

On the use of guards for logics with data



Gabriele Puppis

LaBRI / CNRS

based on joint works with

Thomas Colcombet
and Clemens Ley

What is this talk about?

Data languages: sets of finite words over an **infinite alphabet** **invariant under permutations** of the letters.

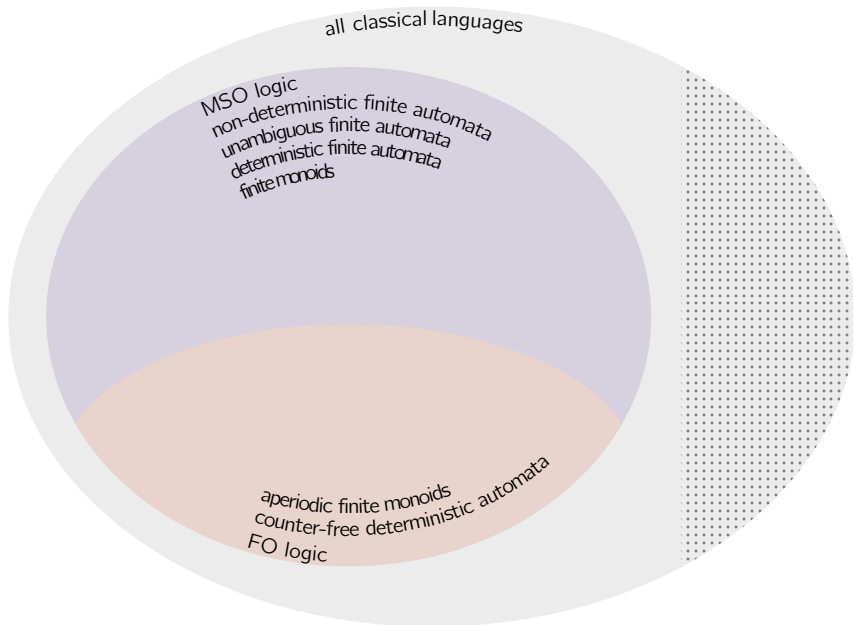
Lifting of some classical results to data languages, e.g.

- translations between monoids and logic
- characterization of first-order definability

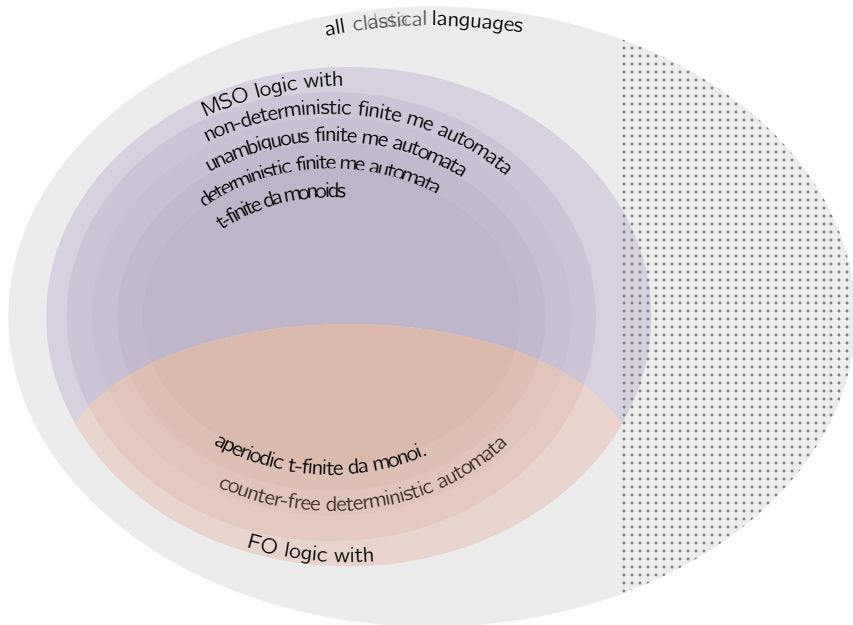
Example

$$\begin{aligned} L &= \{w \in \mathbb{N}^* : \text{at most 2 distinct values in } w\} \\ &= \{\varepsilon, \text{red}, \text{green}, \text{red red green}, \text{grey grey purple}, \text{green purple green green}, \dots\} \\ &= \{\varepsilon, \text{white}, \text{white white white}, \text{white white white white}, \dots\} \end{aligned}$$

