

Regular Languages of Infinite Words

DANA Fisman
Ben-Gurion University

based on joint works with

Dana Angluin, Timos Antonopoulos, Udi Boker,
Nevin George, Yaara Shoval

Learning Langs of Finite Words

Active Learning

$w \in L?$
Yes / No

Passive Learning

w	aba	aab
	$aabb$	$abba$
	baa	$baba$

L is a language in a given class of languages \mathbb{L} .

Aut is an acceptor from a given class of acceptors \mathbb{A} representing \mathbb{L} .

Learning Langs of Infinite Words

Active Learning

$abaabbabababababab\bar{a}babab\dots$

$\in L?$

Yes / No

Passive Learning

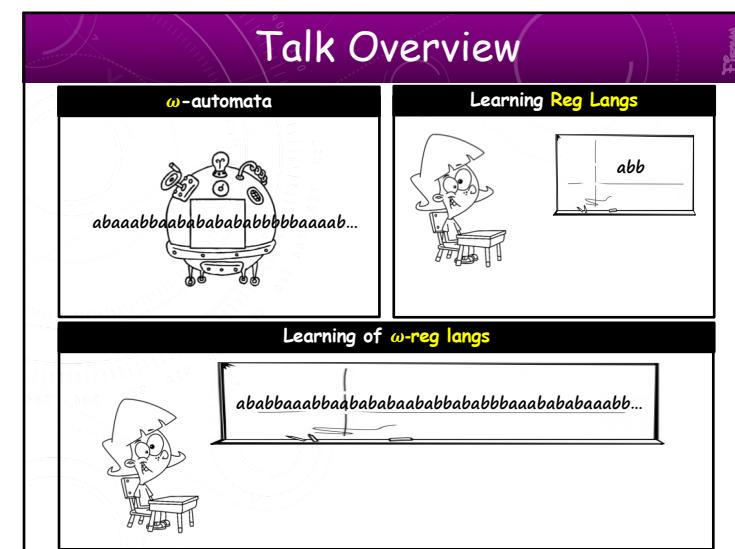
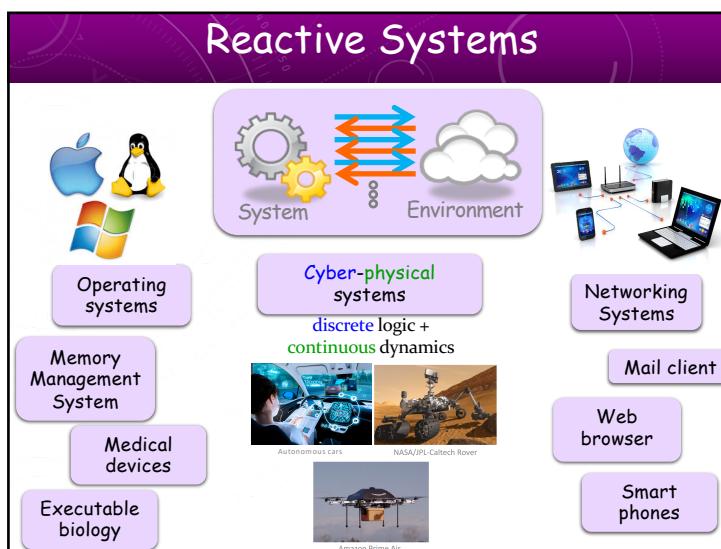
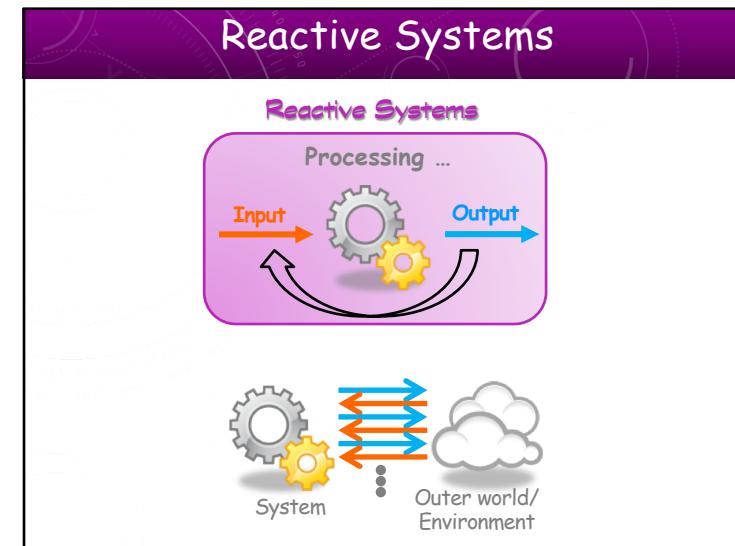
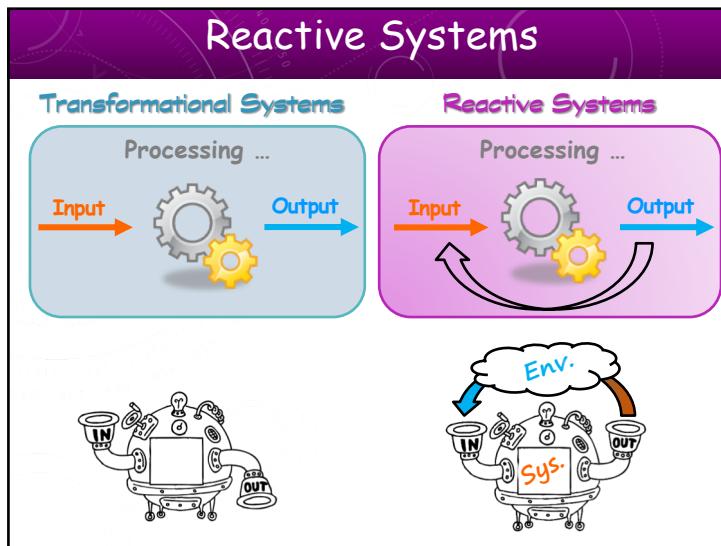
$abaaababbbbababababaa$

