

WORDS, TREES AND AUTOMATA MINIMIZATION



1

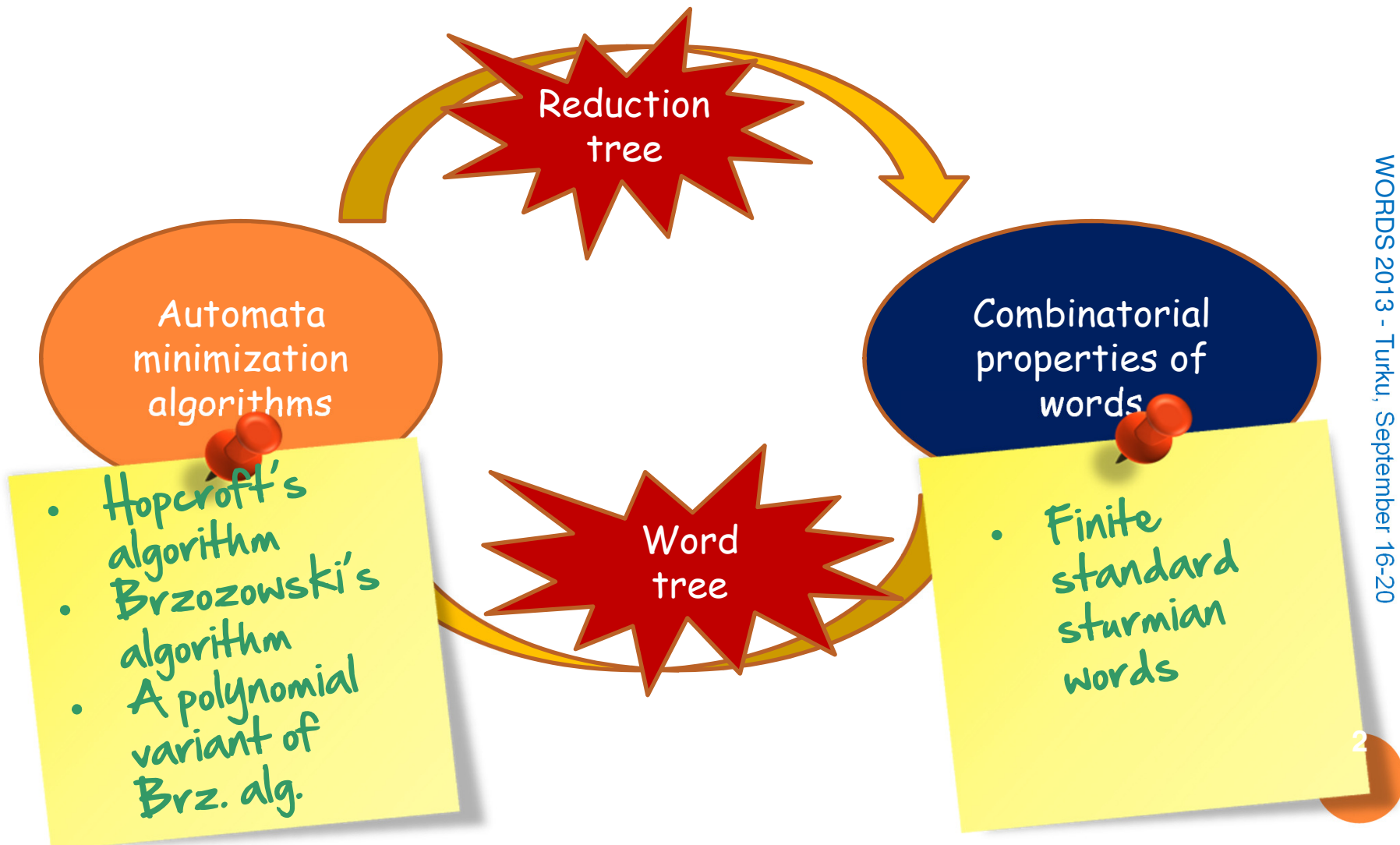
Marinella Sciortino

Joint work with Giusi Castiglione

University of Palermo

WORDS 2013 - Turku, September 16-20

GRAPHICAL SUMMARY



REGULAR LANGUAGES AND FSA

- Regular languages can be represented equivalently as rational expressions or as deterministic finite automata.
- A finite state automaton is denoted by

$$A=(Q,\Sigma,I,\delta,F),$$

where Q is the set of states, Σ is the finite alphabet, I the set of initial states, F the set of final states, δ the transition function from $Q \times \Sigma$ to 2^Q .

- The automaton is deterministic (DFA), denoted by $A=(Q,\Sigma,i,\delta,F)$, if it has a unique initial state and if δ is a function from $Q \times \Sigma$ to Q .
- There is a unique DFA (up to a renaming of the states) with a minimal number of states, called **minimal automaton**, recognizing a regular language L .
- The description of a language with its minimal automaton is important when space considerations matter in implementations of applications, such as in Pattern Matching, Lexical Analysis, Coding Systems, and so on.

NERODE EQUIVALENCE

- The minimal automaton equivalent to a given automaton A (i.e. recognizing the same language) is defined by the right-invariant Nerode equivalence.
- Given a state $p \in Q$, we consider the language
$$L_p(A) = \{v \in \Sigma^* \mid \delta(p, v) \in F\}$$
- The Nerode equivalence is defined as follows:
$$p \sim q \text{ if } L_p(A) = L_q(A)$$
- It is a congruence, i.e. for any $a \in \Sigma$,
$$p \sim q \text{ implies } \delta(p, a) \sim \delta(q, a),$$
and it is the coarsest congruence such that F is union of classes of the congruence.
- The equivalence classes are the states of the minimal automaton.

HOW TO COMPUTE IT?

MINIMIZATION ALGORITHMS

- There are lots of methods which can be used to minimize a finite automaton, i.e. to find the minimal automaton equivalent to a given DFA.
- Some of them operate by successive refinements of a partition of the states:
 - Moore, 1956. Time complexity $O(kn^2)$, $k=|\Sigma|$
 - Hopcroft, 1971. Time complexity $O(kn \log n)$, $k=|\Sigma|$
- The Brzozowski's method (1962) operates by reversal and determinization repeated twice.
Time complexity: exponential worst case, but good performance in practice.
- A taxonomy of finite automata minimization algorithms is given by B. Watson, 1994.
- "Minimization of automata" by Berstel, Boasson, Carton, Fagnot, 2010.
- Polynomial variants of Brzozowski's method have been recently introduced.
- Many papers on experimental comparison between minimization algorithms.

HOPCROFT'S ALGORITHM: THE FASTEST

- ◉ No faster algorithm is known for general deterministic automata
- ◉ In order to obtain the Nerode equivalence, a sequence of refinements of equivalence relations is realized.
- ◉ The sequence starts with the partition in two classes separating final and not final states.
- ◉ It is based on the so-called "smaller half" strategy

