

Encoding Co-Lex Orders of Finite State-Automata in Linear Space

36th Annual Symposium on Combinatorial Pattern Matching

R. Becker¹, N. Cotumaccio², S. Kim¹,
N. Prezza¹, **C. Tosoni**¹

1. Ca' Foscari University of Venice
2. University of Helsinki

Milan, Bicocca University

17/06/2025

The Burrows-Wheeler transform

The **BWT** is a famous reversible string transformation invented by Burrows and Wheeler [1].

- ① **Enhance the compressibility** of strings.
- ② Allows the implementation of **indexes for pattern matching**.

$$\textit{banana\$} \longrightarrow \textit{annb\$aa}$$

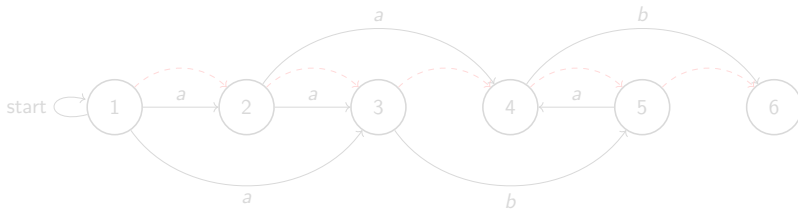
1. Burrows M., Wheeler D.: A block-sorting lossless data compression algorithm. SRS Research Report (1994)

Extending the BWT to NFAs

The BWT has been extended to **nondeterministic finite automata** (NFAs) using **Wheeler orders** [2].

Wheeler order

Total order \leq of nodes **consistent with the strings** reaching the states.



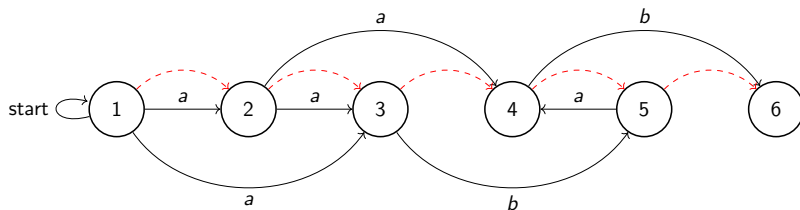
2. Gagie et al.: *Wheeler graphs: A framework for BWT-based data structures*. *Theor. Comput. Sci.* (2017)

Extending the BWT to NFAs

The BWT has been extended to **nondeterministic finite automata** (NFAs) using **Wheeler orders** [2].

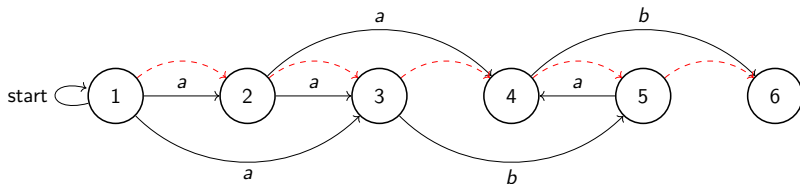
Wheeler order

Total order \leq of nodes **consistent with the strings** reaching the states.



2. Gagie et al.: *Wheeler graphs: A framework for BWT-based data structures*. *Theor. Comput. Sci.* (2017)

Extending the BWT to NFAs



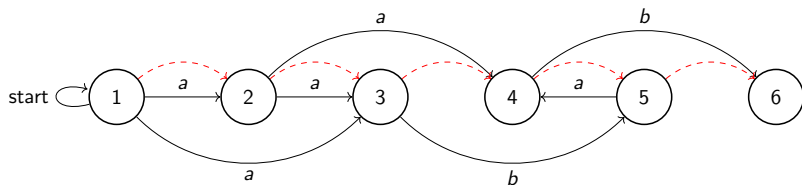
What does “**consistent**” mean?

- **Initial state** is the **first node** of \leq (node u_1)

Strings reaching u_3 and $u_4 \rightarrow l_{u_3} = \{a, aa\}$ and $l_{u_4} = \{aa, aba, aaba\}$

- $u_3 \leq u_4 \implies$ excluding the strings in common $\{aa\}$, strings in l_{u_3} **are smaller than those in** l_{u_4} .

