



Università  
Ca' Foscari  
Venezia

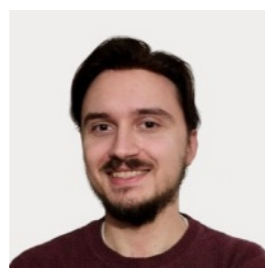
# Random Wheeler Automata

CPM 2024

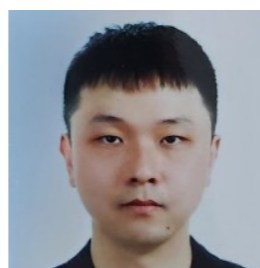
Fukuoka, Japan, 25 June 2024



**Ruben Becker**



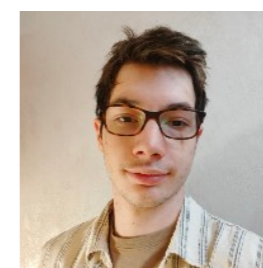
Davide Cenzato



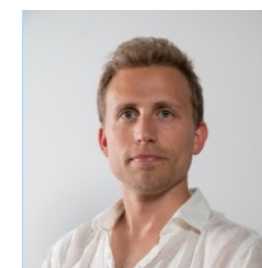
Sung-Hwan Kim



Bojana Kodric



Riccardo Maso



Nicola Prezza

# Setting and Contribution

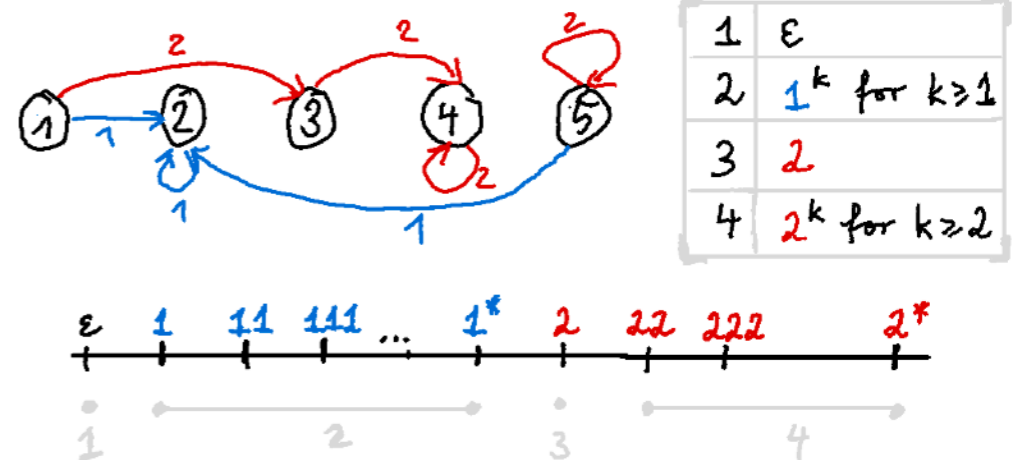
**Setting:** - **Wheeler Automata** constitute class of Automata, s.t. we can

- store them **concisely** ( $O(1)$  space per edge)
- do **pattern matching** on their accepted language

→ generalizing FM index from strings to (some) graphs

- Here: Focus on **DFA**s, then:

**Wheeler** = DFA's whose states can be ordered **consistently** with **colex-order of accepted strings**.

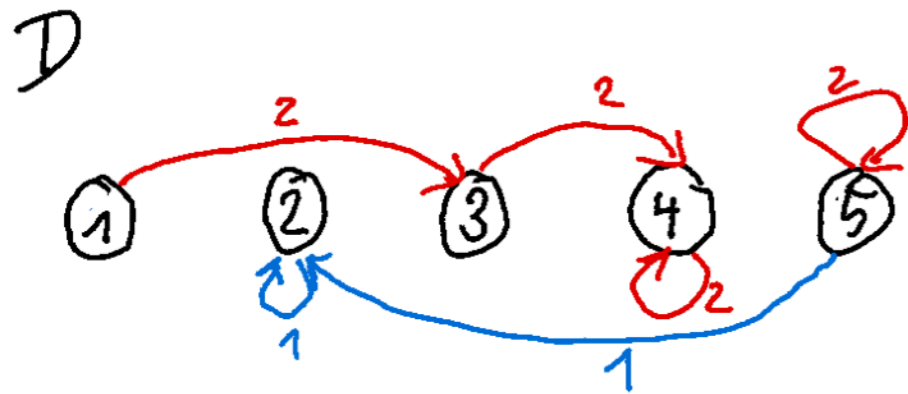


- Plenty of algorithmic / language-theoretic research, but **few available data sets**.

