

General Algorithms for Mining Closed Flexible Patterns

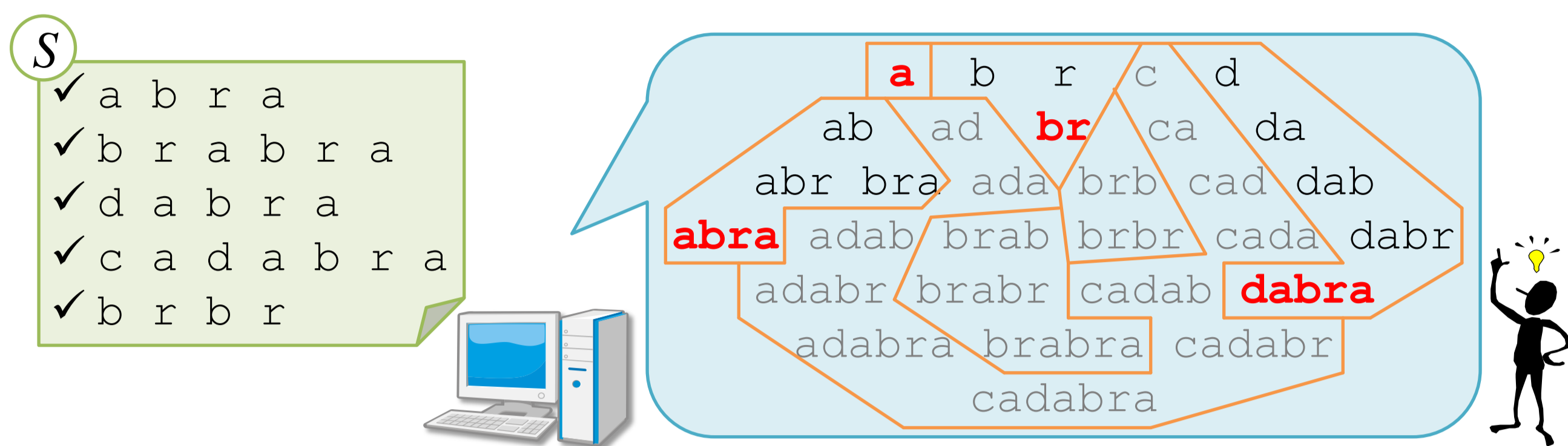
under Various Equivalence Relations

○井 智弘, 荏隈 勇樹, 坂内 英夫, 竹田 正幸 九州大学

頻出飽和文字列パターン列挙問題

入力: テキスト集合 S , 頻度の閾値 σ , 同値関係 \equiv_S
出力: S 中の σ 個以上の文字列に出現する飽和パターンを列挙

例 $\sigma = 2$ のときの頻出飽和部分文字列パターン列挙



同値関係を導入するメリット

1. 簡潔で見やすい出力
2. うまく計算すると探索・出力コスト減

Flexible パターン

$*_1 w_1 *_2 w_2 *_3 \dots *_k w_k *_k$ $k \geq 1$ $* \notin \Sigma$: ギャップ文字
 $w_i \in \Sigma^+$: i 番目のセグメント

パターン P $* a b * e * r a *$

埋め込み $\theta = (cd, r, r, bra)$ を適用

$P\theta$ $c d a b r e r r a b r a$

埋め込み $\theta = (c*, r, e, b*a)$ を適用

$P\theta$ $c * a b r e r a b * a$

パターンの出現区間, 極小出現区間

パターン P $* a b * e * r a *$
テキスト T $y a b x a b e x r a x r a y$
 $\theta = (x, xab, x, xrax)$
 $\theta = (x, xab, xrax, x)$
 $\theta = (xabx, e, x, xrax)$
 $\theta = (xabx, e, xrax, x)$

出現区間 $(P, T) = \{[2, 10], [2, 13], [5, 10], [5, 13]\}$
極小出現区間 $(P, T) = \{[5, 10]\}$

