

Weak Factor Automata: Comparing (Failure) Oracles and Storacles

Loek Cleophas, Derrick G. Kourie, and Bruce W. Watson

FASTAR Research Group, University of Pretoria and Stellenbosch University,
South Africa

`{loek,derrick,bruce}@fastar.org`

`http://www.fastar.org`

Prague Stringology Conference, 2–4 September 2013



(Weak) factor automata

- ▶ Factor automaton (DAWG)
 - ▶ Accepts factors of keyword
 - ▶ Used for efficient backward pattern matching
 - ▶ Used as index
- ▶ Weak factor automata
 - ▶ (Small) over-approximation
... to save space
 - ▶ OK for pattern matching
 - ▶ May be OK for indexing

Contributions

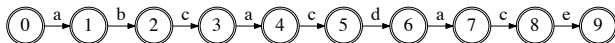
- ▶ New weak factor automata constructions
 - ▶ Based on factor oracle and factor storage
 - ▶ Using failure transitions
 - ▶ Failure factor oracle
 - ▶ Failure factor storage
- ▶ Empirical size comparison
 - ▶ On generated strings of lengths 4 – 9
 - ▶ On English word list (lengths 4 – 28)

Factor oracle construction algorithm

- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}
- 5: **for** i from 2 to m **do**
- 6: Let the longest path from state 0 that spells a prefix of $p_i \dots p_m$ end in state j and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$)
- 7: **if** $k \neq m$ **then**
- 8: Build a new transition $j \xrightarrow{p_{k+1}} k + 1$

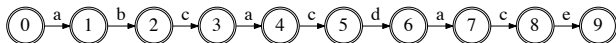
Factor oracle construction algorithm

- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}



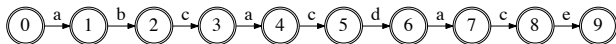
Factor oracle construction algorithm

- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}
- 5: **for** i from 2 to m **do**
- 6: Let the longest path from state 0 that spells a prefix of $p_i \dots p_m$ end in state j and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$)



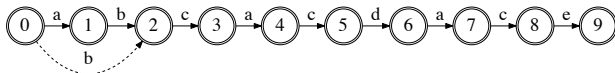
Factor oracle construction algorithm

- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}
- 5: **for** i from 2 to m **do**
- 6: Let the longest path from state 0 that spells a prefix of $p_i \dots p_m$ end in state j and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$)
- 7: **if** $k \neq m$ **then**
- 8: Build a new transition $j \xrightarrow{p_{k+1}} k + 1$



Factor oracle construction algorithm

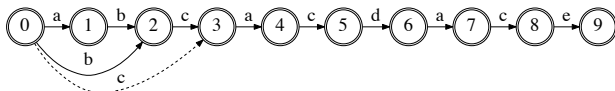
- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}
- 5: **for** i from 2 to m **do**
- 6: Let the longest path from state 0 that spells a prefix of $p_i \dots p_m$ end in state j and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$)
- 7: **if** $k \neq m$ **then**
- 8: Build a new transition $j \xrightarrow{p_{k+1}} k + 1$



Suffix bcadace

Factor oracle construction algorithm

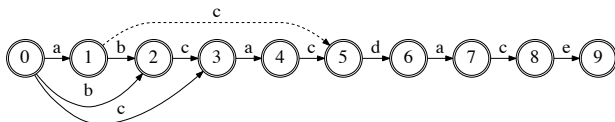
- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}
- 5: **for** i from 2 to m **do**
- 6: Let the longest path from state 0 that spells a prefix of $p_i \dots p_m$ end in state j and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$)
- 7: **if** $k \neq m$ **then**
- 8: Build a new transition $j \xrightarrow{p_{k+1}} k + 1$



Suffix caccace

Factor oracle construction algorithm

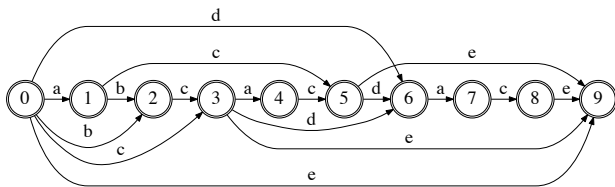
- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}
- 5: **for** i from 2 to m **do**
- 6: Let the longest path from state 0 that spells a prefix of $p_i \dots p_m$ end in state j and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$)
- 7: **if** $k \neq m$ **then**
- 8: Build a new transition $j \xrightarrow{p_{k+1}} k + 1$



