#### 2024/7/24 文字列ゼミ

# CDAWGによる LZ78部分文字列圧縮

柴田 紘希(九州大学)

#### 本研究の概要

研究でやったこと

CDAWGを用いて、LZ78分解の部分文字列圧縮を

省領域で高速に行う手法を提案した

- CDAWG: 文字列の省領域索引
- LZ78分解: 文字列の圧縮表現のひとつ
- 部分文字列圧縮:事前にテキスト文字列が与えられ、

<u>部分文字列の圧縮結果</u>を返すクエリを高速に処理する問題

準備: LZ78分解と部分文字列圧縮問題

- 文字列の圧縮手法の一つ
- 文字列 T を特定のアルゴリズムにしたがって

$$T = F_0 F_1 \dots F_f$$
 に分解する(ただし、 $F_0$  は空文字列)

$$T = abbabaaab = F_0 F_1 \dots F_f$$

$$F = (\varepsilon, a, b, ba, baa, ab)$$

- $F_i$  ( $i \ge 1$ ) を計算するアルゴリズムは以下の通り:
  - $T' = T[1 + |F_0| + |F_1| + \dots + |F_{i-1}| \dots n]$  とし、T' の接頭辞となる最長の $F_j$  (j < i) を求める
  - $F_i = F_j T'[|F_j| + 1]$  とする

$$T = abbabaaab = F_0 F_1 \dots F_f \ (F_0 = \varepsilon)$$
 $T' = abbabaaab$ 
 $F = (\varepsilon)$ 

- $F_i$  ( $i \ge 1$ ) を計算するアルゴリズムは以下の通り:
  - $T' = T[1 + |F_0| + |F_1| + \dots + |F_{i-1}| \dots n]$  とし、T' の接頭辞となる最長の $F_j$  (j < i) を求める
  - $F_i = F_j T'[|F_j| + 1]$  とする

$$T' = abbabaaab$$
  
 $F = (\epsilon, a)$ 

- $F_i$  ( $i \ge 1$ ) を計算するアルゴリズムは以下の通り:
  - $T' = T[1 + |F_0| + |F_1| + \dots + |F_{i-1}| \dots n]$  とし、T' の接頭辞となる最長の $F_j$  (j < i) を求める
  - $F_i = F_j T'[|F_j| + 1]$  とする

$$T' = abbabaaab$$
  
 $F = (\epsilon, a, b)$ 

- $F_i$  ( $i \ge 1$ ) を計算するアルゴリズムは以下の通り:
  - $T' = T[1 + |F_0| + |F_1| + \dots + |F_{i-1}| \dots n]$  とし、T' の接頭辞となる最長の $F_j$  (j < i) を求める
  - $F_i = F_j T'[|F_j| + 1]$  とする

$$T' = abbabaaab$$
  
 $F = (\epsilon, a, b, ba)$ 

- $F_i$  ( $i \ge 1$ ) を計算するアルゴリズムは以下の通り:
  - $T' = T[1 + |F_0| + |F_1| + \dots + |F_{i-1}| \dots n]$  とし、T' の接頭辞となる最長の $F_j$  (j < i) を求める
  - $F_i = F_j T'[|F_j| + 1]$  とする

$$T' = abbabaaab$$
  
 $F = (\varepsilon, a, b, ba, baa)$ 

- $F_i$  ( $i \ge 1$ ) を計算するアルゴリズムは以下の通り:
  - $T' = T[1 + |F_0| + |F_1| + \dots + |F_{i-1}| \dots n]$  とし、T' の接頭辞となる最長の $F_j$  (j < i) を求める
  - $F_i = F_j T'[|F_j| + 1]$  とする

$$T' = abbabaaab$$
  
 $F = (\varepsilon, a, b, ba, baa, ab)$ 

- $F_i$  ( $i \ge 1$ ) を計算するアルゴリズムは以下の通り:
  - $T' = T[1 + |F_0| + |F_1| + \dots + |F_{i-1}| \dots n]$  とし、T' の接頭辞となる最長の $F_j$  (j < i) を求める
  - $F_i = F_j T'[|F_j| + 1]$  とする

$$T' = abbabaaab$$
  
 $F = (\varepsilon, a, b, ba, baa, ab)$ 

%もし  $T' = F_i$  となったら  $F_i = F_i$ とする

- Factor  $F_i$  は整数  $j_i$  と文字 $c_i$ を用いて  $(j_i, c_i)$  と表現できる
  - (j<sub>i</sub>, c<sub>i</sub>)の列 F' から元の文字列 T を復元できるため、この列が 圧縮表現になる!

```
T = abbabaaab

F = (\epsilon, a, b, ba, baa, ab)

F' = ((0, a), (0, b), (2, a), (3, a), (a, b))
```

#### 部分文字列圧縮問題

- 長さ *n* のテキスト *T* が<u>事前に</u>与えられる
  - T についての索引構造を事前に構築しておいてよい
- 以下のクエリを処理:
  - 入力: 整数  $l, r (1 \le l \le r \le n)$
  - 出力: T[l..r] の圧縮表現

T = abbabaaab, (l, r) = (2, 7)

 $\rightarrow T[l..r]$  のLZ78分解は (b, ba, baa)

#### LZ78分解の部分文字列圧縮

LZ78分解の部分文字列圧縮 [Köppl, '21]

LZ78分解の部分文字列圧縮は、1クエリあたり

 $O(z_{l,r})$  time • O(n) space

で解くことができる

% n = |T| で、 $z_{l,r}$  は T[l...r] のLZ78分解のfactor数

■ 接尾辞木(Suffix Tree)を用いた手法

本研究はこの手法の省領域化

%ワードサイズ  $\Omega(n)$  のword-RAM modelを仮定、空間計算量はワード数で表記

#### 発表の流れ

- 1. 接尾辞木によるLZ78分解の計算
- 2. 1. の部分文字列圧縮版
- 3. CDAWGによるLZ78分解の計算

4. 3. の部分文字列圧縮版

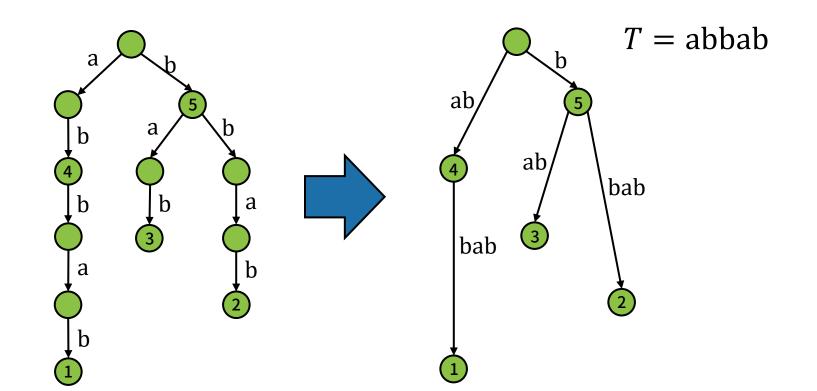
本研究の貢献

## 既存手法: 接尾辞木を使ったLZ78分解

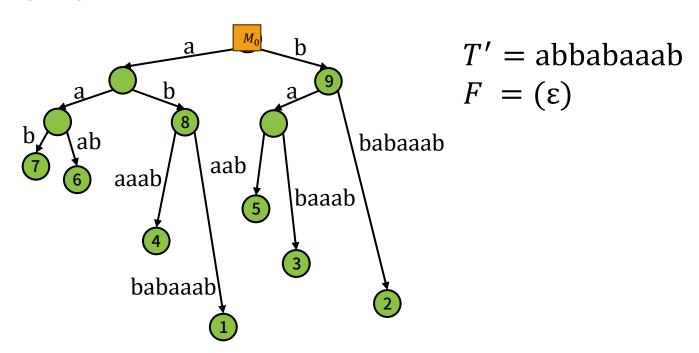
+部分文字列圧縮

# 接尾辞木 (Suffix Tree)

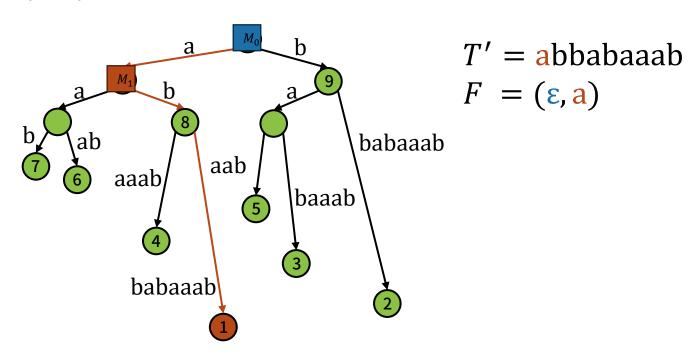
文字列の接尾辞集合を表すトライ木から、接尾辞を表さない出次数 1のノードをひとまとめにすることで得られる辺ラベル付き木



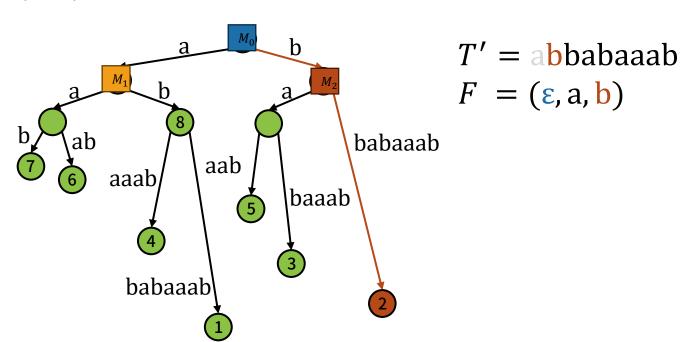
- $F_i$  ( $i \ge 1$ ) を計算するアルゴリズムは以下の通り:
  - 1.  $T' = T[1 + |F_0| + |F_1| + \cdots, |F_{i-1}| ... n]$  を表す接尾辞木上のパスを見つける
  - 2. パス上に存在する中で最深のマーク $M_i$ を探す
  - $F_i = F_j T'[|F_j| + 1]$  として、接尾辞木中の $F_i$  に対応する地点にマーク $M_i$  をつける