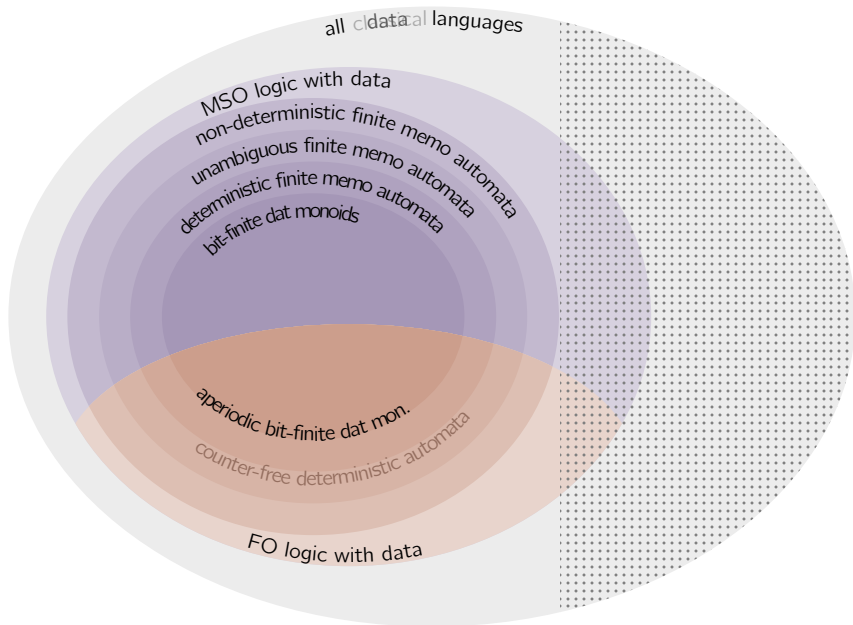
Logics, automata, and monoids over finite alphabets:



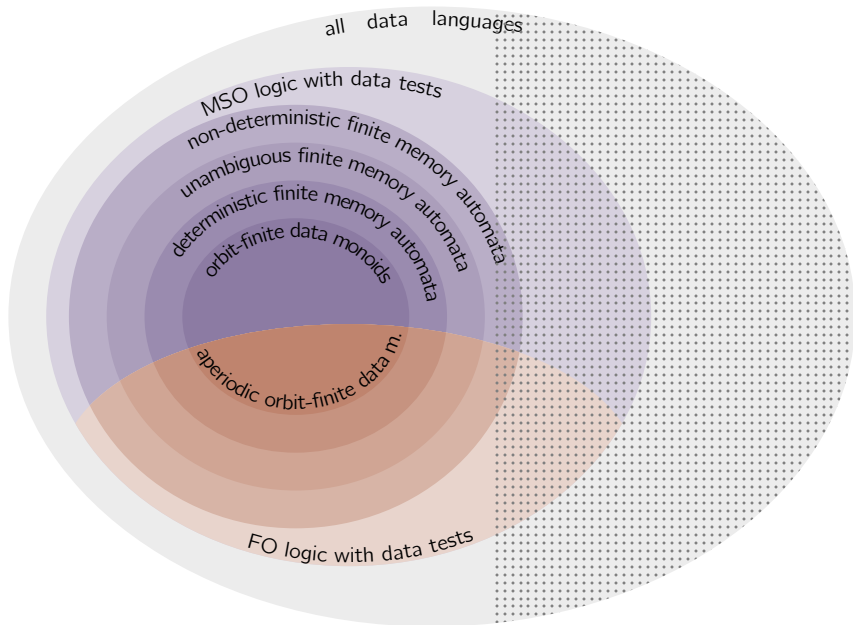
Logics, automata, and monoids over infinite alphabets:



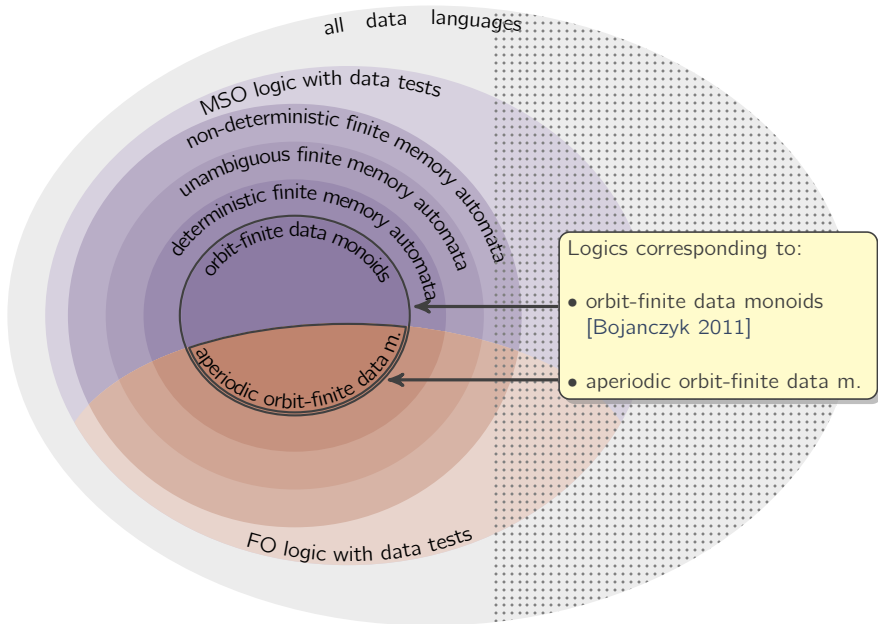
Logics, automata, and monoids over infinite alphabets:



Logics, automata, and monoids over infinite alphabets:



Logics, automata, and monoids over infinite alphabets:



Finite alphabet:

- **MSO logic**

$$\begin{aligned} \exists X. \text{ first} \in X \wedge \text{ last} \notin X \\ \wedge \forall y. y \in X \leftrightarrow y + 1 \notin X \end{aligned}$$

Infinite alphabet:

- **MSO logic with data tests**

$$\begin{aligned} \exists X. \text{ first} \in X \wedge \text{ last} \notin X \\ \wedge \forall y. y \in X \leftrightarrow y + 1 \notin X \\ \wedge \forall y, z \in X \rightarrow y \sim z \end{aligned}$$

Finite alphabet:

- **MSO logic**

$$\begin{aligned} \exists X. \text{ first} \in X \wedge \text{ last} \notin X \\ \wedge \forall y. y \in X \leftrightarrow y + 1 \notin X \end{aligned}$$

- **Finite Automata**



Infinite alphabet:

- **MSO logic with data tests**

$$\begin{aligned} \exists X. \text{ first} \in X \wedge \text{ last} \notin X \\ \wedge \forall y. y \in X \leftrightarrow y + 1 \notin X \\ \wedge \forall y, z \in X \rightarrow y \sim z \end{aligned}$$

- **Finite Memory Automata**



Finite alphabet:

- MSO logic

$$\begin{aligned} \exists X. \text{ first} \in X \wedge \text{ last} \notin X \\ \wedge \forall y. y \in X \leftrightarrow y + 1 \notin X \end{aligned}$$

- Finite Automata



- Finite monoids

$$s \cdot t = s$$

$$t \cdot s = t$$

Infinite alphabet:

- MSO logic with data tests

$$\begin{aligned} \exists X. \text{ first} \in X \wedge \text{ last} \notin X \\ \wedge \forall y. y \in X \leftrightarrow y + 1 \notin X \\ \wedge \forall y, z \in X \rightarrow y \sim z \end{aligned}$$

- Finite Memory Automata

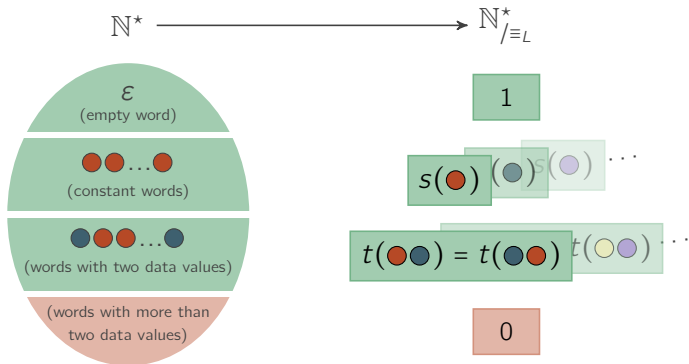


- Orbit-finite Data Monoids

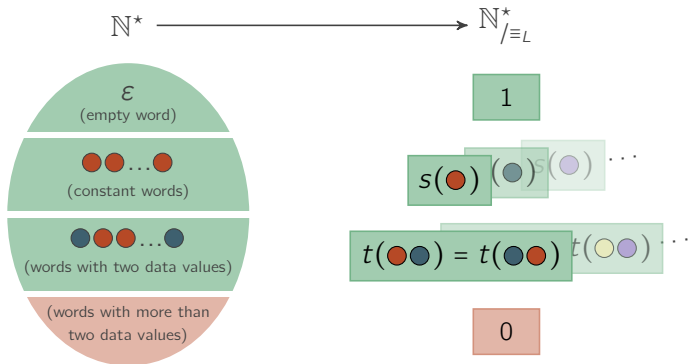
$$s(\text{red } \bullet \text{ green } \bullet) \cdot t(\text{grey } \bullet) = s(\text{red } \bullet \text{ grey } \bullet)$$

$$t(\text{grey } \bullet) \cdot s(\text{red } \bullet \text{ green } \bullet) = t(\text{grey } \bullet)$$

