

New Challenges for Automata Learning

Property-Directed Verification of Recurrent Neural Networks

Benedikt Bollig

CNRS, LSV, ENS Paris-Saclay, Université Paris-Saclay

joint work with **LeaRNNify** team:

Benoît Barbot, Alain Finkel, Serge Haddad, Igor Khmelnitsky, Daniel Neider,
Martin Leucker, Rajarshi Roy, Lina Ye

DATAIA Workshop "Safety & AI"
September 23, 2020
CentraleSupélec

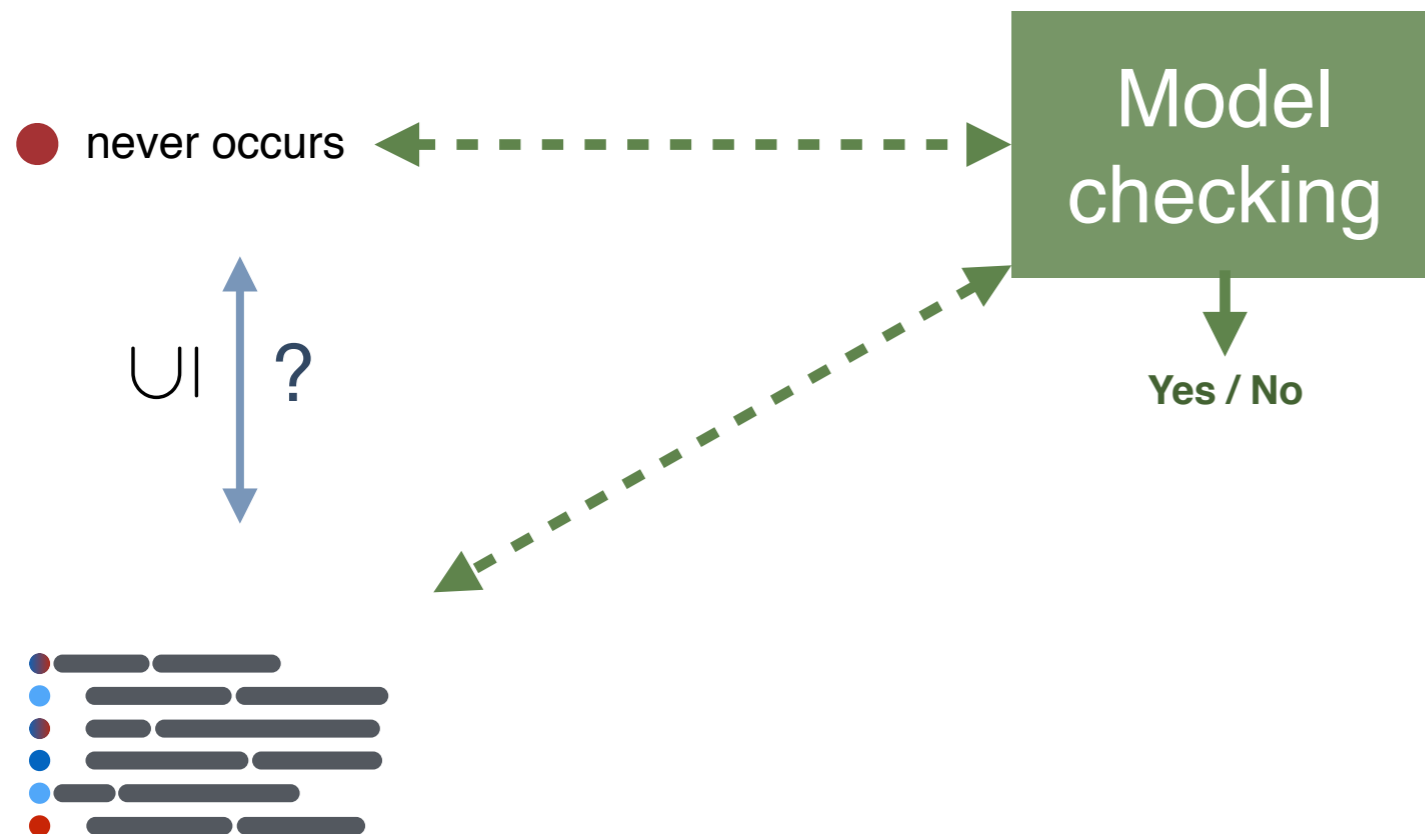