Extending the BWT to NFAs



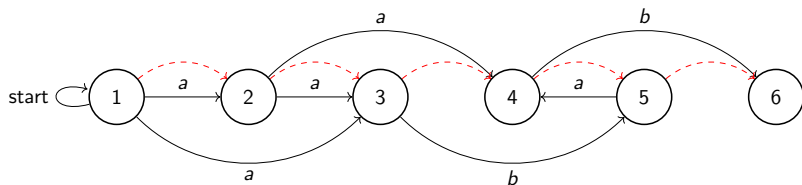
What does “**consistent**” mean?

- **Initial state** is the **first node** of \leq (node u_1)

Strings reaching u_3 and $u_4 \rightarrow I_{u_3} = \{a, aa\}$ and $I_{u_4} = \{aa, aba, aaba\}$

- $u_3 \leq u_4 \implies$ excluding the strings in common $\{aa\}$, strings in I_{u_3} **are smaller than those in I_{u_4} .**

Extending the BWT to NFAs



What does “**consistent**” mean?

- **Initial state** is the **first node** of \leq (node u_1)

Strings reaching u_3 and $u_4 \rightarrow I_{u_3} = \{a, aa\}$ and $I_{u_4} = \{aa, aba, aaba\}$

- $u_3 \leq u_4 \implies$ excluding the strings in common $\{aa\}$, strings in I_{u_3} **are smaller than those in** I_{u_4} .

Wheeler orders: pros and cons

Pros

- NFAs admitting a Wheeler order can be **efficiently compressed and indexed**.

Cons

- Most NFAs do not admit a Wheeler order.
- Recognizing if this is the case is an **NP-complete problem** [3]!

3. Gibney, Thankachan: *On the Hardness and Inapproximability of Recognizing Wheeler Graphs*. ESA. (2019)

Wheeler orders: pros and cons

Pros

- NFAs admitting a Wheeler order can be **efficiently compressed and indexed**.

Cons

- Most NFAs do not admit a Wheeler order.
- Recognizing if this is the case is an **NP-complete problem** [3]!

3. Gibney, Thankachan: *On the Hardness and Inapproximability of Recognizing Wheeler Graphs*. ESA. (2019)

Wheeler orders: pros and cons

Pros

- NFAs admitting a Wheeler order can be **efficiently compressed and indexed**.

Cons

- Most NFAs do not admit a Wheeler order.
- Recognizing if this is the case is an **NP-complete problem** [3]!

3. Gibney, Thankachan: *On the Hardness and Inapproximability of Recognizing Wheeler Graphs*. ESA. (2019)

CFS orders

To address these issues → **CFS orders!** [4]

Coarsest Forward-Stable co-lex (CFS) orders

A CFS order \leq_{FS} is a **partial preorder** on an NFA's states **consistent with the strings reaching them**.

State-Of-The-Art for **indexing** and **compressing** NFAs!

- **Exists** and is **unique** for each NFA.
- **Polynomial time** to compute it.

4. Becker et al.: *Indexing Finite-State Automata Using Forward-Stable Partitions*. SPIRE (2024)

CFS orders

To address these issues → **CFS orders!** [4]

Coarsest Forward-Stable co-lex (CFS) orders

A CFS order \leq_{FS} is a **partial preorder** on an NFA's states **consistent with the strings reaching them**.

State-Of-The-Art for **indexing** and **compressing** NFAs!

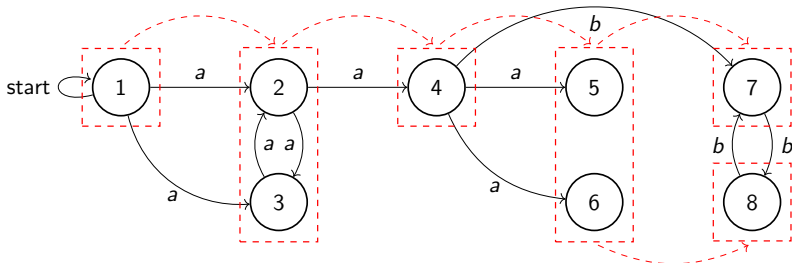
- **Exists** and is **unique** for each NFA.
- **Polynomial time** to compute it.

