

Learning Automata with Name Allocation

1st July 2025

Florian Frank, Stefan Milius, Jurriaan Rot
and Henning Urbat

Research Seminar

Chair for Computer Science 8 (Theoretical Computer Science)
Friedrich-Alexander-Universität Erlangen-Nürnberg




T.CS



Friedrich-Alexander-Universität
Faculty of Engineering

» Data languages are formal languages over an infinite alphabet.



A : admissible user IDs for
a server (\rightsquigarrow *infinite set*)

» Data languages are formal languages over an infinite alphabet.

\mathbb{A} : admissible user IDs for
a server (\rightsquigarrow infinite set)

$$\mathcal{L} = \left\{ a_1 \cdots a_n \in \mathbb{A}^* : \left(a_1 = a_n \wedge \forall 1 < i < n. a_1 \neq a_i \right) \right\}$$

'first and last user coincide and differ from any other user'

» Data languages are formal languages over an infinite alphabet.

Standard model: Register Automata

\mathbb{A} : admissible user IDs for
a server (\rightsquigarrow infinite set)

$$\mathcal{L} = \left\{ a_1 \cdots a_n \in \mathbb{A}^* : \left(a_1 = a_n \wedge \forall 1 < i < n. a_1 \neq a_i \right) \right\}$$

'first and last user coincide and differ from any other user'

» Data languages are formal languages over an infinite alphabet.

Standard model: Register Automata \rightsquigarrow **unfeasible for model checking (undecidable inclusion)**

\mathbb{A} : admissible user IDs for
a server (\rightsquigarrow *infinite set*)

$$\mathcal{L} = \left\{ a_1 \cdots a_n \in \mathbb{A}^* : \left(a_1 = a_n \wedge \forall 1 < i < n. a_1 \neq a_i \right) \right\}$$

'first and last user coincide and differ from any other user'

» Data languages are formal languages over an infinite alphabet.

Standard model: Register Automata \rightsquigarrow **unfeasible for model checking (undecidable inclusion)**

» To gain decidability, we must accept restrictions in their expressivity.

\mathbb{A} : admissible user IDs for
a server (\rightsquigarrow *infinite set*)

$$\mathcal{L} = \left\{ a_1 \cdots a_n \in \mathbb{A}^* : \left(a_1 = a_n \wedge \forall 1 < i < n. a_1 \neq a_i \right) \right\}$$

'first and last user coincide and differ from any other user'

» Data languages are formal languages over an infinite alphabet.

Standard model: Register Automata \rightsquigarrow **unfeasible for model checking (undecidable inclusion)**

» To gain decidability, we must accept restrictions in their expressivity.

\mathbb{A} : admissible user IDs for
a server (\rightsquigarrow *infinite set*)

$$\mathcal{L} = \left\{ a_1 \cdots a_n \in \mathbb{A}^* : \left(a_1 = a_n \wedge \forall 1 < i < n. a_1 \neq a_i \right) \right\}$$

'first and last user coincide and differ from any other user'

Result

Schröder, Kozen, Milius, Wißmann '17

(Specific) **languages expressible by binding signatures** and their automata have decidable inclusion problems. \longleftarrow What are 'words with binders'?

» Data languages are formal languages over an infinite alphabet.

Standard model: Register Automata \rightsquigarrow **unfeasible for model checking (undecidable inclusion)**

» To gain decidability, we must accept restrictions in their expressivity.

\mathbb{A} : admissible user IDs for a server (\rightsquigarrow *infinite set*)

$$\mathcal{L} = \left\{ a_1 \cdots a_n \in \mathbb{A}^* : \left(\begin{array}{l} a_1 = a_n \wedge \\ \forall 1 < i < n. a_1 \neq a_i \end{array} \right) \right\}$$

'first and last user coincide and differ from any other user'

$\lambda a. (\lambda b.)^* a$
(using shadowing)



Result

Schröder, Kozen, Milius, Wißmann '17

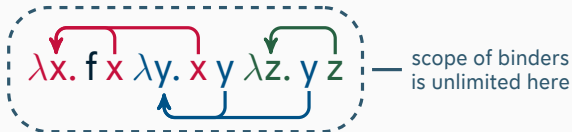
(Specific) **languages expressible by binding signatures** and their automata have decidable inclusion problems.

What are 'words with binders'?