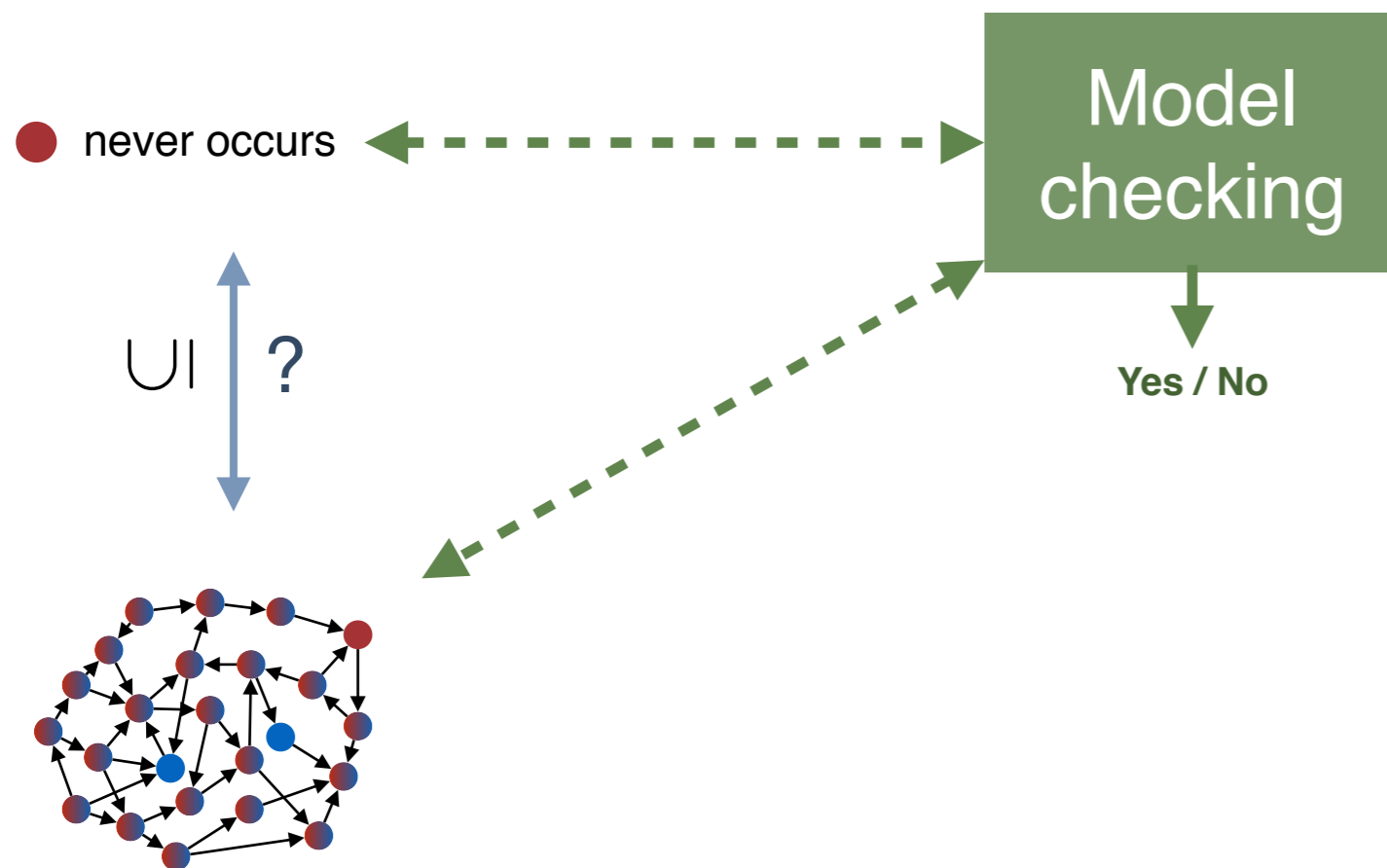
● never occurs

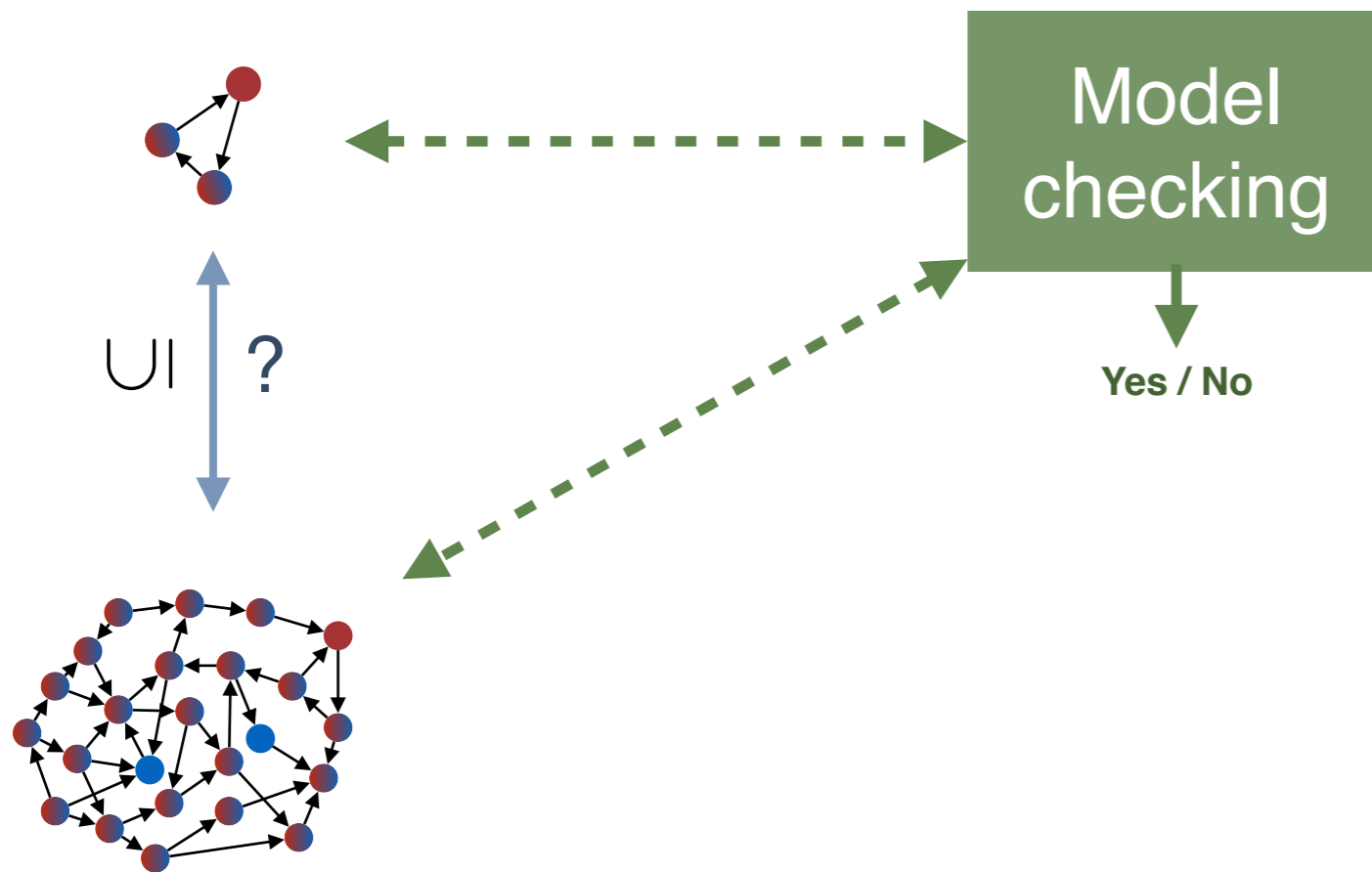


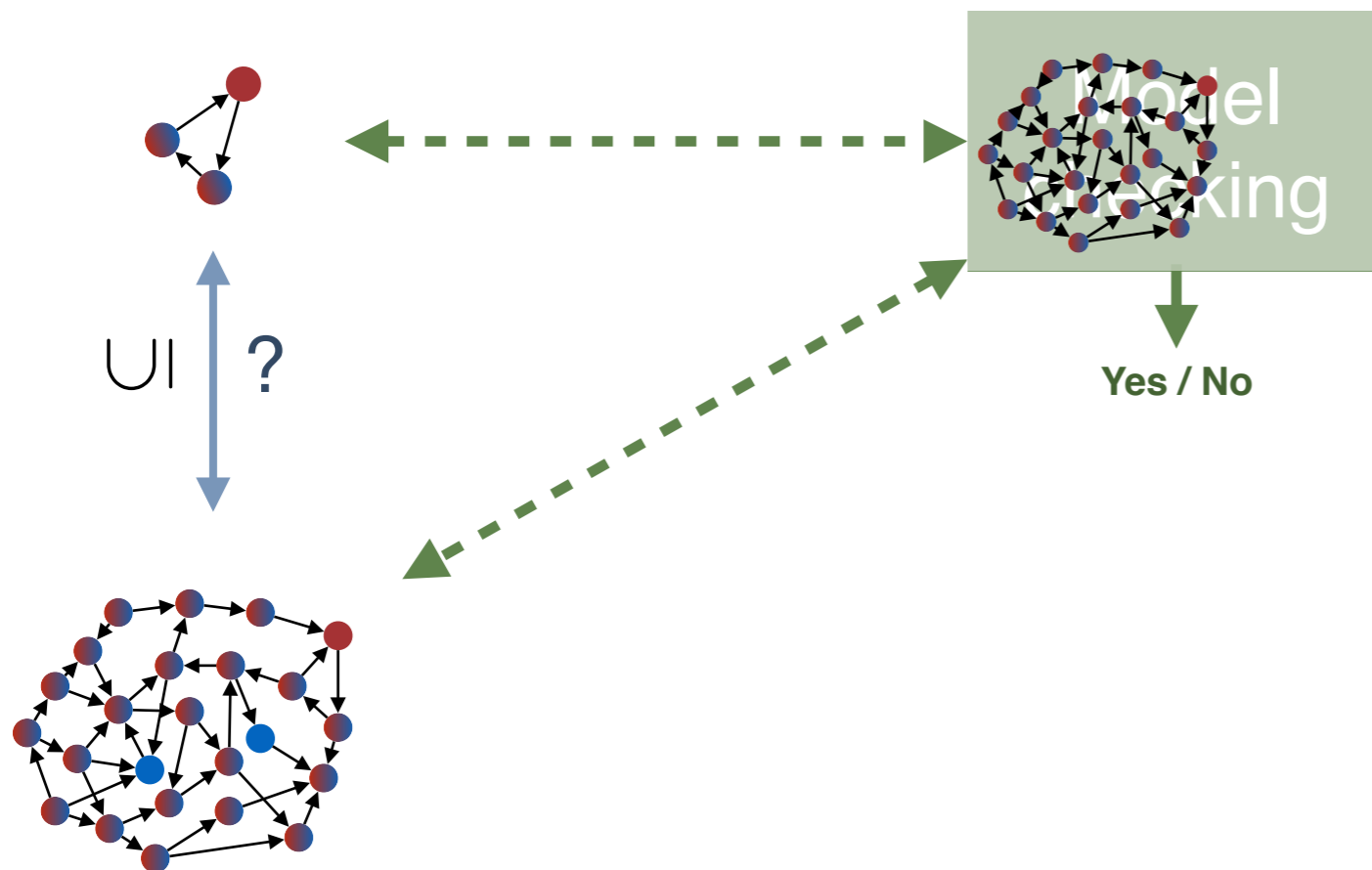
● never occurs

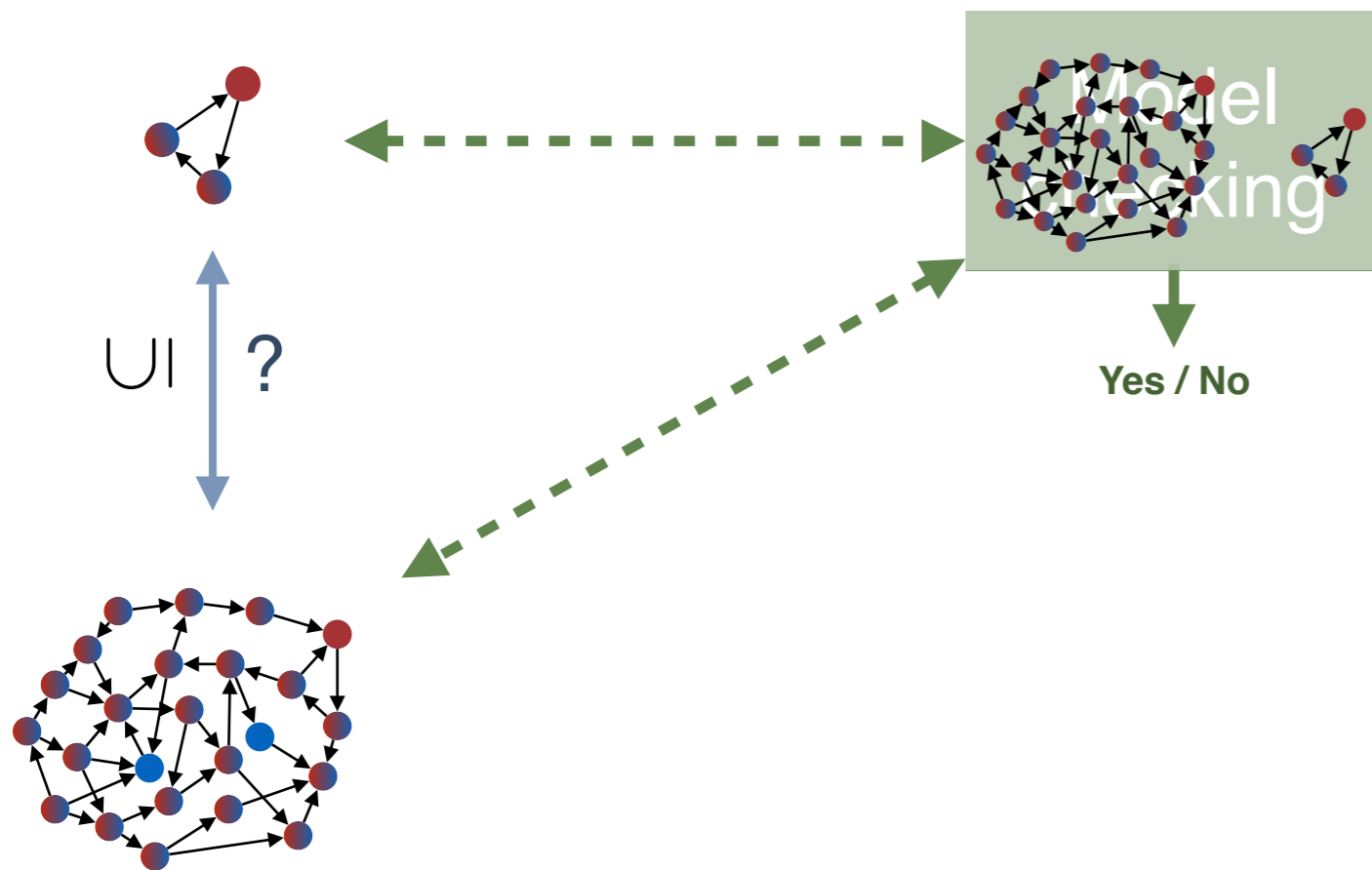


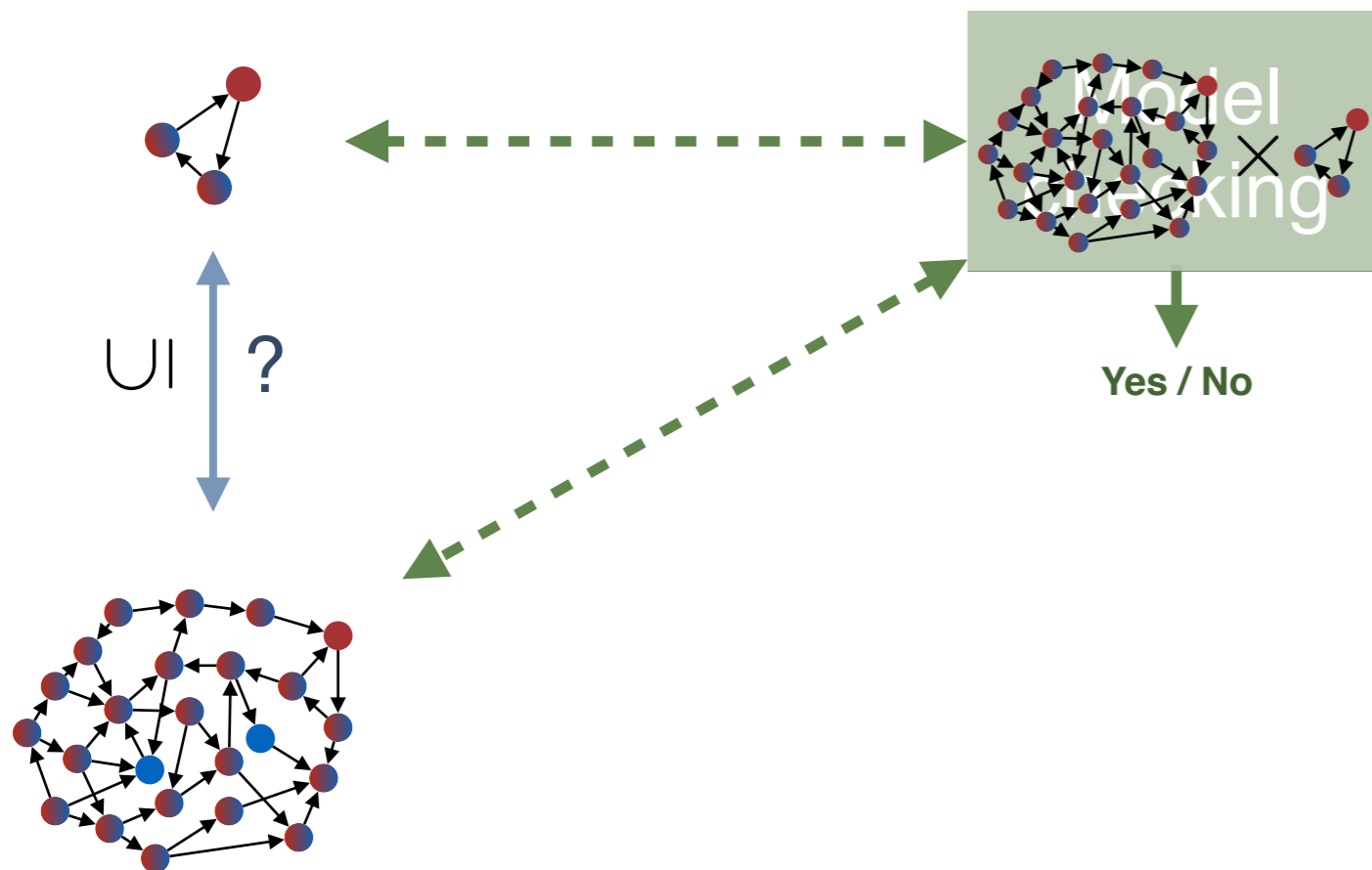


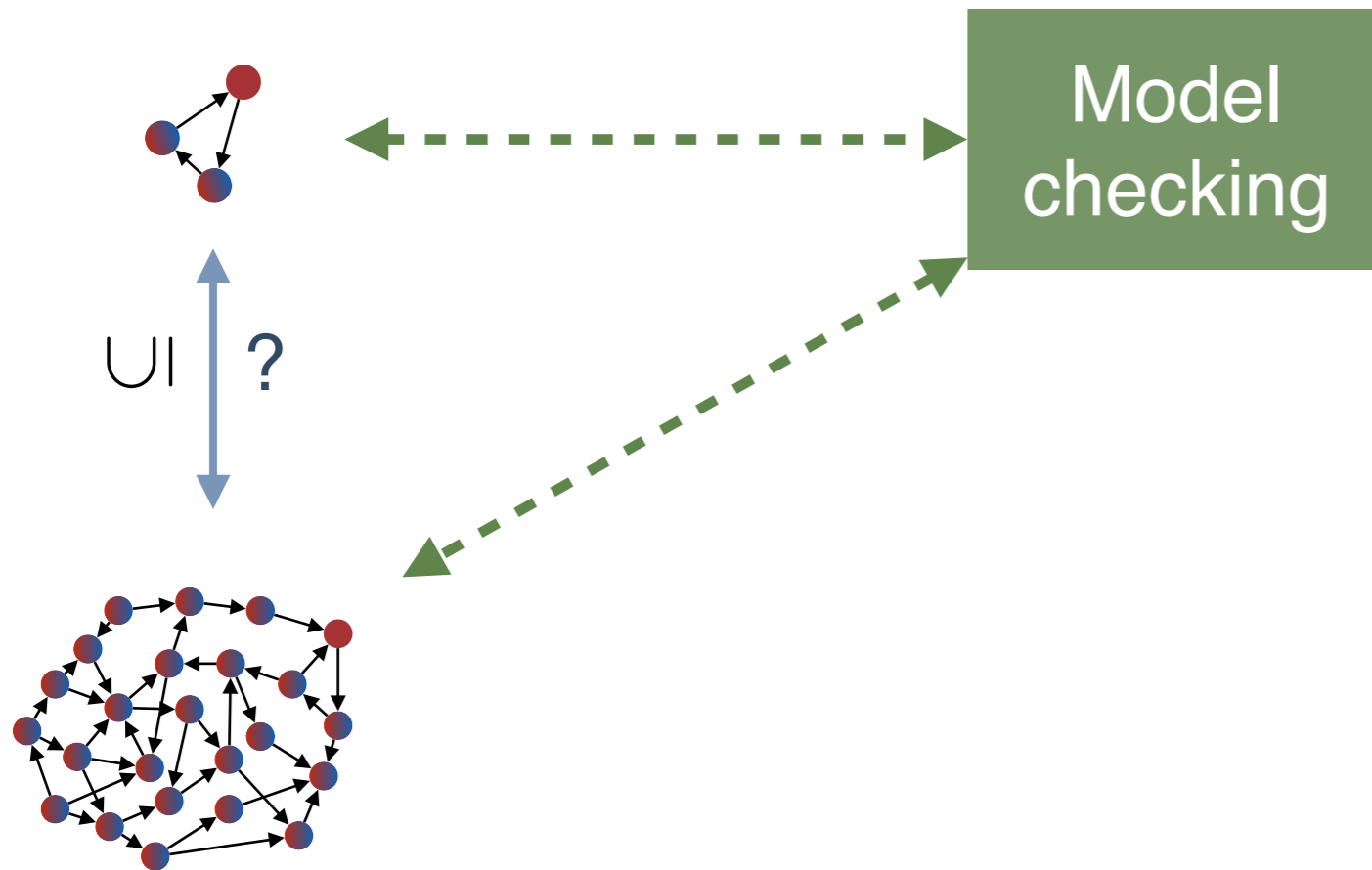






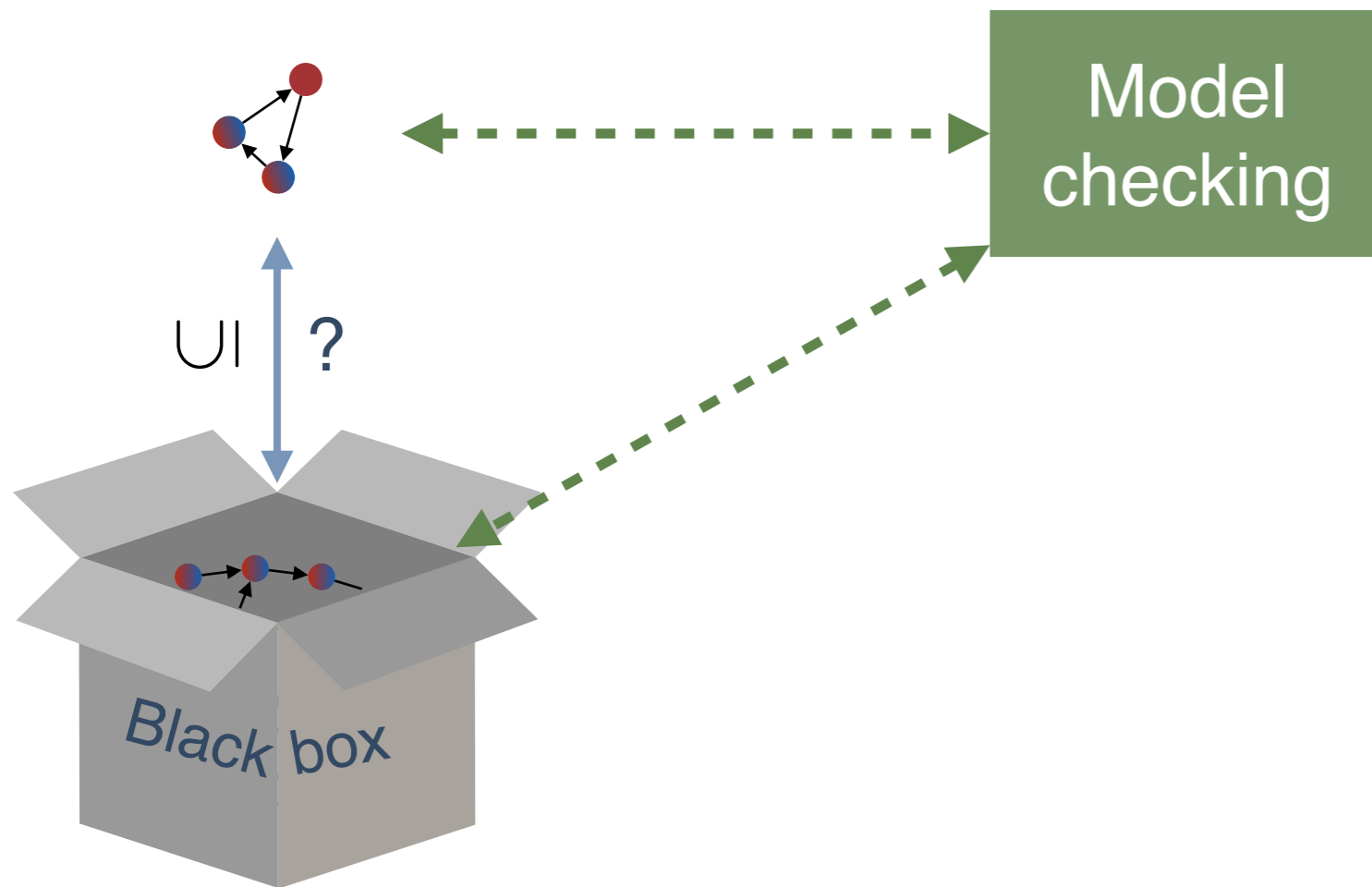






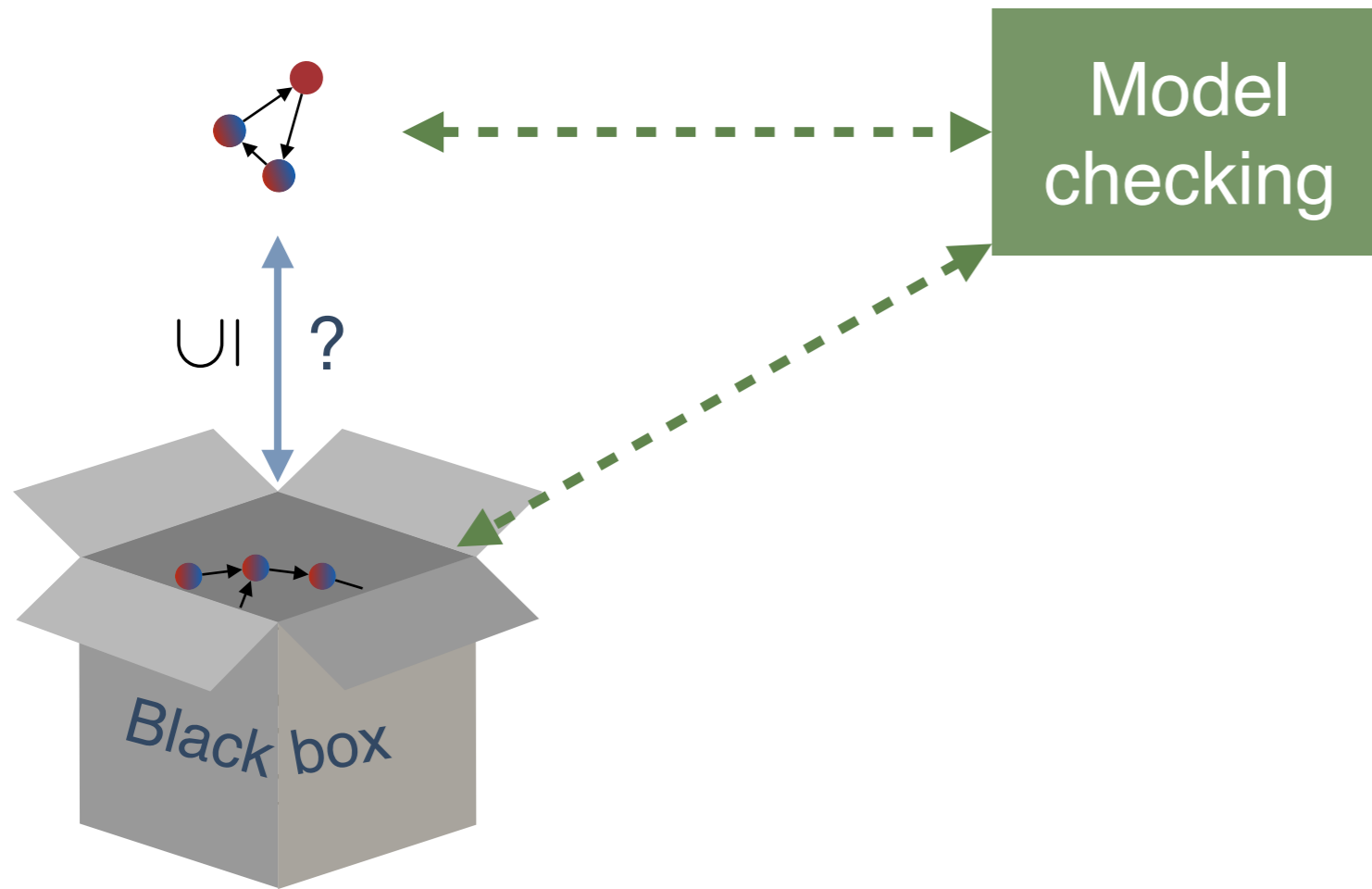
Model Checking