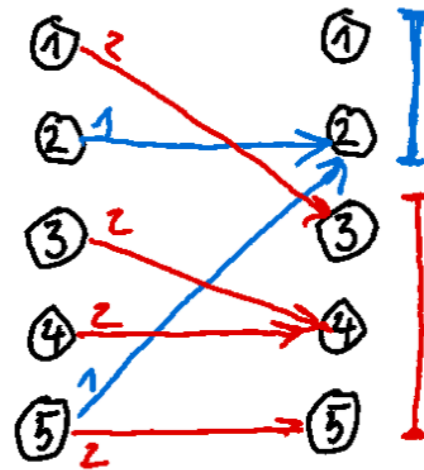
**Contribution:** Algorithm to generate **uniform** Wheeler DFA from  $\mathcal{D}_{n,m,v}$  in **constant space** & expected **linear time** all Wheeler DFA's with  $n$  nodes,  $m$  edges, alphabet size  $v$ .

# Wheeler DFAs

Def.: Wheeler DFA  $D = (Q, \delta, \Sigma)$  is s.t.  $\prec$  satisfies



"Bipartite Representation"



- (i) ranges | |
- (ii) edges of same color do not cross

$O =$

	1	2
1		
2		
3		
4		
5		

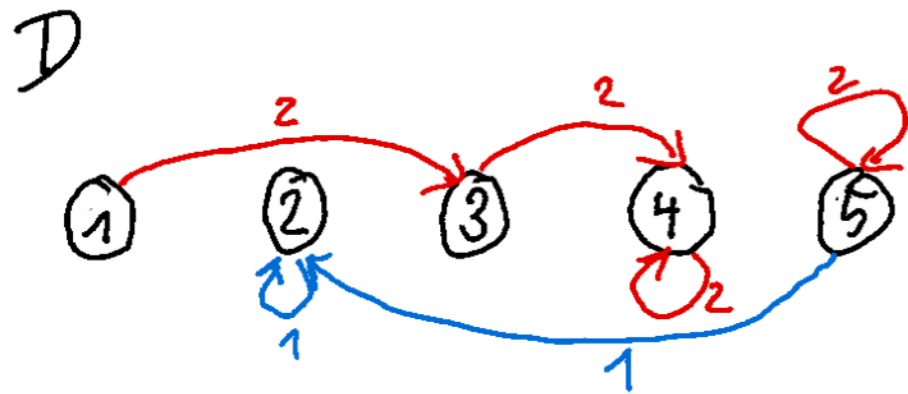
$I =$

--	--	--	--	--	--	--

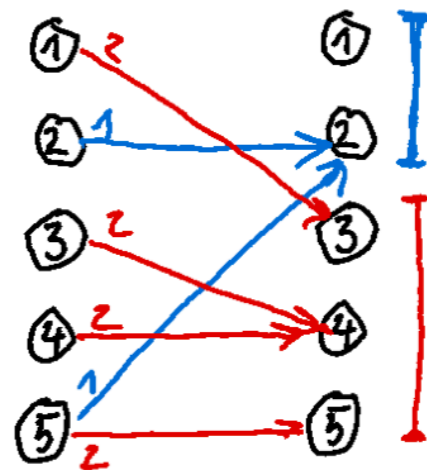
# Wheeler DFAs

Def.: **Wheeler DFA**  $D = (Q, \delta, \Sigma)$  is s.t.  $\prec$  satisfies

- (i)  $u' = \delta_a(u), v' = \delta_{a'}(v), a < a' \Rightarrow u' < v'$
- (ii)  $u' = \delta_a(u) \neq v' = \delta_a(v), u < v \Rightarrow u' < v'$



"Bipartite Representation"