パターン上の半順序

パターン $P \neq Q$ に対して,
 $P\theta = Q$ なる埋め込みが存在するとき
 Q は P より具体的であるという

パターン P $* a b * e * r a *$

埋め込み $\theta = (*, x*b, x, xr*)$ を適用

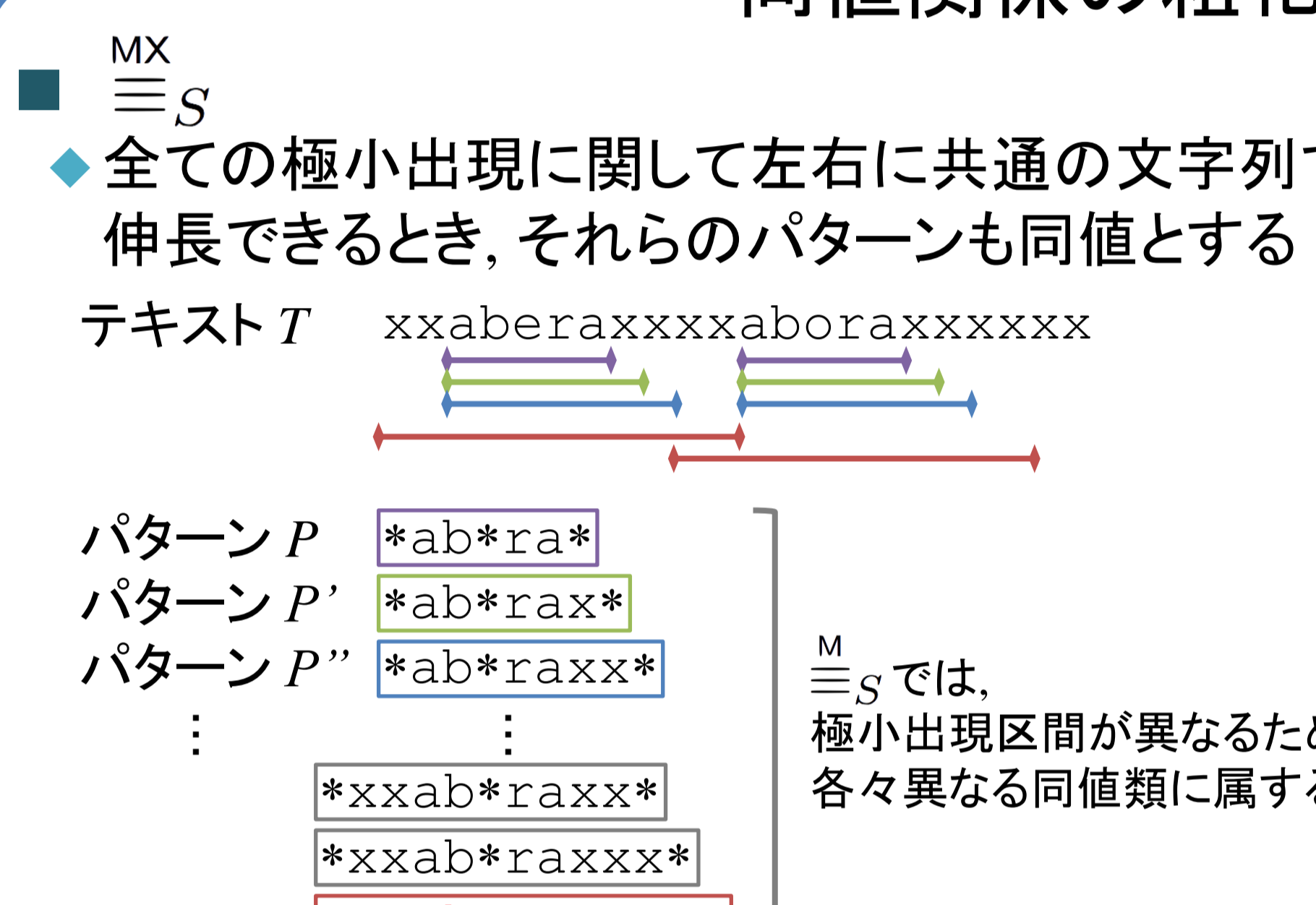
パターン Q $* a b x * b e x r a x r *$

Flexible パターン上の様々な同値関係に対して統一的な多項式時間遅延アルゴリズムを提案

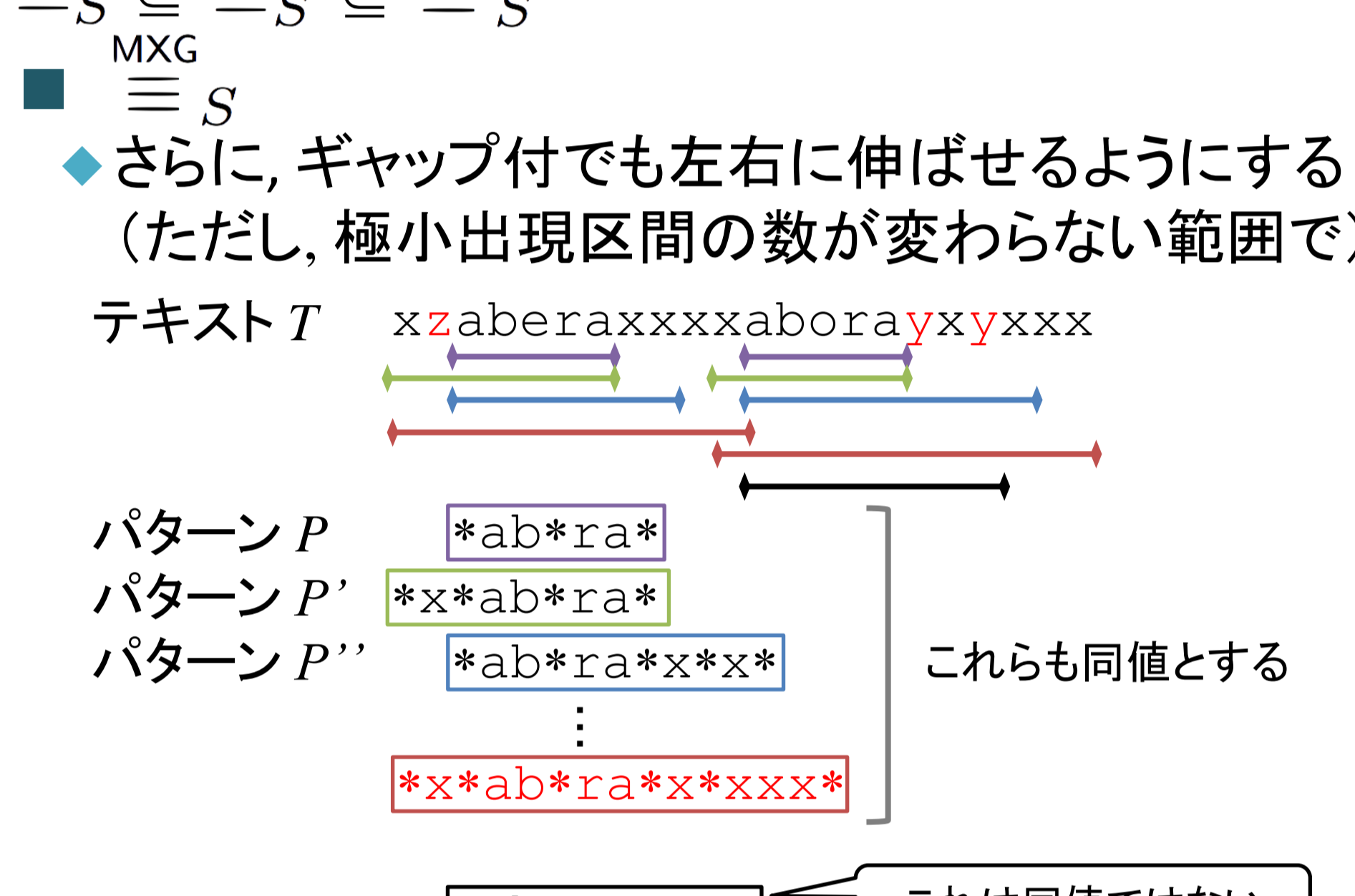
同値関係	説明	遅延計算量
$\overset{I}{\equiv}_S$	出現区間の一致	$O(\Sigma \ S\ d)$
$\overset{IX}{\equiv}_S$	$\overset{I}{\equiv}_S$ の拡張	$O(\Sigma \ S\ d)$
$\overset{M}{\equiv}_S$	極小出現区間の一致	$O(\Sigma \ S\ d)$
$\overset{MX}{\equiv}_S$	$\overset{M}{\equiv}_S$ の拡張1	$O(\Sigma \ S\ d)$
$\overset{MXG}{\equiv}_S$	$\overset{M}{\equiv}_S$ の拡張2	$O(\Sigma \ S\ ^2 d)$
$\overset{E}{\equiv}_S$	出現の終了位置集合の一致	$O(\Sigma \ S\)$
$\overset{EX}{\equiv}_S$	$\overset{E}{\equiv}_S$ の拡張	$O(\Sigma \ S\)$
$\overset{B}{\equiv}_S$	出現の開始位置集合の一致 (有村&宇野)	$O(\Sigma \ S\)$

