

Model Checking for String Problems

Milka Hutagalung Martin Lange

School of Electr. Eng. and Comp. Sc., University of Kassel, Germany

9th Int. Computer Science Symposium in Russia, CSR'14
09/06/2014

- 1 Model Checking
- 2 String Problems
- 3 The Polyadic μ -Calculus
- 4 An Algorithm for LCS

Model Checking – The Big Business

The graphic is a yellow cloud-like shape with a black outline. Inside, several topics are listed in blue text, and two book covers are shown. The topics are arranged in a roughly circular pattern around a central diagram. The topics include: temporal logic, automata, probabilistic systems, SAT/SMT, NUSMV, real-time systems, SPIN, quantitative m.c., hybrid systems, UPPAAL, BDDs, on-the-fly m.c., fixpoints, and bounded model checking. The central diagram is a red square with a white network of nodes and edges. Below the diagram is the book cover for 'Principles of Model Checking' by Charles Burch and David Long. To the left of the diagram is the book cover for 'Model Checking' by Edward M. Clarke, Jr., Dana D. Long, and Steven A. Jha. To the right of the diagram is the book cover for '25 Years of Model Checking' edited by Michael Huth and Thomas Wehr.

temporal logic

automata

probabilistic systems

SAT/SMT

NUSMV

real-time systems

SPIN

quantitative m.c.

hybrid systems

UPPAAL

BDDs

on-the-fly m.c.

fixpoints

bounded model checking

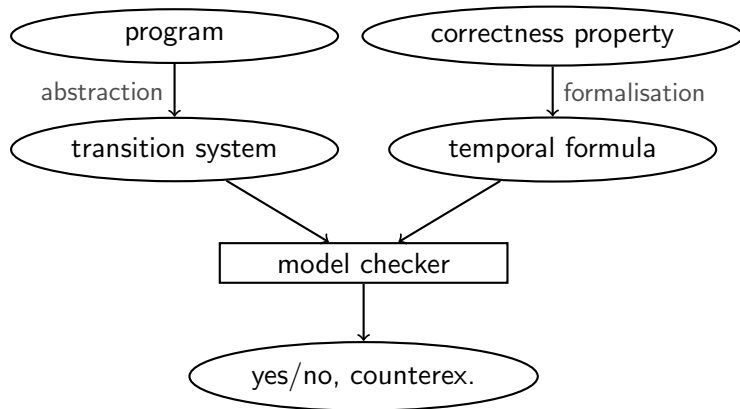
Model Checking
Edward M. Clarke, Jr., Dana D. Long, and Steven A. Jha

Principles of Model Checking
Charles Burch and David Long

25 Years of Model Checking
History, Achievements, Perspectives
Michael Huth and Thomas Wehr

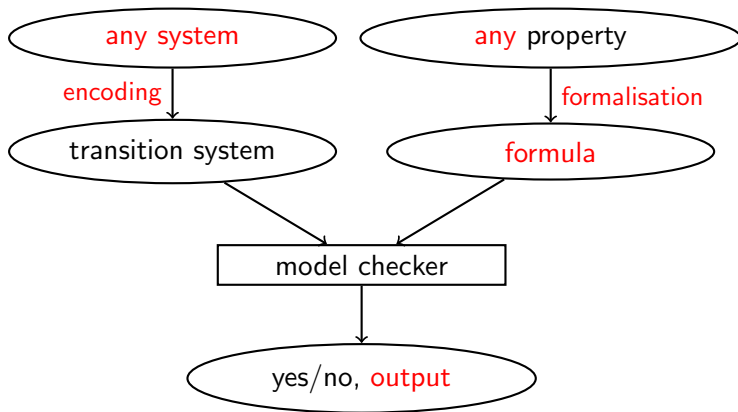
Model Checking – Be Generous

model checking as a method in program verification



Model Checking – Be Generous

model checking as a **generic** method for **problem solving**



Model Checking for Decision Problems

how to use model checking technology to solve **decision** problem P

- 1 take suitable **logic** \mathcal{L} with decidable model checking problem
- 2 find suitable **encoding** enc for instances of P
- 3 find suitable **defining formula** φ_P