- (i) ranges I
- (ii) edges of same color do not cross

$O =$

	1	2
1	0	
2	1	
3	0	
4	0	
5	1	

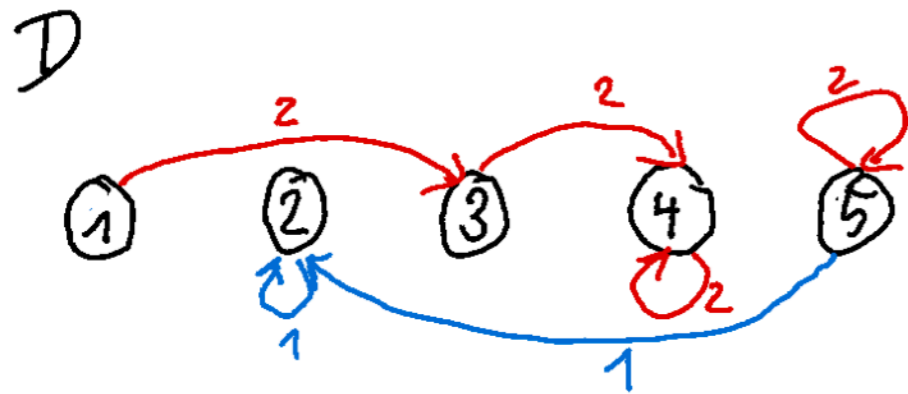
$I =$

	2					
1	1	0				

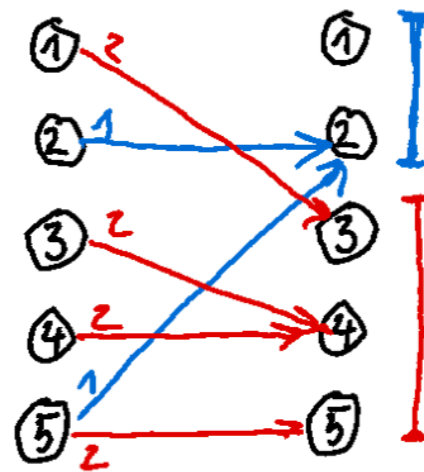
# Wheeler DFAs

Def.: **Wheeler DFA**  $D = (Q, \delta, \Sigma)$  is s.t.  $\prec$  satisfies

- (i)  $u' = \delta_a(u), v' = \delta_{a'}(v), a < a' \implies u' < v'$
- (ii)  $u' = \delta_a(u) \neq v' = \delta_a(v), u < v \implies u' < v'$



"Bipartite Representation"



- (i) ranges | |
- (ii) edges of same color do not cross

$O =$

	1	2
1	0	1
2	1	
3	0	
4	0	
5	1	

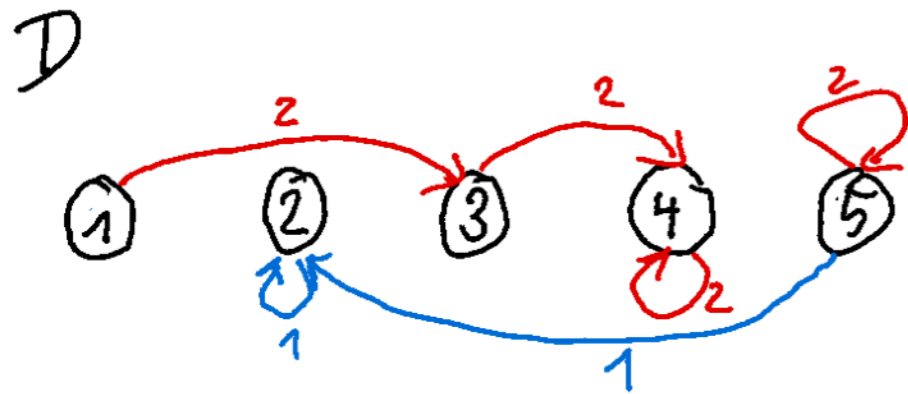
$I =$

	2	3				
1	1	0	1			

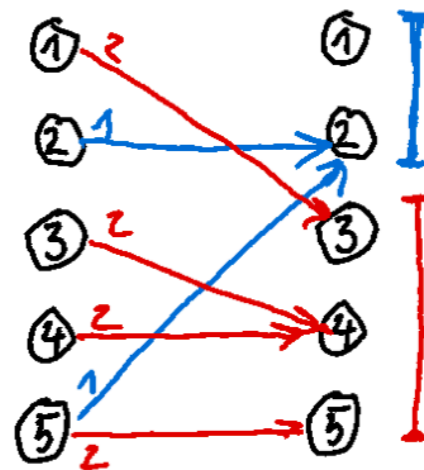
# Wheeler DFAs

Def.: **Wheeler DFA**  $D = (Q, \delta, \Sigma)$  is s.t.  $\prec$  satisfies

- (i)  $u' = \delta_a(u), v' = \delta_{a'}(v), a < a' \implies u' < v'$
- (ii)  $u' = \delta_a(u) \neq v' = \delta_a(v), u < v \implies u' < v'$



"Bipartite Representation"