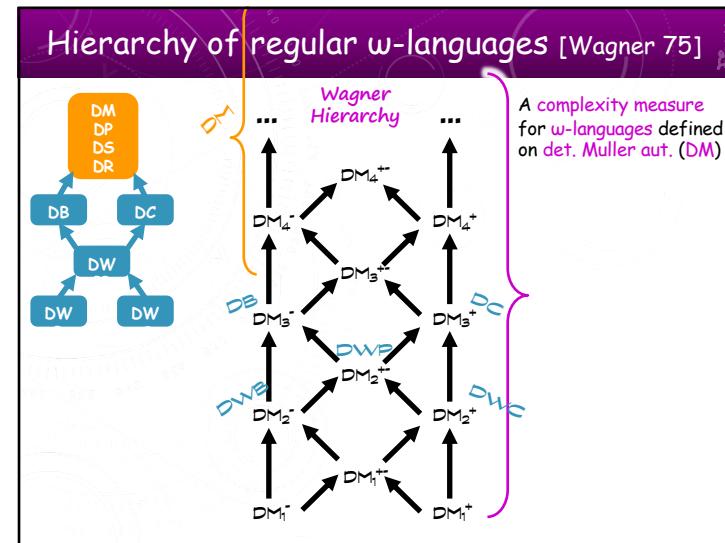
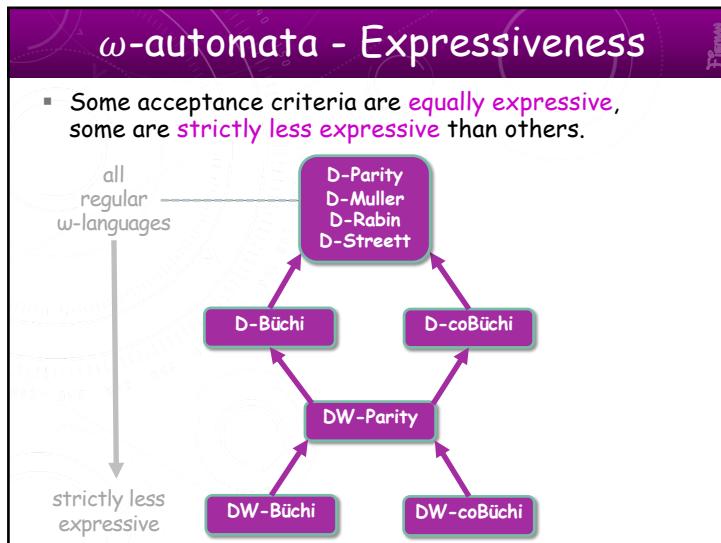
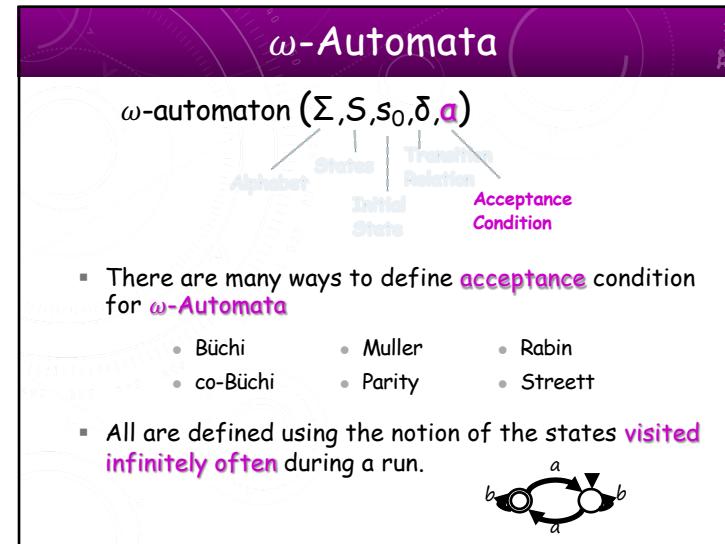
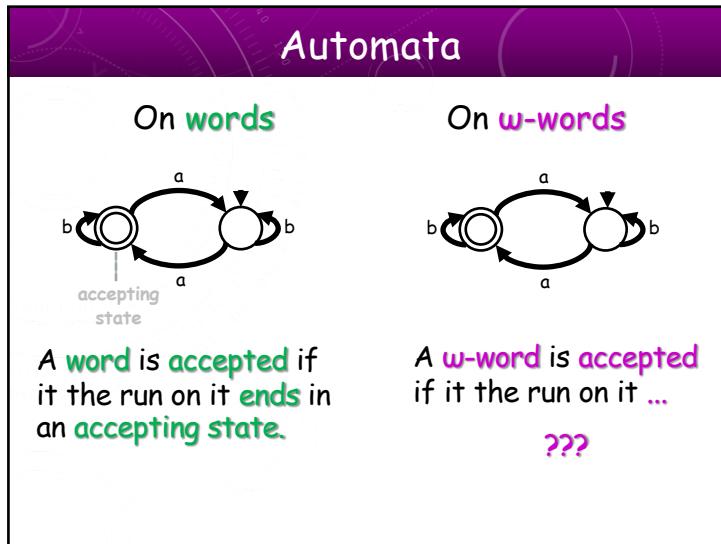
$abaaaabaaaabbabababab$

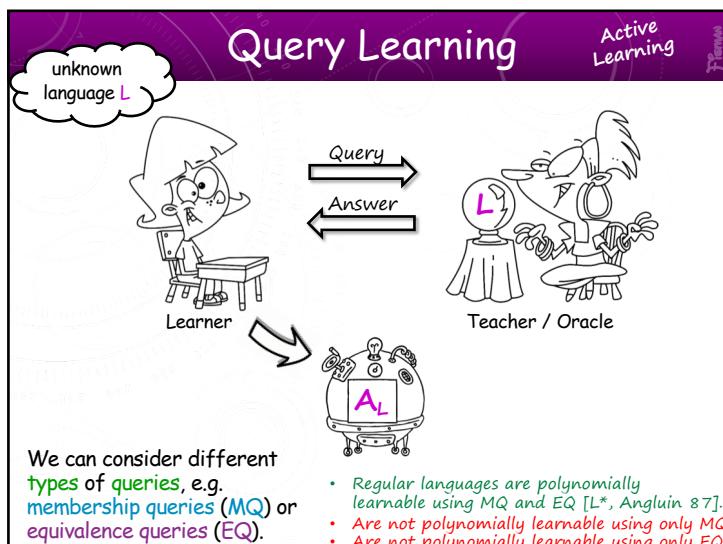
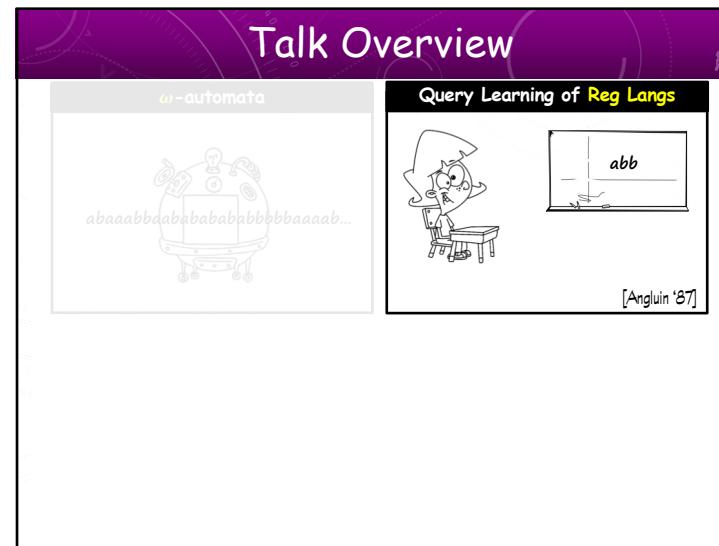
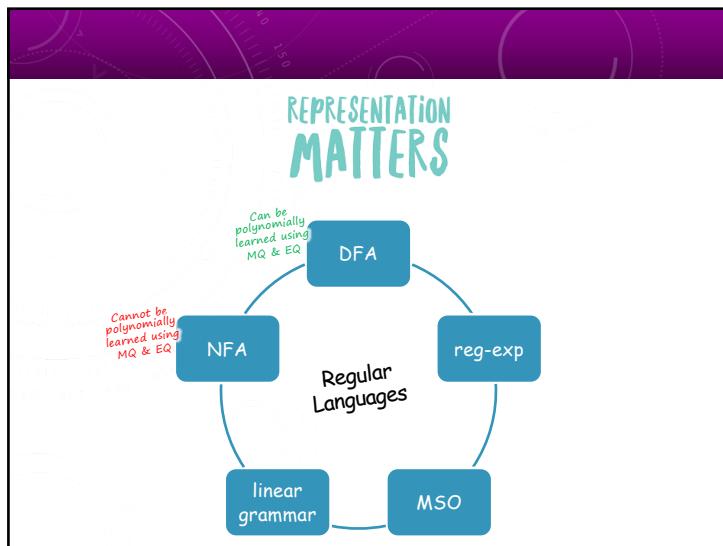
What are ω -automata?