EXECUTION OF HOPCROFT'S ALGORITHM

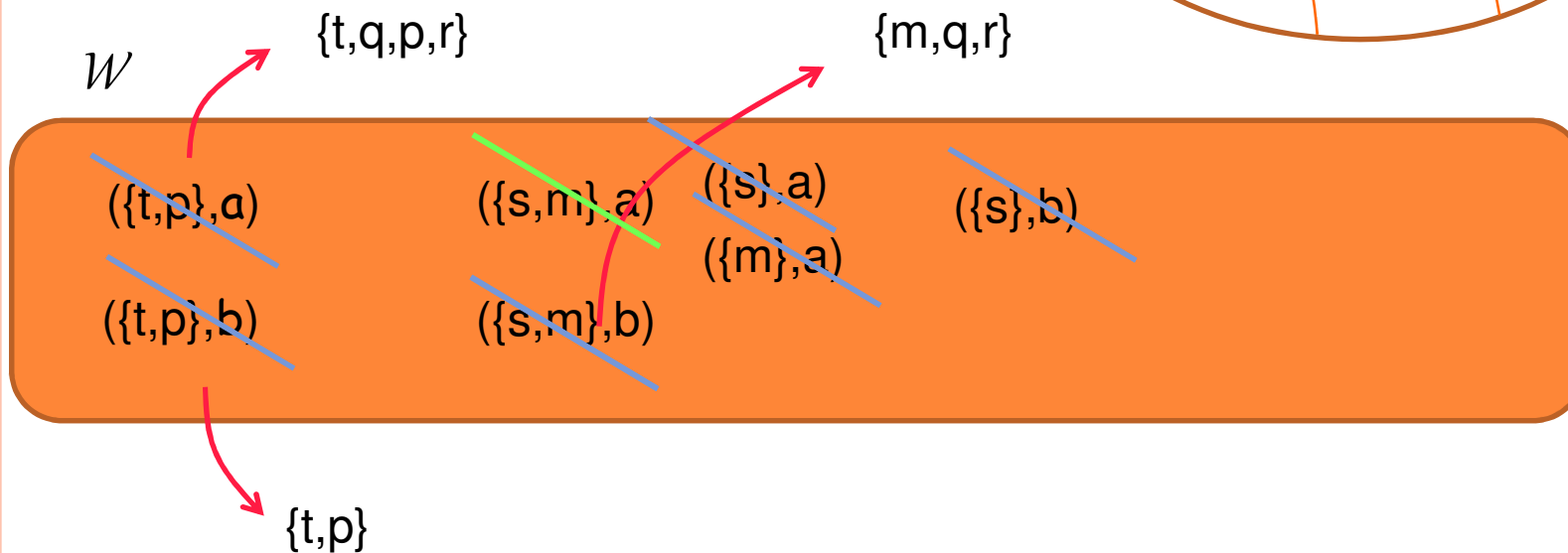
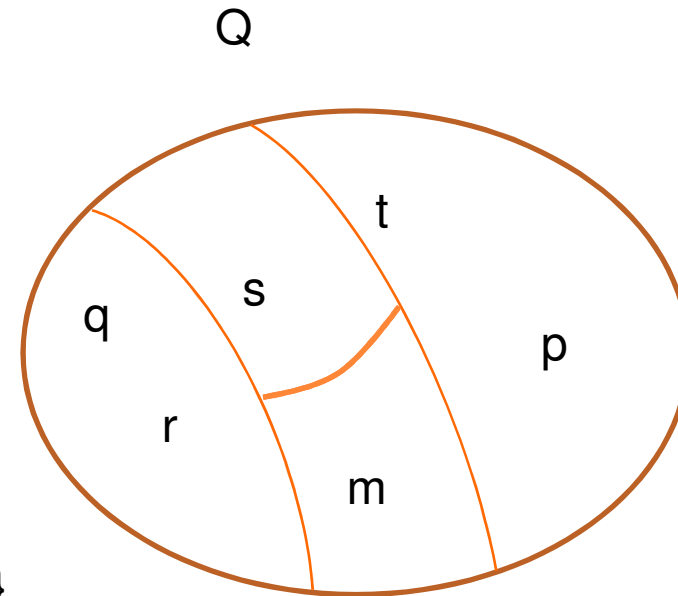
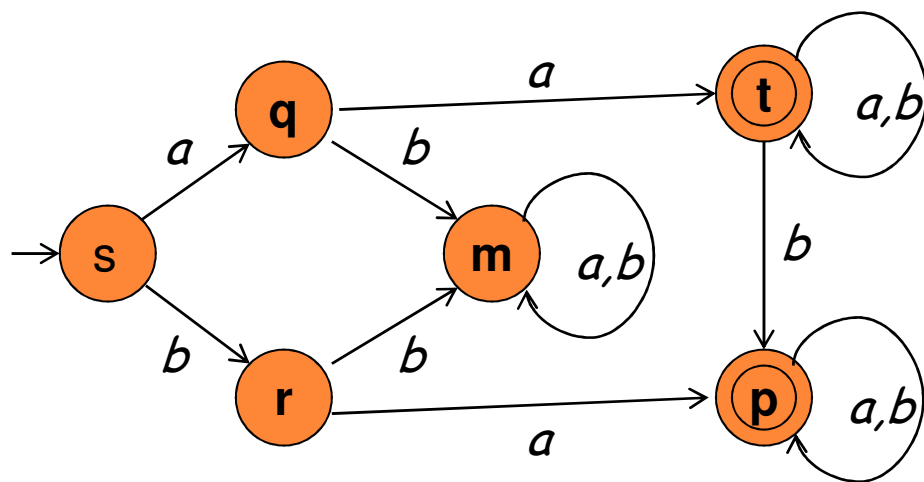
Let Π a partition of Q :

- A pair (C, a) where $C \in \Pi$ and $a \in \Sigma$ is a **splitter** if there exists a class P in Π such that
$$\emptyset \subsetneq \delta^{-1}_a(C) \cap P \subsetneq P$$
- The class P is split in $\delta^{-1}_a(C) \cap P$ and $P \setminus \delta^{-1}_a(C)$
- The smallest set coupled with some symbol of the alphabet goes into the **waiting set**.
- At each step a pair from the waiting set is extracted and the algorithm stops when the waiting set is empty.

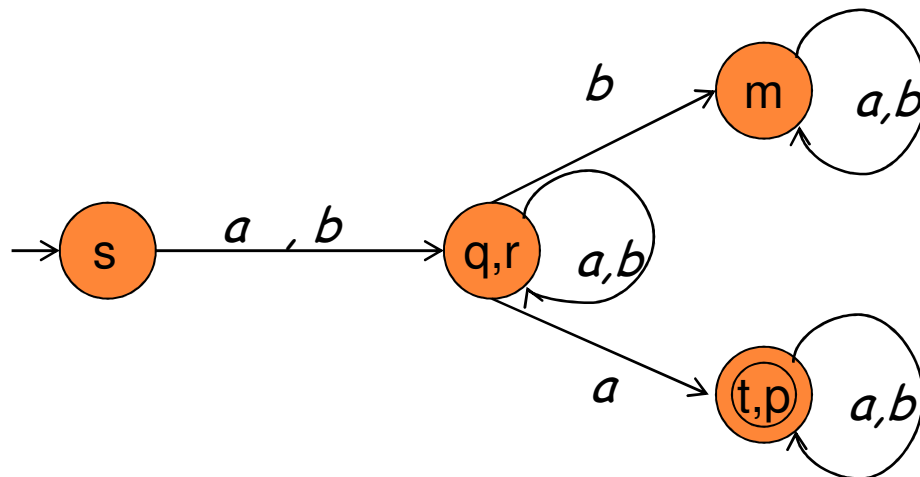
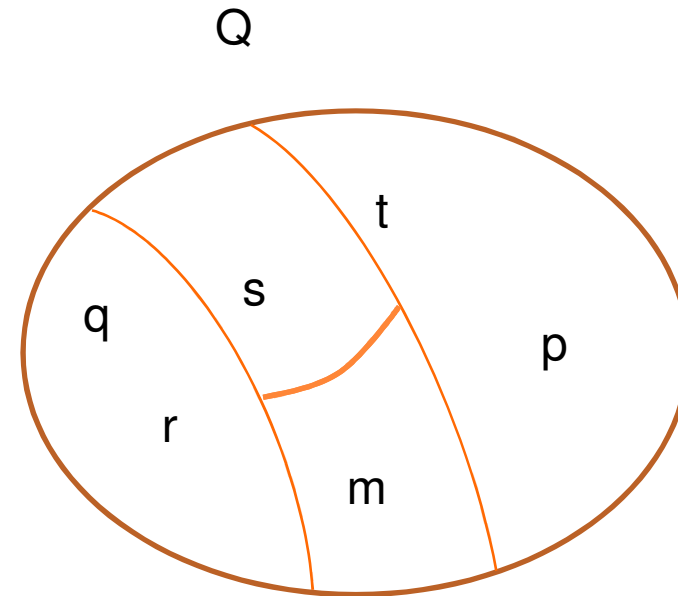
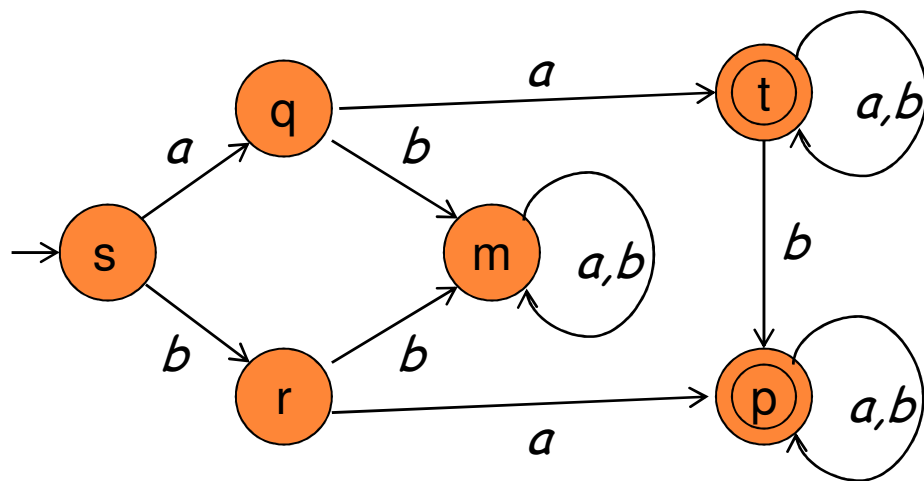
Let $A = (Q, \Sigma, i, \delta, F)$ be a DFA.