[Clarke-Emerson, Queille-Sifakis '80s]



Model Checking

[Clarke-Emerson, Queille-Sifakis '80s]

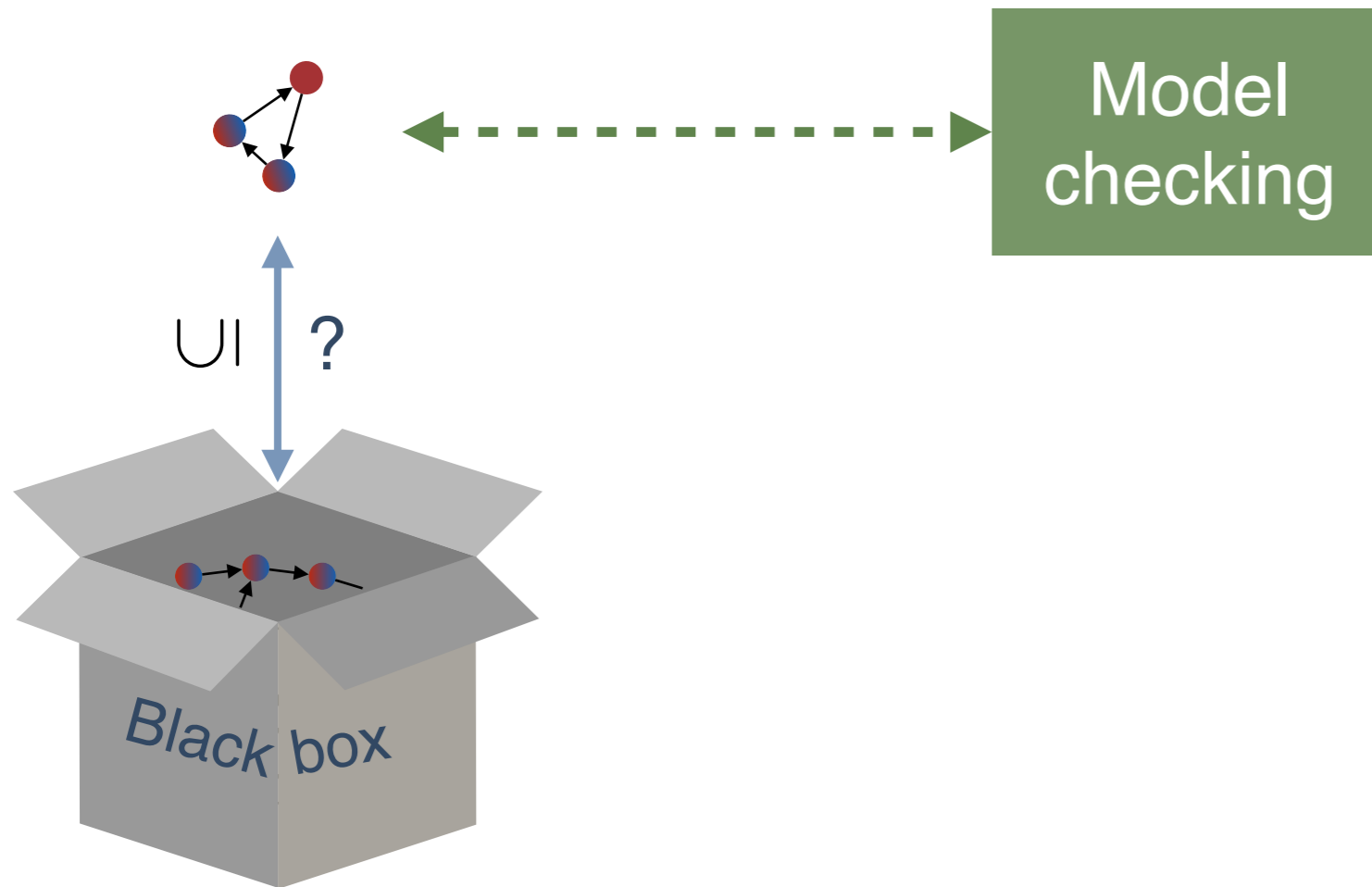


- legacy software
- code not open source
- third-party software
- embedded systems



Model Checking

[Clarke-Emerson, Queille-Sifakis '80s]

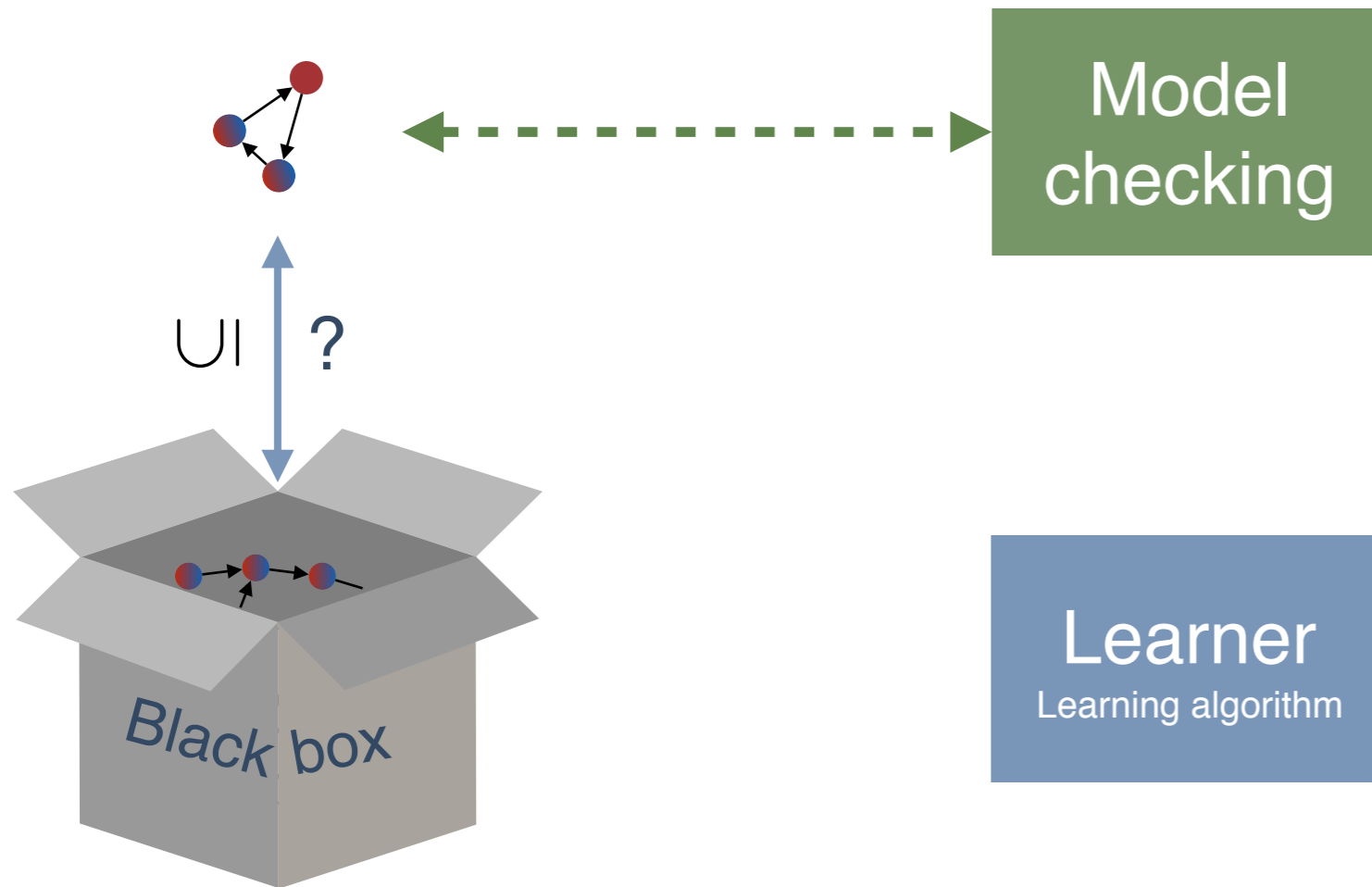


- legacy software
- code not open source
- third-party software
- embedded systems



