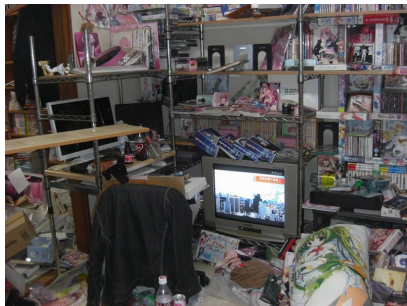# The *rational* construction of a (Wheeler) DFA

Giovanni Manzini, <u>Alberto Policriti</u>, Nicola Prezza, and Brian Riccardi

# Order is important ...

# Order is important ...

# Order is important ...



... and some orders are more important than others (e.g. the order of $\mathbb{Q}$)

## Automata and (Co-lexicographic) Order

$\mathcal{A}$ *is input-consistent*: $(\forall u, v \in Q)(\delta(u, a_1) = \delta(v, a_2) \rightarrow a_1 = a_2)$.

$$\delta(q) = q' \text{ stands for } \delta(q, \lambda(q')) = q'.$$

### Definition

A Wheeler DFA (WDFA) $\mathcal{A} = (Q, s, \delta, F, \prec)$ is such that $(Q, \prec)$ is a total order with $s$ as minimum, and letting $v_1 = \delta(u_1)$, and $v_2 = \delta(u_2)$:

  i $v_1 \prec v_2 \Rightarrow \lambda(v_1) \leqslant \lambda(v_2)$;

  ii $(\lambda(v_1) = \lambda(v_2) \wedge v_1 \prec v_2) \Rightarrow u_1 \prec u_2$.

(GAGIE-MANZINI-SIRÉN. TCS 2017) (ALANKO-D'AGOSTINO-P.-PREZZA. INF. AND COMP. 2021)

(COTUMACCIO-D'AGOSTINO-P.-PREZZA. JACM. 2023)

## Automata and (Co-lexicographic) Order

$\mathcal{A}$ *is input-consistent*: $(\forall u, v \in Q)(\delta(u, a_1) = \delta(v, a_2) \rightarrow a_1 = a_2)$.

$$\delta(q) = q' \text{ stands for } \delta(q, \lambda(q')) = q'.$$

### Definition
A Wheeler DFA (WDFA) $\mathcal{A} = (Q, s, \delta, F, \prec)$ is such that $(Q, \prec)$ is a total order with $s$ as minimum, and letting $v_1 = \delta(u_1)$, and $v_2 = \delta(u_2)$:

   i $v_1 \prec v_2 \Rightarrow \lambda(v_1) \leqslant \lambda(v_2)$;

   ii $(\lambda(v_1) = \lambda(v_2) \wedge v_1 \prec v_2) \Rightarrow u_1 \prec u_2$.

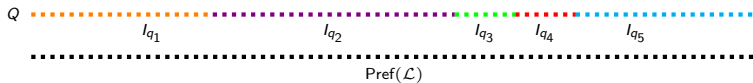(GAGIE-MANZINI-SIRÉN. TCS 2017) (ALANKO-D'AGOSTINO-P.-PREZZA. INF. AND COMP. 2021)

(COTUMACCIO-D'AGOSTINO-P.-PREZZA. JACM. 2023)

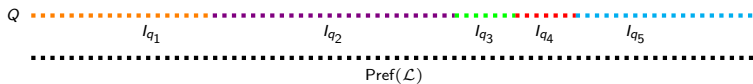$\Sigma = \{a_1, a_2, \ldots, a_\sigma\}$ is ordered and W(i)-W(ii) extend the ordering to strings reaching states

<div align="center">co-lexicographically</div>

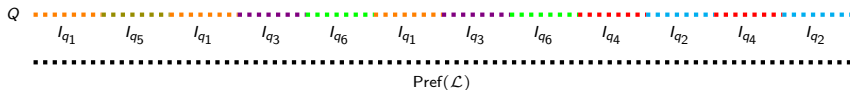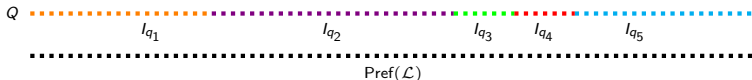align to the right and compare right-to-left (cbaabba < aaacba)