- 1: $\Pi \leftarrow \{F, F^c\}$
- 2: **for** all $a \in \Sigma$ **do**
- 3: $\text{Add}((\min(F, F^c), a), W)$
- 4: **while** $W \neq \emptyset$ **do**
- 5: $(C, a) \leftarrow \text{TakeSome}(W)$
- 6: **for** each $P \in \Pi$ which is split by (C, a) **do**
- 7: $P', P'' \leftarrow (C, a) \mid P$
- 8: Replace P by P' and P'' in Π
- 9: **for** all $b \in \Sigma$ **do**
- 10: **if** $(P, b) \in W$ **then**
- 11: Replace (P, b) by (P', b) and (P'', b) in W
- 12: **else**
- 13: $\text{Add}((\min(P', P''), b), W)$

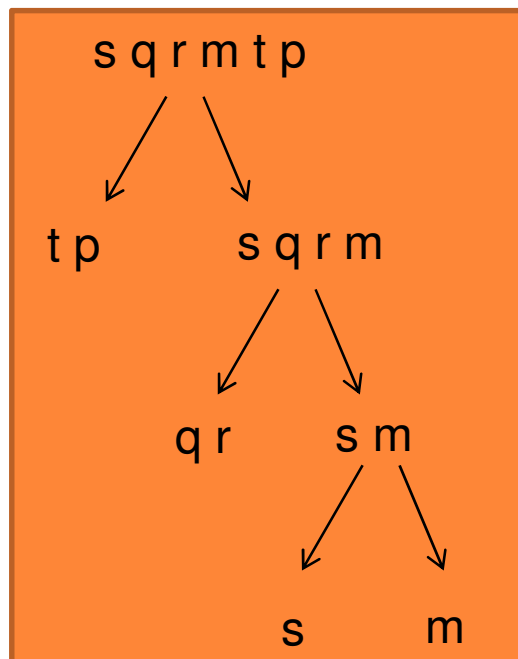
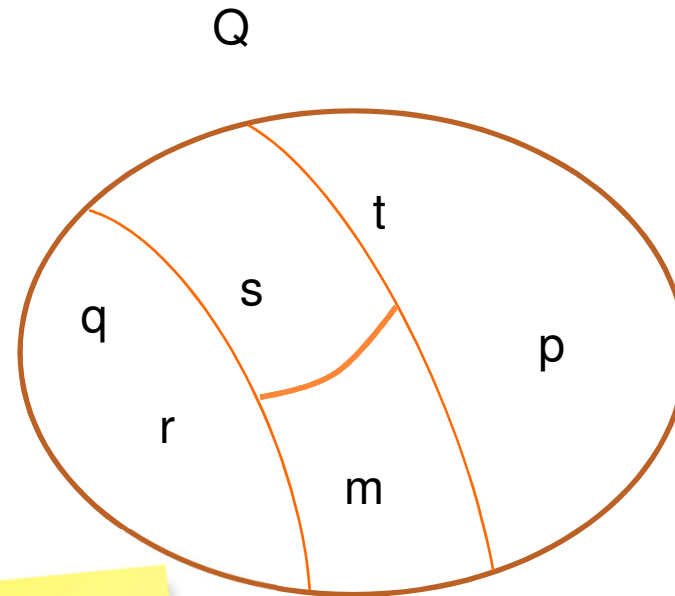
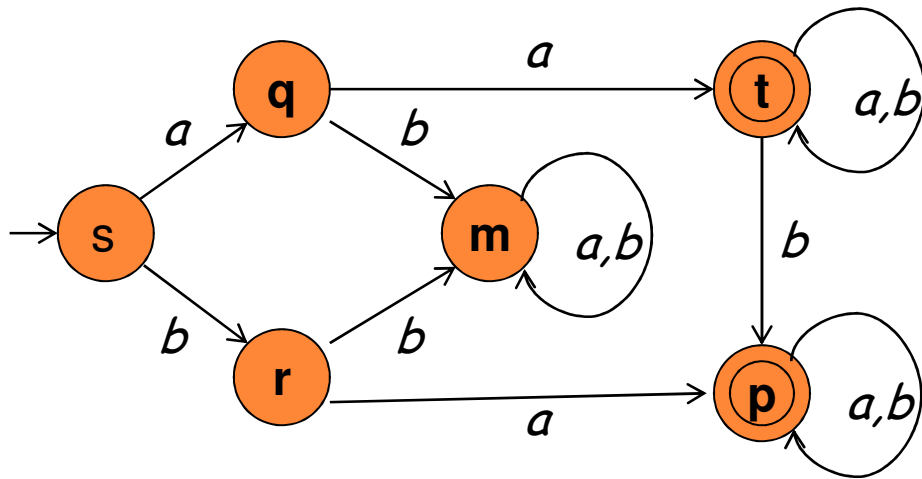
HOPCROFT'S ALGORITHM: AN EXAMPLE



HOPCROFT'S ALGORITHM: AN EXAMPLE



HOPCROFT'S ALGORITHM: THE DERIVATION TREE



The Derivation tree describes the refinement process

GENERAL FEATURES OF THE ALGORITHM

- The refinement process leading to the final partition could be not unique.
- For a given automaton, there could be different executions with different time complexity.

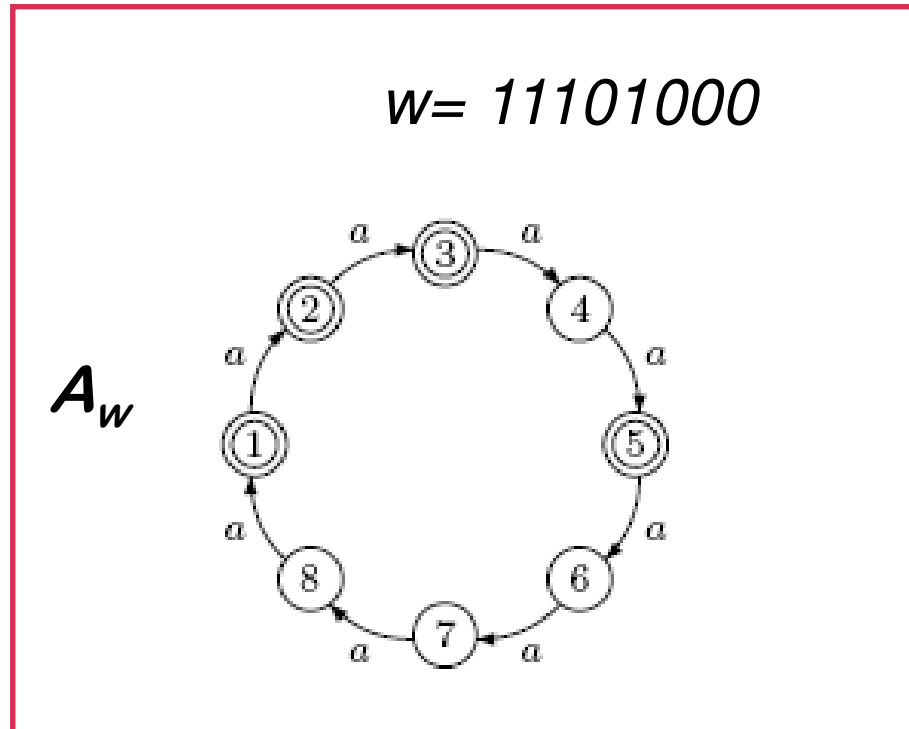
Question: Is the time complexity tight?

[Berstel, Carton 2004]

There exist cyclic unary automata where you need $\Omega(n \log n)$ time if you are "unlucky".

They are related to De Bruijn words.

CYCLIC UNARY AUTOMATA



w is a word
 $a_1 a_2 \dots a_n$ over $A = \{0, 1\}$.

$A_w = (Q, \Sigma, \delta, F)$ the cyclic automaton associated to w :
 $Q = \{1, 2, \dots, n\}$, $\Sigma = \{a\}$

$\delta(i, a) = (i+1)$, $i < n$
 $\delta(n, a) = 1$

$F = \{i \in Q \mid a_i = 1\}$

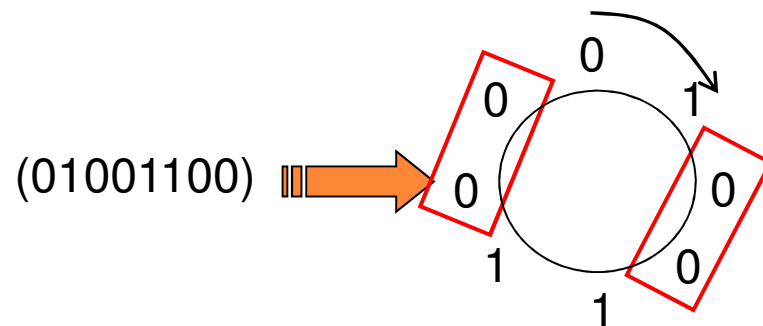
A natural question is:
do there exist infinite families of automata for which
each execution of Hopcroft's algorithm runs in
 $O(n \log n)$?

