

複数テキスト索引構造の オンライン構築

高木 拓也 (富士通研究所)

○ 稲永 俊介 (九州大学)

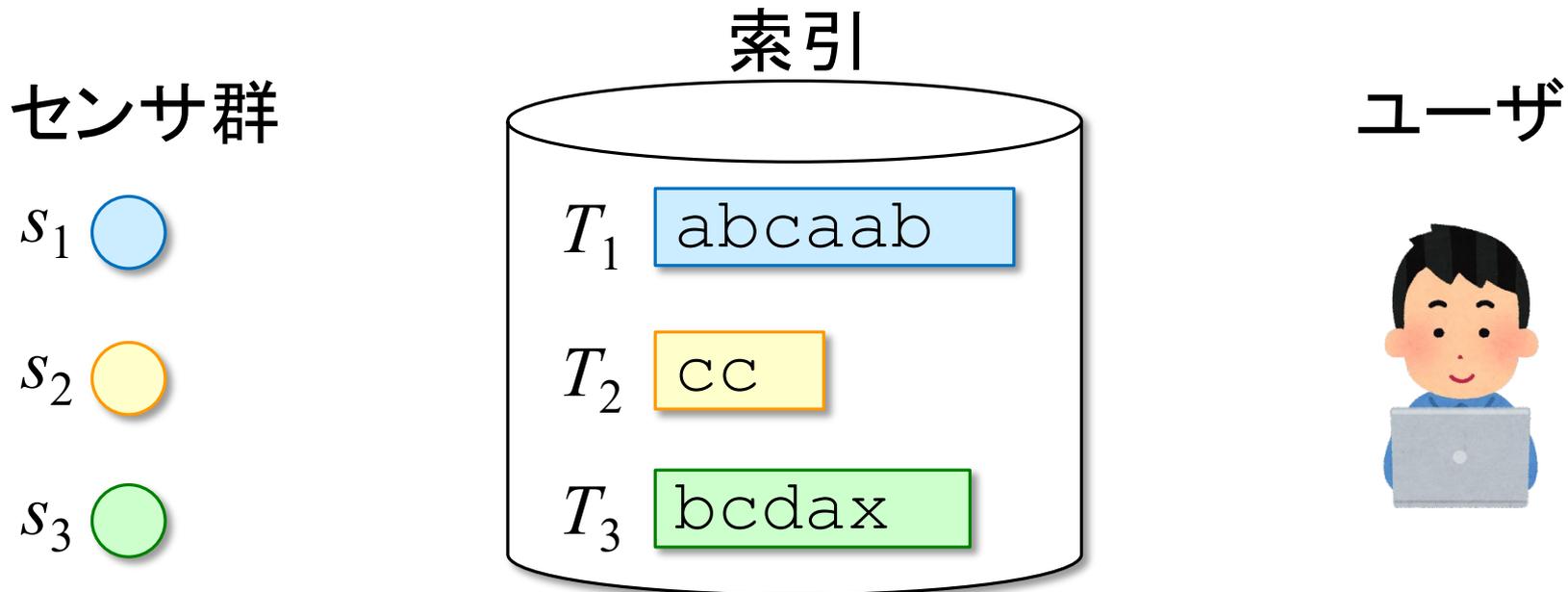
有村 博紀 (北海道大学)

Dany Breslauer (無所属)

Diptarama Hendrian (東北大学)

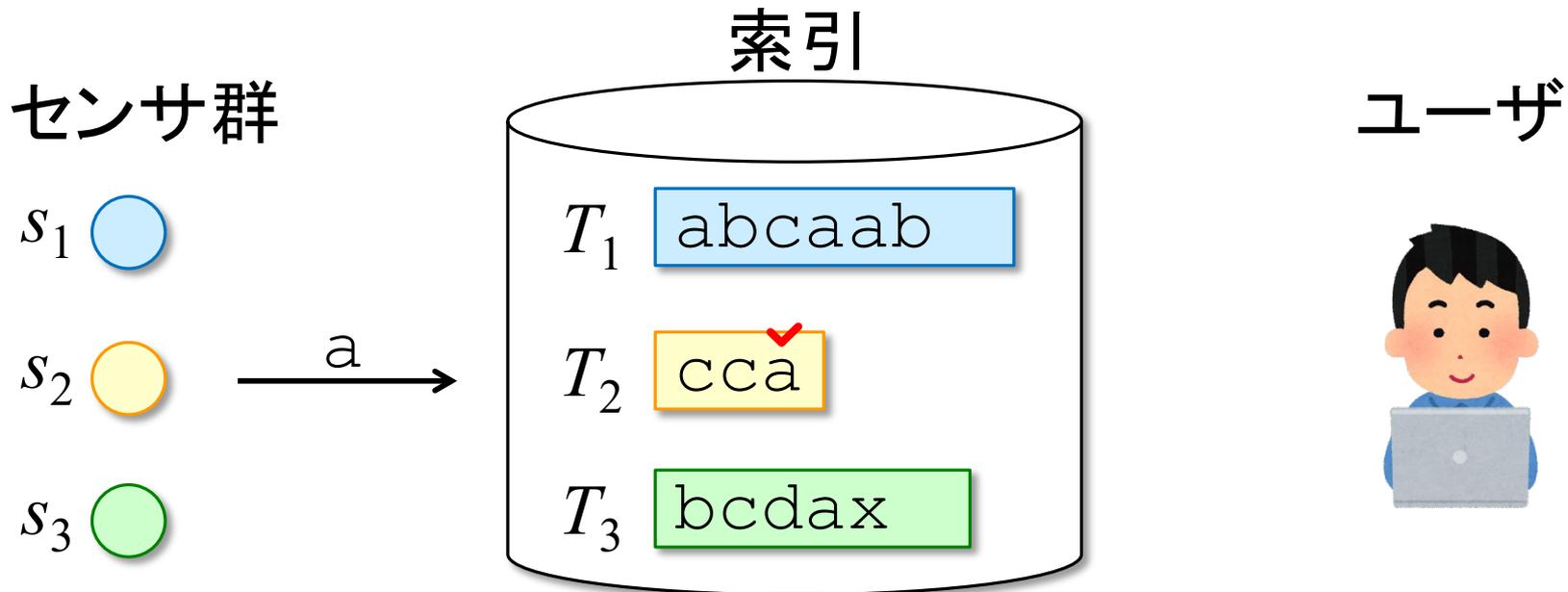
オンライン複数文字列の索引

- 問題: 末尾に文字の追加を許す複数の文字列に対する索引を高速に構築せよ.
- 動機: マルチストリームデータ処理
 - ◆ 各種センサ, トラジェクトリ, Twitter, などなど



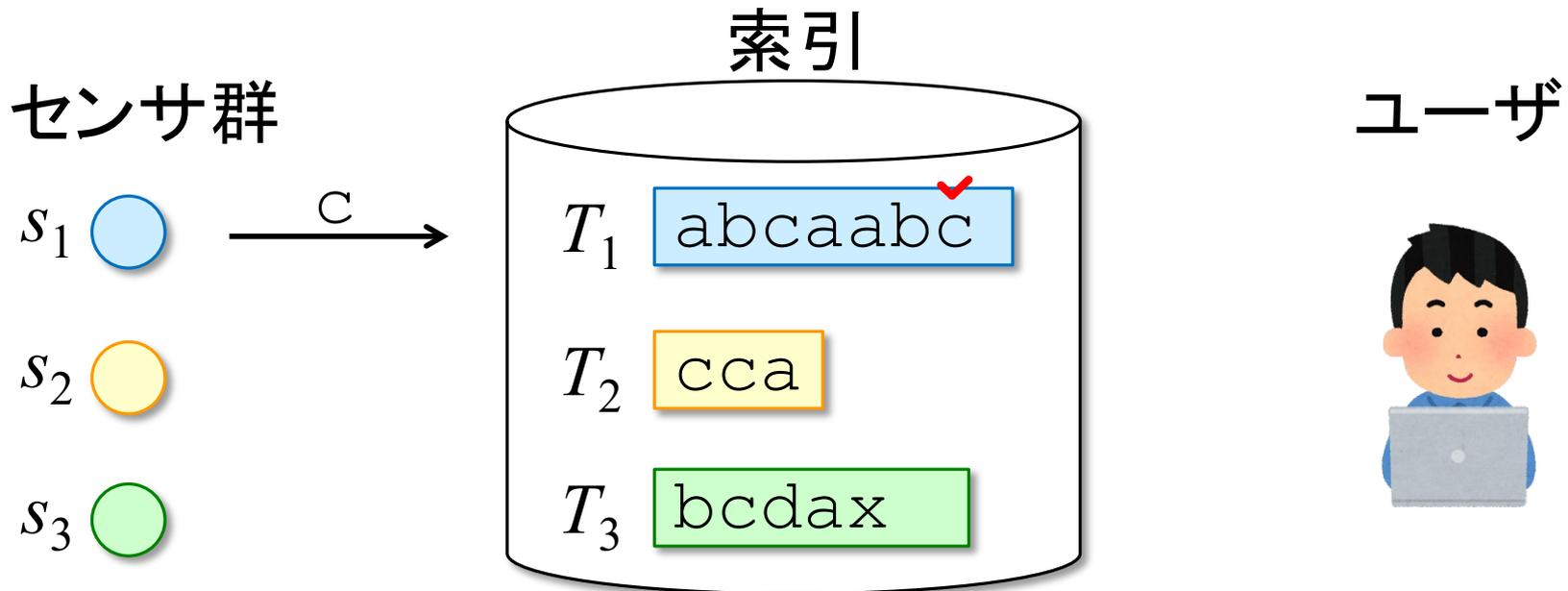
オンライン複数文字列の索引

- 問題: 末尾に文字の追加を許す複数の文字列に対する索引を高速に構築せよ.
- 動機: マルチストリームデータ処理
 - ◆ 各種センサ, トラジェクトリ, Twitter, などなど