Suffix acadce

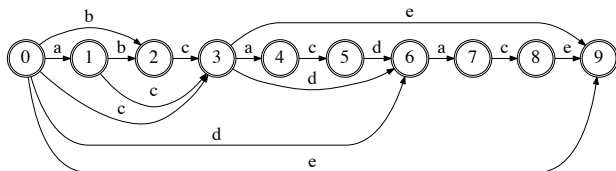
Factor oracle

... etc. leads to



Factor storacle

- ▶ Modification to $O(m^2)$ FO construction
- ▶ **shortest forward transition oracle**
 - ▶ ... keeping it homogeneous
 - ▶ Accidental...
 - ▶ Example smaller than FO



- ▶ Not smaller than FO in general (*Cleophas & Watson 2012*)
... usually *slightly larger*
- ▶ Conjecture m to $3m$ transitions

Factor storage construction algorithm

- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}
- 5: **for** i from 2 to m **do**
- 6: Let the longest path from state 0 that spells a prefix of $p_i \dots p_m$ end in state j and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$)
- 7: **while** $k \neq m$ **do**
- 8: Let the first state from state j onward that has an incoming transition on p_{k+1} be state l ($j < l \leq k + 1$)
- 9: Build a new transition $j \xrightarrow{p_{k+1}} l$
- 10: Let the longest path from state 0 that spells a prefix of $p_i \dots p_m$ end in state j and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$)

Factor storage construction algorithm

- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}

Factor storage construction algorithm

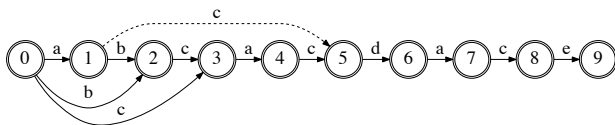
- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}
- 5: **for** i from 2 to m **do**
- 6: Let the longest path from state 0 that spells a prefix of $p_i \dots p_m$ end in state j and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$)
- 8: Let the first state from state j onward that has an incoming transition on p_{k+1} be state l ($j < l \leq k + 1$)
- 9: Build a new transition $j \xrightarrow{p_{k+1}} l$

Factor storage construction algorithm

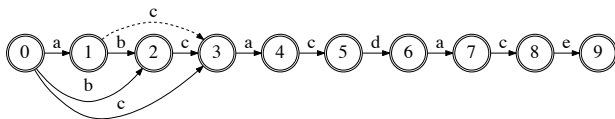
- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}
- 5: **for** i from 2 to m **do**
- 6: Let the longest path from state 0 that spells a prefix of $p_i \dots p_m$ end in state j and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$)
- 7: **while** $k \neq m$ **do**
- 8: Let the first state from state j onward that has an incoming transition on p_{k+1} be state l ($j < l \leq k + 1$)
- 9: Build a new transition $j \xrightarrow{p_{k+1}} l$
- 10: Let the longest path from state 0 that spells a prefix of $p_i \dots p_m$ end in state j and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$)

Factor storage construction example

Recall factor oracle case:

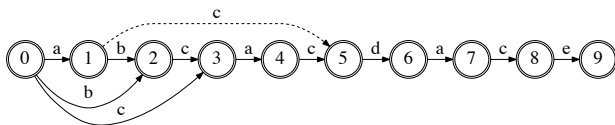


Factor storage case:

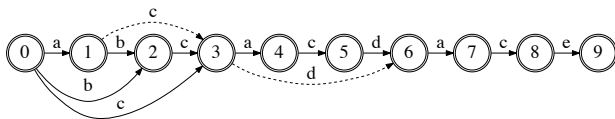


Factor storage construction example

Recall factor oracle case:



Factor storage case:



Failure transitions

- ▶ Allow failure transitions in addition to symbol ones
- ▶ Save space, but more transition use...
- ▶ Not new
 - ▶ Aho-Corasick
 - ▶ Generalized by *Crochemore & Hancart 1997*
- ▶ First general *DFA* \rightarrow *F DFA* algorithm by *Kourie et al. 2012*
 - ▶ Intermediate lattice construction
... keeping state set fixed
- ▶ *Björklund et al. 2013*
 - ▶ Complexity...
... even if keeping state set fixed
 - ▶ Algorithm to reach $\frac{2}{3}$ of optimal savings
- ▶ Both *ex post facto*...