$\|S\|$: テキスト長の総和
 d : セグメント数の最大値
計算領域はいずれも $O(\|S\|d)$

同値関係の粗化

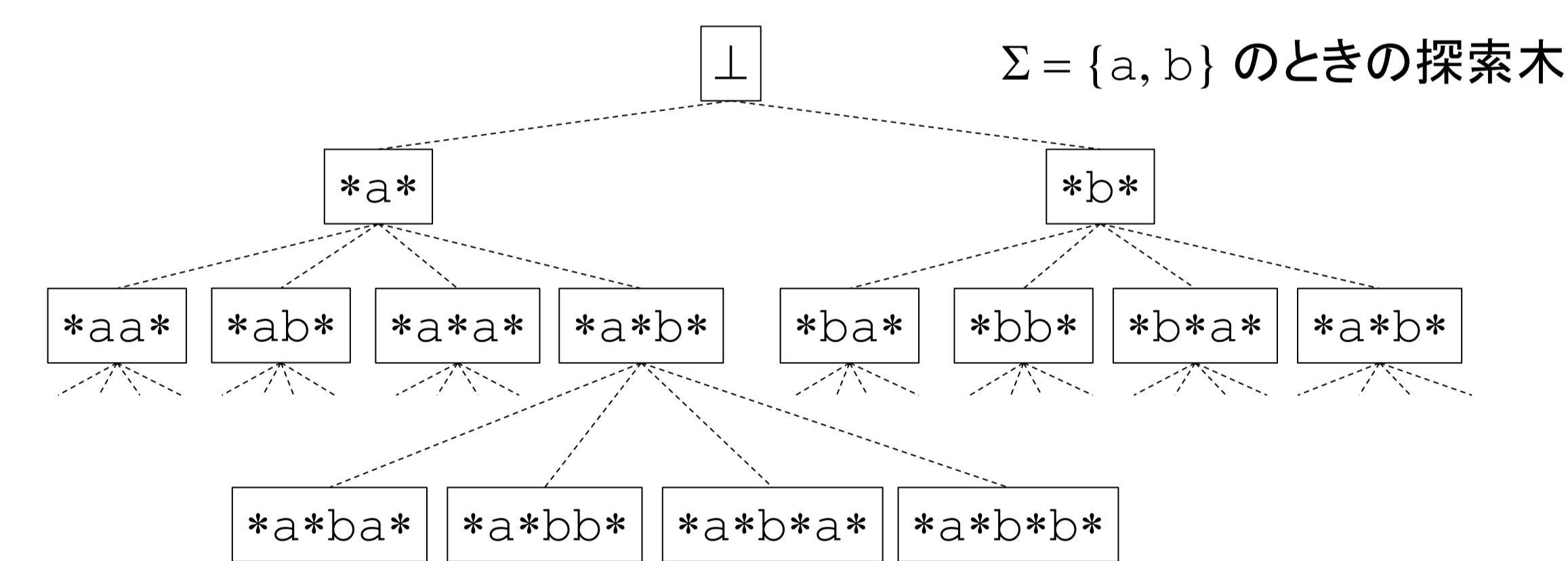


同値関係の粗化 (続き)