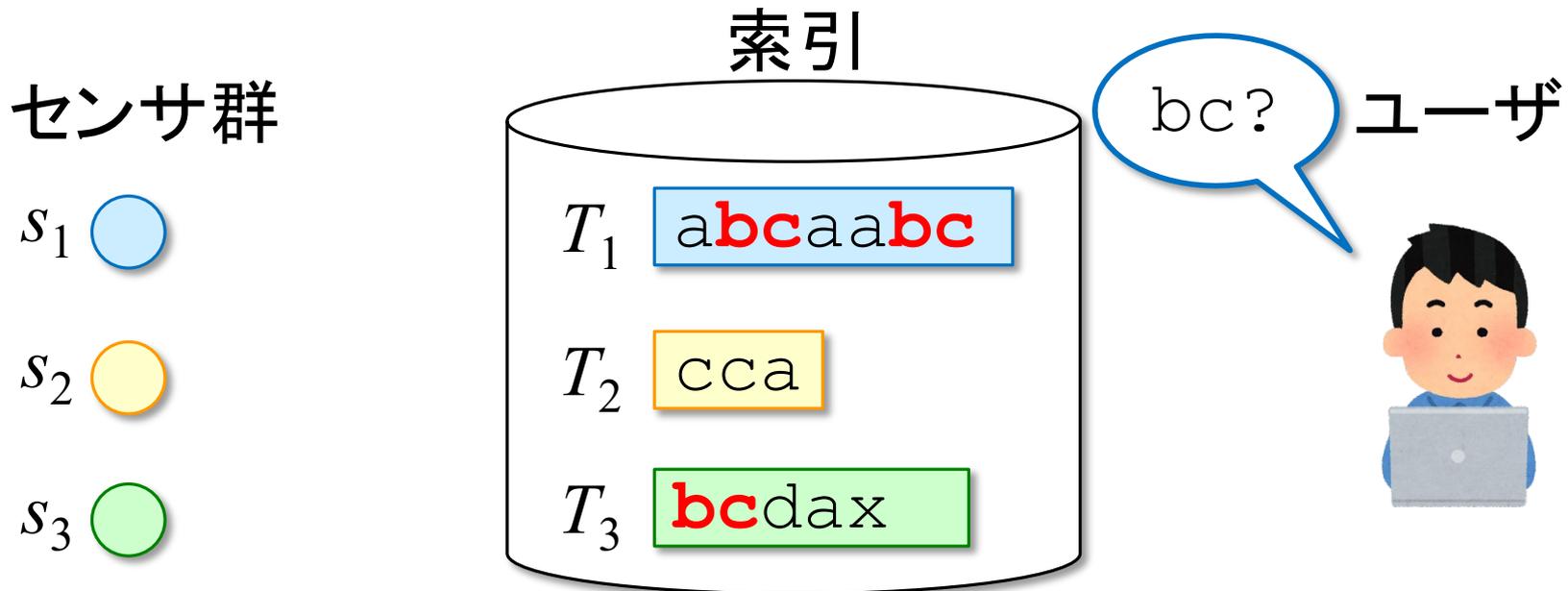
オンライン複数文字列の索引

- 問題: 末尾に文字の追加を許す複数の文字列に対する索引を高速に構築せよ.
- 動機: マルチストリームデータ処理
 - ◆ 各種センサ, トラジェクトリ, Twitter, などなど



オンライン複数文字列の索引

- 問題: 末尾に文字の追加を許す複数の文字列に対する索引を高速に構築せよ.
- 動機: マルチストリームデータ処理
 - ◆ 各種センサ, トラジェクトリ, Twitter, などなど



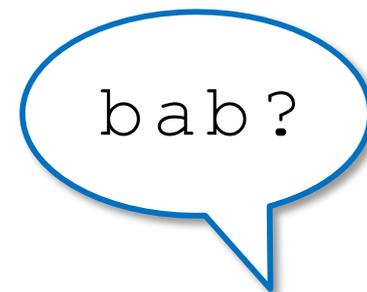
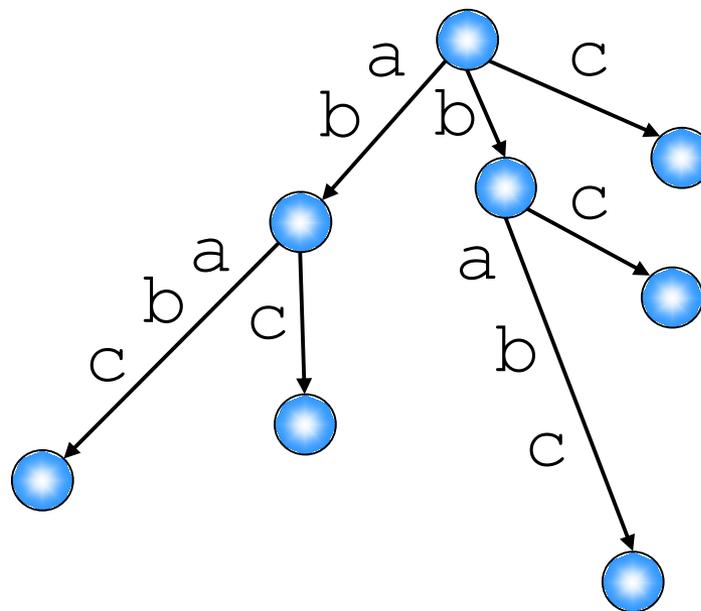
接尾辞木 (Suffix Tree)

- 入力文字列 T のすべての接尾辞を表現する根つき木 (オートマトンの亜種)

$T = \text{ababc}$

T の接尾辞

- ababc
- babc
- abc
- bc
- c



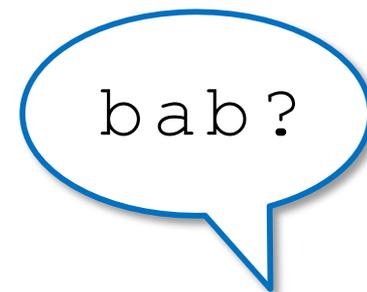
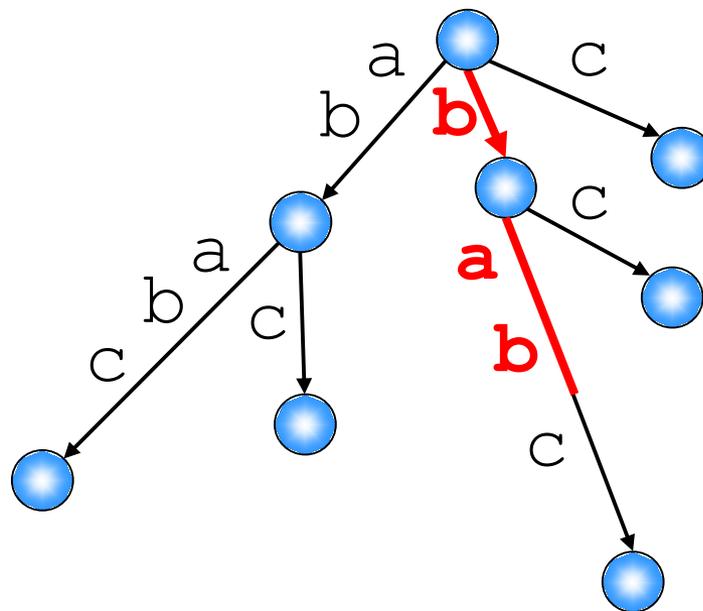
接尾辞木 (Suffix Tree)

- 入力文字列 T のすべての接尾辞を表現する根つき木 (オートマトンの亜種)

$T = a**bab**c$

T の接尾辞

- ababc
- bab**c
- abc
- bc
- c



複数文字列に対する接尾辞木

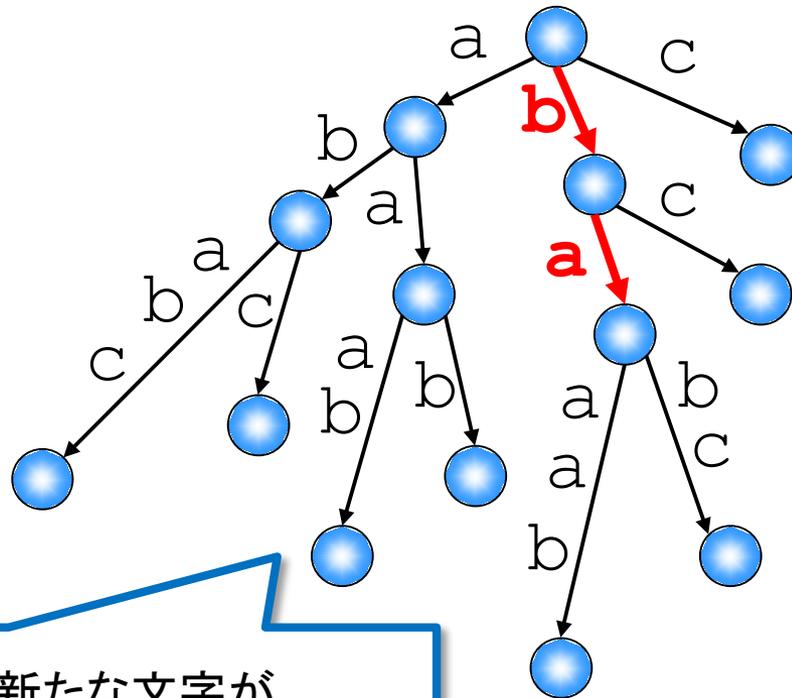
- 複数の入力文字列 T_1, \dots, T_K のすべての接尾辞を表現する根つき木

複数文字列の共通パターン発見

$T_1 = a$ **ba** bc

$T_2 =$ **ba** aab

$K = 2$ の場合



ba?



各文字列 T_i の末尾/先頭に新たな文字が追加される毎に、接尾辞木を高速に更新したい

単一文字列の場合

【本研究の出発点】
複数文字列の場合は？

右左オンライン構築 [Weiner 1973]

文字を先頭に追加していく単一文字列の接尾辞木を
 $O(n \log \sigma)$ 時間・ $O(n)$ 領域で構築できる。

左右オンライン構築 [Ukkonen 1995]

文字を末尾に追加していく単一文字列の接尾辞木を
 $O(n \log \sigma)$ 時間・ $O(n)$ 領域で構築できる。

n : 文字列の長さ, σ : 文字種類数

※ 比較モデルでは $O(n \log \sigma)$ 時間が最適
← ソートの下界 $\Omega(n \log \sigma)$ が適用される

研究の初期段階

2015年5月～

右左オンライン(先頭に文字追加)については,
Weiner のアルゴリズムをそのまま
複数文字列に拡張できそうですね



高木君



有村先生



稲永

研究の初期段階

イケそうですね



高木君

イケてますね



有村先生



稲永

研究の初期段階

左右オンライン(末尾に文字追加)については、
Ukkonen のアルゴリズムをそのまま
複数文字列に拡張するのは難しそうですね？



高木君



有村先生



稲永

研究の初期段階

そこをなんとか



高木君



有村先生



稲永

研究の初期段階

右左オンラインの接尾辞木を並行して構築することで、
左右オンラインの接尾辞木を構築できそうです



高木君



有村先生



稲永

研究の初期段階

イケそうですね

イケてますね



高木君



有村先生



稲永

本研究にまつわるヒストリー

2015/5	研究開始(高木, 稲永, 有村)
2016/7	文字列処理に関する国際会議 CPM で発表

CPM 2016 の論文の主張

主張1 (右左オンライン構築)

Weiner のアルゴリズムをそのまま適用することで、先頭に文字が追加されていく**複数文字列**の接尾辞木を $O(N \log \sigma)$ 時間・ $O(N)$ 領域で構築できる。

主張2 (左右オンライン構築)

主張1の結果と、Ukkonen のアルゴリズムを用いて、末尾に文字が追加されていく**複数文字列**の接尾辞木を $O(N \log \sigma)$ 時間・ $O(N)$ 領域で構築できる。

N : 複数文字列の長さの総和, σ : 文字種類数

Dany からメールが届く 2017年1月



Dany Breslauer

あなた方の CPM の論文を読みました。
残念ながら**間違ってる**と思います。

Mjsk!?



高木君



有村先生



稲永

Dany からメールが届く 2017年1月



Dany Breslauer

マジです。
複数文字列のオンライン処理特有の
問題を見逃していると思います。

やってもた...



高木君



有村先生



稲永

我々の CPM 2016 論文は

完全に
間違い!!

主張1 (右左オンライン構築)

Weiner のアルゴリズムをそのまま適用することで、先頭に文字が追加されていく**複数文字列**の接尾辞木を $O(N \log \sigma)$ 時間 $O(N)$ 領域で構築できる。

主張2 (左右オンライン構築)

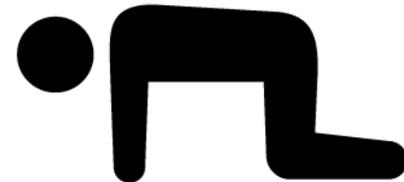
主張1の結果と、Ukkonen のアルゴリズムを用いて、末尾に文字が追加されていく**複数文字列**の接尾辞木を $O(n \log \sigma)$ 時間 $O(n)$ 領域で構築できる。

N : 複数文字列の長さの総和, σ : 文字種類数

本研究にまつわるヒストリー

2015/5	研究開始(高木, 稲永, 有村)
2016/7	文字列処理に関する国際会議 CPM で発表
2017/1	Dany から3人にメールが届き, CPM の論文に致命的なエラーが見つかる

全然イケてない...