- (i) ranges I
- (ii) edges of same color do not cross

$O =$

	1	2
1	0	1
2	1	0
3	0	1
4	0	1
5	1	

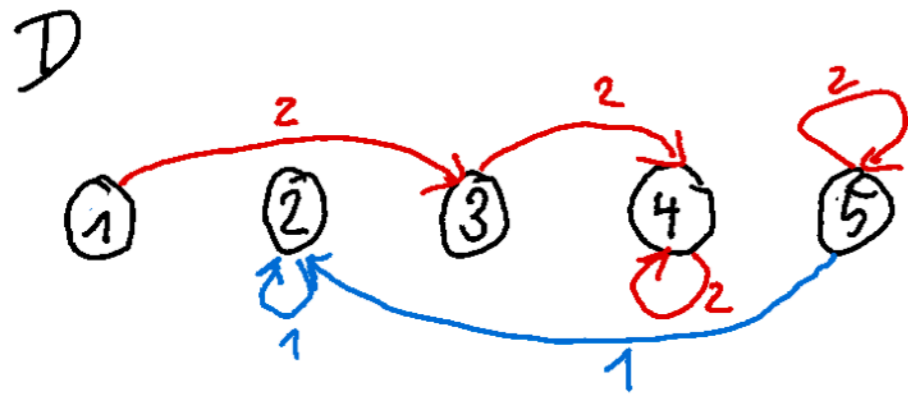
$I =$

	2	3	4		
1	1	0	1	1	0

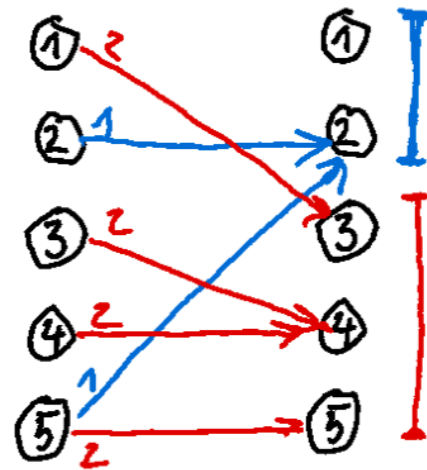
# Wheeler DFAs

Def.: **Wheeler DFA**  $D = (Q, \delta, \Sigma)$  is s.t.  $\prec$  satisfies

- (i)  $u' = \delta_a(u), v' = \delta_{a'}(v), a < a' \Rightarrow u' < v'$
- (ii)  $u' = \delta_a(u) \neq v' = \delta_a(v), u < v \Rightarrow u' < v'$



"Bipartite Representation"