$$x \in P \quad \text{iff} \quad enc(x) \models \varphi_P$$

Model Checking for Decision Problems

how to use model checking technology to solve **decision** problem P

- 1 take suitable **logic** \mathcal{L} with decidable model checking problem
- 2 find suitable **encoding** enc for instances of P
- 3 find suitable **defining formula** φ_P

$$x \in P \quad \text{iff} \quad enc(x) \models \varphi_P$$

- 4 **optimise** model checking algorithm for \mathcal{L} into algorithm for P using **partial evaluation**

String Problems

input: $w_1, \dots, w_n \in \Sigma^*$

compute

- ▶ longest common substring (LCST):
 v such that $w_i = \dots v \dots$ for all $i = 1, \dots, n$

String Problems

input: $w_1, \dots, w_n \in \Sigma^*$

compute

- ▶ longest common substring (LCST):
 v such that $w_i = \dots v \dots$ for all $i = 1, \dots, n$
- ▶ shortest common superstring (SCST):
 v such that $v = \dots w_i \dots$ for all $i = 1, \dots, n$

String Problems

input: $w_1, \dots, w_n \in \Sigma^*$

compute

- ▶ longest common substring (LCST):
 v such that $w_i = \dots v \dots$ for all $i = 1, \dots, n$
- ▶ shortest common superstring (SCST):
 v such that $v = \dots w_i \dots$ for all $i = 1, \dots, n$
- ▶ longest common subsequence (LCSS):
 $a_1 \dots a_k$ such that $w_i = \dots a_1 \dots a_2 \dots a_k \dots$ for all $i = 1, \dots, n$

String Problems

input: $w_1, \dots, w_n \in \Sigma^*$

compute

- ▶ **longest common substring (LCST):**
 v such that $w_i = \dots v \dots$ for all $i = 1, \dots, n$
- ▶ **shortest common superstring (SCST):**
 v such that $v = \dots w_i \dots$ for all $i = 1, \dots, n$
- ▶ **longest common subsequence (LCSS):**
 $a_1 \dots a_k$ such that $w_i = \dots a_1 \dots a_2 \dots a_k \dots$ for all $i = 1, \dots, n$

Algorithmic Solutions to String Problems

input: $w_1, \dots, w_n \in \Sigma^*$, assuming $|w_i| \leq m$

- ▶ dynamic programming in $\mathcal{O}(m^n)$
- ▶ suffix tree algorithm in $\mathcal{O}(mn)$

String Problems via Computational Logic

goal: take suitable logic \mathcal{L}

define these string problems in \mathcal{L} on suitable input representation

use **model checking** and **partial evaluation** to design algorithms for LCST, SCST, LCSS

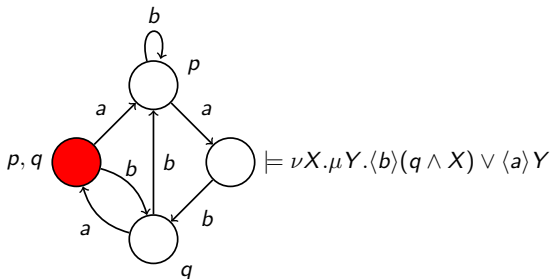
slight mismatch: LCST, SCST, LCSS are **computation** problems