# WDFA: states are intervals of strings
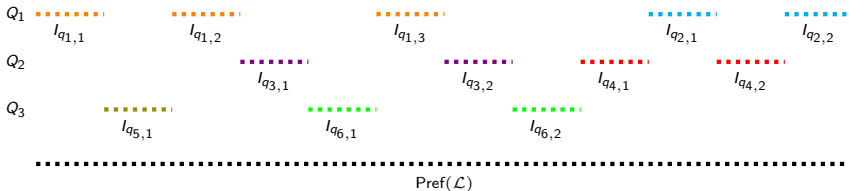
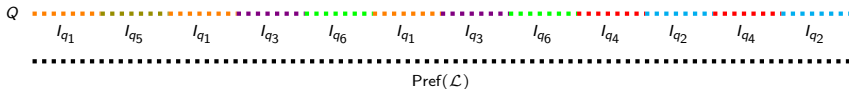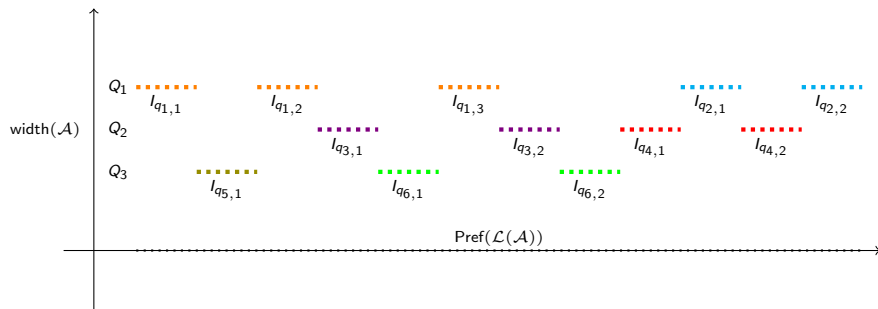## WDFA: states are intervals of strings



## DFA

## WDFA: states are intervals of strings

$Q$

$I_{q_1}$   $I_{q_2}$   $I_{q_3}$   $I_{q_4}$   $I_{q_5}$

$\text{Pref}(\mathcal{L})$

## DFA

$Q$

$I_{q_1}$  $I_{q_5}$  $I_{q_1}$  $I_{q_3}$  $I_{q_6}$  $I_{q_1}$  $I_{q_3}$  $I_{q_6}$  $I_{q_4}$  $I_{q_2}$  $I_{q_4}$  $I_{q_2}$

$\text{Pref}(\mathcal{L})$

$Q_1$

$I_{q_{1,1}}$   $I_{q_{1,2}}$   $I_{q_{1,3}}$   $I_{q_{2,1}}$   $I_{q_{2,2}}$

$Q_2$

$I_{q_{3,1}}$   $I_{q_{3,2}}$   $I_{q_{4,1}}$   $I_{q_{4,2}}$

$Q_3$

$I_{q_{5,1}}$   $I_{q_{6,1}}$   $I_{q_{6,2}}$
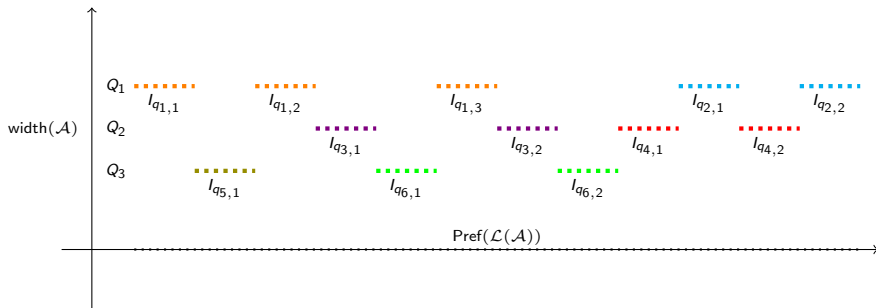
$\text{Pref}(\mathcal{L})$