- (i) ranges I
- (ii) edges of same color do not cross

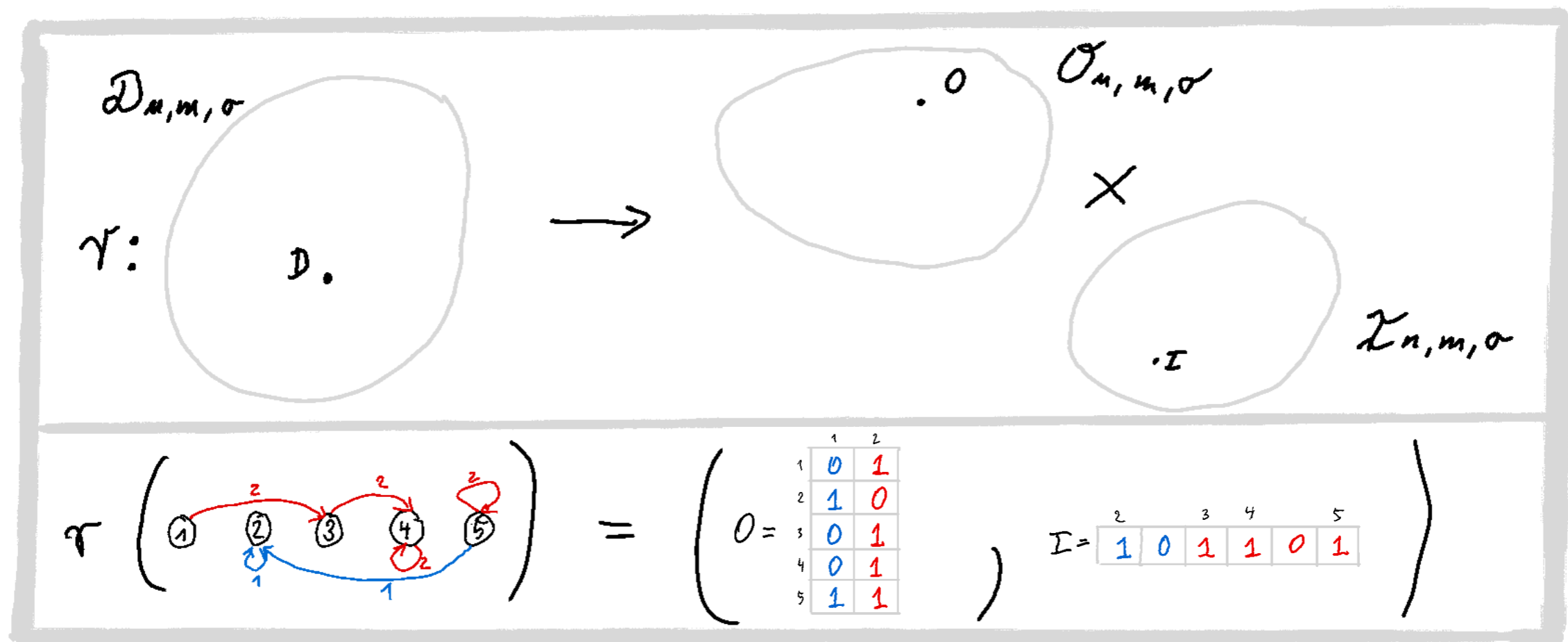
$O =$

	1	2
1	0	1
2	1	0
3	0	1
4	0	1
5	1	1

$I =$

	2	3	4	5
1	1	0	1	1
2	0	1	1	0
3	1	0	1	1
4	0	1	1	0
5	1	1	0	1

# Sampling Wheeler DFAs

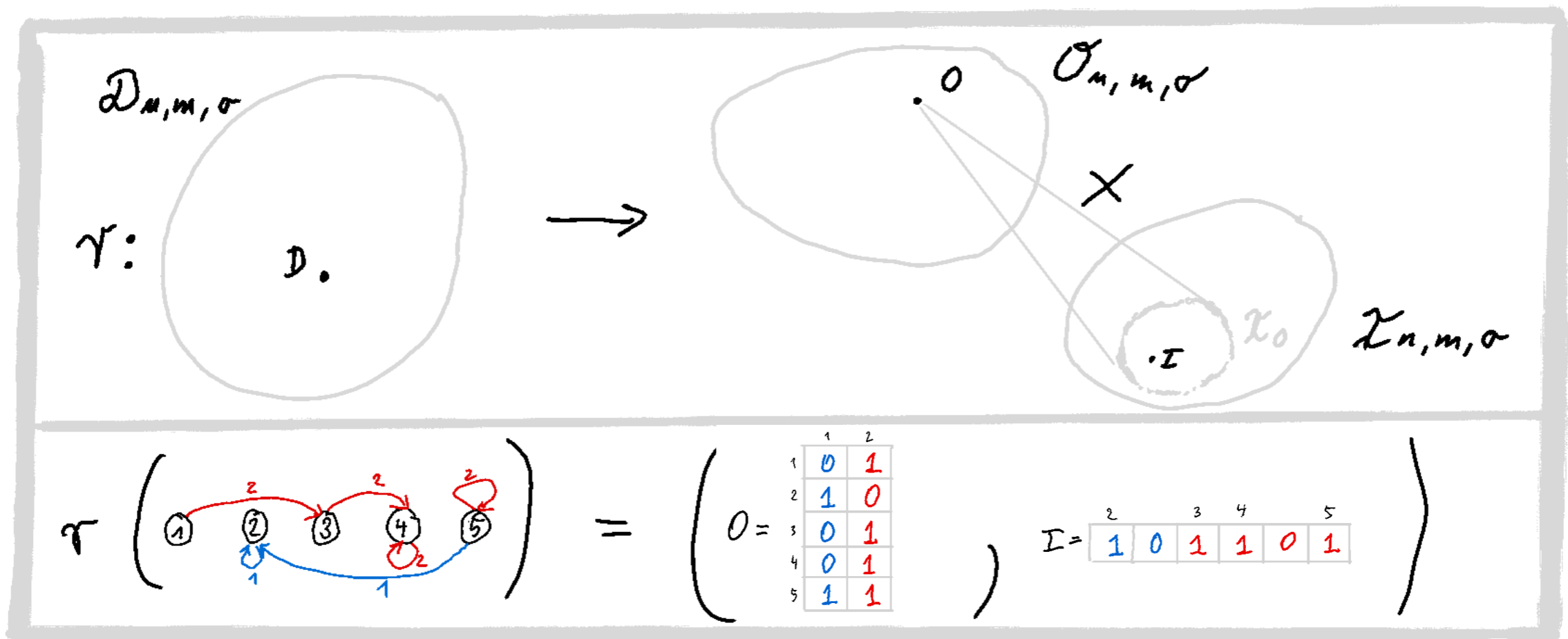


**Idea:** Sample from  $D_{n,m,\sigma}$  by sampling from  $O_{n,m,\sigma} \times I_{n,m,\sigma}$

- $\pi$  is injective  $\checkmark$  [a pair corresponds to at most one DFA]
- $\pi$  is not surjective  $\times$  [not every pair corresponds to a DFA]