The Modal μ -Calculus

syntax

$$\varphi ::= q \mid \bar{q} \mid \varphi \wedge \varphi \mid \varphi \vee \varphi \mid \langle a \rangle \varphi \mid [a] \varphi \mid X \mid \mu X. \varphi \mid \nu X. \varphi$$

interpreted in **states** of a transition system



The Polyadic μ -Calculus

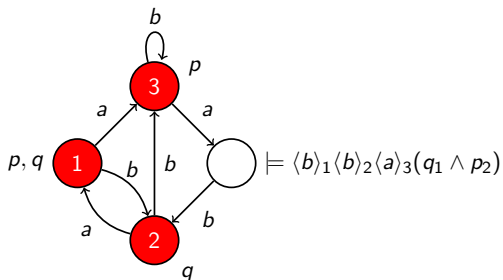
syntax

[Andersen'94, Otto'99]

$$\varphi ::= q_i \mid \bar{q}_i \mid \varphi \wedge \varphi \mid \varphi \vee \varphi \mid \langle a \rangle_i \varphi \mid [a]_i \varphi \mid X \mid \mu X. \varphi \mid \nu X. \varphi \mid \{\bar{i} \leftarrow \bar{j}\} \varphi$$

with $1 \leq i \leq d$

interpreted in d -tuples of states of a transition system



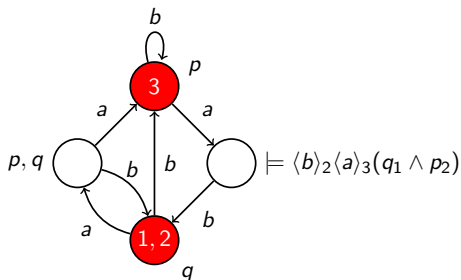
The μ -Calculus

syntax

$$\varphi ::= q_i \mid \bar{q}_i \mid \varphi \wedge \varphi \mid \varphi \vee \varphi \mid \langle a \rangle_i \varphi \mid [a]_i \varphi \mid X \mid \mu X. \varphi \mid \nu X. \varphi \mid \{\bar{i} \leftarrow \bar{j}\} \varphi$$

with $1 \leq i \leq d$

interpreted in d -tuples of states of a transition system



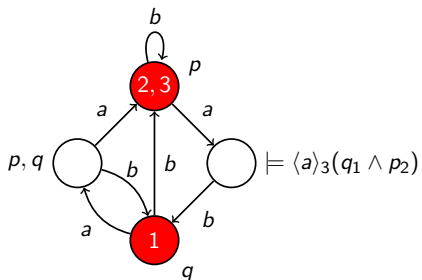
The μ -Calculus

syntax

$$\varphi ::= q_i \mid \bar{q}_i \mid \varphi \wedge \varphi \mid \varphi \vee \varphi \mid \langle a \rangle_i \varphi \mid [a]_i \varphi \mid X \mid \mu X. \varphi \mid \nu X. \varphi \mid \{\bar{i} \leftarrow \bar{j}\} \varphi$$

with $1 \leq i \leq d$

interpreted in d -tuples of states of a transition system



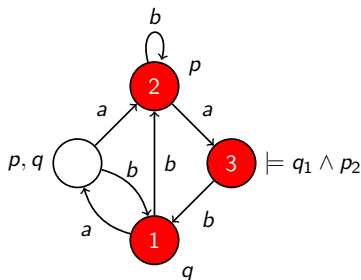
The μ -Calculus

syntax

$$\varphi ::= q_i \mid \bar{q}_i \mid \varphi \wedge \varphi \mid \varphi \vee \varphi \mid \langle a \rangle_i \varphi \mid [a]_i \varphi \mid X \mid \mu X. \varphi \mid \nu X. \varphi \mid \{\bar{i} \leftarrow \bar{j}\} \varphi$$

with $1 \leq i \leq d$

interpreted in d -tuples of states of a transition system



The Polyadic μ -Calculus

Fact: $\mathcal{L}_\mu = \mathcal{L}_\mu^1 \leq \mathcal{L}_\mu^2 \leq \dots \leq \mathcal{L}_\mu^\omega$

standard example: \mathcal{L}_μ^2 can define **bisimilarity**

$$\nu X. \left(\bigwedge_q \overline{q_1} \vee q_2 \right) \wedge \left(\bigwedge_a [a]_1 \langle a \rangle_2 X \right) \wedge \{(1, 2) \leftarrow (2, 1)\} X$$

other simulations definable

[Andersen'94, L./Lozes/Vargas'12]