How do we cope with infinite words (ω -words)?

Why infinite words?







The idea behind many learning alg.s.

- Start with the smallest automaton adhering to data
- Add states when you realize you must
- If $xv \in L$ and $yv \notin L$ then x and y must lead to different states
 - The word v is termed a distinguishing suffix.
- We say that $x \not\sim_L y$ (as defined next).
- The relation \sim_L plays a key role in learning languages !!!

The Syntactic Right Congruence

Let $x, y \in \Sigma^*$.

Let $L \subseteq \Sigma^*$:

$$x \sim_L y \text{ iff } \forall z \in \Sigma^*. xz \in L \Leftrightarrow yz \in L$$

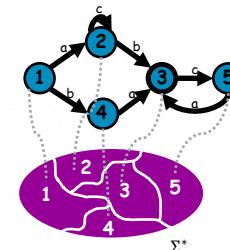
Let $L \subseteq \Sigma^\omega$:

$$x \sim_L y \text{ iff } \forall z \in \Sigma^\omega. xz \in L \Leftrightarrow yz \in L$$

Myhill-Nerode THM

For every regular language L the number of equivalence classes under \sim_L is finite.

The minimal DFA has one state for every equivalence class of L .



The matrix of a language

Let $\Sigma = \{a, b\}$, $L \subseteq \Sigma^*$:

Enumeration
of all strings

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

$bb \notin L$

$aab \in L$

Enumeration
of all strings

The matrix of a language

Let $\Sigma = \{a, b\}$, $L \subseteq \Sigma^*$:

Enumeration
of all strings

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

If $\epsilon \sim_L ab$ then their rows must be the same!

Enumeration
of all strings

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

How does L^* work?

- It maintains an **observation table**
- Where it tries to **find all different rows** of the infinite **table**.
- It asks **MQ** to **gradually extends** the table
- When the table meets certain criteria ("is **closed**") it is possible to **extract** from it an **automaton** and ask **EQ**.
- Smart processing of a **counterexample** to an **EQ** will reveal at **least one new row** to the table.
- If **n** is the **number of states** of the **smallest DFA** for **L**
- Then when the **table** is **closed** and has as **n** **distinct rows** the next **EQ** will be the **smallest DFA**.

Talk Overview

ω -automata

abaaaabbababababbbbbaaaab...

Learning Reg Langs

abb

Learning of ω -reg langs

ababbaaabbababababaababbababbbaaababaaaabb...

Coping with ω -words

Is **abcccccabdebbaaaabcdaa...** in **L**?

prefixes
ultimately periodic words

Learner

Teacher

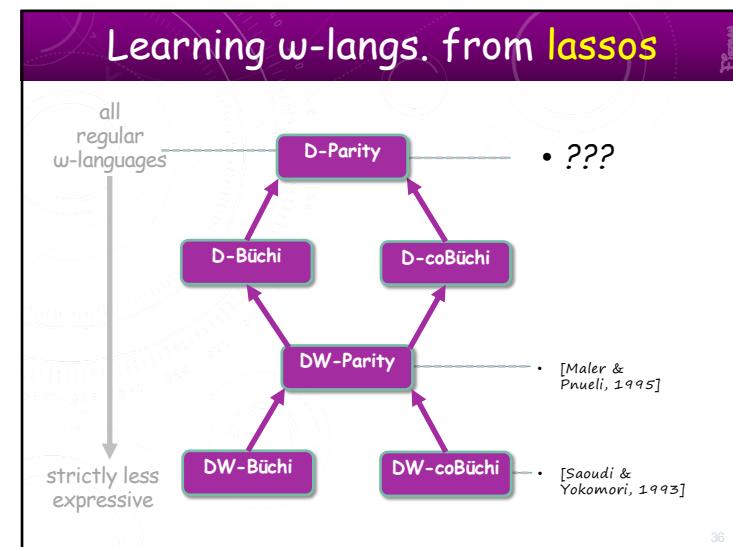
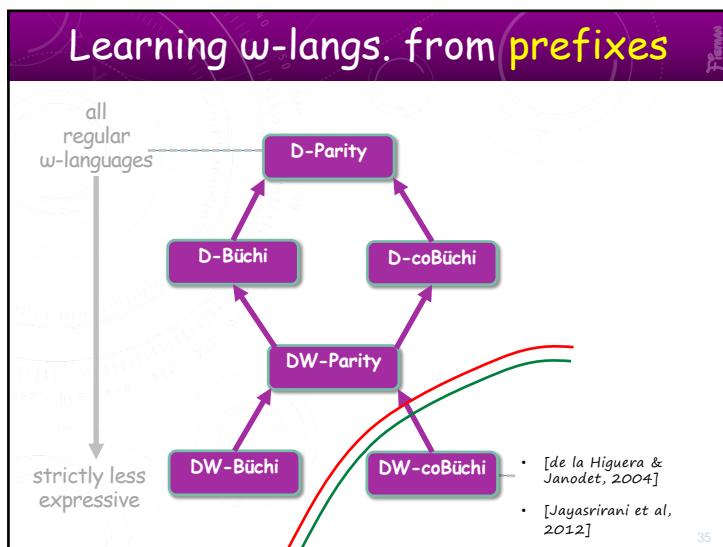
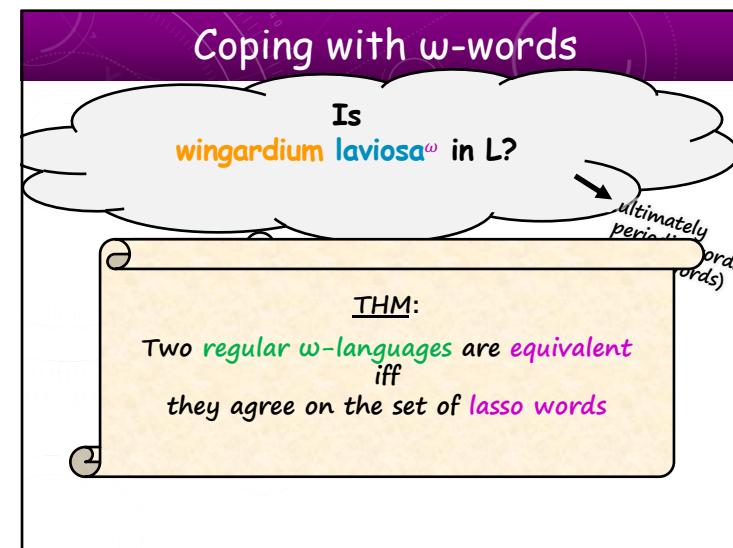
Coping with ω -words

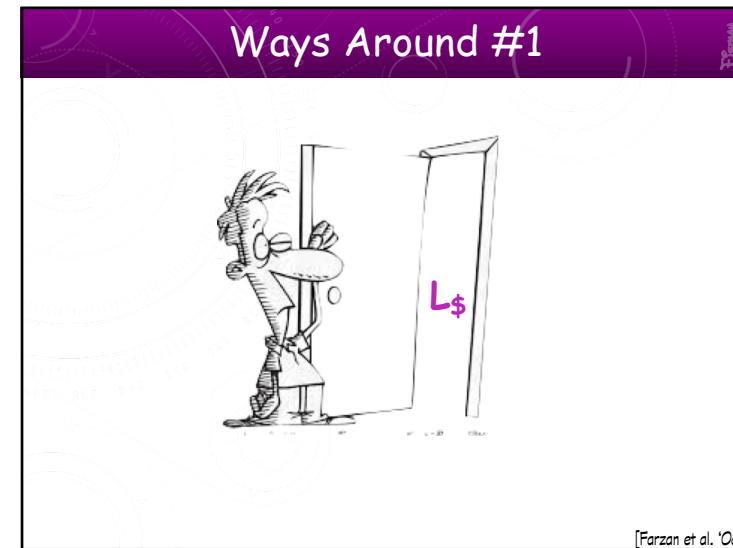
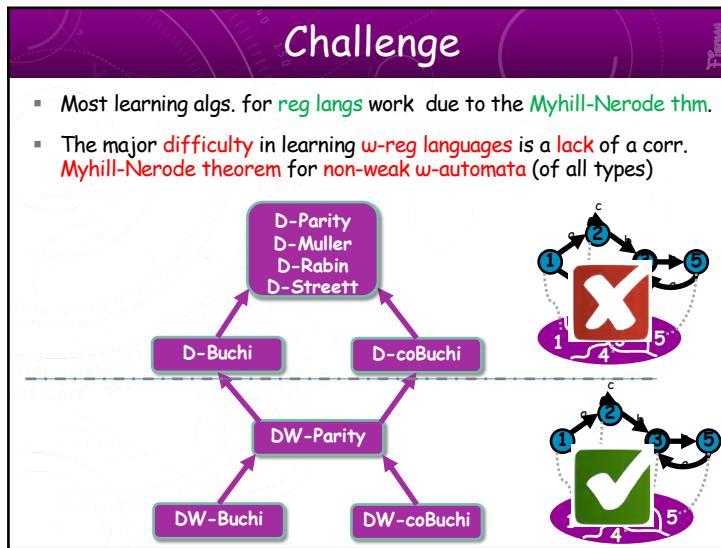
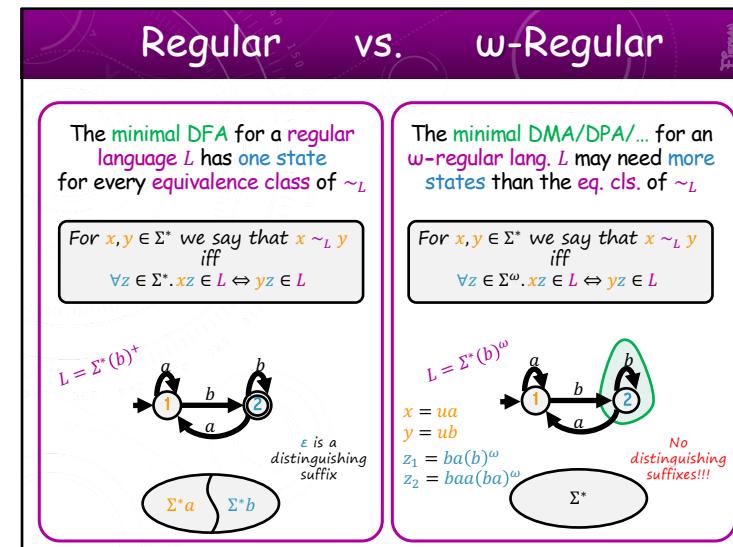
Is **wingardium laviosa ^{ω}** in **L**?

ultimately periodic words

Learner

Teacher





$L_{\$}$

- It follows from

THM:

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

that an ω -reg. lang. L can be represented by the reg. lang $L_{\$} = \{ u\$v \mid uv^\omega \in L \}$ [Calbrix, Nivat & Podolski 1993]

- [Farzan et al. '08] proposed to learn N-Büchi using L^* and a DFA for $L_{\$}$
- Unfortunately an N-Büchi with n states \Rightarrow DFA for $L_{\$}$ with $2^n + 3^{n^2}$ states [Kuperberg, Pinault & Pous 2019]

Ways Around #2

[Angluin & F., '14]

Finding a Myhill-Nerode Rep.

- Look for an alternative representation or reg- ω langs, that does have a Myhill-Nerode Theorem.
- In the quest for a representation for ω -langs for which a Myhill-Nerode THM exists [Maler and Staiger 1997] introduced

Family of Right Congruences (FORC)

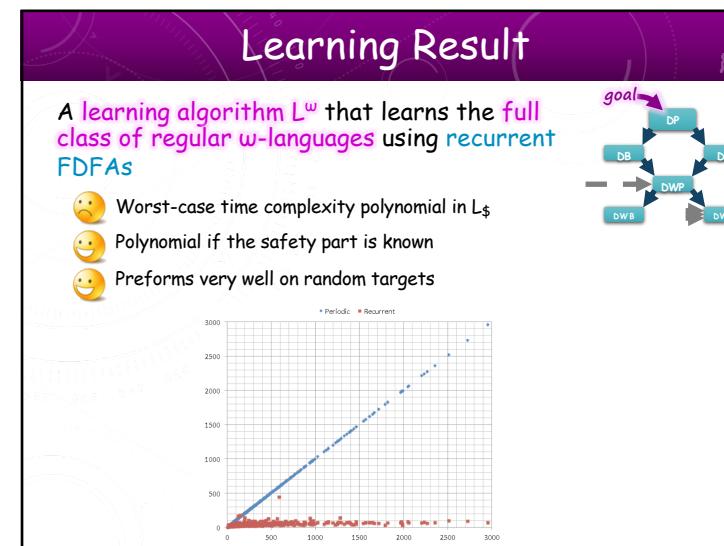
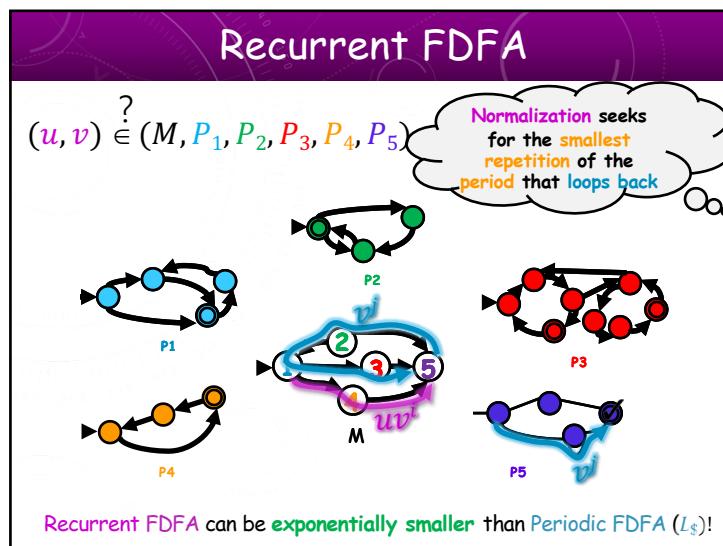
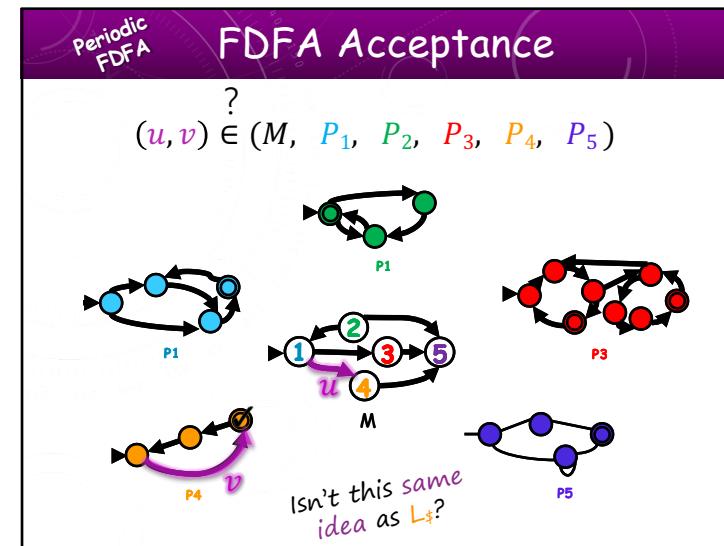
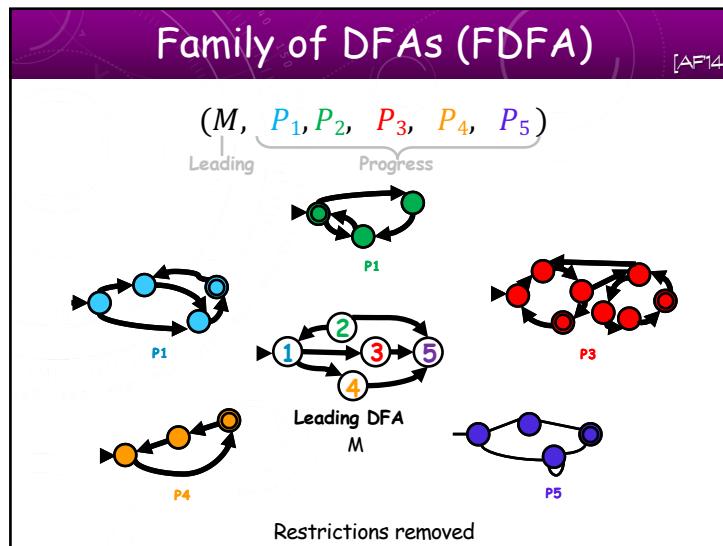
Family of Right Congruences [MS97]

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

Leading Progress

Leading Right Congruence \sim

Plus some restrictions

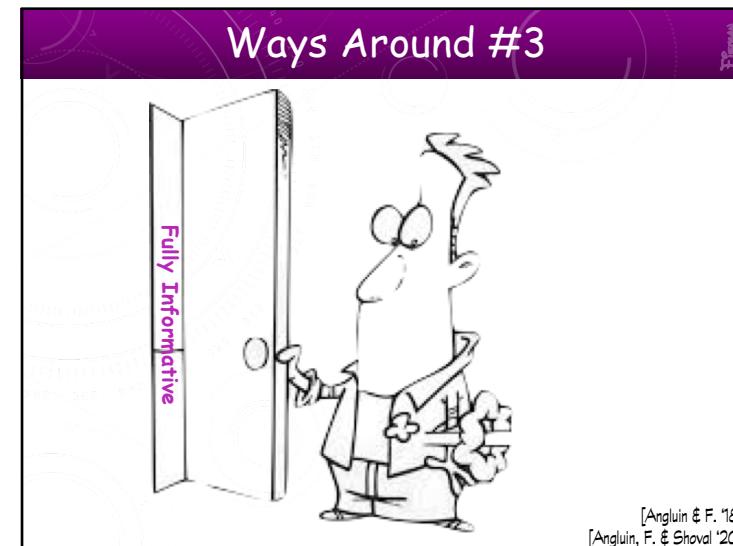


FDFAs as Acceptors of ω -Langs

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

DPA → FDFA → NBA

union, intersection, complementation
emptiness, inclusion, equivalence



Which Sub-Classes have a MN Thm?

What do we mean by this?

D-Parity
D-Parity
D-Muller
D-Rabin
D-Streett

D-Parity
 $L \quad \times L' \quad L'''$
 $\checkmark L''$

D-Büchi
DW-Parity
DW-Büchi
DW-coBüchi

D-coBüchi

Are there langs.
not in
DW-Parity
with a MN-Thm?

How informative is \sim_L ?

Σ^*

If \sim_L has one equivalence class we say it is zero-informative (or trivial)

If \exists det. aut with 1-1 correspondence to \sim_L we say it is fully-informative

Questions

- Are there other subsets of ω -regular langs. that have a fully informative right congruence?
- Can we fully characterize languages with a trivial right congruence?

Characterization Result

[AF18]

- A regular ω -language L has a trivial right congr. iff L is prefix-independent.
- If L has a trivial right congr., then L is a liveness property, but the other direction is not necessarily true.
- The class of languages with a trivial right congr. is closed under intersection, union and complementation.
- A regular ω -language has a trivial right congr. iff $L = \Sigma^*(R_1^\omega + R_2^\omega + \dots + R_k^\omega)$ for some regular languages R_1, R_2, \dots, R_k .

Notations

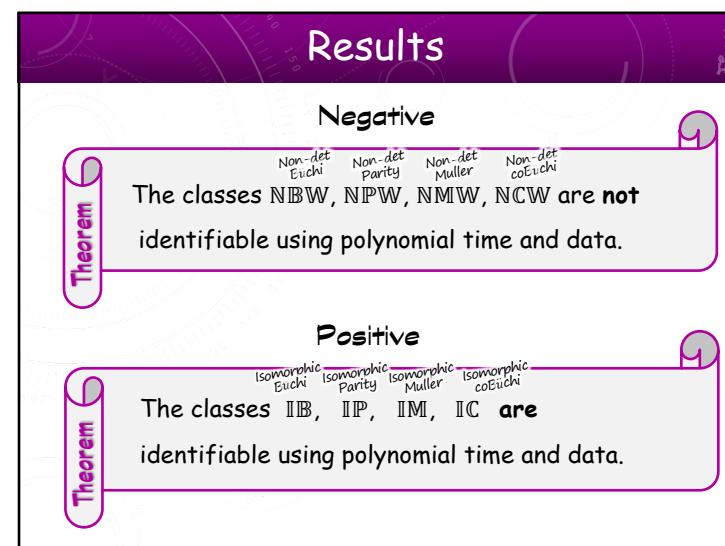
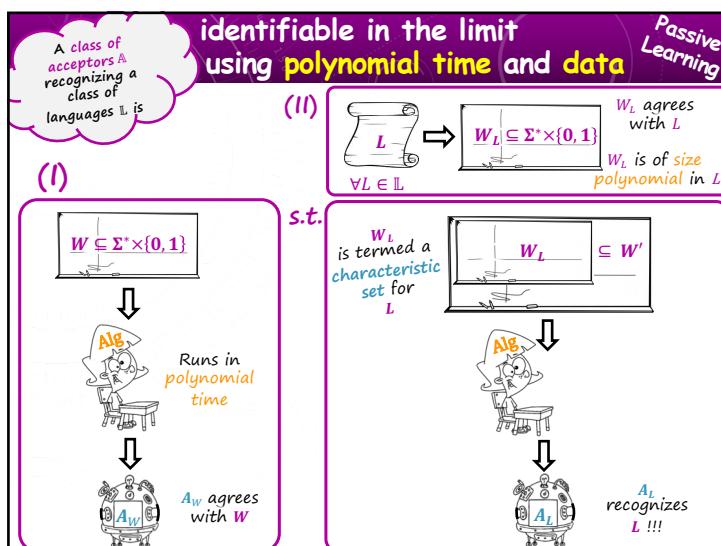
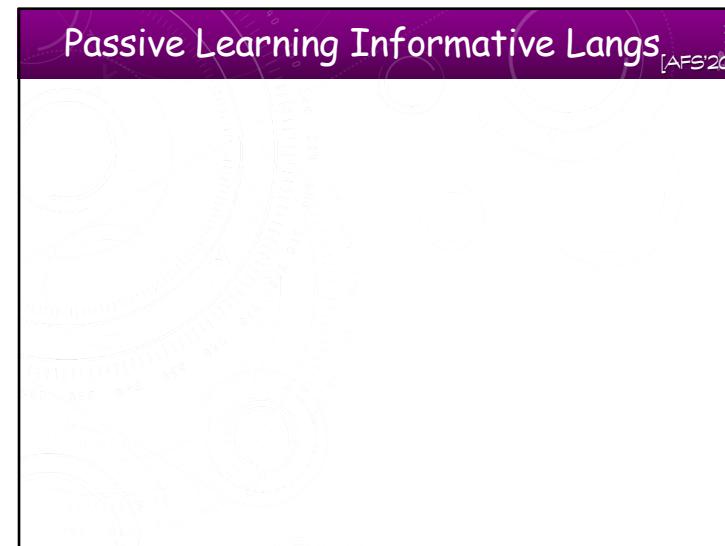
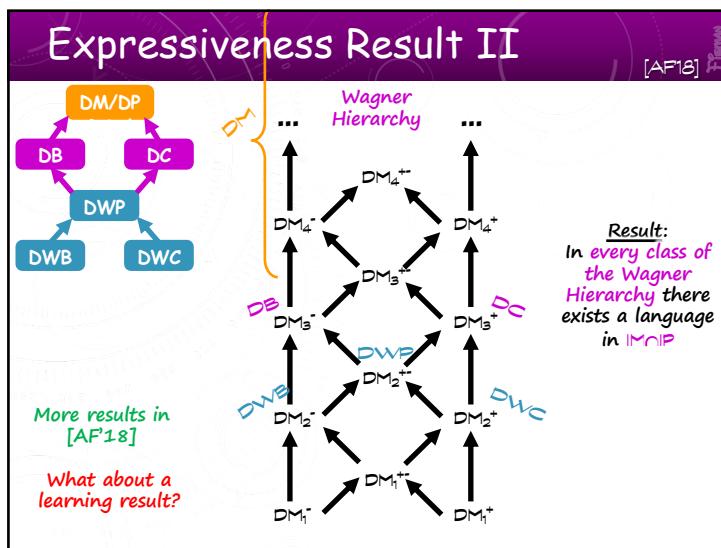
IB, IP, IM, IC

The "isomorphic classes"

