

Path-Equivalent Removals of ε -transitions in a Genomic Weighted Finite Automaton

Mathieu Giraud, Philippe Veber, and Dominique Lavenier

IRISA / CNRS / Université de Rennes 1
35042 Rennes Cedex, France
{mgiraud, pveber, lavenier}@irisa.fr

Abstract. Weighted finite automata (WFA) are used with accelerating hardware to scan large genomic banks. Hardwiring such automata raise surface area and clock frequency constraints, requiring efficient ε -transitions-removal techniques. In this paper, we present new bounds on the number of new transitions for several ε -transitions-removal problems. We study the case of acyclic WFA. We introduce a new problem, the partial removal of ε -transitions while accepting short chains of ε -transitions.

1 Introduction

Weighted Finite Automata (WFA) are used to find occurrences of biological patterns in genomic databases containing tens of gigabytes of data. Biological patterns can be seen as regular or weighted expressions over the 20-letter amino acid alphabet. They may represent the signature of a protein family, the features of a domain or the specific location of an active site. The usual length ranges of the patterns are from a few amino acids to a few tenths.

WFA can be efficiently hardwired onto reconfigurable architectures (namely FPGA components) to speed up the search of biological patterns, reducing computation time from hours to minutes [1]. Today, with the exponential growth of genomic data, the hardware WFA alternative offers an interesting approach compared to pure software implementation.

Hardware speed comes from the ability to compute all WFA states simultaneously. Actually, genomic data (input string) are processed on-the-fly, and the performance of a hardwired WFA is mainly determined by the input data rate. Thus, the processing time becomes independent of the WFA size, and is only dictated by the time for accessing all the items of the database.

This scheme is valid as long as the WFA fits into FPGA components. Unfortunately, biological patterns may require consequent reconfigurable resources, particularly when insertion/deletion errors are considered. In that case, insertions are modeled by cyclic transitions and deletions by ε -transitions. Resulting WFA are thus much larger in terms of the number of transitions. From a hardware point of view, the resources are directly related to the number of transitions to hardware. Hence, finding equivalent automata with less transitions is highly beneficial.

Beside the automaton size, a direct hardware implementation of ε -transitions is not realistic. Fig. 1 exemplifies the hardware mapping of a WFA with ε -transitions. Paths with ε -transitions are represented by dotted lines: they systematically bypass state registers. The main consequence is that a long *critical path* (dashed line) is created from the input to the output. The critical path is defined as the longest path between two registers, and determines the maximum clock frequency of the circuit. The longer the path, the lower the frequency. Hence, to keep a reasonable working frequency, the critical path needs to be broken into smaller parts by removing some ε -transitions.

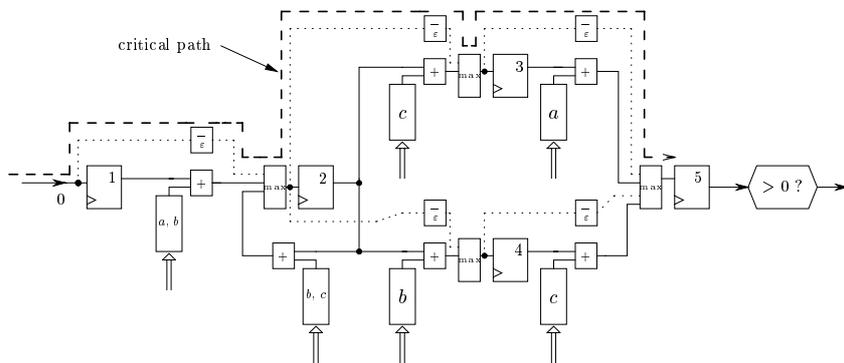


Fig. 1. Hardwiring a WFA with 5 regular transitions doubled with ε -transitions [2]. A critical path runs through the whole automaton.

