



Safety Proof Synthesis for Regular Transition Systems

Chih-Duo Hong

FLOLAC 2025



Transition system

- A *transition system* is a triple (S, I, T) , where
 - S is the set of states
 - $I \subseteq S$ is the set of initial states
 - $T \subseteq S \times S$ is the set of transitions
- A *trace* of (S, I, T) is a sequence $\sigma = \sigma_0 \sigma_1 \sigma_2 \dots \in S^\omega$ such that
 - $\sigma_0 \in I$
 - for all $i \geq 0$, $(\sigma_i, \sigma_{i+1}) \in T$
- That is, a trace is a finite/infinite sequence of consecutive transitions starting from an initial state.

Safety property

- Safety properties are concerned with the assurance that certain undesirable behaviors will never occur in a system
- Typical safety properties of software:
 1. **Division by zero:** A program will never divide a number by zero
 2. **Null dereference:** A program will never dereference a null or uninitialized pointer
 3. **Data race:** A shared variable will never be updated simultaneously
- Safety of a transition system
 - Does every trace never reach a bad state?
- Model checking a safety property
 - Yes + proof
 - No + counterexample (a system trace that reaches a bad state)

Liveness property

- Liveness properties are concerned with the assurance that certain desirable behaviors will eventually occur in a system
- Typical liveness properties of software:
 1. **Termination:** The program eventually terminates
 2. **Response:** The system eventually responds to an input event
 3. **Non-Starvation:** The thread/process is scheduled for execution infinitely often
- Liveness of a transition system
 - Does every trace eventually reach a good state?
- Model checking a liveness property
 - Yes + proof
 - No + counterexample (a system trace that never reaches a good state)

Liveness property

- Liveness properties are concerned with the assurance that certain desirable behaviors will eventually occur in a system
- Typical liveness properties of software:
 1. **Termination:** The program eventually terminates
 2. **Response:** The system eventually responds to an input event
 3. **Non-Starvation:** The thread/process is scheduled for execution infinitely often
- Liveness of a transition system
 - Does every trace eventually reach a good state?
- Liveness or safety?
 - **Timely response:** The system should respond within a given time frame
 - **Progress:** At least one thread/process is running at any point of time

Symbolic transition system

- We usually specify and reason about a transition system using a *symbolic representation*
- In this lecture, we will introduce two symbolic representations for infinite-state transition systems:
 1. Logical formulas (over a background theory)
 2. Regular languages

Symbolic transition system

- We usually specify and reason about a transition system using a *symbolic representation*
- In this lecture, we will introduce two symbolic representations for infinite-state transition systems:
 - 1. Logical formulas (over a background theory)**
 2. Regular languages

Formulas as symbolic representation

- A *symbolic transition system* is a tuple (V, I, T) , where
 - V is a set of variables,
 - I is a formula over variables V
 - T is a formula over variables $V \cup V'$
(E.g., $i' = i + 1$ is a formula over $\{i\} \cup \{i'\}$ that increments i by 1)
- A *state* $\sigma \in A$ is a type-consistent assignment to V
- A *trace* of (V, I, T) is a sequence $\sigma = \sigma_0 \sigma_1 \sigma_2 \dots \in A^\omega$, where
 - $\sigma_0 \models I(V)$
 - $\sigma_i, \sigma'_{i+1} \models T(V, V')$ for all $i \geq 0$

Example: the Collatz transition system

- Consider the following operation on a natural number:
 - If the number is even, divide it by two. ($x' = x/2$)
 - If the number is odd, triple it and add one. ($x' = 3x + 1$)
- Applying this operation to a number repeatedly will generate a sequence, for example: $21 \rightarrow 64 \rightarrow 32 \rightarrow 16 \rightarrow 8 \rightarrow 4 \rightarrow 2 \rightarrow 1$
- The corresponding symbolic transition system is (V, I, T) , where

$$V := \{x\}$$

$$I := (x \geq 1)$$

$$T := (\exists k. x = 2k \wedge x' = k) \vee (\exists k. x = 2k + 1 \wedge x' = 3x + 1)$$

Example: the Collatz transition system

- Consider the following operation on a natural number:
 - If the number is even, divide it by two. ($x' = x/2$)
 - If the number is odd, triple it and add one. ($x' = 3x + 1$)
- Applying this operation to a number repeatedly will generate a sequence, for example: $21 \rightarrow 64 \rightarrow 32 \rightarrow 16 \rightarrow 8 \rightarrow 4 \rightarrow 2 \rightarrow 1$

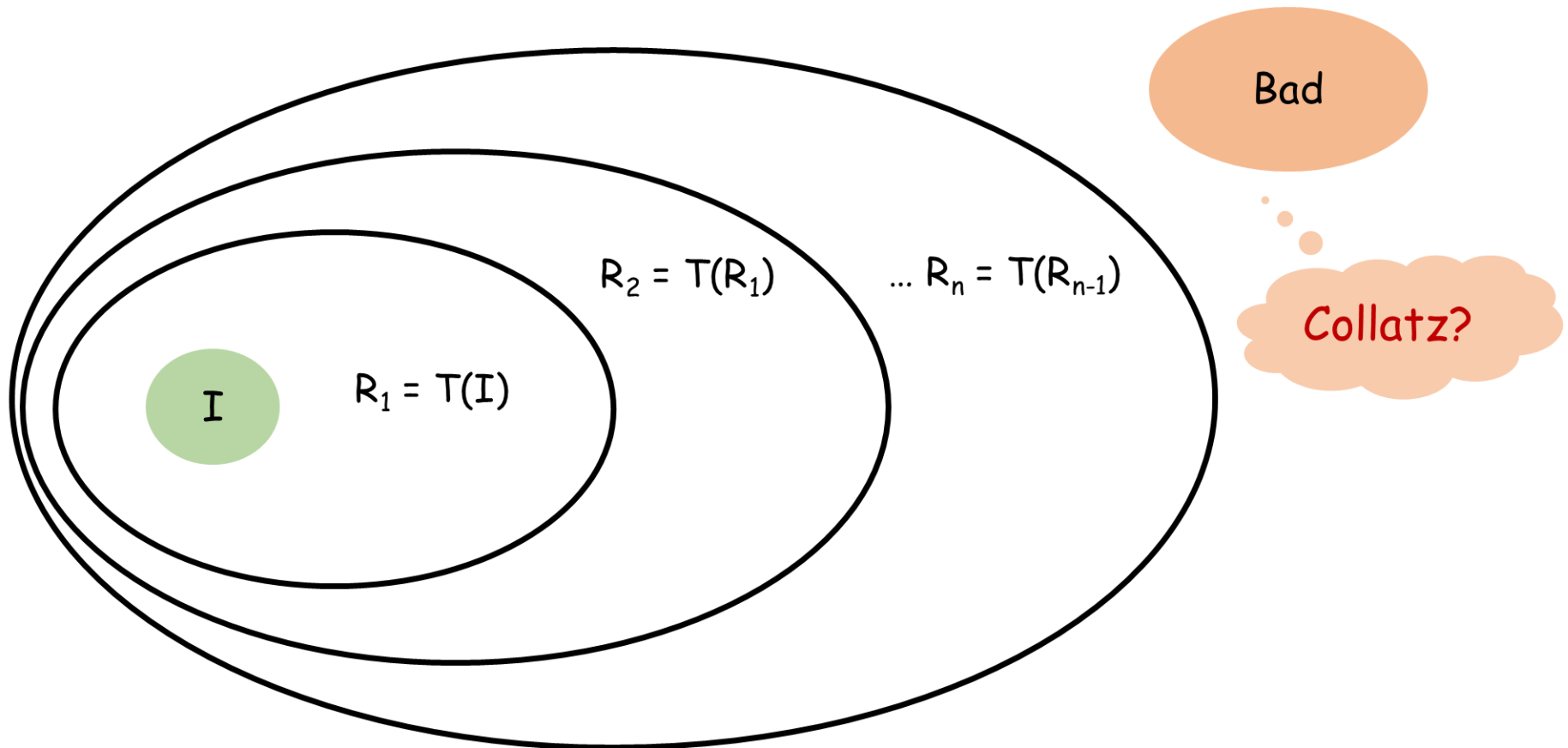
An example safety property:

“Every sequence starting from a power of 2 will reach no odd numbers but 1.”

An example liveness property: (The Collatz conjecture)

“Every sequence will eventually reach 1.”