チーム結成!!

2017年2~3月



高木君



有村先生



稲永

一緒に問題を
解決しましょう!



Dany



Diptarama君

Weiner アルゴリズムの下界

定理 1

Weiner のアルゴリズムをそのまま
複数文字列に適用すると $\Omega(N^{1.5})$ 時間かかる.

ウソの主張 1 (再掲)

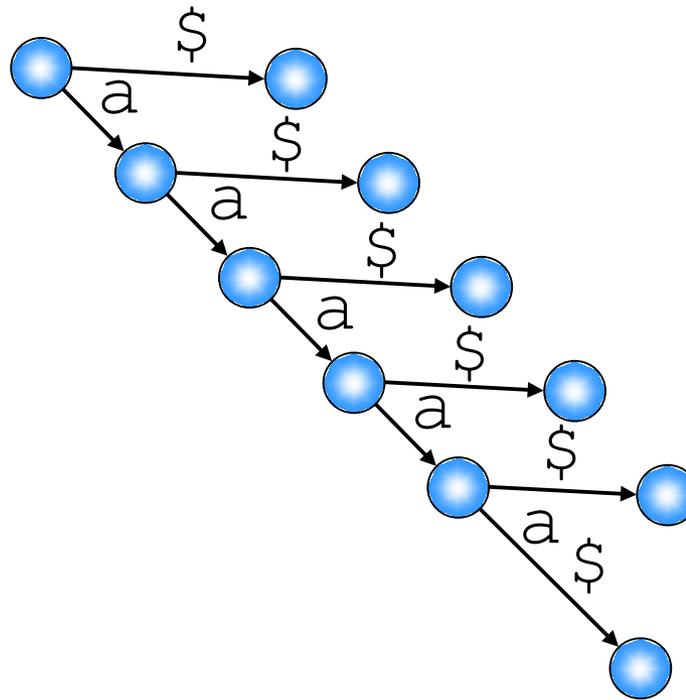
Weiner のアルゴリズムをそのまま適用することで、
先頭に文字が追加されていく **複数文字列**の接尾辞木を
 $O(N \log \sigma)$ 時間 $O(N)$ 領域で構築できる.

N : 複数文字列の長さの総和, σ : 文字種類数
→ $\sigma \leq N$ が常に成り立つ

$\Omega(N^{1.5})$ 時間の証明

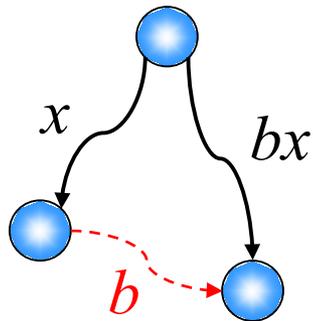
aaaaa\$ T_1
aaaa\$ T_2
aaa\$ T_3
aa\$ T_4
a\$ T_5

$\Omega(N^{1.5})$ 時間の証明

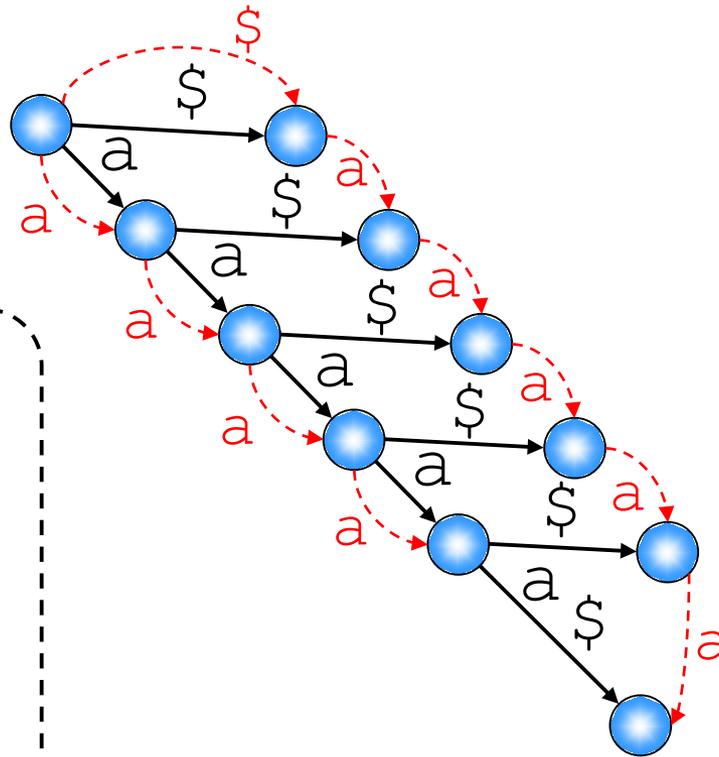


- aaaaa\$ T_1
- aaaa\$ T_2
- aaa\$ T_3
- aa\$ T_4
- a\$ T_5

$\Omega(N^{1.5})$ 時間の証明



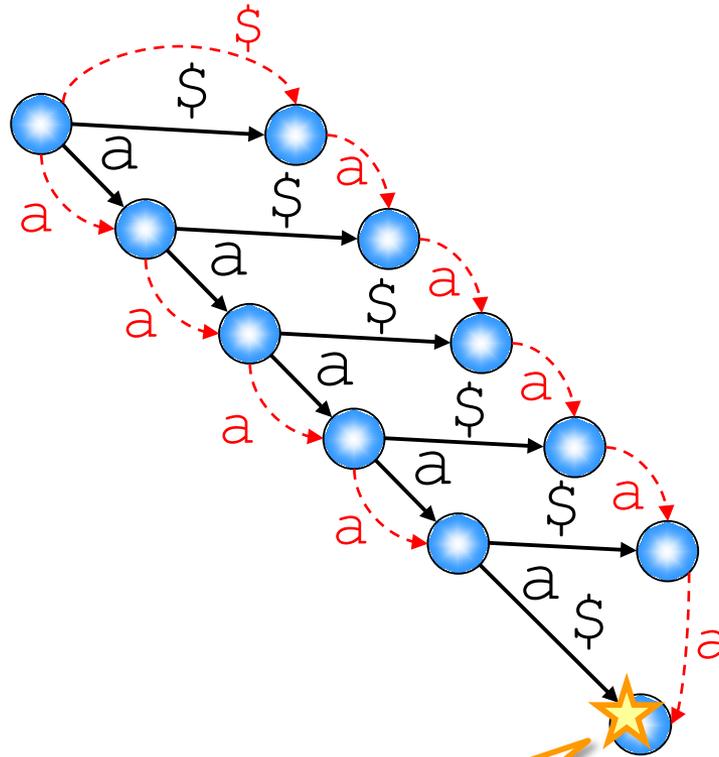
x は文字列, b は文字とする.
 x と bx を表す頂点があるとき,
 b でラベル付けされたリンクを
 x から bx に向かって張る.