Consider the Myhill-Nerode equivalence \equiv_L for the data language $L = \{w \in \mathbb{N}^* : \text{at most 2 distinct values in } w\}$



Consider the Myhill-Nerode equivalence \equiv_L for the data language $L = \{w \in \mathbb{N}^* : \text{at most 2 distinct values in } w\}$



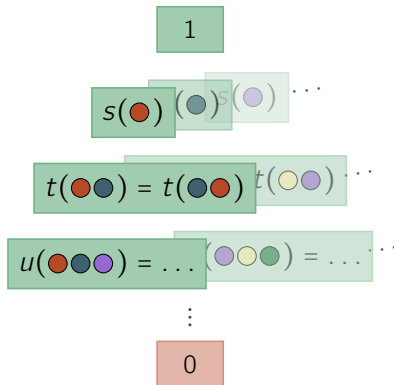
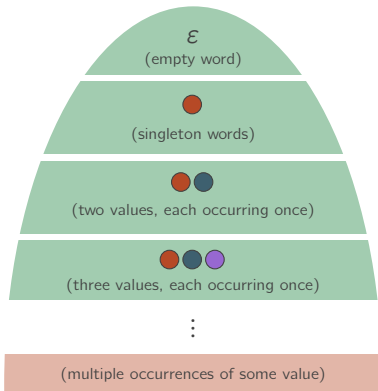
This gives a data monoid with products defined by equations like

$$s(\bullet) \cdot s(\bullet) = t(\bullet\bullet) \quad t(\bullet\bullet) \cdot s(\bullet) = 0 \quad \dots$$

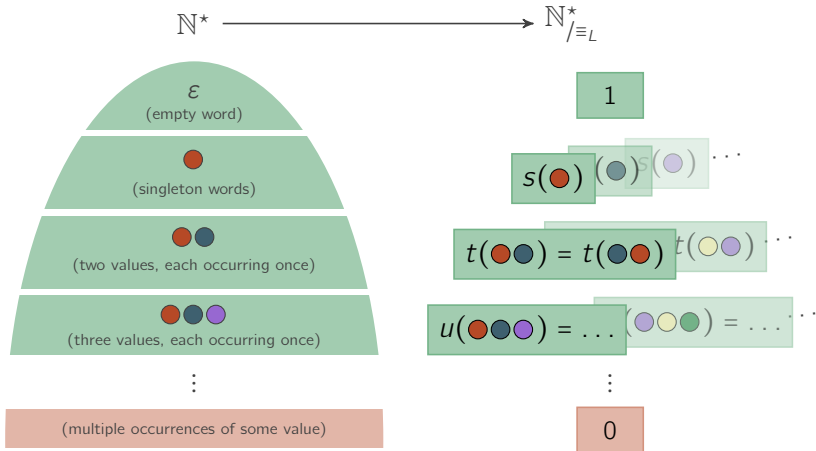
👉 The data monoid is **orbit-finite** because it has **finitely many elements up to renamings**

Consider the Myhill-Nerode equivalence \equiv_L for the data language $L = \{w \in \mathbb{N}^* : w \text{ contains at most one occurrence of each value}\}$

$$\mathbb{N}^* \longrightarrow \mathbb{N}^*/\equiv_L$$



Consider the Myhill-Nerode equivalence \equiv_L for the data language $L = \{w \in \mathbb{N}^* : w \text{ contains at most one occurrence of each value}\}$



👉 This syntactic data monoid is **not orbit-finite** (unbounded “memory” \rightarrow infinitely many orbits!)

Languages recognized by orbit-finite data monoids:

✓ Closed under all boolean operations

✓ *First value = last value*



✓ *At least two distinct values*



✓ *At least three distinct values*



✗ *First value reappears later*

✗ *Some value appears twice*

✗ *Every value appears at most once*

Languages recognized by orbit-finite data monoids:

✓ Closed under all boolean operations

✓ *First value = last value*



$$\exists x, y. x \sim y \wedge x = \text{first} \wedge y = \text{last}$$

✓ *At least two distinct values*



✓ *At least three distinct values*



✗ *First value reappears later*

✗ *Some value appears twice*

✗ *Every value appears at most once*

Languages recognized by orbit-finite data monoids:

✓ Closed under all boolean operations

✓ *First value = last value*



$$\exists x, y. x \sim y \wedge x = \text{first} \wedge y = \text{last}$$

✓ *At least two distinct values*



$$\exists x, y. x \not\sim y \wedge y = x+1$$

✓ *At least three distinct values*



✗ *First value reappears later*

✗ *Some value appears twice*

✗ *Every value appears at most once*

Languages recognized by orbit-finite data monoids:

✓ Closed under all boolean operations

✓ *First value = last value*



$$\exists x, y. x \sim y \wedge x = \text{first} \wedge y = \text{last}$$

✓ *At least two distinct values*



$$\exists x, y. x \neq y \wedge y = x+1$$

✓ *At least three distinct values*



$$\exists x, y. x \neq y \wedge (x \neq x+1) \wedge (y-1 \neq y) \\ \wedge \nexists z \in [x+1, y-2]. (z \neq z+1)$$

✗ *First value reappears later*

✗ *Some value appears twice*

✗ *Every value appears at most once*

Languages recognized by orbit-finite data monoids:

✓ Closed under all boolean operations

✓ *First value = last value*



$$\exists x, y. x \sim y \wedge x = \text{first} \wedge y = \text{last}$$

✓ *At least two distinct values*



$$\exists x, y. x \neq y \wedge y = x + 1$$

✓ *At least three distinct values*



$$\exists x, y. x \neq y \wedge (x \neq x + 1) \wedge (y - 1 \neq y) \\ \wedge \nexists z \in [x + 1, y - 2]. (z \neq z + 1)$$

✗ *First value reappears later*

✗ *Some value appears twice*

✗ *Every value appears at most once*

All data tests are
guarded by formulas
defining **bijections**!

Definition

Rigidly guarded MSO[~] is the fragment of MSO with data tests defined by the following grammar:

$$\varphi \mapsto x < y \mid x \in Y \mid \neg\varphi \mid \varphi \wedge \varphi \mid \exists x. \varphi \mid \exists Y. \varphi \mid \\ x \sim y \wedge \alpha(x, y)$$

where $\alpha(x, y)$ is generated by the same grammar and is **rigid** i.e. in every word, it **determines x from y and vice versa**.

Definition

Rigidly guarded MSO[~] is the fragment of MSO with data tests defined by the following grammar:

$$\varphi \mapsto x < y \mid x \in Y \mid \neg\varphi \mid \varphi \wedge \varphi \mid \exists x. \varphi \mid \exists Y. \varphi \mid \\ x \sim y \wedge \alpha(x, y)$$

where $\alpha(x, y)$ is generated by the same grammar and is **rigid** i.e. in every word, it **determines x from y and vice versa**.

 Rigidity can be checked, or enforced syntactically like in

$$\alpha^{\min}(x, y) = \alpha(x, y) \wedge \nexists x, y'. [x', y'] \not\subseteq [x, y] \wedge \alpha(x', y')$$

 We can use shorthands like

$$\alpha(x, y) \wedge x \not\sim y = \alpha(x, y) \wedge \neg(\alpha(x, y) \wedge x \sim y) \\ \alpha(x, y) \rightarrow x \sim y = \neg\alpha(x, y) \vee (\alpha(x, y) \wedge x \sim y)$$

Theorem 1

Rigidly guarded MSO^{\sim}
||
Orbit-finite data monoids.

...and as in the Schützenberger-McNaughton-Papert's theorem:

Theorem 2

Rigidly guarded \mathbf{FO}^{\sim}
||
Aperiodic orbit-finite data monoids.

From logic to monoids: induction and closure properties

- ① Negation: easy, by definition of recognizability
- ② Conjunction: product of orbit-finite data monoids

From logic to monoids: induction and closure properties

- 1 Negation: easy, by definition of recognizability
- 2 Conjunction: product of orbit-finite data monoids
- 3 Existential quantification: **powerset construction**

i.e. given $h: \underbrace{(\mathbb{N} \times \mathbb{B})^* \rightarrow M}_{\text{for a formula } \varphi(X)}$, construct $h': \underbrace{\mathbb{N}^* \rightarrow \wp(M)}_{\text{for the quantified formula } \exists X. \varphi(X)}$:

- elements of $\wp(M)$ are sets of elements of M
- product is naturally defined by $S \cdot T = \{s \cdot t \mid s \in S, t \in T\}$

From logic to monoids: induction and closure properties

- 1 Negation: easy, by definition of recognizability
- 2 Conjunction: product of orbit-finite data monoids
- 3 Existential quantification: **powerset construction**

i.e. given $h : (\mathbb{N} \times \mathbb{B})^* \rightarrow M$, construct $h' : \mathbb{N}^* \rightarrow \wp(M)$:

for a formula $\varphi(X)$

for the quantified formula $\exists X. \varphi(X)$

- elements of $\wp(M)$ are sets of **pairwise orbit-distinct** elements of M
- product is naturally defined by $S \cdot T = \{s \cdot t \mid s \in S, t \in T\}$
- **stronger invariant (projectability)** that forbids the following case