Forward reachability analysis



$$T(A) := \{ s' : s \in A \text{ and } (s, s') \in T \}$$

Inductive invariant

A set of states Inv is an **inductive invariant** if it satisfies the following three conditions:

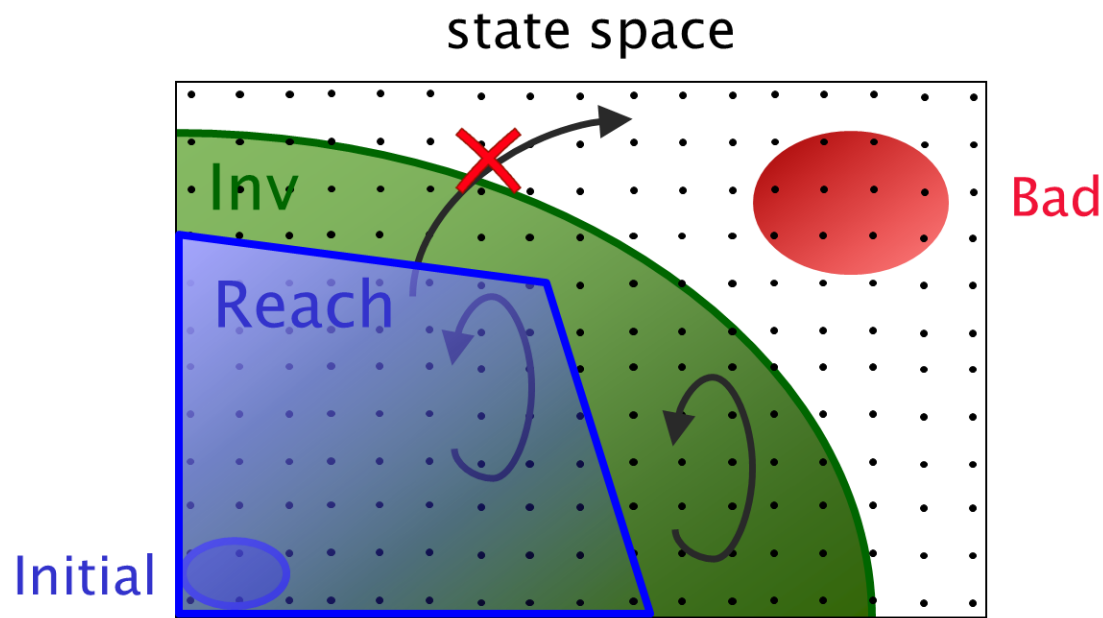
- Initiation: $I \subseteq Inv$
- Consecution: $T(Inv) \subseteq Inv$
- Safety: $Inv \cap B = \emptyset$

When I , F , B , Inv are expressed in formulas, these conditions are equivalent to

$$\begin{aligned} I(V) &\Rightarrow Inv(V) \\ Inv(V) \wedge T(V, V') &\Rightarrow Inv(V') \\ Inv(V) &\Rightarrow \neg B(V) \end{aligned}$$

Inductive invariant (cont'd)

- Initiation: $I \subseteq Inv$
- Consecution: $T(Inv) \subseteq Inv$
- Safety: $Inv \cap B = \emptyset$



A system is safe iff it has an inductive invariant

Example: inductive invariant

- Consider a symbolic transition system (V, I, T) , where

$$V := \{x, y\}$$

$$I := x = 1 \wedge y = 1$$

$$T := (x' = x + y) \wedge (y' = y + x)$$

- We want to prove the safety property $P := y \geq 1$.

Example: inductive invariant (cont'd)

$P := y \geq 1$ is not an inductive invariant

- $I \Rightarrow P$:

- $(x = 1 \wedge y = 1) \Rightarrow y \geq 1$

- But $P \wedge T \not\Rightarrow P'$:

- $y \geq 1 \wedge (x' = x + y \wedge y' = x + y) \not\Rightarrow y' \geq 1$

$V := \{x, y\}$

$I := x = 1 \wedge y = 1$

$T := (x' = x + y) \wedge (y' = y + x)$

$P := y \geq 1$

$I(V) \Rightarrow \text{Inv}(V)$

$\text{Inv}(V) \wedge T(V, V') \Rightarrow \text{Inv}(V')$

$\text{Inv}(V) \Rightarrow \neg B(V)$

Example: inductive invariant (cont'd)

$P := y \geq 1$ is not an inductive invariant

- $I \Rightarrow P$:

- $(x = 1 \wedge y = 1) \Rightarrow y \geq 1$

- But $P \wedge T \not\Rightarrow P'$:

- $y \geq 1 \wedge (x' = x + y \wedge y' = x + y) \not\Rightarrow y' \geq 1$

$$V := \{x, y\}$$

$$I := x = 1 \wedge y = 1$$

$$T := (x' = x + y) \wedge (y' = y + x)$$

$$P := y \geq 1$$

Consider $\text{Inv} := x \geq 0 \wedge y \geq 1$

- $(x = 1 \wedge y = 1) \Rightarrow x \geq 0 \wedge y \geq 1$

- $x \geq 0 \wedge y \geq 1 \wedge (x' = x + y \wedge y' = x + y) \Rightarrow x' \geq 0 \wedge y' \geq 1$

- $x \geq 0 \wedge y \geq 1 \Rightarrow y \geq 1$

Property proved!

$$I(V) \Rightarrow \text{Inv}(V)$$

$$\text{Inv}(V) \wedge T(V, V') \Rightarrow \text{Inv}(V')$$

$$\text{Inv}(V) \Rightarrow \neg B(V)$$

Example: inductive invariant (cont'd)

Induction hypothesis

$P := y \geq 1$ is not a Base case invariant

$V := \{x, y\}$

$I := x = 1 \wedge y = 1$

$T := (x' = x + y) \wedge (y' = y + x)$

$P := y \geq 1$

• $I \Rightarrow P$:

– $(x = 1 \wedge y = 1)$ Induction step

• But $P \wedge T \not\Rightarrow P'$:

– $y \geq 1 \wedge (x' = x + y \wedge y' = x + y)$

Strengthening the induction hypothesis

Consider $\text{Inv} := x \geq 0 \wedge y \geq 1$

– $(x = 1 \wedge y = 1) \Rightarrow x \geq 0 \wedge y \geq 1$

– $x \geq 0 \wedge y \geq 1 \wedge (x' = x + y \wedge y' = x + y) \Rightarrow x' \geq 0 \wedge y' \geq 1$

– $x \geq 0 \wedge y \geq 1 \Rightarrow y \geq 1$

Property proved!

Exercise

- Consider a transition system (V, I, T) such that for some $m \geq 0$,

$$V := \{x, y\}$$

$$I := x = 0 \wedge y = m$$

$$T := (x < m \Rightarrow (x' = x + 1 \wedge y' = y - 1)) \\ \wedge (x \geq m \Rightarrow (x' = x \wedge y' = y))$$

- Consider the safety property $P := 0 \leq y \leq m$.
- Is P an inductive invariant? If not, how to strengthen it?

Symbolic transition system

- We usually specify and reason about a transition system using a *symbolic representation*
- In this lecture, we introduce two common symbolic representation for infinite transition systems:
 1. Logical formulas (over a background theory)
 - 2. Regular languages**

Regular language as symbolic representation

- For a **finite alphabet** Σ and a **padding symbol** $\#$.
- A **pair** of finite strings (u, v) can be represented by a string

$$u \otimes v := \begin{bmatrix} u_1 \\ v_1 \end{bmatrix} \cdots \begin{bmatrix} u_n \\ v_n \end{bmatrix} \in u\#^* \times v\#^*,$$

where $n = \max(|u|, |v|)$.

- A **binary relation** R over strings in Σ^* can be represented by language (i.e. a set of strings)

$$\{u \otimes v : (u, v) \in R\} \subseteq \Sigma^*\#^* \times \Sigma^*\#^*.$$

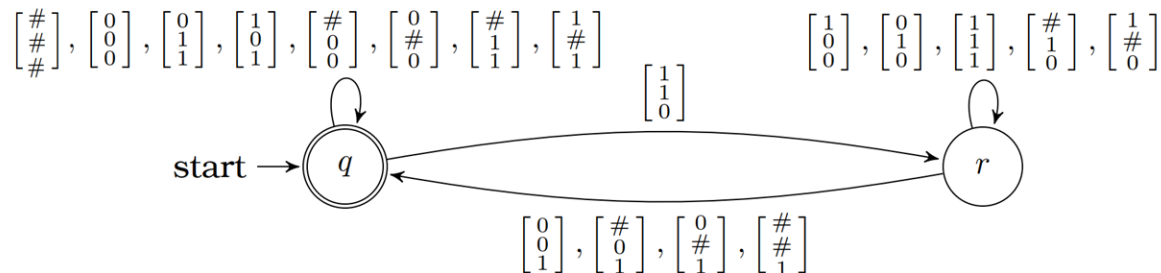
The same idea can be extended to general n -ary relations.