CIRCULAR WORDS

$A=\{0,1\}$ binary alphabet.

- $u, v \in A^*$ **conjugate**: $\exists z, w$ in A^* s.t. $v=zw$ and $u=wz$
- conjugacy is an equivalence relation:
let $w \in A^*$, by (w) we denote the class of all the conjugates of the word w and we call it **circular word**;
- a **factor** of (w) is a factor of ww of length not greater than $|w|$;
- u is a **special factor** of (w) if both $u0$ and $u1$ are factors of (w) .

Example:



00 is a special factor

CIRCULAR STANDARD WORDS

Recall the notion of **standard word**:

$d_1, d_2, \dots, d_n, \dots$ a sequence of natural integers $d_1 \geq 0, d_i > 0$
 $\{s_n\}_{n \geq 0}$ sequence of words over $A = \{0, 1\}$ where

$$s_1 = 0, s_{n+1} = s_n^{d_n} s_{n-1} \text{ for } n \geq 1$$

Each finite word s_n is called **standard word**.

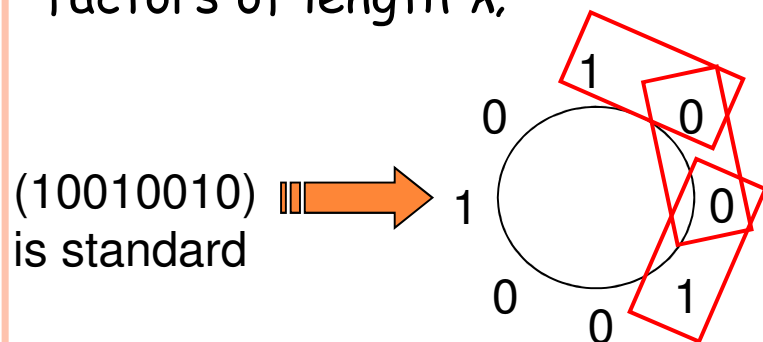
$(d_1, d_2, \dots, d_n, \dots)$ **directive sequence**

$(1, 1, \dots, 1, \dots)$ corresponds to Fibonacci words.

(w) is a **circular standard word** if w is a standard word.

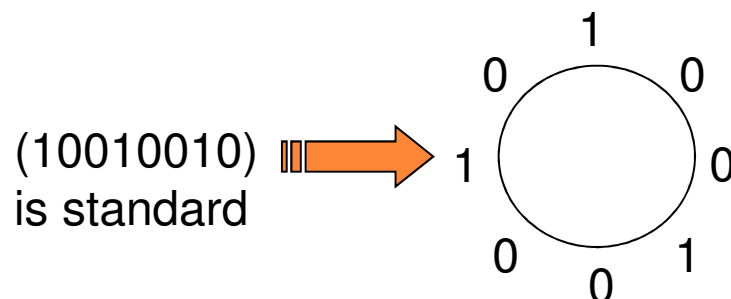
CIRCULAR STANDARD WORDS

Proposition (Borel, Reutenauer). Let v be a word of length $n \geq 2$.
 (v) is a circular standard word iff for $k=0, \dots, n-1$, (v) has exactly $k+1$ factors of length k ;



For $k=2$, there are 3 distinct factors of length 2

Proposition. Let v be a word of length n .
 (v) is standard iff $\forall k=0, \dots, n-2$, (v) has a unique special factor of length k .



Special factors:

| | |
|-------|---------------|
| $k=0$ | ε |
| $k=1$ | 0 |
| $k=2$ | 10 |
| $k=3$ | 010 |
| $k=4$ | 0010 |
| $k=5$ | 10010 |
| $k=6$ | 010010 |

SOME ANSWERS ON UNARY AUTOMATA

- [Castiglione, Restivo, S. 2008]
 - There exist cyclic unary automata for which the execution of Hopcroft's algorithm is unique. They are associated to circular standard words.
 - There exist unary automata where you need always $\Omega(n \log n)$ time. They are associated to circular Fibonacci words.
- [Berstel, Boasson, Carton, 2009]

The same holds for cyclic unary automata associated to all circular standard words having directive sequences for which the sequence of geometric means is bounded.

SOME ANSWERS ON UNARY AUTOMATA

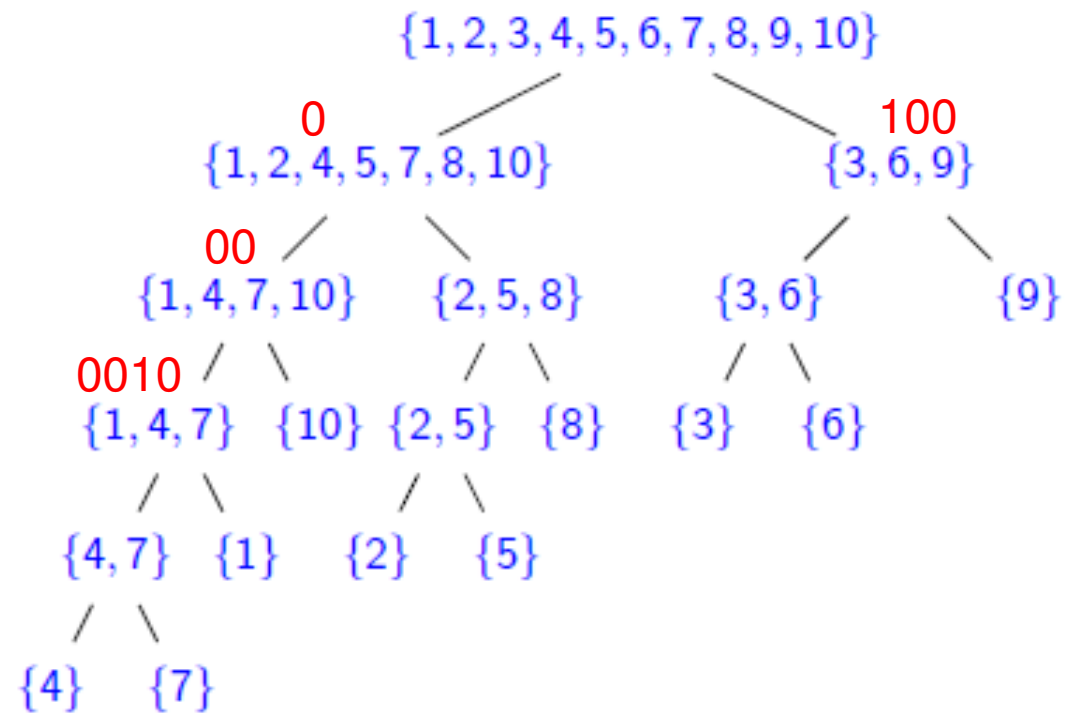
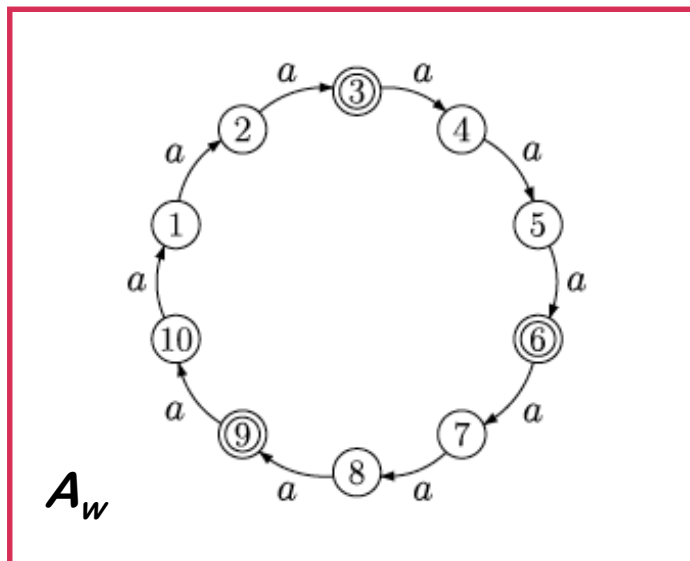
- [Castiglione, Restivo, S. 2008]
 - There exist cyclic unary automata for which the execution of Hopcroft's algorithm is unique. They are associated to circular Fibonacci words.
 - There exist unary automata with $\Omega(n \log n)$ time. They are associated to Fibonacci words.
- [Berstel, Boasson, Carton, 2009]

The same holds for cyclic unary automata. To all circular standard words having a bounded sequence of geometric means is bounded.