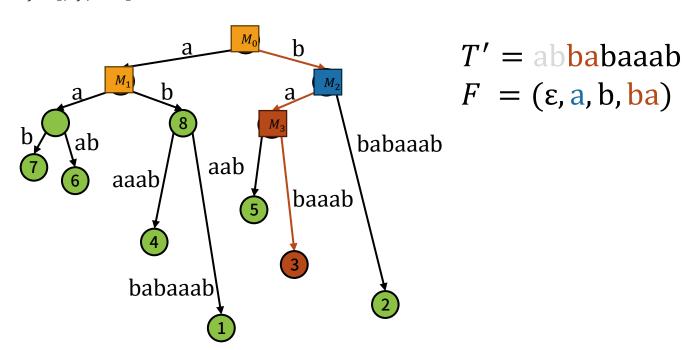
- $F_i$   $(i \ge 1)$  を計算するアルゴリズムは以下の通り:
  - 1.  $T' = T[1 + |F_0| + |F_1| + \cdots, |F_{i-1}| \cdot \cdot n]$  を表す接尾辞木上のパスを見つける
  - 2. パス上に存在する中で最深のマーク $M_i$ を探す
  - $F_i = F_j T'[|F_j| + 1]$  として、接尾辞木中の $F_i$  に対応する地点にマーク $M_i$  をつける



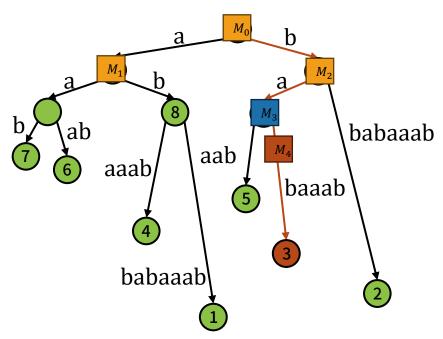
- $F_i$   $(i \ge 1)$  を計算するアルゴリズムは以下の通り:
  - 1.  $T' = T[1 + |F_0| + |F_1| + \cdots, |F_{i-1}| \cdot \cdot n]$  を表す接尾辞木上のパスを見つける
  - 2. パス上に存在する中で最深のマーク $M_j$ を探す
  - $F_i = F_j T'[|F_j| + 1]$  として、接尾辞木中の $F_i$  に対応する地点にマーク $M_i$  をつける



- $F_i$   $(i \ge 1)$  を計算するアルゴリズムは以下の通り:
  - 1.  $T' = T[1 + |F_0| + |F_1| + \cdots, |F_{i-1}| ... n]$  を表す接尾辞木上のパスを見つける
  - 2. パス上に存在する中で最深のマーク $M_i$ を探す
  - $F_i = F_i T'[|F_i| + 1]$  として、接尾辞木中の $F_i$  に対応する地点にマーク $M_i$  をつける



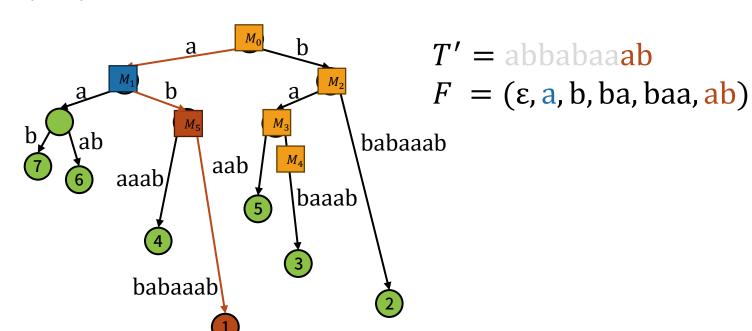
- $F_i$   $(i \ge 1)$  を計算するアルゴリズムは以下の通り:
  - 1.  $T' = T[1 + |F_0| + |F_1| + \cdots, |F_{i-1}| ... n]$  を表す接尾辞木上のパスを見つける
  - 2. パス上に存在する中で最深のマーク $M_i$ を探す
  - $F_i = F_j T'[|F_j| + 1]$  として、接尾辞木中の $F_i$  に対応する地点にマーク $M_i$  をつける



T' = abbabaaab

 $F = (\varepsilon, a, b, ba, baa)$ 

- $F_i$  ( $i \ge 1$ ) を計算するアルゴリズムは以下の通り:
  - 1.  $T' = T[1 + |F_0| + |F_1| + \cdots, |F_{i-1}| ...n]$  を表す接尾辞木上のパスを見つける
  - 2. パス上に存在する中で最深のマーク $M_i$ を探す
  - $F_i = F_j T'[|F_j| + 1]$  として、接尾辞木中の $F_i$  に対応する地点にマーク $M_i$  をつける



#### アルゴリズム中で用いるテクニック

- T'を表すパス上の最深マークを見つける・マークをつける
  - ▲ Lowest Marked Ancestor Problemであり、O(1) 時間で処理 できる [Westbrook, 1992]
- *F<sub>i</sub>* に対応する接尾辞木中の位置を得る
  - ▲ 接尾辞木上のWeighted Level Ancestor Problemであり、O(1) 時間で求まる [Gawrychowski et al., 2014, Bellazougui et al., 2021]

#### 時間•空間計算量

前処理も O(n) 時間

- **接尾辞木・LMA索引・WLA索引は全て** O(n) spaceなため、全体でも O(n) space
- マーク地点の取得・マークが定数時間で行えるため、factorあたり O(1) 時間
- $\Rightarrow$  全体の計算量は 空間 O(n) ・時間 O(z)

※ z は TのLZ分解のfactor数

# 接尾辞木によるLZ78分解の部分文字列圧縮

部分文字列圧縮への対応: とても簡単

■ 始点 *l* の対応: 前述のアルゴリズムの *T'* を

$$T' = T[l + |F_0| + |F_1| + \dots + |F_{i-1}| \dots n]$$
 に書き換えれば良い

■ 終点 r の対応: 通常通り計算して、 T[l..r] の範囲をはみ出したら末尾のfactorを削れば良い

計算量:  $O(z_{l,r})$  time •  $O(n+z_{l,r})$  space

※索引自体に <math>O(n) space必要で、作業領域が  $O(z_{l,r})$  space

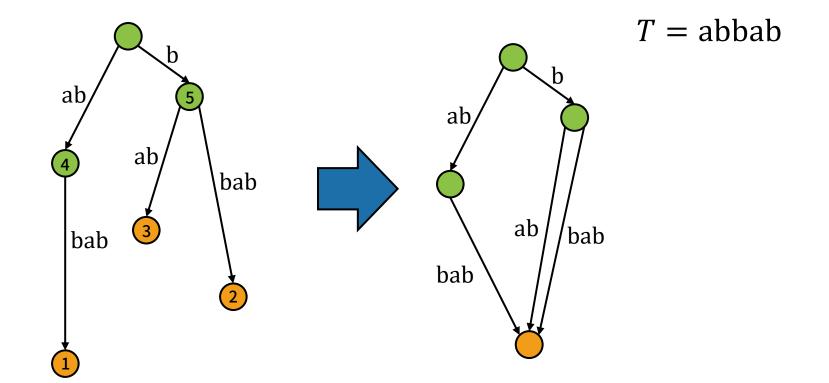
 $x \ge z_{l,r}$  は T[l...r] のLZ分解のfactor数

# 提案手法: CDAWGを使ったLZ78分解

+部分文字列圧縮

#### **CDAWG (Compacted Directed Acyclic Word Graph)**

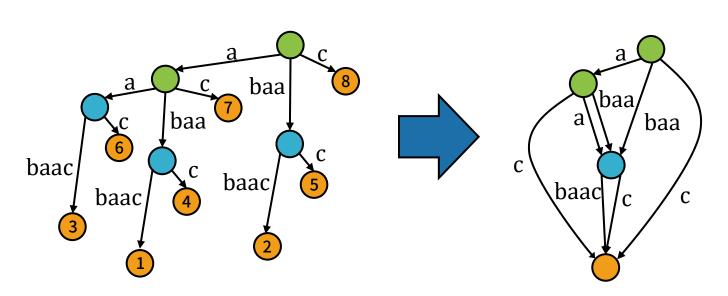
接尾辞木の<u>同型な部分木をひとまとめにする</u>ことで得られる 辺ラベル付きDAG



#### **CDAWG (Compacted Directed Acyclic Word Graph)**

接尾辞木の<u>同型な部分木をひとまとめにする</u>ことで得られる 辺ラベル付きDAG

T = abaabaac



#### CDAWGによる文字列圧縮

- e:CDAWGの辺の数
  - 接尾辞木の辺数は 2n-1 以下なので、  $e \le 2n-1$

性質: 文字列が繰り返し構造を持つと e は小さくなり、  $e \in \Theta(\log n)$  であるような文字列も存在

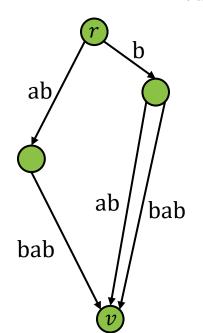


O(e) spaceで索引を作ることができれば、元文字列より省領域な圧縮索引になる!