4. Becker et al.: *Indexing Finite-State Automata Using Forward-Stable Partitions*. SPIRE (2024)

CFS orders

Definition of CFS order \leq_{FS}

- 1) Compute the **quotient NFA** defined by **coarsest forward stable partition**.
- 2) For every **transitions** $v \xrightarrow{a} u, v' \xrightarrow{a'} u'$ in the **quotient NFA**:
 - $a < a' \implies u < u'$
 - $(a = a') \wedge (v < v') \implies u \leq u'$



Complexity of CFS Orders

The CFS order can be computed in **polynomial time**, more precisely in **$O(m^2)$** time, with **m** being the number of **transitions in the NFA**.

- Unfeasible for **big data applications** (i.e., to index **pagename graphs**)

Main open problem: devise a (near) **linear time algorithm** :D

Complexity of CFS Orders for DFAs

For the special case of **DFAs**, a **$O(m \log n)$** time algorithm to compute CFS orders **is already known**. [5]

This algorithm was preceded by a **$O(n)$ linear space representation** to encode them. [6]

$m \rightarrow$ Number of transitions.

$n \rightarrow$ Number of states.

5. Becker et al.: *Sorting Finite Automata via Partition Refinement*. *ESA*. (2023)

6. Kim et al.: *Faster Prefix-Sorting Algorithms for Deterministic Finite Automata*. *CPM*. (2023)

Complexity of CFS Orders for DFAs

For the special case of **DFAs**, a **$O(m \log n)$** time algorithm to compute CFS orders **is already known**. [5]

This algorithm was preceded by a **$O(n)$ linear space representation** to encode them. [6]

$m \rightarrow$ Number of transitions.

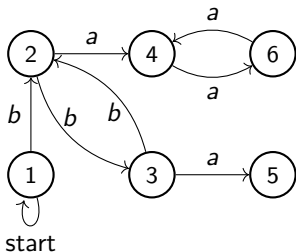
$n \rightarrow$ Number of states.

5. Becker et al.: *Sorting Finite Automata via Partition Refinement*. *ESA*. (2023)

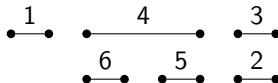
6. Kim et al.: *Faster Prefix-Sorting Algorithms for Deterministic Finite Automata*. *CPM*. (2023)

Complexity of CFS Orders for DFAs

Indeed, **in DFAs** the CFS order \leq_{FS} can be **encoded using intervals**



Interval order of \leq_{FS}



However, this is **not possible in the case of nondeterminism**

Our main result

Theorem

There exists an **$O(n)$ representation** of \leq_{FS}

- Arbitrary partial preorders on V ($|V| = n$) require **$\Omega(n^2)$ bits to be represented** (binary matrix)
- This special class can be encoded in just linear space

Step towards a subquadratic algorithm?

Our main result

Theorem

There exists an **$O(n)$ representation** of \leq_{FS}

- Arbitrary partial preorders on V ($|V| = n$) require **$\Omega(n^2)$ bits to be represented** (binary matrix)
- This special class can be encoded in just linear space

Step towards a subquadratic algorithm?

Our main result

Theorem

There exists an **$O(n)$ representation** of \leq_{FS}

- Arbitrary partial preorders on V ($|V| = n$) require **$\Omega(n^2)$ bits to be represented** (binary matrix)
- This special class can be encoded in just linear space

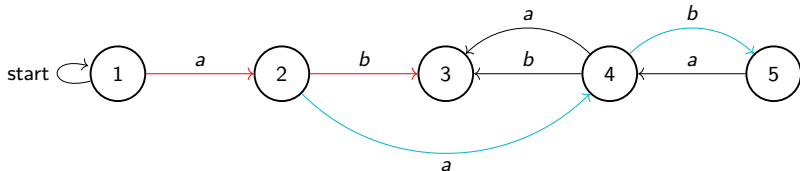
Step towards a subquadratic algorithm?