IDEA: Splits are related to the occurrences of special factors in the circular word

DERIVATION TREE

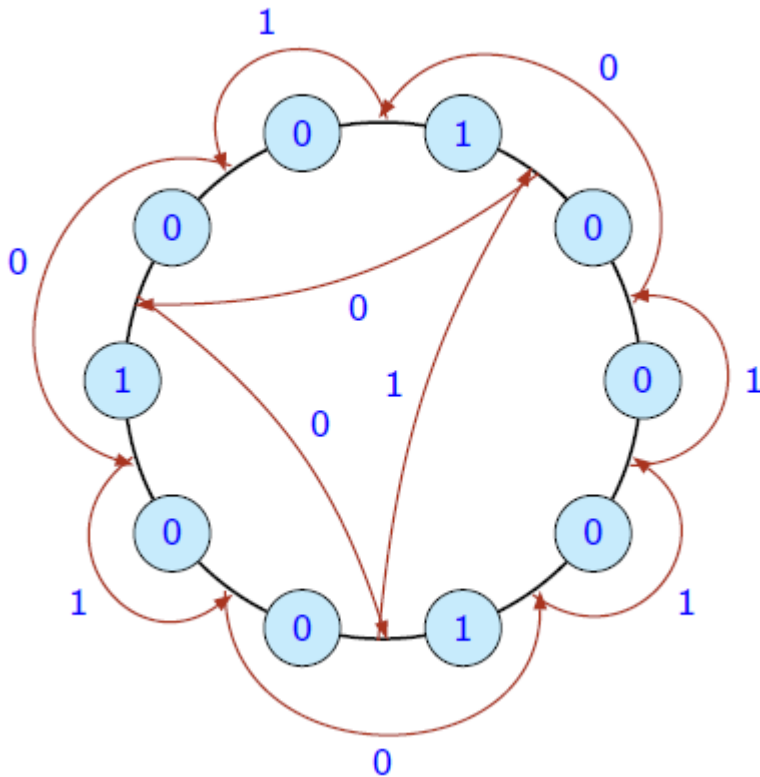
$w = 0010010010$



And so on...

FACTORIZATION OF CIRCULAR STANDARD WORDS

$w = 0010010010$



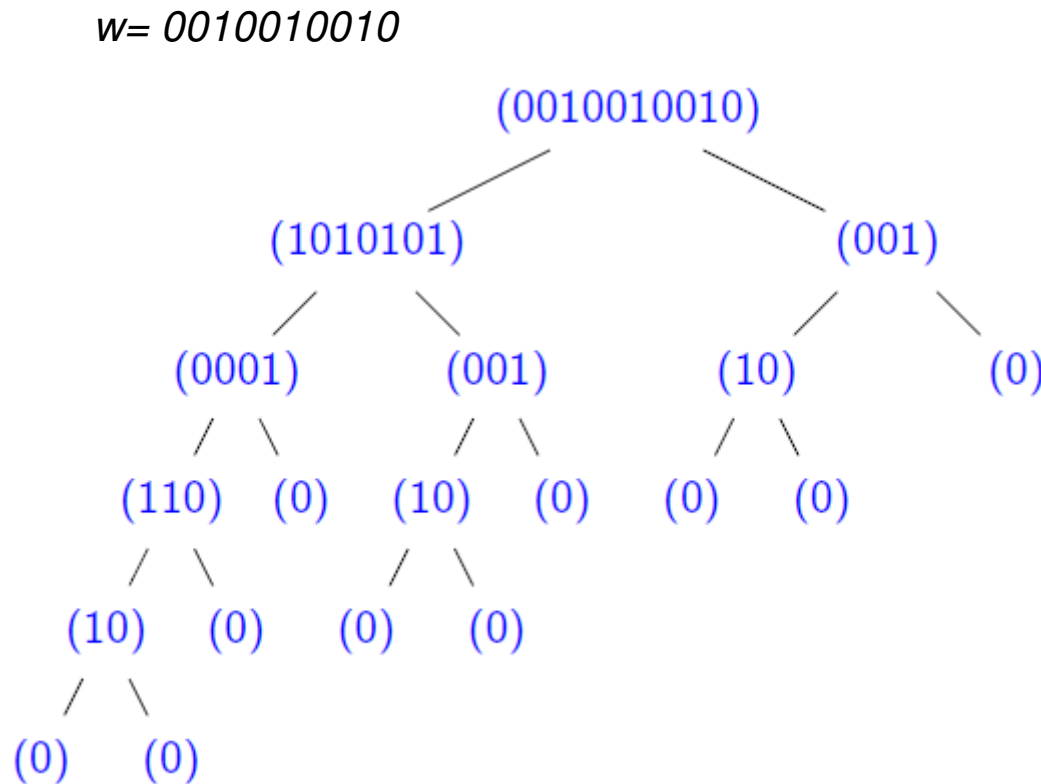
- Each circular standard word can be uniquely "circularly" factored in $\{0, 01\}$ and $\{10^p, 10^{p+1}\}$.
- By using such encoding we obtain again circular standard words.
- Example:

$w = 0010010010 =$

- $0|01|0|01|0|01|0|$
 $\rightarrow 1010101$ denoted by $L(w)$
- $00|100|100|10$
 $\rightarrow 100$ denoted by $R(w)$

REDUCTION TREE OF A CIRCULAR STANDARD WORD

- if $(w)=(0)$ or $(w)=(1)$, $\tau(w)$ is a single node with label (0) ;
- if $|w| > 1$, $\tau(w)$ is a tree with root labeled by (w) having respectively as left and right subtrees $\tau(L(w))$ and $\tau(R(w))$.

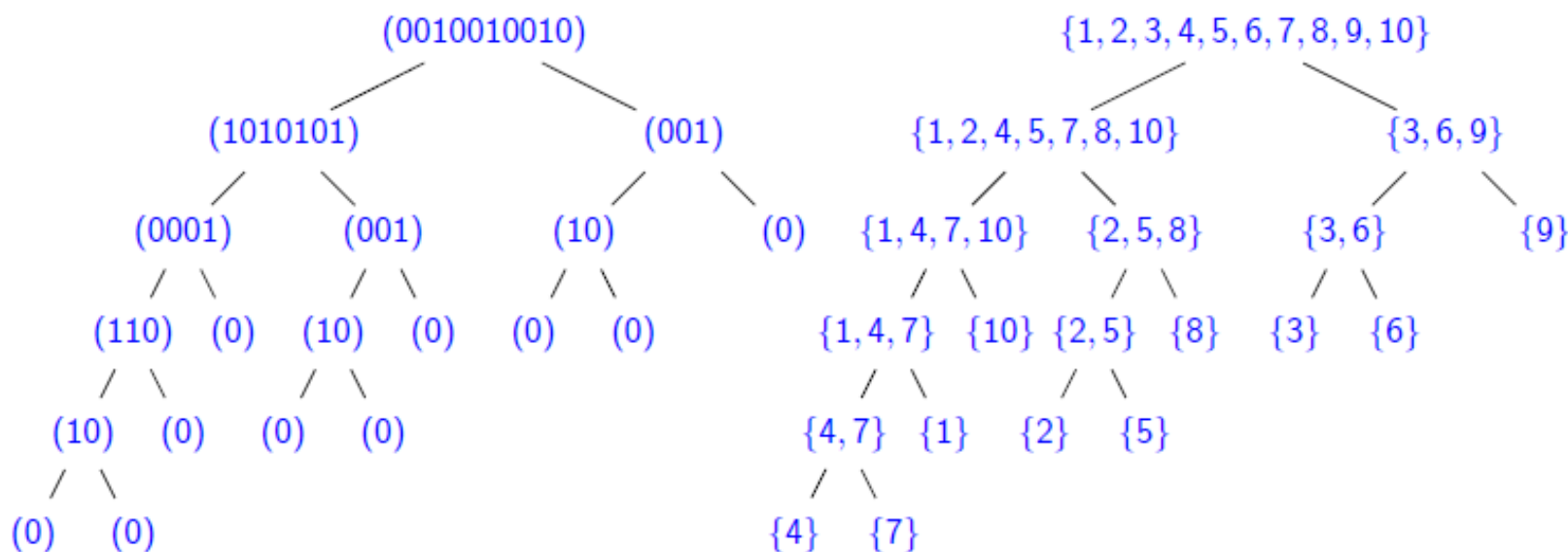


- THEOREM** (Castiglione, Restivo, S. 2009): Let (w) and (v) be circular standard words, then $\tau(w)=\tau(v)$ iff $(w)=(v)$ up to exchanging 0 with 1.

ISOMORPHISM BETWEEN TREES

- THEOREM** (Castiglione, Restivo, S. 2009):
 If (w) is a circular standard word then the
 derivation tree of A_w and the reduction tree are
 isomorphic.