Model Checking

[Clarke-Emerson, Queille-Sifakis '80s]

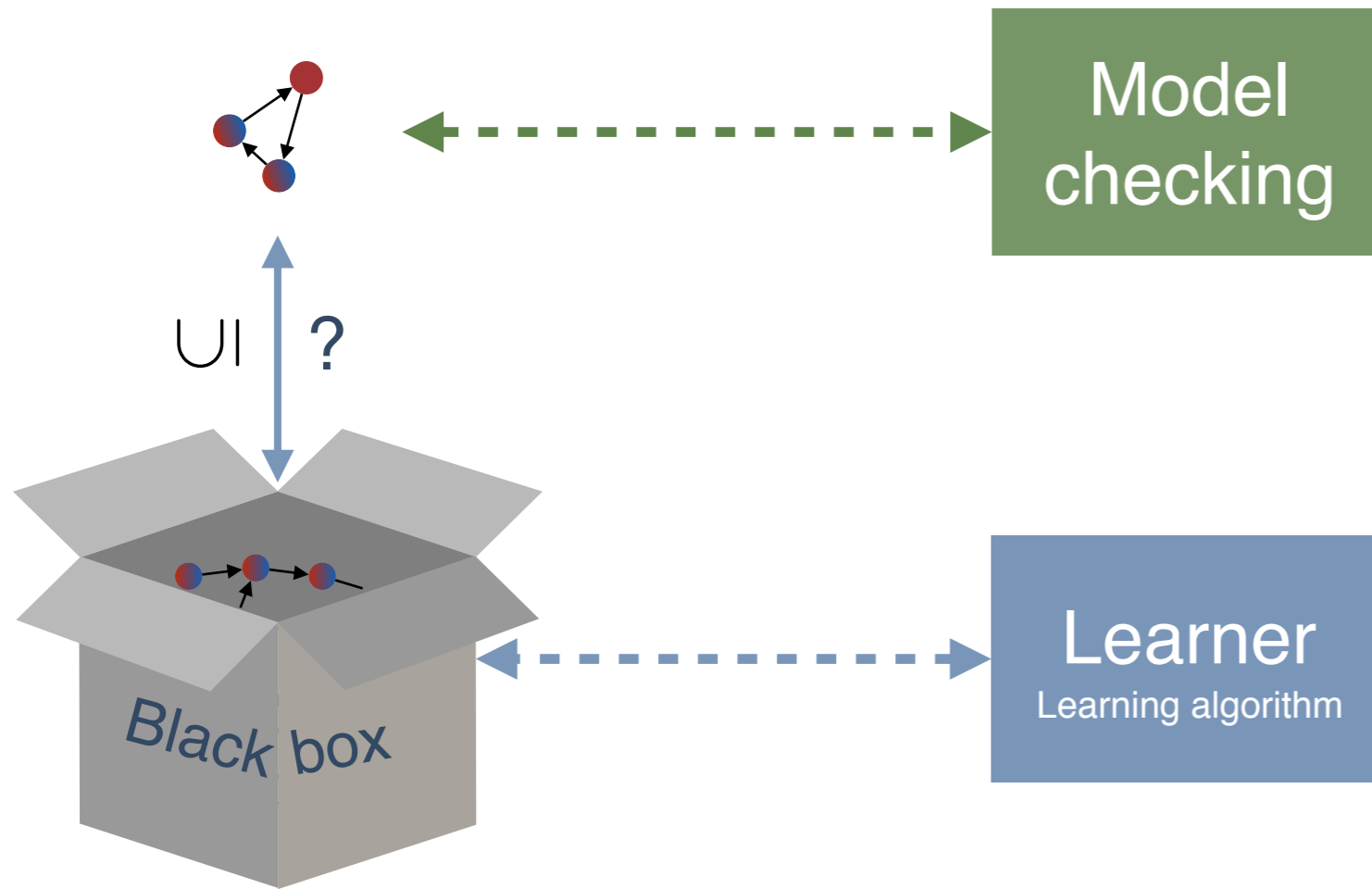


- legacy software
- code not open source
- third-party software
- embedded systems



Model Checking

[Clarke-Emerson, Queille-Sifakis '80s]

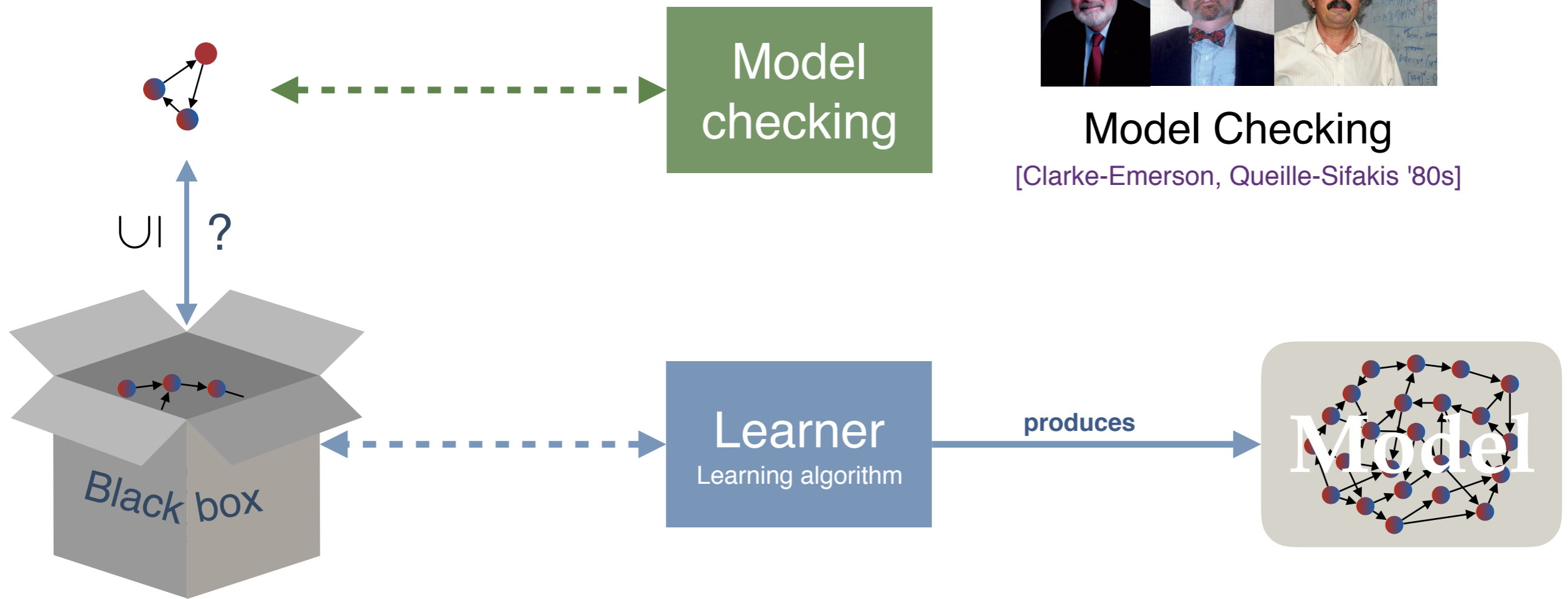


- legacy software
- code not open source
- third-party software
- embedded systems



Model Checking

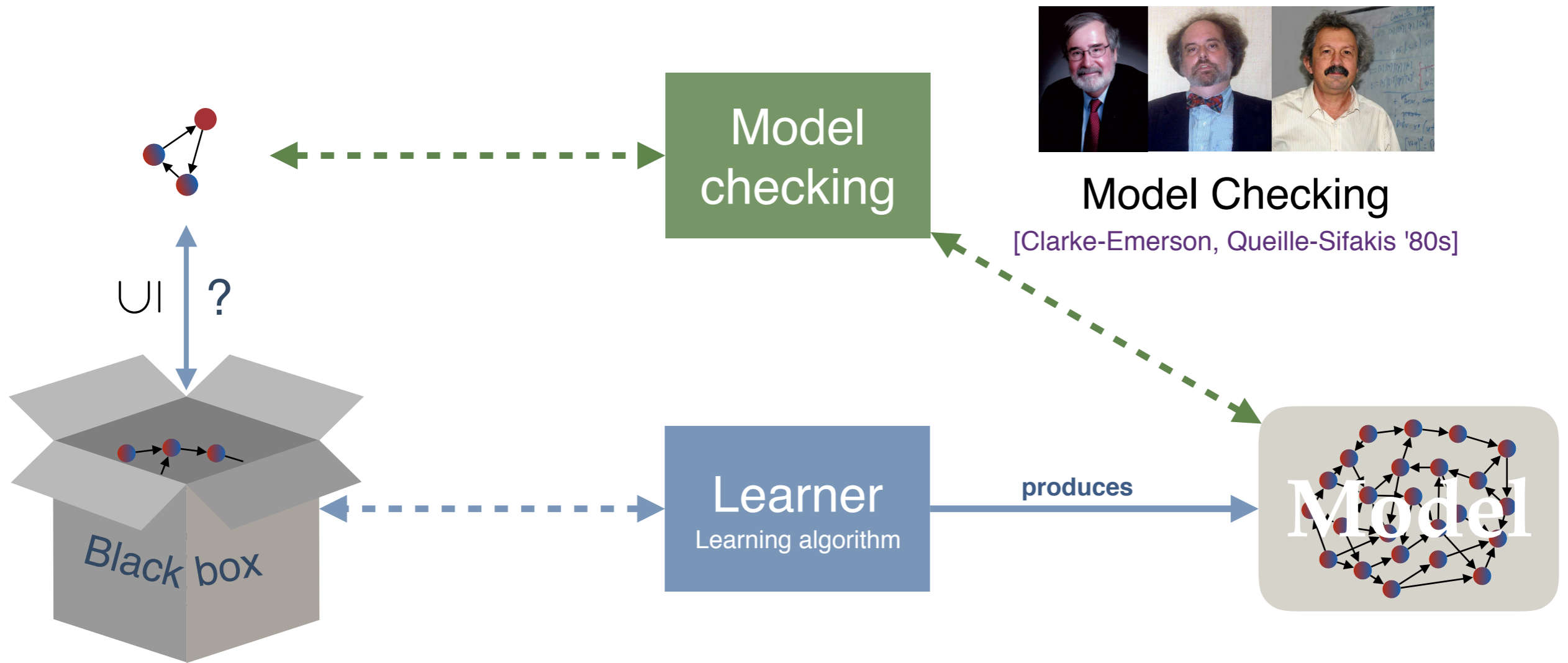
[Clarke-Emerson, Queille-Sifakis '80s]



Model Checking

[Clarke-Emerson, Queille-Sifakis '80s]

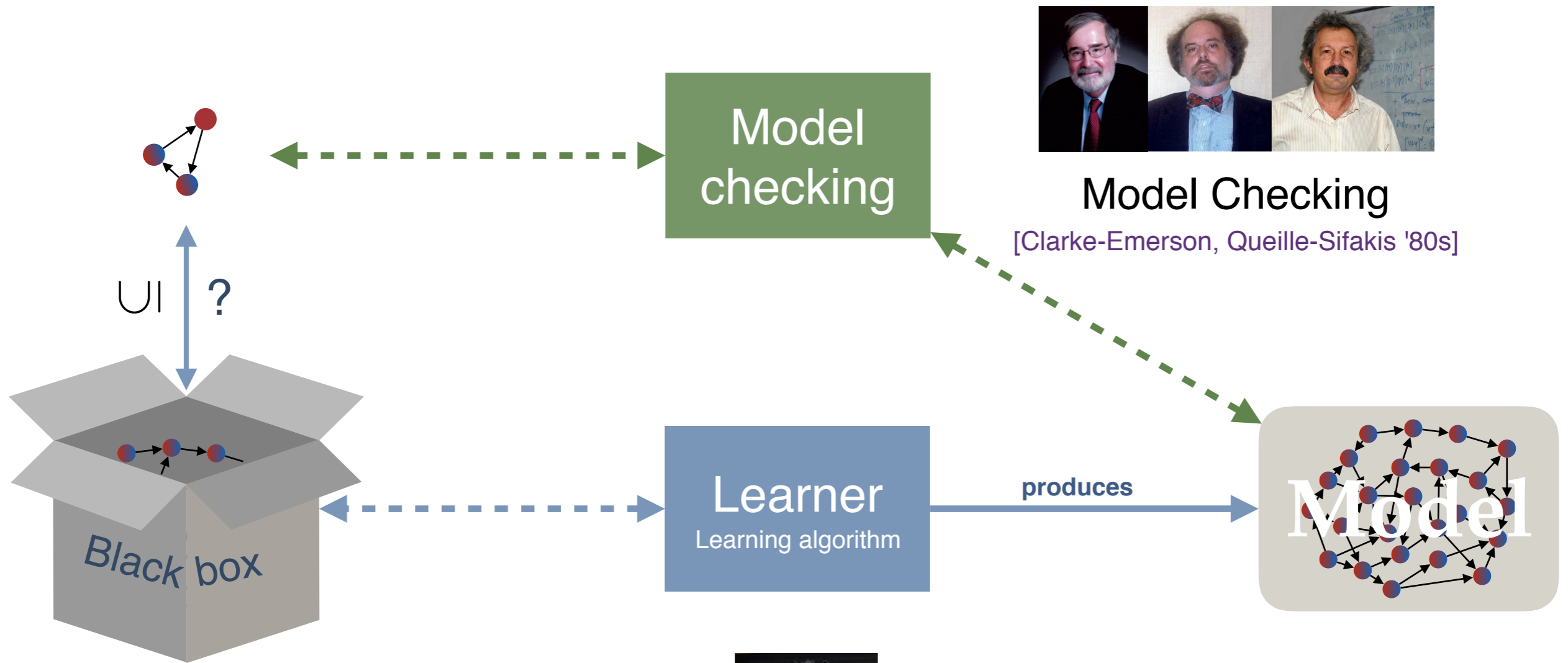
- legacy software
- code not open source
- third-party software
- embedded systems