次の親子関係で定義される探索木を深さ優先探索する

- 任意の $x \in \Sigma$ に対して, $*x*$ の親は \perp
- 任意の $x \in \Sigma, w_1, \dots, w_k \in \Sigma^+$ に対して,
 $*w_1 \dots w_k x*$ と $*w_1 \dots w_k *$ の親は $*w_1 \dots w_k *$



任意のパターン P に対して以下が成り立つ

- 頻度チェック**
 P が頻出でないなら, P の子孫も頻出でない (枝刈り)
- 左端飽和チェック**
左端のギャップ文字を具体的にして P と同値なパターンがあれば P と P の子孫は飽和パターンではない (枝刈り)
- 内部飽和チェック**
左右端以外のギャップ文字を具体的にして P と同値なパターンがあれば P と P の子孫は飽和パターンではない (枝刈り)
- 右端飽和チェック**
右端のギャップ文字を具体的にして P と同値なパターンがあるとき P は飽和パターンではないが P の子孫に頻度が等しく飽和なパターンが必ず存在する

和歌データを用いた実験

algorithm/equiv	$\sigma = 2$ (0.13%)		$\sigma = 4$ (0.26%)		$\sigma = 8$ (0.52%)	
	patterns	seconds	patterns	seconds	patterns	seconds
GenCloFlex $\overset{I}{\equiv}_S$	85,665,856	4084.4	24,941,537	2088.3	4,634,071	823.54
GenCloFlex+ $\overset{I}{\equiv}_S$	85,665,856	3983.7	24,941,537	2065.8	4,634,071	817.12
GenCloFlex+ $\overset{IX}{\equiv}_S$	80,622,811	3932.2	24,661,456	2063.0	4,622,063	816.77
GenCloFlex $\overset{M}{\equiv}_S$	83,042,793	3993.8	24,596,215	2064.2	4,601,106	819.16
GenCloFlex+ $\overset{M}{\equiv}_S$	83,042,793	3872.2	24,596,215	2040.3	4,601,106	810.88
GenCloFlex+ $\overset{MX}{\equiv}_S$	77,392,851	3819.8	24,287,323	2035.1	4,588,408	811.97
GenCloFlex+ $\overset{MXG}{\equiv}_S$	55,131,507	3509.8	23,508,993	2032.7	4,580,550	814.46
GenCloFlex $\overset{E}{\equiv}_S$	67,210,837	3567.4	24,253,028	2055.3	4,609,610	825.44
GenCloFlex+ $\overset{E}{\equiv}_S$	67,210,837	3500.8	24,253,028	2034.5	4,609,610	817.87
GenCloFlex+ $\overset{EX}{\equiv}_S$	65,242,280	3497.0	24,064,249	2030.3	4,601,551	806.16

山家集 (1552 首) に対して, 頻度の閾値 σ を 2, 4, 8 に設定し, 各同値関係における頻出飽和パターンの数と計算時間を比較 (画面表示はオフ)

「御裳濯和歌集序文」と「八代集の各首」の間で, $\overset{M}{\equiv}_S, \overset{MX}{\equiv}_S, \overset{MXG}{\equiv}_S$ における共通飽和パターンを列挙し結果を比較

御裳濯和歌集序文
やまとうたはあしはらのなかつくによりおこり, いなためすかのまよりひるまれり, あとをたれたまへるしんみやうもこのみちをむねとし, あらははれてたまへるふつものことわざをすてたまはず, おほよそわかくにのふうそくとしてむかしよりまいたにゆることなし, つらつらよせんをたつぬれはすてにははたいになり, ところどころのうらきしふそのかすまたおほし, かねにしたふころすすみて, (1241 文字)

- 最初のセグメントから最後のセグメントの距離が 38 の出現のみを考える
- セグメントの長さの総和が 15 以上のパターンのみ表示する