Regular language as symbolic representation (cont'd)

- A **regular set** (over strings) is a set of strings that can be encoded by a regular language.
- A **regular relation** (over strings) is a relation that can be encoded by a regular language.
- Operations over regular sets and relations correspond to automata manipulations over finite-state automata.
 - Intersection, Union, Difference, Complement
 - Composition, Inverse, Image, Pre-image

Regular language as symbolic representation (cont'd)

- Many data types can be encoded as finite strings over a finite alphabet:
 - Mathematical integer / rational numbers?
 - Machine integer / floating point numbers?
- Many data operations can be represented by regular relations.
- For example, the relation $R = \{(x, y, z) : x + y = z\}$ can be defined by intersecting $(\{0,1\}^*\#^*)^3$ with the language of



Regular language as symbolic representation (cont'd)

Every relation definable in $\text{FO}(\mathbb{N}, 0, 1, +, <)$ (aka. Presburger arithmetic) can be represented by a regular language!

- How do we perform logical operations using automata manipulations?
 - And, Or, Not
 - Exist, Forall

$$\text{project } W: \begin{matrix} \cancel{w} \\ z: \begin{bmatrix} \ominus \\ 1 \end{bmatrix} \end{matrix} \mapsto Z: [1]$$

Regular language as symbolic representation (cont'd)

Exercises

- Let S be a regular set of natural numbers.
 - How do you express the set $2S := \{2n : n \in S\}$?
 - How do you express kS for general $k \in \mathbb{N}$?
- How do you model-check the following statement?
 - $S(x) \wedge 2S(x)$ is satisfiable
 - $S(0) \Rightarrow 2S(0)$ is true
 - $S(x) \Rightarrow 2S(2x)$ is valid
 - $\forall x. S(x) \Rightarrow 2S(2x)$ is true

Regular language as symbolic representation (cont'd)

Every relation definable in $\text{FO}(\mathbb{N}, 0, 1, +, <)$ (aka. Presburger arithmetic) can be represented by a regular language!

- How do we perform logical operations using automata manipulations?
 - And, Or, Not
 - Exist, Forall

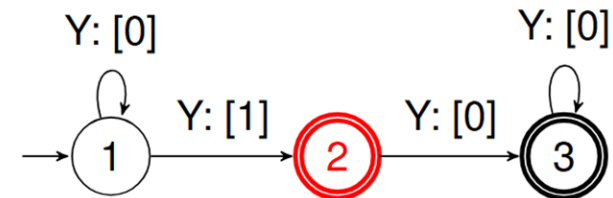
$$\text{project } W: \begin{matrix} \cancel{w} & \oplus \\ z & 1 \end{matrix} \mapsto Z: [1]$$

- **Many theories are effectively expressible by regular languages (or variants thereof). These theories can be checked by automata operations and are therefore decidable.**

Regular language as symbolic representation (cont'd)

Issues with projection (existential quantification) in the binary encoding

- After removing of the tracks not all models would be accepted(problem with 0-padding)
- So after projection, we need to adjust the final states by **saturation**, i.e., pump the final states with all states backward reachable with 0
- Consider $B(y) := \exists x. A(x, y)$



Regular language as symbolic representation (cont'd)

Issues with projection (existential quantification) in the binary encoding

- After removing of the tracks not all models would be accepted (problem with 0-padding)
- So after projection, we need to adjust the final states by **saturation**, i.e., pump the final states with all states backward reachable with 0
- Consider $B(y) := \exists x. A(x, y)$
- Illustration:

Regular language as symbolic representation (cont'd)

Every relation definable in $\text{FO}(\mathbb{N}, 0, 1, +, <)$ (aka. Presburger arithmetic) can be represented by a regular language!

- How do we perform logical operations using automata manipulations?
 - Validity and Satisfiability
 - And, Or, and Not
 - Exist and Forall

Many theories are effectively expressible by regular languages (or variants thereof). These theories can be checked by automata operations and are therefore decidable.

Regular transition system

A **regular transition system** (RTS) is a triple (Σ, I, T) , where $I \subseteq \Sigma^*$ is a regular set, and $T \subseteq \Sigma^* \times \Sigma^*$ is a regular relation.

- Each **state** is a finite string over the alphabet Σ
- The set of **initial states** is a regular set $\llbracket I \rrbracket \subseteq \Sigma^*$
- The **transition relation** is regular relation $\llbracket T \rrbracket \subseteq \Sigma^* \times \Sigma^*$

Regular transition system (cont'd)

A **regular transition system** (RTS) is a triple (Σ, I, T) , where $I \subseteq \Sigma^*$ is a regular set, and $T \subseteq \Sigma^* \times \Sigma^*$ is a regular relation.

Example The Collatz problem can be modeled by an RTS.

Recall that Presburger-definable relations can be encoded in regular languages.

$$V := \{x\}$$

$$I := (x \geq 1)$$

$$T := (\exists k. x = 2k \wedge x' = k) \vee (\exists k. x = 2k + 1 \wedge x' = 3x + 1)$$

Regular transition system (cont'd)

A **regular transition system** (RTS) is a triple (Σ, I, T) , where $I \subseteq \Sigma^*$ is a regular set, and $T \subseteq \Sigma^* \times \Sigma^*$ is a regular relation.

Example The Collatz problem can be modeled by an RTS.

Recall that Presburger-definable relations can be encoded in regular languages.

Example Turing machine's configuration graph can be modeled by an RTS.

A TM with a two-sided tape can be simulated by a TM with a one-sided tape.

Each TM transition modifies at most one tape position, and hence is regular.

Safety of regular transition systems

Fix an RTS (Σ, I, T) . Let B be a regular representation of a set of bad states.

- The RTS (Σ, I, T) is *safe* if B cannot be reached from I with a finite number of steps through the transition T
- A **safety proof** is a regular language P satisfying
 - $I \subseteq P$
 - $P \cap B = \emptyset$
 - $T(P) \subseteq P$
- A regular transition system is safe iff it has a safety proof

Example: the Collatz transition system

The Collatz system applies the following operation on natural numbers:

- If the number is even, divide it by two.
- If the number is odd, triple it and add one.

We can specify the Collatz transition system as an RTS (Σ, I, T) by encoding natural numbers in binary with the least significant bit first without trailing zeros.

Consider the safety property: “Every sequence starting from a power of 2 will reach no odd numbers but 1.”

We set $I := 0^*1\#^*$ as the initial states and $B := 1(0 + 1)(0 + 1)^*\#^*$ as the bad states. Observe that $T(\underbrace{0 \cdots 0}_{n \text{ zeros}}1\#^*) = \underbrace{0 \cdots 0}_{n-1 \text{ zeros}}1\#^*$ for each $n \geq 1$.

We therefore have $I \cap B = \emptyset$ and $T(I) \subseteq I$. Namely, I is itself a safety proof.

Regular model checking

- The **regular model checking** problem is to find a *regular* safety proof for a regular transition system
- A RTS may not have a regular proof even if it is safe!
 - A safe RTS always has a safety proof, but it can happen that all proofs are not regular. Example: Turing machines
- For some subclass of RTSs, a regular proof is guaranteed to exist when the system is safe
 - For example, the set of reachable states is regular for RTSs like Petri nets, pushdown systems, and lossy-channel systems
 - Such systems have a regular safety proof whenever they are safe
- If a system have a regular safety proof iff it is safe, then it is decidable to check safety of the system.

Regular model checking (cont'd)

- If an RTS has a regular proof when it is safe, then safety checking of the system is decidable. Idea: launch two procedures as follows at the same time

Procedure A:

```
while true do
   $i := 1$ 
  let  $A_i$  be the  $i$ -th DFA, and let  $P := L_A$ 
  if  $I \subseteq P$  and  $P \cap B = \emptyset$  and  $T(P) \subseteq P$  then
    terminate and report “safe”
   $i := i + 1$ 
```

Procedure B:

```
while true do
   $i := 0$ 
  if  $B$  is reachable from  $I$  in  $i$  step then
    terminate and report “unsafe”
   $i := i + 1$ 
```

- Eventually one of the two procedures will terminate!

Learning proofs for regular model checking

- In the rest of this lecture, we will look at two methods to find a regular proof for an RTS:
 - SAT-based learning
 - L^* -based learning
- The SAT-based method is less scalable (i.e. it is not effective when all regular proofs are large). However, it has the same termination guarantee as brute-force enumeration.
- The L^* -based method is more scalable and is capable of finding very large regular proofs in practice. However, it is not guaranteed to find a regular proof even if one exists.