Our idea

- ▶ Introduce failure transitions *during construction*
- ▶ We call the resulting automata *Failure Factor (St)Oracles*
- ▶ Complexity as for non-failure cases
- ▶ Idea: instead of constructing $j \xrightarrow{a} k + 1 \dots$
... construct $j \xrightarrow{fail} k$, from which $k \xrightarrow{a} k + 1$ exists
- ▶ Potential problem...

Our idea

- ▶ Introduce failure transitions *during construction*
- ▶ We call the resulting automata *Failure Factor (St)Oracles*
- ▶ Complexity as for non-failure cases
- ▶ Idea: instead of constructing $j \xrightarrow{a} k + 1 \dots$
... construct $j \xrightarrow{fail} k$, from which $k \xrightarrow{a} k + 1$ exists
- ▶ Potential problem...
- ▶ Using sequence of existing (failure, symbol) transitions may end up in $j > k$

Potential problem...

- ▶ Using sequence of existing (failure, symbol) transitions may end up in $j > k$
 - ... potential for *backward* failure transition
 - ... hence for *cycle*
 - ... hence for *failure cycle*
 - ... hence for *divergent failure cycle*
 - ... leading to live-lock in construction or use of automaton
- ▶ Solution: do not construct backward *failure* transition
 - ... instead create non-forward *symbol* transition
 - ... still potential for cycle (but manageable)

Failure factor oracle construction algorithm

- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}
- 5: **for** i from 2 to m **do**
- 6: Let the longest recognized prefix of $p_i \dots p_m$ be recognized in state j' and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$), and let the longest failure transition path from j' end in state j
- 7: **if** $k \neq m$ **then**
- 8: **if** $k > j$ **then**
- 9: Build a new failure transition $j \xrightarrow{fail} k$
- 10: **else**
- 11: Build a new symbol transition $j \xrightarrow{p_{k+1}} k + 1$

Failure factor oracle construction algorithm

- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}

Failure factor oracle construction algorithm

- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}
- 5: **for** i from 2 to m **do**
- 6: Let the longest recognized prefix of $p_i \dots p_m$ be recognized in state j' and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$), and let the longest failure transition path from j' end in state j

Failure factor oracle construction algorithm

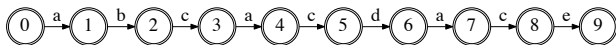
- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}
- 5: **for** i from 2 to m **do**
- 6: Let the longest recognized prefix of $p_i \dots p_m$ be recognized in state j' and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$), and let the longest failure transition path from j' end in state j
- 7: **if** $k \neq m$ **then**
- 8: **if** $k > j$ **then**
- 9: Build a new failure transition $j \xrightarrow{fail} k$
- 10: **else**
- 11: Build a new symbol transition $j \xrightarrow{p_{k+1}} k + 1$

Failure factor oracle construction algorithm

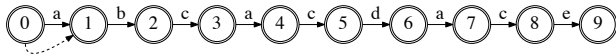
- 1: **for** i from 0 to m **do**
- 2: Create a new final state i
- 3: **for** i from 0 to $m - 1$ **do**
- 4: Create a new transition from i to $i + 1$ on symbol p_{i+1}
- 5: **for** i from 2 to m **do**
- 6: Let the longest recognized prefix of $p_i \dots p_m$ be recognized in state j' and spell out $p_i \dots p_k$ ($i - 1 \leq k \leq m$), and let the longest failure transition path from j' end in state j
- 7: **if** $k \neq m$ **then**
- 8: **if** $k > j$ **then**
- 9: Build a new failure transition $j \xrightarrow{\text{fail}} k$
- 10: **else**
- 11: Build a new symbol transition $j \xrightarrow{p_{k+1}} k + 1$

Failure factor storage construction algorithm similar

Failure factor oracle construction example

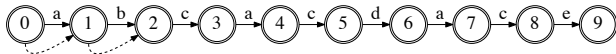


Failure factor oracle construction example



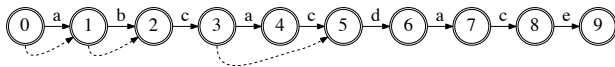
Suffix bcacdace

Failure factor oracle construction example



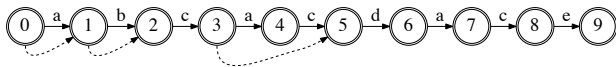
Suffix cascade

Failure factor oracle construction example



Suffix acdace

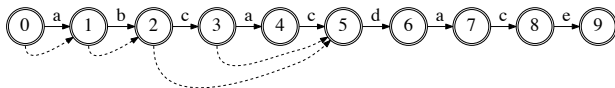
Failure factor oracle construction example