## ... coordinates

## ... coordinates



- ▸ (substring closure) membership in $\mathcal{L}(\mathcal{A})$ in $O(\text{width}(\mathcal{A})^2)$ per matched character;
- ▸ any NFA $\mathcal{N}$ is equivalent to a DFA $\mathcal{D}$ with at most $2^{\text{width}(\mathcal{N})}(|\mathcal{N}| - \text{width}(\mathcal{N}) + 1)$ states;
- ▸ any automaton of width $p$ can be encoded in $O(\log p + \log |\Sigma|)$ bits per transition.

(COTUMACCIO-D'AGOSTINO-P.-PREZZA. JACM. 2023)

**Intermezzo: Path coherence**

### Definition (Path coherence)

Given a NFA $\mathcal{A}$ and an order $\prec$ over the set of its states, we say that $\mathcal{A}$ is *path coherent* iff strings send intervals (of states) into intervals:

$$\forall q \leq q' \ \forall \alpha \ \exists p \leq p' \ \left( [q, q'] \overset{\alpha}{\leadsto} [p, p'] \right).$$

## The rational embedding of strings

### Definition (**The Rational Embedding of $\Sigma^*$**)

The *Rational Embedding* of $\Sigma^*$ is the map $\mathfrak{q} : \Sigma^* \to \mathbb{Q}[0,1)$ such that, for any $\alpha = \alpha_1 \ldots \alpha_m \in \Sigma^*$:

$$\mathfrak{q}(\alpha) = \sum_{i=1}^{m} \alpha_i \cdot (\sigma + 2)^{-(m-i+1)}.$$

### Example

start $\rightarrow$ ( # ) $\rightarrow$ ( 4 ) $\rightarrow$ ( 7 ) $\rightarrow$ ( 7 ) $\rightarrow$ ( 5 ) $\rightarrow$ ( 3 )    $\mathfrak{q}(\alpha) = \mathfrak{q}(47753) = 0.35774$

## The rational embedding of strings

### Definition (**The Rational Embedding of $\Sigma^*$**)

The *Rational Embedding* of $\Sigma^*$ is the map $\mathfrak{q} : \Sigma^* \to \mathbb{Q}[0,1)$ such that, for any $\alpha = \alpha_1 \ldots \alpha_m \in \Sigma^*$:

$$\mathfrak{q}(\alpha) = \sum_{i=1}^{m} \alpha_i \cdot (\sigma + 2)^{-(m-i+1)}.$$

### Example



$\mathfrak{q}(\alpha) = \mathfrak{q}(47753) = 0.35774$

property:
$\boldsymbol{\alpha} < \boldsymbol{\beta}$ (in co-lex order) $\Leftrightarrow \mathfrak{q}(\boldsymbol{\alpha}) < \mathfrak{q}(\boldsymbol{\beta})$ (as rational numbers).

# The rational embedding of DFAs

### Definition

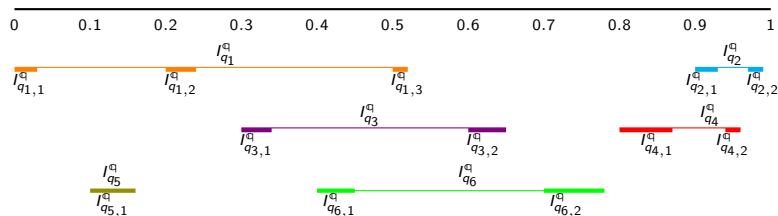$I_{\mathbb{Q}[0,1)}$ the collection of convex sets of rationals in $\mathbb{Q}[0,1)$:

### Definition (**The Rational Embedding of a DFA**)

The *Rational Embedding* of $\mathcal{A} = (Q, s, \delta, F)$ is the map
$I^{\mathbb{q}} : Q \to I_{\mathbb{Q}[0,1)}$ defined as follows: for any $q \in Q$,

$$I^{\mathbb{q}}(q) = \bigcap \{J \in I_{\mathbb{Q}[0,1)} \mid (\forall \alpha \in I_q)(\mathbb{q}(\alpha) \in J)\}.$$

$I^{\mathbb{q}}(q)$ is the convex closure (hull) of $I_q$. (Notation: $I^{\mathbb{q}}_q = I^{\mathbb{q}}(q)$)

## The rational embedding of DFAs



Determinism: $q \neq q' \Rightarrow I_q \cap I_{q'} = \varnothing$.
But it might be that $q \neq q' \wedge I_q^{\mathbb{Q}} \cap I_{q'}^{\mathbb{Q}} \neq \varnothing$.