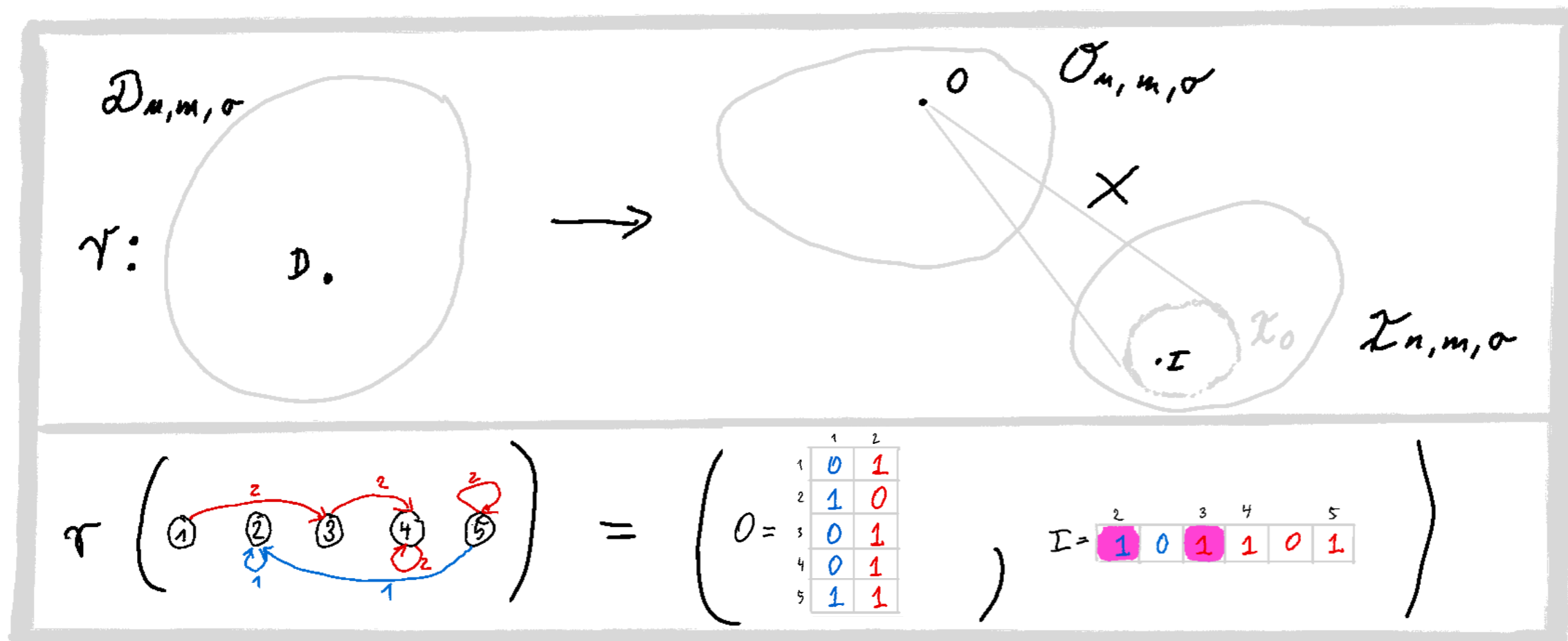


# Sampling Wheeler DFAs ctd.



$$\mathcal{X}_0 := \left\{ I \in \mathcal{X}_{n,m,\sigma} : I_{1 + \sum_{k=1}^{j-1} \|0_k\|_1} = 1 \quad \forall j \in [\sigma] \right\}$$

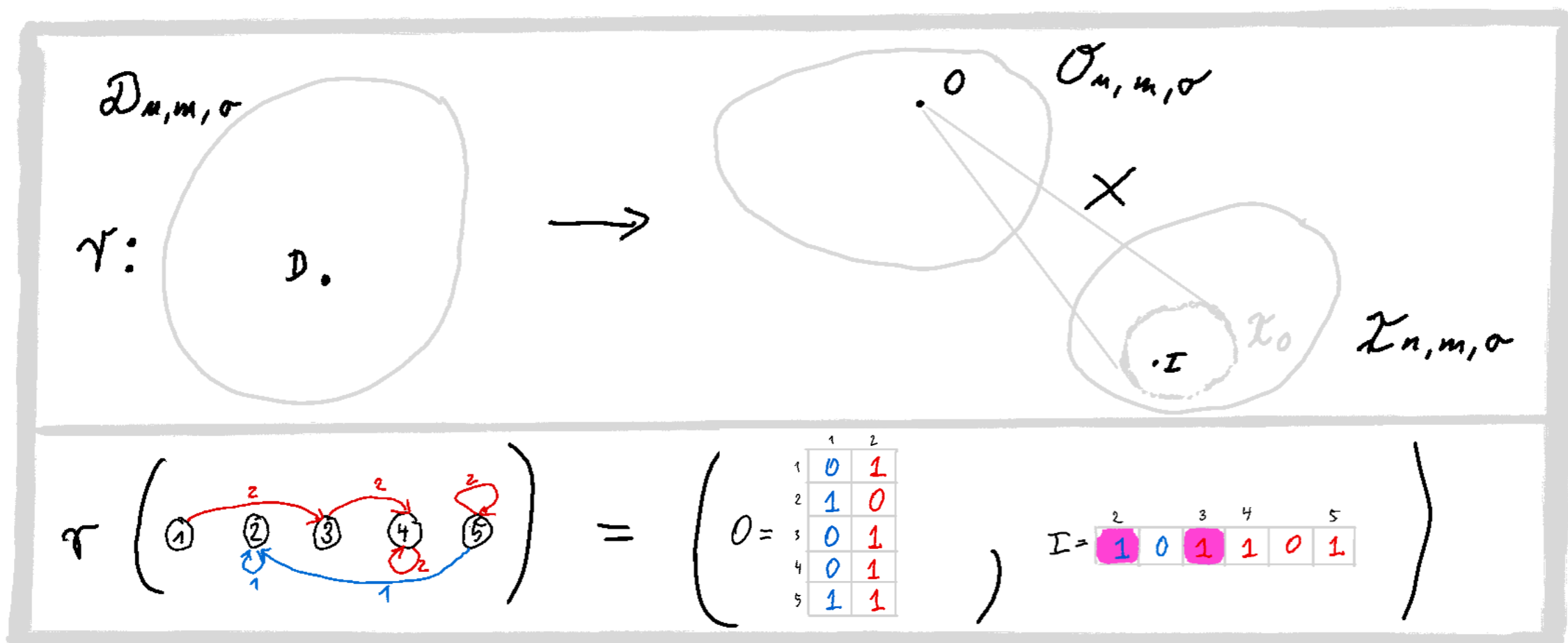
# Sampling Wheeler DFAs ctd.