### 接尾辞木→CDAWG で置き換える際の課題

先行研究のアルゴリズムをCDAWGでシミュレートしたい → 課題あり

- 1. CDAWGでは祖先が一意に定まらない
- 2. CDAWGでは1つの頂点が複数文字列を表す
- ⇒ 別のアプローチを用いて解決する



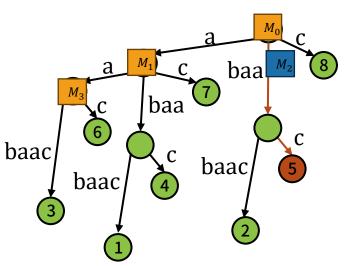
$$T = abbab$$
  
 $S(v) = \{abbab, bbab, bab\}$ 

v の入次数は3で、祖先は2頂点ある

 $X \otimes S(v)$ : CDAWGのroot  $r \to$  頂点 v のパスが表す文字列の集合

先行研究のアルゴリズム中で扱うMarked Ancestor は以下のようなもの:

- **1.** 頂点 v にマークをつける
- 2. 葉  $\ell$  を選び、 $\ell$  から根へのパス上で最も深いマークを見つける



- % 辺 (u,v) へのマークは頂点 v へのマークとしてよい
- ※一般には接尾辞を表す頂点が葉とは限らないが、 Tの末尾にuniqueな文字を追加することで 葉であることを保証できる

アルゴリズム中で扱うMarked Ancestor Problemは以下のようなもの:

- 1. 頂点 v にマークをつける
- 2. 葉  $\ell$  を選び、 $\ell$  から根へのパス上で最も深いマークを見つける

辺の探索順序は辞書順にしておく

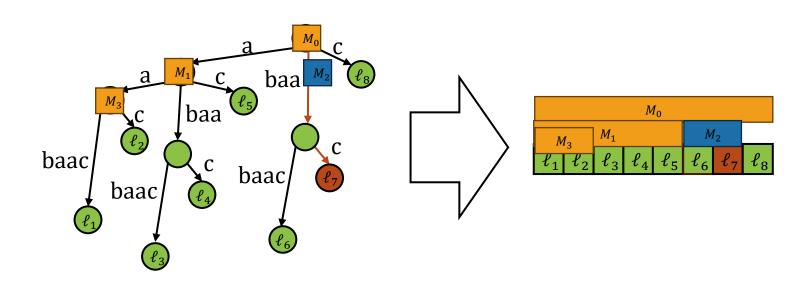
木の葉を行きがけ順に並べて  $\ell_1,\ell_2,...,\ell_n$  とすることで、上の2つの処理は以下の2つの処理に置き換えられる:

- 1. v の子孫である葉  $\ell_a$ ,  $\ell_{a+1}$ , ...,  $\ell_b$  に重み depth(v) でマークを付ける
- 2.  $\ell_k$  についたマークのうち、重みが最大のものを求める

※depth(v):頂点 v の深さ

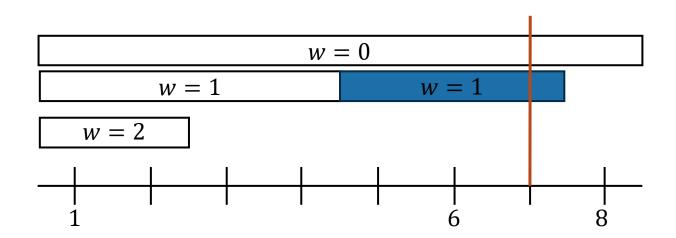
木の葉を行きがけ順(辺順序は辞書順)に並べて  $\ell_1,\ell_2,...,\ell_n$  とする

- 1. v の子孫である葉  $\ell_a$ ,  $\ell_{a+1}$ , ...,  $\ell_b$  に重み depth(v) でマークを付ける
- 2.  $\ell_k$  についたマークのうち、重みが最大のものを求める



問題をさらに抽象化すると、以下の問題に帰着できる:

- 1. 重み w の区間 [a, b] を集合に追加
- 2. 整数 k を包含するような区間のうち、重みが最大のものを求める



※頂点に対応する [a,b] や葉に対応する k は高速に求まると仮定

問題をさらに抽象化すると、以下の問題に帰着できる:

- 1. 重み w の区間 [a, b] を集合に追加
- 2. 整数 k を包含するような区間のうち、重みが最大のものを求める

この問題は1次元のstabbing-max problemであり、以下の条件を満たすようなデータ構造が存在 [Nekrich, 2011]

- 各クエリの時間計算量: ならし O(log n)
- 空間計算量: *0*(*m*) words

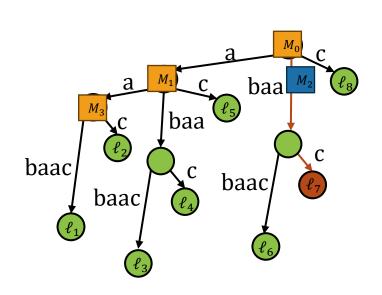
※ m:区間の個数

# T[l...r] に対応する区間の取得

以下の2つを高速に処理する必要がある:

- 1. (l,r) が与えられたとき、T[l..r] に対応する葉の区間 [a,b] を求める
- 2. i が与えられたとき、T[i...n] に対応する葉  $\ell_k$  を求める

葉を辞書順に並べると、これらは接尾辞配列上の操作に対応



				SA[i]	T[SA[i], n]
,	$M_3$	$M_1$	$M_0$	3	aabaac
				6	aac
				1	abaabaac
				4	abaac
				7	ac
		$M_2$		2	baabaac
				<b>-</b> 5	baac
				8	С

## T[l...r] に対応する区間の取得

以下の性質を満たすCDAWG-basedの圧縮索引が存在する

[Bealazzougui and Cunial, CPM 2017]

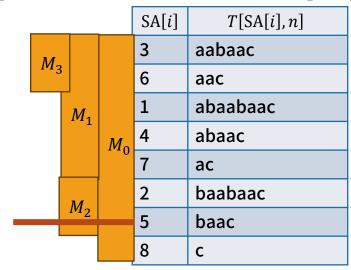
- T[i], ISA[i]の取得:  $O(\log n)$  time
- *T*[*l*..*r*] に対応する区間の取得: *O*(log *n*) time
- 空間計算量: *O(e)* space

%ISA: 接尾辞配列 SA の逆配列 (ISA[SA[i]] = i)

#### 以下のような手続きの繰り返しでLZ78のfactorを計算できる:

- 1.  $T' = T[1 + |F_0| + |F_1| + \dots + |F_{i-1}| \dots n]$  に対応する SA 上の位置  $k = \text{ISA}[1 + |F_0| + |F_1| + \dots + |F_{i-1}|]$  を探す
- 2. k を包含する重み w が最大の区間 [a, b] を求める
- 3.  $F_i = F_i T'[w+1]$  とする
- 4.  $F_i$  に対応する SA 上の区間 [a,b] を探し、重みw+1 の区間 [a,b] を追加

$$T' = abaabaac$$
  
 $F = (\varepsilon, a, b, aa)$ 



#### 以下のような手続きの繰り返しでLZ78のfactorを計算できる:

- 1.  $T' = T[1 + |F_0| + |F_1| + \dots + |F_{i-1}| \dots n]$  に対応する SA 上の位置  $k = \text{ISA}[1 + |F_0| + |F_1| + \dots + |F_{i-1}|]$  を探す
- 2. k を包含する重み w が最大の区間 [a, b] を求める
- 3.  $F_i = F_i T'[w+1]$  とする
- 4.  $F_i$  に対応する SA 上の区間 [a,b] を探し、重みw+1 の区間 [a,b] を追加

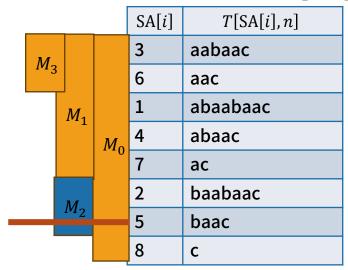
$$T' = abaabaac$$
  
 $F = (\varepsilon, a, b, aa)$ 

_				<u> </u>	
				SA[i]	T[SA[i], n]
	1/	;	$M_0$	3	aabaac
	$M_3$			6	aac
		$M_1$		1	abaabaac
		1		4	abaac
	_			7	ac
		11/1		2	baabaac
		$M_2$		5	baac
				8	С
				١٥	C

#### 以下のような手続きの繰り返しでLZ78のfactorを計算できる:

- 1.  $T' = T[1 + |F_0| + |F_1| + \dots + |F_{i-1}| \dots n]$  に対応する SA 上の位置  $k = \text{ISA}[1 + |F_0| + |F_1| + \dots + |F_{i-1}|]$  を探す
- 2. k を包含する重み w が最大の区間 [a, b] を求める
- 3.  $F_i = F_i T'[w+1]$  とする
- 4.  $F_i$  に対応する SA 上の区間 [a,b] を探し、重みw+1 の区間 [a,b] を追加

$$T' = abaabaac$$
  
 $F = (\varepsilon, a, b, aa, ba)$ 



#### 以下のような手続きの繰り返しでLZ78のfactorを計算できる:

- 1.  $T' = T[1 + |F_0| + |F_1| + \cdots, |F_{i-1}| ...n]$  に対応する SA 上の位置  $k = \text{ISA}[1 + |F_0| + |F_1| + \cdots + |F_{i-1}|]$  を探す
- 2. *k* を包含する重み *w* が最大の区間 [*a*, *b*] を求める
- 3.  $F_i = F_i T'[w+1]$  とする
- 4.  $F_i$  に対応する SA 上の区間 [a,b] を探し、重み w+1 の区間 [a,b] を追加

#### 計算量:

- 時間: ならし  $O(\log n)$  / factor  $\rightarrow$  全体で  $O(z \log n)$  time
- 空間: *O*(*e* + *z*) space
  - *e*(CDAWGの辺数)はCDAWG-basedの索引の分
  - z(TのLZ分解のfactor数)はstabbing-maxのためのデータ構造の分

factorの個数分だけ区間を保存する

## LZ78部分文字列圧縮への対応

接尾辞木の場合と同様に、部分文字列圧縮に対応できる

- 始点lの対応:前述のアルゴリズムのT'を $T' = T[\frac{l}{l} + |F_0| + |F_1| + \dots + |F_{i-1}| \dots n]$  に書き換えれば良い
- 終点 r の対応: 通常通り計算して、 T[l..r] の範囲を はみ出したら末尾のfactorを削れば良い

接尾辞木と同じ

計算量:  $O(z_{l,r}\log n)$  time •  $O(e+z_{l,r})$  チウェリの度にstabbing-maxのためのデータ構造を構築する

※索引自体は O(e) space、作業領域が  $O(z_{l,r})$  space

 $\divideontimes z_{l,r}$  は T[l...r] のLZ分解のfactor数

## まとめ・展望

#### LZ78分解の部分文字列圧縮を以下の計算量で行う手法を提案:

■ 時間:  $O(z_{l,r}\log n)$ 

■ 空間: *O*(*e*)

既存手法( $O(z_{l,r})$  time, O(n) space)からの省領域化を達成

#### 展望

- 他の問題も O(e) sizeの索引で解けないか
  - CDAWG上のパスクエリ等も圧縮領域で行えるため、応用は幅広い
- CDAWG以外の索引・他の部分文字列圧縮への一般化
  - 1文字アクセス・ISA・SA範囲の3処理が行える索引ならCDAWGの代わりに使える
  - 他の部分文字列圧縮もこの枠組みで解けるかもしれない(一部は既知)

# 補足: CDAWGを用いた索引の紹介

[Bealazzougui and Cunial, CPM 2017]