Learning proofs for regular model checking

- In the rest of this lecture, we will look at two methods to find a regular proof for an RTS:
 - SAT-based learning
 - L^* -based learning
- The SAT-based method is less scalable (i.e. it is not effective when all regular proofs are large). However, it has the same termination guarantee as brute-force enumeration.
- The L^* -based method is more scalable and is capable of finding very large regular proofs in practice. However, it is not guaranteed to find a regular proof even if one exists.

SAT-based learning for safety proofs

Fix a regular system (Σ, I, T) and a set of bad states B

For each $n \geq 1$, we construct a Boolean formula Φ_n such that a model of Φ_n corresponds to a DFA A of n states and vice versa

SAT-based learning of regular proofs:

```
 $n := 1, C := \emptyset$   
while true do  
  construct  $\Phi_n$   
  while  $\Phi_n \wedge \Phi_C$  has a model  $\alpha$  do  
    construct a DFA  $A$  from  $\alpha$   
    if  $L_A$  is a safety proof then  
      return  $A$   
    let  $cex$  be a witness of the violation  
     $C := C \cup \{cex\}$   
   $n := n + 1$ 
```

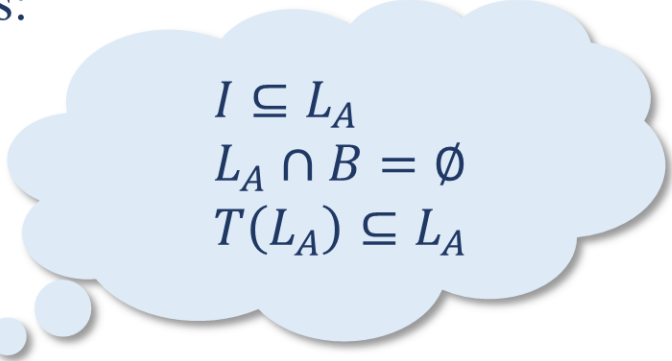
SAT-based learning for safety proofs

Fix a regular system (Σ, I, T) and a set of bad states B

For each $n \geq 1$, we construct a Boolean formula Φ_n such that a model of Φ_n corresponds to a DFA A of n states and vice versa

SAT-based learning of regular proofs:

```
 $n := 1, C := \emptyset$   
while true do  
  construct  $\Phi_n$   
  while  $\Phi_n \wedge \Phi_C$  has a model  $\alpha$  do  
    construct a DFA  $A$  from  $\alpha$   
    if  $L_A$  is a safety proof then  
      return  $A$   
    let  $cex$  be a witness of the violation  
     $C := C \cup \{cex\}$   
   $n := n + 1$ 
```


$$\begin{aligned} I &\subseteq L_A \\ L_A \cap B &= \emptyset \\ T(L_A) &\subseteq L_A \end{aligned}$$

SAT-based learning for safety proofs

Fix a regular system (Σ, I, T) and a set of bad states B

For each $n \geq 1$, we construct a Boolean formula Φ_n such that a model of Φ_n corresponds to a DFA A

SAT-based learning

```
n := 1, C := ∅
while true do
  construct  $\Phi_n$ 
  while  $\Phi_n \wedge \Phi_C$  has a model  $\alpha$  do
    construct a DFA  $A$  from  $\alpha$ 
    if  $L_A$  is a safety proof then
      return  $A$ 
    let  $cex$  be a witness of the violation
     $C := C \cup \{cex\}$ 
  n := n + 1
```

$\alpha \models \Phi_C$ iff for all $c \in C$,
 c is not a witness of A_α
violating the proof rules

SAT encoding of DFA

Encoding of a DFA $(\Sigma, S, s_0, \delta, F)$

- Given Σ and S , it suffices to fix s_0 and define only δ and F .
- For each $i, j \in S$ and $a \in \Sigma$, we define a Boolean variable $t_{i,a,j}$ such that “ $t_{i,a,j}$ is true” corresponds to “ $\delta(i, a) = j$ ”.
- For each $i \in S$, we define a Boolean variable f_i such that “ f_i is true” corresponds to “ $i \in F$ ”.
- We use the following constraint to ensure that the DFA is deterministic and complete:

$$\left(\bigwedge_{i,j,k \in S, j \neq k, a \in \Sigma} \neg(t_{i,a,j} \wedge t_{i,a,k}) \right) \wedge \left(\bigwedge_{i \in S, a \in \Sigma} \bigvee_{j \in S} t_{i,a,j} \right)$$

SAT encoding of DFA (cont'd)

Encoding of a DFA $(\Sigma, S, s_0, \delta, F)$

- For each $n \geq 1$, we define a propositional formula $\phi_{\text{DFA}}^n(\bar{t}, \bar{f})$ as

$$\left(\bigwedge_{1 \leq i, j, k \leq n, j \neq k, a \in \Sigma} \neg(t_{i,a,j} \wedge t_{i,a,k}) \right) \wedge \left(\bigwedge_{1 \leq i \leq n, a \in \Sigma} \bigvee_{1 \leq j \leq n} t_{i,a,j} \right)$$

with free variables

$$\{ t_{i,a,j} : 1 \leq i, j \leq n, a \in \Sigma \} \text{ and } \{ f_i : 1 \leq i \leq n \}.$$

- Any $\alpha \models \phi_{\text{DFA}}^n(\bar{t}, \bar{f})$ corresponds to a DFA $A_\alpha := (\Sigma, S, s_0, \delta, F)$:
 - $S = \{1, \dots, n\}$, $s_0 = 1$
 - For $i \in S$ and $a \in \Sigma$, $\delta(i, a) = j$ iff $\alpha(t_{i,a,j}) = \text{true}$
 - $F = \{i : \alpha(f_i) = \text{true}\}$

Counterexample refinement

SAT-based learning of regular proofs:

```
 $n := 1, C := \emptyset$   
while true  
  construct  $\Phi_n$   
  while  $\Phi_n \wedge \Phi_C$  has a model  $\alpha$   
    construct a DFA  $A$  from  $\alpha$   
    if  $L_A$  is a safety proof then  
      return  $A$   
    let cex be a witness of the violation  
     $C := C \cup \{cex\}$   
 $n := n + 1$ 
```

Counterexample refinement (cont'd)

- **Positive counterexample**

- A positive cex is a string supposed to be accepted by A .
- We obtain a positive cex $w \in I \setminus L_A$ when $I \not\subseteq L_A$.

$$I \subseteq L_A$$

$$L_A \cap B = \emptyset$$

$$T(L_A) \subseteq L_A$$

- **Negative counterexample**

- A negative cex is a string not supposed to be accepted by A .
- We obtain a negative cex $w \in L_A \cap B$ when $L_A \cap B \neq \emptyset$.

- **Implication counterexample**

- An implication cex is a pair of strings (w, w') such that “ w is in L_A ” implies “ w' is in L_A ”
- We obtain an implication cex when $T(L_A) \not\subseteq L_A$. In such case, we can find a pair of strings (w, w') such that $w \in L_A$ and $w' \in T(w) \setminus L_A$.

SAT encoding of positive counterexample

Encoding the membership of a string

- Suppose we got a positive counterexample w
- We give a formula ϕ_w^n such that

if $\alpha \models \phi_{\text{DFA}}^n \wedge \phi_w^n$, then A_α accepts w

- We introduce variables $\{v_{k,i} : 0 \leq k \leq |w|, 1 \leq i \leq n\}$ and let

$$\begin{aligned} \phi_w^n := & v_{0,1} \wedge \left(\bigwedge_{1 \leq k \leq |w|} \bigvee_{1 \leq i \leq n} v_{k,i} \right) \wedge \left(\bigwedge_{1 \leq i \leq n} (v_{|w|,i} \Rightarrow f_i) \right) \\ & \wedge \left(\bigwedge_{1 \leq k \leq |w|} \bigwedge_{1 \leq i,j \leq n} (v_{k-1,i} \wedge v_{k,j} \Rightarrow t_{i,a_k,j}) \right) \end{aligned}$$

Intuitively, $\alpha(v_{k,i}) = \text{true}$ iff the DFA A_α reaches state i after reading the prefix $a_1 \cdots a_k$ of the string w .

