

領域効率のよい文法圧縮
のための可変長符号
(SPIRE 2012 発表予定)

高畠嘉将(九工大)

田部井靖生(JST ERATO)

坂本比呂志(九工大)

CFG(文脈自由文法)による圧縮

$w = \{bbacbbacbbacbbac\}$

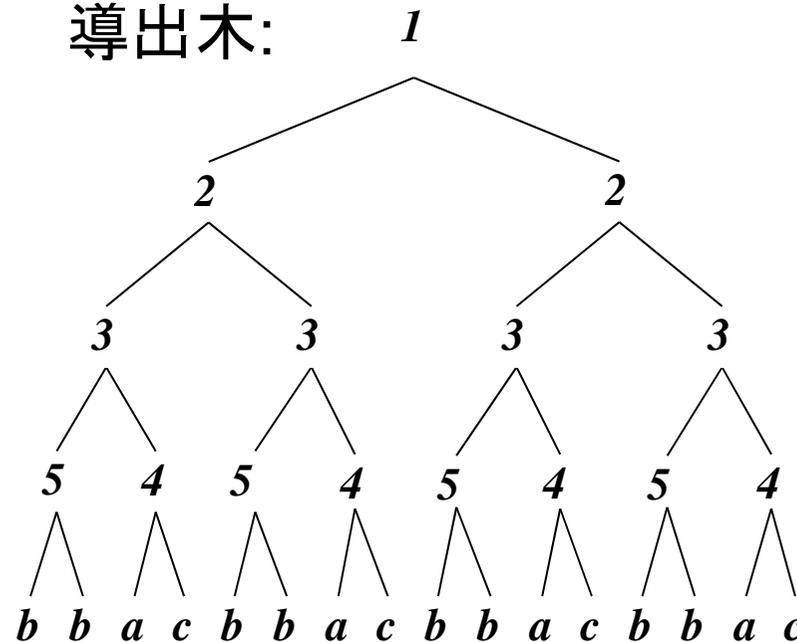
↓ CFG

辞書

Z	X	Y
1	2	2
2	3	3
3	5	4
4	a	c
5	b	b

=

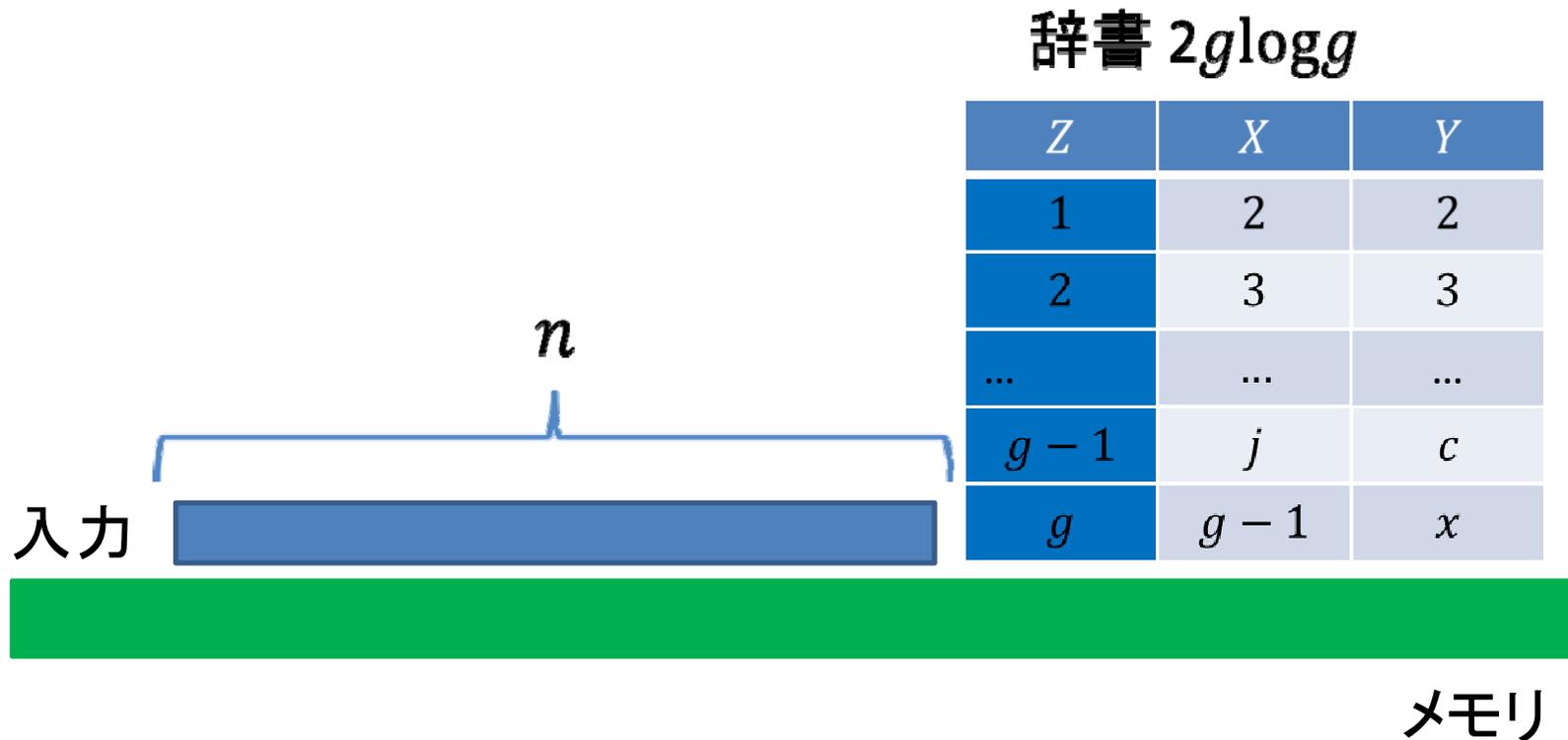
導出木:



最適な辞書を構築するのは、**NP-困難** => 近似アルゴリズム

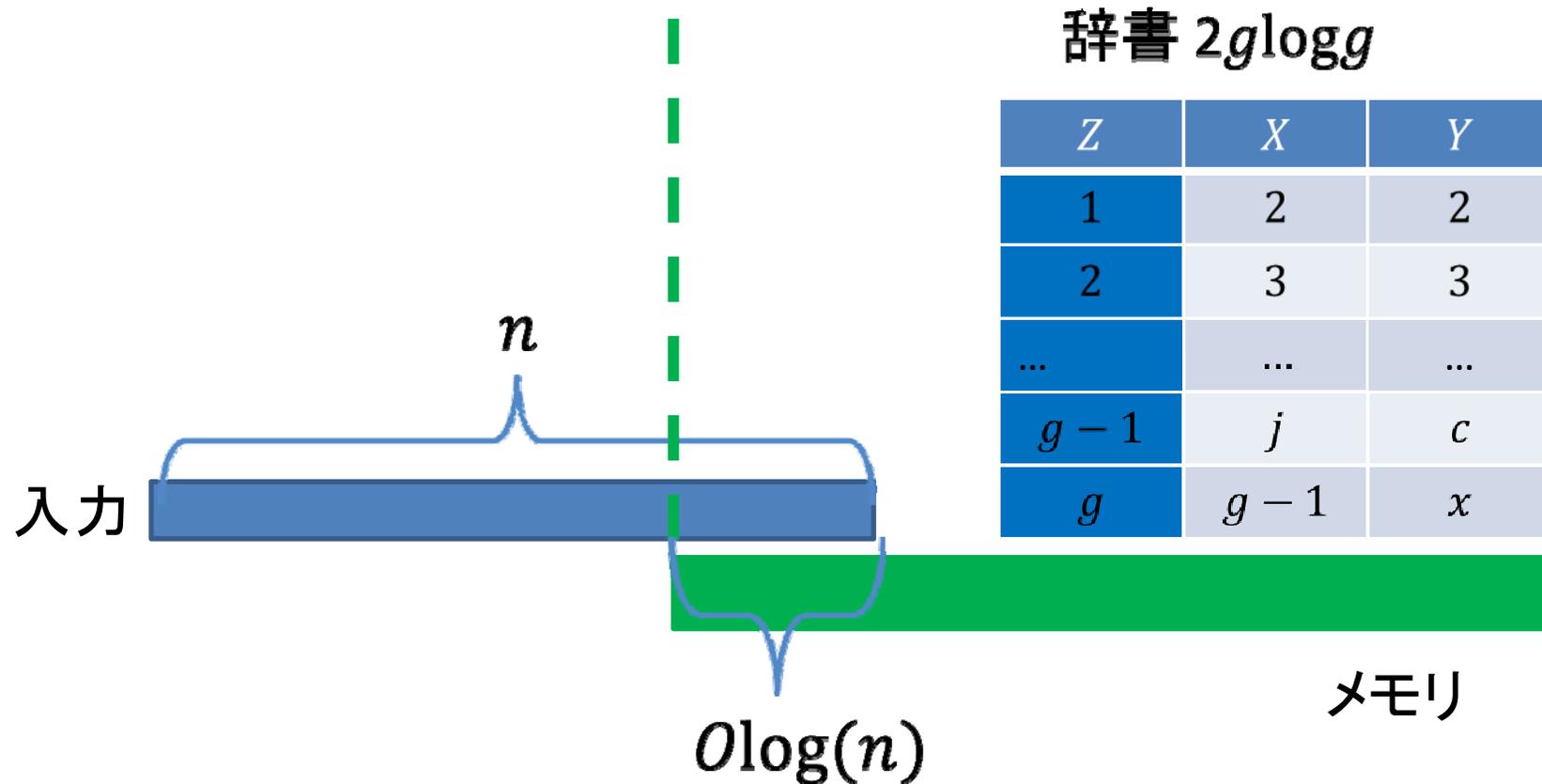
最近の研究

ESP-index(offline)[Maruyama,Kishiue,Nakahara,Sakamoto et al.,SPIRE 2011]

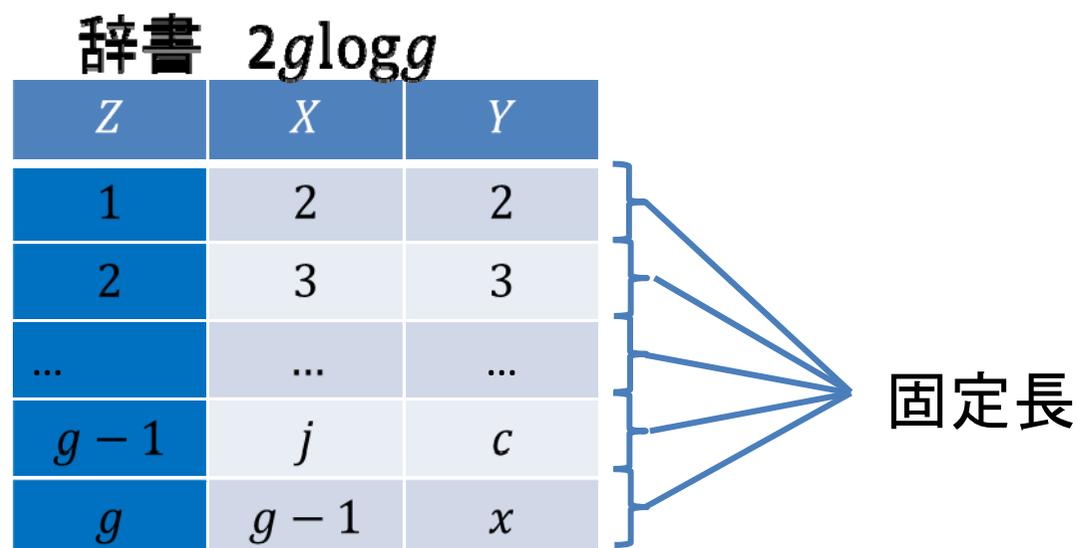


最近の研究

LCA-online [Maruyama, Takeda, Sakamoto et al., Algorithms 2012]



