFG: 87.08

Parsing the WSJ Penn Treebank

Setup: grammar refinement (split symbols into subsymbols)

Training on one section:



Training on 20 sections: ($K = 16$)

Standard PCFG: 86.23

HDP-PCFG: 87.08

Results:

- HDP-PCFG overfits less than standard PCFG
- If have large amounts of data, HDP-PCFG \cong standard PCFG