SAT encoding of negative counterexample

Encoding the non-membership of a string

- Suppose we got a negative counterexample w
- We give a formula ψ_w^n such that

if $\alpha \models \phi_{\text{DFA}}^n \wedge \psi_w^n$, then A_α **does not** accept w

- We introduce variables $\{u_{k,i} : 0 \leq k \leq |w|, 1 \leq i \leq n\}$ and let

$$\psi_w^n := u_{0,1} \wedge \left(\bigwedge_{1 \leq k \leq |w|} \bigvee_{1 \leq i \leq n} u_{k,i} \right) \wedge \left(\bigwedge_{1 \leq i \leq n} (u_{|w|,i} \Rightarrow \neg f_i) \right) \\ \wedge \left(\bigwedge_{1 \leq k \leq |w|} \bigwedge_{1 \leq i,j \leq n} (u_{k-1,i} \wedge u_{k,j} \Rightarrow t_{i,a_k,j}) \right)$$

Intuitively, $\alpha(u_{k,i}) = \text{true}$ iff the DFA A_α reaches state i after reading the prefix $a_1 \cdots a_k$ of the string w .

SAT encoding of negative counterexample

Encoding the non-membership of a string

- Suppose we got a negative counterexample w
- We give a formula ψ_w^n such that

if $\alpha \models \psi_w^n$ works only if A_α is not accept w
a complete DFA

- We introduce variables $u_{k,i}$, $1 \leq k \leq |w|$, $1 \leq i \leq n$ and let

$$\psi_w^n := u_{0,1} \wedge \left(\bigwedge_{1 \leq k \leq |w|} \bigvee_{1 \leq i \leq n} u_{k,i} \right) \wedge \left(\bigwedge_{1 \leq i \leq n} (u_{|w|,i} \Rightarrow \neg f_i) \right) \\ \wedge \left(\bigwedge_{1 \leq k \leq |w|} \bigwedge_{1 \leq i, j \leq n} (u_{k-1,i} \wedge u_{k,j} \Rightarrow t_{i,a_k,j}) \right)$$

Intuitively, $\alpha(u_{k,i}) = \text{true}$ iff the DFA A_α reaches state i after reading the prefix $a_1 \cdots a_k$ of the string w .

SAT-based learning for safety proofs

SAT-based learning with counterexample refinement:

$n := 1$, $Pos := \emptyset$, $Neg := \emptyset$, $Imp := \emptyset$

while true do

 while $\phi_{DFA}^n \wedge \Gamma_n$ has a satisfying assignment α do

 construct a DFA A from α

 if A is a safety proof then

 return A

 add a new counterexample to either Pos , Neg , or Imp

$n := n + 1$

$$\Gamma_n := \left(\bigwedge_{w \in Pos} \phi_w^n \right) \wedge \left(\bigwedge_{w \in Neg} \psi_w^n \right) \wedge \left(\bigwedge_{(w,v) \in Imp} \psi_w^n \vee \phi_v^n \right)$$

Learning proofs for regular model checking

- In the rest of this lecture, we will look at two methods to find a regular proof for an RTS:
 - SAT-based learning
 - L^* -based learning
- The SAT-based method is less scalable (i.e. it is not effective when all regular proofs are large). However, it has the same termination guarantee as the brute-force enumeration.
- The L^* -based method is more scalable and can find very large regular proofs in practice. However, it is not guaranteed to find a regular proof even if one exists.

Myhill-Nerode Theorem

Given a language $L \subseteq \Sigma^*$, we can define an equivalence relation \equiv_L over Σ^* such that $x \equiv_L y$ if and only if

$$\forall z \in \Sigma^*, xz \in L \Leftrightarrow yz \in L.$$

We will call \equiv_L the *Nerode congruence*.

Fact $x \not\equiv_L y$ if and only if there exists $z \in \Sigma^*$ such that either $xz \in L \wedge yz \notin L$, or $xz \notin L \wedge yz \in L$.

In such case, we say z is a *distinguishing string* for x and y .

Myhill-Nerode Theorem (cont'd)

Given a language $L \subseteq \Sigma^*$, we can define an equivalence relation \equiv_L over Σ^* such that $x \equiv_L y$ if and only if

$$\forall z \in \Sigma^*, xz \in L \Leftrightarrow yz \in L.$$

We will call \equiv_L the *Nerode congruence*.

Example 1

Consider $\Sigma := \{a, b\}$ and $L := (aa)^*$.

Is $\varepsilon \equiv_L aa$?

Is $a \equiv_L aa$?

Is $ab \equiv_L ba$?

What are the equivalence classes induced by \equiv_L ?

Myhill-Nerode Theorem (cont'd)

Given a language $L \subseteq \Sigma^*$, we can define an equivalence relation \equiv_L over Σ^* such that $x \equiv_L y$ if and only if

$$\forall z \in \Sigma^*, xz \in L \Leftrightarrow yz \in L.$$

We will call \equiv_L the *Nerode congruence*.

Example 1

Consider $\Sigma := \{a, b\}$ and $L := (aa)^*$.

Is $\varepsilon \equiv_L aa$?

Is $a \equiv_L aa$?

Is $ab \equiv_L ba$?

What are the equivalence classes induced by \equiv_L ?

Myhill-Nerode Theorem (cont'd)

Given a language $L \subseteq \Sigma^*$, we can define an equivalence relation \equiv_L over Σ^* such that $x \equiv_L y$ if and only if

$$\forall z \in \Sigma^*, xz \in L \Leftrightarrow yz \in L.$$

We will call \equiv_L the *Nerode congruence*.

Example 2

Consider $\Sigma := \{a, b\}$ and $L := \{a^n b^n : n \geq 0\}$.

What are the equivalence classes induced by \equiv_L ?

Myhill-Nerode Theorem (cont'd)

Given a language $L \subseteq \Sigma^*$, we can define an equivalence relation \equiv_L over Σ^* such that $x \equiv_L y$ if and only if

$$\forall z \in \Sigma^*, xz \in L \Leftrightarrow yz \in L.$$

We will call \equiv_L the *Nerode congruence*.

Myhill-Nerode Theorem

L is regular iff \equiv_L induces a finite number of equivalence classes.

Key observation

When L is regular, the set of the equivalence classes is isomorphic to the set of states of the minimal DFA that recognizes L .

Nerode congruence vs DFA

When L is regular, the set of equivalence classes induced by \equiv_L is isomorphic to the set of states of the minimal DFA that recognizes L .

Nerode congruence vs DFA (cont'd)

When L is regular, the set of equivalence classes induced by \equiv_L is isomorphic to the set of states of the minimal DFA that recognizes L .

DFA to equivalence classes

Suppose $A := (\Sigma, s_0, S, \delta, F)$ is the minimal DFA recognizing L .

Let $L_s \subseteq \Sigma^*$ be the language accepted by $A_s := (\Sigma, s_0, S, \delta, \{s\})$.

Then $\{L_s : s \in S\}$ is the set of equivalence classes induced by \equiv_L .

Nerode congruence vs DFA (cont'd)

When L is regular, the set of equivalence classes induced by \equiv_L is isomorphic to the set of states of the minimal DFA that recognizes L .

DFA to equivalence classes

Suppose $A := (\Sigma, s_0, S, \delta, F)$ is the minimal DFA recognizing L .

Let $L_s \subseteq \Sigma^*$ be the language accepted by $A_s := (\Sigma, s_0, S, \delta, \{s\})$.

Then $\{L_s : s \in S\}$ is the set of equivalence classes induced by \equiv_L .

$\{L_s : s \in S\}$ forms a partitioning of Σ^* (by determinism of A)

If $x, y \in L_s$ for some $s \in S$, then $x \equiv_L y$ (by def. of L_s and \equiv_L)

If $x \equiv_L y$, then $x, y \in L_s$ for some $s \in S$ (by minimality of A)

Nerode congruence vs DFA (cont'd)

When L is regular, the set of equivalence classes induced by \equiv_L is isomorphic to the set of states of the minimal DFA that recognizes L .

DFA to equivalence classes

Suppose $A := (\Sigma, s_0, S, \delta, F)$ is the minimal DFA recognizing L .

Let $L_s \subseteq \Sigma^*$ be the language accepted by $A_s := (\Sigma, s_0, S, \delta, \{s\})$.

Then $\{L_s : s \in S\}$ is the set of equivalence classes induced by \equiv_L .

Equivalence classes to DFA

Let $\{[x]_L : x \in \Sigma^*\}$ be the set of equivalence classes induced by \equiv_L .

Define an automaton $A_L := (\Sigma, s_0, S, \delta, F)$ as follows:

$$s_0 := [\varepsilon]_L$$

$$S := \{[x]_L : x \in \Sigma^*\}$$

$$\delta := \{([x]_L, a, [xa]_L) : x \in \Sigma^*, a \in \Sigma\}$$

$$F := \{[x]_L : x \in L\}$$

Nerode congruence vs DFA (cont'd)

When L is regular, the set of equivalence classes induced by \equiv_L is isomorphic to the set of states of the minimal DFA that recognizes L .