Model Checking


[Clarke-Emerson, Queille-Sifakis '80s]

- legacy software
- code not open source
- third-party software
- embedded systems



Model Checking

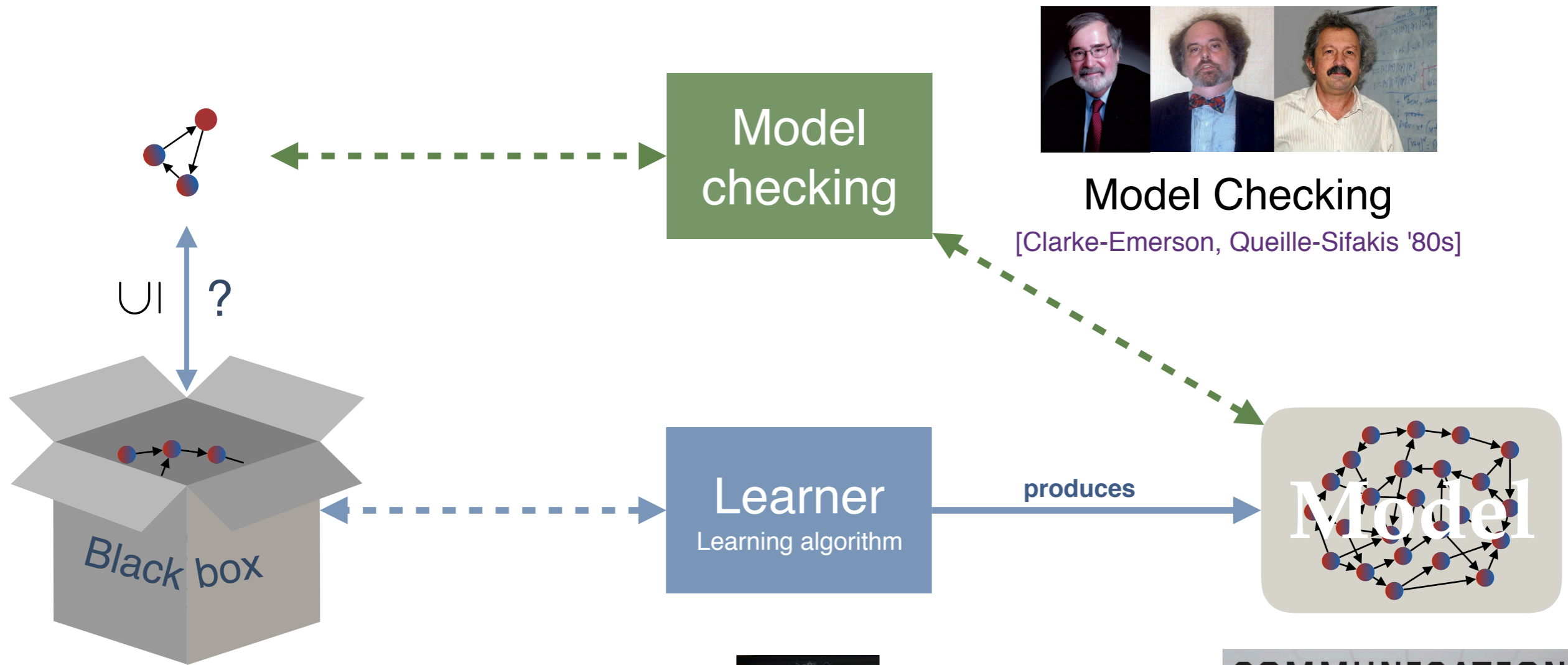
[Clarke-Emerson, Queille-Sifakis '80s]



Inference/Model Learning

[Moore '50s, Angluin '80s]

- legacy software
- code not open source
- third-party software
- embedded systems



Model Checking

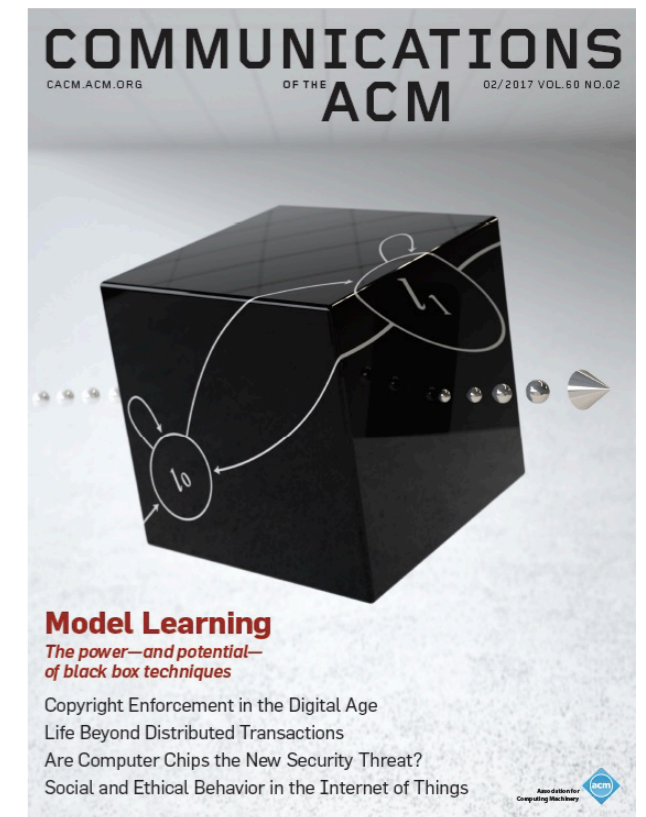
[Clarke-Emerson, Queille-Sifakis '80s]

- legacy software
- code not open source
- third-party software
- embedded systems



Inference/Model Learning

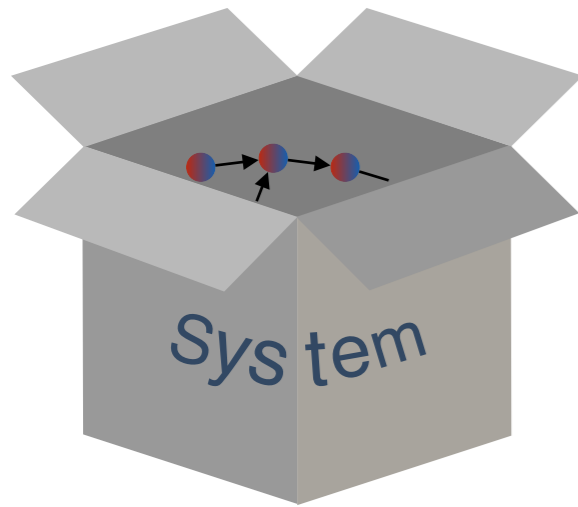
[Moore '50s, Angluin '80s]



[F. Vaandrager 2017]

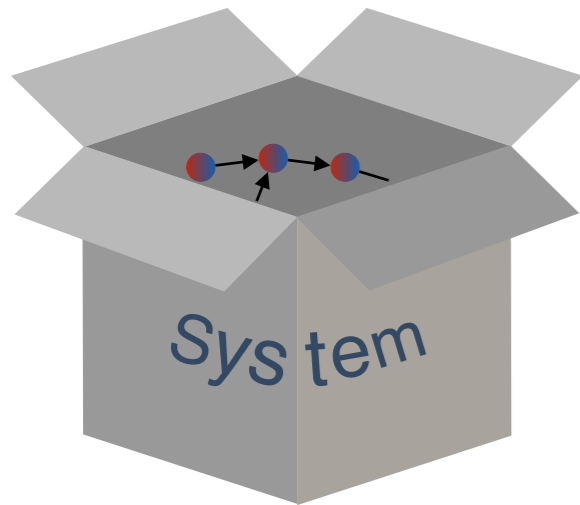
Angluin's L^* Algorithm

[Angluin 1987]



Angluin's L^* Algorithm

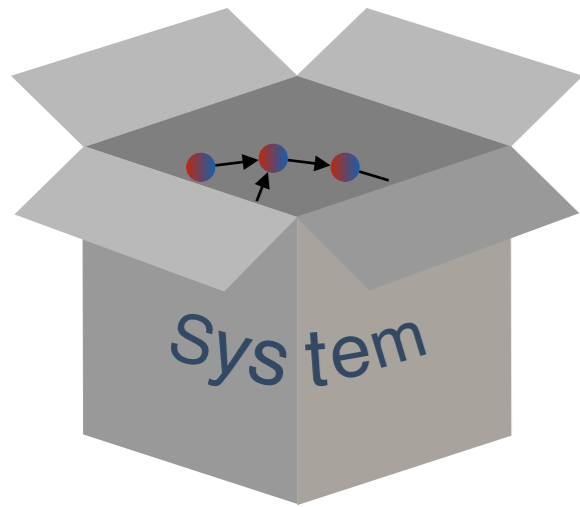
[Angluin 1987]



$$L = (a + b)^* b(a + b)$$

Angluin's L^* Algorithm

[Angluin 1987]



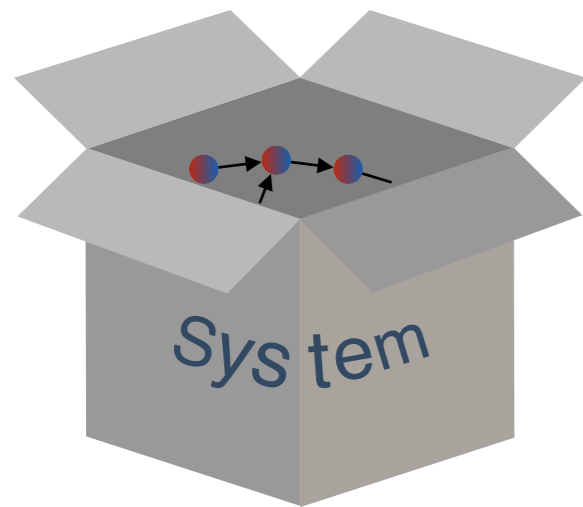
$$L = (a + b)^* b(a + b)$$

$$f_L : \Sigma^* \rightarrow \{0, 1\}$$

$$f_L(w) = 1 \text{ iff } w \in L$$

Angluin's L^* Algorithm

[Angluin 1987]



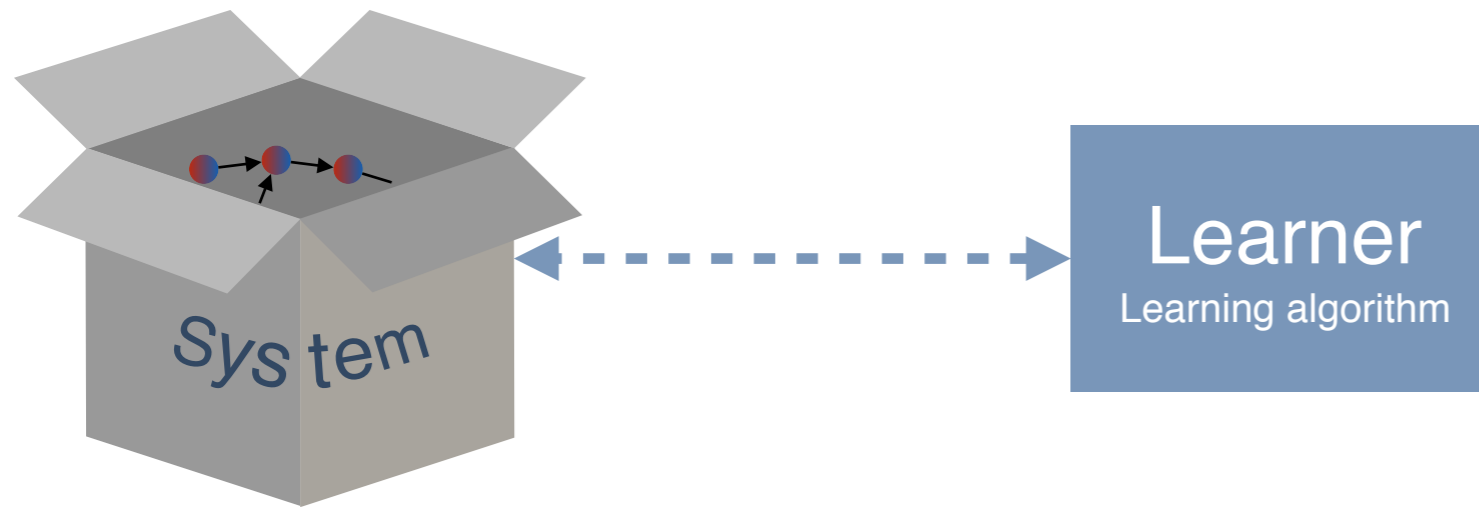
$$L = (a + b)^* b(a + b)$$

$$f_L : \Sigma^* \rightarrow \{0, 1\}$$

$$f_L(w) = 1 \text{ iff } w \in L$$

Angluin's L^* Algorithm

[Angluin 1987]



