

**THE COMPLEXITY AND APPLICATION
OF SYNTACTIC PATTERN
RECOGNITION USING FINITE
INDUCTIVE STRINGS**

**Elijah Myers
Paul S. Fisher
Keith Irwin
Jinsuk Baek
Joao Setubal**

FINITE INDUCTIVE STRINGS

- A finite inductive string is one that:
 - Consists of symbols selected from a finite alphabet of possible symbols
 - Eventually terminates and is therefore finite in length
 - Can be characterized using inductive reasoning (i.e. observations of the specific occurrences of symbols can be applied generally to define the entire string)

FACTORING

- Given a sample string: aacacgacgt
 - Append the start symbol: Saacacgacgt
 - Create a storage structure (ruling) that uniquely identifies each symbol

| Ruling | |
|---------------------|----------------------|
| $S \rightarrow a$ | $cac \rightarrow g$ |
| $Sa \rightarrow a$ | $cacg \rightarrow a$ |
| $aa \rightarrow c$ | $ga \rightarrow c$ |
| $aac \rightarrow a$ | $gac \rightarrow g$ |
| $ca \rightarrow c$ | $gacg \rightarrow t$ |

FACTORING

- Given the residual string: agagt
 - Append the start symbol: Sagagt
 - Create a second ruling that uniquely identifies each symbol

| Level 0 Ruling | |
|--------------------|--------------------|
| $S \rightarrow a$ | $ca \rightarrow c$ |
| $Sa \rightarrow a$ | $ga \rightarrow c$ |
| $aa \rightarrow c$ | |

| Level 1 Ruling | |
|-------------------|---------------------|
| $S \rightarrow a$ | $ag \rightarrow a$ |
| $a \rightarrow g$ | $gag \rightarrow t$ |

FACTORING

- Given the residual string: t
 - Append the start symbol: St
 - Create a third ruling that uniquely identifies each symbol

| Level 0 Ruling | |
|--------------------|--------------------|
| $S \rightarrow a$ | $ca \rightarrow c$ |
| $Sa \rightarrow a$ | $ga \rightarrow c$ |
| $aa \rightarrow c$ | |

| Level 1 Ruling | |
|-------------------|--------------------|
| $S \rightarrow a$ | $ag \rightarrow a$ |
| $a \rightarrow g$ | |

| Level 2 Ruling | |
|-------------------|--|
| $S \rightarrow t$ | |

FOLLOWING

- Given a sample string: aacacgacat
 - Append the start symbol: Saacacgacat
 - Apply the first ruling to identify those elements that do not conform to the ruling

| Level 0 Ruling | |
|--------------------|--------------------|
| $S \rightarrow a$ | $ca \rightarrow c$ |
| $Sa \rightarrow a$ | $ga \rightarrow c$ |
| $aa \rightarrow c$ | |

FOLLOWING

- Given the residual string: agacat
 - Append the start symbol: Sagacat
 - Apply the second ruling to identify those elements that do not conform to the ruling

| Level 1 Ruling | |
|-------------------|--------------------|
| $S \rightarrow a$ | $ag \rightarrow a$ |
| $a \rightarrow g$ | |

FOLLOWING

- Given the residual string: cat
 - Append the start symbol: S**cat**
 - Apply the third ruling to identify those elements that do not conform to the ruling

Level 2 Ruling

S → t

- With all rulings applied, the final residual “cat” is the portion of the string that did not match the pattern of the factored string

PERFORMANCE

○ Definition of Terms

- b – number of symbols in the alphabet (i.e. base)
- IB – inductive base
- N – length of the examined string (number of characters)
- L – number of levels factored or followed
- R – number of rules in a given ruling level