Suffix acdace

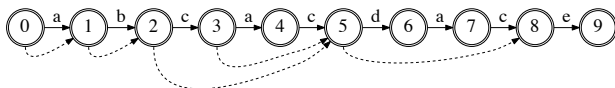
Suffix cdace

Failure factor oracle construction example



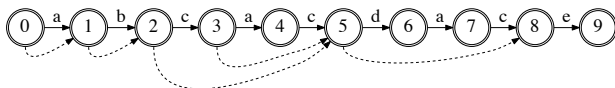
Suffix dace

Failure factor oracle construction example



Suffix ace

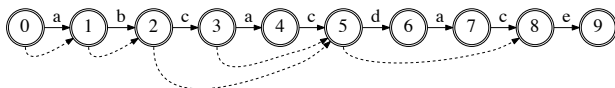
Failure factor oracle construction example



Suffix ace

Suffix ce

Failure factor oracle construction example



Suffix ace

Suffix ce

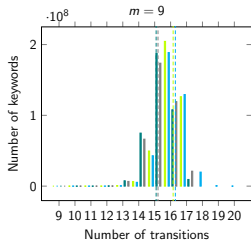
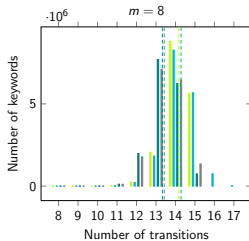
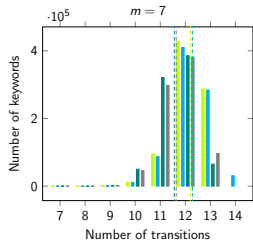
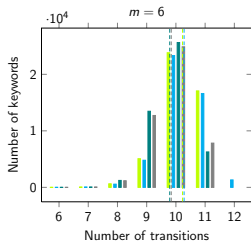
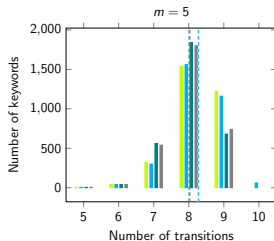
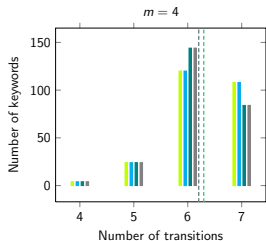
Suffix e

Empirical results

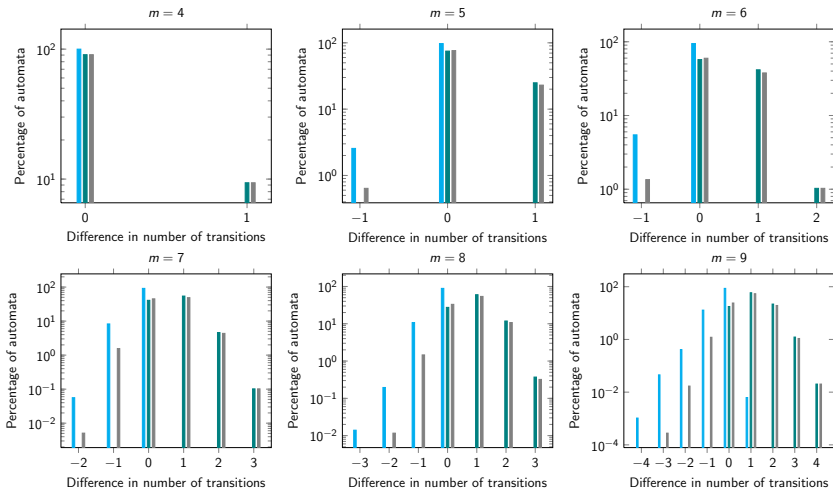
Two data sets

- ▶ Generated strings: all strings of length $m \in [4, 9]$ over alphabet of size m
- ▶ English words

Generated strings—number of transitions



Generated strings—difference in #transitions of FO vs. ...