$$h(w, X) = s(\bullet) \quad \text{and} \quad h(w, X') = s(\bullet)$$

From logic to monoids: induction and closure properties

- 1 Negation: easy, by definition of recognizability
- 2 Conjunction: product of orbit-finite data monoids
- 3 Existential quantification: **powerset construction**

i.e. given $h : (\mathbb{N} \times \mathbb{B})^* \rightarrow M$, construct $h' : \mathbb{N}^* \rightarrow \wp(M)$:

for a formula $\varphi(X)$

for the quantified formula $\exists X. \varphi(X)$

- elements of $\wp(M)$ are sets of **pairwise orbit-distinct** elements of M
- product is naturally defined by $S \cdot T = \{s \cdot t \mid s \in S, t \in T\}$
- **stronger invariant (projectability)** that forbids the following case

$$h(w, X) = s(\bullet) \quad \text{and} \quad h(w, X') = s(\bullet)$$

- 4 Rigidly guarded data tests: product with **non-projectable** morphism

i.e. given $h : (\mathbb{N} \times \mathbb{B} \times \mathbb{B})^* \rightarrow M$, construct $h' : (\mathbb{N} \times \mathbb{B} \times \mathbb{B})^* \rightarrow M'$

for a rigid guard $\alpha(x, y)$

for the data test $\alpha(x, y) \wedge x \sim y$

From monoids to logic: generalization of Schützenberger's proof

Given a morphism $h : \mathbb{N}^* \rightarrow M$, logically define the language $h^{-1}(s)$ by induction on the size of the **infix-closed set** $s^\uparrow = \{t \mid s \in M \cdot t \cdot M\}$.

From monoids to logic: generalization of Schützenberger's proof

Given a morphism $h : \mathbb{N}^* \rightarrow M$, logically define the language $h^{-1}(s)$ by induction on the size of the **infix-closed set** $s^\uparrow = \{t \mid s \in M \cdot t \cdot M\}$.

Key ingredients in the classical aperiodic case:

$$\textcircled{1} s = (s \cdot M) \cap (M \cdot s) \cap s^\uparrow$$

From monoids to logic: generalization of Schützenberger's proof

Given a morphism $h : \mathbb{N}^* \rightarrow M$, logically define the language $h^{-1}(s)$ by induction on the size of the **infix-closed set** $s^\uparrow = \{t \mid s \in M \cdot t \cdot M\}$.

Key ingredients in the classical aperiodic case:

- ① $s = (s \cdot M) \cap (M \cdot s) \cap s^\uparrow$
- ② $h(w) \in (s \cdot M) \leftrightarrow \exists$ prefix $u \cdot a$ of w such that $h(u) \cdot h(a) \in (s \cdot M)$
(w.l.o.g. let u be **maximal** such that $h(u)^\uparrow \not\subseteq s^\uparrow$)

From monoids to logic: generalization of Schützenberger's proof

Given a morphism $h : \mathbb{N}^* \rightarrow M$, logically define the language $h^{-1}(s)$ by induction on the size of the **infix-closed set** $s^\uparrow = \{t \mid s \in M \cdot t \cdot M\}$.

Key ingredients in the classical aperiodic case:

$$\textcircled{1} \quad s = (s \cdot M) \cap (M \cdot s) \cap s^\uparrow$$

$$\textcircled{2} \quad h(w) \in (s \cdot M) \leftrightarrow \exists \text{ prefix } u \cdot a \text{ of } w \text{ such that } h(u) \cdot h(a) \in (s \cdot M) \\ \text{(w.l.o.g. let } u \text{ be } \mathbf{maximal} \text{ such that } h(u)^\uparrow \not\subseteq s^\uparrow)$$

$$\textcircled{3} \quad h(w) \in s^\uparrow \leftrightarrow \forall \text{ infixes } a \cdot u \cdot b \text{ of } w, h(a) \cdot h(u) \cdot h(b) \in s^\uparrow \\ \text{(w.l.o.g. let } u \text{ be } \mathbf{maximal} \text{ such that } h(u)^\uparrow \not\subseteq s^\uparrow)$$

From monoids to logic: generalization of Schützenberger's proof

Given a morphism $h : \mathbb{N}^* \rightarrow M$, logically define the language $h^{-1}(s)$ by induction on the size of the **infix-closed set** $s^\uparrow = \{t \mid s \in M \cdot t \cdot M\}$.