$$\mathcal{L}_0 := \{ I \in \mathcal{X}_{n,m,\sigma} : I_{1 + \sum_{k=1}^{j-1} \|O_k\|_1} = 1 \quad \forall j \in [0] \}$$

$$\mathcal{R}_{n,m,\sigma} := \{ (D, I) : D \in \mathcal{O}_{n,m,\sigma} \text{ and } I \in \mathcal{L}_0 \}$$

# Sampling Wheeler DFAs ctd.



$$\mathcal{I}_0 := \{ I \in \mathcal{I}_{n,m,\sigma} : I_{1 + \sum_{k=1}^{j-1} \|O_k\|_1} = 1 \quad \forall j \in [0] \}$$

$$\mathcal{R}_{n,m,\sigma} := \{ (O, I) : O \in \mathcal{O}_{n,m,\sigma} \text{ and } I \in \mathcal{I}_0 \}$$

$\gamma: \mathcal{D}_{n,m,\sigma} \rightarrow \mathcal{R}_{n,m,\sigma}$  is **bijective**!

→ **sample from  $\mathcal{D}_{n,m,\sigma}$  by sampling from  $\mathcal{R}_{n,m,\sigma}$**   
 ↳ uniformity follows from  $|\mathcal{I}_0| = |\tilde{\mathcal{I}}_0| \quad \forall O, O' \in \mathcal{O}_{n,m,\sigma}$

# Sampling from $\mathcal{R}_{n,m,\sigma}$

$$\mathcal{R}_{n,m,\sigma} := \{(O, I) : O \in \mathcal{O}_{n,m,\sigma} \text{ and } I \in \mathcal{I}_0\},$$

where  $\mathcal{I}_0 := \{I \in \mathcal{I}_{n,m,\sigma} : I_{1 + \sum_{k=1}^{j-1} \|O_k\|_1} = 1 \ \forall j \in [\sigma]\}$

## sample\_O

repeat

$O := \text{reshape}_{n,\sigma}(\text{shuffle}(1^m 0^{n\sigma-m}))$

until  $\|O_j\| \geq 1$  for all  $j \in [\sigma]$

return  $O$

## sample\_I(O)

mask :=  $1 \#^{\|O_1\|-1} 1 \#^{\|O_2\|-1} \dots 1 \#^{\|O_\sigma\|-1}$

$I := \text{fill}(\text{mask}, \text{shuffle}(1^{n-\sigma-1} 0^{m-n+1}))$

return  $I$

## build-D(O, I)

$\delta := \emptyset, i := 1, v := 1$

for  $j = 1, 2, \dots, \sigma$  do

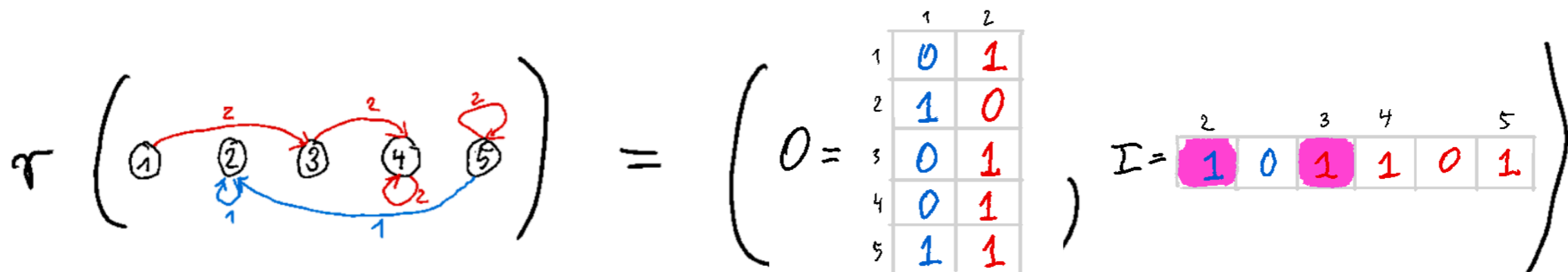
    for  $u = 1, 2, \dots, n$  do

        if  $O_{u,j} = 1$  then

            if  $I_i = 1$  then  $v := v + 1$

$\delta := \delta \cup \{((u, j), v)\}$

$i := i + 1$



# Sampling from $\mathcal{R}_{n,m,\sigma}$

$$\mathcal{R}_{n,m,\sigma} := \{ (O, I) : O \in \mathcal{O}_{n,m,\sigma} \text{ and } I \in \mathcal{I}_0 \},$$

$$\text{where } \mathcal{I}_0 := \{ I \in \mathcal{I}_{n,m,\sigma} : I_{1 + \sum_{k=1}^{j-1} \|O_k\|_1} = 1 \quad \forall j \in [\sigma] \}$$

sample\_O

repeat

$O := \text{reshape}_{n,\sigma}(\text{shuffle}(1^m 0^{n\sigma-m}))$

until  $\|O_j\| \geq 1$  for all  $j \in [\sigma]$

return  $O$

sample\_I(O)

mask :=  $1 \#^{\|O_1\|-1} 1 \#^{\|O_2\|-1} \dots 1 \#^{\|O_\sigma\|-1}$

$I := \text{fill}(\text{mask}, \text{shuffle}(1^{n-\sigma-1} 0^{m-n+1}))$

return  $I$

build-D(O,I)