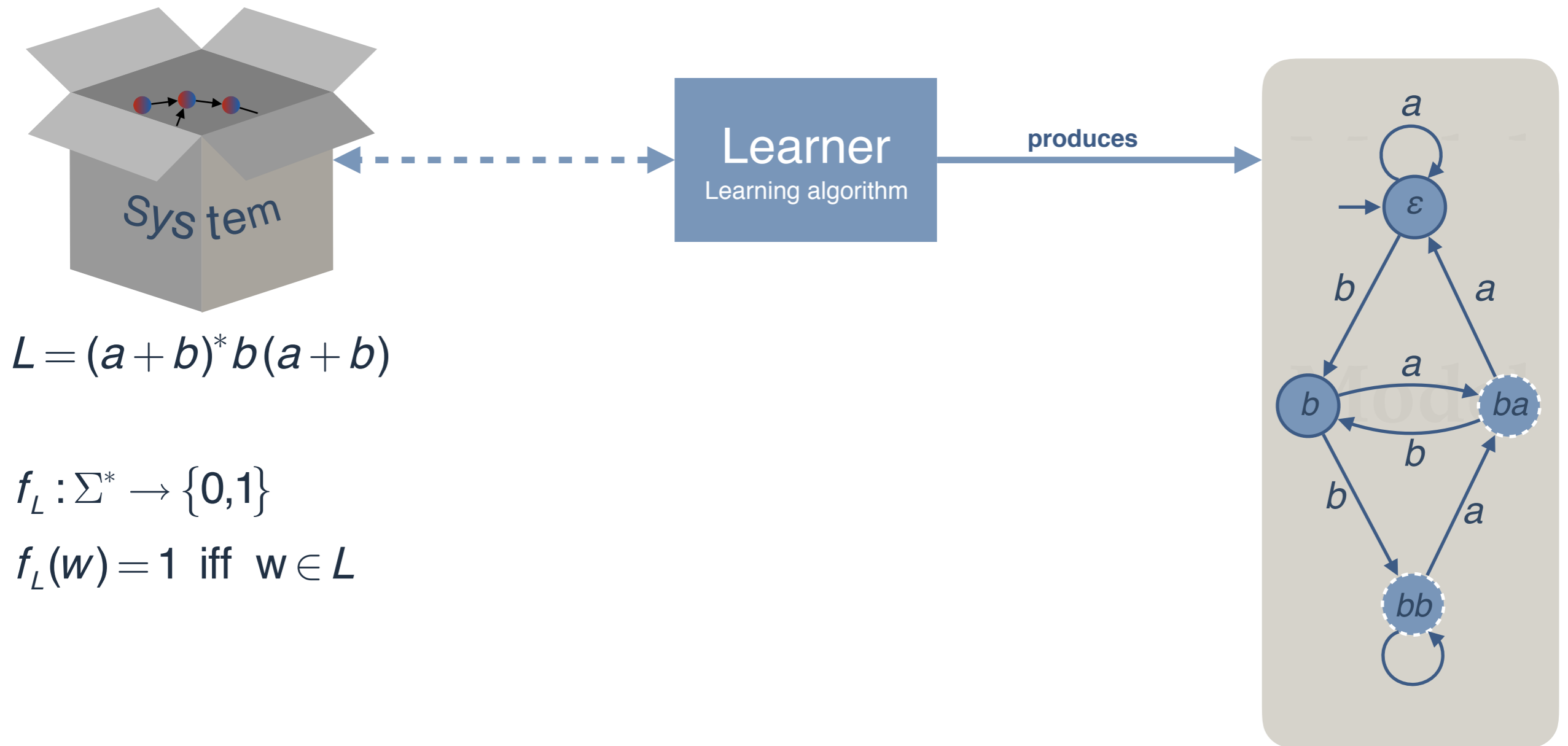
$$L = (a + b)^* b(a + b)$$

$$f_L : \Sigma^* \rightarrow \{0,1\}$$

$$f_L(w) = 1 \text{ iff } w \in L$$

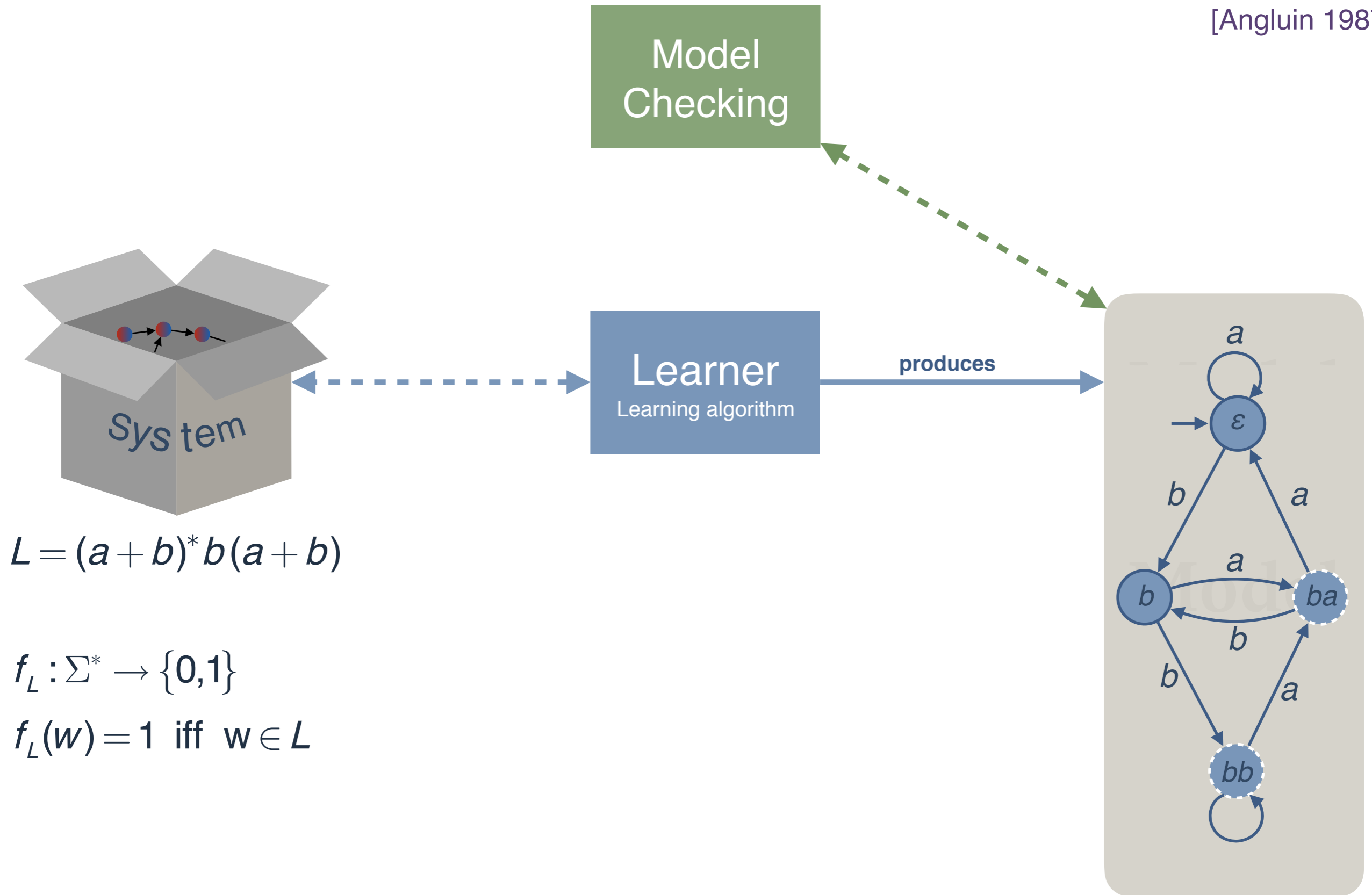
Angluin's L* Algorithm

[Angluin 1987]



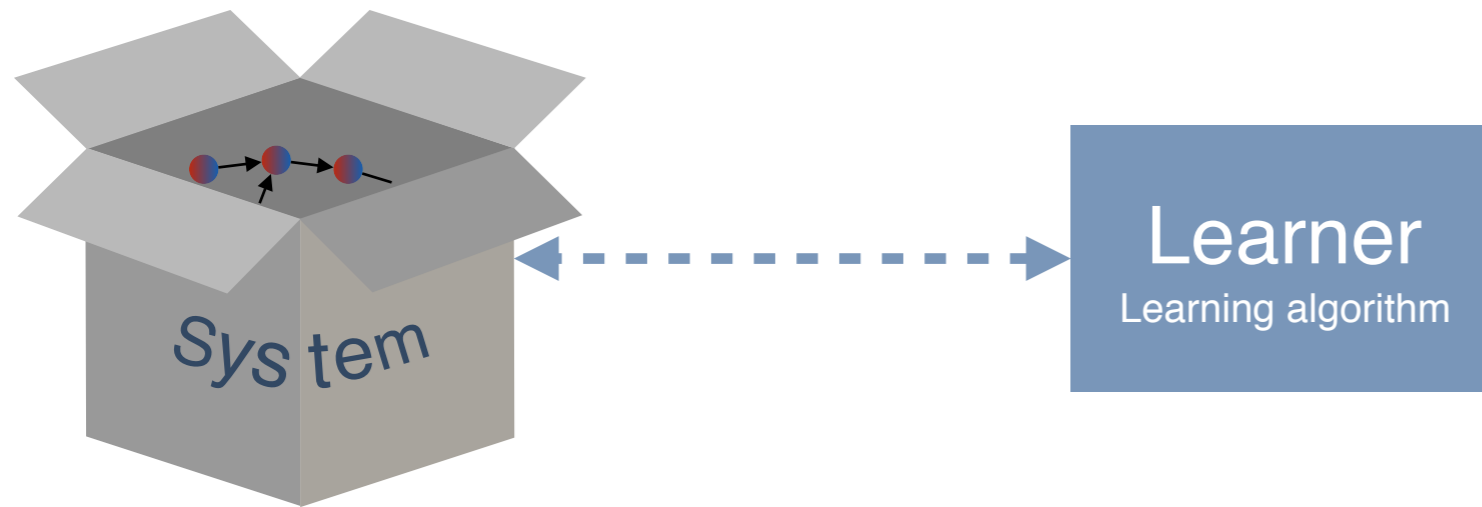
Angluin's L* Algorithm

[Angluin 1987]



Angluin's L^* Algorithm

[Angluin 1987]



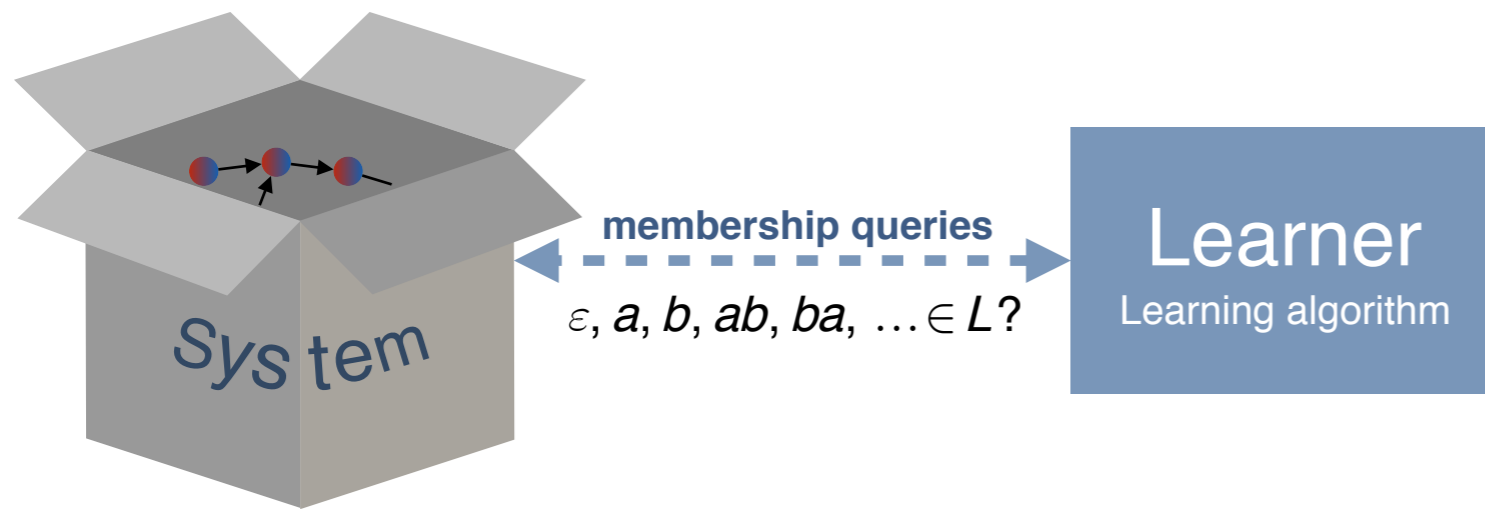
$$L = (a + b)^* b (a + b)$$

$$f_L : \Sigma^* \rightarrow \{0,1\}$$

$$f_L(w) = 1 \text{ iff } w \in L$$

Angluin's L^* Algorithm

[Angluin 1987]