Characterizing \leq_{FS}

Preceding pairs [7]

Let \mathcal{A} be an NFA and Q its set of states.

Consider $(w, z), (u, v) \in Q \times Q$, we say that (w, z) **precedes** (u, v) , denoted $(w, z) \Rightarrow (u, v)$, if there exist α such that $w \xrightarrow{\alpha} u$ and $z \xrightarrow{\alpha} v$



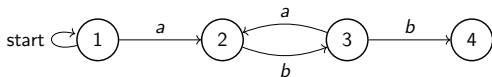
Since $1 \xrightarrow{a} 2 \xrightarrow{b} 3$ and $2 \xrightarrow{a} 4 \xrightarrow{b} 5$, we have $(1, 2) \Rightarrow (3, 5)$.

7. Cotumaccio N.: *Graphs can be succinctly indexed for pattern matching in $O(|E|^2 + |V|^{5/2})$ time.* DCC. (2022)

Characterizing \leq_{FS}

input-consistency: Incoming transitions are **labeled with same character**.

We denote with $\lambda(u)$ the label of the incoming edges of u .



Ex. $\lambda(3) = b$

Corollary

For every states $u, v \in Q$;

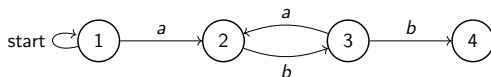
$u \leq_{FS} v \iff$ for each **preceding pair** $(w, z) \Rightarrow (u, v)$ we have $\lambda(w) \leq \lambda(z)$.

Ex. $\neg(4 \leq_{FS} 3)$ because $(3, 2) \Rightarrow (4, 3)$ and $\lambda(3) > \lambda(2)$

Characterizing \leq_{FS}

input-consistency: Incoming transitions are **labeled with same character**.

We denote with $\lambda(u)$ the label of the incoming edges of u .



Ex. $\lambda(3) = b$

Corollary

For every states $u, v \in Q$;

$u \leq_{FS} v \iff$ for each **preceding pair** $(w, z) \Rightarrow (u, v)$ we have $\lambda(w) \leq \lambda(z)$.

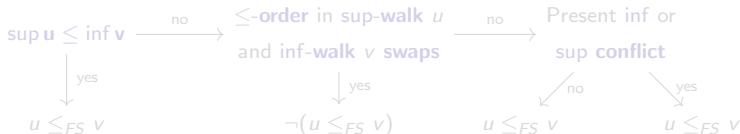
Ex. $\neg(4 \leq_{FS} 3)$ because $(3, 2) \Rightarrow (4, 3)$ and $\lambda(3) > \lambda(2)$

Overview

Encoding \leq_{FS} : $O(n)$ space

- **Linear Extension** \leq
- Lexicographically smallest (largest) string **inf** u (**sup** u) and its respective walk **Inf-walk** u (**Sup-walk** u)
- **Inf** (**Sup**) **Conflicts** of u

If $u < v$:

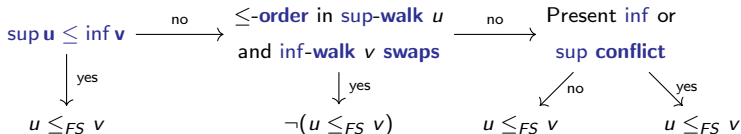


Overview

Encoding \leq_{FS} : $O(n)$ space

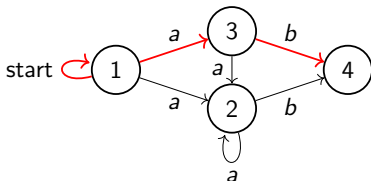
- **Linear Extension** \leq
- Lexicographically smallest (largest) string **inf** u (**sup** u) and its respective walk **Inf-walk** u (**Sup-walk** u)
- **Inf** (**Sup**) **Conflicts** of u

If $u < v$:

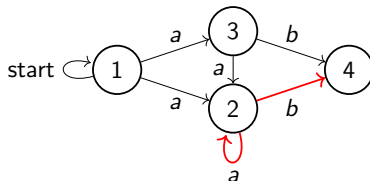


Infima/Suprema

The **Infima** (**Suprema**) are the **smallest** (**largest**) strings reaching the states from the initial state. **We can**: compute them in $O(m \log n)$ [5] stores them in $O(n)$ [6]



$\text{inf } 4 = ba\#\#\#\# \dots = ba\#\omega$
an **infimum walk**: $1 \xrightarrow{\#} 1 \xrightarrow{a} 3 \xrightarrow{b} 4$



$\text{sup } 4 = baaaaa \dots = ba\omega$
a **supremum walk**: $2 \xrightarrow{a} 2 \xrightarrow{b} 4$

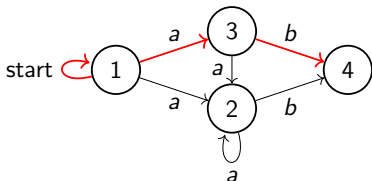
Lemma

For every two states u, v , $\text{sup } u \leq \text{inf } v \implies u \leq_{FS} v$

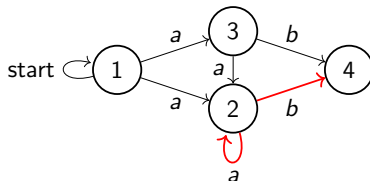
5. Sorting Finite Automata via Partition Refinement. 6. Faster Prefix-Sorting Algorithms for DFAs.

Infima/Suprema

The **Infima** (**Suprema**) are the **smallest** (**largest**) strings reaching the states from the initial state. **We can**: compute them in $O(m \log n)$ [5] stores them in $O(n)$ [6]



$\text{inf } 4 = ba\#\#\#\# \dots = ba\#\omega$
 an **infimum walk**: $1 \xrightarrow{\#} 1 \xrightarrow{a} 3 \xrightarrow{b} 4$



$\text{sup } 4 = baaaaa \dots = ba\omega$
 a **supremum walk**: $2 \xrightarrow{a} 2 \xrightarrow{b} 4$

Lemma

For every two states u, v , $\text{sup } u \leq \text{inf } v \implies u \leq_{FS} v$

5. Sorting Finite Automata via Partition Refinement. 6. Faster Prefix-Sorting Algorithms for DFAs.

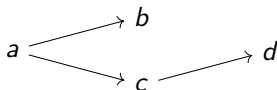
Linear Extension

Second part of our encoding: a linear extension! ($O(n)$ space)

Linear Extension

A total order \leq is a **linear extension** of \leq_{FS} if:

$$u \leq_{FS} v \implies u \leq v$$



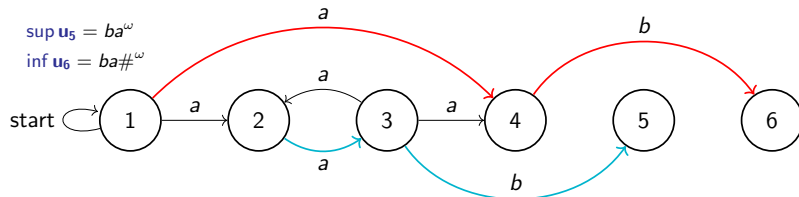
possible **partial order**



possible **linear extension**

The linear extension allows us **to reconstruct** \leq_{FS} when $\sup u > \inf v$.

Inf/Sup-Walks Cross



In this case the **linear extension** is $1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5 \rightarrow 6$.