The classical method removing ε -transitions in automata use the ε -closure of every state [3,4]. Recently, for WFA, Mohri proposed a generic algorithm with a smallest distance method [5]. A certain condition must be checked to ensure that the weights are well-defined in cycles.

These algorithms can raise the number of transitions from n to $\mathcal{O}(n^2)$. The resulting automaton can be minimized [6], but for large automata, such a limit makes the hardware implementation impossible. As an example, in [7], we experienced an 80-state automaton for discovering olfactory receptor genes in the dog genome. On this automaton, the classical ε -transitions-removal algorithms produce more than 3100 new transitions. This number reaches the limit of today FPGA's technology and prevents larger automata from being hardwired.

Hromkovic proposed a study for ε -transitions in finite automata [8]. There are rational expressions of size $\mathcal{O}(n)$ such that every ε -free recognizing automaton has a size $\Omega(n \log n)$. Lifshits raised this bound to $\Omega(n \log^2 n / \log \log n)$ [9]. Other works optimized the creation time of those automata [10].

In this paper we study the *development* of WFA: we double every transition with an ε -transition, and we study the number of new transitions created when removing the ε -transitions. We previously proposed a first study for linear-shaped automata: in this case, we designed an optimal method that produces automata with $\Theta(n \log n)$ new transitions [2].

The rest of the paper is organized as follows. Section 2 provides WFA background. Then, in Section 3, we study the development of WFA for acyclic automata. Section 4 presents a new problem driven by the hardware constraints: the removal while accepting short ε -chains. The final section concludes with experimental results and perspectives.

2 Background

2.1 WFA and Pattern Matching

Definition 1. A *Weighted Finite Automaton (WFA)* is a 5-tuple $\mathcal{A} = (Q, \Sigma, \Delta, I, F)$, where Q is a finite set of states, Σ a finite alphabet, $\Delta \subset Q \times Q \times (\delta : \Sigma \text{eps} \mapsto \mathbb{K})$ a finite transition table, $I \subseteq Q$ and $F \subseteq Q$ the sets of initial and final states.

The number of transitions of the WFA is $|\Delta|$. For each transition $\tau = (q, q', \delta) \in \Delta$, we denote by $i[\tau] = q$ its initial state, $f[\tau] = q'$ its final state, and $\delta[\tau] = \delta$ its weight function. A WFA without ε -transitions is a WFA such that $\delta(\varepsilon) = -\infty$ for every transition (q, q', δ) . Now we define paths as consecutive labeled transitions:

Definition 2. A *path* $\pi = (\tau_1, \alpha_1) \dots (\tau_k, \alpha_k) \in (\Delta \times (\Sigma \text{eps}))^*$ in a WFA \mathcal{A} is a succession of pairs of transitions and characters where the transitions $\tau_1 \dots \tau_k$ are consecutive transitions, that is $f[\tau_i] = i[\tau_{i+1}]$ for $i = 1 \dots k-1$, and where the characters α_i are in Σeps . The *label* of π is the word $\alpha_1 \dots \alpha_k$.

The weight function δ can be extended to paths: for a path $\pi = (\tau_1, \alpha_1) \dots (\tau_k, \alpha_k)$, we define $\delta(\pi) = \delta[\tau_1](\alpha_1) + \dots + \delta[\tau_k](\alpha_k)$. Weights on words used in pattern matching are computed as weights on paths between some initial and final states.

2.2 Path-Equivalence

Now we give a definition of our ε -transition-removal problem. We define it as finding a new automaton with a special kind of equivalence, the path-equivalence, which requires that some paths (the closed paths, see below) have a superior path in the corresponding automaton.

Definition 3. One path π is *superior* to another one π' if both paths have the same label, the same initial state and the same final state, and if $\delta(\pi) \geq \delta(\pi')$.

Definition 4. A path $\pi = (\tau_1, \alpha_1) \dots (\tau_k, \alpha_k)$ is *left-closed* if it begins with an initial state ($i[\tau_1] \in I$) or if its first character α_1 is different than ε . Similarly, a path is *right-closed* if $f[\tau_k] \in F$ or $\alpha_k \neq \varepsilon$. A path is *closed* if it is closed at both sides.