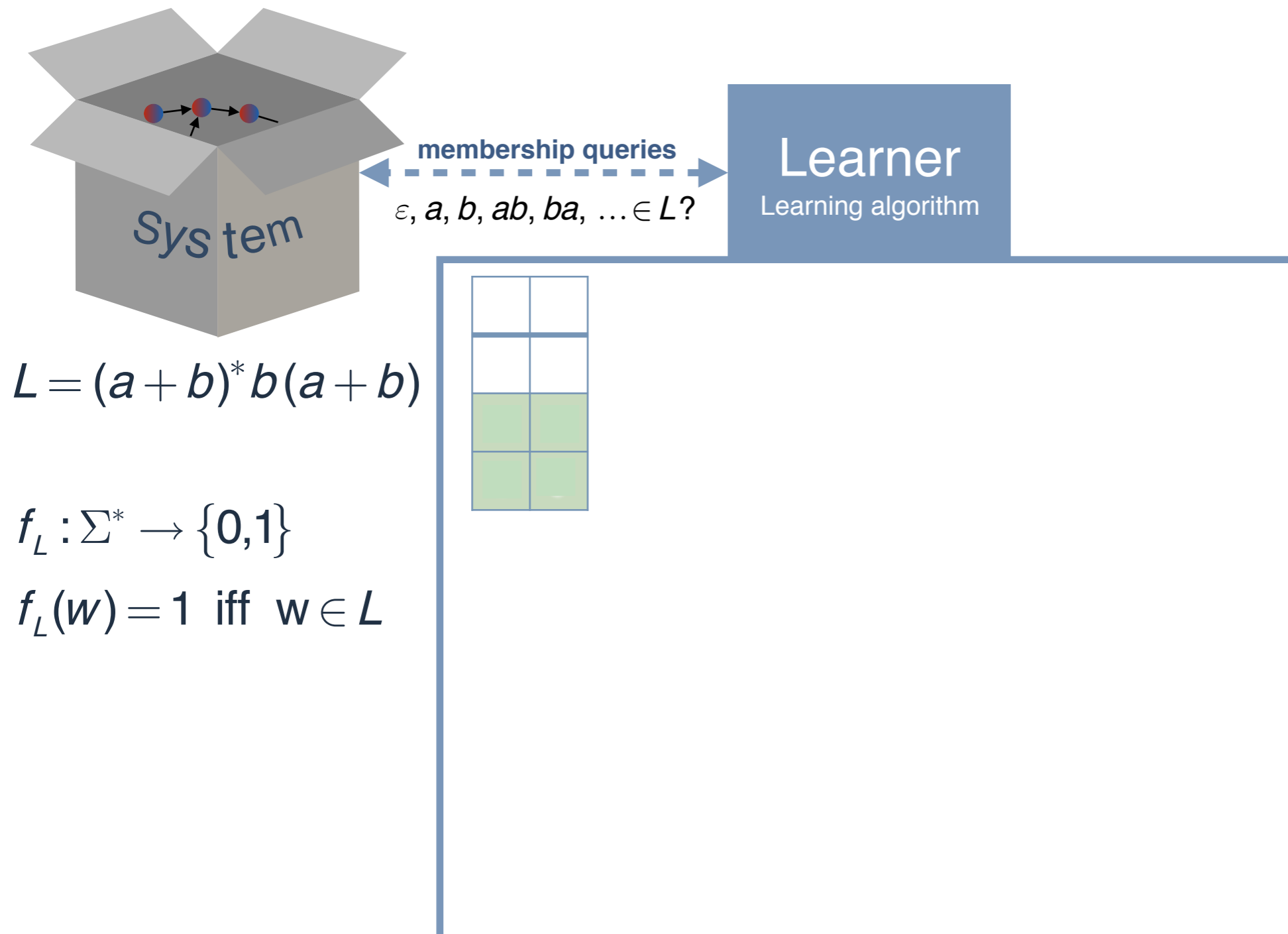
$$L = (a + b)^* b(a + b)$$

$$f_L : \Sigma^* \rightarrow \{0,1\}$$

$$f_L(w) = 1 \text{ iff } w \in L$$

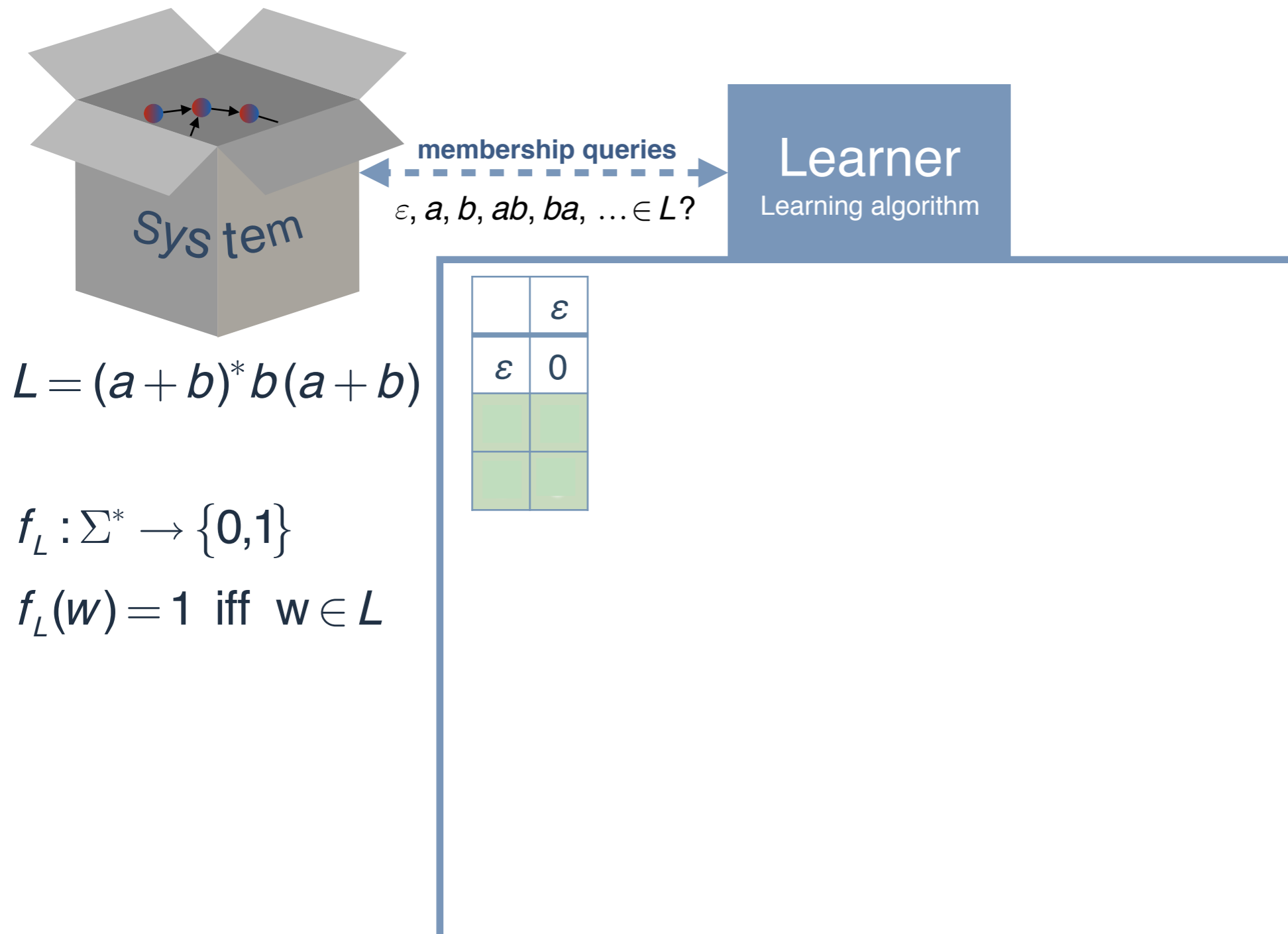
Angluin's L* Algorithm

[Angluin 1987]



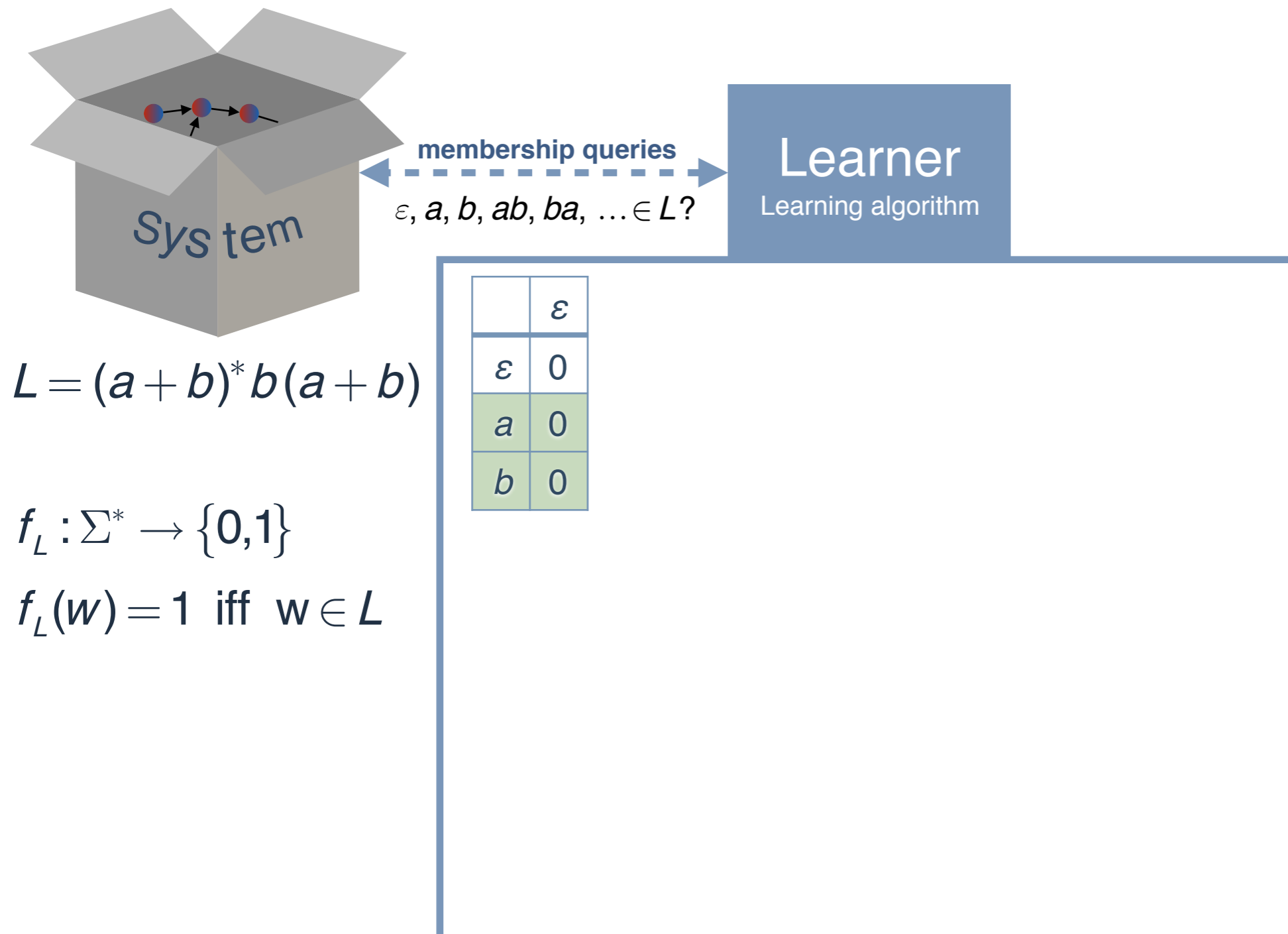
Angluin's L^* Algorithm

[Angluin 1987]



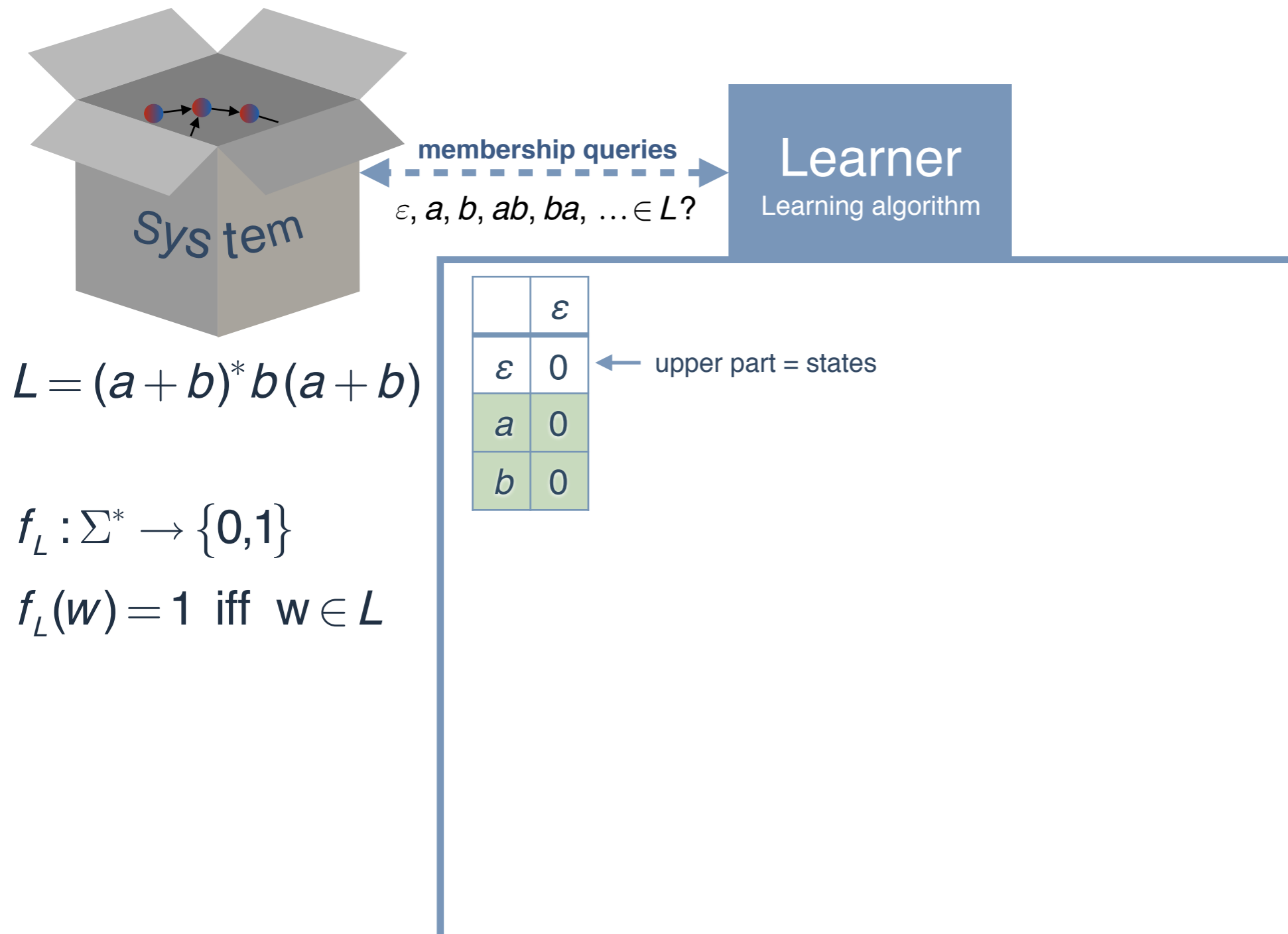
Angluin's L* Algorithm

[Angluin 1987]