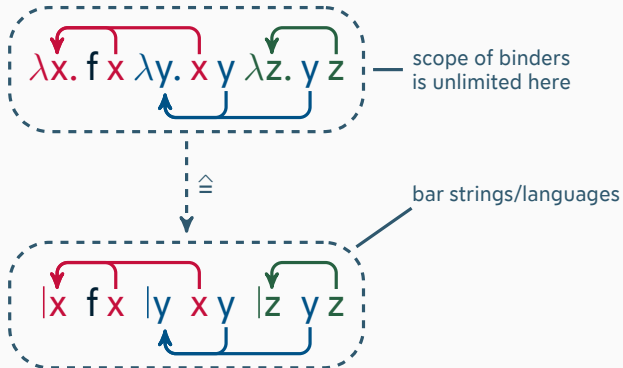
- » We consider data languages with explicit binders, which we see with λ -terms without parenthesis:

$\lambda x. f x \lambda y. x y \lambda z. y z$ — scope of binders is unlimited here

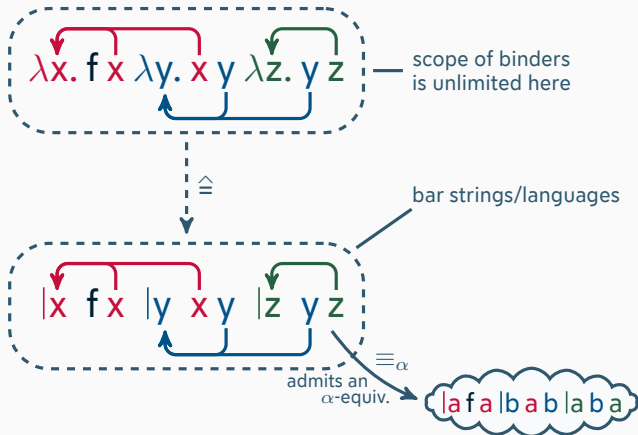
- » We consider data languages with explicit binders, which we see with λ -terms without parenthesis:



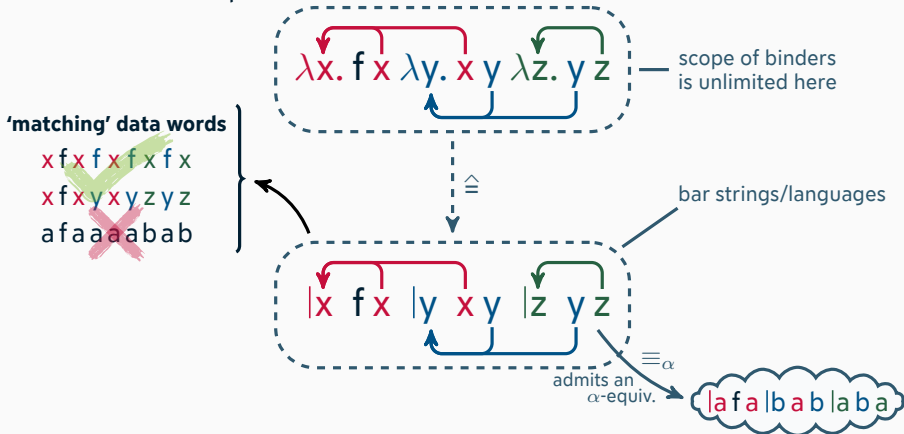
- » We consider data languages with explicit binders, which we see with λ -terms without parenthesis:



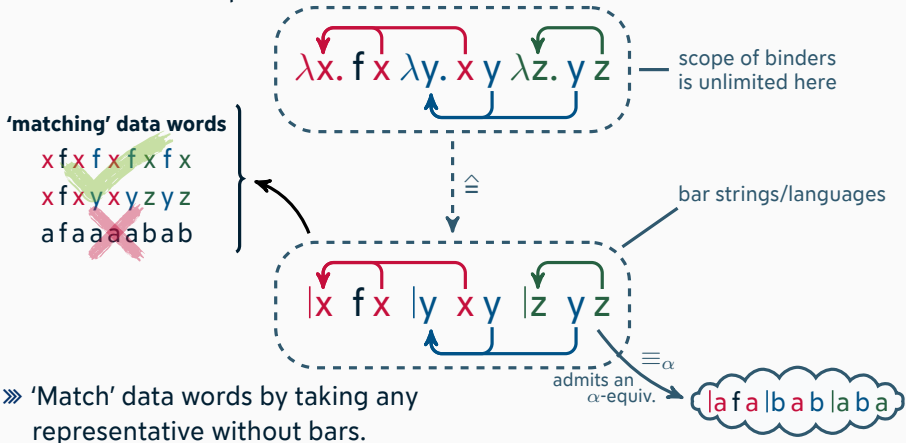
- » We consider data languages with explicit binders, which we see with λ -terms without parenthesis:



- » We consider data languages with explicit binders, which we see with λ -terms without parenthesis:

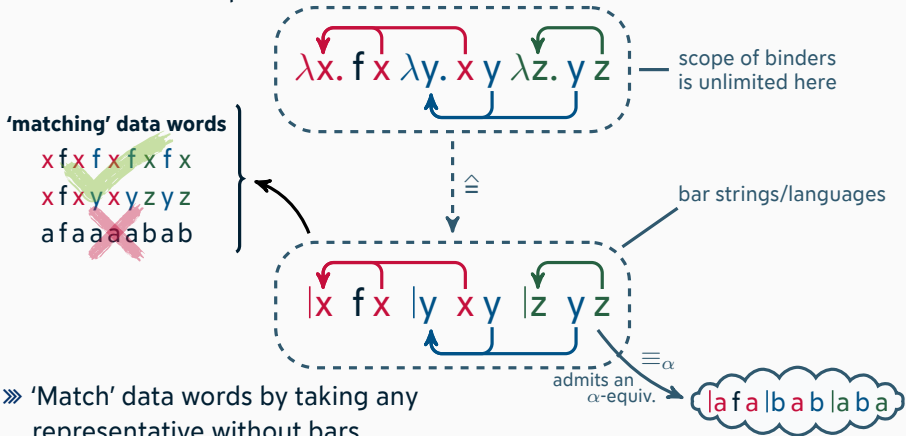


- » We consider data languages with explicit binders, which we see with λ -terms without parenthesis:



- » 'Match' data words by taking any representative without bars.

- » We consider data languages with explicit binders, which we see with λ -terms without parenthesis:



- » 'Match' data words by taking any representative without bars.
- » **Our Problem:** How do we get the necessary models for model checking from a black-box system?

» We consider classical automata over *finite* subalphabets $\overline{A}_0 \subseteq_f \overline{A}$:

Definition (Bar DFA)

A bar DFA \mathcal{A} is a DFA over a finite alphabet $\overline{A}_0 \subseteq_f \overline{A}$.

Its *bar language* $L_\alpha(\mathcal{A}) = \{w \in \overline{A}^* : w \equiv_\alpha w' \in L(\mathcal{A})\}$ consists of all representatives of its α -equivalence classes.

automaton is *closed* iff $L(\mathcal{A}) = L_\alpha(\mathcal{A})$

- » We consider classical automata over *finite* subalphabets $\overline{A}_0 \subseteq_f \overline{A}$:

Definition (Bar DFA)

A bar DFA \mathcal{A} is a DFA over a finite alphabet $\overline{A}_0 \subseteq_f \overline{A}$.

Its *bar language* $L_\alpha(\mathcal{A}) = \{w \in \overline{A}^* : w \equiv_\alpha w' \in L(\mathcal{A})\}$ consists of all representatives of its α -equivalence classes.

automaton is *closed* iff $L(\mathcal{A}) = L_\alpha(\mathcal{A})$

- » Correspond precisely to Schröder et al.'s nominal automata.
(Schröder, Kozen, Milius, Wißmann '17)

- » We consider classical automata over *finite* subalphabets $\overline{A}_0 \subseteq_f \overline{A}$:

Definition (Bar DFA)

A bar DFA \mathcal{A} is a DFA over a finite alphabet $\overline{A}_0 \subseteq_f \overline{A}$.

Its *bar language* $L_\alpha(\mathcal{A}) = \{w \in \overline{A}^* : w \equiv_\alpha w' \in L(\mathcal{A})\}$ consists of all representatives of its α -equivalence classes.

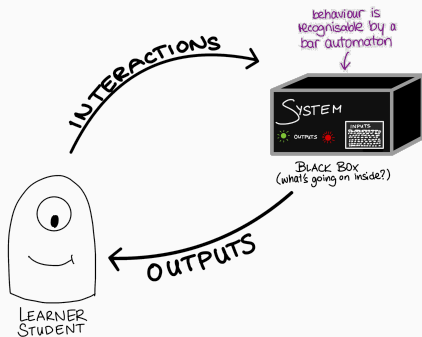
automaton is *closed* iff $L(\mathcal{A}) = L_\alpha(\mathcal{A})$

- » Correspond precisely to Schröder et al.'s nominal automata.
(Schröder, Kozen, Milius, Wißmann '17)
- » Expressivity (data languages):
subclass of register automata



Assumption: Black-box system recognises a bar language L_T

» **Task:** Infer an automaton behaving 'identically' to the black-box system.

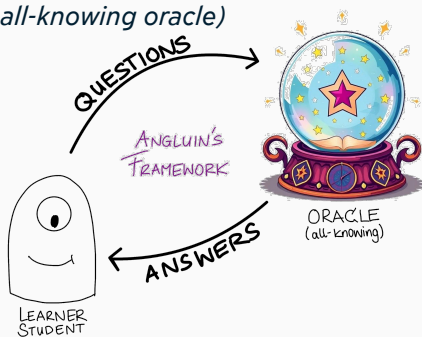




Assumption: Black-box system recognises a bar language L_T

- » **Task:** Infer an automaton behaving 'identically' to the black-box system.
- » Take Angluin's framework:
(*replace the black-box system by an all-knowing oracle*)

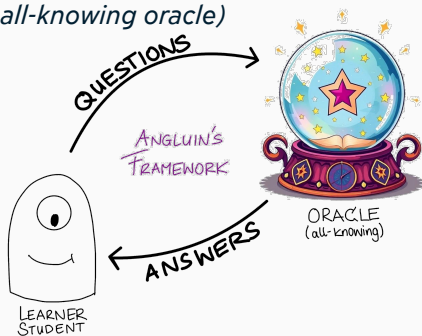
(Angluin '87)





Assumption: Black-box system recognises a bar language L_T

- » **Task:** Infer an automaton behaving 'identically' to the black-box system.
- » Take Angluin's framework: (Angluin '87)
(*replace the black-box system by an all-knowing oracle*)
- » Two kinds of queries:



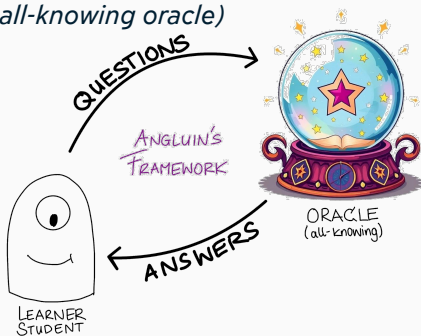


Assumption: Black-box system recognises a bar language L_T

- » **Task:** Infer an automaton behaving 'identically' to the black-box system.
- » Take Angluin's framework: (Angluin '87)
(*replace the black-box system by an all-knowing oracle*)
- » Two kinds of queries:

MQ

Given a bar string $w \in \overline{A}^*$,
is $w \in L_T$?





Assumption: Black-box system recognises a bar language L_T

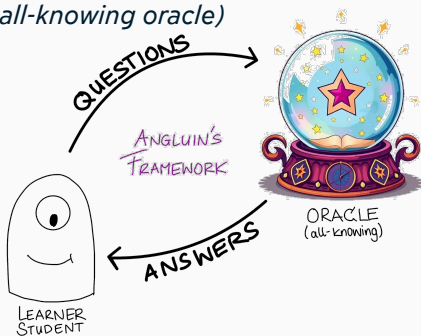
- » **Task:** Infer an automaton behaving 'identically' to the black-box system.
- » Take Angluin's framework: (Angluin '87)
(*replace the black-box system by an all-knowing oracle*)
- » Two kinds of queries:

MQ

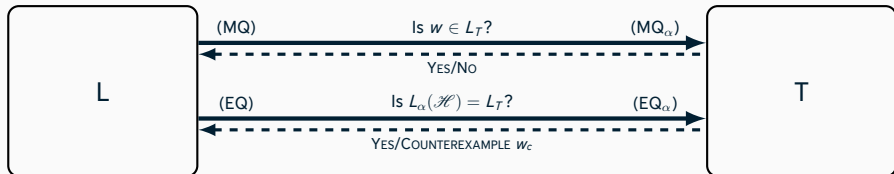
Given a bar string $w \in \bar{A}^*$,
is $w \in L_T$?

EQ

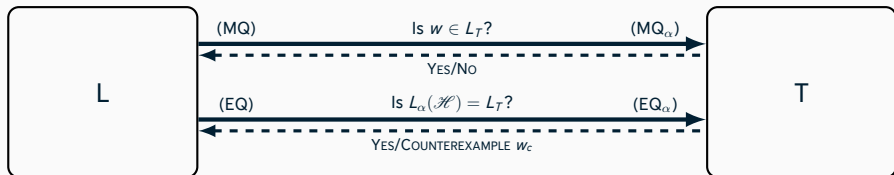
Given a bar automaton \mathcal{H} ,
is $L_\alpha(\mathcal{H}) = L_T$?



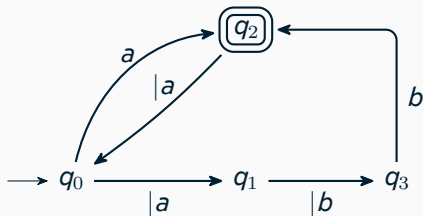
Question: Can't we just use classical learning algorithms for bar automata (as they are just DFA with additional semantics)?



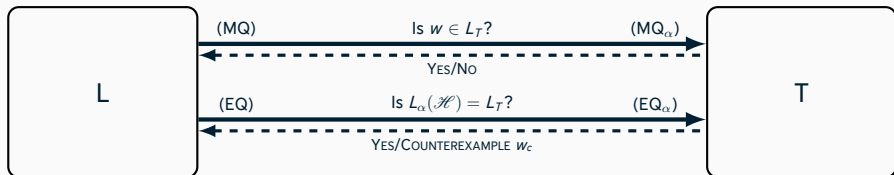
Question: Can't we just use classical learning algorithms for bar automata (as they are just DFA with additional semantics)?



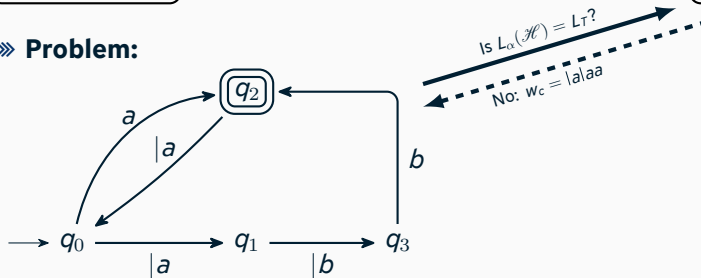
» **Problem:**



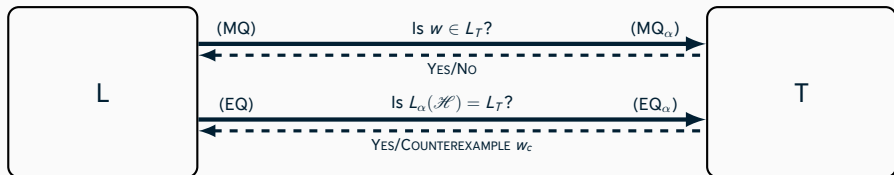
Question: Can't we just use classical learning algorithms for bar automata (as they are just DFA with additional semantics)?



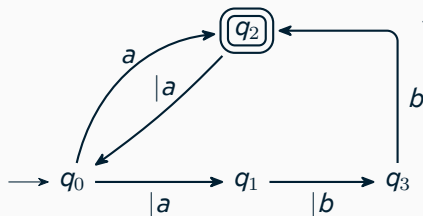
» **Problem:**



Question: Can't we just use classical learning algorithms for bar automata (as they are just DFA with additional semantics)?



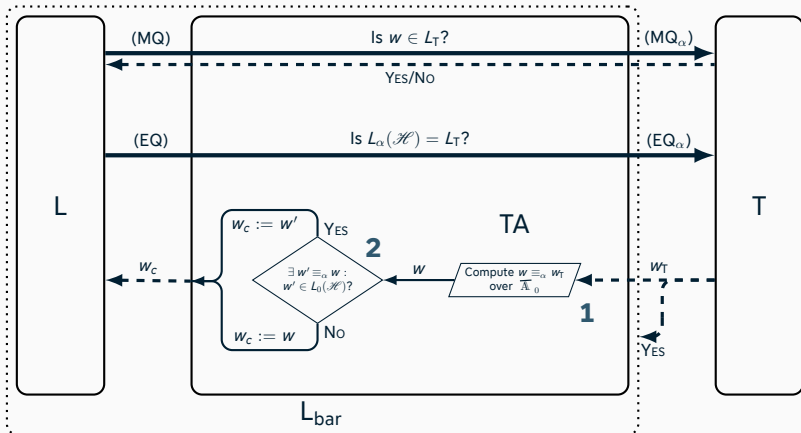
» **Problem:**



Is $L_\alpha(\mathcal{H}) = L_T$?
No: $w_c = |a|a$

w_c is **not** a counterexample in the classical sense.
(but $|a|bb \equiv w_c$ would be)

» Indeed, this is the only problem with the previous approach:





Given some bar word $w_T \in \overline{\mathbb{A}}^*$, compute $w \equiv_\alpha w_T$ over $\overline{\mathbb{A}}_0^*$



Given some bar word $w_T \in \overline{\mathbb{A}}^*$, compute $w \equiv_\alpha w_T$ over $\overline{\mathbb{A}}_0^*$

- » With λ -terms, checking if two terms are α -equivalent essentially boils down to computing a De Bruijn representation.



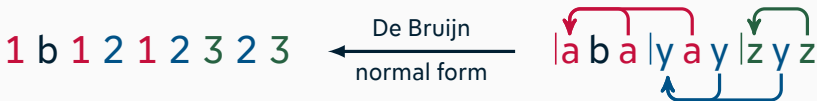
Given some bar word $w_T \in \overline{\mathbb{A}}^*$, compute $w \equiv_\alpha w_T$ over $\overline{\mathbb{A}}_0^*$

- » With λ -terms, checking if two terms are α -equivalent essentially boils down to computing a De Bruijn representation.
(\rightsquigarrow apply De Bruijn representations to bar strings)



Given some bar word $w_T \in \overline{\mathbb{A}}^*$, compute $w \equiv_\alpha w_T$ over $\overline{\mathbb{A}}_0^*$

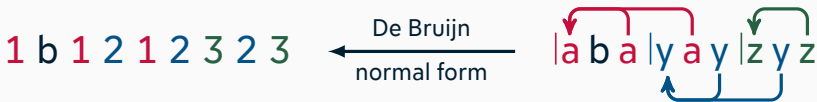
- » With λ -terms, checking if two terms are α -equivalent essentially boils down to computing a De Bruijn representation.
 (\rightsquigarrow apply De Bruijn representations to bar strings)





Given some bar word $w_T \in \overline{\mathbb{A}}^*$, compute $w \equiv_\alpha w_T$ over $\overline{\mathbb{A}}_0^*$

- » With λ -terms, checking if two terms are α -equivalent essentially boils down to computing a De Bruijn representation.
 (\rightsquigarrow apply De Bruijn representations to bar strings)



- » These normal forms are *unique* (per equivalence class) and computable in *polynomial time* (with linear-logarithmic space).



Given some $w \in \overline{\mathbb{A}}_0$, is there a $w' \equiv_\alpha w$ with $w' \in L(\mathcal{H})$?



Given some $w \in \overline{\mathbb{A}}_0$, is there a $w' \equiv_\alpha w$ with $w' \in L(\mathcal{H})$?

» Easy non-deterministic algorithm (in NP):



Given some $w \in \overline{\mathbb{A}}_0$, is there a $w' \equiv_{\alpha} w$ with $w' \in L(\mathcal{H})$?

» Easy non-deterministic algorithm (in NP):

Algorithm

- 1) Guess a bar string of equal length;
- 2) Check for α -equivalence.



Given some $w \in \overline{\mathbb{A}}_0$, is there a $w' \equiv_{\alpha} w$ with $w' \in L(\mathcal{H})$?

» Easy non-deterministic algorithm (in NP):

Algorithm

- 1) Guess a bar string of equal length;
- 2) Check for α -equivalence.

» Is a deterministic poly-time algorithm possible?



Given some $w \in \overline{\mathbb{A}}_0$, is there a $w' \equiv_{\alpha} w$ with $w' \in L(\mathcal{H})$?

» Easy non-deterministic algorithm (in NP):

Algorithm

- 1) Guess a bar string of equal length;
- 2) Check for α -equivalence.

» Is a deterministic poly-time algorithm possible?

not exactly ...



Given some $w \in \overline{\mathbb{A}}_0$, is there a $w' \equiv_{\alpha} w$ with $w' \in L(\mathcal{H})$?

» Easy non-deterministic algorithm (in NP):

Algorithm

- 1) Guess a bar string of equal length;
- 2) Check for α -equivalence.

» Is a deterministic poly-time algorithm possible? not exactly ...
» In general: NP-completeness (via the Hamilton cycle problem) ...



Given some $w \in \overline{\mathbb{A}}_0$, is there a $w' \equiv_\alpha w$ with $w' \in L(\mathcal{H})$?

» Easy non-deterministic algorithm (in NP):

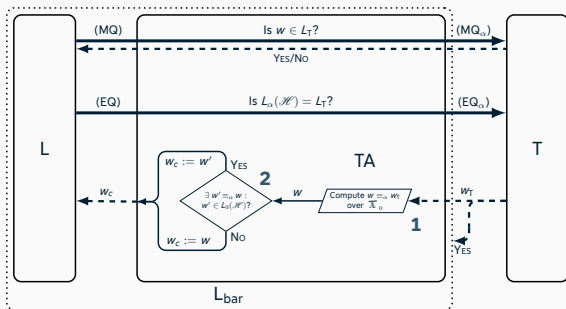
Algorithm

- 1) Guess a bar string of equal length;
- 2) Check for α -equivalence.

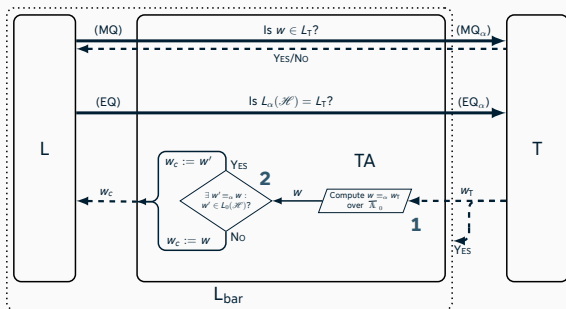
» Is a deterministic poly-time algorithm possible? not exactly ...

- » In general: NP-completeness (via the Hamilton cycle problem) ...
- » ... but for fixed alphabets in deterministic poly-time.
by comparing with the *closed* bar automaton.

» By solving both problems, correctness of our approach is shown.



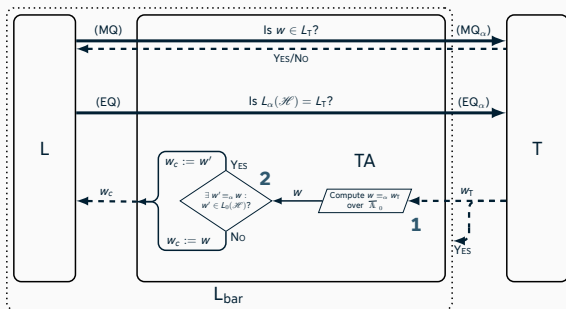
» By solving both problems, correctness of our approach is shown.



Query Complexity

L_{bar} asks at most as many queries as L would.

» By solving both problems, correctness of our approach is shown.



Query Complexity

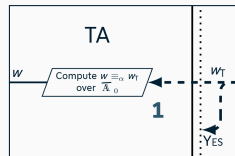
L_{bar} asks at most as many queries as L would.

» Implicit Assumption: L_{bar} knows the *number of registers* (size of alphabet) needed for L_T .

Idea

If L_{bar} gets stuck, just restart anew (with an extended alphabet).

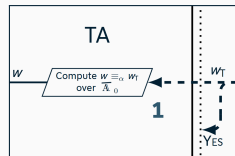
- 1) Start with the empty alphabet $\overline{\mathbb{A}}_0 = \emptyset$.



Idea

If L_{bar} gets stuck, just restart anew (with an extended alphabet).

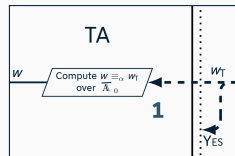
- 1) Start with the empty alphabet $\overline{\mathbb{A}}_0 = \emptyset$.
- 2) **Repeat:**



Idea

If L_{bar} gets stuck, just restart anew (with an extended alphabet).

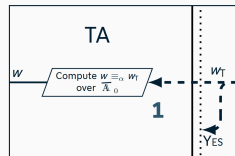
- 1) Start with the empty alphabet $\overline{\mathbb{A}}_0 = \emptyset$.
- 2) **Repeat:**
 - a) If a **correct** hypothesis is found, **terminate**.



Idea

If L_{bar} gets stuck, just restart anew (with an extended alphabet).

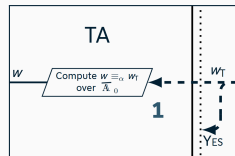
- 1) Start with the empty alphabet $\overline{\mathbb{A}}_0 = \emptyset$.
- 2) **Repeat:**
 - a) If a **correct** hypothesis is found, **terminate**.
 - b) If T delivers a counterexample w_T and step **1** works, **continue**.



Idea

If L_{bar} gets stuck, just restart anew (with an extended alphabet).

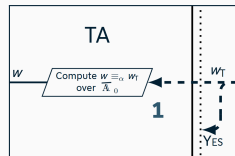
- 1) Start with the empty alphabet $\overline{A}_0 = \emptyset$.
- 2) **Repeat:**
 - a) If a **correct** hypothesis is found, **terminate**.
 - b) If T delivers a counterexample w_T and step **1** **works, continue**.
 - c) **Else:** Extend \overline{A}_0 to $\overline{A}_0 \subseteq \overline{A}'_0$ and try again.



Idea

If L_{bar} gets stuck, just restart anew (with an extended alphabet).

- 1) Start with the empty alphabet $\bar{A}_0 = \emptyset$.
- 2) **Repeat:**
 - a) If a **correct** hypothesis is found, **terminate**.
 - b) If T delivers a counterexample w_T and step **1** **works, continue**.
 - c) **Else:** Extend \bar{A}_0 to $\bar{A}_0 \subseteq \bar{A}'_0$ and try again.

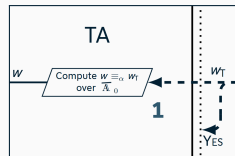


Choose extension via De Bruijn normal form.

Idea

If L_{bar} gets stuck, just restart anew (with an extended alphabet).

- 1) Start with the empty alphabet $\bar{A}_0 = \emptyset$.
- 2) **Repeat:**
 - a) If a **correct** hypothesis is found, **terminate**.
 - b) If T delivers a counterexample w_T and step **1** **works, continue**.
 - c) **Else:** Extend \bar{A}_0 to $\bar{A}_0 \subseteq \bar{A}'_0$ and try again.



Choose extension via De Bruijn normal form.

Theorem (*Extension Complexity*)

L_{bar} can infer a bar automaton using the minimal alphabet with at most as many queries as L would need (summed over all smaller cardinalities).

- » Next to finite word languages, we extended our approach to bar ω - and bar tree languages:

- » Next to finite word languages, we extended our approach to bar ω - and bar tree languages:

Alphatic Trees

As expected, everything works out analogously:

- » Normal forms are computed branchwise (and thus still in *poly-time*).

- » Next to finite word languages, we extended our approach to bar ω - and bar tree languages:

Alphatic Trees

As expected, everything works out analogously:

- » Normal forms are computed branchwise (and thus still in *poly-time*).
- » Finding accepted bar trees is NP-complete, but also in para. poly-time.

- » Next to finite word languages, we extended our approach to bar ω - and bar tree languages:

Alphatic Trees

As expected, everything works out analogously:

- » Normal forms are computed branchwise (and thus still in *poly-time*).
- » Finding accepted bar trees is NP-complete, but also in para. poly-time.
- » Handling unknown alphabets works completely identical.

- » Next to finite word languages, we extended our approach to bar ω - and bar tree languages:

Alphatic Trees

As expected, everything works out analogously:

- » Normal forms are computed branchwise (and thus still in *poly-time*).
- » Finding accepted bar trees is NP-complete, but also in para. poly-time.
- » Handling unknown alphabets works completely identical.

Challenges in the infinite word case

- » Next to finite word languages, we extended our approach to bar ω - and bar tree languages:

Alphatic Trees

As expected, everything works out analogously:

- » Normal forms are computed branchwise (and thus still in *poly-time*).
- » Finding accepted bar trees is NP-complete, but also in para. poly-time.
- » Handling unknown alphabets works completely identical.

Challenges in the infinite word case

- » There is no suitable normal form for infinite strings.

(example at blackboard)

- » Next to finite word languages, we extended our approach to bar ω - and bar tree languages:

Alphatic Trees

As expected, everything works out analogously:

- » Normal forms are computed branchwise (and thus still in *poly-time*).
- » Finding accepted bar trees is NP-complete, but also in para. poly-time.
- » Handling unknown alphabets works completely identical.

Challenges in the infinite word case

- » There is no suitable normal form for infinite strings.
(example at blackboard)
- » Checking α -equivalence is still in poly-time, but guessing seemingly impossible.

- » Next to finite word languages, we extended our approach to bar ω - and bar tree languages:

Alphatic Trees

As expected, everything works out analogously:

- » Normal forms are computed branchwise (and thus still in *poly-time*).
- » Finding accepted bar trees is NP-complete, but also in para. poly-time.
- » Handling unknown alphabets works completely identical.

Challenges in the infinite word case

- » There is no suitable normal form for infinite strings.
(example at blackboard)
- » Checking α -equivalence is still in poly-time, but guessing seemingly impossible.
- » Expensive computation of closures is important!

- » Learnability of various kinds of bar languages in Angluin's framework.
- » Our learner is in terms of query complexity as optimal as the underlying learner.
- » Introduced efficient procedures for checking α -equivalence.



Future Work

- » Is guessing possible for bar ω -languages? Can the computation of closures be removed?
- » Can this approach be extended to efficiently learn data languages?
- » What to do about conformance testing?

Questions?



Friedrich-Alexander-Universität
Faculty of Engineering

-  Angluin, Dana. **'Learning regular sets from queries and counterexamples'**. *Information and Computation* 75.2 (Nov. 1987), pp. 87–106. ISSN: 0890-5401. DOI: 10.1016/0890-5401(87)90052-6. URL: [https://doi.org/10.1016/0890-5401\(87\)90052-6](https://doi.org/10.1016/0890-5401(87)90052-6).
-  Schröder, Lutz, Dexter Kozen, Stefan Milius, Thorsten Wißmann. **'Nominal Automata with Name Binding'**. *Proc. 20th International Conference on Foundations of Software Science and Computation Structures, (FOSSACS 2017)*. Vol. 10203. Lect. Notes Comput. Sci. 2017, pp. 124–142.