DFA to equivalence classes

Suppose $A := (\Sigma, s_0, S, \delta, F)$ is the minimal DFA recognizing L .

Let $L_s \subseteq \Sigma^*$ be the language accepted by $A_s := (\Sigma, s_0, S, \delta, \{s\})$.

Then $\{L_s : s \in S\}$ is the set of equivalence classes induced by \equiv_L .

Equivalence classes to DFA

Let $\{[x]_L : x \in \Sigma^*\}$ be the set of equivalence classes induced by \equiv_L .

Define an automaton $A_L := (\Sigma, s_0, S, \delta, F)$ as follows:

$$s_0 := [\varepsilon]_L$$

$$S := \{[x]_L : x \in \Sigma^*\}$$

$$\delta := \{([x]_L, a, [xa]_L) : x \in \Sigma^*, a \in \Sigma\}$$

$$F := \{[x]_L : x \in L\}$$

A_L is finite

A_L is deterministic

A_L is minimal

Nerode congruence vs DFA (cont'd)

Let $L \subseteq \Sigma^*$ be a language. Define an automaton $A_L := (\Sigma, s_0, S, \delta, F)$ as:

$$s_0 := [\varepsilon]_L$$

$$S := \{[x]_L : x \in \Sigma^*\}$$

$$\delta := \{([x]_L, a, [xa]_L) : x \in \Sigma^*, a \in \Sigma\}$$

$$F := \{[x]_L : x \in L\}$$

Key observation When L is regular, **there is a finite $D \subseteq \Sigma^*$** such that A_L is isomorphic to a DFA $A_{D,L} := (\Sigma, s_0, S, \delta, F)$ defined as follows:

$$s_0 := [\varepsilon]_{D,L}$$

$$S := \{[x]_{D,L} : x \in \Sigma^*\}$$

$$\delta := \{([x]_{D,L}, a, [xa]_{D,L}) : x \in \Sigma^*, a \in \Sigma\}$$

$$F := \{[x]_{D,L} : x \in L\}$$

Here, $y \in [x]_{D,L}$ iff x, y cannot be distinguished by any string in D w.r.t. L .

Nerode congruence vs DFA (cont'd)

Let $L \subseteq \Sigma^*$ be a language. Define an automaton $A_L := (\Sigma, s_0, S, \delta, F)$ as:

$$s_0 := [\varepsilon]_L$$

$$S := \{[x]_L : x \in \Sigma^*\}$$

$$\delta := \{([x]_L, a, [xa]_L) : x \in \Sigma^*, a \in \Sigma\}$$

$$F := \{[x]_L : x \in L\}$$

Key observation When L is regular, **there is a finite $D \subseteq \Sigma^*$** such that A_L is isomorphic to a DFA $A_{D,L} := (\Sigma, s_0, S, \delta, F)$ defined as follows:

$$s_0 := [\varepsilon]_{D,L}$$

$$S := \{[x]_{D,L} : x \in \Sigma^*\}$$

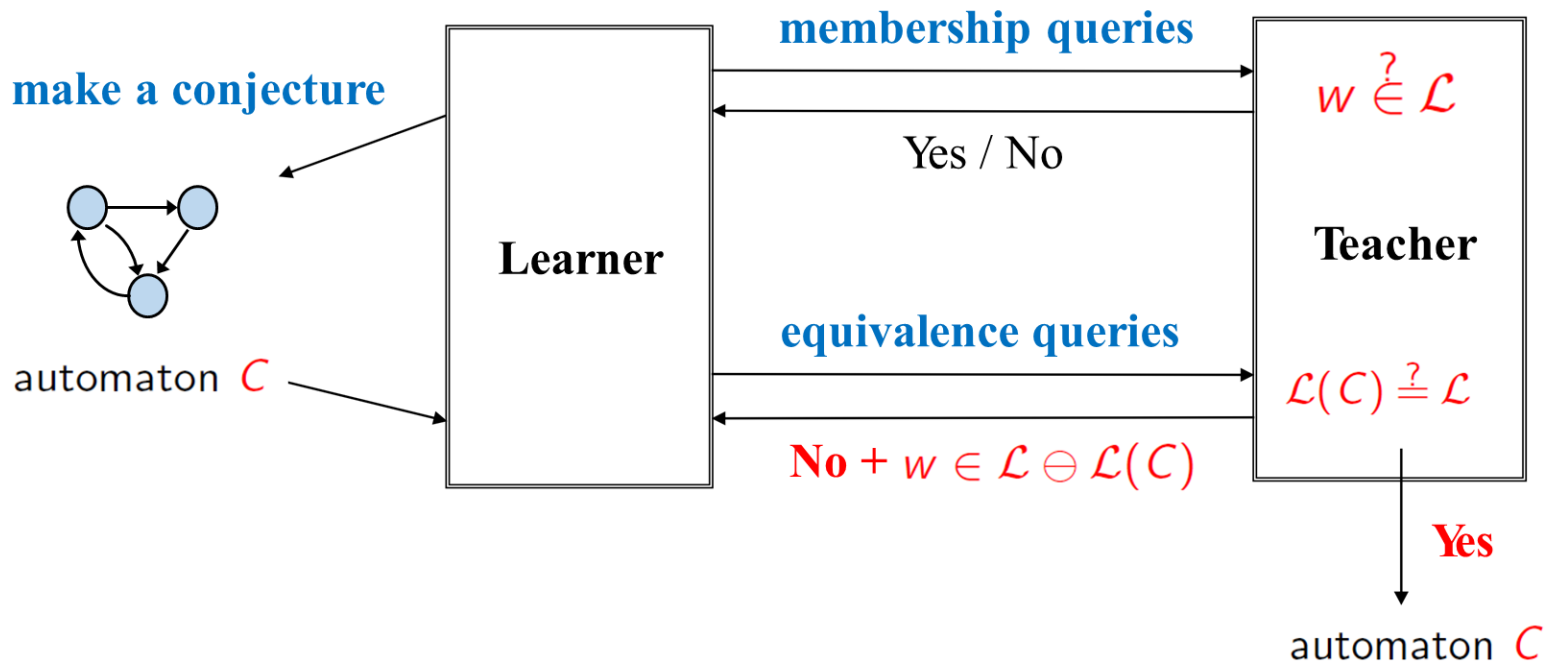
$$\delta := \{([x]_{D,L}, a, [xa]_{D,L}) : x \in \Sigma^*, a \in \Sigma\}$$

$$F := \{[x]_{D,L} : x \in L\}$$

Question: what's the relationship between $A_{D,L}$ and A_L in general?

L^* automata learning algorithm

L^* was proposed by Dana Angluin in 1987 and later improved by Rivest and Schapire in 1993. We will introduce R&S's version in this lecture.



L^* automata learning algorithm

L^* learning remains a popular approach in the AI era...

[L[^]LM: Learning Automata from Examples using Natural Language Oracles](#)

[M Vazquez-Chanlatte, K Elmaaroufi, SJ Witwicki...](#) - arXiv preprint arXiv ..., 2024 - arxiv.org

[Automata extraction from transformers](#)

[Y Zhang, Z Wei, M Sun](#) - arXiv preprint arXiv:2406.05564, 2024 - arxiv.org

[Extracting automata from recurrent neural networks using queries and counterexamples \(extended version\)](#)

[G Weiss, Y Goldberg, E Yahav](#) - Machine Learning, 2024 - Springer

[Stability analysis of various symbolic rule extraction methods from recurrent neural network](#)

[N Dave, D Kifer, CL Giles, A Mali](#) - arXiv preprint arXiv:2402.02627, 2024 - arxiv.org

[Learning minimal automata with recurrent neural networks](#)

[BK Aichernig, S König, C Mateis, A Pferscher...](#) - Software and Systems ..., 2024 – Springer

[Deepdfa: Automata learning through neural probabilistic relaxations](#)

[E Umili, R Capobianco](#) - ECAI 2024, 2024 - ebooks.iospress.nl

[Learning automata models of operator activity for human digital twin construction](#)

[F Caltabiano](#) - 2024 - politesi.polimi.it

L^* automata learning algorithm (cont'd)

Goal: learn a minimal DFA $A := (\Sigma, s_0, S, \delta, F)$ for a language L such that

$$s_0 := [\varepsilon]_L$$

$$S := \{[x]_L : x \in \Sigma^*\}$$

$$\delta := \{([x]_L, a, [xa]_L) : x \in \Sigma^*, a \in \Sigma\}$$

$$F := \{[x]_L : x \in L\}$$

L^* automata learning algorithm (cont'd)