aaaaa\$ T_1
 aaaa\$ T_2
 aaa\$ T_3
 aa\$ T_4
 a\$ T_5

←文字列を左に伸ばすときに便利

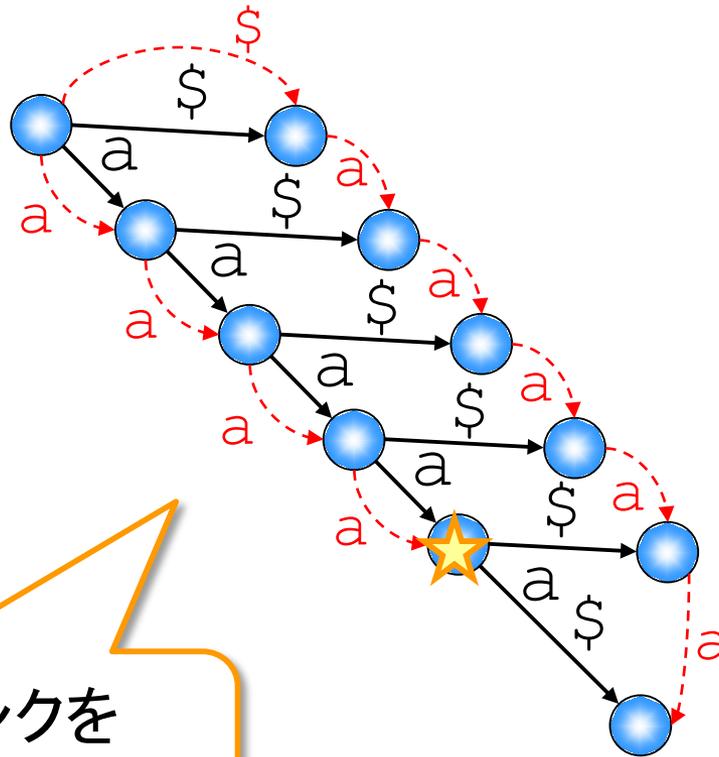
$\Omega(N^{1.5})$ 時間の証明



caaaaaa\$ T_1
▲ aaaaa\$ T_2
aaa\$ T_3
aa\$ T_4
a\$ T_5

文字追加前の T_1 を
表す葉からスタート

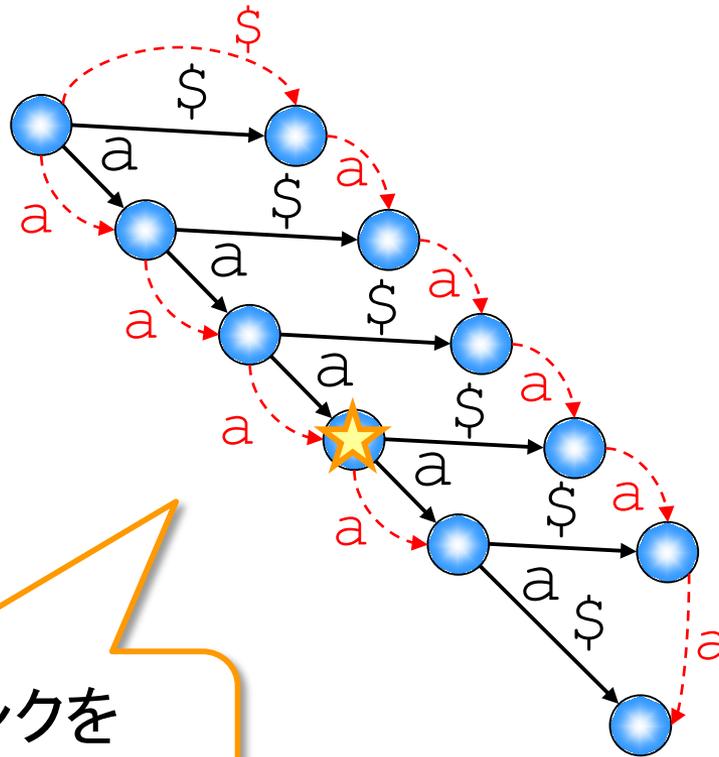
$\Omega(N^{1.5})$ 時間の証明



caaaaa\$ T_1
▲ aaaa\$ T_2
aaa\$ T_3
aa\$ T_4
a\$ T_5

追加文字 c のリンクを持つ最も深い頂点を探す

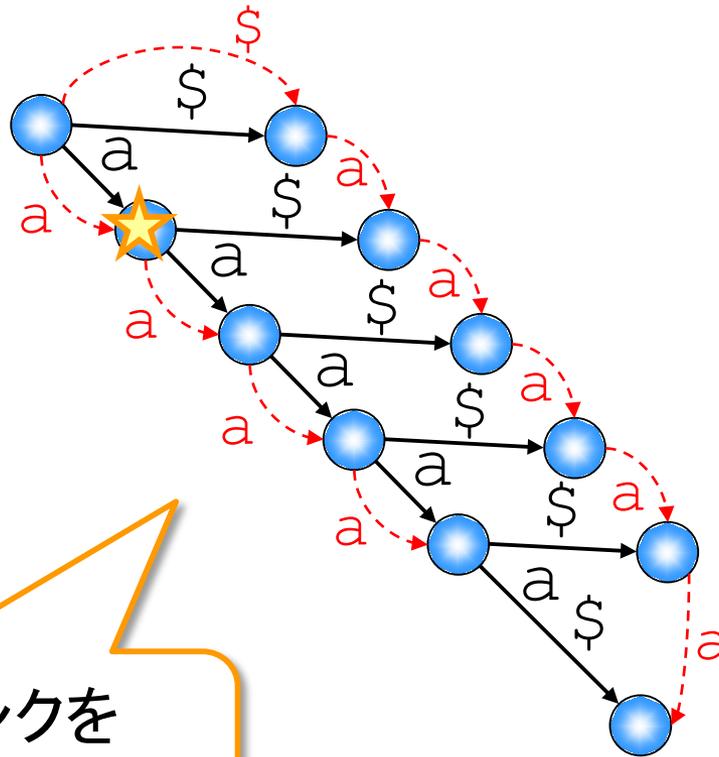
$\Omega(N^{1.5})$ 時間の証明



caaaaa\$ T_1
▲ aaaa\$ T_2
aaa\$ T_3
aa\$ T_4
a\$ T_5

追加文字 c のリンクを
持つ最も深い頂点を探す

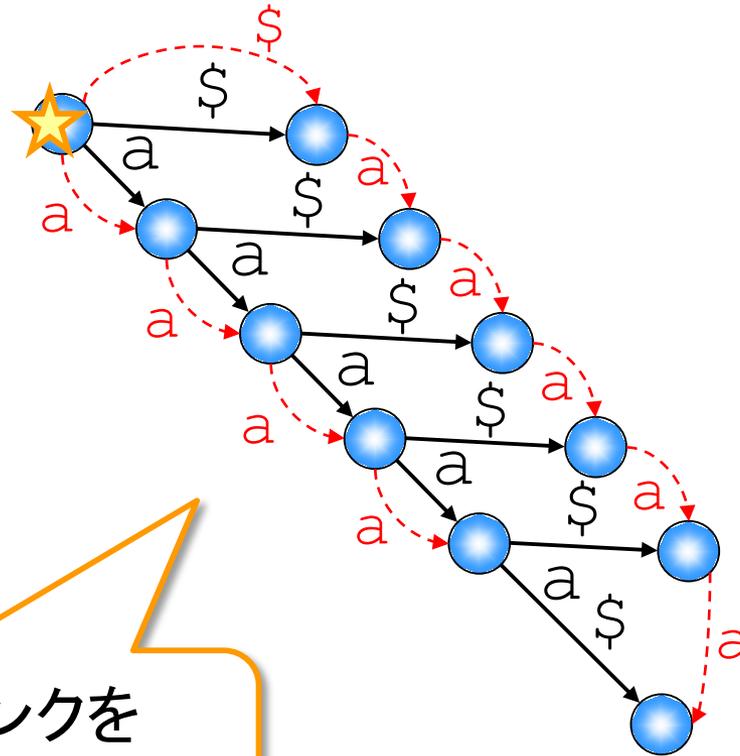
$\Omega(N^{1.5})$ 時間の証明



caaaaa\$ T_1
▲ aaaa\$ T_2
aaa\$ T_3
aa\$ T_4
a\$ T_5

追加文字 c のリンクを
持つ最も深い頂点を探す

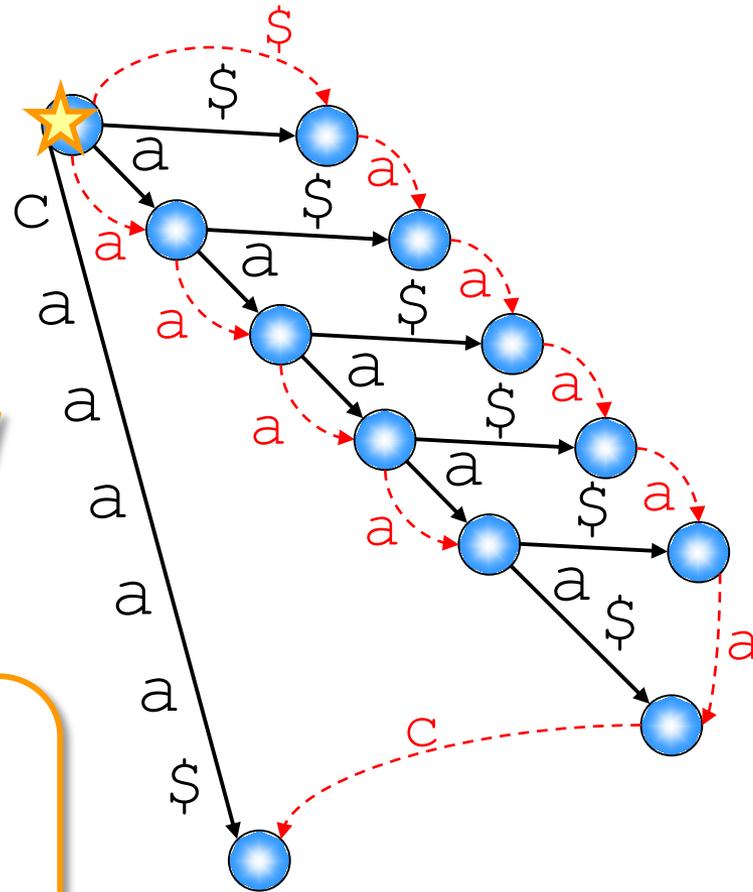
$\Omega(N^{1.5})$ 時間の証明



- caaaaa\$ T_1
- aaaaa\$ T_2
- aaa\$ T_3
- aa\$ T_4
- a\$ T_5

追加文字 c のリンクを
持つ最も深い頂点を探す

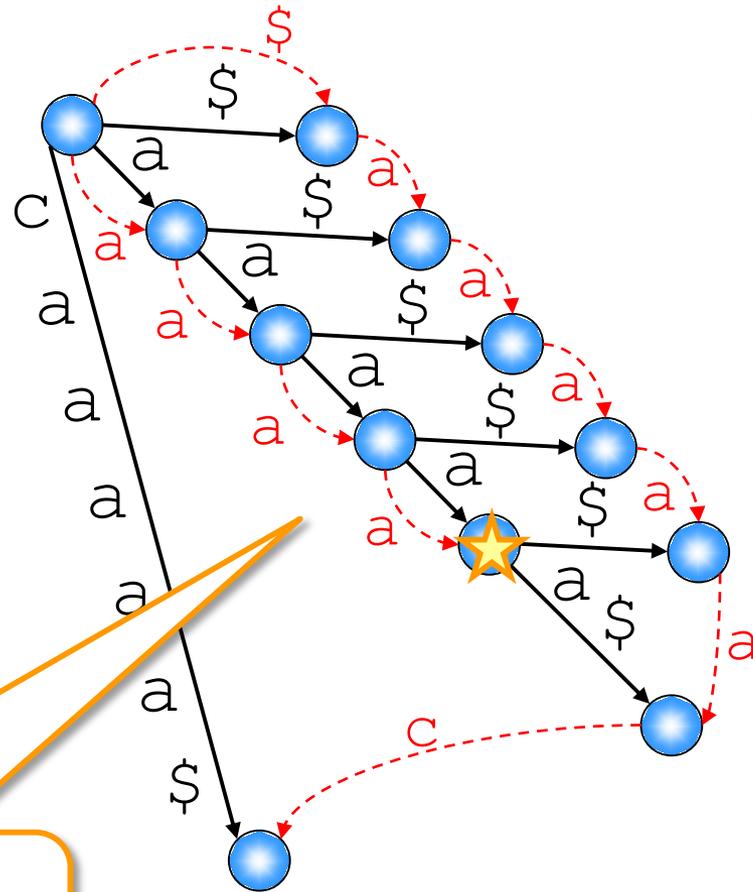
$\Omega(N^{1.5})$ 時間の証明



$c a a a a a \$$ T_1
 $\triangle a a a a \$$ T_2
 $a a a \$$ T_3
 $a a \$$ T_4
 $a \$$ T_5