## CDAWGを用いた索引: 定理

定理([Belazzougui and Cunial, CPM 2017])

以下の処理を  $O(\log n)$  時間で行える

O(e) 領域のデータ構造が存在する:

- **1.** SA[i] の計算
- **2.** ISA[*i*] の計算
- 3. *T*[*i*] の計算
- 4. T[l...r] に対応する SA 上の区間の取得

# 説明の流れ

実は前述の4クエリのうち SA[i] 以外は、以下の方法で求められる

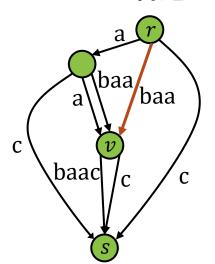
- 1. T の長さ k の接尾辞を表すパスを  $O(\log n)$  時間で求める
  - ullet パスの長さは  $\mathit{O}(e)$  になってしまうが、上手い表現(後述)で求める
- 2. パス上で何らかの処理を行ってクエリの答えを求める

計算時間を気にしない手法から紹介

そのため、まずは長さkの接尾辞を表すパスの取得方法を説明

## 準備: CDAWGに関する諸定義

- root r: 入次数が 0 の頂点(空文字列を表す頂点)
- sink s:出次数が 0 の頂点(接尾辞を表す頂点)
- $\blacksquare$  R(v), S(v): r-v パス, v-s パスの個数
  - 便宜上 R(s) = S(s) = 1 とする
- CDAWGの各辺は(始点,終点,辺ラベルの先頭文字,辺ラベル長)で表現



T = abaabaac

$$R(v) = 3$$

$$S(v) = 2$$

<del>赤色の辺</del>は (*r*, *v*, b, 3) のように表現

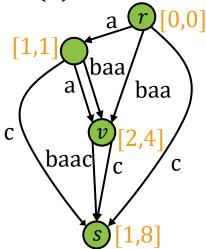
※以降ではTの末尾にuniqueな文字があると仮定

## 準備: CDAWGの性質

- str(v): CDAWGのroot  $r \rightarrow 頂点 v$  のパスが表す文字列の集合
  - |str(v)| = R(v) が成り立つ
- long(v):str(v) 中で最長の文字列

性質: str(v) は long(v) の先頭を 0,1,...,R(v)-1 文字削った接尾辞の集合

 $\rightarrow str(v)$  に含まれる文字列の長さの集合は区間  $[min_v, max_v]$  として表せる



$$T = abaabaac$$
  
 $str(v) = \{abaa, baa, ba\}$   
 $[min_v, max_v] = [2,4]$ 

# T[l,n] を表すパスの計算

v = s の状態からroot方向に、

u と v を結ぶ、辺ラベルの先頭文字が c で長さ  $\ell$  の辺

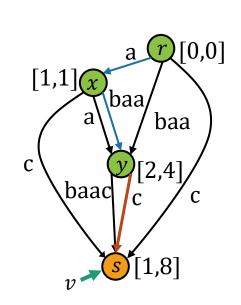
「長さkのr-v パスにおける最後の辺 $(u,v,c,\ell)$ 」を見つける処理を繰り返す

$$l = 3$$
,  $T[l, n] = abaabaac$ 

長さk=5のパスを見つけたい

最後の辺は (y, s, c, 1)

**」**辺の見つけ方はこれから説明



# T[l,n] を表すパスの計算

v = s の状態からroot方向に、

u と v を結ぶ、辺ラベルの先頭文字が c で長さ  $\ell$  の辺

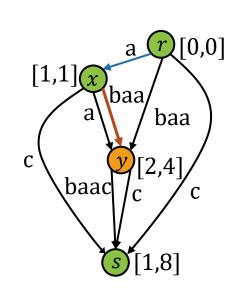
「長さkのr-v パスにおける最後の辺 $(u,v,c,\ell)$ 」を見つける処理を繰り返す

$$l = 3$$
,  $T[l, n] = abaabaac$ 

長さk=4のパスを見つけたい

最後の辺は (x, y, b, 3)

辺の見つけ方はこれから説明



# T[l,n] を表すパスの計算

v = s の状態からroot方向に、

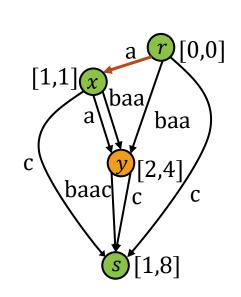
u と v を結ぶ、辺ラベルの先頭文字が c で長さ  $\ell$  の辺

「長さkのr-vパスにおける最後の辺 $(u,v,c,\ell)$ 」を見つける処理を 繰り返す

$$l = 3$$
,  $T[l, n] = abaabaac$ 

長さk=1のパスを見つけたい

**最後の辺は** (*r*, *x*, a, 1) 辺の見つけ方はこれから説明



# 最後の辺 $(u, v, c, \ell)$ の特定

長さkのr-vパスPが存在すると仮定する

このとき、P から最後の辺  $(u, v, c, \ell)$  を取り除いた r-u パスの長さは  $k-\ell$ 

 $\rightarrow k - \ell_i \in [\min_{u_i}, \max_{u_i}]$  である辺  $(u_i, v, c_i, \ell_i)$  が条件を満たす辺の候補

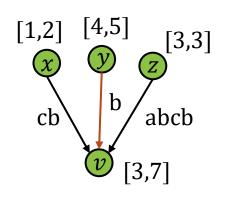
長さ $k - \ell_i$ のr - uパスが存在する

#### 右図の例で長さ5のパスを探す場合:

$$(x, v, c, 2)$$
:  $5 - 2 = 3 \notin [1,2]$ 

$$(y, v, b, 1)$$
:  $5 - 1 = 4 \in [4,5]$ 

(z, v, a, 4):5 - 4 = 1 ∉ [3,3] より、辿るべき辺は (y, v, b, 1)



# 最後の辺 $(u, v, c, \ell)$ の特定

長さkのr-vパスPが存在すると仮定する

このとき、P から最後の辺  $(u, v, c, \ell)$  を取り除いた r-w パスの長さは  $k-\ell$ 

 $ightarrow k - \ell_i \in \left[\min_{u_i}, \max_{u_i}\right]$  である辺  $(u_i, v, c_i, \ell_i)$  が条件を満たす辺の候補 長さ  $k - \ell_i$ の r - u パスが存在する

重要な性質:  $k - \ell_i \in S(u_i)$  を満たす辺  $(u_i, v, c_i, \ell_i)$  は必ず存在し、一意に定まる

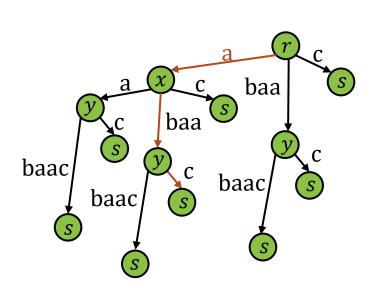
証明: 仮定より長さkのr-vパスは必ず存在する。r-vパスが表す文字列の集合は接尾辞集合になるため、同じ長さのr-vパスは一意に定まる。

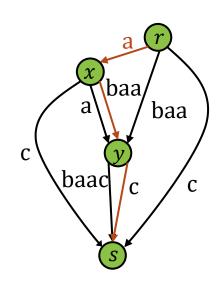
条件を満たす辺がない・複数存在するとパスの存在性・一意性と矛盾(証明終)

→ 条件を満たす辺を使ったroot方向への移動を繰り返せばパスが計算できる!

# T[l,n] を表すパスからの T[l] の計算

- 1. 前述の方法で T[l,n] に対応するパスを得る
- 2. パス中でrootと接続している辺が  $(u, v, c, \ell)$  のとき、 T[l] = c
  - ※ CDAWG索引の各辺は(始点,終点,辺ラベルの1文字目,長さ)の組



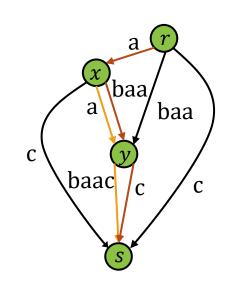


# T[l,n] を表すパスからの ISA[l] の計算

辺  $e = (u, v, c, \ell)$  について、u を始点とするような、辺ラベルが c より小さい辺の集合を L(e) とする

このとき、 $ISA[l] = 1 + \sum_{e \in P} \sum_{(u,v,c,\ell) \in L(e)} S(v)$  が成り立つ

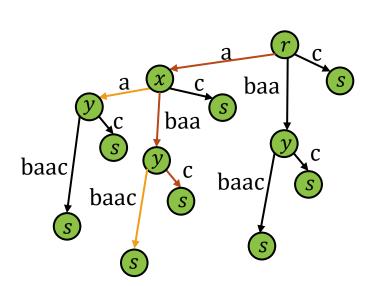
$$l = 3$$
,  $T[l, n] = abaabaac$   
 $L((x, y, b, 3)) = \{(x, y, a, 1)\}$   
 $L(y, s, c, 1)) = \{(y, s, b, 4)\}$   
 $ISA[l] = 1 + S(y) + S(s) = 1 + 2 + 1 = 4$ 

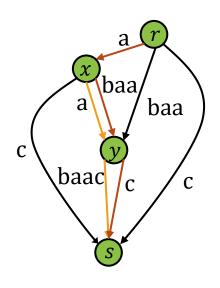


# T[l,n] を表すパスからの ISA[l] の計算

辺  $e = (u, v, c, \ell)$  について、u を始点とするような、辺ラベルが c より小さい辺の集合を L(e) とする

このとき、 $ISA[l] = 1 + \sum_{e \in P} \sum_{(u,v,c,\ell) \in L(e)} S(v)$  が成り立つ



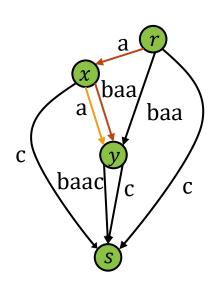


# 内部頂点に対応する SA 上の区間の計算

T[l...r] に対応する SA 上の区間 [a,b] は以下のように求められる:

- 1. T[l..n] を表すパスの深さr-l+1 地点までを表すパスP' を得る
- 2. ISA[l] の計算と同じ方法で、SA 上の区間の左端 a を求める
- 3. 区間の右端 b を b = a + S(v) 1 で求める ただし、v はパス P' の終端点

$$(l,r) = (3,6),$$
  $T[l,r] = abaabaac$   
 $E(x) = \{(x, y, a, 1)\}$   
 $a = 1 + S(y) = 3$   
 $b = a + S(y) - 1 = 4$ 

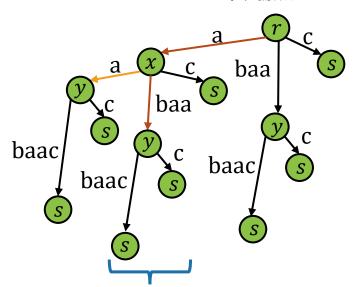


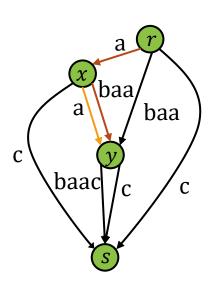
# 内部頂点に対応する SA 上の区間の計算

T[l..n] を表すパスの深さr-l+1 地点までを表すパスP' を得る

- 1. ISA[l] の計算と同じ方法で、SA 上の区間の左端 a を求める
- 2. 区間の右端 b を b = a + S(v) 1 で求める

ただし、v はパス P' の終端点





# 処理の高速化

説明では計算量について触れなかったが、この手法はとても遅い (愚直に実装するとO(e) 時間かかる)

以降では、この処理を2段階に分けて高速化する

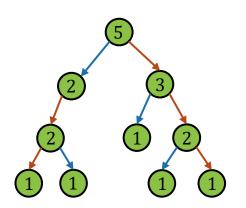
- 1. Symmetric Centroid Path Decomposition + 二分探索:  $O(\log^2 n)$
- 2. ↑の手法 + Interval Biased Search Tree:  $O(\log n)$

# **Heavy Path Decomposition (HPD)**

HPD: 木の辺をHeavy EdgeからなるパスとLight Edgeに分解する手法

各頂点の子方向の辺を、以下の二種類に分類する

- **Heavy Edge:** S(v) が最大の辺(複数ある場合はどれか一つ)
  - S(v) は v から葉へのパスの数の総和
- Light Edge: Heavy Edge以外の辺



## HPD: 定理1

始点を共有する辺のうち R(v) が最大の辺はHeavy Edge、そうでない辺は Light Edge

 $\lfloor \log_2 S(r) \rfloor = m$  とおくと、以下の定理が成り立つ。

定理1: Light Edgeの個数

任意のパス中のLight Edgeの個数は m 以下

#### 証明:

- 1. u の子の集合を C(u) とおくと、  $S(u) = \sum_{v \in C(u)} S(v)$
- 2. u-v にLight Edgeがあるとき、始点を共有するHeavy Edge u-v' が必ず存在
- 3.  $S(u) \ge S(v) + S(v'), S(v) \le S(v')$  より、 $S(u) \ge 2S(v)$  が成り立つ
- 4. Light Edgeを通るたびに S(v) が半分以下になるため、パス中のLight Edgeの個数は  $\lfloor \log_2 S(r) \rfloor = m$  以下

## HPD: 定理2

始点を共有する辺のうち R(v) が最大の辺はHeavy Edge、そうでない 辺はLight Edge

定理2: Heavy Path

同じ頂点を始点・終点とするHeavy Edgeの組は存在しない (→Heavy Edgeのみからなるグラフはパスの集合になる)

#### 証明(ほぼ自明):

- 1. 木なので、終点を共有するHeavy Edgeはない
- 2. Heavy Edgeは始点ごとに1つなため、始点を共有するHeavy Edgeもない

# Symmetric Centroid Path Decomposition (SCPD)

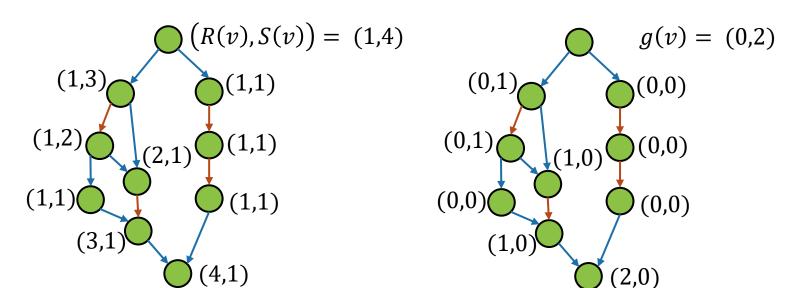
SCPD: root/sinkが1つのDAGに対するHeavy Path Decomposition

root, sinkからのパスの数を2進表記したときの最上位bitの位置

$$g(v) = (\lfloor \log_2 R(v) \rfloor, \lfloor \log_2 S(v) \rfloor)$$
 **L 5**

DAGにおいてu-vを結ぶ辺のうち、

g(u) = g(v) を満たす辺をHeavy Edge、そうでない辺をLight Edgeと呼ぶ



# Symmetric Centroid Path Decomposition (SCPD)

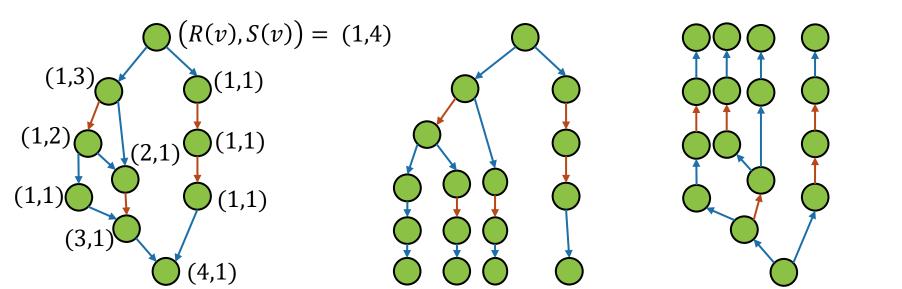
SCPD: root/sinkが1つのDAGに対するHeavy Path Decomposition

root, sinkからのパスの数を2進表記したときの最上位bitの位置

$$g(v) = (\lfloor \log_2 R(v) \rfloor, \lfloor \log_2 S(v) \rfloor)$$
 とする

DAGにおいてu-vを結ぶ辺のうち、

g(u) = g(v) を満たす辺をHeavy Edge、そうでない辺をLight Edgeと呼ぶ



## SCPD: 定理1

g(u) = g(v) を満たす辺はHeavy Edge、そうでない辺はLight Edge  $\lfloor \log_2 R(s) \rfloor = \lfloor \log_2 S(r) \rfloor = m$  とおくと、以下の定理が成り立つ。

定理1: Light Edgeの個数

任意のr-s パス中のLight Edgeの個数は2m 以下

#### 証明:

g(r) = (0, m), g(s) = (m, 0) が成り立つ。 任意の辺 (u, v) について  $R(u) \le R(v), S(u) \ge S(v)$  であることから、 $g(u) \ne g(v)$  である場合は $\lfloor \log_2 R(u) \rfloor < \lfloor \log_2 R(v) \rfloor$  または  $\lfloor \log_2 S(u) \rfloor > \lfloor \log_2 S(v) \rfloor$  となる。

Light Edgeを通るたびに  $\lfloor \log_2 R(v) \rfloor$  が増えるか  $\lfloor \log_2 S(v) \rfloor$  が減り、パスの両端は g(r)=(0,m),g(s)=(m,0) なため、Light Edgeの個数は 2m 以下。(証明終)

## SCPD: 定理2

g(u) = g(v) を満たす辺はHeavy Edge、そうでない辺はLight Edge 以下の定理が成り立つ。

定理2: Heavy Path

同じ頂点を始点・終点とするHeavy Edgeの組は存在しない (→Heavy Edgeのみからなるグラフはパスの集合になる)

#### 証明:

同じ頂点を始点とするHeavy Edgeの組 (u,v), (u,v') が存在すると仮定する。 Heavy Edgeの定義より  $\lfloor \log_2 S(u) \rfloor = \lfloor \log_2 S(v) \rfloor = \lfloor \log_2 S(v') \rfloor$  だが、  $\lfloor \log_2 S(u) \rfloor \geq \lfloor \log_2 S(v) + \log_2 S(v') \rfloor = \lfloor \log_2 S(v) \rfloor + 1$  となるため、矛盾。 終点が一致する場合も  $\lfloor \log_2 R(v) \rfloor$  を用いて同様の議論ができる。(証明終)

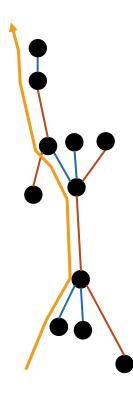
# 探したいパスとHeavy Pathの関係

探したいパスはいくつかのHeavy Path(Heavy Edgeからなる極大パス)とLight Edgeに分けられる

CDAWGでは R(s) = S(r) = n なので、 パス上のLight Edgeは  $2\log n$  個以下(定理1)

Light Edgeの個数が  $O(\log n)$  なため、Heavy Pathの個数も  $O(\log n)$  個

- $lacksymbol{\square}$  パスをHeavy Path上の一部分 or Light Edgeの列の形で $O(\log n)$  wordsで表現できる
- Heavy PathとLight Edgeを  $O(\log n)$  時間で探索できれば、全体の計算量を  $O(\log^2 n)$  時間にできる



# Heavy Path上の二分探索

Heavy Path(Heavy Edgeからなる極大パス)上の探索を高速化したい

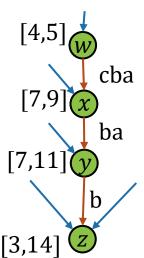
#### 定義:

- d(a,b): a-b を結ぶHeavy Pathの長さ(パス文字列の長さ)
  - ただし、a,b は同じHeavy Pathに属するとする

<u>長さkのr-vパスP</u>とvを含むHeavy Pathの共通部分を探すことを考える求めたいパスPがHeavy Path中のu-vパスを含むかは、

 $k - d(u, v) \in [\min_{u}, \max_{u}]$  によって判定できる

右図の例で長さ9のr-zパスとの共通部分を探す場合



# Heavy Path上の二分探索

Heavy Path(Heavy Edgeからなる極大パス)上の探索を高速化したい

#### 定義:

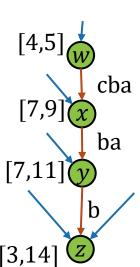
- d(a,b): a-b を結ぶHeavy Pathの長さ(パス文字列の長さ)
  - ただし、a,b は同じHeavy Pathに属するとする