The Polyadic μ -Calculus

Fact: $\mathcal{L}_\mu = \mathcal{L}_\mu^1 \leq \mathcal{L}_\mu^2 \leq \dots \leq \mathcal{L}_\mu^\omega$

standard example: \mathcal{L}_μ^2 can define **bisimilarity**

$$\nu X. \left(\bigwedge_q \overline{q_1} \vee q_2 \right) \wedge \left(\bigwedge_a [a]_1 \langle a \rangle_2 X \right) \wedge \{(1, 2) \leftarrow (2, 1)\} X$$

other simulations definable

[Andersen'94, L./Lozes/Vargas'12]

Thm.: $\mathcal{L}_\mu^\omega = \text{PTIME}/\sim$

[Otto'99]

Thm.: \mathcal{L}_μ^2 satisfiability is undecidable

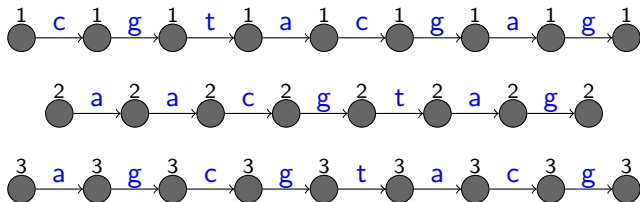
[Otto'99]

Thm.: \mathcal{L}_μ^d model checking in $\mathcal{O}((en^k)^d)$

[L./Lozes'12]

Defining String Problems in $\mathcal{L}_{\mu}^{\omega}$

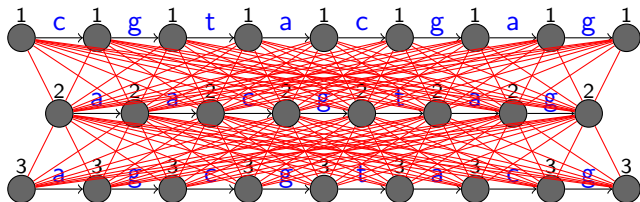
consider input w_1, \dots, w_n as transition system, e.g.



$$\Phi_{\text{LCST}} = \nu X. \left(\bigwedge_{i=1}^3 i_i \right) \wedge \bigvee_{b \in \{a, c, g, t\}} \langle b \rangle_1 \langle b \rangle_2 \langle b \rangle_3 X$$

Defining String Problems in $\mathcal{L}_{\mu}^{\omega}$

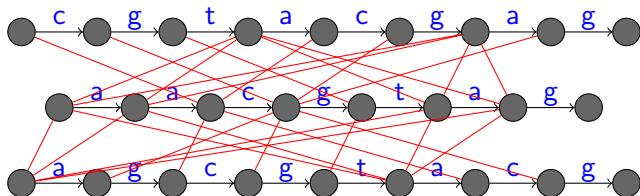
consider input w_1, \dots, w_n as transition system, e.g.



$$\Phi_{\text{LCST}} = \nu X. \left(\bigwedge_{i=1}^3 i_i \right) \wedge \bigvee_{b \in \{a, c, g, t\}} \langle b \rangle_1 \langle b \rangle_2 \langle b \rangle_3 X$$

Defining String Problems in $\mathcal{L}_{\mu}^{\omega}$

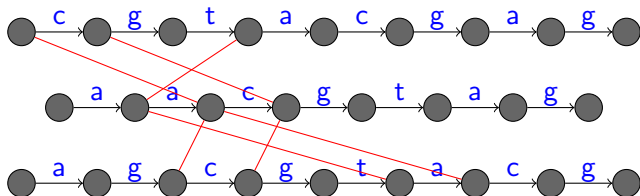
consider input w_1, \dots, w_n as transition system, e.g.