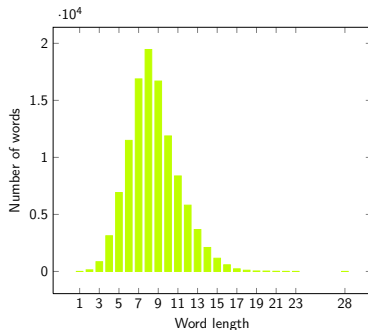


■ FO vs. FSto
 ■ FO vs. FFO
 ■ FO vs. FFSto



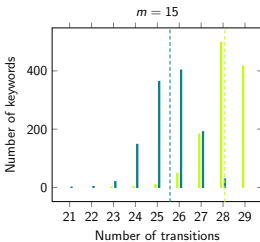
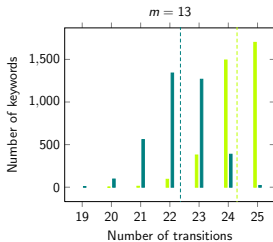
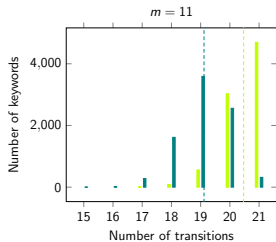
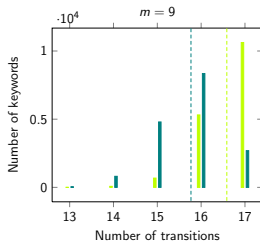
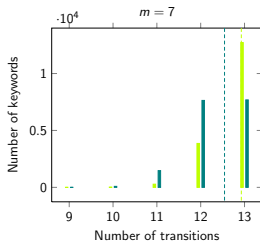
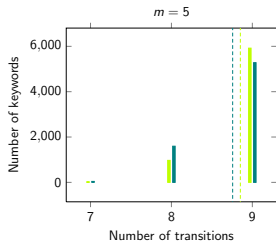
Empirical results on English words

- ▶ English word list from <http://www.sil.org/linguistics/wordlists/english>.

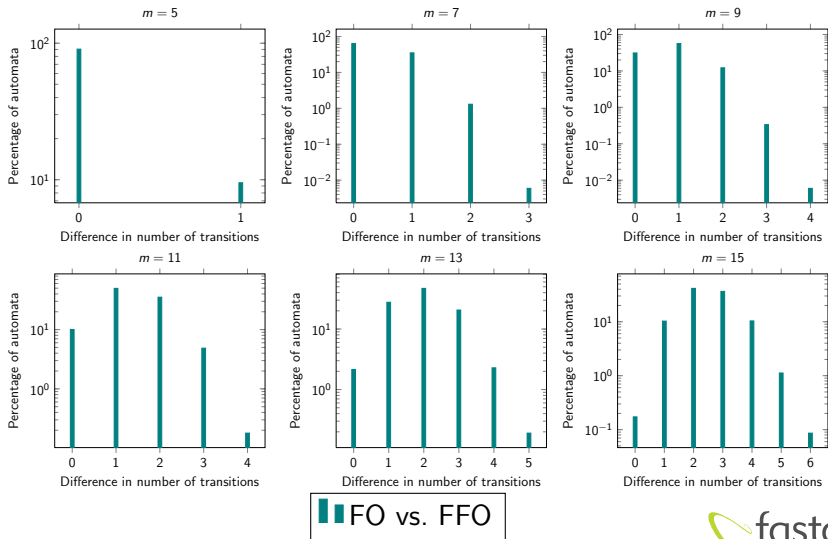


... disregarding words p with $|p| < 4$

English words—number of transitions



English words—difference in #transitions of FO vs. FFO



Concluding remarks

- ▶ Failure versions save ca. 2-10% on #transitions
... possibly more on space
- ▶ Open questions
 - ▶ Upper bounds on number of transitions
 - ▶ Languages
 - ▶ Comparison to general super automata
 - ▶ Comparison to general FDFA construction algorithm
- ▶ Performance when using automata
... recent work: FFO for DNA strings of lengths 4 – 2048
 - ▶ Savings of 8 – 10% for lengths 16 – 2048
 - ▶ Also rudimentary processing; runtime increases 34 – 88%

References

- ▶ Allauzen, Crochemore & Raffinot, *Factor oracle: a new structure for pattern matching*, SOFSEM 1999.
- ▶ Björklund, Björklund & Zechner, *Compact representation of finite automata with failure transitions*, Umea Univ. TR 2013/011.
- ▶ Cleophas & Watson, *On Factor Storacles: an Alternative to Factor Oracles?*, Festschrift for Bořivoj Melichar, 2012.
- ▶ Cleophas, Zwaan & Watson, *Constructing Factor Oracles*, PSC 2003 and JALC special issue 2005.
- ▶ Crochemore & Hancart, *Automata for matching patterns*, pp. 399–462 in Handbook of Formal Languages, 1997.
- ▶ Kourie, Watson, Cleophas & Venter, *Failure Deterministic Finite Automata*, PSC 2012.