$w = 0010010010$



TIGHTNESS PROBLEM: UNARY CASE?

REMARK:

The unary case is a very special case because automata minimization can be achieved also in linear time when the alphabet has only one letter [Paige, Tarjan, Bonic, 1985]

Do there exist infinite families of binary automata for which each execution of Hopcroft's algorithm runs in $O(n \log n)$?

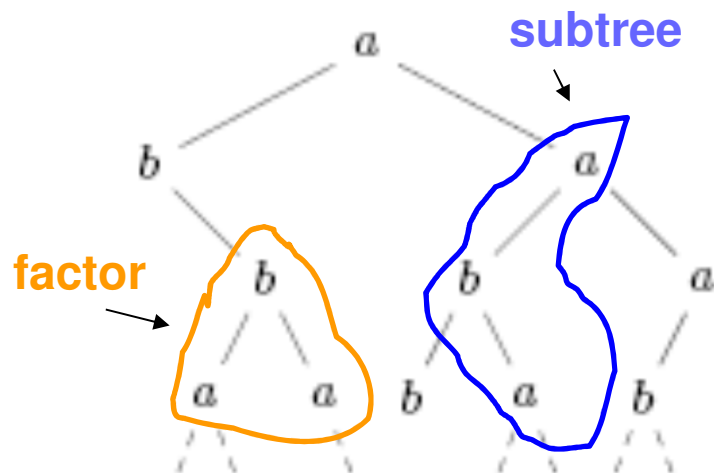
LABELLED BINARY TREES

Let $\Sigma = \{0,1\}$ (alphabet for the shape) and $A = \{a,b\}$ (alphabet for the labels) be two binary alphabets.

A binary labelled tree is a map $\tau: \Sigma^* \rightarrow A$ whose domain $\text{dom}(\tau)$ is a prefix-closed subset of Σ^* .

EXAMPLE

If $\text{dom}(t)$ is finite, the tree is finite, else the tree is infinite.



Height of a finite tree: the maximal length plus 1 of the elements of $\text{dom}(t)$.

Complete infinite tree:
the domain is Σ^* .

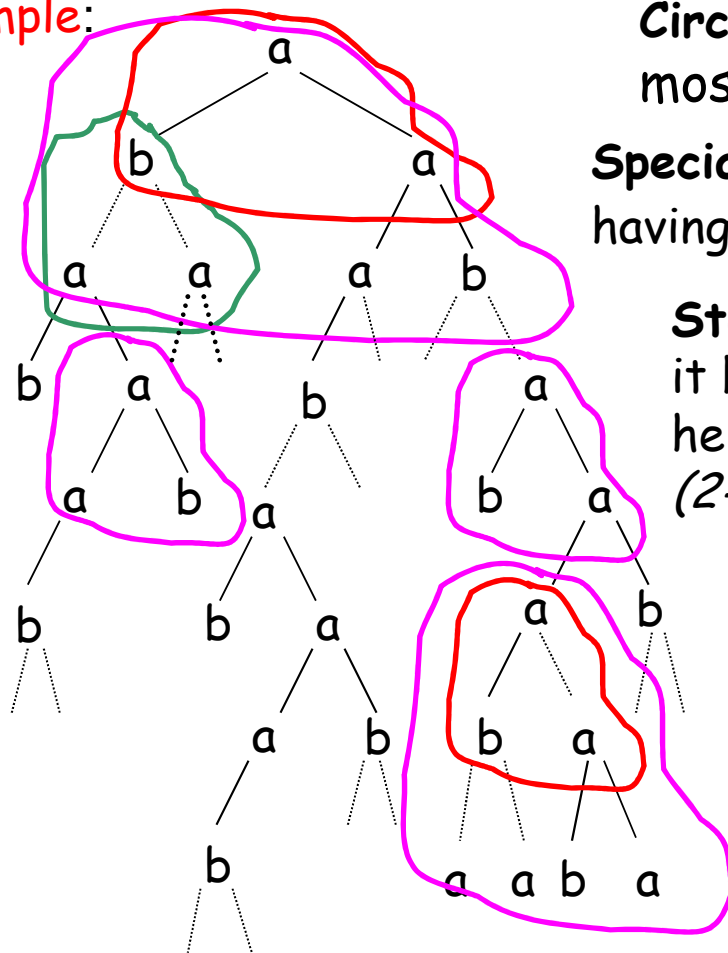
Complete finite tree of height n :
the domain is Σ^{n-1} .

Factor: a finite complete subtree

STANDARD TREES

Given a finite tree τ , we denote by τ^ω the complete infinite tree obtained by recursively concatenating τ to each node of the frontier.

Example:



Circular factor: a factor of τ^ω of height at most $h(\tau)$

Special circular factor: a circular factor having at least two complete extensions in τ^ω

Standard tree: if for each $0 \leq h \leq h(\tau)-2$ it has exactly one special circular factor of height h and it has exactly 2 extensions (*2-special circular factor*).

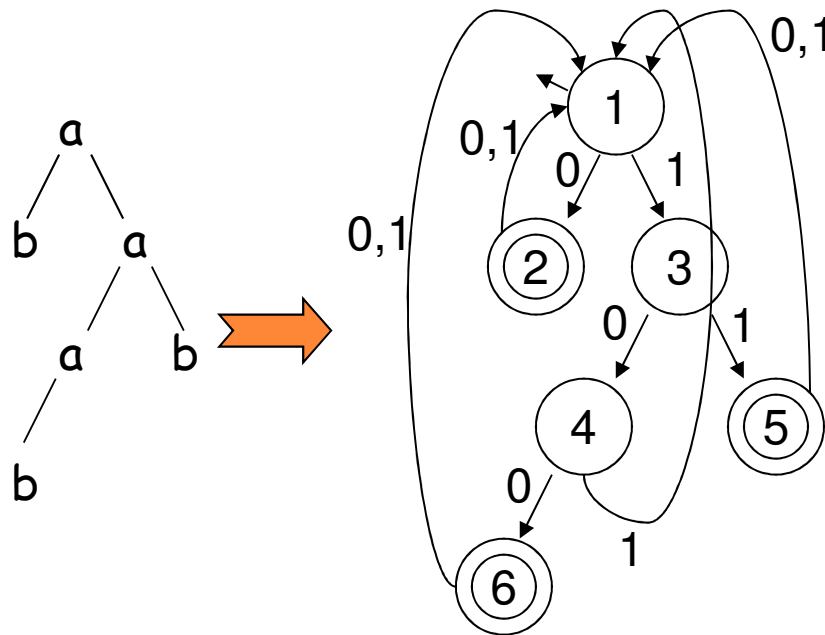
$h=0$ Empty tree

$h=1$ a is 2-special

$h=2$ $\begin{array}{c} a \\ / \quad \backslash \\ b \quad a \end{array}$ is 2-special

TREE-LIKE AUTOMATA

- Given a finite binary labeled tree τ we can uniquely associate a **tree-like automaton** A_τ having τ as skeleton.



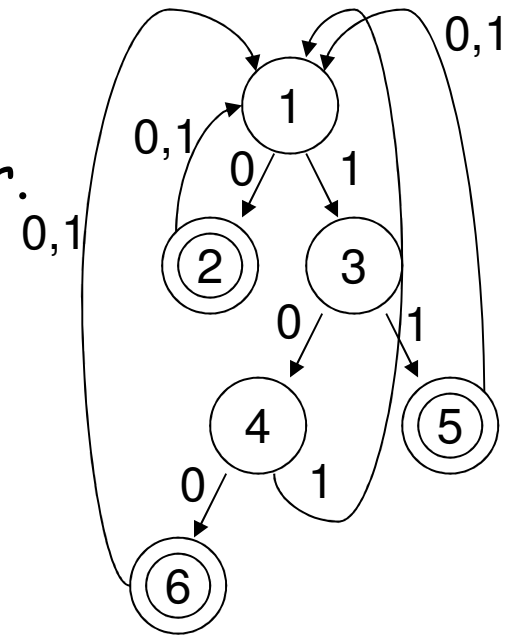
- For each missing edge we add a transition to the root of the tree.
- Moreover, the root is the initial state and the states corresponding to nodes labelled by a (resp. b) are non-final (resp. final) states.

HOPCROFT'S ALGORITHM ON STANDARD TREE-LIKE AUTOMATA

A **standard tree-like automaton** is a tree-like automaton $A_\tau = (Q, \Sigma, i, \delta, F)$ associated to a standard tree τ .

We denote by Q_σ the set of states where the tree σ occurs as a circular factor.