<u>長さkのr-vパスP</u>とvを含むHeavy Pathの共通部分を探すことを考える求めたいパスPがHeavy Path中のu-vパスを含むかは、

 $k - d(u, v) \in [\min_{u}, \max_{u}]$  によって判定できる

Heavy Pathとの共通部分を求める計算量:

- u を決めたときの判定は定数時間で可能 各頂点 u について、d(u,h) と  $\min_u, \max_u$  を保存しておく (ただし、h は u を含むHeavy Pathの終点) d(u,v) = d(u,h) - d(v,h) より、d(u,v) は定数時間で計算できる
- Heavy Path上で二分探索すると、Heavy Path毎に  $O(\log n)$  time



# Light Edgeの二分探索

Light Edgeの探索を高速化したい

長さkのr-vパスPが辺 $(u,v,c,\ell)$ を含むかは $k-\ell \in [\min_u, \max_u]$ で判定できる

 $r-u_i-v$  パスが表す文字列の長さの最小値

v を終点として持つ各辺  $(u_i, v, c_i, \ell_i)$  を  $\min_{u_i} + \ell_i$  でソートしておくことで、

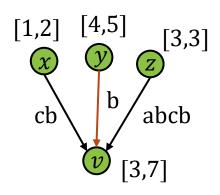
二分探索によって  $O(\log n)$  時間で目的のLight Edgeを特定できる

$$(x, v, c, 2)$$
:  $\min_{u} + \ell_{i} = 1 + 2 = 3$ 

$$(y, v, b, 1)$$
:  $\min_{u} + \ell_{i} = 4 + 1 = 5$ 

$$(z, v, a, 4)$$
:  $\min_{u} + \ell_{i} = 2 + 4 = 7$ 

長さが6のパスを探す場合、 $\min_u + \ell_i$ が6以下で最大の辺(y,v,b,1)が答え



# 現時点の計算量

- Heavy Edgeの探索: <u>Heavy Path毎に</u>  $O(\log n)$  時間
- Light Edgeの探索: <u>辺1つの特定に</u>  $O(\log n)$  時間

CDAWGでは R(s) = S(r) = n なので、 パス上のLight Edgeは  $2\log n$  個以下(定理1)

r-s パスは高々  $O(\log n)$  個しかLight Edgeを持たず、Heavy Pathの個数も  $O(\log n)$  となるため、計算量は合計で  $O(\log^2 n)$ 

 $\rightarrow O(\log n)$  に改善する手法を紹介

## **Biased Search Tree**

- 重み付きの集合  $S = \{(a_1, w_1), ..., (a_k, w_k)\}$  上の検索を高速に行うデータ構造
  - $a_i \ t + (a_1 < \dots < a_k)$
  - w<sub>i</sub> は要素重み (w<sub>i</sub>≥ 1)
- 要素  $(a_i, w_i)$  を検索する計算量は  $O\left(1 + \log \frac{w}{w_i}\right)$ 
  - W は要素重みの総和

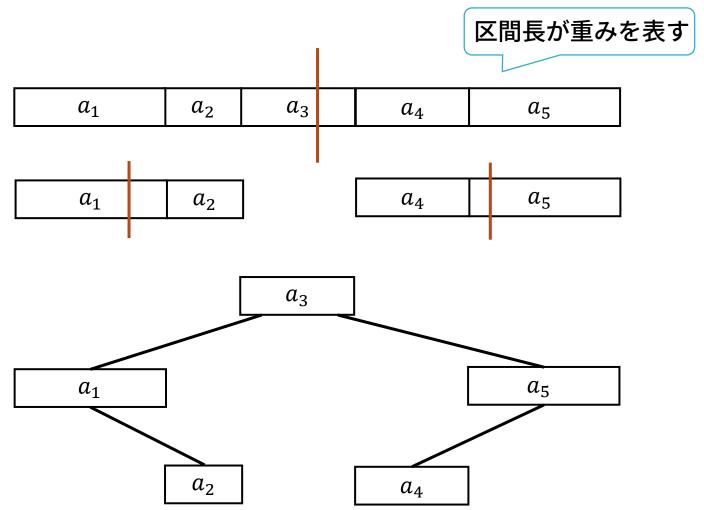
#### Biased Search Treeの構造

- $S = \{(a_1, w_1), ..., (a_k, w_k)\}$  について、要素  $(a_i, w_i)$  の累積重み  $z_i$  を  $z_i = \sum_{j=1}^i w_j$  で定義
- 集合 S のBST: 以下のように再帰的に定義
  - $S = \{(a_i, x_i)\}$  のとき:  $a_i$  のみからなる根のみの木
  - |S| ≥ 2 のとき: 以下のように定義

重み和の半分で分割

- lack 根頂点: 累積重みが $rac{W}{2}$ 以上となる中で最小の要素  $(a_i,x_i)$  の値  $a_i$
- ▲ 左の子:  $\{(a_1, w_1), ..., (a_{i-1}, w_{i-1})\}$  のBST
- ▲ 右の子:  $\{(a_{i+1}, w_{i+1}), ..., (a_k, w_k)\}$  のBST

#### Biased Search Treeの例



#### Biased Search Treeでの検索

- 検索: 通常の二分探索木と同様に比較しながら潜れば良い
- 要素  $(a_i, w_i)$  を検索する計算量は  $O\left(1 + \log \frac{W}{w_i}\right)$ 
  - 木を子方向へ移動する度、部分木の重み和は半分になる
    - lack 深さ d の頂点の重み和は $rac{W}{2^d}$ 以下
  - ullet 探索を続けると、 $a_i$  を表す頂点に必ず到達
  - *a<sub>i</sub>* を表す頂点以下の部分木重み和は必ず *w<sub>i</sub>* 以上
  - $\rightarrow w_i \leq$ 部分木重み和  $\leq \frac{W}{2^d}$  より、  $d \leq \log_2 \frac{W}{w_i}$

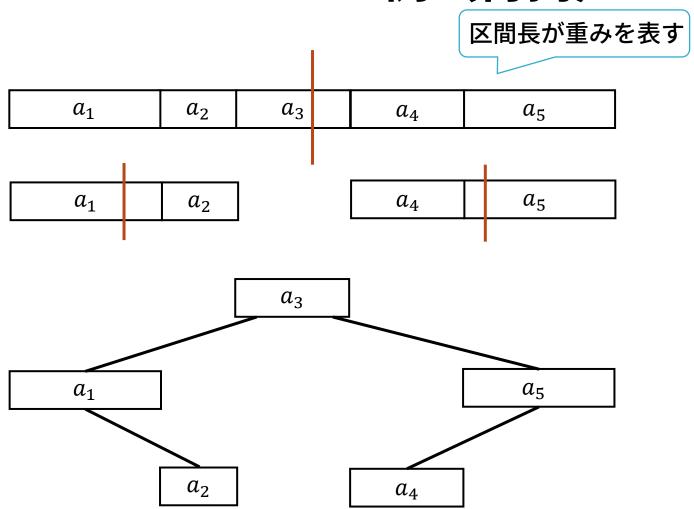
### Interval Biased Search Tree (IBST)

IBST: Biased Search Treeを区間の集合に対応させた亜種

- 集合  $S = \{(a_1, [l_1, r_1]) \dots, (a_k, [l_k, r_k])\}$  上の検索を高速に行うデータ構造
  - 区間  $[l_i,r_i]$  に値  $a_i$  が紐づいている
  - ullet  $[l_1,r_1],...,[l_k,r_k]$  はある区間の分割である必要がある
- Biased Search Treeを区間に対応させただけなので詳細は省略
- 区間長の和がWで探したい区間の長さが $w_i$ のとき、

計算量は
$$O\left(1 + \log \frac{W}{w_i}\right)$$

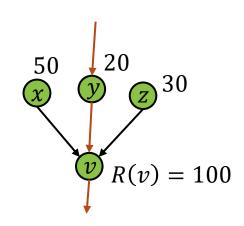
# Biased Search Treeの例(再掲)



#### IBSTを使うアイデア

- CDAWGのsink→root方向の探索時の性質をおさらい
  - 探索中では R(v) は広義単調減少で、 R(s) = n, R(r) = 1
  - v 終点の辺の始点の多重集合を  $S_v$  とすると、  $R(v) = \sum_{u \in S_v} R(u)$

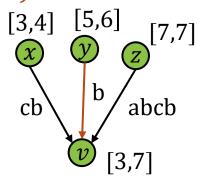
- 上の性質とIBSTを用いて、
  - $O(N \setminus A \cup B)$  でパスを得る手法を説明
    - SCPDを使わない手法
    - 後でSCPDと組み合わせて  $O(\log n)$  に改善



#### IBSTによる高速化

- 全ての頂点について、その頂点に入る辺をIBSTで管理
  - - ▲ 区間は辺  $(u,v,c,\ell)$ を通るような r-u-v パスの長さの集合を表す
    - $\blacktriangle$  区間長は R(u) で、頂点 v のIBSTの区間長の総和は R(v)

- このときのパス検索の計算量:  $O(\mathcal{N}$ ス長 +  $\log n$ )
  - 次ページで詳しく説明



#### IBSTによる高速化

■ このときのパス検索の計算量: *O*(パス長 + log *n*)

$$r - s$$
 パス  $v_1, ..., v_t$  を考える  $(v_1 = r, v_t = s)$ 

 $v_i$  のIBSTを用いて  $v_{i-1}$  を求める計算量は  $O\left(1 + \log \frac{R(v_i)}{R(v_{i-1})}\right)$ 

したがって、計算量の総和は

$$\sum_{i=t}^{i=2} O\left(1 + \log \frac{R(v_i)}{R(v_{i-1})}\right)$$
 間の項が打ち消せる 
$$= \sum_{i=t}^{i=2} O(1 + \log R(v_i) - \log R(v_{i-1}))$$
 
$$= O(t + \log R(v_t) - \log R(v_1))$$
 
$$= O(t + \log n)$$

#### さらなる改善

- 現在のパス検索の計算量: *O*(パス長 + log *n*)
  - 0(パス長) の部分がネック
- ボトルネックの原因: 辺を一本ずつ遡っていること
  - → Heavy Pathを高速に探索したい