Goal: learn a minimal DFA $A := (\Sigma, s_0, S, \delta, F)$ and $D \subseteq \Sigma^*$ for L such that

$$s_0 := [\varepsilon]_{D,L}$$

$$S := \{[x]_{D,L} : x \in \Sigma^*\}$$

$$\delta := \{([x]_{D,L}, a, [xa]_{D,L}) : x \in \Sigma^*, a \in \Sigma\}$$

$$F := \{[x]_{D,L} : x \in L\}$$

L^* automata learning algorithm (cont'd)

Goal: learn a minimal DFA $A := (\Sigma, s_0, S, \delta, F)$ and $D \subseteq \Sigma^*$ for L such that

$$s_0 := [\varepsilon]_{D,L}$$

$$S := \{[x]_{D,L} : x \in \Sigma^*\}$$

$$\delta := \{([x]_{D,L}, a, [xa]_{D,L}) : x \in \Sigma^*, a \in \Sigma\}$$

$$F := \{[x]_{D,L} : x \in L\}$$

The learner maintains an **observation table**:

		u_1	\dots	u_m
Each w is a candidate representative of state $[w]_{D,L}$	w_1	$w_1 u_1 \in? L$	\dots	$w_1 u_m \in? L$
	\vdots	\vdots		
	w_n			
Successors of the representatives: $[w]_{D,L} \xrightarrow{a} [wa]_{D,L}$	$w_1 a_1$	$w_1 a_1 u_1 \in? L$		
	\vdots	\vdots		
	$w_n a_k$			

$D := \{u_1, \dots, u_m\}$ is a set of distinguishing strings for the representatives w_1, \dots, w_n

L^* algorithm: the initial table

Fix $\Sigma := \{a, b\}$ and suppose that the target language is $L := (ab + aab)^*$.

The learner creates an initial table:

	ε
ε	
a	
b	

In the initial table, the column is indexed by ε , while the rows are indexed by $\{\varepsilon\} \cup \Sigma$.

L^* algorithm: the initial table

Fix $\Sigma := \{a, b\}$ and suppose that the target language is $L := (ab + aab)^*$.

The learner creates an initial table:

	ε
ε	T
a	F
b	F

The learner then fills the table by making membership queries.

Now we know that the state $[\varepsilon]_L$ differs from its successors $[a]_L$ and $[b]_L$.

We extend the table by adding a (or b) to the state space.

L^* algorithm: extending the table

Fix $\Sigma := \{a, b\}$ and suppose that the target language is $L := (ab + aab)^*$.

After extending the state space with a , we obtain the table

	ε
ε	T
a	F
b	F
aa	
ab	

The learner then extends the table with the successors of a, b and fills the table by making membership queries.

L^* algorithm: extending the table

Fix $\Sigma := \{a, b\}$ and suppose that the target language is $L := (ab + aab)^*$.

After extending the state space with a , we obtain the table

	ε
ε	T
a	F
b	F
aa	F
ab	T

Now every successor class has a representative in the table with respect to the current set of distinguishing strings.

We say that the table is *closed*.

L^* algorithm: making a conjecture

Fix $\Sigma := \{a, b\}$ and suppose that the target language is $L := (ab + aab)^*$.

The table is closed. We construct a DFA $A_{D,L}$ from the table with $D = \{\varepsilon\}$.

	ε
ε	T
a	F
b	F
aa	F
ab	T

$$s_0 := [\varepsilon]_{D,L}$$

$$S := \{[x]_{D,L} : x \in \Sigma^*\}$$

$$\delta := \{([x]_{D,L}, a, [xa]_{D,L}) : x \in \Sigma^*, a \in \Sigma\}$$

$$F := \{[x]_{D,L} : x \in L\}$$

L^* algorithm: making a conjecture

Fix $\Sigma := \{a, b\}$ and suppose that the target language is $L := (ab + aab)^*$.

	ε
ε	T
a	F
b	F
aa	F
ab	T

The learner then makes an equivalence query $Eq(A_D)$ to the teacher.

The teacher replies “No” and provides a counterexample $w \in L_{A_D} \ominus L$.

Then this string w contains a suffix that is a valid distinguishing string.

L^* algorithm: making a conjecture

Fix $\Sigma := \{a, b\}$ and suppose that the target language is $L := (ab + aab)^*$.

	ε
ε	T
a	F
b	F
aa	F
ab	T

The learner then makes an equivalence query $Eq(A_D)$ to the teacher.

The teacher replies “No” and provides a counterexample $w \in L_{A_D} \oplus L$.

Suppose that the teacher returns bb . Then b is a distinguishing string for at least two D -equivalent states.

L^* algorithm: the 2nd iteration

Fix $\Sigma := \{a, b\}$. Suppose that the language to learn is $L := (ab + aab)^*$.

	ε	b
ε	T	F
a	F	T
b	F	F
aa	F	T
ab	T	F

We include b in the state space and extend the table accordingly.

The representatives **a** and **b** are separated by the new distinguishing string!

L^* algorithm: the 2nd iteration

Fix $\Sigma := \{a, b\}$. Suppose that the language to learn is $L := (ab + aab)^*$.

	ε	b
ε	T	F
a	F	T
b	F	F
aa	F	T
ab	T	F
ba	F	F
bb	F	F

The table is now closed. The learner makes an equivalence query to the teacher. The teacher replies “No”, and we obtain a new distinguishing string ab after analyzing the counterexample.

L^* algorithm: the 3rd iteration

Fix $\Sigma := \{a, b\}$. Suppose that the language to learn is $L := (ab + aab)^*$.

	ε	b	ab
ε	T	F	T
a	F	T	T
b	F	F	F
aa	F	T	F
ab	T	F	T
ba	F	F	F
bb	F	F	F
aaa	F	F	F
aab	T	F	T

The learner successfully learns a minimal DFA A for L in the 3rd iteration.

L^* algorithm: counterexample analysis

Claim If the teacher returns a counterexample $w \in L(A) \ominus L$ for an equivalence query $Eq(A)$, then one can make $\log|w|$ membership queries to find a string that distinguish two states of A .

Recall that a hypothesis automaton $A := (\Sigma, s_0, S, \delta, F)$ is defined as

$$s_0 := [\varepsilon]_{D,L}$$

$$S := \{[x]_{D,L} : x \in \Sigma^*\}$$

$$\delta := \{([x]_{D,L}, a, [xa]_{D,L}) : x \in \Sigma^*, a \in \Sigma\}$$

$$F := \{[x]_{D,L} : x \in L\}$$

Consider a counterexample string $w := a_1 \dots a_m$. Then A will reach state $[a_1 \dots a_k]_{D,L}$ after reading the prefix $a_1 \dots a_k$ of w .

If $w \in L(A) \ominus L$, then there exists $1 \leq k \leq m$ such that $a_{k+1} \dots a_m$ is a distinguishing string for some $x, y \in [a_1 \dots a_k]_{D,L}$.

We can locate this k using binary search with $\log|w|$ membership queries. Adding $a_{k+1} \dots a_m$ to D will identify at least one new state.

L^* algorithm: complexity

Complexity result of L^*

If the minimal DFA recognizing the target language has n states, then

1. The learner needs at most n equivalence queries
2. The learner needs $O(|\Sigma|n^2 + n \log m)$ membership queries

where m is the maximum size of counterexample returned by the teacher.

L^* -based learning for safety proofs

We introduce below how to use the L^* algorithm to learn a safety proof for a regular transition system (Σ, I, T) .

- We need a target language for L^* . We cannot use the proof to learn as the target language since safety proof is not unique.
- Instead, we set (the language representation of) the reachable states $T^*(I)$ as the target language.
- Recall that $T^*(I)$ is unique, and is a proof when the system is safe.
- We will design a teacher for L^* such that when the system is safe and $T^*(I)$ is regular, the learner is guaranteed to find a proof.

L^* -based learning for safety proofs (cont'd)

- We set the reachable states $T^*(I)$ as the target language.
- We will design a teacher for L^* such that when the system is safe and $T^*(I)$ is regular, the learner is guaranteed to find a regular proof.
- **Resolving Mem(w):**
 - $w \in T^*(I)$ iff w is reachable from I .
- **Resolving Eq(A):**

It suffices to check the proof rules for safety:

 - $I \subseteq L_A$
 - $L_A \cap B = \emptyset$
 - $T(L_A) \subseteq L_A$

L^* -based learning: resolving equivalence query

- We check the proof rules for safety to resolve $\text{Eq}(A)$:
 - $I \subseteq L_A$
 - $L_A \cap B = \emptyset$
 - $T(L_A) \subseteq L_A$
- If any of the checks fails:
 - $I \not\subseteq L_A$: any $w \in I \setminus L_A$ is a **positive cex**
 - $L_A \cap B \neq \emptyset$: any $w \in L_A \cap B$ is a **negative cex**
 - $T(L_A) \not\subseteq L_A$: there is $w \in L_A$ and $T(w) \setminus L_A \neq \emptyset$.