Definition 5. Two WFA $\mathcal{A} = (Q, \Sigma, \Delta, I, F)$ and $\mathcal{A}' = (Q, \Sigma, \Delta', I, F)$ are *path-equivalent* if every closed path in \mathcal{A} labeled by a word $w \neq \varepsilon$ has a superior path in \mathcal{A}' and reciprocally.

Basically, the path-equivalence states that the two automata simulate each other through their paths. Usual algorithms that remove the ε -transitions such as [4] or [5] produce path-equivalent automata.

2.3 Development of an Automaton

Given a WFA without ε -transitions $\mathcal{A} = (Q, \Sigma, \Delta, I, F)$ and a deletion cost c_ε , we define \mathcal{A}_ε as the WFA in which all transitions of \mathcal{A} are doubled by ε -transitions. More precisely, every transition $(q, q', \delta) \in \Delta$ is extended with $\delta(\varepsilon) = c_\varepsilon$.

Definition 6. *Given a WFA \mathcal{A} , any WFA \mathcal{A}' is a development of \mathcal{A} if \mathcal{A}' is path-equivalent to \mathcal{A}_ε and if has no ε -transitions. We say that \mathcal{A}' is developed from \mathcal{A} if \mathcal{A}' is a development of \mathcal{A} .*

To be efficiently harwired, a WFA needs to be developed with *as few new transitions as we can*. In the general case, the ε -transitions-removal from an automaton with n transitions gives an automaton with $\mathcal{O}(n^2)$ new transitions. In [2], we studied the case of linear-shaped automata. We designed an optimal method that produces automata with $\Theta(n \log n)$ new transitions.

3 Removal in Acyclic Automata

Here we use the results on linear-shaped WFA to analyze the number of new transitions in the developments of some more generic automata. To ensure that the weights are well defined, automata with cycles require special constraints [5]. The section 3.1 considers *acyclic automata* with n states : we give an upper bound to develop such automata. The section 3.2 extends the result to automata with cycles, but with no cycles on ε -transitions. Such automata are common in biological applications (Fig. 2).

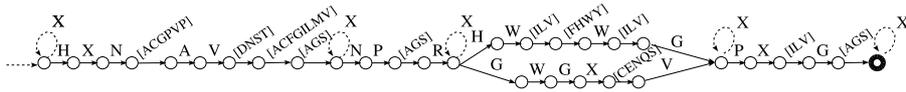


Fig. 2. Detail of a genomic automaton recognizing MIP membrane proteins [11]. The complete automaton has more than 300 transitions. Except for some insertion transitions (X), this automaton is acyclic.

3.1 Acyclic Automata

Definition 7. *A WFA $\mathcal{A} = (Q, \Sigma, \Delta, I, F)$ is acyclic if its graph has no cycle. The states of an acyclic WFA can be numbered q_1, q_2, \dots, q_n such than there is no backward transition $(q_i, q_j, \delta_{i,j})$ with $i \geq j$.*

We call such a WFA a *numbered automaton*. The following algorithm develops a numbered automaton with n states from the development of two sub-automatons obtained by cutting the automaton at a state q_2 .

Algorithm 1. Development of a numbered WFA

Input: a numbered WFA with n states $\mathcal{A} = (Q, \Sigma, \Delta, I, F)$, an integer $z \in [2, n - 1]$, a cost c_ε

Let \mathcal{C} be the set of all cut transitions $(q_i, q_j, \delta_{i,j})$ with $i < z < j$