For instance if $\sigma = \begin{array}{c} a \\ b \quad a \end{array}$ then $Q_\sigma = \{1, 4\}$



WORDS 2013 - Turku, September 16-20

THEOREM (Castiglione, Restivo, S. 2010): Let A_τ be a standard tree-like automaton. The sequence of the refinement process $\Pi_1, \Pi_2, \dots, \Pi_m$ is uniquely determined, $m = h(\tau) - 2$ and $\Pi_k = \{Q_\sigma \mid \sigma \text{ is a circular factor of } \tau \text{ with } h(\sigma) = k\}$ and $|\Pi_k| = k + 1$

IDEA OF THE PROOF

The proof is by induction on k .

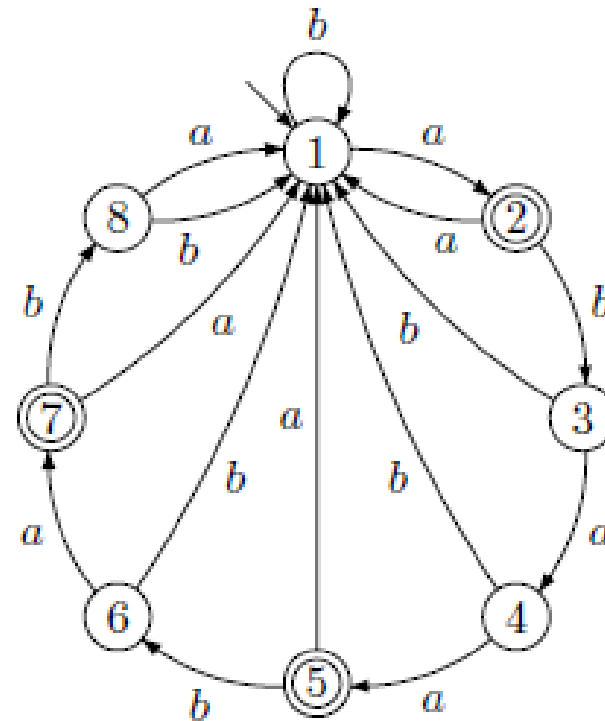
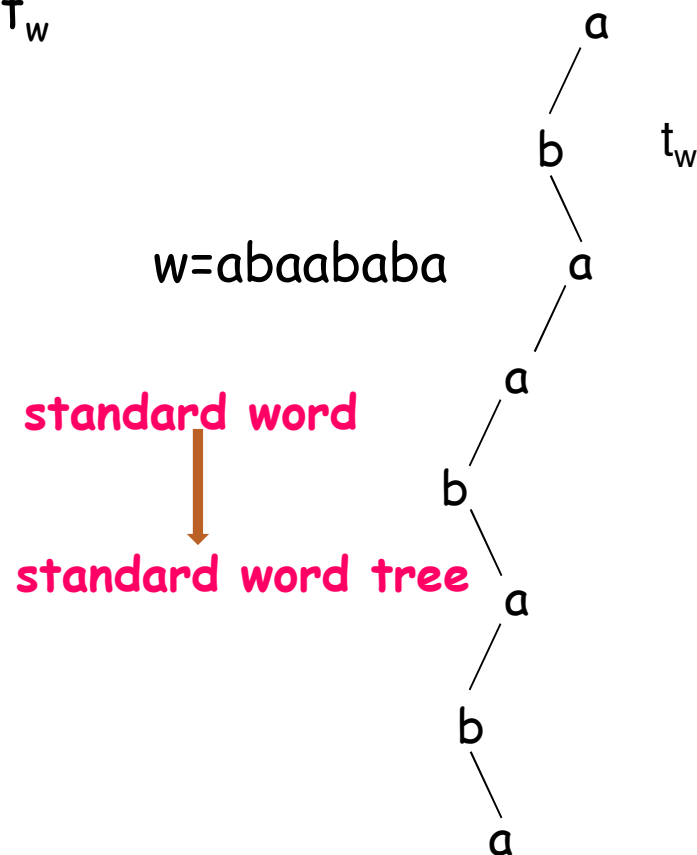
It uses the following result:

Proposition: Let A_τ be a standard tree-like automaton. Let Q_σ and Q_γ be classes of a partition of Q . If (Q_γ, x) splits Q_σ , for some $x \in \Sigma$, with $h(\gamma) = h(\sigma)$, then σ is a 2-special circular factor of τ . The resulting classes are $Q_{\sigma'}$ and $Q_{\sigma''}$, where σ' and σ'' are the only two possible extensions of σ in τ .

We proved that at each step, the splitter in the waiting set satisfies the hypothesis of the proposition.

WORD TREES AND WORD AUTOMATA

Given a word w , a **word tree** t_w is a binary labeled tree in which the alphabet of the shape and the labels coincide. Moreover, $\text{dom}(t_w)$ is the set of prefixes of w . The **word automaton** A_w is the tree-like automaton associated to t_w



TIGHTNESS OF HOPCROFT'S ALGORITHM

We consider the executions of Hopcroft's algorithm on binary word automata.

THEOREM (Castiglione, Restivo, S. 2010): There exists an infinite family of binary automata for which each execution of Hopcroft's algorithm runs in time $\Theta(n \log n)$.

The bound is tight on word automata associated to standard word trees constructed with Fibonacci words.

IN THE PROOF

We prove and use the following:

Proposition: Let τ_w be a standard word tree.

Each execution of Hopcroft's algorithm on $A\tau_w$ has time complexity C_H satisfying

$$\sum_{\sigma \in \text{sp}(\tau_w)} \min(|Q_{\sigma'}|, |Q_{\sigma''}|) + n - 1 \leq C_H(A\tau_w) \leq 2 \times \sum_{\sigma \in \text{sp}(\tau_w)} \min(|Q_{\sigma'}|, |Q_{\sigma''}|)$$

where with $\text{sp}(\tau_w)$ we denote the set of 2-special circular factor of τ_w .

REMARK: If we consider the circular standard word $w = a^n b$, all the executions run in time $\Theta(n)$. If word automata associated to fibonacci words, the time complexity is $\Theta(n \log n)$.

IN THE PROOF

We prove and use the following:

Proposition: Let τ_w be a standard word tree.

Each execution of Hopcroft's algorithm on $A\tau_w$ has time complexity C_H satisfying

$$\sum_{\sigma \in \text{sp}(\tau_w)} \min(|Q_{\sigma'}|, |Q_{\sigma''}|) + n - 1 \leq C_H(A\tau_w) \leq 2 \times \sum_{\sigma \in \text{sp}(\tau_w)} \min(|Q_{\sigma'}|, |Q_{\sigma''}|)$$

where with $\text{sp}(\tau_w)$ we denote the set of 2-children of τ_w .

REMARK

In case of standard words, the circular word automaton is $\Theta(n \log n)$. If word automaton is $\Theta(n \log n)$.

In case of standard words, time complexity does not depend from implementation of the waiting set!

For de Bruijn words, a stack implementation runs in $\Theta(n)$, a queue implementation in $\Theta(n \log n)$.

SOME RESULTS ON AVERAGE

[David 2012]

- For the uniform distribution on complete deterministic automata, the average time complexity of Moore's state minimization algorithm is $O(n \log \log n)$.
- There is a family of implementations of Hopcroft's state minimization algorithm are always faster than Moore's algorithm.

[Bassino, David, Sportiello 2012]

- For the uniform distribution on complete deterministic accessible automata, there exists a family of implementations of Hopcroft's state minimization algorithm whose average complexity is $O(n \log \log n)$.

WHAT ABOUT OTHER MINIMIZATION ALGORITHMS?

- Word automata associated to Fibonacci words are a challenging class for other minimization algorithms.

BRZOWSKI'S METHOD

Given the automaton $A = (Q, \Sigma, \delta, q_0, F)$ recognizing the language L

- With $d(A)$ we denote the deterministic automaton equivalent to A (obtained by the subset construction).
- With $r(A)$ we denote the reverse of A .

The minimal automaton recognizing L is $d(r(d(r(A))))$

TIME COMPLEXITY OF BRZOWSKI'S ALGORITHM

- Since the reverse of an automaton can be computed in linear time with respect to the number of states and transitions, the critical part is the accessible subset construction.
- Time complexity is exponential in the worst case.
- The sum of the size of all accessible states in the subset construction $d(r(A))$, is a lower bound of the time complexity of the accessible subset construction and then a lower bound for the algorithm.

BRZOWSKI'S ALGORITHM AND WORD AUTOMATA

- [Castiglione, Nicaud, S. 2011]
Brzowski's algorithm has a $\Omega(n^2)$ time complexity for word automata constructed by Fibonacci words.

OTHER RESULTS

- De Felice, Nicaud: Brzozowski Algorithm is Generically Super-Polynomial for Deterministic Automata, 2013.
- Brzozowski, Tamm: Minimal Nondeterministic Finite Automata and Atoms of Regular Languages, 2013
- Brzozowski, Tamm: Quotient Complexities of Atoms of Regular Languages, 2012.
- Vazquez de Parga, García, López: A polynomial double reversal minimization algorithm for deterministic finite automata, 2013
- Vazquez de Parga, García, López: DFA minimization: from Brzozowski to Hopcroft, 2013

WHAT ABOUT THE VARIANTS?