Key ingredients in the classical aperiodic case:

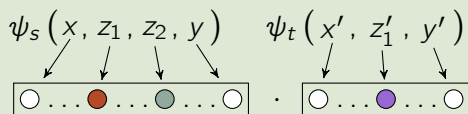
- ① $s = (s \cdot M) \cap (M \cdot s) \cap s^\uparrow$
- ② $h(w) \in (s \cdot M) \leftrightarrow \exists$ prefix $u \cdot a$ of w such that $h(u) \cdot h(a) \in (s \cdot M)$
(w.l.o.g. let u be **maximal** such that $h(u)^\uparrow \not\subseteq s^\uparrow$)
- ③ $h(w) \in s^\uparrow \leftrightarrow \forall$ infixes $a \cdot u \cdot b$ of w , $h(a) \cdot h(u) \cdot h(b) \in s^\uparrow$
(w.l.o.g. let u be **maximal** such that $h(u)^\uparrow \not\subseteq s^\uparrow$)

Additional difficulties with data monoids:

- ⊕ products depend on data values
- ⊕ data tests must be performed under rigid guards
- ⊕ “data groups” must be considered to keep track of data values

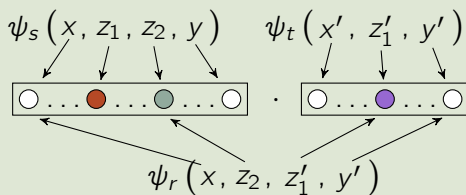
Example

Consider the product $s(\bullet\bullet) \cdot t(\bullet) = r(\bullet\bullet)$ of two elements



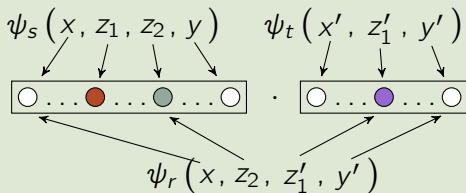
Example

Consider the product $s(\bullet \bullet) \cdot t(\bullet) = r(\bullet \bullet)$ of two elements



Example

Consider the product $s(\bullet\bullet) \cdot t(\bullet) = r(\bullet\bullet)$ of two elements



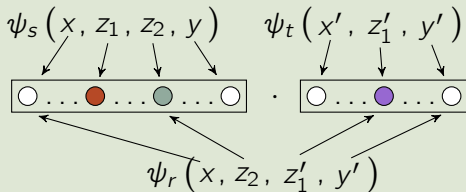
Stronger inductive hypothesis

Given a morphism $h : \mathbb{N}^* \rightarrow M$ and an **orbit** $s(\circ, \dots, \circ)$, one can construct the following objects by induction on s^\dagger :

- 1 a formula $\varphi_s^\dagger(x, y)$ that defines the infixes $w[x, y] \in h^{-1}(s^\dagger)$

Example

Consider the product $s(\bullet \bullet) \cdot t(\bullet) = r(\bullet \bullet)$ of two elements



Stronger inductive hypothesis

Given a morphism $h : \mathbb{N}^* \rightarrow M$ and an **orbit** $s(\circ, \dots, \circ)$, one can construct the following objects by induction on s^\dagger :

- 1 a formula $\varphi_s^\dagger(x, y)$ that defines the infixes $w[x, y] \in h^{-1}(s^\dagger)$
- 2 for each rigid guard $\alpha(x, y)$ that entails $\varphi_s^\dagger(x, y)$, a rigid formula $\psi_s^\alpha(x, z_1, \dots, z_k, y)$ such that

$$w \models \psi_s^\alpha(x, z_1, \dots, z_k, y) \rightarrow h(w[x, y]) = s(w(z_1), \dots, w(z_k))$$

From monoids to logic: non-aperiodic case

- 1 Check that $h(w) \in (M \cdot s) \cap (s \cdot M) \cap s^\uparrow$

From monoids to logic: non-aperiodic case

1 Check that $h(w) \in (M \cdot s) \cap (s \cdot M) \cap s^\uparrow$

2 Factorize $w = \overbrace{u_1 a_1 \leftarrow v_1 \rightarrow b_1}^{w_1} \overbrace{u_2 a_2 \leftarrow v_2 \rightarrow b_2}^{w_2} \dots$

with $h(u_i)^\uparrow \not\subseteq s^\uparrow$ and v_i **maximal** such that $h(v_i)^\uparrow \not\subseteq s^\uparrow$

From monoids to logic: non-aperiodic case

- 1 Check that $h(w) \in (M \cdot s) \cap (s \cdot M) \cap s^\dagger$

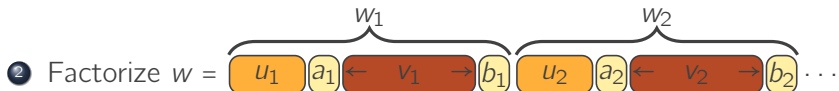
- 2 Factorize $w =$ \dots
The diagram shows a sequence of boxes representing characters in a word w . The boxes are arranged in a row: u_1 (orange), a_1 (yellow), v_1 (dark red), b_1 (yellow), u_2 (orange), a_2 (yellow), v_2 (dark red), b_2 (yellow), followed by an ellipsis. Brackets above the boxes group them into two segments: w_1 covers u_1, a_1, v_1, b_1 and w_2 covers u_2, a_2, v_2, b_2 . Inside the dark red boxes v_1 and v_2 , there are arrows pointing outwards: \leftarrow on the left and \rightarrow on the right.