From (ALANKO-D'AGOSTINO-P.-PREZZA. INF. AND COMP. 2021)[THEOREM 4.3]:

$$\mathcal{A} \text{ is Wheeler iff } q \neq q' \Rightarrow I_q^{\mathbb{Q}} \cap I_{q'}^{\mathbb{Q}} = \varnothing.$$

**The rational embedding of DFAs**

### Example

There are *"solid"* collections of rational embeddings of strings accepted by a DFA: Σ* (non-denumerable set of accumulation points)

## The rational embedding of DFAs

### Example

There are *"solid"* collections of rational embeddings of strings accepted by a DFA: $\Sigma^*$ (non-denumerable set of accumulation points)

### Example

Missing largest and smallest digit $\Rightarrow$ (in general) there is a successor but <u>no</u> predecessor.

$$0.7348 < \cdots < 0.73488 < \cdots < 0.734888 < \cdots \; 0.7345 < 0.73451$$

## The rational embedding of DFAs

### Example
There are *"solid"* collections of rational embeddings of strings accepted by a DFA: $\Sigma^*$ (non-denumerable set of accumulation points)

### Example
Missing largest and smallest digit $\Rightarrow$ (in general) there is a successor but <u>no</u> predecessor.

$$0.7348 < \cdots < 0.73488 < \cdots < 0.734888 < \cdots 0.7345 < 0.73451$$

Question: are $I_q^{\mathbb{q}}$'s left and right limits $(\ell_q, r_q)$ always in $\mathbb{Q}$?
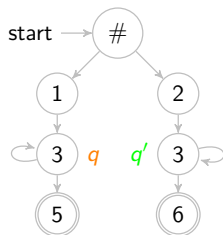
## Ordering <u>states</u>: Entanglement

### Definition

$Q' \subseteq Q$ is entangled if there exists a <u>monotone sequence</u> $(\alpha_i)_{i \in \mathbb{N}}$ in $Pref(\mathcal{L}(\mathcal{D}))$ such that:

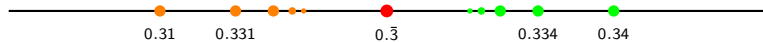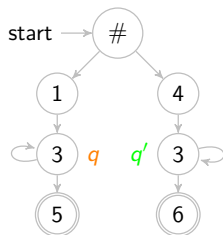$$\forall u' \in Q' \; \delta(s, \alpha_i) = u' \text{ for infinitely many } i's$$

### Definition
$Q' \subseteq Q$ is entangled if there exists a <u>monotone sequence</u> $(\alpha_i)_{i \in \mathbb{N}}$ in $Pref(\mathcal{L}(\mathcal{D}))$ such that:

$$\forall u' \in Q' \ \delta(s, \alpha_i) = u' \text{ for infinitely many } i's$$

### Lemma
*If a value $x$ is a left-accumulation point (resp. right-accumulation point) for both the sets $I_q^{\mathfrak{q}}$ and $I_{q'}^{\mathfrak{q}}$ then $q$ and $q'$ are entangled.*