$$\Phi_{\text{LCST}} = \nu^1 X. \left(\bigwedge_{i=1}^3 i_i \right) \wedge \bigvee_{b \in \{a, c, g, t\}} \langle b \rangle_1 \langle b \rangle_2 \langle b \rangle_3 X$$

Defining String Problems in $\mathcal{L}_{\mu}^{\omega}$

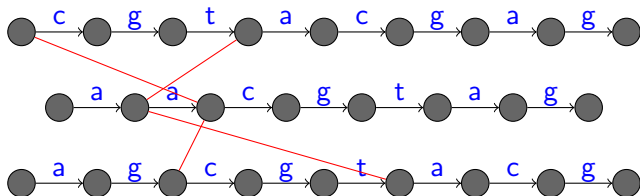
consider input w_1, \dots, w_n as transition system, e.g.



$$\Phi_{\text{LCST}} = \nu^2 X. \left(\bigwedge_{i=1} i_j \right) \wedge \bigvee_{b \in \{a, c, g, t\}} \langle b \rangle_1 \langle b \rangle_2 \langle b \rangle_3 X$$

Defining String Problems in $\mathcal{L}_{\mu}^{\omega}$

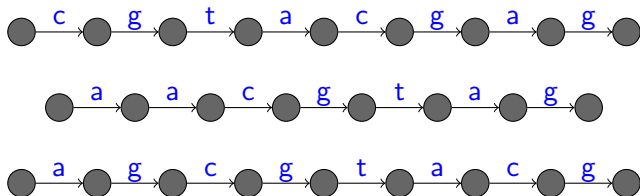
consider input w_1, \dots, w_n as transition system, e.g.



$$\Phi_{\text{LCST}} = \nu^3 X. \left(\bigwedge_{i=1}^3 i_j \right) \wedge \bigvee_{b \in \{a, c, g, t\}} \langle b \rangle_1 \langle b \rangle_2 \langle b \rangle_3 X$$

Defining String Problems in $\mathcal{L}_{\mu}^{\omega}$

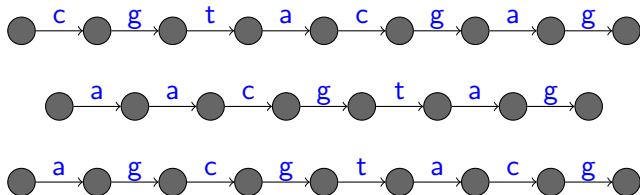
consider input w_1, \dots, w_n as transition system, e.g.



$$\Phi_{\text{LCST}} = \nu^4 X. \left(\bigwedge_{i=1}^4 i_j \right) \wedge \bigvee_{b \in \{a, c, g, t\}} \langle b \rangle_1 \langle b \rangle_2 \langle b \rangle_3 X$$

Defining String Problems in $\mathcal{L}_{\mu}^{\omega}$

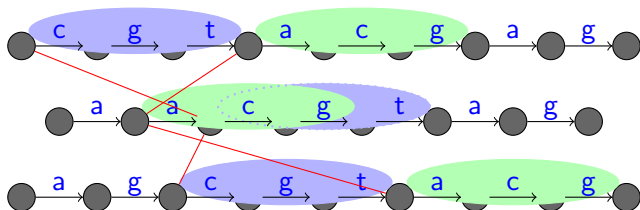
consider input w_1, \dots, w_n as transition system, e.g.



$$\Phi_{\text{LCST}} = \nu X. \left(\bigwedge_{i=1} i_i \right) \wedge \bigvee_{b \in \{a, c, g, t\}} \langle b \rangle_1 \langle b \rangle_2 \langle b \rangle_3 X \equiv \perp$$

Defining String Problems in $\mathcal{L}_{\mu}^{\omega}$

consider input w_1, \dots, w_n as transition system, e.g.