辞書のサイズの問題

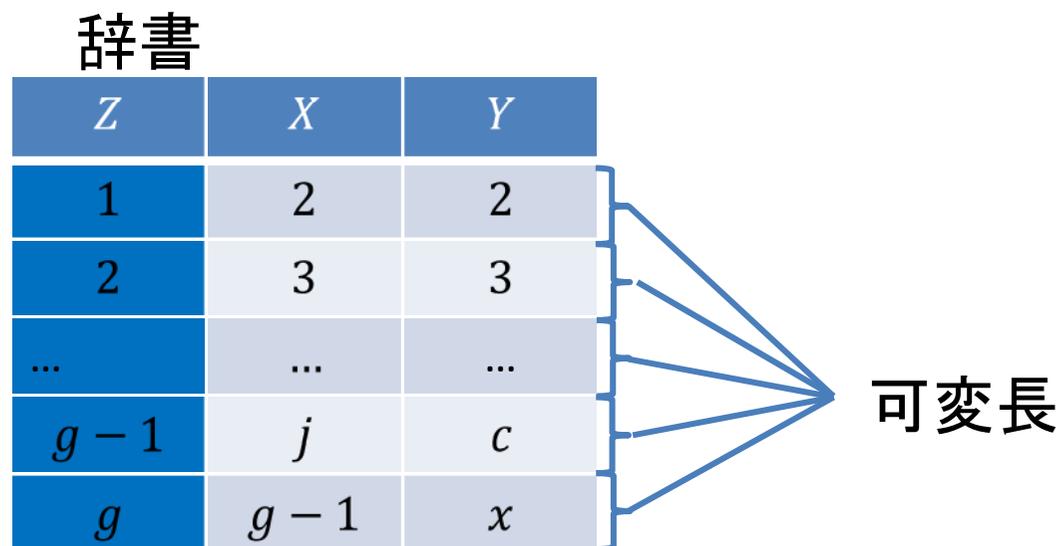


入力によっては辞書のサイズが大きい

メモリ使用量
wikipedia

入力	辞書
5533MB	9401MB

本研究



- 動的に可変長で構築
=>メモリの省スペース化を実現

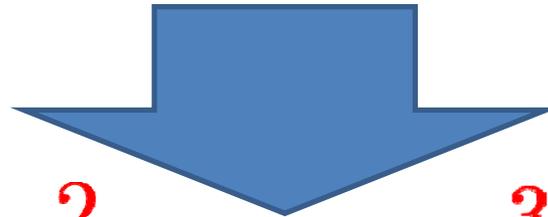
$$2g \log g \rightarrow \frac{3}{2} g \log g + \frac{15}{2} g$$