# Examples
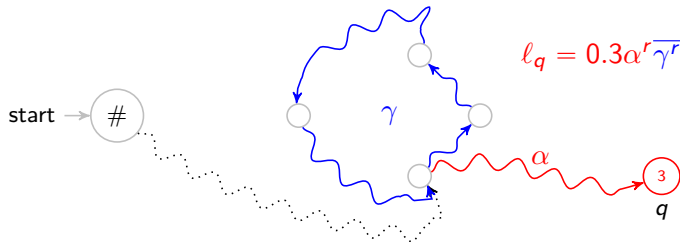
# Examples

## Finding Left and Right Limits

**Lemma**
*Let $\mathcal{L} = \mathcal{L}(\mathcal{D})$, with $\mathcal{L}$ Wheeler and $\mathcal{D}$ either minimum or Wheeler. For any $q \in Q$ we have:*

$$\ell_q = 0.a_{q,1}\cdots a_{q,h}\overline{a_{q,h+1}\cdots a_{q,h+j}},$$

*with $h + j \leqslant |Q|$, and $j > 0$ if and only if $\ell_q \notin I_q^{\mathfrak{q}}$.*

## Finding Left and Right Limits

### Lemma

*Let $\mathcal{L} = \mathcal{L}(\mathcal{D})$, with $\mathcal{L}$ Wheeler and $\mathcal{D}$ either minimum or Wheeler. For any $q \in Q$ we have:*

$$\ell_q = 0.a_{q,1} \cdots a_{q,h}\overline{a_{q,h+1} \cdots a_{q,h+j}},$$
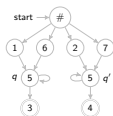
*with $h + j \leqslant |Q|$, and $j > 0$ if and only if $\ell_q \notin I_q^{\mathbb{q}}$.*

### Proof's idea

walk backward from $q$ and find $\alpha$ whose rational embedding $\mathbb{q}(\alpha)$ is the smallest among those reaching $q$.



$\ell_q = 0.3\alpha^r\overline{\gamma^r}$

## Finding Left and Right Limits

### Lemma

Let $\mathcal{L} = \mathcal{L}(\mathcal{D})$, with $\mathcal{L}$ Wheeler and $\mathcal{D}$ either minimum or Wheeler. For any $q \in Q$ we have:

$$\ell_q = 0.a_{q,1} \cdots a_{q,h}\overline{a_{q,h+1} \cdots a_{q,h+j}},$$

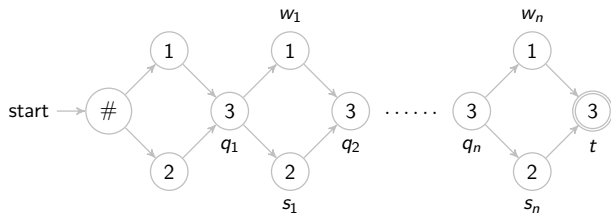with $h + j \leqslant |Q|$, and $j > 0$ if and only if $\ell_q \notin I_q^{\mathbb{Q}}$.

### Theorem

If $\mathcal{L} = \mathcal{L}(\mathcal{D}) = \mathcal{L}(\mathcal{D}_w)$, with $\mathcal{L}$ Wheeler and $\mathcal{D}$ either minimum or Wheeler, then for all $q \in Q$, we have $\ell_q, r_q \in \mathbb{Q}$.