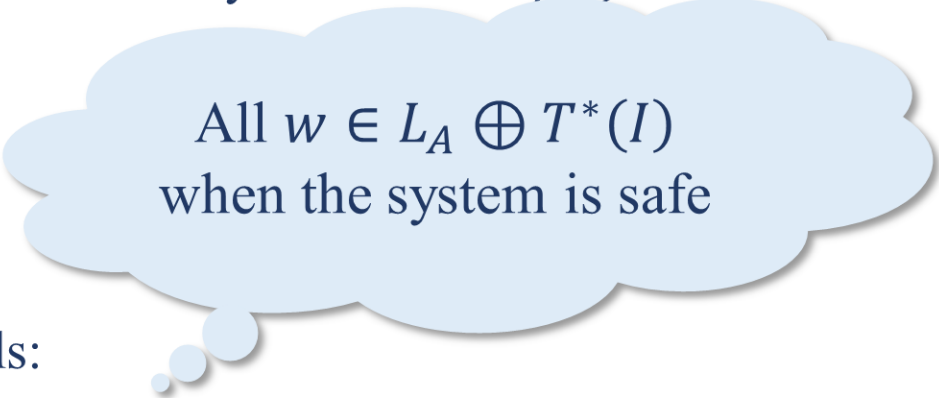
If $\text{Mem}(w)$ is “no”, then $w \notin T^*(I)$ and thus is a **negative cex**

If $\text{Mem}(w)$ is “yes”, then any $w \in T(w) \setminus L_A$ is a **positive cex**

L^* -based learning: resolving equivalence query

- We check the proof rules for safety to resolve $\text{Eq}(A)$:

- $I \subseteq L_A$
- $L_A \cap B = \emptyset$
- $T(L_A) \subseteq L_A$



All $w \in L_A \oplus T^*(I)$
when the system is safe

- If any of the checks fails:

- $I \not\subseteq L_A$: any $w \in I \setminus L_A$ is a positive cex
- $L_A \cap B \neq \emptyset$: any $w \in L_A \cap B$ is a negative cex
- $T(L_A) \not\subseteq L_A$: there is $w \in L_A$ and $T(w) \setminus L_A \neq \emptyset$.

If $\text{Mem}(w)$ is “no”, then $w \notin T^*(I)$ and thus is a negative cex

If $\text{Mem}(w)$ is “yes”, then any $w \in T(w) \setminus L_A$ is a positive cex

Active learning algorithms for DFAs (up to 2015)

	Algorithm	Publication
Angluins et al. 1987	Angluin's L^*	Learning regular sets from queries and counterexamples
Rivest and Schapire 1993	R & S 's Algorithm	Inference of Finite Automata Using Homing Sequences
Kearns and Vazirani 1994	K & V 's Algorithm	An introduction to computational learning theory
Parekh et al. 1997	ID and IID	A polynomial time incremental algorithm for regular grammar inference
Denis et al. 2001	DeLeTe2	Learning regular languages using RFSAs
Bongard et al. 2005	Estimation- Exploration	Active Coevolutionary Learning of Deterministic Finite Automata
Isberner et al. 2014	The TTT Algorithm	The TTT Algorithm: A Redundancy-Free Approach to Active Automata Learning
Volpato et al. 2015	LearnLTS	Approximate Active Learning of Nondeterministic Input Output Transition Systems

Example: Proving protocol correctness using L^* -learning

The Israeli-Jalfon's leader election protocol

1. Processes $1, \dots, n$ are organized in a ring
2. At the beginning, at least *two* processes hold a token
3. At each step, a process can pass its token to the right or left
4. When a process receives two tokens, it discards one of them

Safety condition: there is at least one token in the ring.

Example: Proving protocol correctness using L^* -learning

The Israeli-Jalfon's leader election protocol

1. Processes $1, \dots, n$ are organized in a ring
2. At the beginning, at least *two* processes hold a token
3. At each step, a process can pass its token to the right or left
4. When a process receives two tokens, it discards one of them

Safety condition: there is at least one token in the ring.

We model the protocol with an RTS (Σ, I, T) and bad states B , where

$$I: (1 + 0)^* 1 (1 + 0)^* 1 (1 + 0)^*$$

$$T: id^* \begin{bmatrix} 1 \\ 0 \end{bmatrix} \left(\begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right) id^* + id^* \left(\begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right) \begin{bmatrix} 1 \\ 0 \end{bmatrix} id^* + \left(\begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right) id^* \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} id^* \left(\begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right)$$

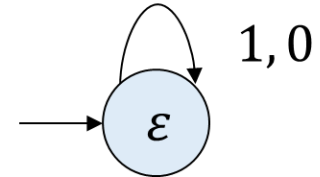
$$B: 0^*$$

$$id := \left(\begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right)$$

L^* -based learning: an example (cont'd)

$I : (1 + 0)^* 1 (1 + 0)^* 1 (1 + 0)^*$

$B : 0^*$



The first closed table:

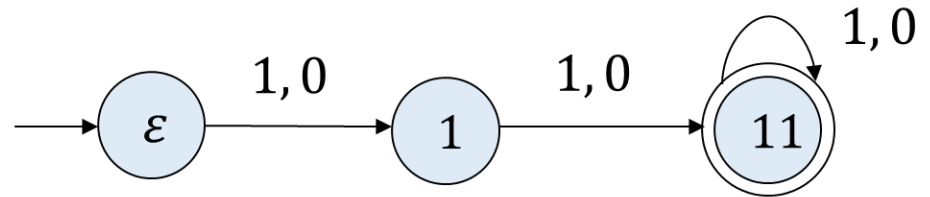
	ϵ
ϵ	F
1	F
0	F

Counterexample: $11 \in I \setminus L_A$. Add a new distinguishing string 1.

L^* -based learning: an example (cont'd)

$I : (1 + 0)^* 1 (1 + 0)^* 1 (1 + 0)^*$

$B : 0^*$



The second closed table:

	ϵ	1
ϵ	F	F
1	F	T
11	T	T
0	F	T
10	T	T
111	T	T
110	T	T

Counterexample: $000 \in L_A \cap B$. Add a new distinguishing string 0.

L^* -based learning: an example (cont'd)

$I : (1 + 0)^* 1 (1 + 0)^* 1 (1 + 0)^*$

$B : 0^*$

The third closed table leads to a regular proof. What is the DFA?

	ϵ	1	0
ϵ	F	F	F
1	F	T	T
0	F	T	F
11	T	T	T
10	T	T	T
01	T	T	T
00	F	T	F
111	T	T	T
110	T	T	T

Tutorial: Specifying and verifying regular properties

The MU Puzzle

You start with the string “MI”. Apply any of the following rules to the current string w:

1. If the last letter is I, append U. Example: $MI \rightarrow MIU$
2. Duplicate the substring after M. Example: $MIU \rightarrow MIUIU$
3. Replace substring III with U. Example: $MUIIIU \rightarrow MUUU$
4. Remove substring UU. Example: $MUUU \rightarrow MU$

Question: can you generate MU?

- 1 Start: "MI"
- 2 $MI \rightarrow MIU$
- 3 $MIU \rightarrow MIUIU$
- 4 $MUIIUU \rightarrow MUUUU$
- 5 $MUUUU \rightarrow MU$
- 6 Bad: "MU"

1. $\forall s. s = \text{"MI"} \rightarrow \text{Inv}(s)$
2. $\forall s, s', x. \text{Inv}(s) \wedge s = x \cdot \text{"I"} \wedge s' = s \cdot \text{"U"} \rightarrow \text{Inv}(s')$
3. $\forall s, s', x. \text{Inv}(s) \wedge s = \text{"M"} \cdot x \wedge s' = \text{"M"} \cdot x \cdot x \rightarrow \text{Inv}(s')$
4. $\forall s, s', x, y. \text{Inv}(s) \wedge s = x \cdot \text{"III"} \cdot y \wedge s' = x \cdot \text{"U"} \cdot y \rightarrow \text{Inv}(s')$
5. $\forall s, s', x, y. \text{Inv}(s) \wedge s = x \cdot \text{"UU"} \cdot y \wedge s' = x \cdot y \rightarrow \text{Inv}(s')$
6. $\forall s. \text{Inv}(s) \wedge s = \text{"MU"} \rightarrow \perp$