- IB is the class of languages L where there is a 1-1 relation between \sim_L and the number of states in the minimal DBW (Det. Buchi Aut) for L
- IC is ... DCW (Det. coBuchi Aut)
- IP is ... DPW (Det. Parity Aut)
- IM is ... DMW (Det. Muller Aut)

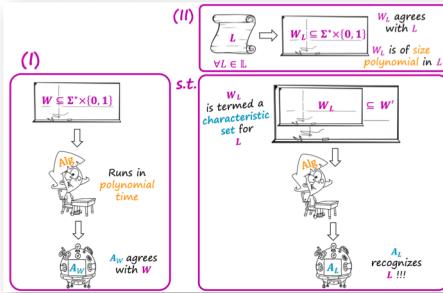
Expressiveness Result I

[AF18]



Polynomial Identification of $\text{IB}, \dots, \text{IM}$

Goal:

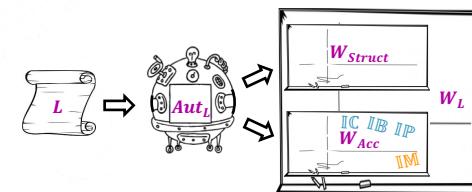


Constructing the characteristic set

Goal:



Scheme:



The Canonical Forest

Theorem

Let $P = (\Sigma, Q, q_0, \delta, \kappa)$ be a DPA.

There exists a **canonical forest** $F^*(P)$ that is

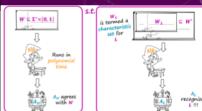
unique up to isomorphism

and has a lot of good properties

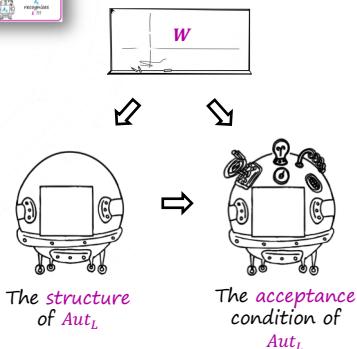
Can be computed
in polynomial time

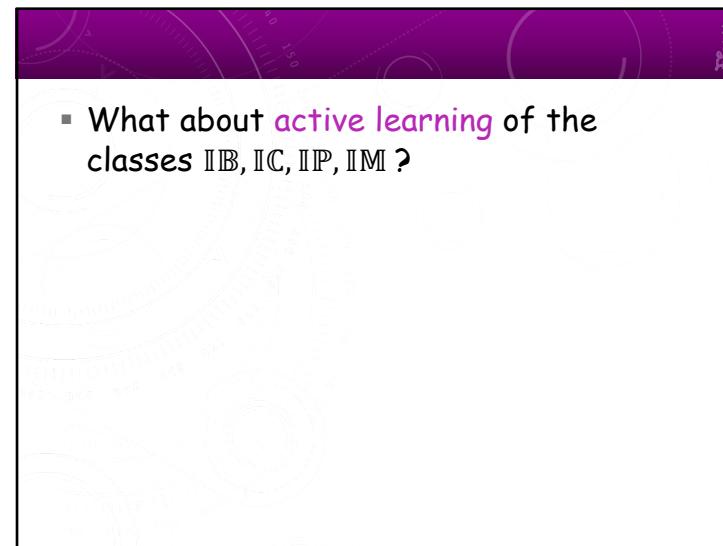
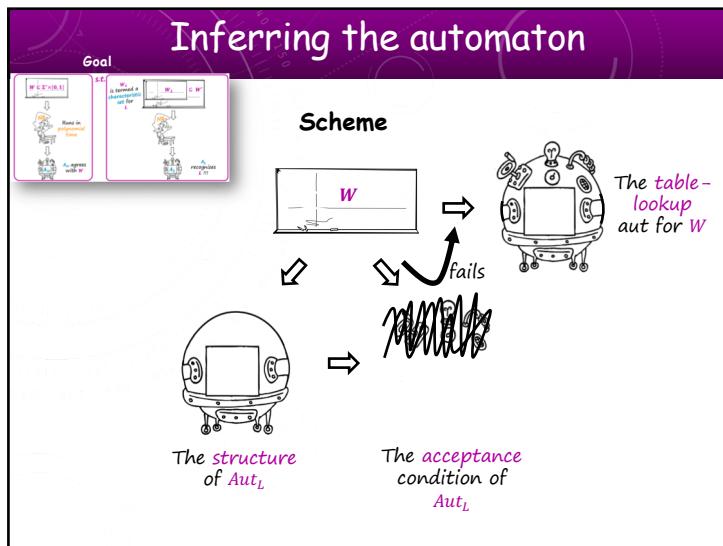
Inferring the automaton

Goal



Scheme





Active Learning Informative Langs

- As hard as learning the deterministic counterparts [Bohn & Löding '21]

Theorem

There is a polynomial time active learning algorithm for IX iff there is a polynomial time active learning algorithm for DX

for $X \in \{\text{B}, \text{C}, \text{P}, \text{M}\}$



What about Non-Det Models?

- Q: Why do we consider only deterministic models ?
- A: It is not hard to show that NBW are not polynomially predictable with MQs and EQs

THM:

Under standard cryptographic assumptions NBAs cannot be polynomially learned using MQs and EQs

- Q: Okay, so what about unambiguous Büchi Automata ?

Unambiguous Automata

DFA	NFA
	
Det. Finite Aut. A single run on each word.	NonDet. Finite Aut. Various runs on each word.

UFA

Unambiguous Finite Aut. Various runs on each word. At most one is accepting

Strongly Unambiguous Büchi Automata

[CM03, BL10]

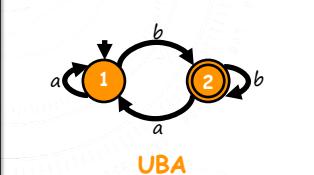
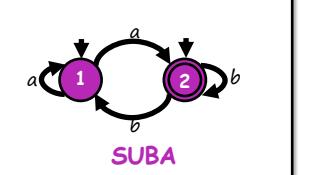
DBA	NBA
	
Det. Büchi A single run on each word.	NonDet. Büchi Various runs on each word.

UBA	SUBA
	
Unambiguous Büchi Various runs on each word. At most one is initial & final, accepting	Strongly Unambiguous Büchi Various runs on each word. At most one is initial & final.

Strongly Unambiguous Büchi Automata

[CM03, BL10]

infinitely many b 's

	
UBA	SUBA

Why SUBAs?