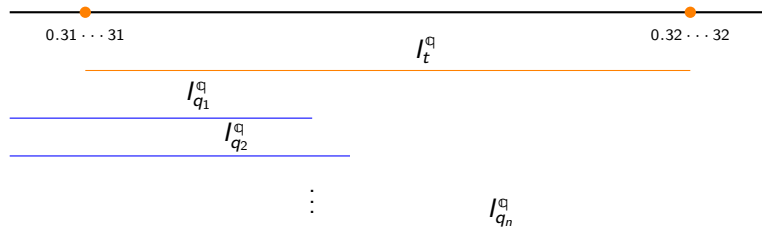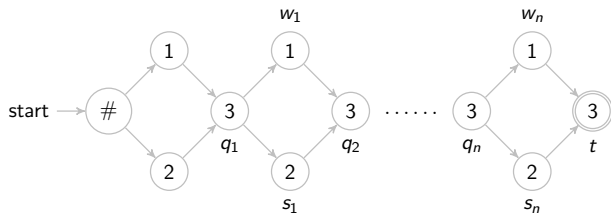More on the algorithmic side on (BECKER-CENZATO-KIM-KODRIC-P.-PREZZA. SPIRE 2023.)
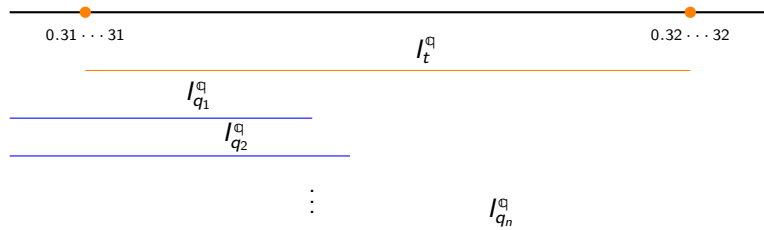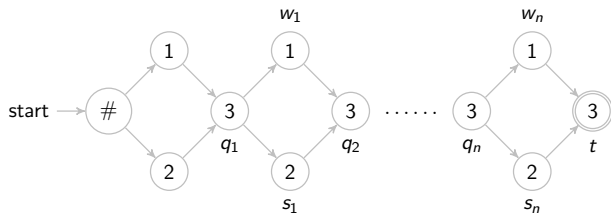More on the entanglement:

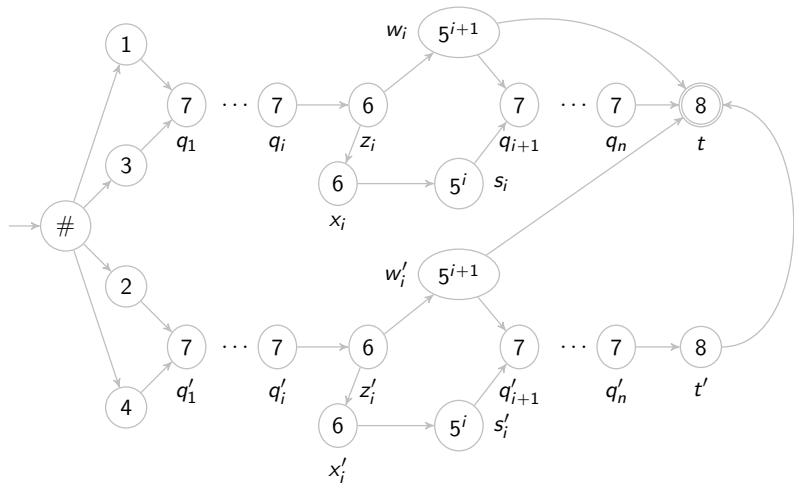# Minimum DFA vs. Minimum Wheeler-DFA
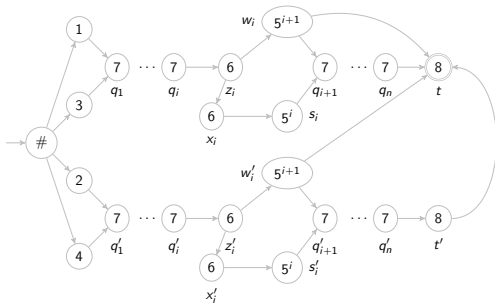
## Minimum DFA vs. Minimum Wheeler-DFA



Question: is the $|Q_w|/|Q|$ related to width($\mathcal{D}$)?

## Lower bound

| State type | Left limit | Right limit |
| --- | --- | --- |
| $s_{i,j}$ | $0.5^j 6675^i 67\ldots$ | $0.5^j 6675^{i-1}667\ldots$ |
| $w_{i,j}$ | $0.5^j 675^i 67\ldots$ | $0.5^j 675^{i-1}667\ldots$ |
| $x_i$ | $0.6675^i 67\ldots$ | $0.6675^{i-1}667\ldots$ |
| $z_i$ | $0.675^i 67\ldots$ | $0.675^{i-1}667\ldots$ |
| $q_i$ | $0.75^i 67\ldots$ | $0.75^{i-1}667\ldots$ |
| $t$ | $0.85^n 67\ldots$ | $0.8875^{n-1}667\ldots$ |
| $t'$ | $0.875^n 67\ldots$ | $0.875^{n-1}667\ldots$ |