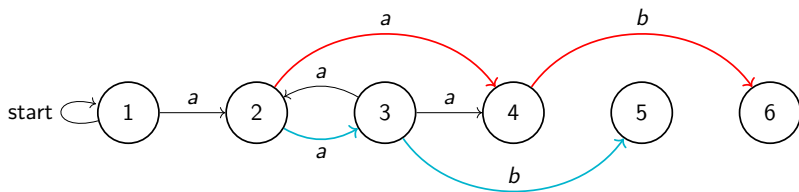
We have $\sup 5 > \inf 6$

Supremum walk to 5: $2 \xrightarrow{a} 3 \xrightarrow{b} 5$

Infimum walk to 6: $1 \xrightarrow{a} 4 \xrightarrow{b} 6$

The order swaps!: $2 > 1 \implies 5 \leq_{FS} 6$ does not hold.

Inf/Sup-Walks Meet

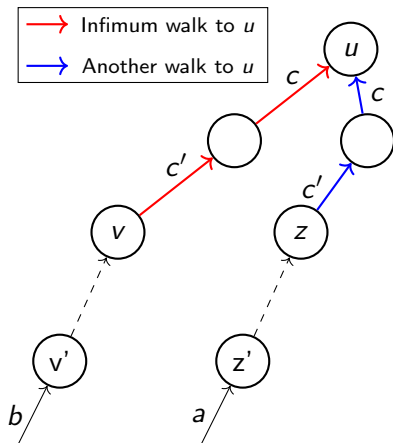


Similar as before, however now the supremum and the infimum **meet at the same node 2**

Special case! We don't know whether $5 \leq_{FS} 6$ holds

We treat this special case separately

Definition Inf/Sup Conflicts



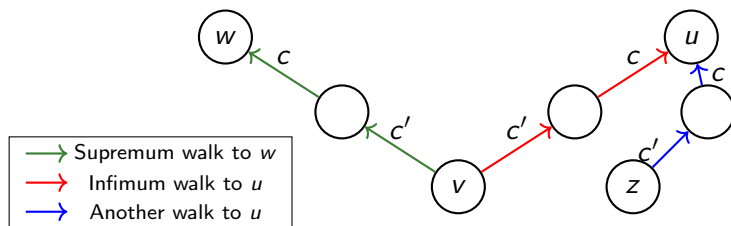
Node v is in **Inf conflict** with u because:

- v is in the **infimum walk** of u
- Node z can reach u with the **same labels $c'c$**
- $v \leq_{FS} z$ **does not hold**

Inf/Sup Conflicts for \leq_{FS}

special case seen before: the **preceding pair** of (v, z) is also a **preceding pair** of $(w, u) \implies w \leq_{FS} u$ **does not hold**

$w \leq_{FS} u \iff v$ not in **sup conflict** with w or in **inf conflict** with u



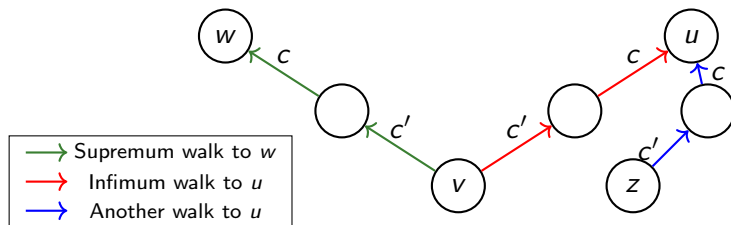
Lemma 22

It is possible to store the **inf/sup conflicts of every state** in $O(n)$ space

Inf/Sup Conflicts for \leq_{FS}

special case seen before: the **preceding pair** of (v, z) is also a **preceding pair** of $(w, u) \implies w \leq_{FS} u$ **does not hold**

$w \leq_{FS} u \iff v$ not in **sup conflict** with w or in **inf conflict** with u



Lemma 22

It is possible to store the **inf/sup conflicts of every state** in **$O(n)$** space

Representation of \leq_{FS}

For every state **u** we save:

- An **infimum walk** to **u**
- A **supremum walk** to **u**
- Its position in the **linear extension**
- Its **inf conflicts**
- Its **sup conflicts**

All of these occupy $O(n)$ total space

Thank you for your attention 😊



Funded by ERC StG “REGINDEX: Compressed indexes for regular languages with applications to computational pan-genomics” grant nr 101039208. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Research Council Executive Agency. Neither the European Union nor the granting authority can be held responsible for them.