PERFORMANCE

○ Factoring

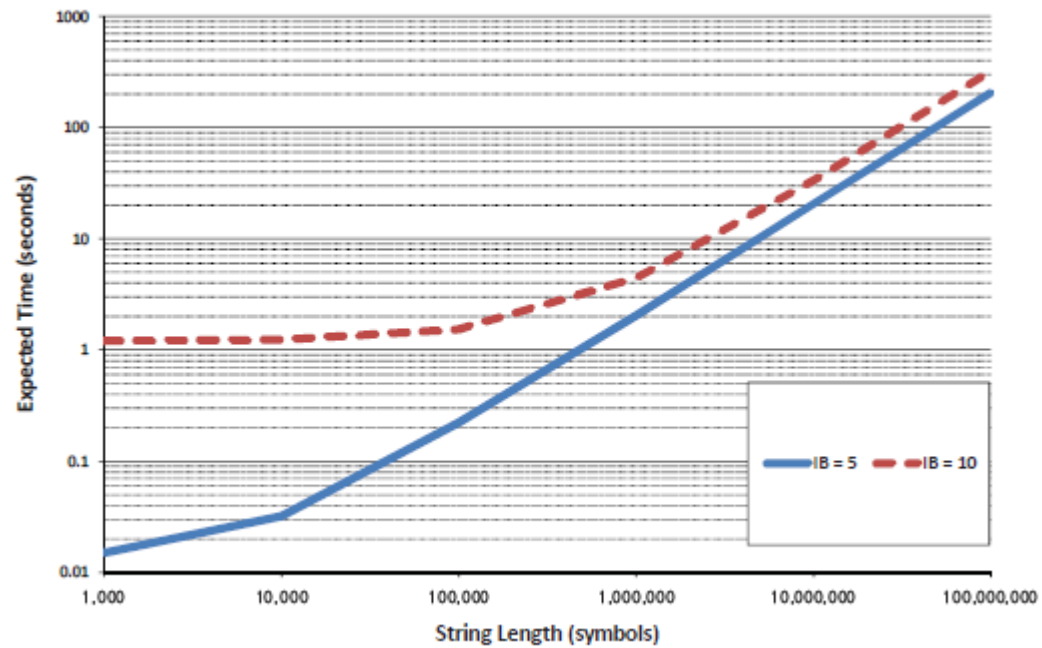
- $O(b + L[b^{\text{IB}} + (\text{IB} - 1)(b^{\text{IB}}) + 3N])$
- $O(b + L[\text{IB}b^{\text{IB}} + 3N])$
- $O(L + LN)$

PERFORMANCE

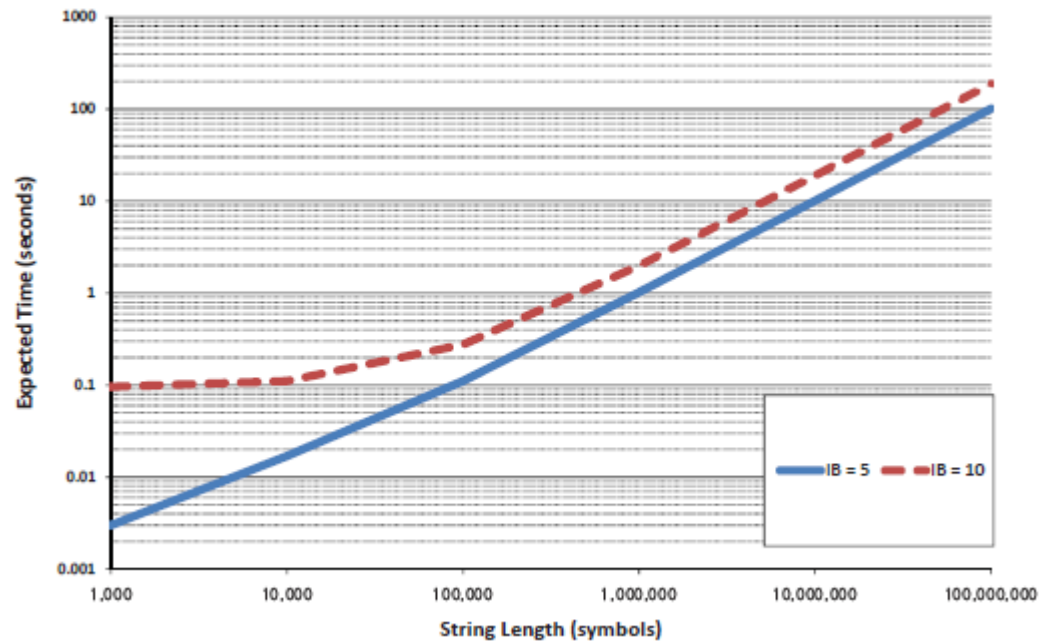
- Following

- $O(b + L[R + 2b^{\text{IB}} + N])$
- $O(b + L[3b^{\text{IB}} + N])$
- $O(L + LN)$

FACTORING RESULTS



FOLLOWING RESULTS



SUMMARY

- Factoring provides a means to discover the underlying patterns in a string of genetic data
- Following provides a means to compare any unknown string to one that has been factored previously to determine convergence/divergence
- Both processes have been determined to perform in linear time with respect to the length of the input strings