**Theorem**
*Let $\mathcal{L} = \mathcal{L}(\mathcal{D}) = \mathcal{L}(\mathcal{D}_w)$, with $\mathcal{L}$ Wheeler, $\mathcal{D}$ minimum, $\mathcal{D}_w$ minimum Wheeler, and let $f(\cdot, \cdot)$ be such that $|\mathcal{D}_w| = O(f(|\mathcal{D}|, width(D)))$. Then, for any $k, p \in \mathbb{N}$,*

$$f(n, p) \notin O(n^k + 2^p).$$

## The arithmetic way

Formally, for the *left* case, we consider the problem of finding the set of all real-valued vectors $x \in \mathbb{R}^Q$ that satisfy the following constraint satisfaction program, that we name $\mathcal{P}_{Left}$:

$$(1) \quad x_s = 0,$$

$$(2) \quad 0 < x_q < 1, \qquad\qquad\qquad\qquad (\forall q \in Q \backslash \{s\})$$

$$(3) \quad (\sigma + 2) \cdot x_q - \lambda(q) = \min \left\{ x_{q'} \mid \delta(q') = q \right\}, \quad (\forall q \in Q \backslash \{s\})$$

## The arithmetic way

Formally, for the *left* case, we consider the problem of finding the set of all real-valued vectors $x \in \mathbb{R}^Q$ that satisfy the following constraint satisfaction program, that we name $\mathcal{P}_{Left}$:

(1) $\quad x_s = 0,$

(2) $\quad 0 < x_q < 1,$ $\hspace{5cm} (\forall q \in Q \backslash \{s\})$

(3) $\quad (\sigma + 2) \cdot x_q - \lambda(q) = \min \{x_{q'} \mid \delta(q') = q\}, \quad (\forall q \in Q \backslash \{s\})$

### Lemma
*Let $\mathcal{L}$ be a Wheeler language, and $\mathcal{D} = (Q, s, \delta, F)$ be either minimum or Wheeler accepting $\mathcal{L}$, and let $\ell \in \mathbb{Q}^Q$ be the vector of left limits. Then, $\ell$ is a solution of $\mathcal{P}_{Left}$.*

### Theorem
*Let $\mathcal{D} = (Q, s, \delta, F)$ be either minimum or Wheeler accepting $\mathcal{L}$ Wheeler, and $\ell \in \mathbb{Q}^Q$ be the vector of left limits. Then, $\mathcal{P}_{Left}$ always admits $\ell$ as its unique solution.*

## The arithmetic way

Consider the following linear program $\mathcal{P}^*_{Left}$:

maximize: $\sum_{q \in Q} x_q,$

subject to: $x_s = 0,$

$0 < x_q < 1, \qquad\qquad\qquad\qquad \forall q \in Q \backslash \{s\},$

$(\sigma + 2) \cdot x_q - \lambda(q) \leqslant x_{q'}, \quad \forall q, q' \in Q \text{ s.t. } \delta(q') = q,$

### Theorem
*Let $\mathcal{D} = (Q, s, \delta, F)$ be either minimum or Wheeler accepting $\mathcal{L}$
Wheeler, and $\ell \in \mathbb{Q}^Q$ be the vector of left limits. Then $\mathcal{P}^*_{Left}$ always
admits $\ell$ as its unique solution.*

## Conclusions and Open problems

- A parameter measuring the DFA vs. WDFA growth (*finite* entanglement)
- On-line splitting of minimum DFA (self-adjusting splitting)
- Optimal *disentanglement*

## Conclusions and Open problems

- ▸ A parameter measuring the DFA vs. WDFA growth (*finite* entanglement)
- ▸ On-line splitting of minimum DFA (self-adjusting splitting)
- ▸ Optimal *disentanglement*

Optimization of the Hasse automaton $\mathcal{H}$ construction

$$\text{width}(\mathcal{L}(\mathcal{H})) = \text{width}(\mathcal{H}) = ent(\mathcal{H})$$

Thank you for your attention.