Let Z be a set of states touching \mathcal{C}

Let $\mathcal{A}_1 = (Q_1 = \{q_1 \dots q_z\} \cup Z, \Sigma, \Delta_1, I, \{q_z\} \cup Z)$

and $\mathcal{A}_2 = (Q_2 = \{q_z \dots q_n\} \cup Z, \Sigma, \Delta_2, \{q_z\} \cup Z, F)$

where the transition tables Δ_1 and Δ_2 are the restrictions of Δ on Q_1 and Q_2

Let \mathcal{A}'_1 and \mathcal{A}'_2 recursively be two developments of \mathcal{A}_1 and \mathcal{A}_2

Let \mathcal{A}' be the concatenation of \mathcal{A}'_1 and \mathcal{A}'_2 : $\mathcal{A}' = (Q, \Sigma, \Delta', I, F)$, $\Delta' = \Delta'_1 \cup \Delta'_2$

For all q_i in Q_1

Add to Δ' the transition (q_i, q, δ'_i) for all final states $q \in F$

with $\delta'_i(\alpha) = \max_{i+1 \leq k \leq n} [(n - i - 1)c_\varepsilon + \delta_k(\alpha)]$

For all q_i in Q_2

Add to Δ' the transition (q, q_i, δ''_i) for all initial states $q \in I$

with $\delta''_i(\alpha) = \max_{1 \leq k \leq i} [(i - 1)c_\varepsilon + \delta_k(\alpha)]$

Output: the WFA $\mathcal{A}' = (Q, \Sigma, \Delta', I, F)$

In the algorithm for linear-shaped WFA (Algorithm 1 in [2]), initial and final states of both sub-automata guarantee that the paths are closed. Here some transitions are *cut* over q_z (Fig. 3). All the paths are closed if one adds to each sub-automaton a set of states Z that *touches* the cut transitions, that is a set Z such that any cut transition starts or ends in Z . Each state in Z is a final state for the left sub-automaton and an initial state for the right one : the sub-automata are overlapping. We have the following property:

Property 1. *The algorithm 1 builds an automaton which is path-equivalent to the initial automaton.*

Proof. We just give the sketch of the proof, which is similar to the case of linear-shaped WFA (Lemma 3 in [2]). $\boxed{\mathcal{A} \mapsto \mathcal{A}'}$ Each closed path of \mathcal{A} not labeled by ε and not completely included in \mathcal{A}_1 or in \mathcal{A}_2 can be written as $\pi_1 \pi_2$, where π_1 and π_2 are closed paths in \mathcal{A}_1 and \mathcal{A}_2 . Any such decomposition leads to a superior closed path in \mathcal{A}' . $\boxed{\mathcal{A}' \mapsto \mathcal{A}}$ Reciprocally, any closed path of \mathcal{A}' either goes through a state $q \in \{q_z\} \cup Z$, or jumps over such a state. In both case a superior closed path of \mathcal{A} can be reconstructed.

Each step of the algorithm adds no more than $|Q_1| \cdot |F| + |Q_2| \cdot |I|$ transitions. To bound this value, we need a bound on $|Z|$.

Definition 8. *Let be a numbered WFA with states $\{q_1, q_2 \dots q_n\}$, and q_z a state. The width κ_z is the number of transitions (q_a, q_b, δ) with $a < z < b$.*

The *maximal width* is $\mathcal{K} = \max_i \kappa_i$: it can be seen as the maximal number of branches in the WFA, except the main branch. On the automaton depicted on