$$\Phi_{\text{LCST}} = \nu^3 X. \left(\bigwedge_{i=1}^3 i_j \right) \wedge \bigvee_{b \in \{a, c, g, t\}} \langle b \rangle_1 \langle b \rangle_2 \langle b \rangle_3 X$$

longest common substrings acg and cgt found in penultimate iteration

- 1 Model Checking
- 2 String Problems
- 3 The Polyadic μ -Calculus
- 4 An Algorithm for LCS

Using Partial Evaluation

observations on model checking Φ_{LCST}

- ▶ deterministic transition relation
- ▶ sets of n -tuples with **exactly one state per input word**
- ▶ operations: **union**, **left-shift-by-one**, **intersection**
- ▶ needed in each step: occurrences of **a, c, g, t**

maintain **occurrence sets** for $b \in \{a, c, g, t\}$

$$\text{Occ}(b) := \{j_i \mid w_i(j) = b\}$$

Using Partial Evaluation

observations on model checking Φ_{LCST}

- ▶ deterministic transition relation
- ▶ sets of n -tuples with **exactly one state per input word**
- ▶ operations: **union**, **left-shift-by-one**, **intersection**
- ▶ needed in each step: occurrences of **a, c, g, t**

maintain **occurrence sets** for $b \in \{a, c, g, t\}$

$$\text{Occ}(b) := \{j_i \mid w_i(j) = b\}$$

extend to occurrences of subwords $w \in \{a, c, g, t\}^+$ by

$$\text{Occ}(bw) := \text{Occ}(b) \cap \{j_i \mid (j+1)_i \in \text{Occ}(w)\}$$

Solutions

Def.: w is **solution** if $\forall i \exists j. j_i \in \text{Occ}(w)$

Lemma: bw is solution $\Rightarrow w$ is solution

Solutions

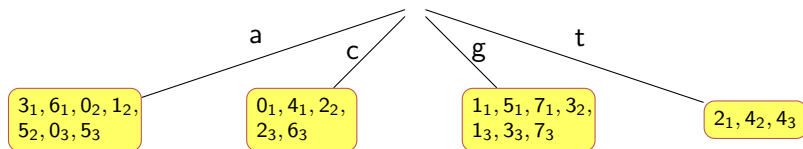
Def.: w is **solution** if $\forall i \exists j. j_i \in Occ(w)$

Lemma: bw is solution $\Rightarrow w$ is solution

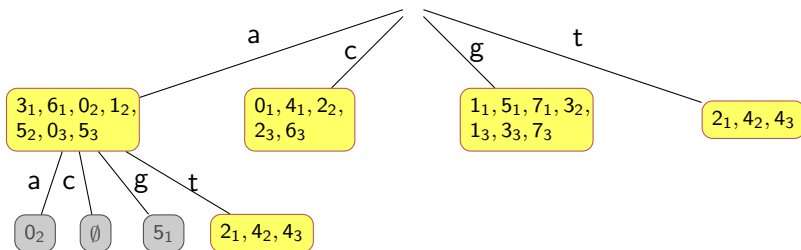
Algorithm LCST

```
procedure LCST( $w_1, \dots, w_n$ )  
   $S \leftarrow \Sigma$   
  compute  $Occ(b)$  for all  $b \in \Sigma$   
  repeat  
     $S' \leftarrow S$   
     $S \leftarrow \{bw \mid w \in S', b \in \Sigma\}$   
    compute  $Occ(bw)$  for all  $bw \in S$   
    eliminate non-solutions from  $S$   
  until  $S = \emptyset$   
  return  $\{(w, Occ(w)) \mid w \in S'\}$   
end procedure
```

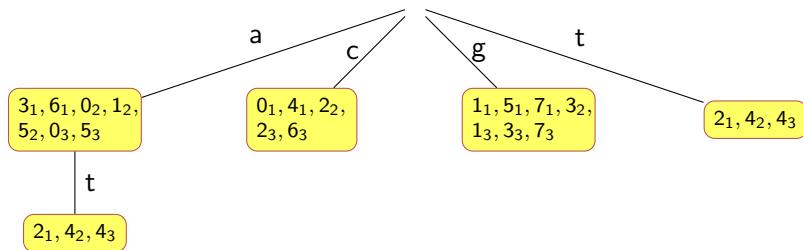
Example (Continued)