根まで行っても c のリンク
 がなかったので、
文字追加後の T_1 を表す
 新たな葉を根から挿入

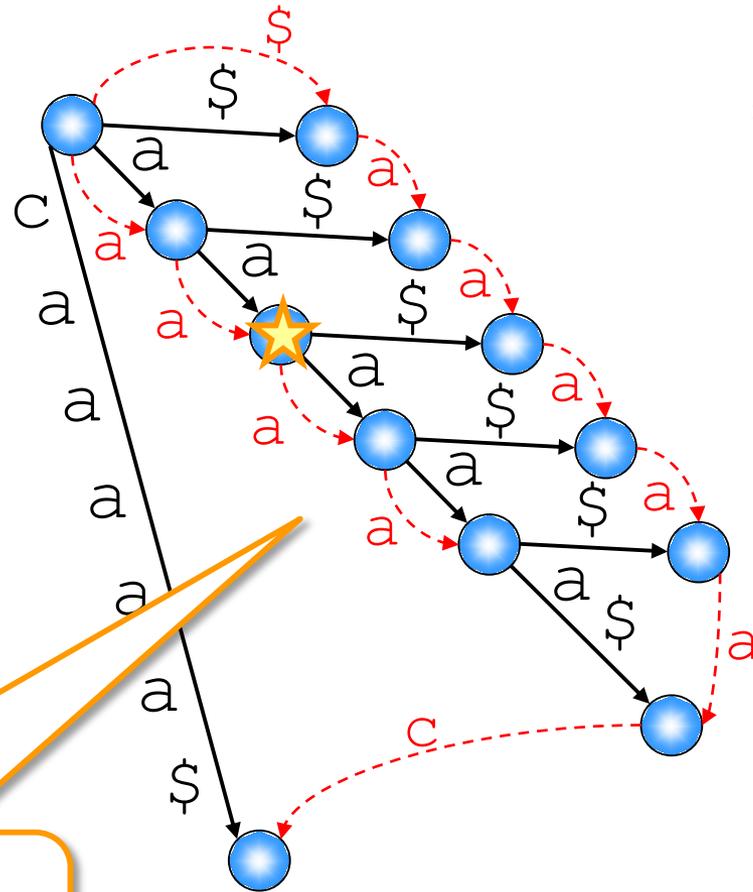
$\Omega(N^{1.5})$ 時間の証明



- caaaaa\$ T_1
- caaaa\$ T_2
- ▲ aaa\$ T_3
- aa\$ T_4
- a\$ T_5

追加文字 c のリンクを持つ最も深い頂点を探す

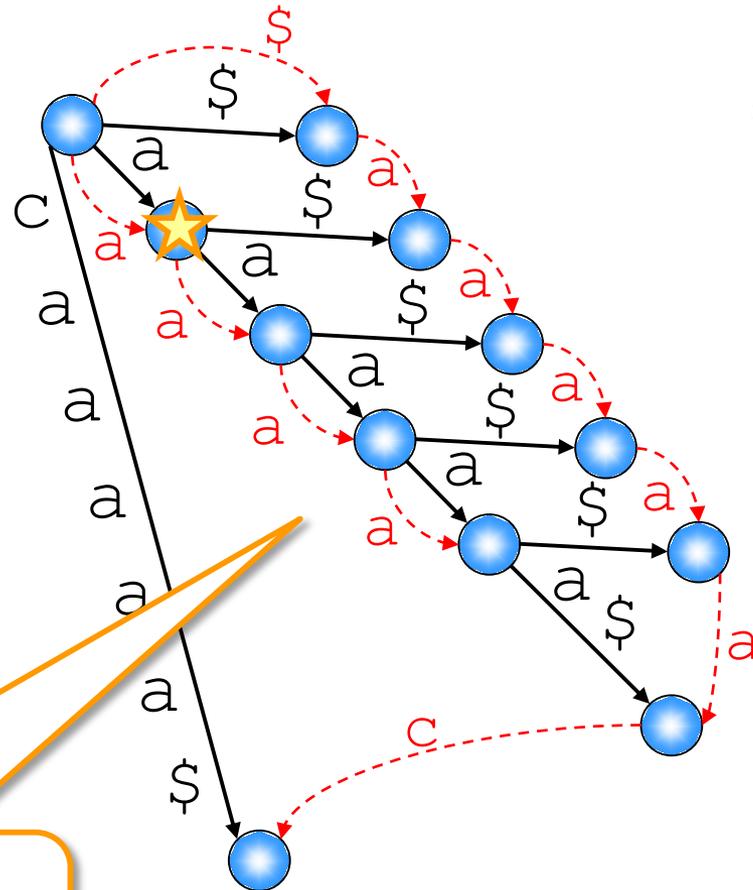
$\Omega(N^{1.5})$ 時間の証明



- caaaaaa\$ T_1
- caaaa\$ T_2
- ▲ aaa\$ T_3
- aa\$ T_4
- a\$ T_5

追加文字 c のリンクを
持つ最も深い頂点を探す

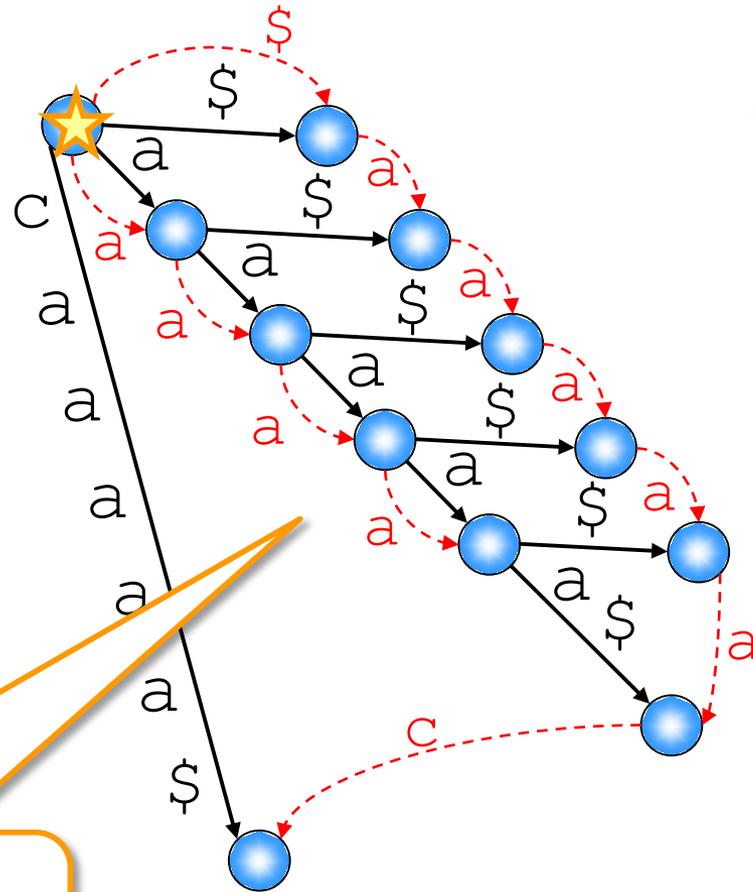
$\Omega(N^{1.5})$ 時間の証明



- caaaaaa\$ T_1
- caaaa\$ T_2
- ▲ aaa\$ T_3
- aa\$ T_4
- a\$ T_5

追加文字 c のリンクを持つ最も深い頂点を探す

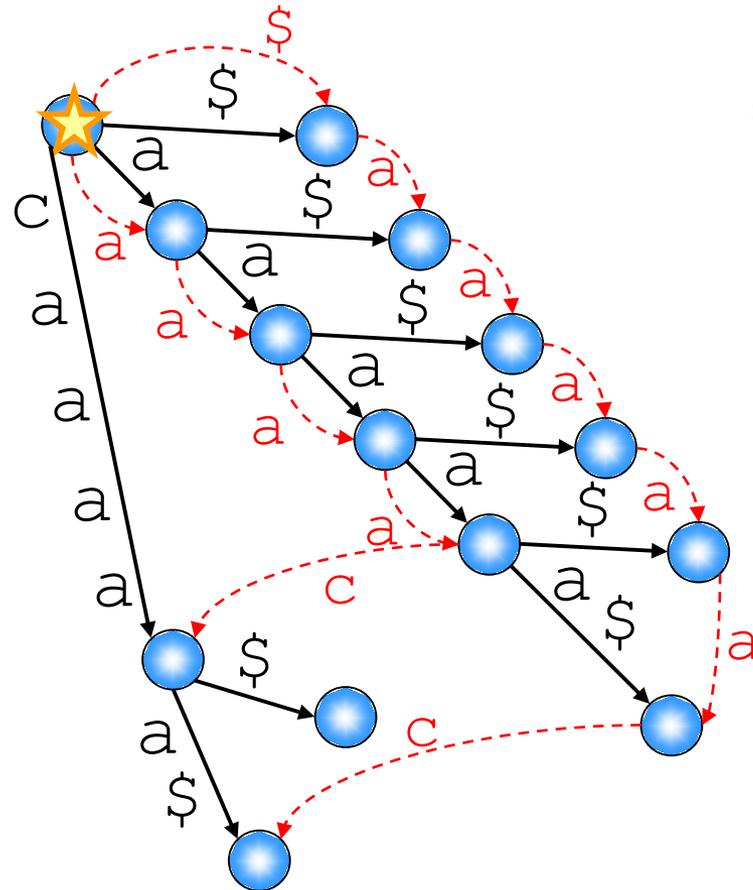
$\Omega(N^{1.5})$ 時間の証明



- caaaaaa\$ T_1
- caaaa\$ T_2
- ▲ aaa\$ T_3
- aa\$ T_4
- a\$ T_5

追加文字 c のリンクを持つ最も深い頂点を探す

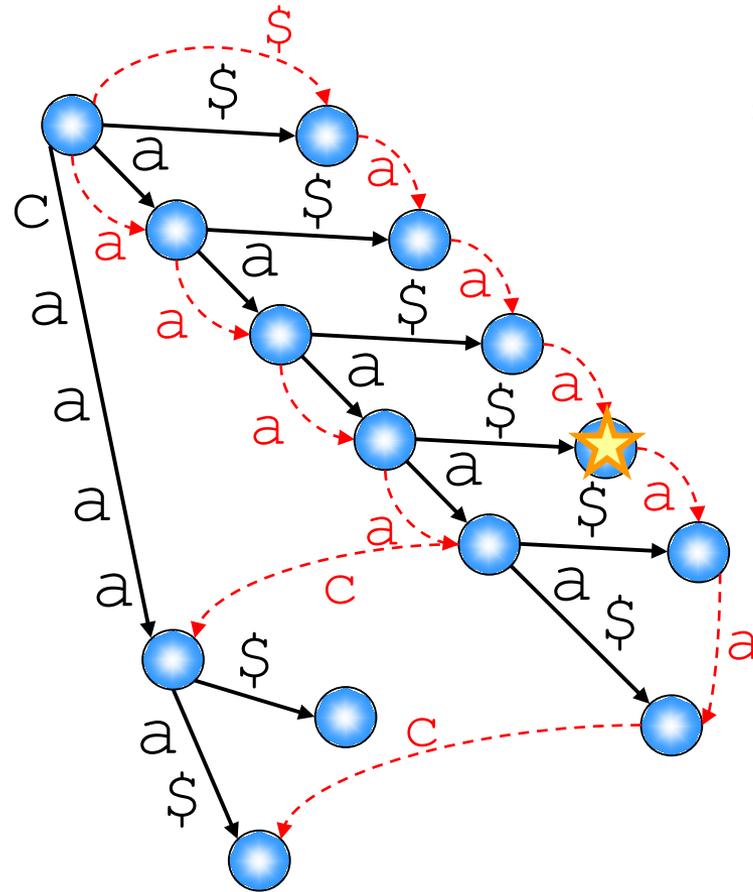
$\Omega(N^{1.5})$ 時間の証明



caaaaa\$ T_1
 caaaa\$ T_2
▲ aaa\$ T_3
 aa\$ T_4
 a\$ T_5

※ このステップに関係ないリンクは一部省略

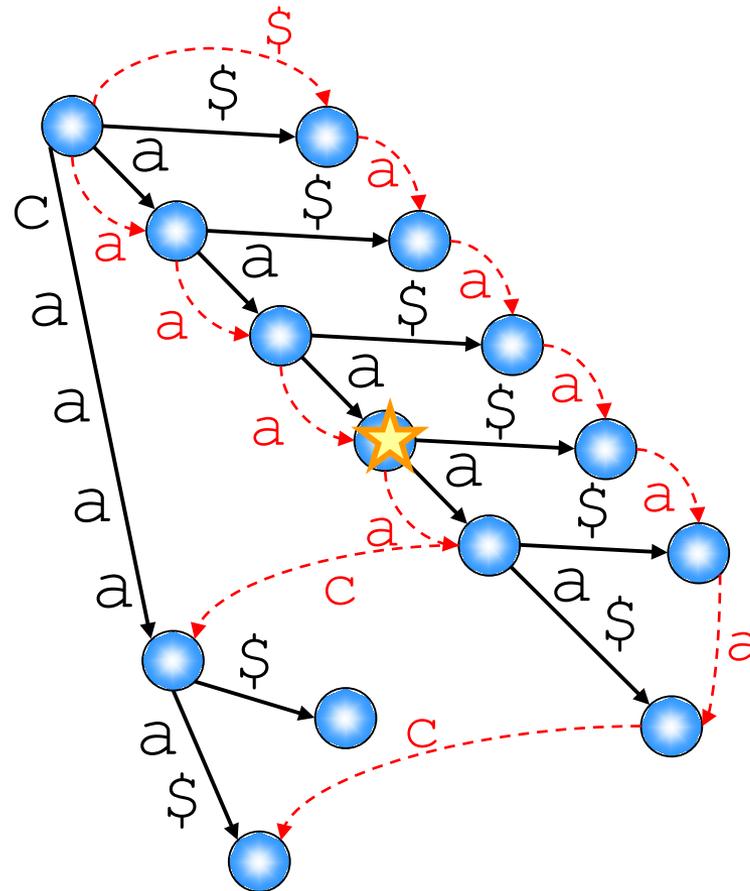
$\Omega(N^{1.5})$ 時間の証明



- caaaaa\$ T_1
- caaaa\$ T_2
- caaaa\$ T_3
- ▲ aa\$ T_4
- a\$ T_5

※ このステップに関係ないリンクは一部省略

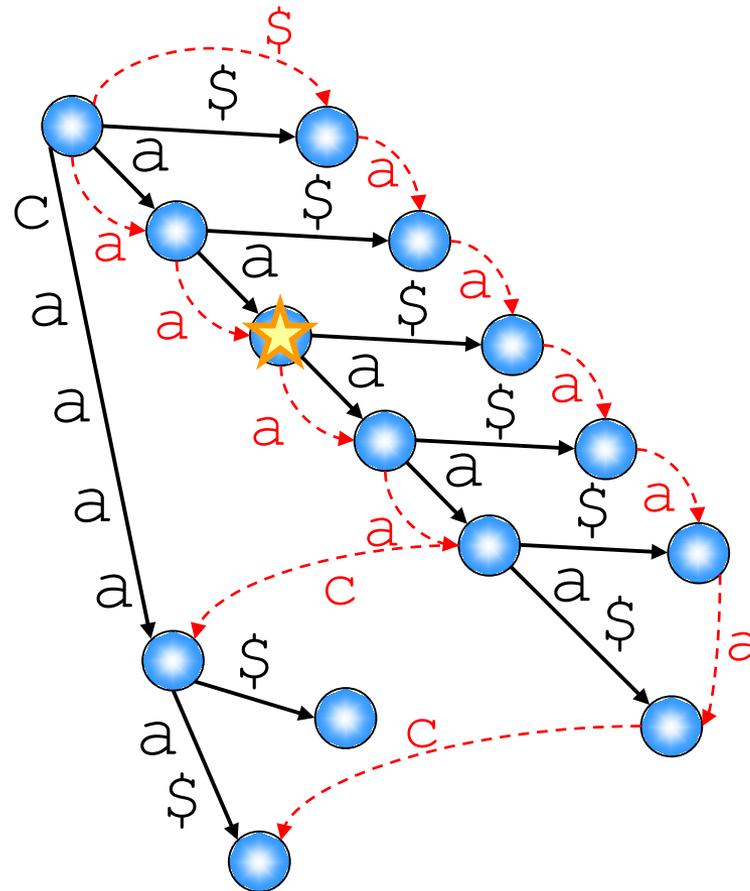
$\Omega(N^{1.5})$ 時間の証明



caaaaa\$ T_1
 caaaa\$ T_2
 caaaa\$ T_3
 ▲ aa\$ T_4
 a\$ T_5

※ このステップに関係ないリンクは一部省略

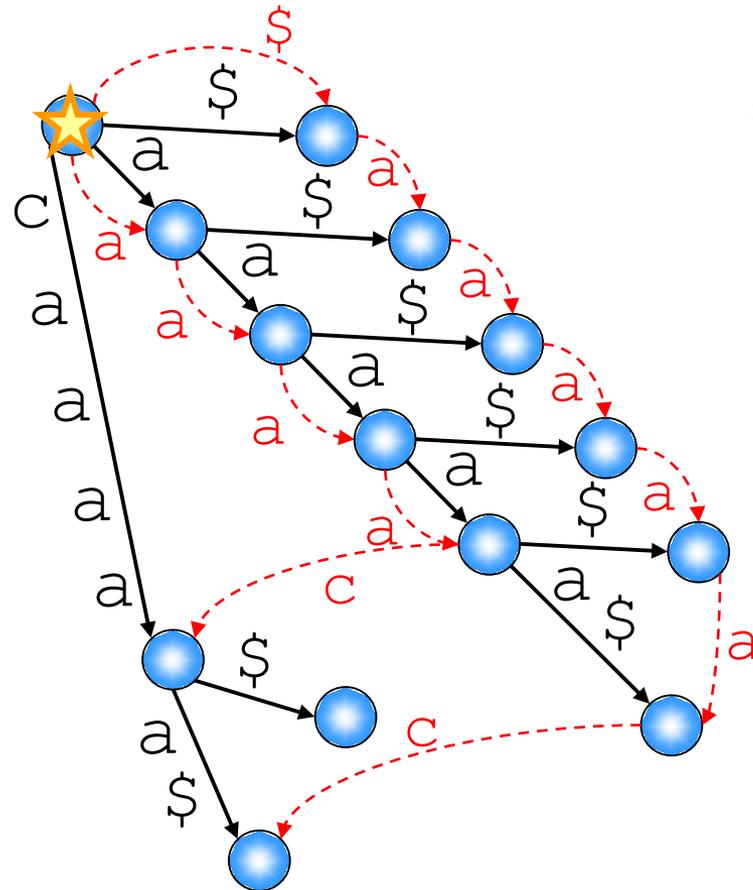
$\Omega(N^{1.5})$ 時間の証明



$c a a a a a \$$ T_1
 $c a a a a \$$ T_2
 $c a a a \$$ T_3
 $\triangle a a \$$ T_4
 $a \$$ T_5

※ このステップに関係ないリンクは一部省略

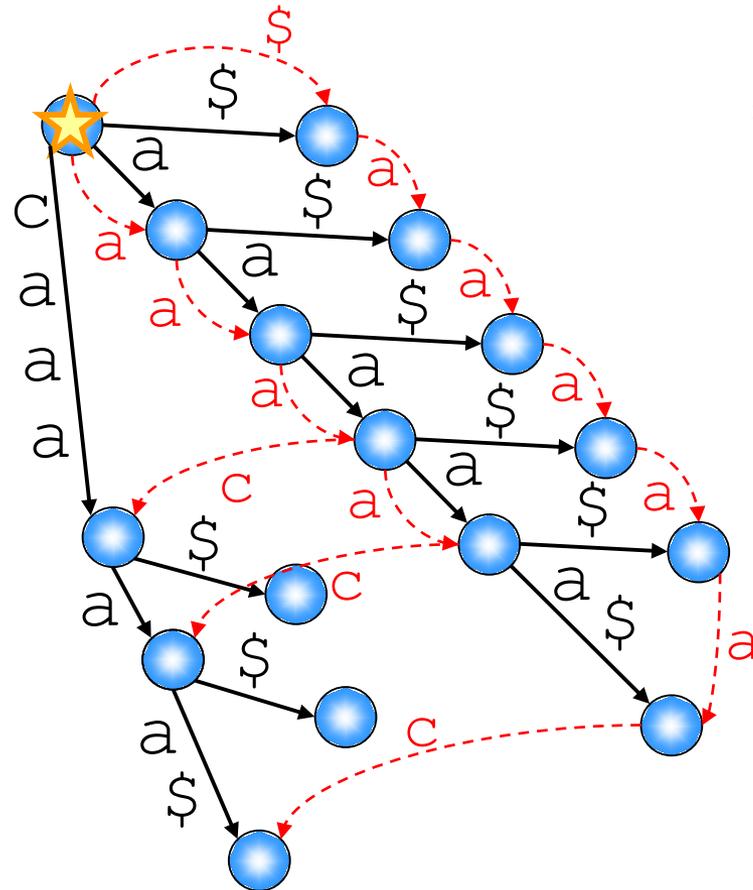
$\Omega(N^{1.5})$ 時間の証明



caaaaa\$ T_1
 caaaa\$ T_2
 caaaa\$ T_3
 ▲ aa\$ T_4
 a\$ T_5

※ このステップに関係ないリンクは一部省略

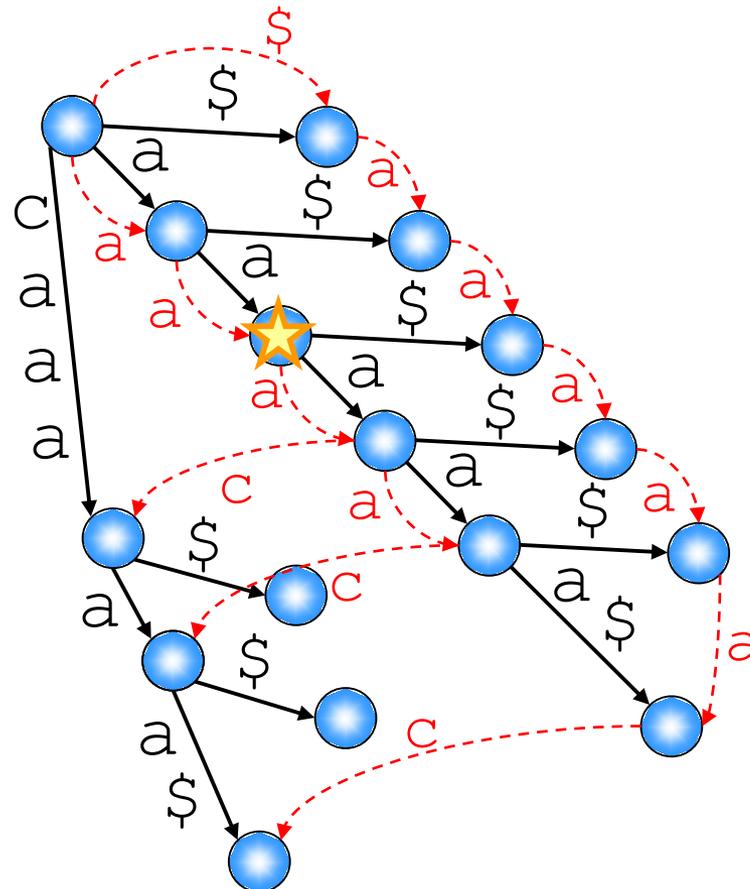
$\Omega(N^{1.5})$ 時間の証明



$caaaaa\$$ T_1
 $caaaa\$$ T_2
 $caaaa\$$ T_3
 \triangle $aaa\$$ T_4
 $a\$$ T_5

※ このステップに関係ないリンクは一部省略

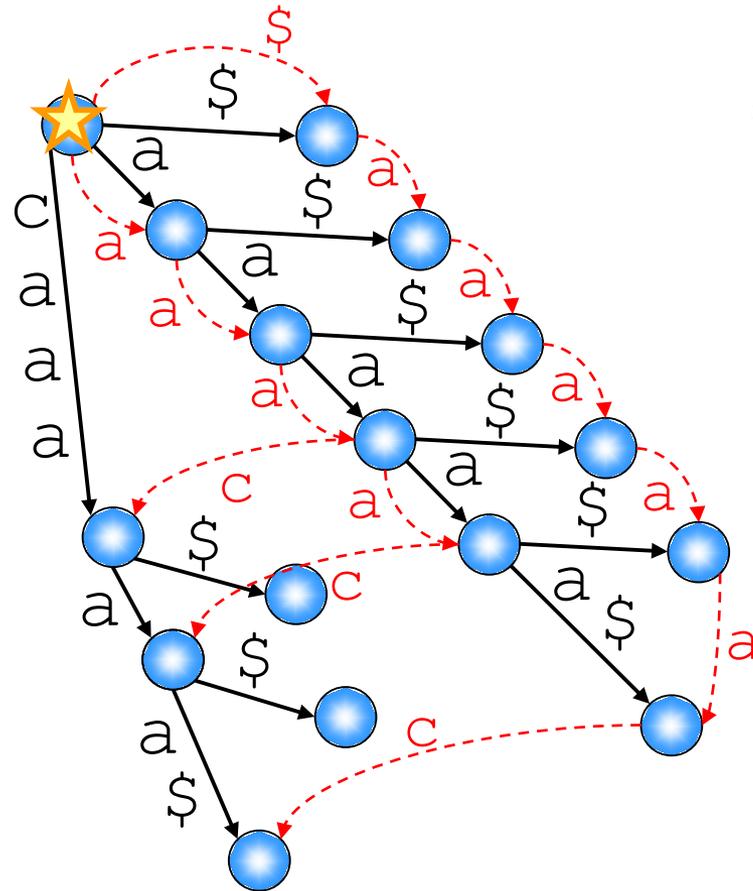
$\Omega(N^{1.5})$ 時間の証明



$caaaaa\$$ T_1
 $caaaa\$$ T_2
 $caaaa\$$ T_3
 \triangle $aa\$$ T_4
 $a\$$ T_5

※ このステップに関係ないリンクは一部省略

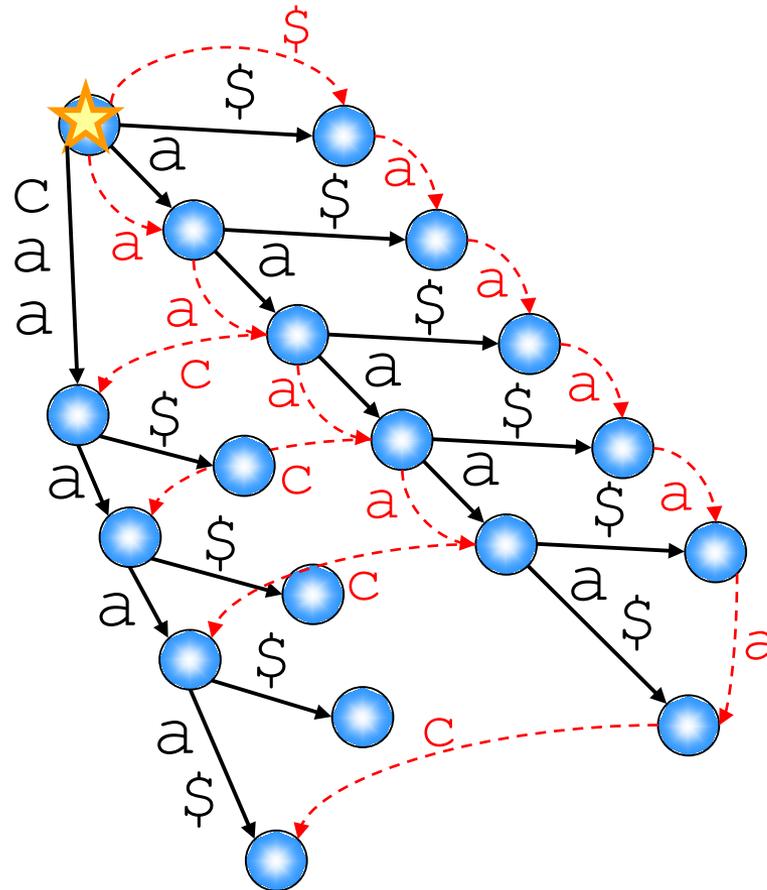
$\Omega(N^{1.5})$ 時間の証明



- caaaaa\$ T_1
- caaaa\$ T_2
- caaaa\$ T_3
- ▲ aa\$ T_4
- a\$ T_5

※ このステップに関係ないリンクは一部省略

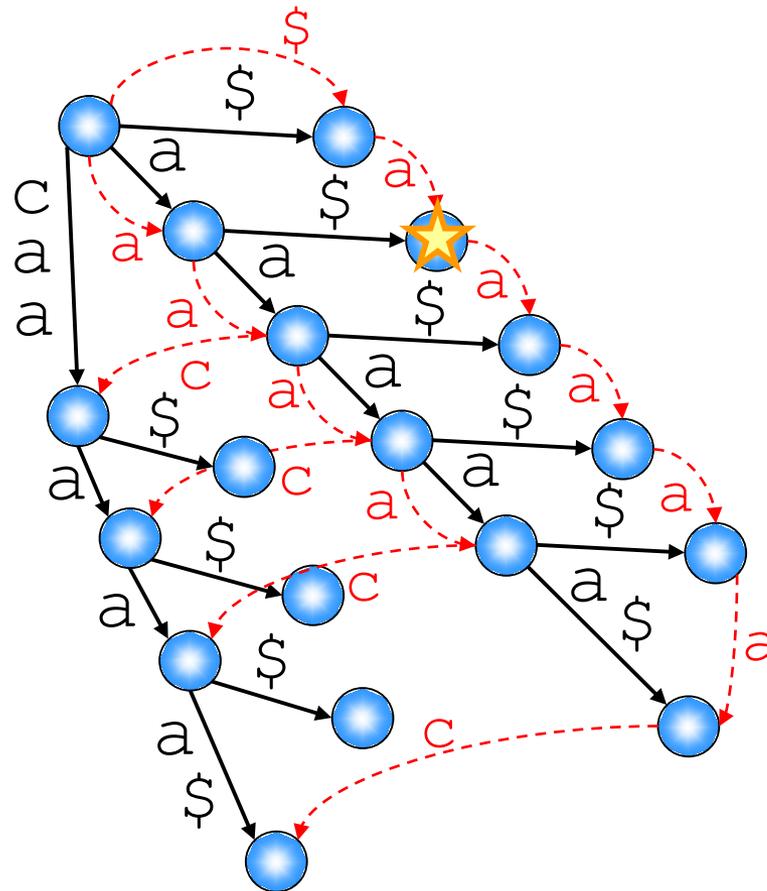
$\Omega(N^{1.5})$ 時間の証明



- caaaaaa\$ T_1
- caaaa\$ T_2
- caaaa\$ T_3
- caa\$ T_4
- ▲ a\$ T_5

※ このステップに関係ないリンクは一部省略

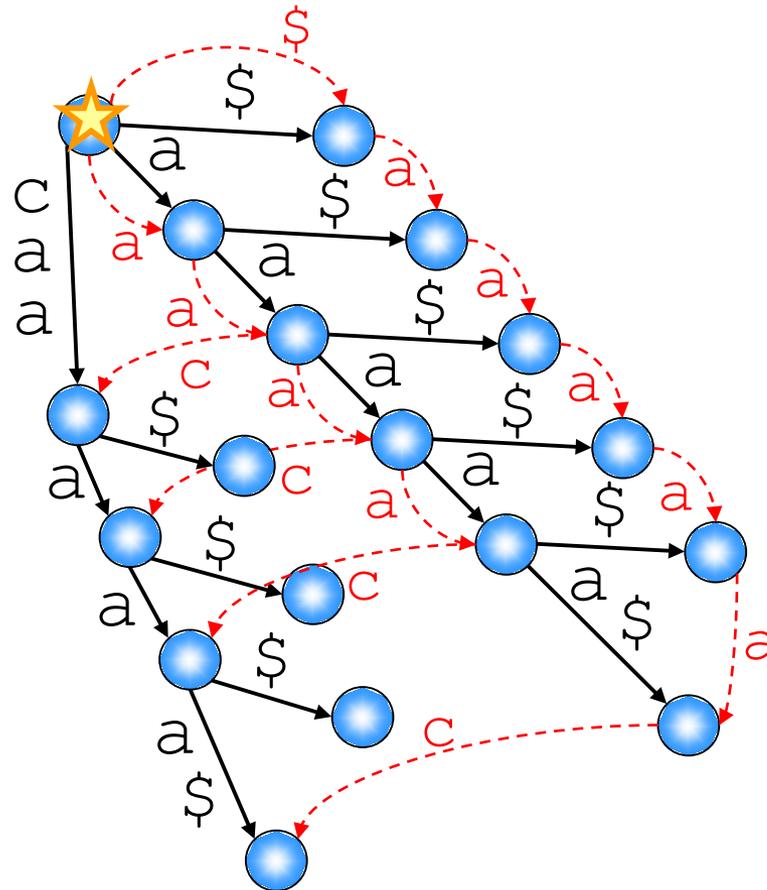
$\Omega(N^{1.5})$ 時間の証明



$caaaaaa\$$ T_1
 $caaaaa\$$ T_2
 $caaaa\$$ T_3
 $caa\$$ T_4
 $ca\$$ T_5
 \triangle

※ このステップに関係ないリンクは一部省略

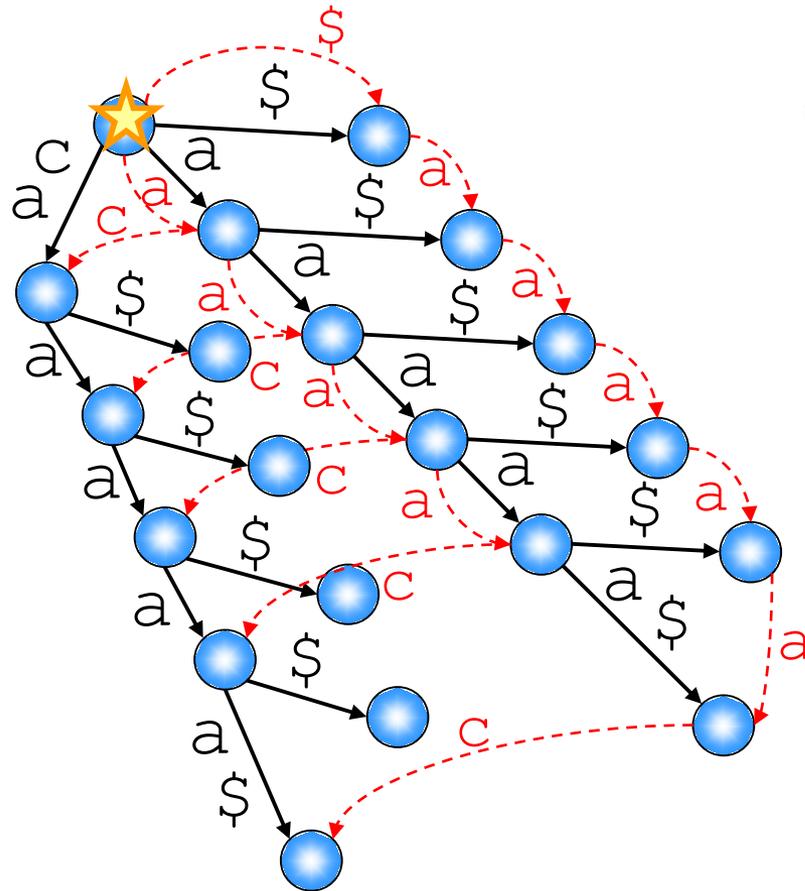
$\Omega(N^{1.5})$ 時間の証明



caaaaa\$ T_1
 caaaa\$ T_2
 caaa\$ T_3
 caa\$ T_4
 ca\$ T_5
 ▲

※ このステップに関係ないリンクは一部省略

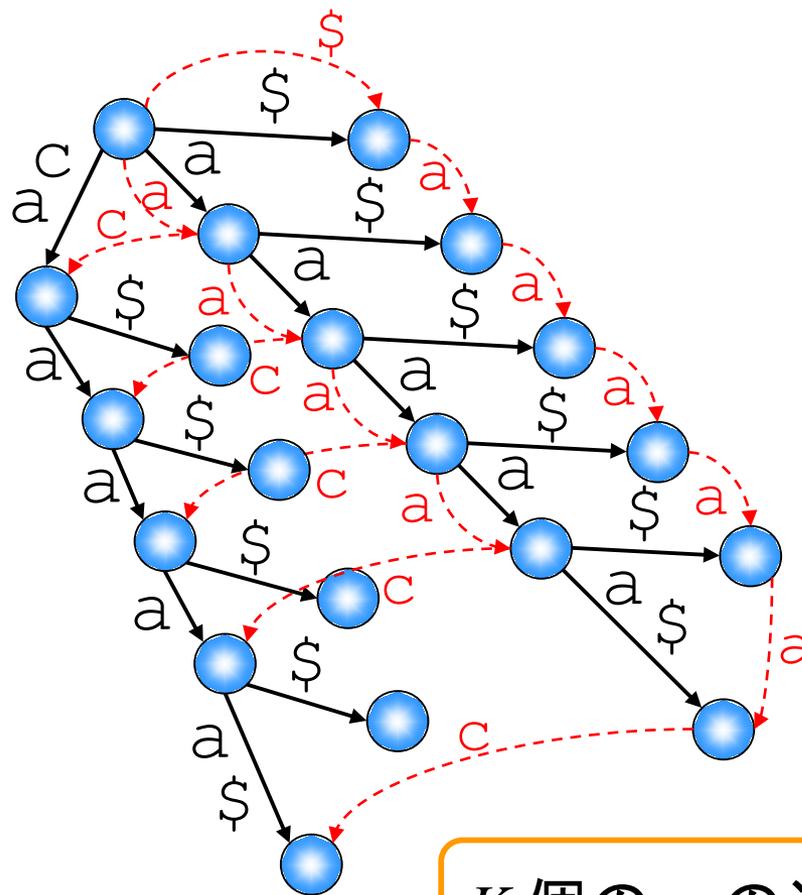
$\Omega(N^{1.5})$ 時間の証明



$caaaaa\$$ T_1
 $caaaa\$$ T_2
 $caaa\$$ T_3
 $caa\$$ T_4
 $ca\$$ T_5
 \triangle

※ このステップに関係ないリンクは一部省略

$\Omega(N^{1.5})$ 時間の証明



$\overbrace{\text{caaaaa}\K
 $\text{caaaa}\$$
 $\text{caaa}\$$
 $\text{caa}\$$
 $\text{ca}\$$

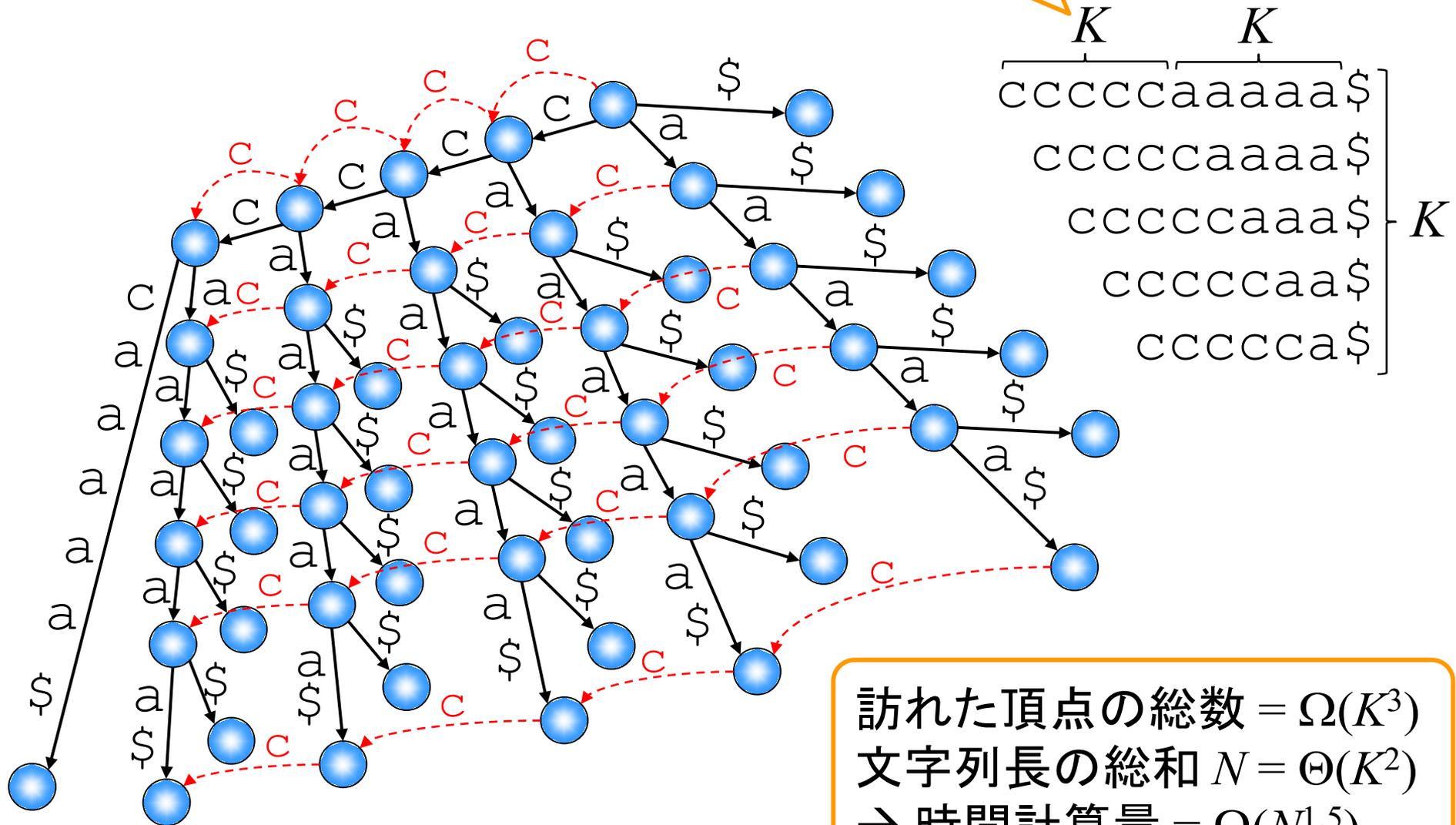
} K

K 個の c の追加に対して
訪れた頂点数は $\Omega(K^2)$

※ このステップに関係ないリンクは一部省略

$\Omega(N^{1.5})$ 時間の証明

先頭への c の追加を
 K ステップ繰り返す



右左オンライン構築～完全版

定理 2

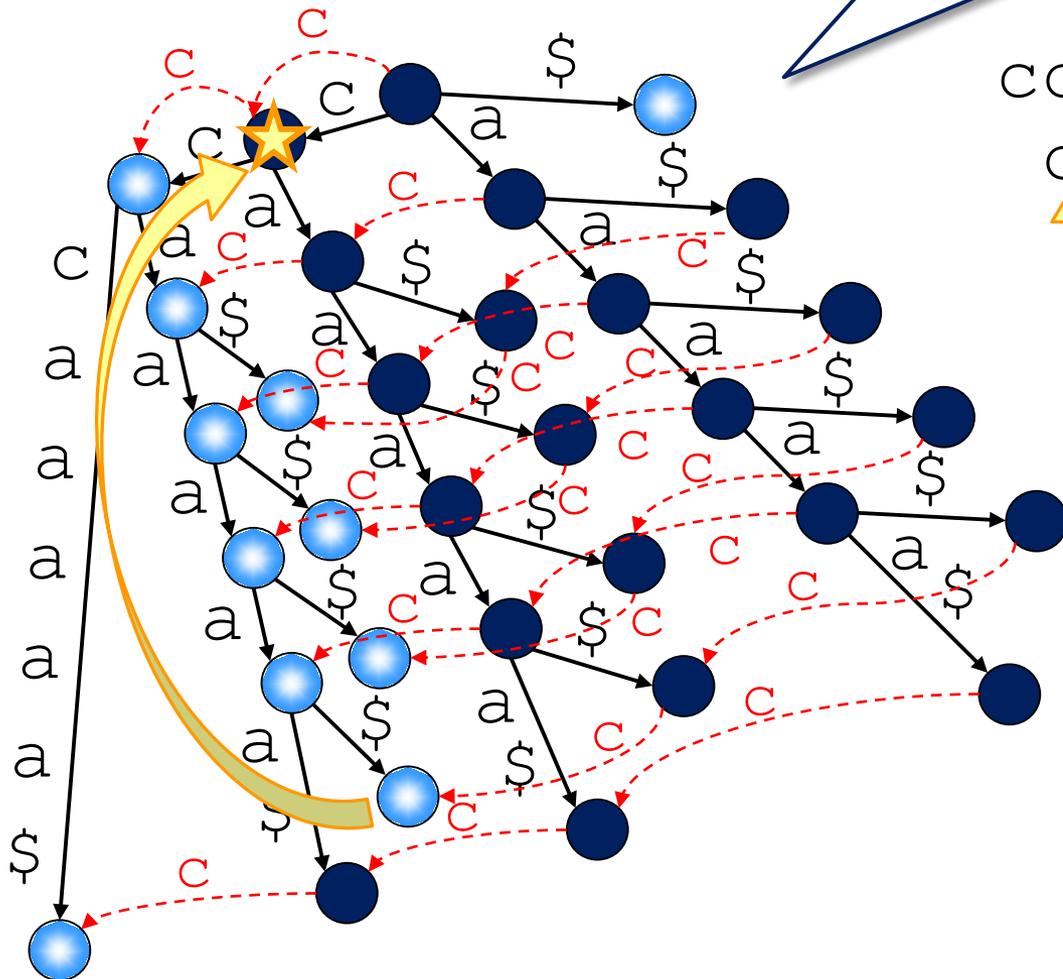
先頭に文字を追加していく複数文字列の接尾辞木を $O(N \log \sigma)$ 時間 $O(N)$ 領域で構築するアルゴリズムが存在する。

【テクニックの概要】

- ◆ 1文字追加につき $\Omega(K)$ 個訪れていた頂点を $O(1)$ 時間でジャンプできればよい。
- ◆ このジャンプを木の最近マーク祖先 (NMA) 問題に帰着。
- ◆ NMA のデータ構造を愚直に用いると $O(\sigma N)$ 領域必要だが、注意深く選んだ頂点集合にのみ用いると $O(N)$ 領域に収まる。

NMAによる高速化

頂点 v がマーク頂点 ●
 $\Leftrightarrow v$ は c のリンクを持つ



- cccaaaaa\$ T_1
- cccaaaa\$ T_2
- ▲ ccaaa\$ T_3
- ccaa\$ T_4
- cca\$ T_5

左右オンライン構築

定理 4

Ukkonen のアルゴリズムをそのまま
複数文字列に適用すると $\Omega(N^2)$ 時間かかる.

定理 5

末尾に文字を追加していく**複数文字列**の接尾辞木を
 $O(N \log \sigma)$ 時間 $O(N)$ 領域で構築するアルゴリズムが
存在する.

※ ただし, 定理5の接尾辞木では葉のラベルを陽に管理しない.
(葉に達する場合のパターン照合は別の方法で行う)

本研究にまつわるヒストリー

2015/5	研究開始(高木, 稲永, 有村)
2016/7	文字列処理に関する国際会議 CPM で発表
2017/1	Dany から3人にメールが届き, CPM の論文に致命的なエラーが見つかる
2017/2~3	Dany と Diptarama が研究に参加
2017/7	すべての問題を解決!
2017/12/20	第1版の執筆が完了し, 共著者にメール
2017/12/21	Dany が急逝
2018/3	Algorithmica に投稿
2019/10	Algorithmica に採択/掲載

Meeting with Dany @北海道大学 (2017/6)



【Dany Breslauer 氏の功績に敬意を表します】