■ これからやること: Heavy Pathの探索をIBSTで高速化

### Heavy Path上の二分探索(再掲)

Heavy Path(Heavy Edgeからなる極大パス)上の探索を高速化したい

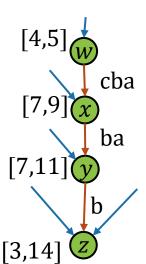
#### 定義:

- lack d(a,b): a-b を結ぶHeavy Pathの長さ(パス文字列の長さ)
  - ただし、a,b は同じHeavy Pathに属するとする

<u>長さkのr-vパスP</u>とvを含むHeavy Pathの共通部分を探すことを考える求めたいパスPがHeavy Path中のu-vパスを含むかは、

 $k - d(u, v) \in [\min_{u}, \max_{u}]$  によって判定できる

右図の例で長さ9のr-zパスとの共通部分を探す場合



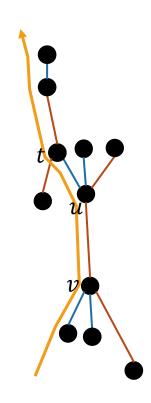
# O(log n) 時間アルゴリズムの概略

頂点 v を含むHeavy Pathと求めたいパスの共通部分を u-v とし、求めたいパス上で u より1つroot側の頂点を t とする

lack u-v 間はHeavy Pathで、t-u 間はLight Edge

このとき、v から t を求める操作が  $O\left(\log \frac{R(v)}{R(t)}\right)$  時間で行える場合、探索全体の計算量が  $O(\log n)$  になる

- Biased Search Treeと同様の計算量評価によって示せる
- ※どのHeavy Pathに属さない頂点も存在するが、 この場合は1頂点からなるHeavy Pathだと考えれば同様に扱える

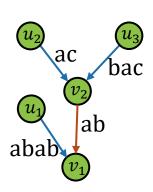


# 頂点・Light Edgeと区間の対応付け

- Heavy Path上の頂点をsink側から $v_1, ..., v_t$ とおく
  - 通常のパスと順序が逆なことに注意
- $lacksymbol{\blacksquare}$  Heavy Path上の頂点  $v_i$  について、

$$I_H(v_i) = \left[\min_{v_i} + d(v_i, v_1), \max_{v_i} + d(v_i, v_1)\right]$$
とする

•  $r - v_i - v_{i-1} - \cdots - v_1$  パスの文字列長の集合



■ Heavy Pathに繋がるLight Edge  $e = (u, v_i, c, \ell)$  について、

$$I_L(e) = \left[\min_{v_i} + d(v_i, v_1) + \ell, \max_{v_i} + d(v_i, v_1) + \ell\right]$$
 とする

• e を通る  $r - u - v_i - v_{i-1} - \cdots - v_1$  パスの文字列長の集合

# 頂点・Light Edgeと区間の対応付け

- - $r v_i v_{i-1} \cdots v_1$  パスの文字列長の集合
- **Light Edge**  $e = (u, v_i, c, \ell)$ について、

$$I_{L}(e) = \left[ \min_{v_i} + d(v_i, v_1) + \ell, \max_{v_i} + d(v_i, v_1) + \ell \right]$$

• e を通る  $r - u - v_i - v_{i-1} - \cdots - v_1$ パスの文字列長の集合

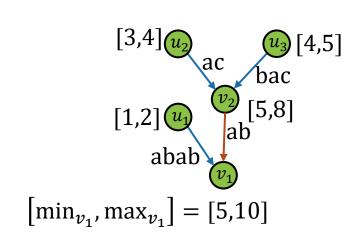
$$I_{H}(v_{1}) = [5,10]$$

$$I_{H}(v_{2}) = [7,10]$$

$$I_{L}((u_{1}, v_{1}, a, 4)) = [5,6]$$

$$I_{L}((u_{2}, v_{2}, a, 3)) = [7,8]$$

$$I_{L}((u_{3}, v_{2}, b, 2)) = [9,10]$$

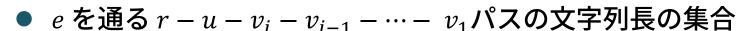


目標: Heavy Path $\rightarrow$ Light Edge の遷移を  $O\left(\log \frac{R(遷移元頂点)}{R(遷移先頂点)}\right)$  に改善したい

# 頂点・Light Edgeと区間の対応付け

- - $r v_i v_{i-1} \cdots v_1$  パスの文字列長の集合
- Heavy Pathに繋がるLight Edge  $e = (u, v_i, c, \ell)$ について、

$$I_{L}(e) = \left[ \min_{v_i} + d(v_i, v_1) + \ell, \max_{v_i} + d(v_i, v_1) + \ell \right]$$



#### ■ 性質: 以下の2つが成り立つ

- $\bullet \quad I_H(v_t) \subseteq I_H(v_{t-1}) \subseteq \cdots \subseteq I_H(v_1)$
- $v_i$  が終点のLight Edgeの集合を  $L(v_i)$  とおいたとき、

 $I_H(v_i)\setminus I_H(v_{i+1})=\bigcup_{X\in L(v_i)}I_L(X)$  であり、右辺の各区間は重ならない i=t なら空集合

# 頂点・Light Edgeと区間の対応付け

- 性質: 以下の2つが成り立つ(再掲)
  - $I_H(v_t) \subseteq I_H(v_{t-1}) \subseteq \cdots \subseteq I_H(v_1)$
  - $v_i$  が終点のLight Edgeの集合を  $L(v_i)$  とおいたとき、

 $I_H(v_i) \setminus I_H(v_{i+1}) = \bigcup_{X \in L(v_i)} I_L(X)$  であり、右辺の各区間は

重なりをもたない

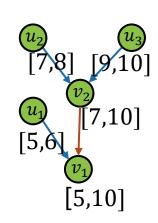
$$I_{H}(v_{1}) = [5,10]$$

$$I_{H}(v_{2}) = [7,10]$$

$$I_{L}((u_{1}, v_{1}, a, 4)) = [5,6]$$

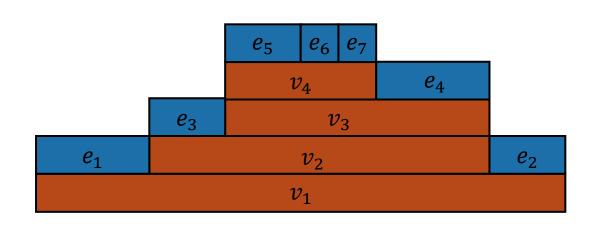
$$I_{L}((u_{2}, v_{2}, a, 3)) = [7,8]$$

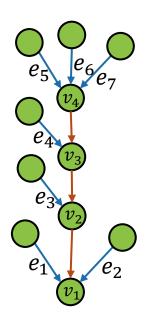
$$I_{L}((u_{3}, v_{2}, b, 2)) = [9,10]$$



# Heavy Path探索のIBST化

- *I<sub>H</sub>, I<sub>L</sub>* の性質から、Heavy Pathの各頂点・Heavy Pathに接続する各Light Edgeが表す区間は、下のように表せる
  - ullet Heavy Pathに接続するLight Edgeが表す区間の和集合が $I_H(v_1)$  と一致する
- Heavy Path毎に、Heavy Pathに接続するLight Edgeの集合をIBSTで保持 することを考える
  - → IBSTを単純に使うだけでは計算量保証ができない





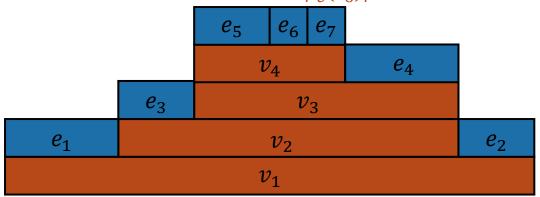
## Heavy Path + IBST の問題点

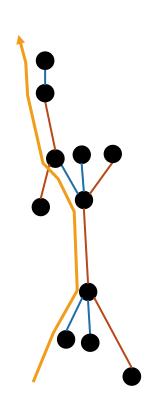
IBST化だけでは計算量は改善しない

原因: Light Edgeで遷移した先がHeavy Edgeの 途中になってしまうことがある

例:  $v_3$  を始点に二分探索をして辺  $e_6$  を探す場合

- $lacksymbol{\blacksquare}$   $O(\log n)$  の達成のために要求される計算量:  $O\left(\log \frac{|I_V(v_3)|}{|I_e(e_6)|}\right)$
- 実際にかかる時間:  $O\left(\log \frac{|I_V(v_1)|}{|I_e(e_6)|}\right)$

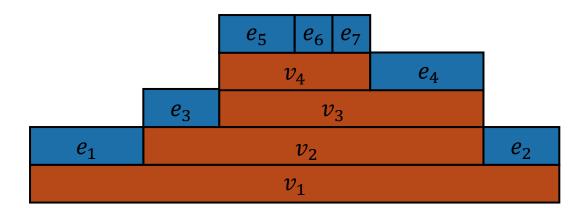




目標: Heavy Path $\rightarrow$ Light Edge の遷移を  $O\left(\log \frac{R($ 遷移元頂点 $)}{R($ 遷移先頂点 $)}$  に改善したい

#### 計算量の削減

- 改善したいこと: Heavy Pathの途中から探索したい
- Heavy Pathの各頂点  $v_i$  にjump pointerを保持して途中から探索することで、計算量を抑えられる
  - Jump pointer: 探索をショートカットできるようなIBSTのノードへのポインタ
  - ullet  $v_3$  にいる状態から  $e_1$  や  $e_3$  を探索しないようにして高速化

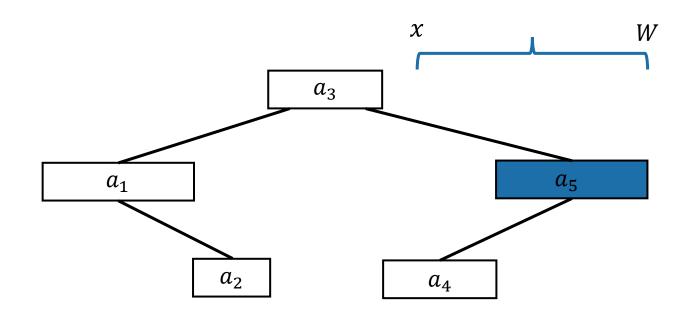


目標: Heavy Path $\rightarrow$ Light Edge の遷移を  $O\left(\log \frac{R(遷移元頂点)}{R(遷移先頂点)}\right)$  に改善したい