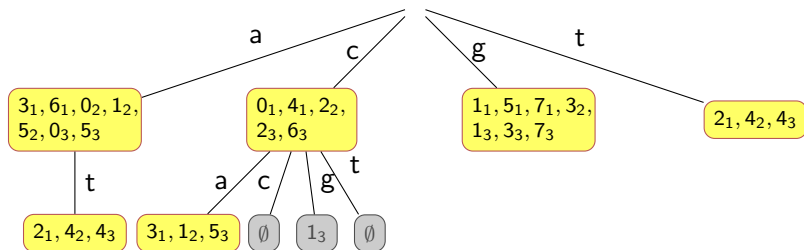
Example (Continued)



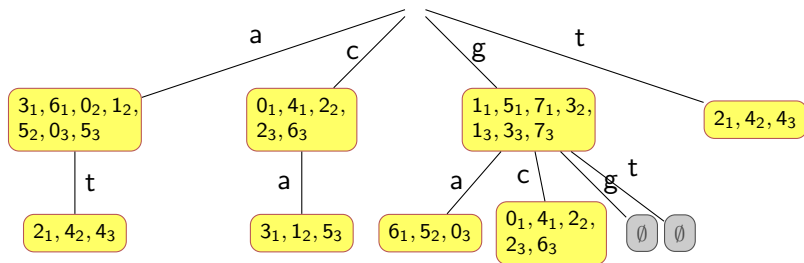
Example (Continued)



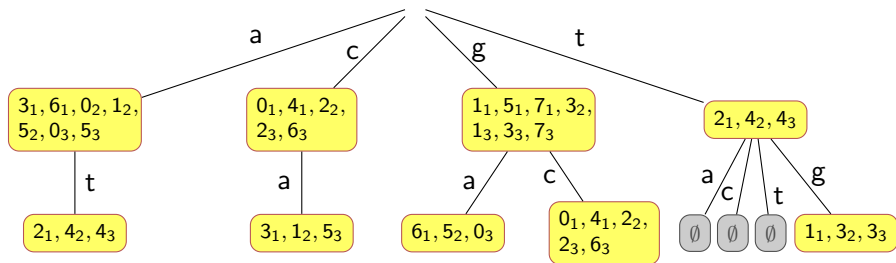
Example (Continued)



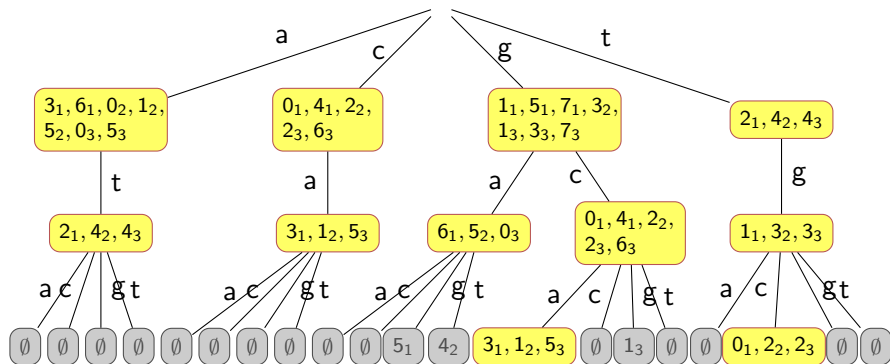
Example (Continued)