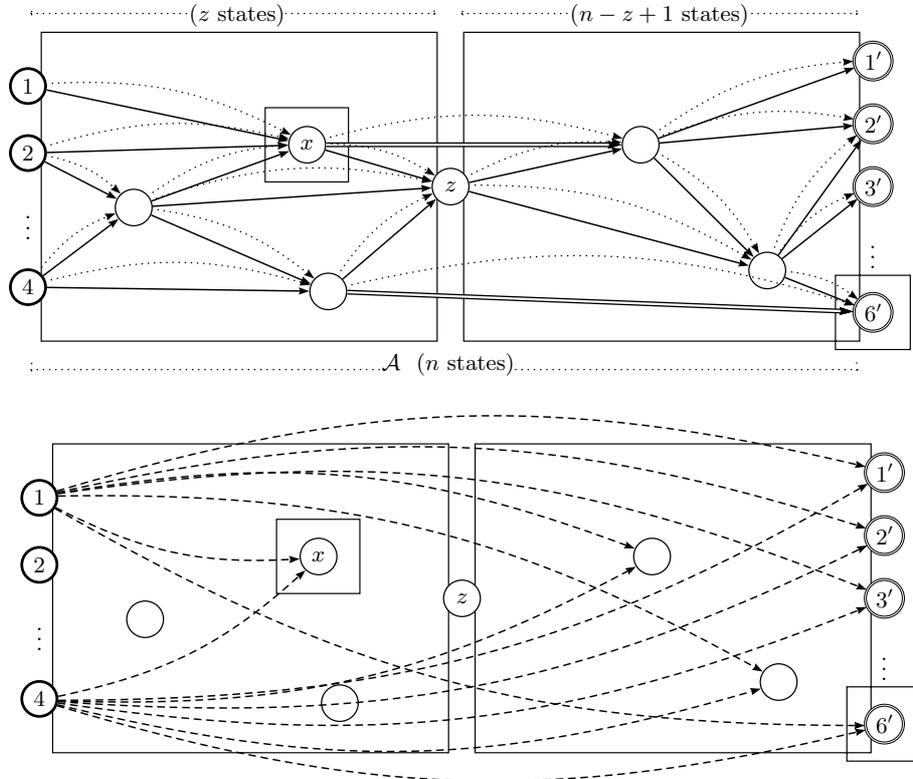


Fig. 3. Algorithm developing a numbered WFA with n states. A state q_z is chosen to split the automaton into two parts with z and $n - z + 1$ states. The two cut transitions are shown in double lines. The set $Z = \{x, 6'\}$ touches every cut transition. This set Z is added to the two parts to give the sub-automata \mathcal{A}_1 and \mathcal{A}_2 . Final states of \mathcal{A}_1 (and initial states of \mathcal{A}_2) are $\{z\} \cup Z$. At the bottom, we add to the developments of the two sub-automata transitions from initial states of \mathcal{A}_1 to all states of \mathcal{A}_2 . With the symmetrical operation, no more than $|Q_1| \cdot |F| + |Q_2| \cdot |I|$ transitions are created.

Fig. 2, we have $\mathcal{K} = 1$ for all numberings. In the general case, the widths depend on the chosen numbering.

At each step, the set Z has no more than \mathcal{K} elements. When applying recursively algorithm 1, the sets I and F will always have no more than $\mathcal{K} + 1$ elements. Then one step of the algorithm adds no more than $(|Q_1| + |Q_2|) \cdot (\mathcal{K} + 1) \leq (n + \mathcal{K}) \cdot (\mathcal{K} + 1)$ transitions. We thus have the following consequence of the property 1:

Property 2. Any numbered WFA with a maximal width \mathcal{K} can be developed with $\mathcal{O}((\mathcal{K} + 1) \cdot n \cdot (\log n + \mathcal{K}))$ transitions.

This coarse bound guarantees that automata with a small maximum width are developed with very few new transitions (Fig. 6). This is sufficient for real-life genomic automata representing biological features. Such automata, hand-crafted or computed by state-merging techniques [11], are compounds of a few linear-shaped parts (Fig. 2).

For the *lower bound*, the generic argument on linear-shaped WFA can be applied to the longest path in the WFA. If this longest path has a size $\ell \leq n$, we have a bound of $\Omega(\ell \log \ell)$.

3.2 ε -acyclic Automata

To extend the previous bounds for automata *with cycles*, we can consider a slightly modified automaton. An *ε -acyclic automaton* is an automaton without cycles of ε -transitions (Fig. 4). As an ε -acyclic automaton has a numbering with no backward ε -transition, the algorithm 1 can still be used. The same bound of $\mathcal{O}((\mathcal{K} + 1) \cdot n \cdot (\log n + \mathcal{K}))$ is obtained (each width κ_i is now the number of ε -transitions cut by the state q_i).