$\delta := \emptyset, i := 1, v := 1$

for  $j = 1, 2, \dots, \sigma$  do

    for  $u = 1, 2, \dots, n$  do

        if  $O_{u,j} = 1$  then

            if  $I_i = 1$  then  $v := v + 1$

$\delta := \delta \cup \{((u, j), v)\}$

$i := i + 1$

# of rejections in sample\_O:

- $\Pr[\text{column } j \text{ empty}] = \prod_{i=1}^m \left( 1 - \frac{n}{n\sigma - (i-1)} \right) \leq \left( 1 - \frac{1}{\sigma} \right)^m = O(1/\sigma)$

- $\Pr[\exists \text{ empty column}] = O(1)$  [union bound]

- $\mathbb{E}[\# \text{ iterations}]$  is constant and # iterations is  $O(\log m)$  w. pr.  $1 - m^{-c}$

# Constant Space Implementation

- build  $D$  accesses  $O$  column- and  $I$  bit-wise
- can generate non-zero entries of  $D$  and  $I$  on the fly in constant space using sequential shuffler  $(N, k)$  of [Shukelyan & Cormode]
- ↳ • generate  $k$  uniform integers from  $[N]$  in ascending order within constant space
  - call with  $N = m\sigma$   $k = m$  for  $D$
  - $N = m - \sigma$   $k = m - \sigma - 1$   $I$
- sequential shuffler  $(N, k)$  is essentially a clever implementation of Knuth's shuffler

**Theorem:** Generate a uniform Wheeler DFA from  $D_{n, m, \sigma}$  in  $O(1)$  space  $O(m)$  exp. time  $O(m \log m)$  time w.k.pr. if  $\sigma \leq \frac{m}{\epsilon n}$ .

**Implementation:** • Generates DFA with  $n = 64 \cdot 10^6$   $m = 8 \cdot 10^9$  in  $\approx 10$  mins  
• Throughput:  $\geq 8 \cdot 10^6$  edges per second

# Add-On: Bound on Number of WDFA's

Our Wheeler DFA representation implies a bound on number of Wheeler DFA's  $\mathcal{D}_{n,\sigma}$  with  $n$  nodes and alphabet size  $\sigma$ .

**Theorem:** For  $\epsilon \in (0, 1/2]$ ,  $n \geq 2/\epsilon$ , and  $\sigma \leq (1-\epsilon) \cdot n$

$$\log |\mathcal{D}_{n,\sigma}| \geq n\sigma + (n-\sigma) \log \sigma - (n + \log \sigma)$$

$$\log |\mathcal{D}_{n,\sigma}| \leq n\sigma + (n-\sigma) \log \sigma + O(n)$$

↳ information-theoretic worst case #bits to encode WDFA from  $\mathcal{D}_{n,\sigma}$

Our representation (opportune encoded using succinct bitvectors) gives an encoding of size  $n\sigma + (n-\sigma) \log \sigma$ , thus being optimal up to an additive  $O(n)$  term.

# Conclusion & Future Work

Theorem: Generate a uniform Wheeler DFA from  $D_{n,m,\sigma}$  in  $O(1)$  space  $O(m)$  exp. time  $O(m \log m)$  time w.k.p. if  $\sigma \leq \frac{m}{\ln m}$ .

Implementation:  
• Generates DFA with  $n = 64 \cdot 10^6$   $m = 8 \cdot 10^9$  in  $\approx 10$  mins  
• Throughput:  $\geq 8 \cdot 10^6$  edges per second

Future Work:  
- Wheeler NFAs  
- Automata of colex width  $p$   
- analyse threshold phenomena on  $D_{n,m,\sigma}$





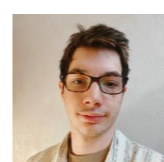
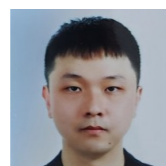
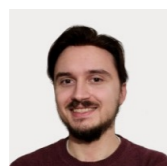
Università  
Ca' Foscari  
Venezia

That's all. Thank you! Questions?

# Random Wheeler Automata

CPM 2024

Fukuoka, Japan, 25 June 2024



**Ruben Becker**, Davide Cenzato, Sung-Hwan Kim, Bojana Kodric, Riccardo Maso, and Nicola Prezza



This work has been funded by the European Union (ERC, REGINDEX, 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. Neither the European Union nor the granting authority can be held responsible for them.