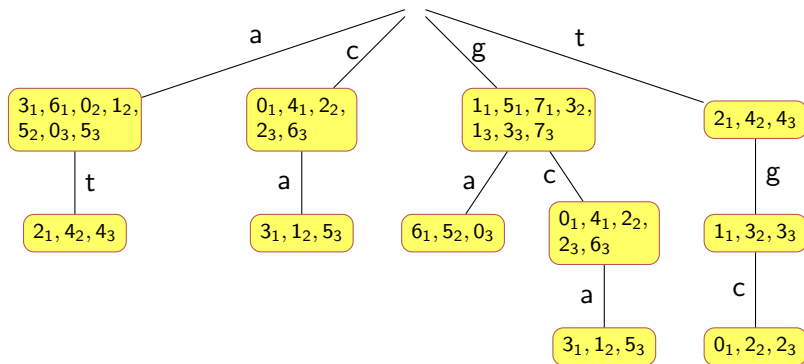
Example (Continued)



Example (Continued)



Example (Continued)



Further Optimisations

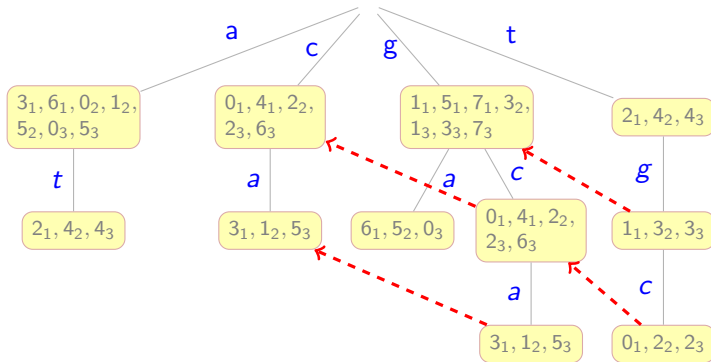
for instance: bw not a solution \Rightarrow bwb' not a solution

\rightsquigarrow point to longest prefix

Further Optimisations

for instance: *bw* not a solution \Rightarrow *bwb'* not a solution

\rightsquigarrow point to longest prefix



Complexity Analysis

input w_1, \dots, w_n with $|w_i| \leq m$

Theorem 1

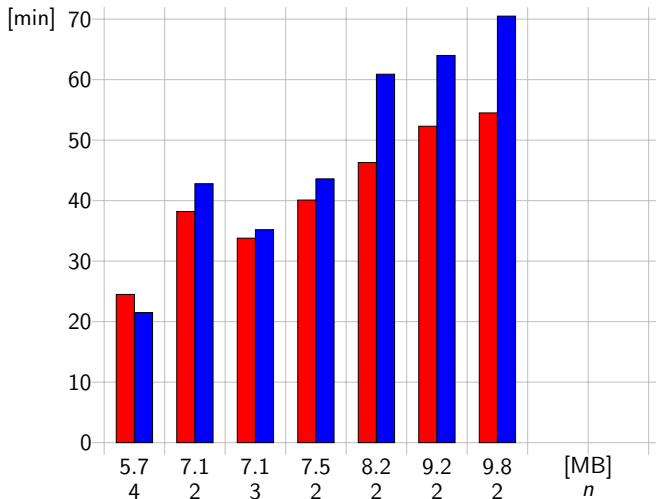
Algorithm LCST can be made to explore at most $\mathcal{O}(nm)$ many occurrence sets.

PROOF IDEA: partition nodes in this graph into three kinds

- ▶ non-extendable to the left $\leq m$
- ▶ extendable to the left by multiple letters $\leq nm$
- ▶ others $\leq m$

Benchmark

LCST against suffix tree algorithm on genome analysis between bacteria / viri



Online Potential

conceptual comparison between LCST and suffix tree algorithm

- ▶ **similar** tree-like data structure \rightsquigarrow similar optimisations possible
- ▶ built **differently**
 - ▶ suffix tree alg. processes **input words** one-by-one
 - ▶ LCST processes **common subwords** one-by-one

rule of thumb: after half the running time

- ▶ suffix tree alg. finds **longest** common substring of **half the inputs**
- ▶ LCST finds **half-longest** common substring of **all inputs**

The End