

Inferring Regular Languages & ω -Languages

DANA FISMAN
Ben-Gurion University

based on joint works with

DANA ANGLUIN, UDI BOKER & SARAH EISENSTAT

Synthesis

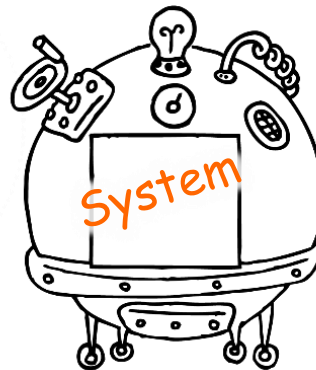
Challenges:

- Hard to characterize using a logical calculus
- Complete bugless spec, really!?



Specification
High Level
What?
Declarative
Ex: temporal logic

Synthesizer



Correct by construction

Implementation
Low Level
How?
Procedural/Executable
Ex: reactive system

A specification scale

examples

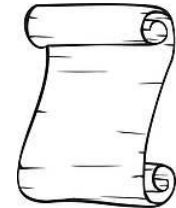
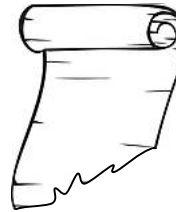
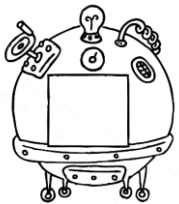
partial
spec

complete
rigorous
mathematical
specification

$\{e_1, e_2, e_3, \dots\}$

Learning

Synthesizer

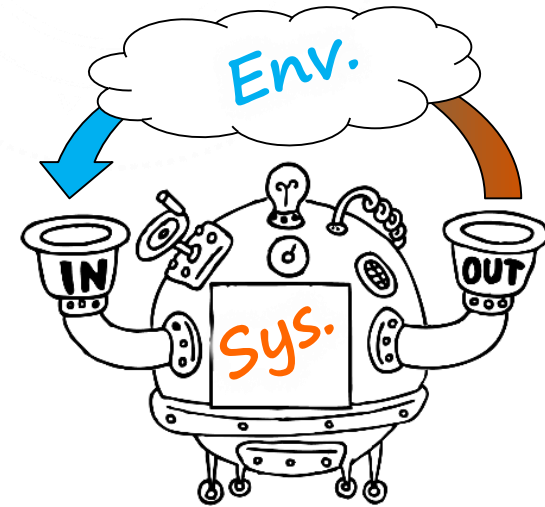


partial
implementation

What kind of examples?

In the context of synthesizing **reactive systems**:

- The **examples** are **words / strings** describing **computations / interfaces**
- The **learned concept** is a set of such examples, presumably a **regular language**.
- For regular languages [Angluin, 1987] suggested **L^*** algorithm.
- **L^*** learns in polynomial time an **unknown regular language** using **membership and equivalence queries**.



L^* - Active Learning with MQ and EQ



Learner



Teacher

Is w in L ?



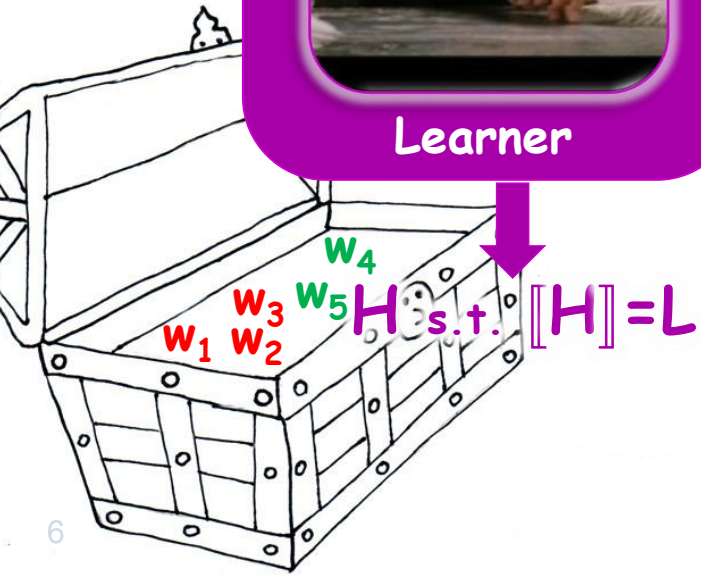
Yes / No



Is $\llbracket H \rrbracket$ same as L ?



Yes / No, c.e: w'



Usages of L^*

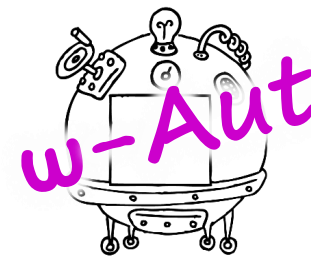
- L^* is an **extremely popular** algorithm. It has applications in many areas including **AI**, **neural networks**, **geometry**, **data mining**, **verification** and **synthesis**.
- Usages of L^* in **verification** and **synthesis** include:
 - Black-box checking [Peled et al.]
 - Assume-guarantee reasoning [Cobleigh et al.]
 - Specification mining [Ammons et al., Gabel et al., ...]
 - Error localization [Chapman et al.]
 - Learning interfaces [Alur et al.]
 - Regular Model Checking [Habermehl & Vonjar]
 - ...

Challenge 1

- L^* learns a regular language of **finite words**. Interesting properties of **reactive systems** e.g. (**liveness** and **fairness**) are **not** expressible by **finite words**.
- Can we extend L^* to L^ω , an alg. that learns regular languages of **infinite words** (**w-words**)?

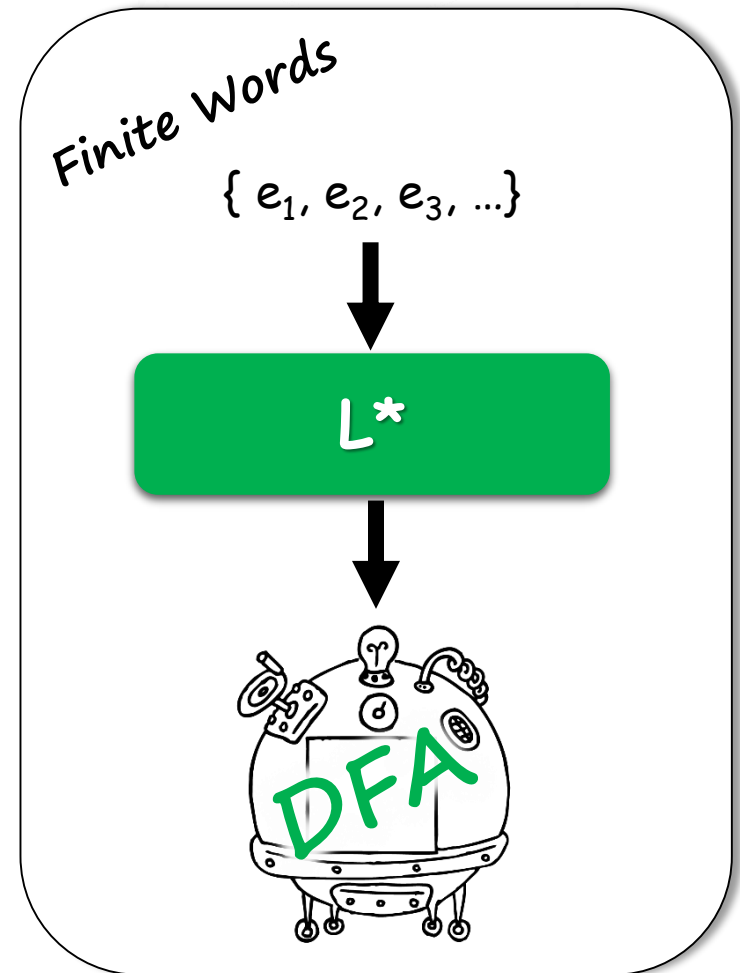
Infinite Words

$\{e_1 = \text{abcdcbaadcacbbccaabcc}$
 $e_2 = \text{bbbcdcaaacbccccccaabcc}$
 $\}$



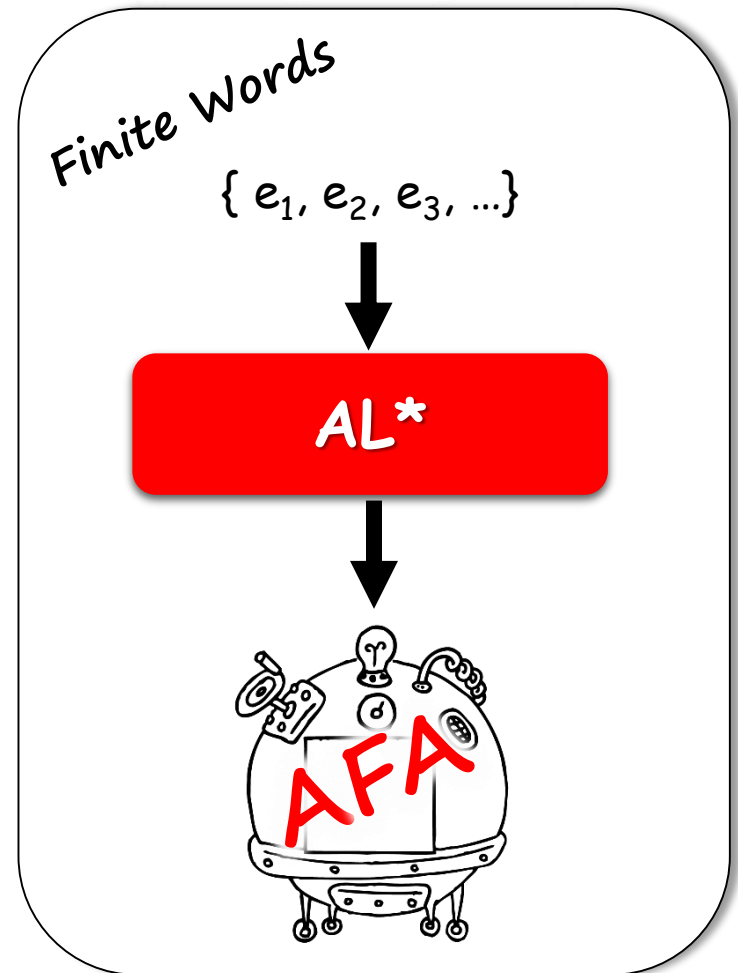
Challenge 2

- L^* produces DFAs (deterministic finite automata), a well behaved representation, yet not a compact one.
- Can we learn more succinct representations, such as non-deterministic finite automata (NFA) or alternating automata (AFA)?

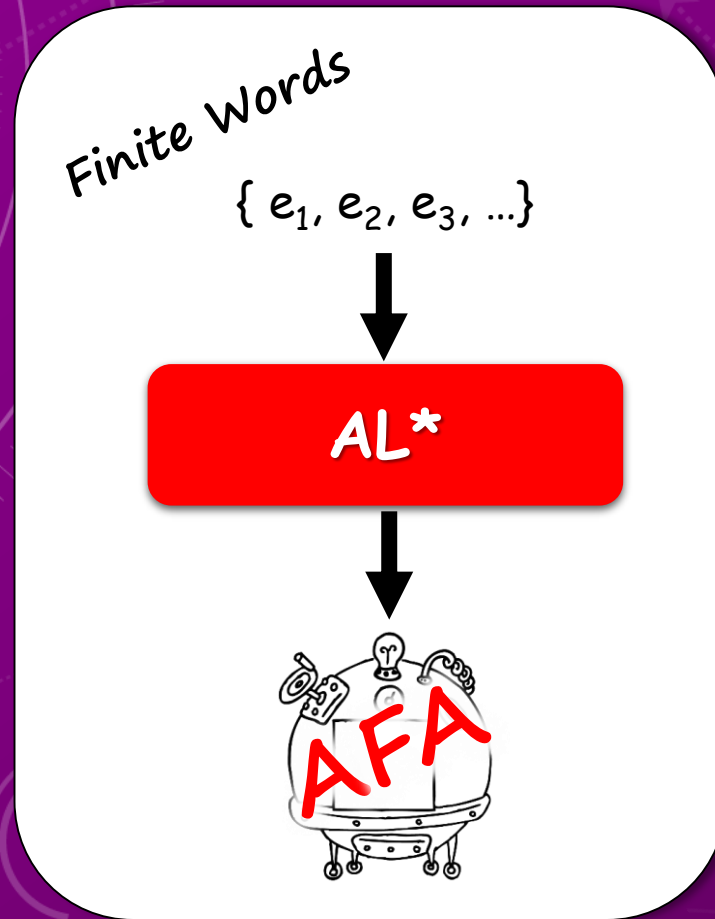


Challenge 2

- L^* produces DFAs (deterministic finite automata), a well behaved representation, yet not a compact one.
- Can we learn more succinct representations, such as non-deterministic finite automata (NFA) or alternating automata (AFA)?



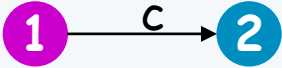
Learning alternating automata



[ANGLUIN, EISENSTAT & FISMAN IJCAI'15]

What are alternating automata?

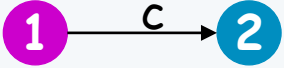
Transition Type	from state	upon reading	to state(s)
Deterministic	s1	c	s2



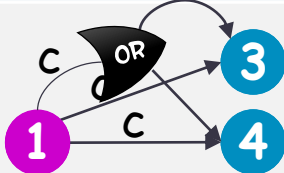
The diagram illustrates a transition from state 1 to state 2. State 1 is represented by a pink circle with the number 1 inside, and state 2 is represented by a blue circle with the number 2 inside. A horizontal arrow points from state 1 to state 2, with the letter 'c' positioned above the arrow.

What are alternating automata?

Transition Type	from state	upon reading	to state(s)
Deterministic	$s1$	c	$s2$
Non-Deterministic	$s1$	c	$s3$ or $s4$



The diagram shows a single transition from state 1 (pink circle) to state 2 (blue circle) labeled with the input symbol 'c'.

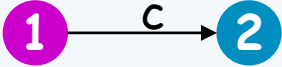
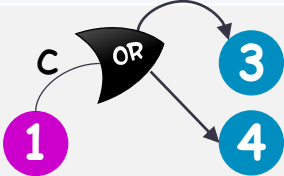
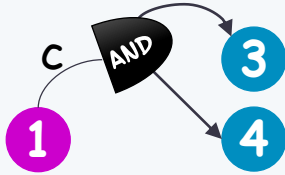
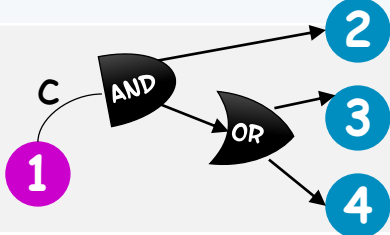


The diagram shows a non-deterministic transition from state 1 (pink circle) to states 3 and 4 (blue circles) on input 'c'. An OR gate symbol is positioned above the transitions to state 3 and 4, indicating that the automaton can choose either state as the next state.

What are alternating automata?

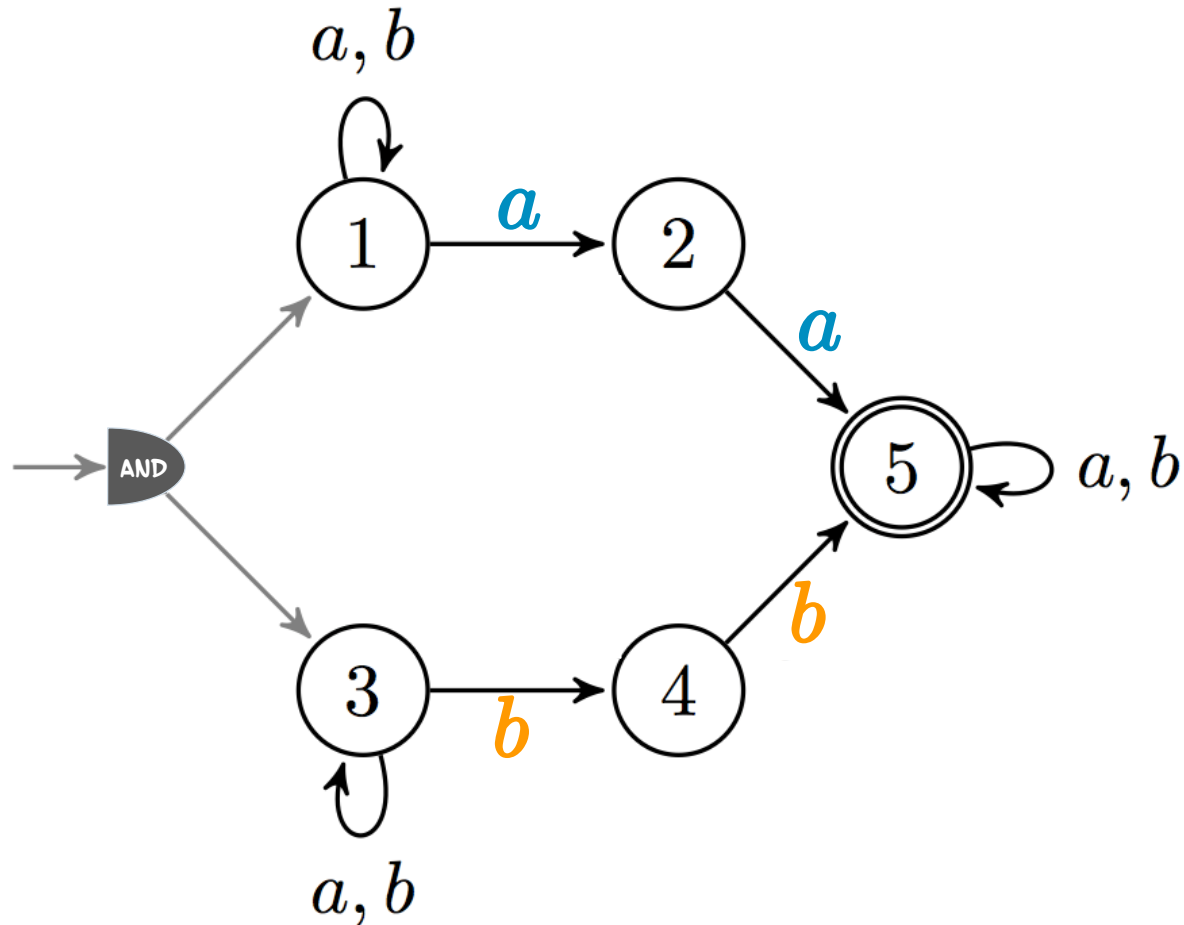
Transition Type	from state	upon reading	to state(s)	
Deterministic	s_1	c	s_2	
Non-Deterministic	s_1	c	s_3 or s_4	
Universal	s_1	c	s_3 and s_4	

What are alternating automata?

Transition Type	from state	upon reading	to state(s)	
Deterministic	s_1	c	s_2	
Non-Deterministic	s_1	c	s_3 or s_4	
Universal	s_1	c	s_3 and s_4	
Alternating	s_1	c	$(s_3$ or $s_4)$ and s_2	

Alternating Automaton - Ex.

$$\Sigma = \{a, b\}$$



Accepts the language $\Sigma^*aa\Sigma^* \cap \Sigma^*bb\Sigma^*$

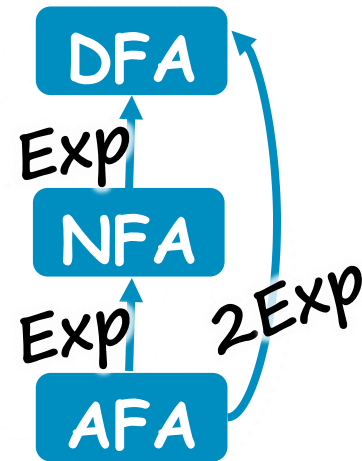
What are they good for?

- AFAs are a **succinct representation**
- The PSL formula

```
always (print-to-both ->  
  ([*], print-a-start, busy[*3..], print-a-end) &  
  ([*], print-b-start, busy[*3..], print-b-end))
```

can be stated by a **12** state **AFA** but the **minimal DFA** requires **115** states.

- **Natural means** to model **conjunctions** and **disjunctions** as well as **existential** and **universal quantification**
- **1-to-1 translations** from **temporal logics**
- Working at the **alternating level** enables **better structured algorithms**, and is the common practice in industry **verification tools**.



Foundation of L^* - Residuality

The **residual** of language L with respect to word u is the set of all words v such that uv in L

$$u^{-1}L = \{ v \mid uv \in L \}$$

Example

$$L = aba^*$$

$$a^{-1}L = ba^*$$

$$ab^{-1}L = a^*$$

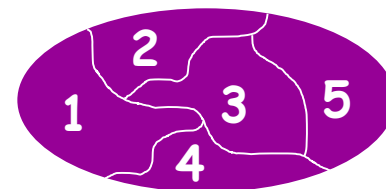
$$abaaa^{-1}L = a^*$$

$$b^{-1}L = \emptyset$$

If $u^{-1}L = v^{-1}L$ we say that $u \sim_L v$.

$$ab \sim_L abaaa$$

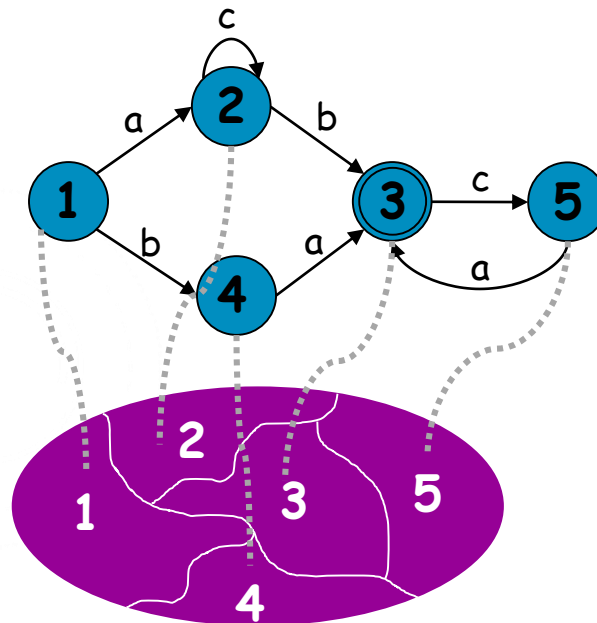
The **residuality index** is the number of equivalence classes of \sim_L



Myhill-Nerode THM

Every regular language L has a **finite number** of residual languages.

The minimal DFA has one state for every residual language of L !!!



Challenge

NFAs and AFAs don't have the residually property, in general.

Residual NFAs

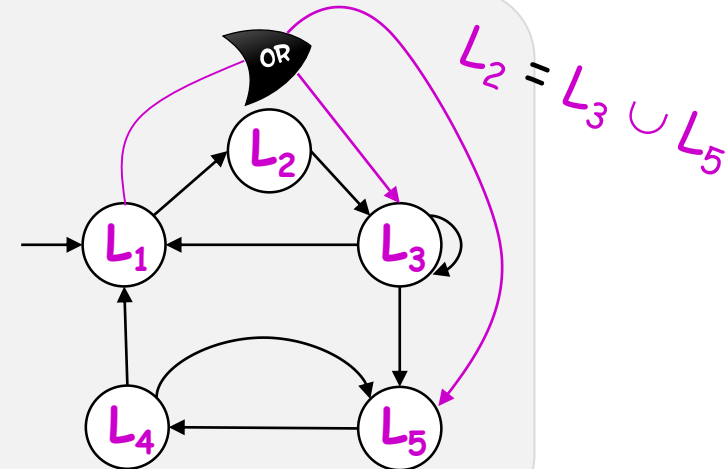
- Dennis et al. [STACS' 01] defined **residual NFAs (NRFA)**
- These are **NFAs** where **each state** corresponds to a **residual language**

Suppose L_1, L_2, \dots, L_n are all the **residual languages** of L

If for some L_i , we have

$$L_i = L_j \cup L_k$$

then we can **remove** the **i^{th} state**,
and use **non-determinism** to
capture it.



Residual NFAs

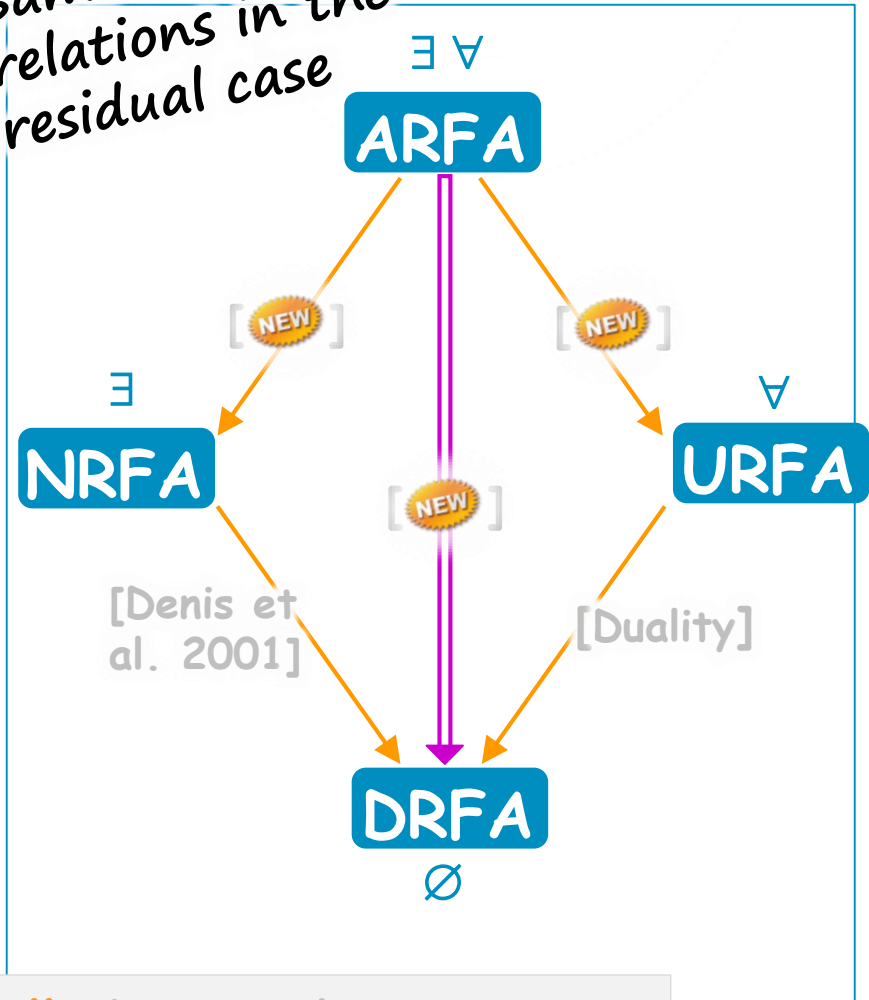
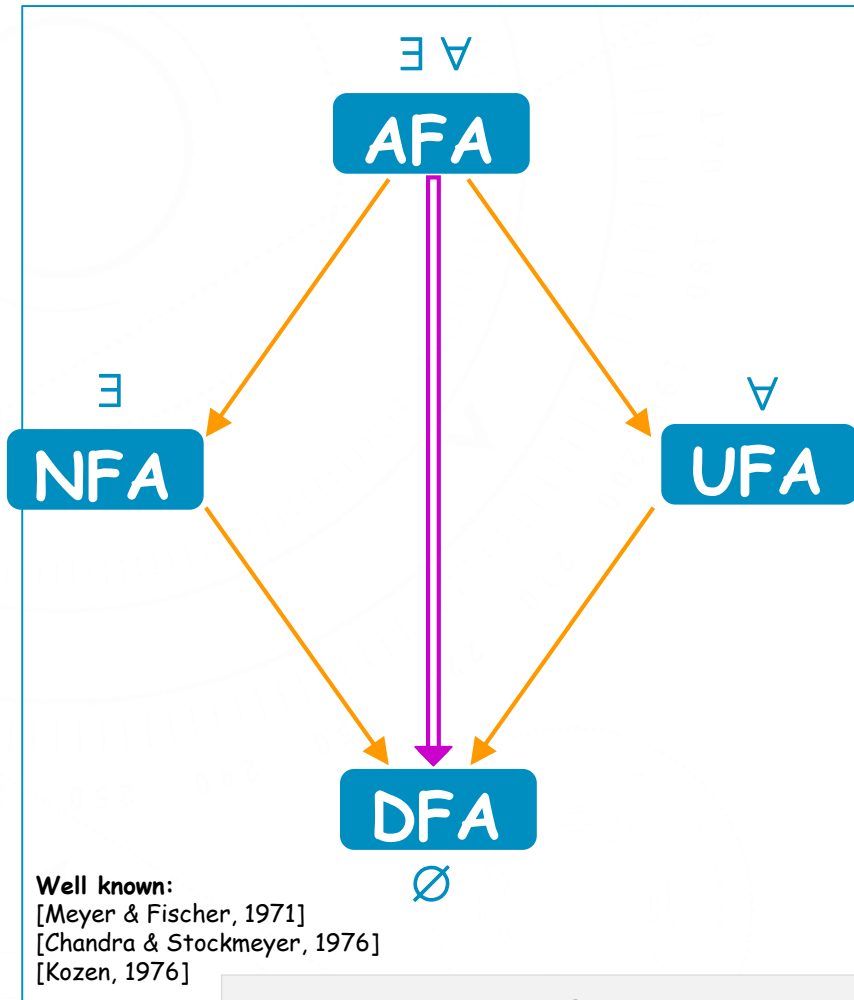
- Dennis et al. showed/provided
 - Every regular language is recognized by a **unique** (canonical) **NRFA** which has a **minimal** number of **states** and a maximal number of transitions.
 - There may be **exponential gaps** between the **minimal DFA**, the **canonical NRFA** and the **minimal NFA**.
- Bollig et al. [IJCAI'09] extended L^* to NL^* (learns NRFA)

Questions

- Can we **extend** the notion of **residually** to **AFAs**?
- Will **exponential gaps** remain?
- Can we define a **canonical one**?
- Can we **learn ARFAs**?

Succinctness

Same size relations in the residual case



- May be exponentially bigger than
- ⇨ May be doubly exponentially bigger than

The learning algorithm

- L^* uses a data structure termed an **observation table**.
- AL^* generalizes NL^* and L^* and the notion of a complete/minimal observation table.
- As shown next...

The table of residual languages

Enumeration
of all strings

ϵ a b aa ab ab bb aaa aab aba abb baa

all the suffixes of
ab that are in L
i.e. $ab^{-1}L$

aa

ab

By Myhill-Nerode the
number of *distinct rows*
is *finite*.

aba

abb

baa

bab

ϵ	0	0	1	1	0	0	0	1	0	1	0	0	1	0
a	0	0	0	1	0	1	0	1	1	0	1	0	0	1
b	1	0	1	1	0	1	1	1	1	0	1	1	1	1
aa	1	1	0	0	1	1	1	0	1	1	1	0	0	0
ab	1	0	0	1	1	0	0	0	1	0	1	0	0	1
ab	0	1	0	0	1	0	1	0	1	1	0	1	0	0
bb	0	1	0	1	1	0	0	1	0	1	1	0	0	1
aaa	0	1	0	1	1	0	1	1	0	1	1	1	1	1
aab	1	0	0	1	1	1	1	0	0	1	1	1	0	0
aba	0	1	1	0	0	0	0	1	1	0	0	0	1	0
abb	1	1	0	0	1	1	0	0	1	0	1	0	1	0
baa	0	0	1	1	0	0	1	0	0	1	1	0	1	0
bab	1	0	0	1	0	1	0	1	0	1	1	1	0	0

The table of residual languages

The number of **distinct columns** is also **finite**.

We call it the **column index**.

Enumeration of all strings

a b aa ab ab bb aaa aab aba abb baa

Enumeration of all strings

	1	0	0	1	1	0	0	0	1	0	1	0	0	1	0
	0	1	0	0	1	0	1	0	1	1	0	1	0	0	1
<i>b</i>	0	1	0	1	1	0	1	1	1	1	0	1	1	1	1
<i>aa</i>	1	1	1	0	0	1	1	1	0	1	1	1	0	0	0
<i>ab</i>	1	0	0	1	1	0	0	0	1	0	1	0	0	1	0
<i>ba</i>	0	1	0	0	1	0	1	0	1	1	0	1	0	0	1
<i>bb</i>	0	0	1	0	1	1	0	0	1	0	1	1	0	0	1
<i>aaa</i>	0	1	1	0	1	1	0	1	1	0	1	1	1	1	1
<i>aab</i>	1	0	0	1	1	1	1	0	0	1	1	1	0	0	0
<i>aba</i>	0	1	1	0	0	0	0	1	1	0	0	0	1	0	1
<i>abb</i>	1	1	0	0	1	1	0	0	1	0	1	0	1	1	0
<i>baa</i>	0	0	1	1	0	0	1	0	0	1	1	0	1	1	0
<i>bab</i>	1	0	0	1	0	1	0	1	0	1	1	1	0	0	0

L* Data Structure

An Observation Table:

Strings:
experiments to
distinguish states

Strings:
candidate
state
represent-
atives

M	e_1	e_2	e_3	e_4	e_5	e_6	...
s_1	1	0	0	1	1		
s_2	0	1	0	0	1		
s_3	0	1	0	1	1		
s_4	1	1	1	0	0		
s_5	1	0	0	1	1		
⋮							

$$M_{i,j} = \begin{cases} 1 & \text{if } s_i e_j \in L \\ 0 & \text{otherwise} \end{cases}$$

Closed Table

An observation table $T = (S, E, M)$
is **closed** w.r.t a subset $B \subseteq S$

S	e_1	e_2	e_3	e_4	e_5	e_6
▶ ϵ	1	0	0	1	1	
▶ a	0	1	0	0	1	
b	1	0	0	1	1	covered
▶ ab	1	0	0	1	1	covered
▶ aa	1	1	1	0	0	
aaa	1	0	0	1	1	covered
aab	0	1	0	0	1	covered

If it satisfies

- 1) Initialization: $\epsilon \in B$
- 2) Consecution: $B\Sigma \subseteq S$
- 3) Coverage: all rows **not** in B are covered by some row in B

The definition of **covers** differs for L^* , NL^* and AL^* .

D-Covered

According to L^*
i.e. when
using DFAs

S	e_1	e_2	e_3	e_4	e_5	e_6
ϵ	1	0	0	1	1	
a	0	1	0	0	1	
b	0	1	0	0	1	
ab	1	0	0	1	1	
aa	1	1	1	0	0	
aaa	1	0	0	1	1	
aab	0	1	0	0	1	

Same
Same
Same
Same

N-Covered

According to NL^*
i.e. when
using NFAs

S	e_1	e_2	e_3	e_4	e_5	e_6
ϵ	1	0	0	1	1	
a	0	1	0	0	1	
b	1	1	0	1	1	
ab	1	1	1	0	1	
aa	1	1	1	0	0	
aaa	1	0	0	1	1	
aab	0	1	0	0	1	

Expressible as bitwise-or
of some rows in B

$$b = (\epsilon \vee a)$$

A-Covered

According to AL^*
i.e. when
using AFAs

S	e_1	e_2	e_3	e_4	e_5	e_6
ϵ	1	0	0	1	1	
a	0	1	0	0	1	
b	0	0	0	0	1	
ab	1	1	1	0	1	
aa	1	1	1	0	0	
aaa	1	0	0	1	1	
aab	0	1	0	0	1	

Expressible as
a monotone combination
of some rows in B

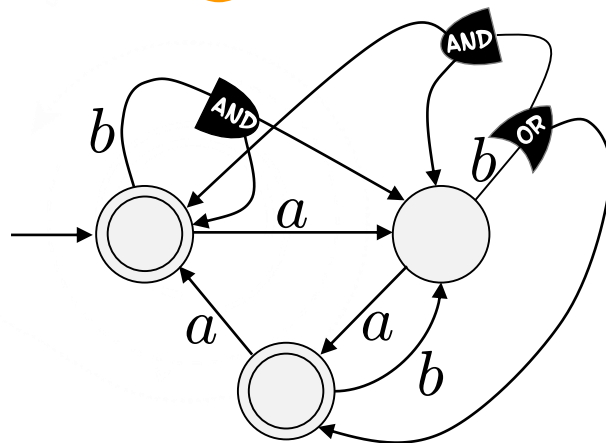
$$b = (\epsilon \wedge a)$$

$$ab = (\epsilon \wedge a) \vee aa$$

From Tables to Automata

Closed and Minimal

S	ϵ	e_2	e_3	e_4	e_5	e_6	
ϵ	1 ✓	0	0	1	1		
a	0	1	0	0	1		
b	0	0	0	0	1		$= (\epsilon \wedge a)$
ab	1	1	1	0	1		$= (\epsilon \wedge a) \vee aa$
aa	1 ✓	1	1	0	0		
aaa	1	0	0	1	1		$= \epsilon$
aab	0	1	0	0	1		$= a$



Need to solve

How to decide

- Is row s a union of rows in B ? *Poly time* [Bollig et al.]
- Is s a monotone combination of rows in B ? *Poly time* [NEW]

Given a set of Boolean vectors S , find a minimal

unique union basis

Poly time [Bollig et al.]

Not unique monotone basis

NP-complete [NEW]

No obvious canonical rep. for ARFAS

Let

$$S = \{0,1\}^3$$

Then both

$$B_1 = \{001, 010, 101\}$$

$$B_2 = \{110, 101, 011\}$$

are minimal monotone bases.

The Learning Alg.

Algorithm 1: XL^* for $X \in \{D, N, U, A\}$

oracles : MQ, EQ

members: Observation table $\mathcal{T} = (S, E, M)$,
Candidate states set P

methods : $IsXClosed$, $IsXMinimal$, X
 InB_X , $XExtractAut$

$S = \langle \epsilon \rangle$, $E = \langle \epsilon \rangle$, $P = \langle \epsilon \rangle$ and M_ϵ .

repeat

$(a_1,$
if a_1

else

if $a_2 = \text{"no"}$ then

$P.RemoveString(s_2)$

else

$A = \mathcal{T}.XExtractAut(P)$

$(a_3, s_3) = EQ(A)$

if $a_3 = \text{"no"}$ then

$\mathcal{T}.XFind\&AddCols(s_3)$

$= \text{"yes"}$

A

Start with

table

If the table is **not closed**, e.g. s_1 is missing,

then

If the table is **not minimal**, e.g. s_2 is redundant then

Ask an **equivalence query**.

If **true**, return.

Otherwise, use the given **counterexample** to find some **columns to add**, and add them.

Start with basis:

$\{\epsilon\}$

THM: Every **counterexample** yields at least one **new column**

Back to finite words

Theorem

The algorithm AL^* returns an AFA for the unknown language after at most

- m equivalence queries
- $O(|\Sigma| mnc)$ membership queries
- $\text{poly}(m, n, c, |\Sigma|)$ time

	L^*	NL^*	AL^*
EQ	n	$O(n^2)$	m
MQ	$O(\Sigma cn^2)$	$O(\Sigma cn^3)$	$O(\Sigma cnm)$

where

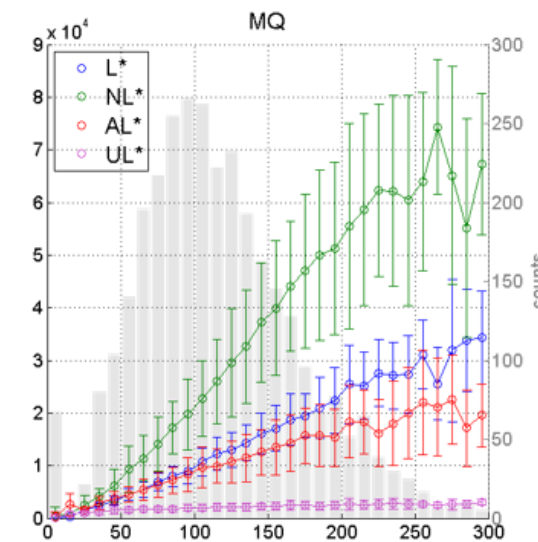
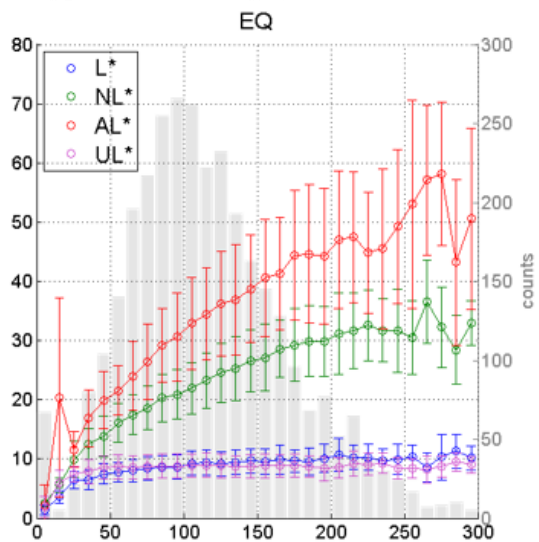
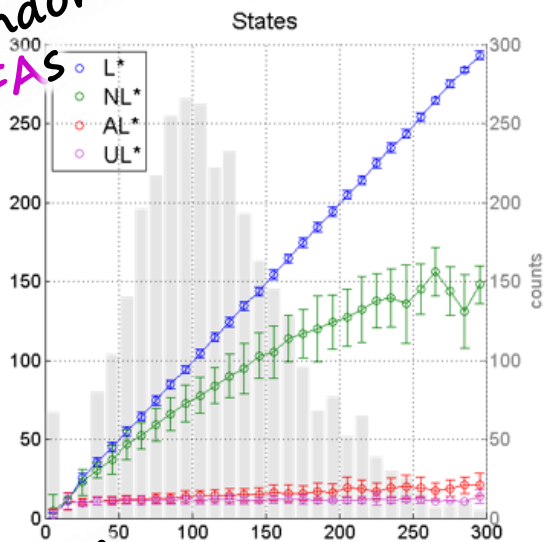
n = row index

m = column index

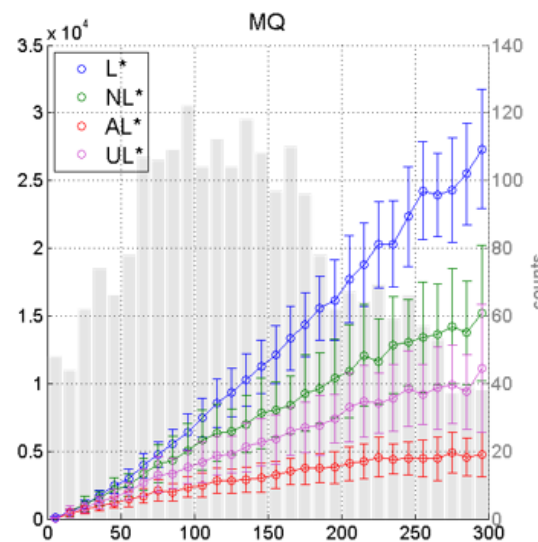
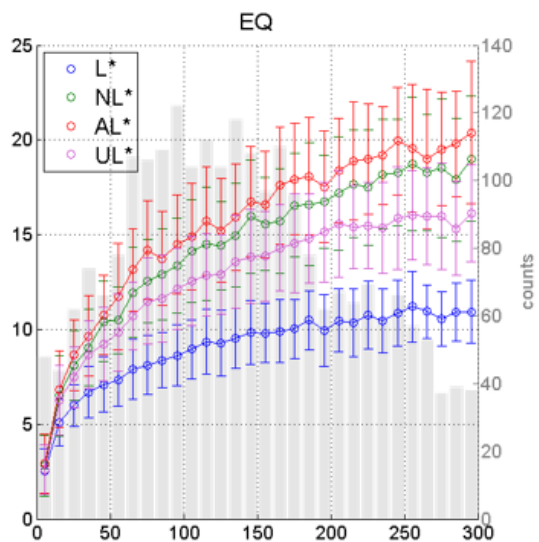
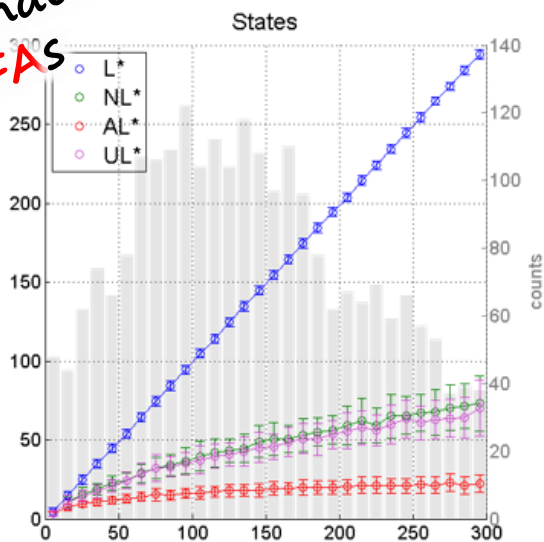
c = length of longest c.e.

Finite words - Empirical results

Random
UFAS



Random
AFAS



Finite words - Empirical results

Rough Summary:

- In terms of **#states** generated,
AL* is always preferable
- In terms of **#MQ**,
xL* outperforms the others when targets are **xFAs**
- In terms of **#EQ**,
L* is always preferable

Open questions & further directions

- Generalization to Boolean Automata ($\wedge \vee \neg$)
- Heuristics combining xL^* s
- Understanding of Residual AFAs
 - Properties of ARFAs
 - **Theorem:** The algorithm AL^* returns an AFA for the unknown language
 - **Conjecture:** The algorithm AL^* returns an ARFA for the unknown language

Coping with w-words

Is
abccccabdebbbaaaabcdaa...
in L?

prefixes

ultimately
periodic words
(Lasso words)



Learner



Teacher

Coping with w-words

Is
wingardium laviosa^w in L?

ultimately
periodic words
(Lasso words)



Learner



Teacher

Coping with w-words

Is

wingardium laviosa^w in L?

wingardium laviosa laviosa laviosa laviosa laviosa ...

ultimately
periodic words
(Lasso words)



Learner



Teacher

Coping with w -words

Is
wingardium **laviosa^w** in L ?

ultimately
periodic
(words
words)

THM:

Two **regular w -languages** are **equivalent**
iff
they agree on the set of **lasso words**

ω -automata

ω -automaton $(\Sigma, S, s_0, \delta, \alpha)$



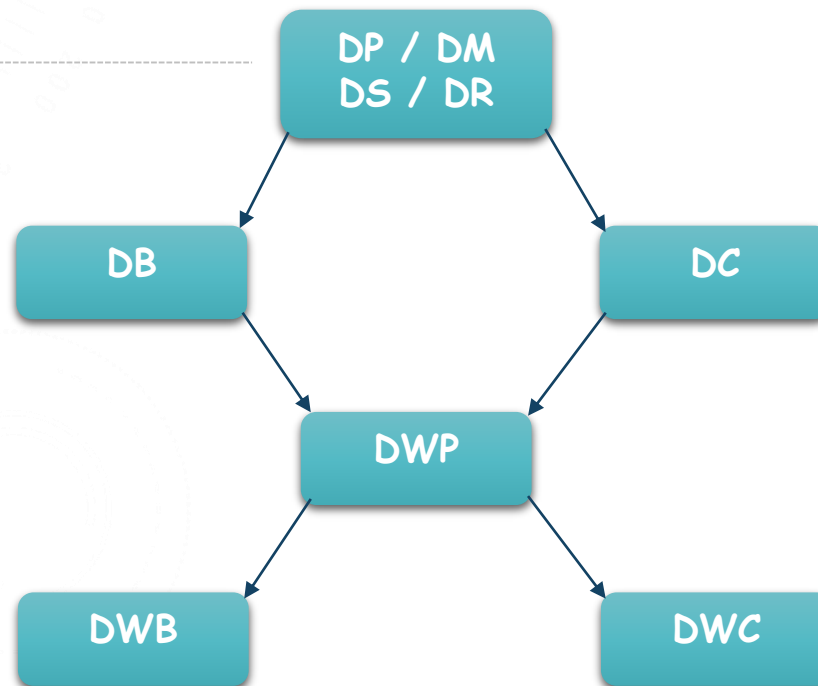
- There are many ways to define **acceptance** condition for **ω -Automata**
 - Büchi
 - Muller
 - Rabin
 - co-Büchi
 - Parity
 - Streett
- Roughly speaking, all are defined using the notion of the states **visited infinitely often** during a run.

w-automata - Expressiveness

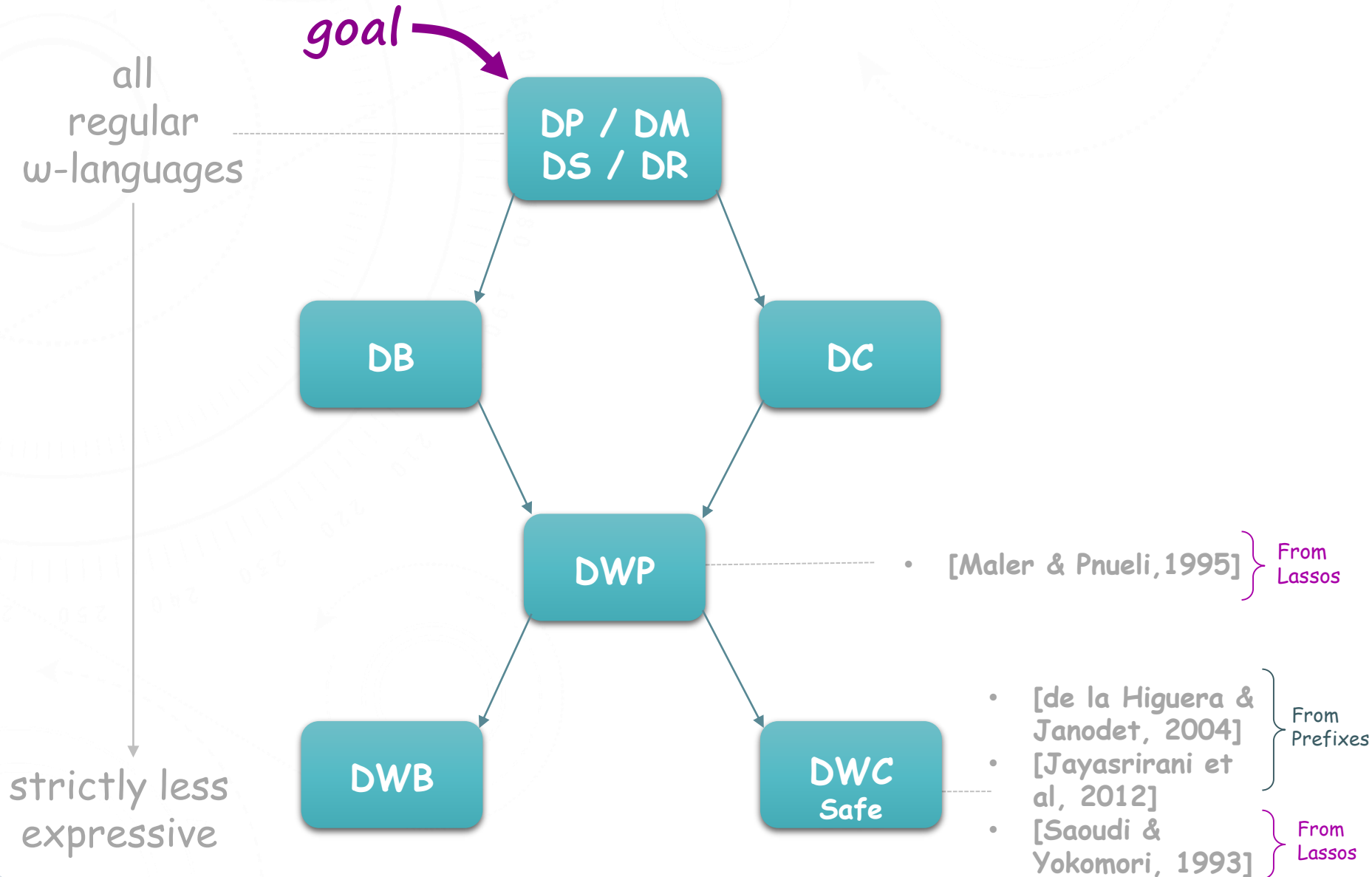
- Some acceptance criteria are **equally expressive**, some are **strictly less expressive** than others.
- Overall picture looks like this:

all
regular
w-languages

strictly less
expressive

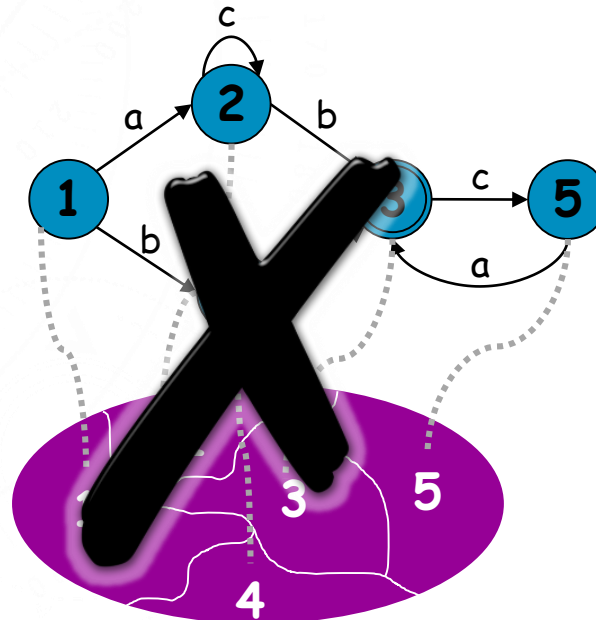


Previous work on learning w-langs.



Challenges

- L^* works due to the **Myhill-Nerode thm.**
- The major **difficulty** in learning w -languages is a **lack** of a corresponding **Myhill-Nerode theorem** for **w -automata** (of all types)



2014

2012

2008

2005

1994

1993

1987 L^*

1962 w -
aut.

Challenges

- It turns out that an **w-regular language** can be represented by a **regular language $L_\$$** of **finite words** [Calbrix, Nivat, Podelski 93]
- And thus one can use **L^*** to learn this representation [Farzan et al. 2008]
- However, this representation is quite **big**:
Büchi with **n** states \Rightarrow **DFA** for **$L_\$$** with **$2^n + 2^{2n^2+n}$**



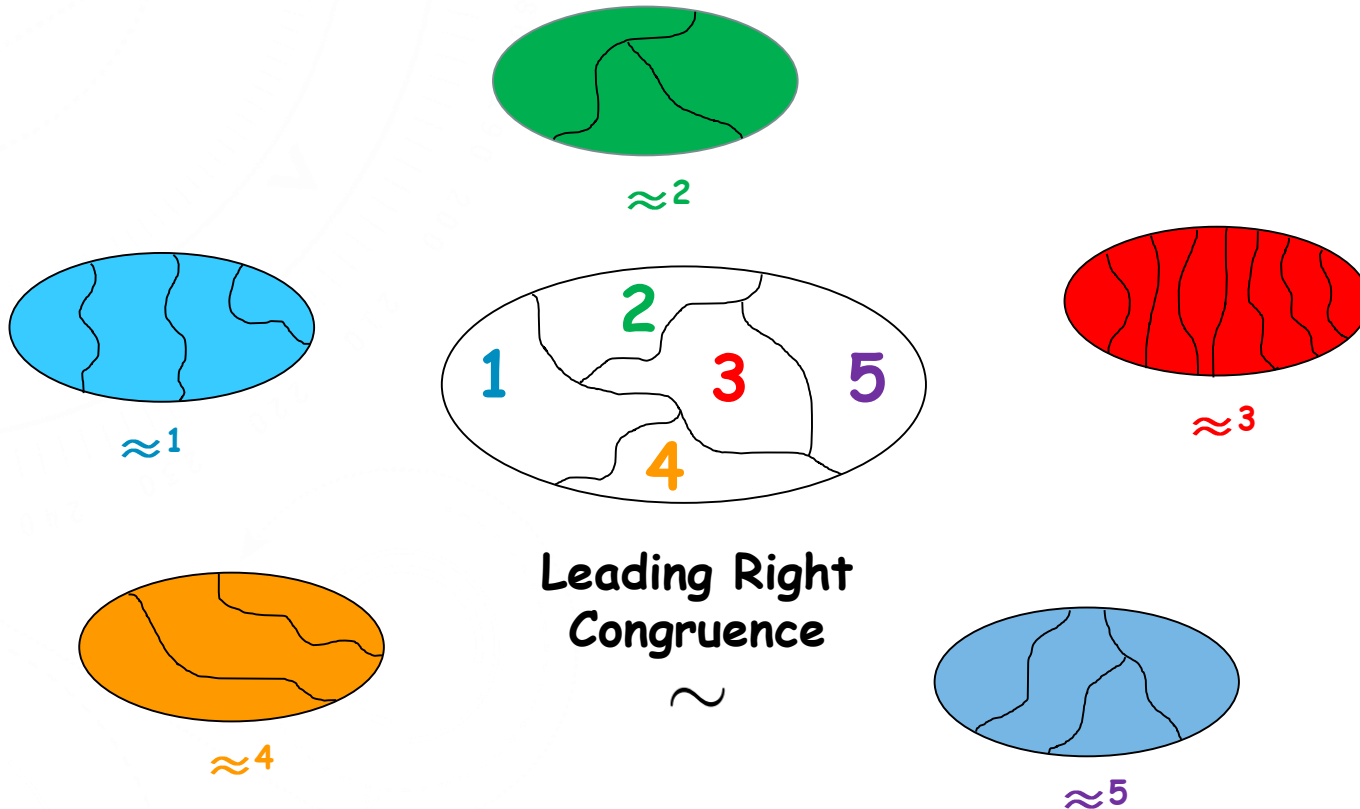
The way out

A new representation: **Family of DFAs** and a new canonical rep **Recurrent FDFAs** based on families of **FORCs** [Maler & Staiger, 95] and the **syntactic FORC** which has a **Myhill-Nerode theorem**

Family of Right Congruences [MS97]

$(\sim, \approx_1, \approx_2, \approx_3, \approx_4, \approx_5)$

Leading Progress



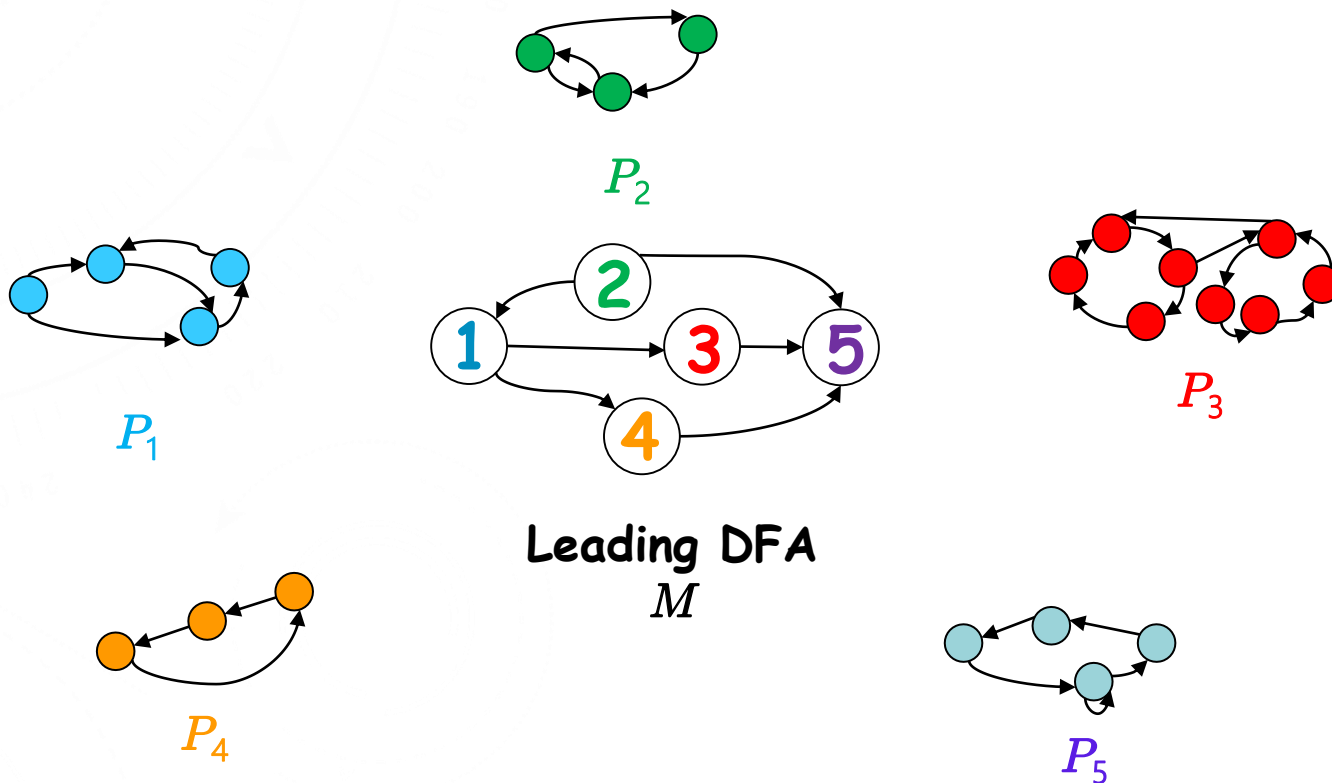
Plus some restriction (details omitted)

Family of DFAs (FDFA)

$(M, P_1, P_2, P_3, P_4, P_5)$

Leading

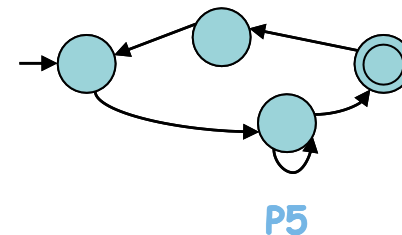
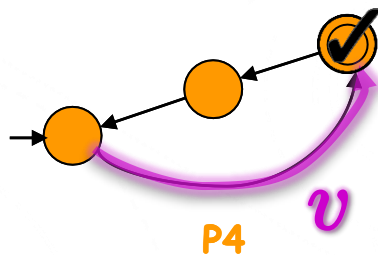
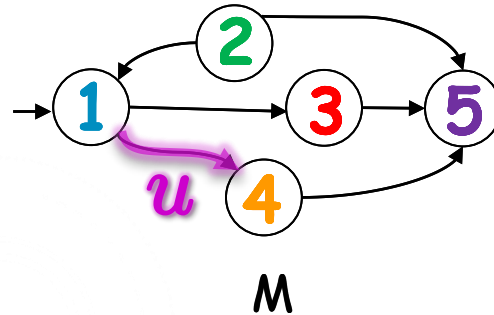
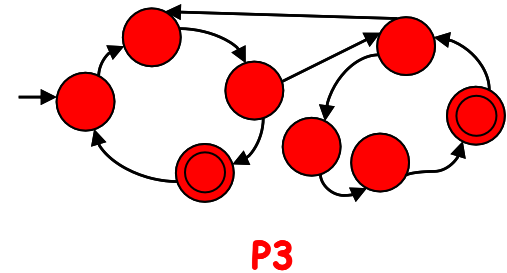
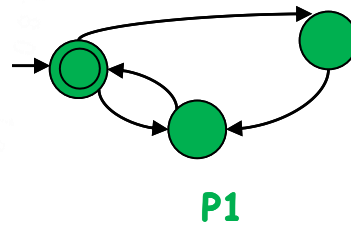
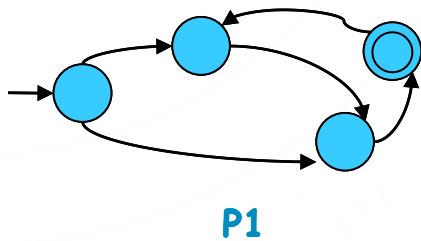
Progress



That restriction is removed

FDFA Acceptance

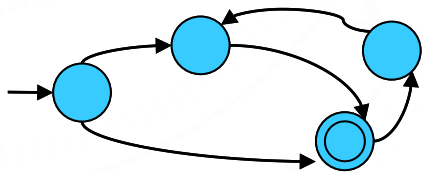
$$(u, v) \stackrel{?}{\in} \left[\overbrace{M, P_1, P_2, P_3, P_4, P_5}^{\mathcal{F}} \right]$$



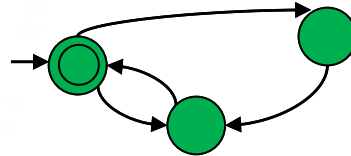
F DFA Normalized Acceptance

$$(u, v) \in \left[M, P_1, P_2, P_3, P_4, P_5 \right] \quad \mathcal{F}$$

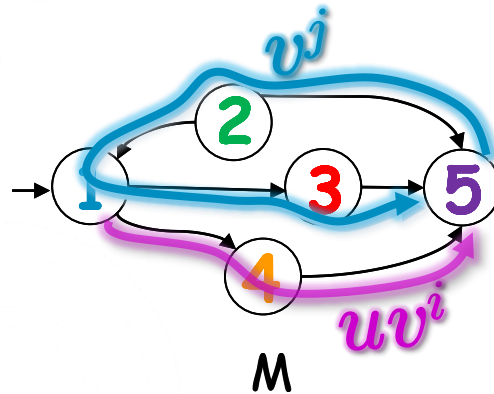
Normalization seeks for the **smallest repetition** of the **period** that **loops back**



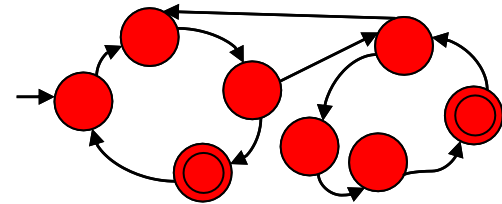
P1



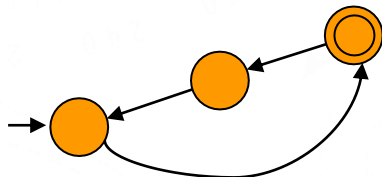
P1



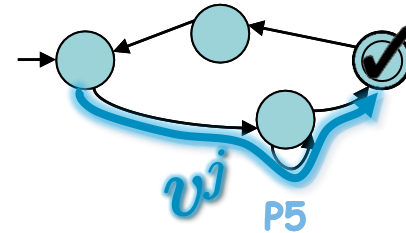
M



P3



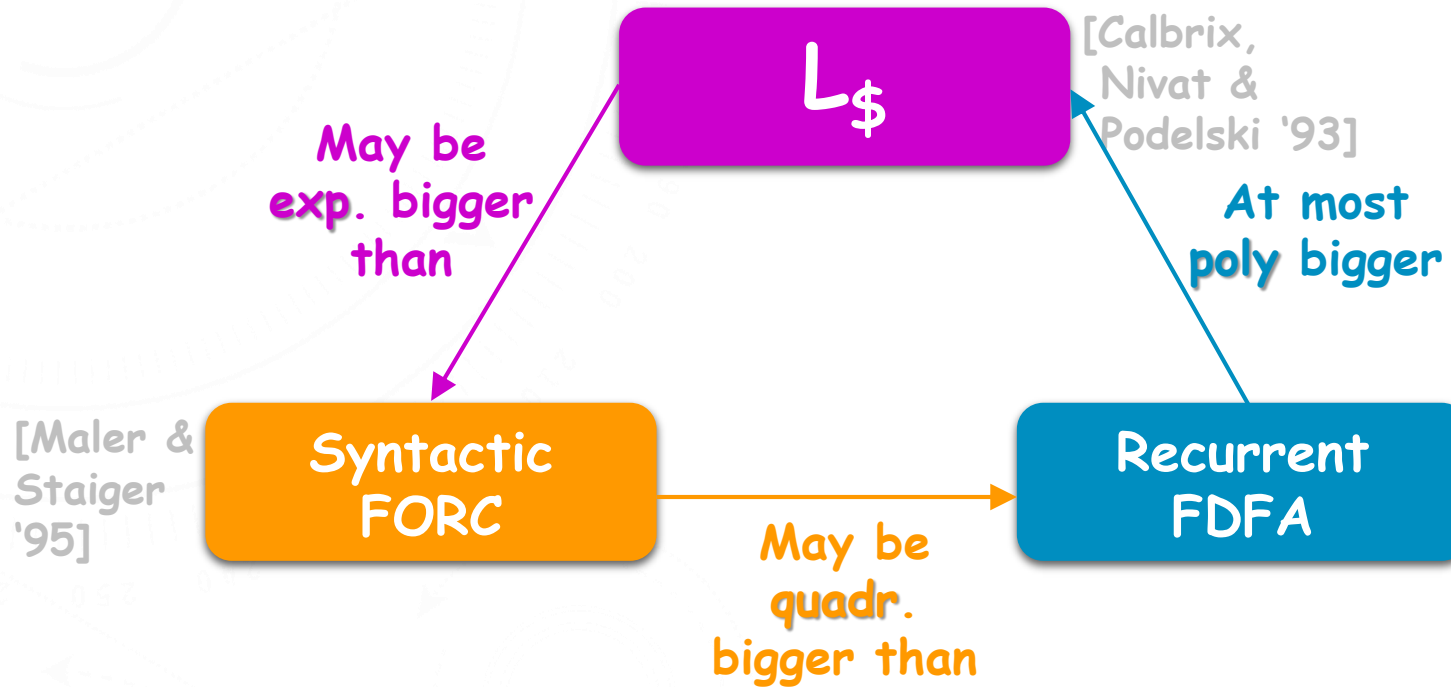
P4



P5

We term **Recurrent FDFA** the FDFA where progress DFA recognize only periods that **loop back**.

Results (1)

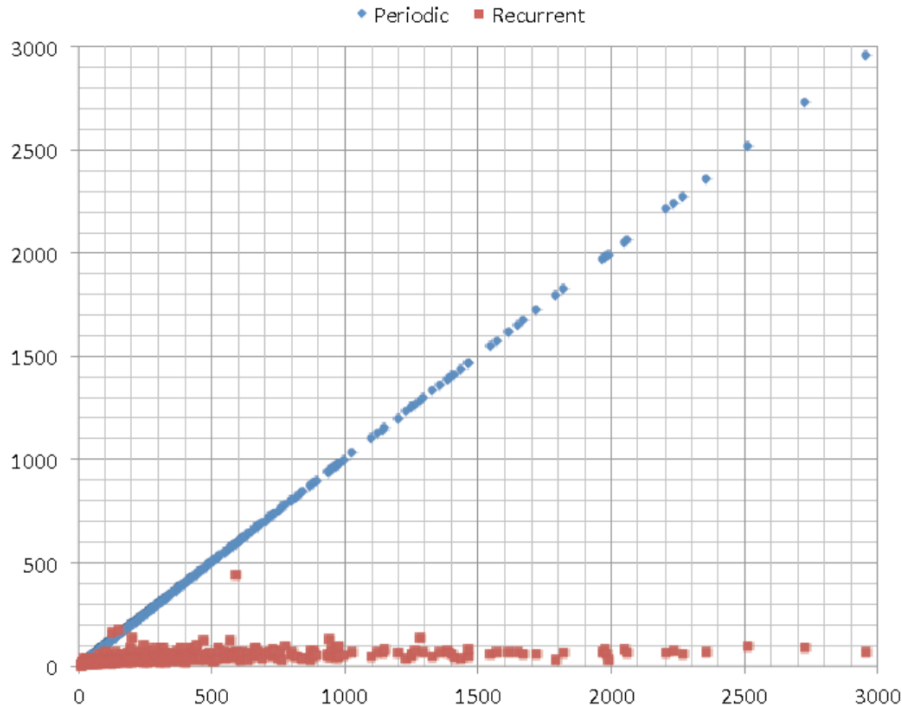
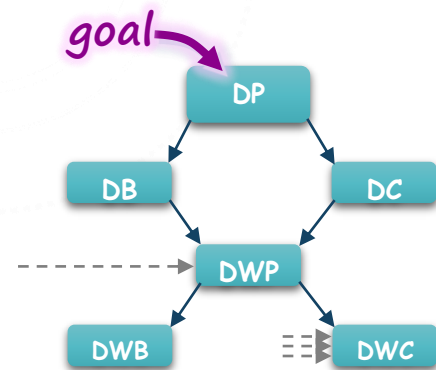


Results (2)

A learning algorithm L^w that learns the full class of regular w -languages using recurrent FDFAs

☹️ Worst-case time complexity polynomial in $L_\$$

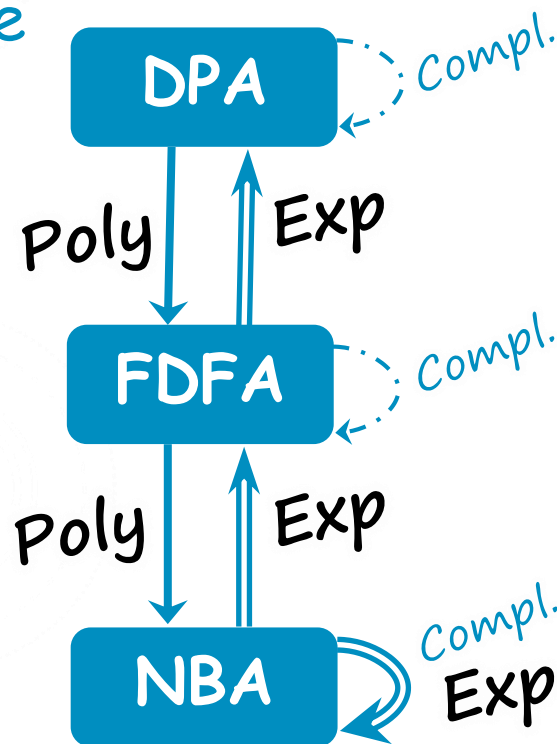
😊 Preforms very well on random targets



FDFAs as Acceptors of w -Langs

[ANGLUIN, BOKER & FISMAN FMCS'16]

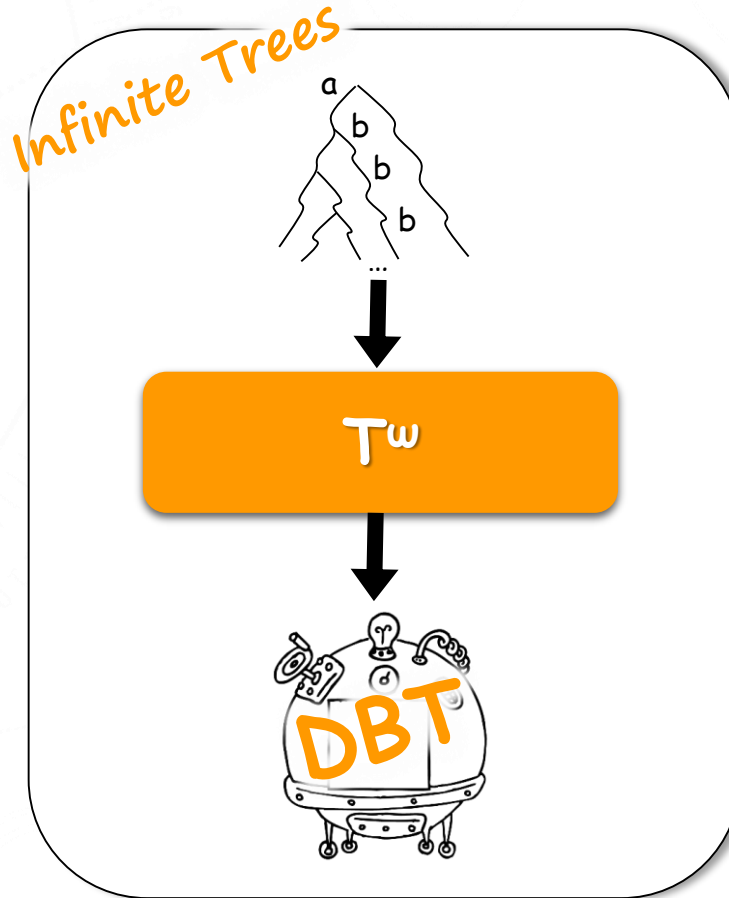
- Have a **Myhill-Nerode** characterization
- **Boolean operations** are in **LOGSPACE**
- **Decision problems** are in **NLOGSPACE**
- **Succinctness-wise**



Some open questions

- Polytime learning of a class of w -Langs more expressive than DWP
- Saturation of FDFA is in PSPACE; currently no lower bound
- Find smaller canonical representations

Further Directions



On going work with
DANA ANGLUIN &
TIMOS ANTONOPOULOS

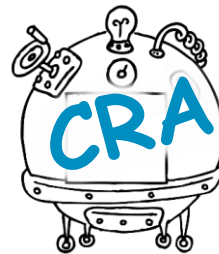
Further Directions

Regular Functions

```
{ f(abb) = 7,  
  f(bbb) = 80,  
  f(aaaaa) = 12  
}
```

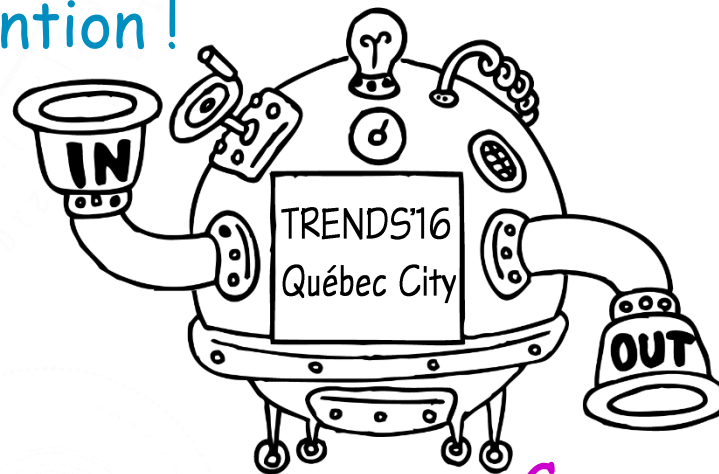


F*



On going work with
RAJEEV ALUR

Merci
beaucoup
pour votre
attention !



Commentaires /
Questions ?