with $h(u_i)^\dagger \not\subseteq s^\dagger$ and v_i **maximal** such that $h(v_i)^\dagger \not\subseteq s^\dagger$

- 3 Guess elements s_1, s_2, \dots (over a bounded data domain) and using i.h. and products, check that each s_i is a **renaming of $h(w_i)$**

From monoids to logic: non-aperiodic case

- 1 Check that $h(w) \in (M \cdot s) \cap (s \cdot M) \cap s^\uparrow$

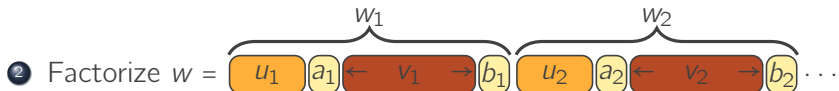


with $h(u_i)^\uparrow \not\subseteq s^\uparrow$ and v_i **maximal** such that $h(v_i)^\uparrow \not\subseteq s^\uparrow$

- 3 Guess elements s_1, s_2, \dots (over a bounded data domain) and using i.h. and products, check that each s_i is a **renaming of $h(w_i)$**
- 4 Check that each pair s_i, s_{i+1} is a **renaming of $h(w_i), h(w_{i+1})$** and, similarly, that s_1, s_n is a **renaming of $h(w_1), h(w_n)$**

From monoids to logic: non-aperiodic case

- 1 Check that $h(w) \in (M \cdot s) \cap (s \cdot M) \cap s^\uparrow$



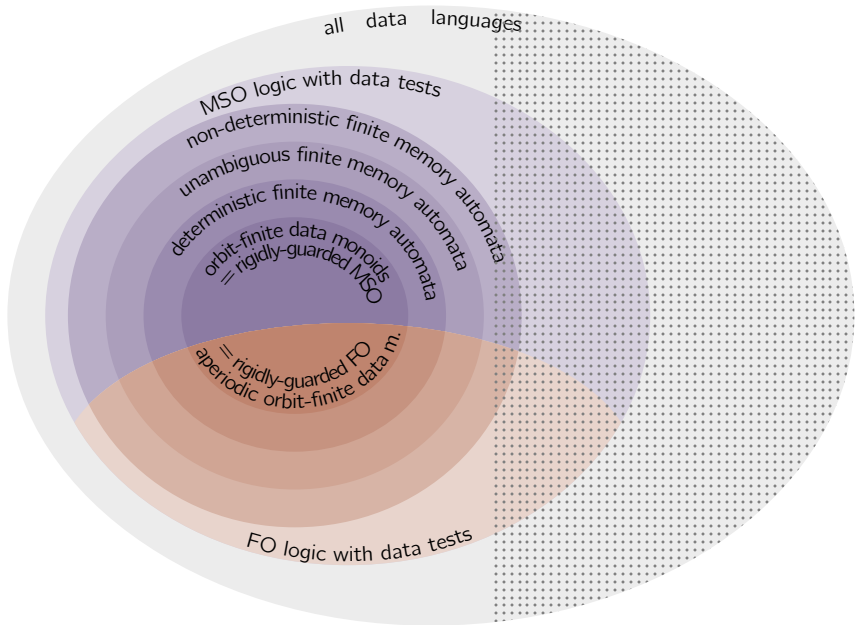
with $h(u_i)^\uparrow \not\subseteq s^\uparrow$ and v_i **maximal** such that $h(v_i)^\uparrow \not\subseteq s^\uparrow$

- 3 Guess elements s_1, s_2, \dots (over a bounded data domain) and using i.h. and products, check that each s_i is a **renaming of $h(w_i)$**
- 4 Check that each pair s_i, s_{i+1} is a **renaming of $h(w_i), h(w_{i+1})$** and, similarly, that s_1, s_n is a **renaming of $h(w_1), h(w_n)$**
- 5 Inductively compute the partial products $s_1 \cdot \dots \cdot s_i$, for $i = 1, \dots, n$:

Lemma

$s_1 \cdot \dots \cdot s_n$ is a renaming of $h(w) = h(w_1) \cdot \dots \cdot h(w_n)$.

Logics, automata, and monoids over infinite alphabets:



Logics, automata, and monoids over infinite alphabets:

