# STAR-FREE LANGUAGES

So far we know:

by Büchi's theorem



WMSO $\overset{\sim}{=}$ REG

FO[<] = (?)

$(aa)^*$

The goal of this part is finding an algebraic characterisation of FO[<]-definable languages.

$\notin$ FO by EF-theorem and winning strategy of Duplicator

## Def (Star-free Languages)

Let $\Sigma$ be an alphabet. The class of <u>star-free languages over $\Sigma$</u> (denoted by $SF_\Sigma$) is the smallest class so that

(1) $\emptyset \in SF_\Sigma$, $\{\varepsilon\} \in SF_\Sigma$, $\{a\} \in SF_\Sigma$ (for each $a \in \Sigma$)

(2) If $L_1$ and $L_2$ are in $SF_\Sigma$ then

- $L_1 \cdot L_2 \in SF_\Sigma$
- $L_1 \cup L_2 \in SF_\Sigma$
- $\overline{L_1} \in SF_\Sigma$ $\longleftarrow$ main difference wrt REG: no kleene star but complement

Intuition: FO[<] allows us to speak about union ($\vee$), finitely many positions at once ($L \cdot L'$) and can complement ($\neg$) but cannot speak of unboundedly many positions at the same time ($L^*$).

1

**Examples** The presence of complement allows us to express languages which can also be expressed using Kleene star (but not all of them)

- $\Sigma^* \in SF_\Sigma$ since $\Sigma^* = \overline{\emptyset}$

- If $L_1, L_2 \in SF_\Sigma$ then
  - $L_1 \cap L_2 = \overline{\overline{L_1} \cup \overline{L_2}} \in SF_\Sigma$
  - $L_1 \setminus L_2 = L_1 \cap \overline{L_2} \in SF_\Sigma$

- Let $D \subseteq \Sigma$. Then $D^* \in SF_\Sigma$ as $D^* = \Sigma^* \setminus (\Sigma^* \cdot (\Sigma \setminus D) \cdot \Sigma^*)$

- The language $(ab)^*$ is star-free because

$$(ab)^* = \Sigma^* \setminus b\Sigma^* \setminus \Sigma^* aa \Sigma^* \setminus \Sigma^* bb \Sigma^* \setminus \Sigma^* a$$

  not starting with b    no consecutive as    no consecutive bs    not ending with a

**Theorem** (McNaughton & Papert '71)   Let $L \subseteq \Sigma^*$.
  1. If $L$ is star-free then $L$ is FO[<]-definable
  2. If $L$ is FO[<]-definable then $L$ is star-free

**Proof of ①** as exercise

**Proof of ②** requires some insights on FO[<] to handle quantifiers

- if $w \in \mathcal{L}(\exists x : \Psi)$ and $qd(\Psi) \leq K$ then
  $w = u\,a\,v$ — position assigned to $x$ to satisfy $\Psi$.
  But then for all $u' \equiv_K u$ and $v' \equiv_K v$, $u'av'$ must also be in $\mathcal{L}(\exists x : \Psi)$

- the equivalence class of $u$, $[u]_{\equiv_K} := \{u' \mid u' \equiv_K u\}$ can be represented as the language of a formula with $qd \leq K$

- $\equiv_K$ as finitely many classes

/2

# Notation

Let $\vec{s} = s_1 \ldots s_m$ and $\vec{x} = x_1 \ldots x_m$.

- $d_m(\varphi(\vec{x})) := \{ (S_v, \vec{s}) \mid S_v, [\vec{s}/\vec{x}] \models \varphi(\vec{x}) \}$

- $\Phi_{k,m} := \{ \varphi(\vec{x}) \in FO[<] \mid qd(\varphi(\vec{x})) \leq k \}$

- $\Phi_k(S_v, \vec{s}) := \{ \varphi \in \overline{\Phi}_{k,m} \mid S_v, [\vec{s}/\vec{x}] \models \varphi \}$

# Lemma

For all $k, m \in \mathbb{N}$ the following holds

① The relation $\equiv_{k,m}$ is an equivalence

② For all $v \in \Sigma^*$ and positions $\vec{s} = s_1 \ldots s_m$,

$$[(S_v, \vec{s})]_{\equiv_{k,m}} = \bigcap_{\varphi \in \Phi_k(S_v, \vec{s})} d_m(\varphi)$$

③ There exists a finite set $\widetilde{\Phi}_{k,m} \subseteq \Phi_{k,m}$ such that

$$\underline{\Phi}_{k,m} = \{ \varphi \mid \exists \varphi' \in \widehat{\Phi}_{k,m} \text{ s.t. } d_m(\varphi) = d_m(\varphi') \}$$

(i.e. $\underline{\Phi}_{k,m}$ is finite up to logical equivalence)

④ The equivalence $\equiv_{k,m}$ has finitely many classes, each characterised by a formula $\varphi_{[(S_v, \vec{s})]_{\equiv_k}}$ such that

$$S_w, [\vec{t}/\vec{x}] \models \varphi_{[(S_v, \vec{s})]_{\equiv_{k,m}}} \iff (S_w, \vec{t}) \in [(S_v, \vec{s})]_{\equiv_{k,m}}$$

# Proof

① easy

② "$\supseteq$" immediate: every structure in $[(S_v, \vec{s})]_{\equiv_{k,m}}$ satisfies the same formulas $\varphi \in \Phi_{k,m}$ as $S_v \vec{s}$ by definition of $\equiv_{k,m}$ so it will be in every $d_m(\varphi)$

"$\subseteq$" exercise

④ We need to show that, up to logical equivalence, there are only finitely many formulas with $n$ variables an $qd \leq k$.

We proceed by induction on $k$.

(k=0) For every quantifier free formula $\varphi(\bar{x}) \in \Phi_0[<]$, we can obtain an equivalent formula $\varphi'(\bar{x})$ in DNF

$$\varphi'(\bar{x}) = \bigvee_{j \in J} \underbrace{\left( \bigwedge_{i \in I} c_{ij} \right)}_{\text{disjunct}}$$

such that no disjunct is repeated and no conjunct in each disjunct is repeated.

Let us count how many different atomic propositions we can write using at most $n$ variables:

$$P_a(x) \rightsquigarrow |\Sigma| n$$

$$x < y \rightsquigarrow n^2$$

Since each $c_{ij}$ is either atomic or a negation of atomic formula we get at most $2(|\Sigma| n + n^2)$ distinct ~~disjuncts~~ $c_{ij}$.

Therefore we have at most $2^{2(|\Sigma| n + n^2)}$ distinct disjuncts and at most $2^{2^{2(|\Sigma| n + n^2)}}$ DNF that we need to put in $\widetilde{\Phi}_{0,n}$ to represent all possible formulas in $\Phi_{0,n}$ up to logical equivalence

(k+1) Induction step is analogous: Count the DNF of formulas with $qd \leq k$ and $n+1$ free variables.

⑤ By ② we have $[(S_v, \bar{s})]_{\equiv_{k,n}} \bigcap_{\varphi \in \Phi_k(S_v, \bar{s})} d_m(\varphi)$ but by ④

we know that for each $\varphi \in \Phi_{k,n}$ there is a formula $\varphi' \in \widetilde{\Phi}_{k,n}$ such that $d(\varphi) = d(\varphi')$

Therefore

$$[(S_v, \bar{s})]_{\equiv_{k,n}} = \bigcap_{\varphi \in \Phi_k(S_v, \bar{s})} d_m(\varphi) = \bigcap_{\varphi \in \widetilde{\Phi}_k \cap \Phi_k(S_v, \bar{s})} d_m(\varphi) = d_m\left( \bigwedge_{\varphi \in \widetilde{\Phi}_k \cap \Phi_k(S_v, \bar{s})} \varphi \right)$$

$\underbrace{\phantom{xxxxxxxxxxxx}}_{\text{finite!}}$

Hence there are $\overset{\text{at most}}{}$ as many classes in $\equiv_{k,n}$ as there are conjunctions of formulas in $\widetilde{\Phi}_k$. ✓ ☐

Now we can proceed with the proof of ② of McNaughter&Papert '71:

Let $\varphi$ be a FO[<]-sentence. Then $\mathcal{L}(\varphi)$ is star-free

Proof We proceed by induction on the quantifier-depth $k$ of closed formulas $\varphi$

(k=0) The only closed formulas with $qd=0$ are, up to logical equivalence,
true and ¬true. $\mathcal{L}(\text{true}) = \Sigma^* = \overline{\emptyset} \in SF_\Sigma$, $\mathcal{L}(\neg\text{true}) = \emptyset \in SF_\Sigma$ ✓

(k+1) Assume formulas of $qd \leq k$ define star-free languages.
A formula of $qd = k+1$ will be a boolean combination of
formulas of the form $\varphi = \exists x: \psi$ with $qd(\psi) \leq k$.
While the boolean connectives can be easily handled by
using the closure properties of $SF_\Sigma$, the existential quantification
requires the following characterisation

$\underline{\text{CLAIM}}$: $\circledast$ $\quad \mathcal{L}(\exists x: \psi) = \bigcup\left\{ [S_u]_{\equiv_k} \cdot a \cdot [S_v]_{\equiv_k} \;\middle|\; S_{uav}, [^{|u|}/_x] \models \right.$

By Lemma ④ the union $\bigcup$ is finite and there are formulas

$\varphi_{[u]_{\equiv_k}}$ and $\varphi_{[v]_{\equiv_k}}$ of $qd \leq k$ such that $\mathcal{L}(\varphi_{[u]_{\equiv_k}}) = [S_u]_{\equiv_k}$
$\qquad\qquad\qquad\qquad$ and $\mathcal{L}(\varphi_{[v]_{\equiv_k}}) = [S_v]_{\equiv_k}$

Since they have $qd \leq k$ we can apply our induction hypothesis to obtain
that $R_u = \mathcal{L}(\varphi_{[u]_{\equiv_k}})$ and $R_v = \mathcal{L}(\varphi_{[v]_{\equiv_k}})$ are star-free languages.

Then $\mathcal{L}(\exists x: \psi) = \bigcup\left\{ R_u \cdot a \cdot R_v \;\middle|\; S_{uav}, [^{|u|}/_x] \models \psi \right\}$

which is a finite union of concatenations of star-free languages,
and hence star-free.
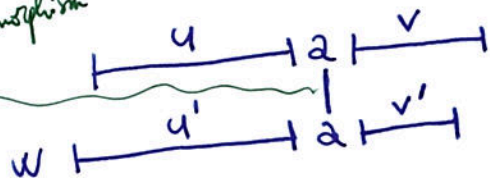
All is left to prove is the CLAIM $\circledast$

⑤

# Proof of ⊛

"⊆" Let $S_w \models \exists x: \Psi$. Then there is a position $i$ with $w(i) = a$
for some $a \in \Sigma$ such that $w = uav$   $|u| = i$ and
$S_{uav}, [^i/_x] \models \Psi$, proving that $w \in [S_u]_{\equiv_k} \cdot a \cdot [S_v]_{\equiv_k}$ and hence in the union

"⊇" Let $w \in [S_u]_{\equiv_k} \cdot a \cdot [S_v]_{\equiv_k}$ for some $u, v \in \Sigma^*$.
Then there are $u'$ and $v'$ such that $S_u \equiv_k S_{u'}$ and $S_v \equiv_k S_{v'}$
and $w = u'av'$. By the EF-theorem, Duplicator wins the games
$G_k((S_u, S_{u'})$ and $G_k(S_v, S_{v'})$.
Consider the game $G_k((S_{uav}, |u|), (S_{u'av'}, |u'|))$
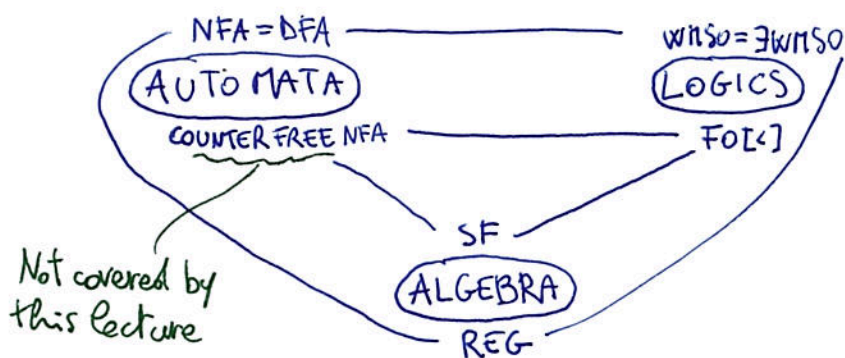
is partial isomorphism



We can win the game by playing
winning strategy of $G_k(S_u, S_{u'})$ in
first half and the one of $G_k(S_v, S_{v'})$
on second half.

Therefore $(S_{u'av'}, |u'|) \equiv_{k,1} (S_{uav}, |u|)$ by EF-theorem.
Since by assumption $S_{uav}, [^{|u|}/_x] \models \Psi$ we get $S_{u'av'}, [^{|u'|}/_x] \models \Psi$
and therefore $S_w \models \exists x: \Psi$ which proves $w \in L(\exists x: \Psi)$
as desired □

---

Overall picture:



Not covered by
this lecture

## Techniques

REG → NFA  Closure properties
NFA → REG  Arden's Lemma (Algebraic view)
COMPLEMENTATION VIA DETERMINISATION
NFA ≡ DFA (Automata view) Powerset
NFA → WMSO  Encoding runs using second order
quantification (logics view)
WMSO → NFA  Encoding ∃x with Σ×and
projection
$(aa)^* \notin FO[<]$  Quantifier depth
and EF-theorem
SF → FO[<] easy
FO[<] → SF by studying $\equiv_k$ classes
and EF-theorem.

6