- Every ω -regular language can be accepted by a SUBA
- Equivalence and containment of SUBAs can be decided in polynomial time.
- SUBAs are backward deterministic
- Intuitively, one starts reading the word from somewhere in the loop and goes backward deterministically
- SUBAs may be exponentially smaller than DBAs and may be exponentially bigger than DMAs

SUBA - Succinctness

DBA $\xrightarrow{\text{exp}}$ SUBA $\xrightarrow{\text{exp}}$ DMA [CM03, BL10]

[AAF'20]

```

graph TD
    DBA -- exp --> SUBA
    SUBA -- exp --> DMA
    DBA -- exp --> FDFA
    DBA -- exp --> DWP
    SUBA -- exp --> DWP
    SUBA -- exp --> L$
    DMA -- exp --> L$ 
    DMA -- exp --> L$^R
    FDFA -- exp --> L$ 
    DWP -- exp --> L$ 
    DWP -- exp --> L$^R
  
```

SUBA - Learning Results

The class of ω -regular languages is polynomially predictable with membership queries when the target language is represented by a SUBA

There is a polynomial time algorithm that learns all regular ω -languages, for target langs. represented by a SUBA using MQ and non-proper EQ

via Mod-² Multiplicity Automata

Multiplicity Automata (MA)

- Multiplicity Automata (MA) are an algebraic generalization of automata, defined wrt. a field \mathbb{F} .
- They associate with every word a value from \mathbb{F} .

$f(w) = 3|w|_a + 1$

($\Sigma, v_I, \{\mu_\sigma\}_{\sigma \in \Sigma}, v_F$)

Alphabet Initial Vector Transition Matrices Final Vector

$w = aaba$

$$\begin{bmatrix} : \\ v_I \\ : \end{bmatrix} \xrightarrow{a} \begin{bmatrix} : \\ \mu_a \\ : \end{bmatrix} \xrightarrow{a} \begin{bmatrix} : \\ \mu_a \\ : \end{bmatrix} \xrightarrow{b} \begin{bmatrix} : \\ \mu_b \\ : \end{bmatrix} \xrightarrow{a} \begin{bmatrix} : \\ \mu_a \\ : \end{bmatrix} \xrightarrow{} \begin{bmatrix} : \\ v_F \\ : \end{bmatrix}^T$$

$$\begin{bmatrix} : \\ v_0 \\ : \end{bmatrix} \quad \begin{bmatrix} : \\ v_1 \\ : \end{bmatrix} \quad \begin{bmatrix} : \\ v_2 \\ : \end{bmatrix} \quad \begin{bmatrix} : \\ v_3 \\ : \end{bmatrix} \quad \begin{bmatrix} : \\ v_4 \\ : \end{bmatrix} \quad c$$

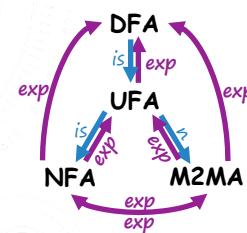
Mod-2 MA (M2MA)

- MA are **Polynomially learnable** with **MQ** and **EQ** [Beimel et al.'00]
- How can we use MA to represent languages?
- M2MA are MA over the field $GF(2) = \{0,1\}$ where sum and product are computed modulo 2.
- We can treat **words associated with 1** as **accepted**, and **words associated with 0** as **rejected**.

Advantages of M2MA

- M2MA are **Polynomially learnable** with **MQ** and **EQ**.
- A **deterministic model**.
- Complementation** incurs only an **additive blowup**.
- Union** and **intersection** incurs only a **quadratic blowup**.
- Equivalence** can be decided in **polynomial time**.

Size Relations



Worst-case size bounds

	DFA	UFA	M2MA	NFA	(G)SUBA	NBA
LTL	$2^{2^{O(n)}}$ via UFA	$2^{O(n)}$ [AFFG]	\leftarrow	\leftarrow	$(2^n, n)$ [BK'08]	$n2^n$ [BK'08]
	\rightarrow	\rightarrow	$2^{\Omega(n)}$ [AFFG]	$2^{\Omega(n)}$ [AFFG]	\rightarrow	$2^{\Omega(n)}$ [BK'08]
DBA	$n + n3^{n^2}$ [KPP'19]	\leftarrow	\leftarrow	\leftarrow	\downarrow	—
	$2^{\Omega(n \log n)}$ [AF'16]	\rightarrow	$2^{\Omega(n)}$ [AFFG]	$2^{\Omega(n)}$ [AFFG]	$2^{\Omega(n)}$ [BL'10]	—
NBA	$2^n + 2^n 3^{n^2}$ [KPP'19]	\leftarrow	\leftarrow	$n + n3^{n^2}$ [KPP'19]	$(12n)^n$ [CM16]	—
	\uparrow	\uparrow	\uparrow	\uparrow	\uparrow	—
SUBA	\uparrow	$n + 2n^2$ [BL'10]	\leftarrow	\leftarrow	—	—
	$2^{\Omega(n)}$ [AAF'20]	\rightarrow	$\Omega(n^2)$ [AFFG]	$\Omega(n^2)$ [AFFG]	—	—

Upper Bound
Lower Bound

Worst-case size bounds						
	DFA	UFA	M2MA	NFA	(G)SUBA	NBA
LTL	$2^{2^{O(n)}}$ via UFA \rightarrow	$2^{O(n)}$ [AFFG] \rightarrow	\leftarrow	\leftarrow	$(2^n, n)$ [AFFG]	$n2^n$ [BK'08] $2^{\Omega(n)}$ [BK'08]
			$2^{\Omega(n)}$ [AFFG]	$2^{\Omega(n)}$ [AFFG]	\rightarrow	
DBA	$n + n3^{n^2}$ [KPP'19]	\leftarrow	\leftarrow	\leftarrow	\downarrow	—
	$2^{\Omega(n \log n)}$ [AF'16]	\rightarrow	$2^{\Omega(n)}$ [AFFG]	$2^{\Omega(n)}$ [AFFG]	$2^{\Omega(n)}$ [BL'10]	—
NBA	$2^n + 2^n 3^{n^2}$ [KPP'19]	\leftarrow	\leftarrow	$n + n3^{n^2}$ [KPP'19]	$(12n)^n$ [CM16]	—
	\uparrow	\uparrow	\uparrow	\uparrow	\uparrow	—
SUBA	\uparrow $2^{\Omega(n)}$ [AAF'20]	$n + 2n^2$ [BL'10] \rightarrow	\leftarrow $\Omega(n^2)$ [AFFG]	\leftarrow $\Omega(n^2)$ [AFFG]	—	—

Recap

ω -regular languages can be learned using MQ and EQ via FDFA_s

- If the safety part of the languages is known, then the liveness part can be polynomially learned
- There are arbitrary complex ω -regular language with a fully informative right congruence
- These languages can be learned in the limit using polynomial time and data
- Actively learning them is not easier than their non-fully informative counterparts
- ω -regular languages can be polynomially learned using MQ and EQ via SUBA_s/M2MA_s