## Jump pointer

定義: jump pointer

区間全体の和集合が [1,W] であるようなIBSTにおいて、 ノード p 以下の部分木の区間の和集合が区間 [x,W] を包 含するような最深のノード  $p_x$  を、位置 x のjump pointerと呼ぶ。



目標: Heavy Path $\rightarrow$ Light Edge の遷移を  $O\left(\log \frac{R(遷移元頂点)}{R(遷移先頂点)}\right)$  に改善したい

# Jump pointerを使った場合の計算量

jump pointer: 部分木が [x, W] を包含する最深ノード  $p_x$  へのポインタ

 $p_x$  から探索を始めて区間  $[l,r] \in [x,W]$ を検索する計算量は  $O\left(1 + \log \frac{W-x+1}{r-l+1}\right)$ 

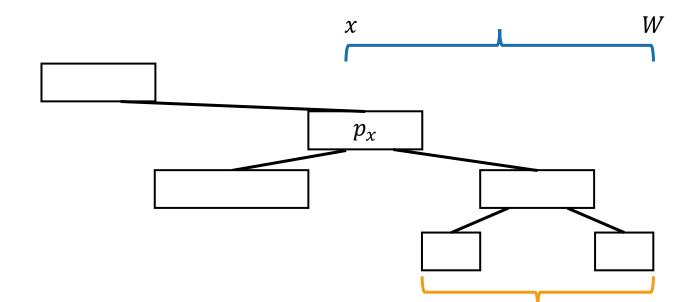
#### ■ 証明:

- 1. ノード  $p_x$  が表す区間は [x, W] と共通部分を持っている
  - lacktriangle そうでないのに部分木が [x, W] を包含するなら、その右の子孫がjump pointerになる
- 2.  $p_x$  の右部分木の任意のノードが表す区間は [x,W] に包含されている
  - $lacktriangle p_x$  が表す区間の右端が x 以上なため、右部分木の任意の区間も x 以上
- 3.  $p_x$ の右部分木の区間長総和がW-x+1以下で、探索する区間長が

$$r-l+1$$
 なため、探索のイテレーション回数は  $O\left(1+\log\frac{W-x+1}{r-l+1}\right)$ 

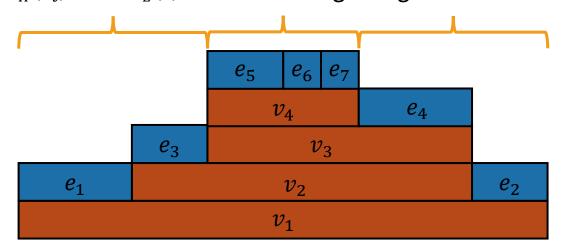
# Jump pointerを使った場合の計算量

- 1. ノード  $p_x$  が表す区間は [x, W] と共通部分を持っている
  - $\blacktriangle$  そうでないのに部分木が [x, W] を包含するなら、その右の子孫がjump pointer
- 2.  $p_x$  の右部分木の任意のノードが表す区間は [x, W] に包含されている
  - $lacktriangle p_x$  が表す区間の右端が x 以上なため、右部分木の任意の区間も x 以上



### Jump pointerを使った探索

- IBSTのjump pointerは区間の後半へのジャンプにしか対応していない
  - 今回は「区間の一部分にジャンプしたい」ため、そのまま使えない
- **IBSTを3種類保持することで解決** 
  - ullet Heavy Pathの始点  $v_t$  に直接繋がるLight Edgeを保持するIBST(図中央)
  - $\max I_L(e) \leq \min I_H(v_t)$  であるようなLight Edge e を保持するIBST(図左側)
  - $\max I_H(v_t) \leq \min I_L(e)$  であるようなLight Edge e を保持するIBST(図右側)



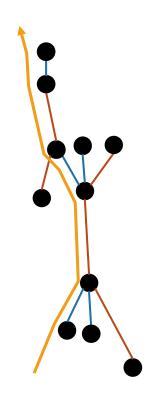
## Jump pointerを使った探索

- IBSTを3種類保持する
  - ullet Heavy Pathの始点  $v_t$  に直接繋がるLight Edgeを保持するIBST
    - lacktriangle jump pointerは使わない(もともと探索区間が  $I_H(v_t)$  の範囲に収まっている)
  - $\max I_L(e) \leq \min I_H(v_t)$  であるようなLight Edge e を保持するIBST
    - ▲ jump pointerで区間の後半にジャンプ
  - $\max I_H(v_t) \leq \min I_L(e)$  であるようなLight Edge e を保持するIBST
    - ▲ jump pointerで区間の前半にジャンプ
- 計算量は  $O\left(\log \frac{R(遷移元頂点)}{R(遷移先頂点)}\right)$

## Heavy Path + IBST手法のまとめ

- Heavy Path→Light Edgeの遷移をIBSTで行う
  - $\bullet$  一回の遷移の計算量は  $O\left(\log \frac{R(遷移元頂点)}{R(遷移先頂点)}\right)$
- 上の操作を  $O(\log n)$  回繰り返してパスを得る
- 全体の計算量は *O*(log *n*)
  - イテレーション回数は  $O(\log n)$
  - 遷移の計算量合計も O(log n)
- よって、*O*(log *n*) 時間でパスを得ることができた

以降、パスから ISA などを求める方法を紹介

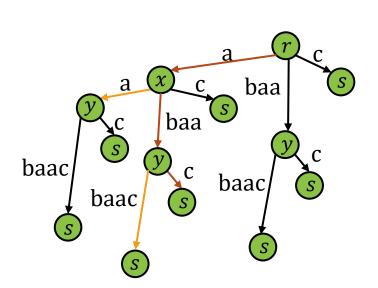


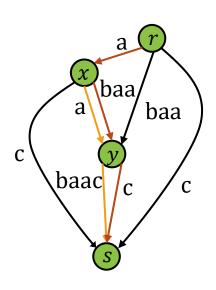
#### T[l,n] を表すパスからの ISA[l] の計算(再掲)

T[l,n] より辞書順で小さい接尾辞の数+1

辺  $e = (u, v, c, \ell)$  について、u を始点とするような、辺ラベルが c より小さい辺の集合を L(e) とする

このとき、 $ISA[l] = 1 + \sum_{e \in P} \sum_{(u,v,c,\ell) \in L(e)} S(v)$  が成り立つ





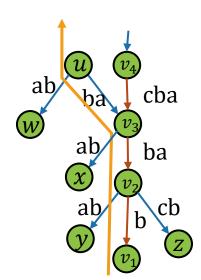
#### T[l,n] を表すパスからの ISA[l] の計算

T[l,n] より辞書順で小さい接尾辞の数+1

■ L(e): 辺  $e = (u, v, c, \ell)$  について、u を始点とするような、辺ラベルが c より小さい辺の集合

求めたいものは、パス上の全ての辺 e における  $\sum_{(u,v,c,\ell)\in L(e)} S(v)$  の総和以下を新たに定義:

- $S_H(v_i) = \sum_{j=1}^{i-1} \sum_{(u,v,c,\ell) \in L(e_{i-1})} S(v)$ 
  - ullet  $v_i$  はHeavy Pathに属する頂点
  - $e_j$  は  $v_{j+1}$  と  $v_j$  を結ぶHeavy Edge
- - e はLight Edge



### T[l,n] を表すパスからの ISA[l] の計算

T[l,n] より辞書順で小さい接尾辞の数+1

- $S_H(v_i) = \sum_{j=1}^{i-1} \sum_{(u,v,c,\ell) \in L(e_{j-1})} S(v)$ 
  - ullet  $v_i$  はHeavy Pathに属する頂点
  - $e_j$  は  $v_{j+1}$  と  $v_j$  を結ぶHeavy Edge
- - e はLight Edge



 $S_H(v_i)$ と  $S_L(e)$ を保存しておくことで、ISA[l] の暫定値を以下のように更新できる

- Heavy Path  $v_i v_j$  をパスに含める場合:  $S_H(v_i) S_H(v_j)$ を加算
- Light Edge e をパスに含める場合:  $S_L(e)$  を加算

ISA[l] の暫定値の更新は定数時間で行えるため、ISA[l] の計算量は  $O(\log n)$ 

### T[l,n] を表すパスからの ISA[l] の計算

T[l,n] より辞書順で小さい接尾辞の数+1

- $S_H(v_i) = \sum_{j=1}^{i-1} \sum_{(u,v,c,\ell) \in L(e_{j-1})} S(v)$ 
  - ullet  $v_i$  はHeavy Pathに属する頂点
  - $e_i$  は  $v_{i+1}$  と  $v_i$  を結ぶHeavy Edge
- - e はLight Edge

T[l,n] に対応する SA 上の区間の計算も 同様の方法で高速に行える(詳細略)

 $S_H(v_i)$  と  $S_L(e)$  を保存しておくことで、ISA[l] の暫定値を以下のように更新できる

- Heavy Path  $v_i v_j$  をパスに含める場合:  $S_H(v_i) S_H(v_j)$ を加算
- Light Edge e をパスに含める場合:  $S_L(e)$  を加算

ISA[l] の暫定値の更新は定数時間で行えるため、ISA[l] の計算量は  $O(\log n)$ 

#### 補足のまとめ

#### 以下のデータ構造について解説しました

定理([Belazzougui and Cunial, CPM 2017]+α)

以下の処理を  $O(\log n)$  時間で行える

O(e) 領域のデータ構造が存在する:

- **1.** SA[*i*] **の計算** 次ページで簡単に説明
- **2.** ISA[*i*] の計算
- **3.** T[i] の計算
- 4. T[l...r] に対応する SA 上の区間の取得

### 補足の補足: SA[i] の計算について

**1.** root→sinkを結ぶ辞書順 *i* 番目のパスを得る

長さでなく辞書順を指定してのパス取得だが、以下の改変で対応できる

- 逆向きではなく順方向(root→sink方向)に探索をするようにする
- 辺をパス長で並べていたところを辞書順で並べるようにする

Heavy Path+IBSTへの対応もでき、計算量は  $O(\log n)$ 

2. パスの長さがkのとき、n-k+1が答え