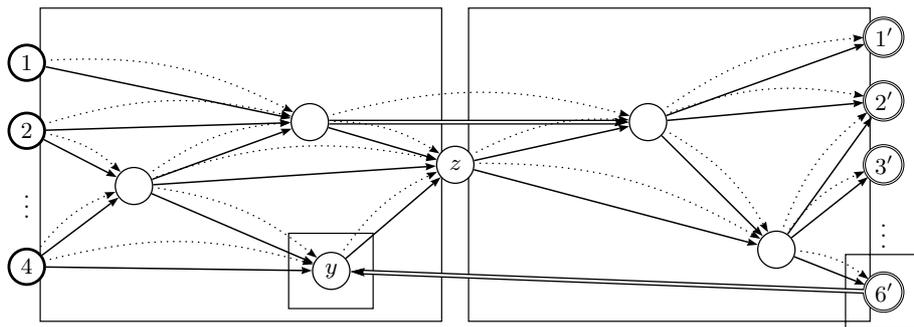


Fig. 4. Unlike the automaton on Fig. 3, this numbered automaton has a backward transition $(6', y, \delta)$. However, that transition is not doubled with an ε -transition.

In real applications, if we have an automaton \mathcal{A} without ε -transitions, we add *some* ε -transitions while keeping the automaton ε -acyclic. This construction is justified when the automaton represents biological structures made of similar units. Those units are separated by sequences that cannot be deleted, as for instance in the case of exon recognition.

4 Removal with Short ε -chains

To further lower the number of new transitions, we can remark that short ε -chains (that is chains of successive ε -transitions) can be actually hardwired with a reasonable critical path (Fig. 5).

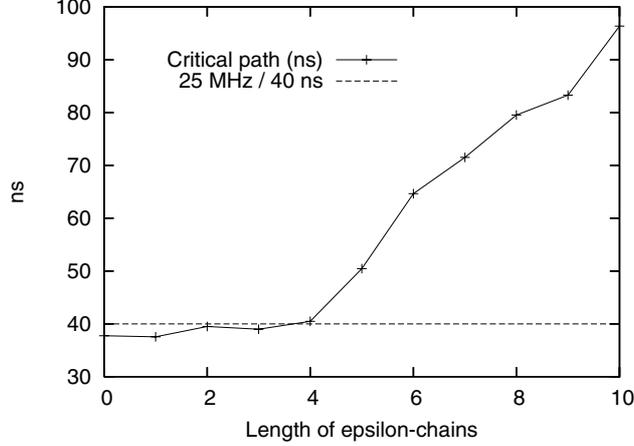


Fig. 5. Critical path for ε -chains with 8-bit weights and a 40 ns (25 MHz) constraint on a WFA with 20 regular transitions and different lengths of ε -chains. Chains of 3 ε -transitions can be hardwired. The FPGA being half-filled (between 48% and 51%), the hardware compiler has a moderate pressure on the different optimisation phases. The critical path should be linear to the length of ε -chains, but the hardware compiler does not further minimize it as soon as it meets the constraint.

Definition 9. Given a WFA \mathcal{A} and $n_\varepsilon \in \mathbb{N}$, any WFA \mathcal{A}'' is a development with short ε -chains of \mathcal{A} if \mathcal{A}'' is path-equivalent to \mathcal{A}_ε and if all ε -chains of \mathcal{A}'' have a length $\leq n_\varepsilon$.

Given a linear-shaped WFA with n states, we can split it into n_ε parts of size $\mathcal{O}(n/n_\varepsilon)$, develop each sub-automaton with $\mathcal{O}(n/n_\varepsilon \cdot \log(n/n_\varepsilon))$ transitions, and finally add an ε -transition that covers each sub-automaton.

Thus we have the following property:

Property 3. A linear-shaped WFA with n states can be developed with short ε -chains with $\mathcal{O}(n \log(n/n_\varepsilon))$ new transitions.

Furthermore, if we restrict that all remaining ε -transitions are *original*, that is, they were present before the removal, the same bound is a lower bound:

Property 4. Given n_ε , any development with short original ε -chains of a linear-shaped WFA with n states has $\Omega(n \log(n/n_\varepsilon))$ new transitions.

The proof, which enumerates some sets in which at least one transition must appear in the automaton, is given in appendix. Although accepting short chains of ε -transitions is a local change, this technique lowers the actual number of new transitions (Fig. 6). The ε -chains can be used in (ε)-acyclic WFA to obtain $\mathcal{O}((\mathcal{K} + 1) \cdot n \cdot (\log(n/n_\varepsilon) + \mathcal{K}))$ new transitions.

Number of states of the initial automaton	20	80	200
Quadratical ε -transitions-removal algorithms	190	3160	19900
Linear-shaped WFA [2]	69	433	1345
Linear-shaped WFA, development with short ε -chains (section 4)			
ε -chains of length $\leq n_\varepsilon = 3$	42	310	1022
ε -chains of length $\leq n_\varepsilon = 5$	30	250	890
(ε -)acyclic WFA (section 3)			
$\mathcal{K} = 1$, best-case	96	522	1555
$\mathcal{K} = 1$, worst-case	141	925	2835
$\mathcal{K} = 2$, best-case	124	612	1766
$\mathcal{K} = 2$, worst-case	176	1292	4096

Fig. 6. Number of new transitions produced while removing ε -transitions on various automata. For acyclic WFA, complexity range from best-case (only one additional branch through the whole WFA) to the worst-case (each cut has a maximal width: the automaton is constantly branching). Even in the worst-case situation, genomic WFA with 80 states and no more than $\mathcal{K} + 2 = 3$ branches can be efficiently hardwired with less than 1300 new transitions.

5 Conclusions and Perspectives

The removal techniques presented in sections 3 and 4 allow larger automata to be hardwired on a given FPGA. For acyclic automata, the best results for a strict application of algorithm 1 would require finding the numbering of the states that minimizes the maximal width \mathcal{K} . In fact, for real automata with a small number of branches as the one in Fig. 2, good solutions are found when cutting at the branching states.

Other studies could find more precise bounds. For acyclic automata, the initial number of transitions could be taken into account. Finally, we plan to study *approximated developments* of automata, in which the resulting automaton would not strictly be path-equivalent to the initial one. In real applications, the cost assigned to deletions prevents sequences with too many ε -transitions from being accepted.

References

1. Giraud, M., Lavenier, D.: Linear encoding scheme for weighted finite automata. In: Ninth International Conference on Implementation and Application of Automata (CIAA 2004), LNCS 3317. (2004)
2. Giraud, M., Lavenier, D.: Dealing with hardware space limits when removing epsilon-transitions in a genomic weighted finite automaton. *Journal of Automata, Languages and Combinatorics* **10** (2005)
3. Aho, A.V., Sethi, R., Ullman, J.D.: *Compilers, Principles, Techniques and Tools*. Addison Wesley (1986)
4. Thompson, K.: Regular expression search algorithm. *Communications of the ACM* **11** (1968) 419–422

5. Mohri, M.: Generic epsilon-removal and input epsilon-normalization algorithms for weighted transducers. *International Journal of Foundations of Computer Science* **13** (2002) 129–143
6. Mohri, M.: Finite-State Transducers in Language and Speech Processing. *Computational Linguistics* **23** (1997) 269–311
7. Quignon, P., Giraud, M., Rimbault, M., Lavigne, P., Tacher, S., Morin, E., Retout, E., Valin, A.S., Lindblad-Toh, K., Nicolas, J., Galibert, F.: The dog and rat olfactory receptor repertoires. *Genome Biology* **6** (2005) R83
8. Hromkovic, J., Seibert, S., Wilke, T.: Translating regular expressions into small epsilon-free nondeterministic finite automata. In: 14th Symposium on Theoretical Aspects of Computer Science (STACS 97), LNCS 1200. (1997) 55–66
9. Lifshits, Y.: A lower bound on the size of ε -free NFA corresponding to a regular expression. *Information Processing Letters* **85** (2003) 293–299
10. Hagenah, C., Muscholl, A.: Computing ε -free NFA from regular expressions in $\mathcal{O}(n \log^2(n))$ time. *Mathematical Foundations of Computer Science* (1998) 277–285
11. Kerbellec, G., Coste, F.: A similar fragments merging approach to learn automata on proteins. In: *Machine Learning: ECML 2005*. (2005)

Appendix

Proof of Property 4

This proof uses a similar technique than the proof of the Lemma 6 in [2], but additional work is done to handle the short ε -chains. The *span* of a transition (q_i, q_j, δ) is $|j - i|$.

Proof. Let \mathcal{A} be a linear-shaped WFA with n states, and \mathcal{A}'' a development with short original ε -chains of \mathcal{A} . Let $\pi = (\tau_{a+1}, \alpha_A)(\tau_{a+2}, \varepsilon) \dots (\tau_{b-1}, \varepsilon)(\tau_b, \alpha_B)$ be a closed path in \mathcal{A} , where α_A and α_B are two characters different from ε . This path has in \mathcal{A}'' a superior path π' that can be written as $\pi' = (\pi'_1, \varepsilon)(\tau_A, \alpha_A)(\pi'_2, \varepsilon)(\tau_B, \alpha_B)(\pi'_3, \varepsilon)$.

As the three paths π'_1, π'_2, π'_3 are original ε -chains, any of them has a span not greater than n_ε transitions, that is $3n_\varepsilon$ globally. Therefore, at least one of the two transitions τ_A and τ_B has a span included in $\{\lceil \frac{k-3n_\varepsilon}{2} \rceil, \dots, k-1\}$ with $k = b - a$ (Figure 7).

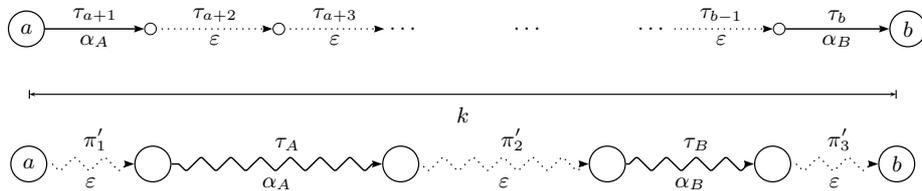


Fig. 7. Proof of the property 4. At least one of the transitions τ_A and τ_B has a span included in $\{\lceil \frac{k-3n_\varepsilon}{2} \rceil, \dots, k-1\}$, where $k = b - a$ is the span of the path $(\tau_A, \alpha_A)(\tau_B, \alpha_B)$.

If we consider all $n-k+1$ pairs (a, b) with the same $k = b-a$, then the WFA \mathcal{A}'' has no less than $(n-k+1)/2$ transitions of span included in $\{\lceil \frac{k-3n_\varepsilon}{2} \rceil, \dots, k-1\}$.

Let (k_i) the sequence defined by $k_{i+1} = 2k_i + 3n_\varepsilon$ and $k_0 = 1$. We have $k_i = 2^i(1 + 3n_\varepsilon) - 3n_\varepsilon$. We consider several k s taking the values of (k_i) from $i = 1$ to the last i such that $k_i \leq n$, that is $i_f = \lfloor \log \frac{n+3n_\varepsilon}{1+3n_\varepsilon} \rfloor = \Theta(\log(n/n_\varepsilon))$.

Then the WFA \mathcal{A}'' has not less than $\sum_{i=1}^{i_f} (n - k_i + 1)/2 = \Theta(n \log(n/n_\varepsilon))$ transitions.