メモリ使用量	入力	辞書(固定長)	辞書(可変長)
wikipedia	5533MB	9401MB	4776MB

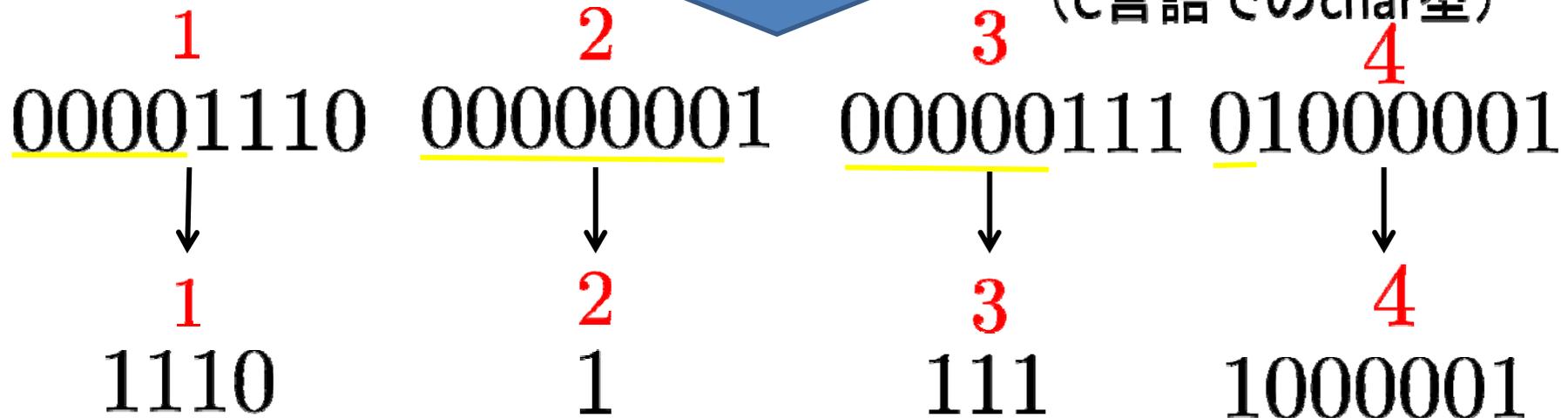
データを可変長で保存

辞書

Z	1	2
X	1:14	3:7
Y	2:1	4:33



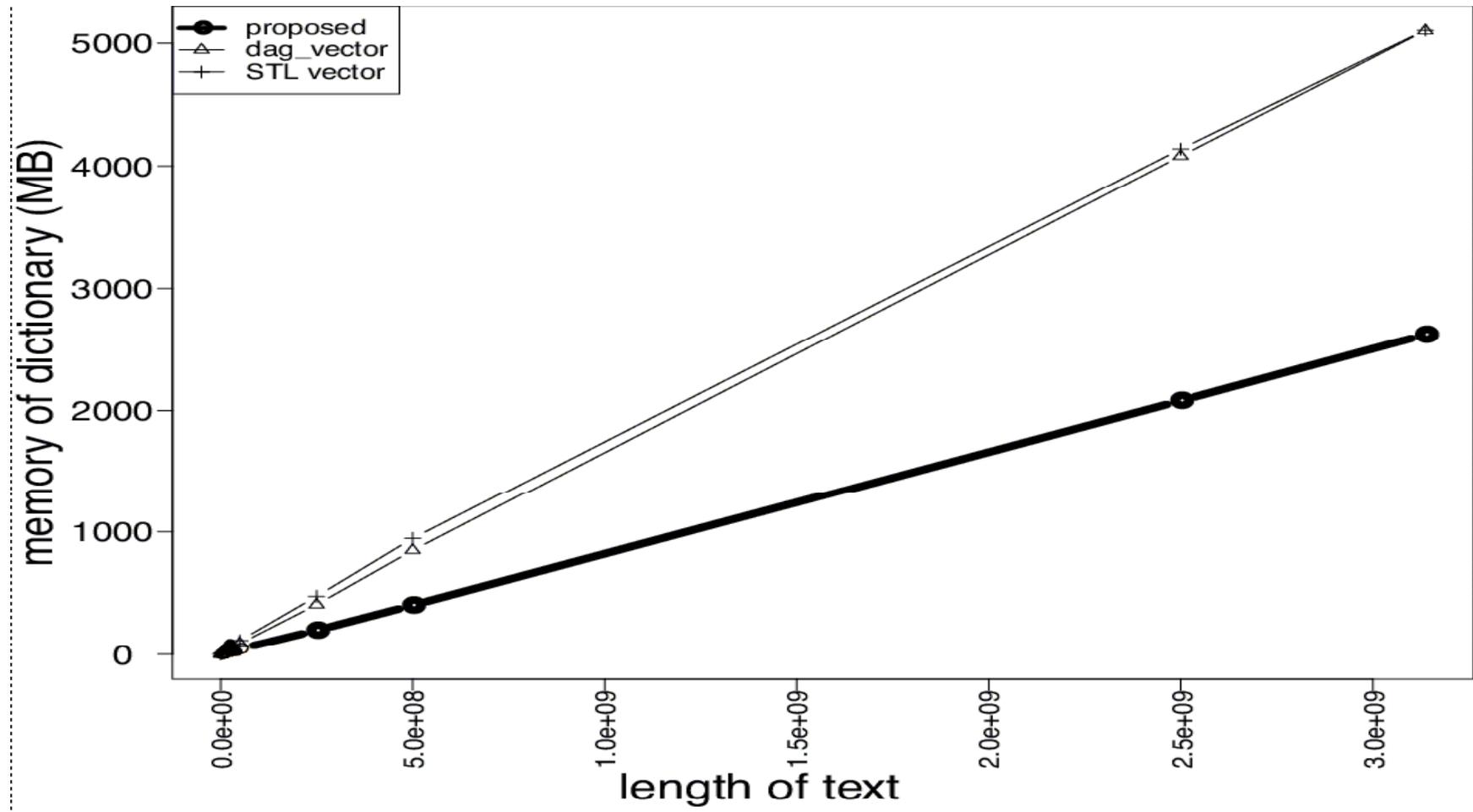
1 byte表現
(C言語でのchar型)



実験環境

- Linux machine
- CPU:8-Core Intel(R) Xeon(R)CPU E7-8837
2.67GHz
- Memory:1TB
- 比較対象
- STL vector: C++のスタンダードテンプレートクラスのvector
- Dag vector:簡潔データ構造を用いたもの

実験(genome)



まとめ

- $\frac{3}{2}g \log g + \frac{15}{2}g$ のメモリスペースで辞書を構築
- 動的に可変長で構築可能