- We consider the recent variant proposed by [Garcia, Lopez, Vazquez de Parga, 2013]

PARTIAL REVERSE DETERMINIZATION

Minimization by PRD ($A = (Q; \Sigma; \delta; q_0; F)$)

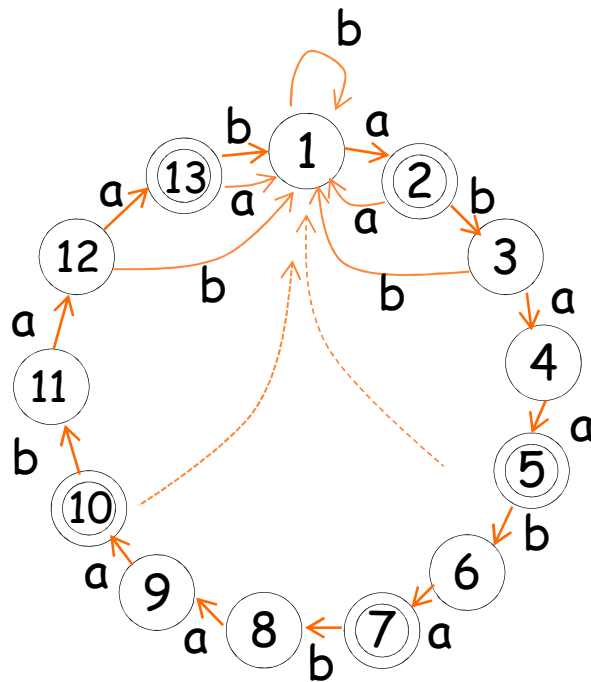
1. $\Pi \leftarrow \{F; Q \setminus F\}$
2. $S = F$
3. $L = \emptyset$
4. **for all** $a \in \Sigma$ **do**
5. $L \leftarrow (S; a)$
6. **while** $L \neq \emptyset$ **do**
7. *choose and delete a pair $(S; a)$*
8. **for all** $B \in \Pi$ **do**
9. **if** B **is split by** $(S; a)$ **then**
10. $B' \leftarrow \delta_a^{-1}(S) \cap B$
11. $B'' \leftarrow B \setminus \delta_a^{-1}(S)$
12. $\Pi \leftarrow \Pi \setminus B \cup \{B' \cup B''\}$
13. **for all** $b \in \Sigma$ **do**
14. $L \leftarrow L \cup \{(\delta_a^{-1}(S), b)\}$



Partial reverse
determinization is
polynomial!

PRD ALGORITHM ON WORD AUTOMATA

- We consider the word automaton associated to fibonacci word f_n
- Example $f_6 = \text{abaababababab}$

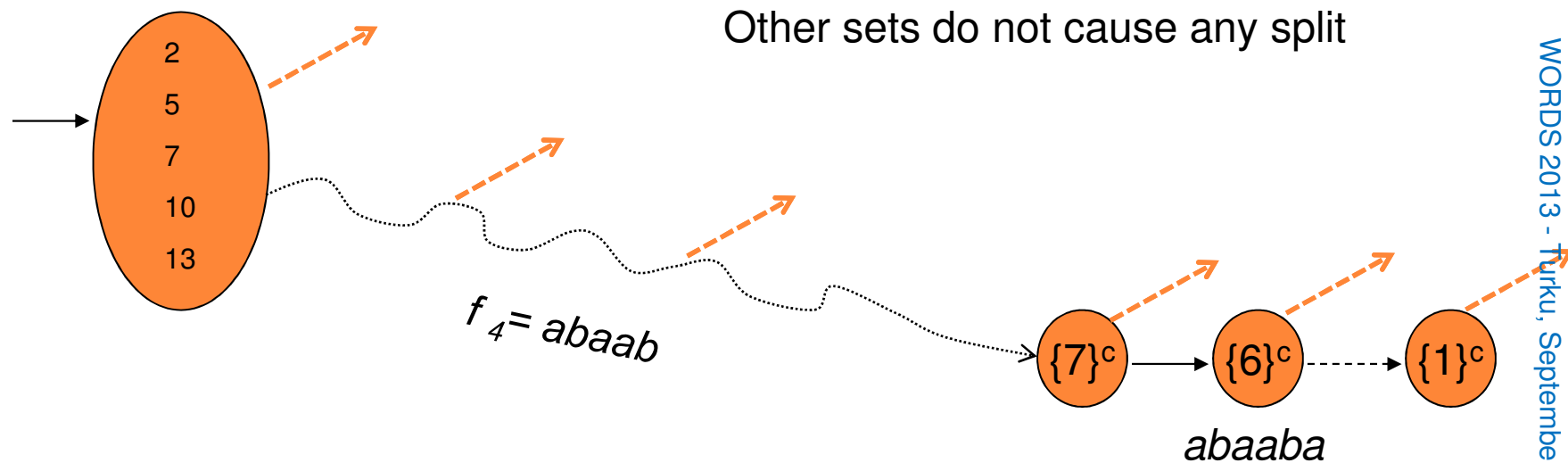


In order to find the time complexity we have to consider the set of states going into L

PRD ON WORD AUTOMATA

If n is even, if we consider a set of δ^{-1} operations by using the letters of f_{n-2} , we obtain sets that are complement of singletons

Ex. $f_6 = \text{abaababaabaab}$

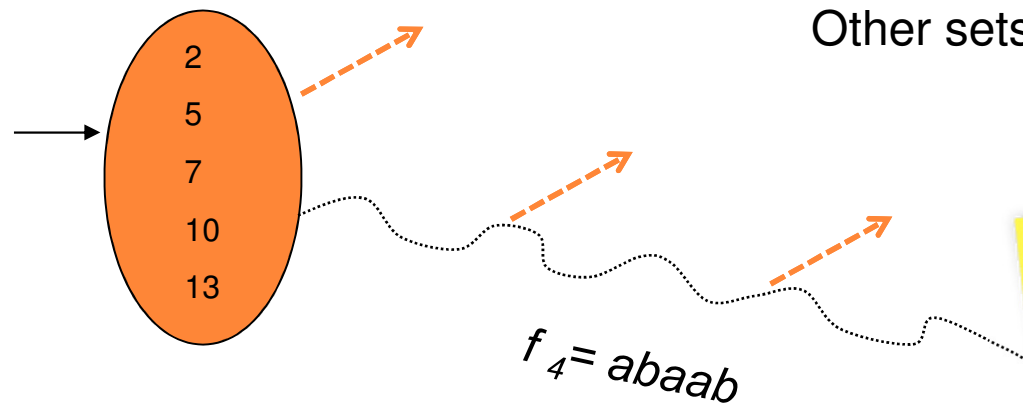


- At each step only a subset causes a split and then goes into the waiting set L .
- The subsets $\{F_{n-1}-1\}^c, \dots, \{1\}^c$ are obtained as inverse image by δ and cause a split.
- The waiting set contains $F_{n-1}-1$ sets of cardinality F_n-1

PRD ON WORD AUTOMATA

If n is even, if we consider a set of δ^{-1} operations by using the letters of f_{n-2} , we obtain sets that are complement of singletons

Ex. $f_6 = \text{abaababaabaab}$



Analogous reasoning can be done when n is odd.

- At each step only a subset causes a split and the waiting set L .
- The subsets $\{F_{n-1}-1\}^c, \dots, \{1\}^c$ are obtained as i and cause a split.
- The waiting set contains $F_{n-1}-1$ sets of cardinality F_n-1

TIME COMPLEXITY FOR STANDARD WORD AUTOMATA

- THEOREM: Let A_w the word automaton associated to a circular standard word. The minimal automaton is obtained by PRD algorithm by a sequence of $\delta^{-1}_{b_1}, \delta^{-1}_{b_2}, \dots, \delta^{-1}_{b_{|w|-2}}$ where $b_1b_2\dots b_{|w|-2}$ is the prefix of length $|w|-2$ of w .
- THEOREM: Let A_{f_n} the word automaton associated to the n -th Fibonacci word f_n with F_n states. Then

$$C_{\text{PRD}}(F_n) = \Theta(F_n^2)$$

They represent the worst case of the algorithm!

PROBLEMS AND PERSPECTIVES

- Are Fibonacci word automata challenging for the minimization process? Is there any better minimization algorithm for DFA?
- Characterizing the words for which the correspondent word automata are always difficult to be minimized.

THANKS FOR YOUR ATTENTION...!!!



WORDS 2003
The first time
I was in Turku!