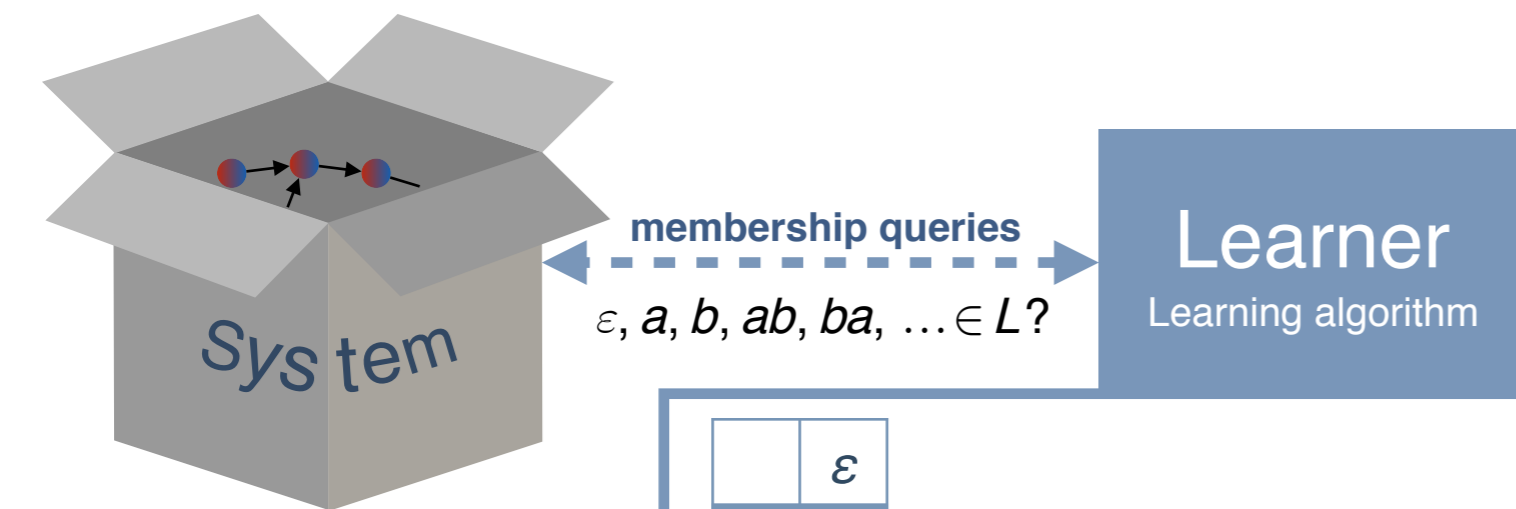
Angluin's L* Algorithm

[Angluin 1987]



Angluin's L* Algorithm

[Angluin 1987]



$$L = (a + b)^* b(a + b)$$

$$f_L : \Sigma^* \rightarrow \{0,1\}$$

$$f_L(w) = 1 \text{ iff } w \in L$$

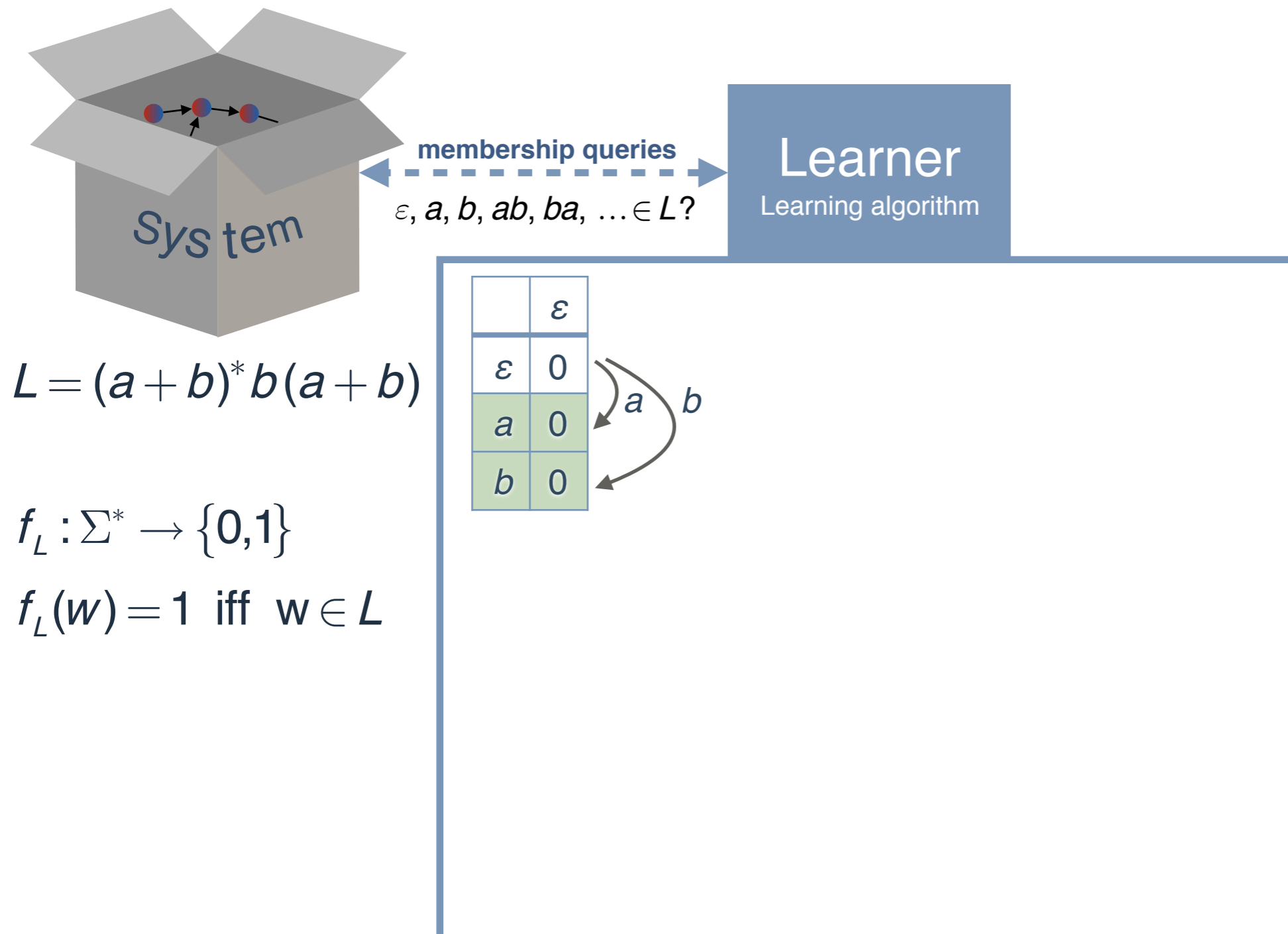
| | ϵ |
|------------|------------|
| ϵ | 0 |
| a | 0 |
| b | 0 |

← upper part = states

↘ lower part = transitions

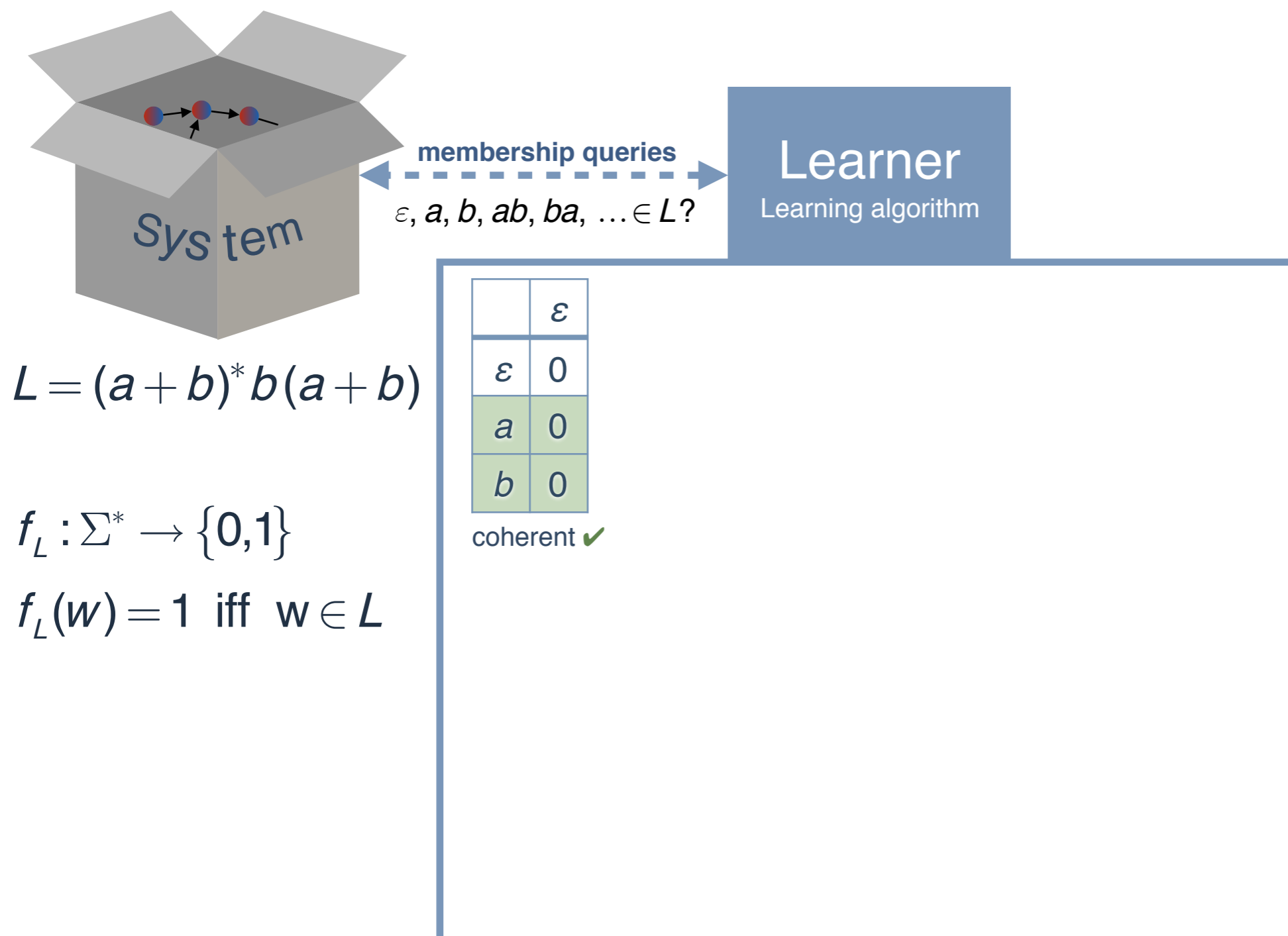
Angluin's L* Algorithm

[Angluin 1987]



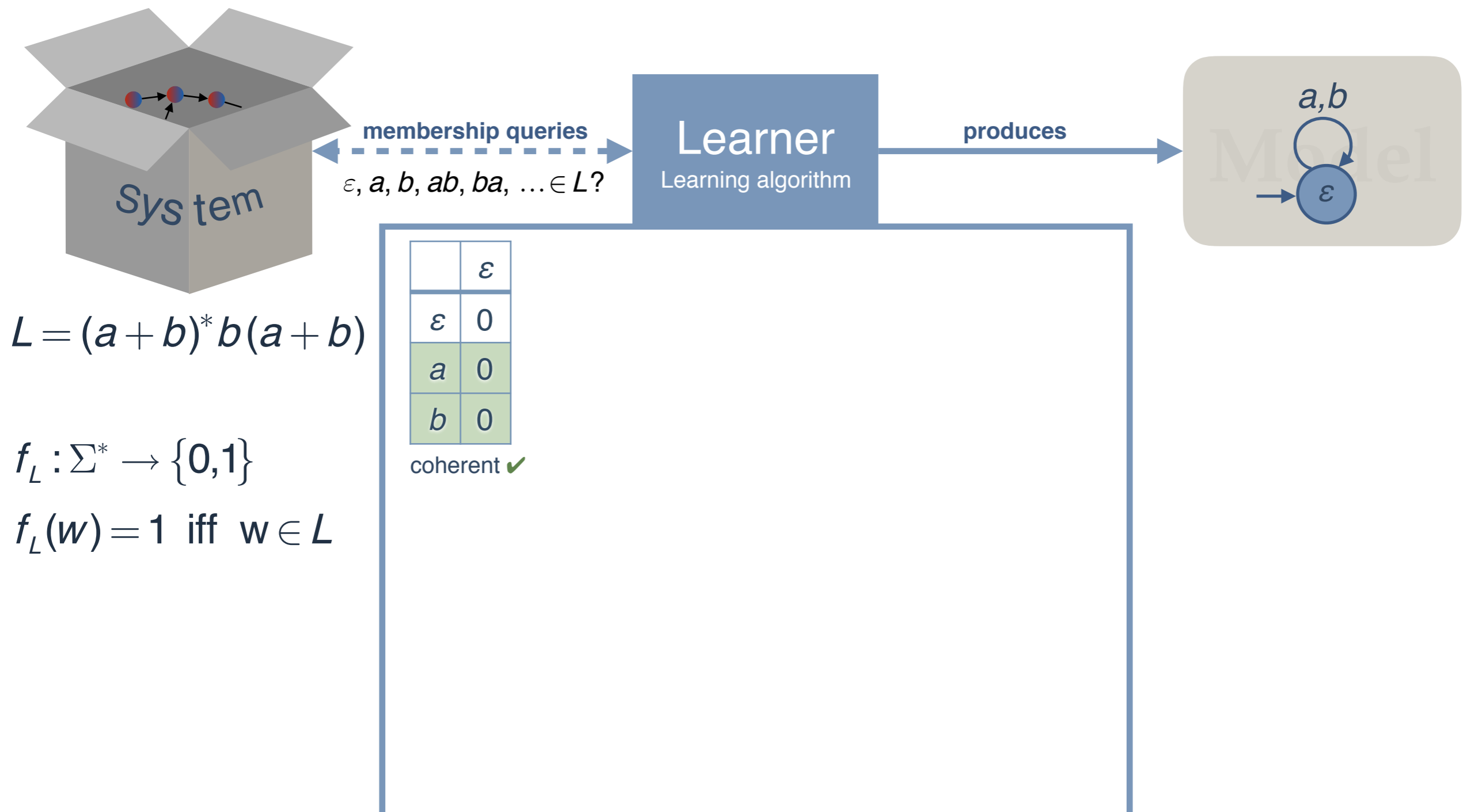
Angluin's L* Algorithm

[Angluin 1987]



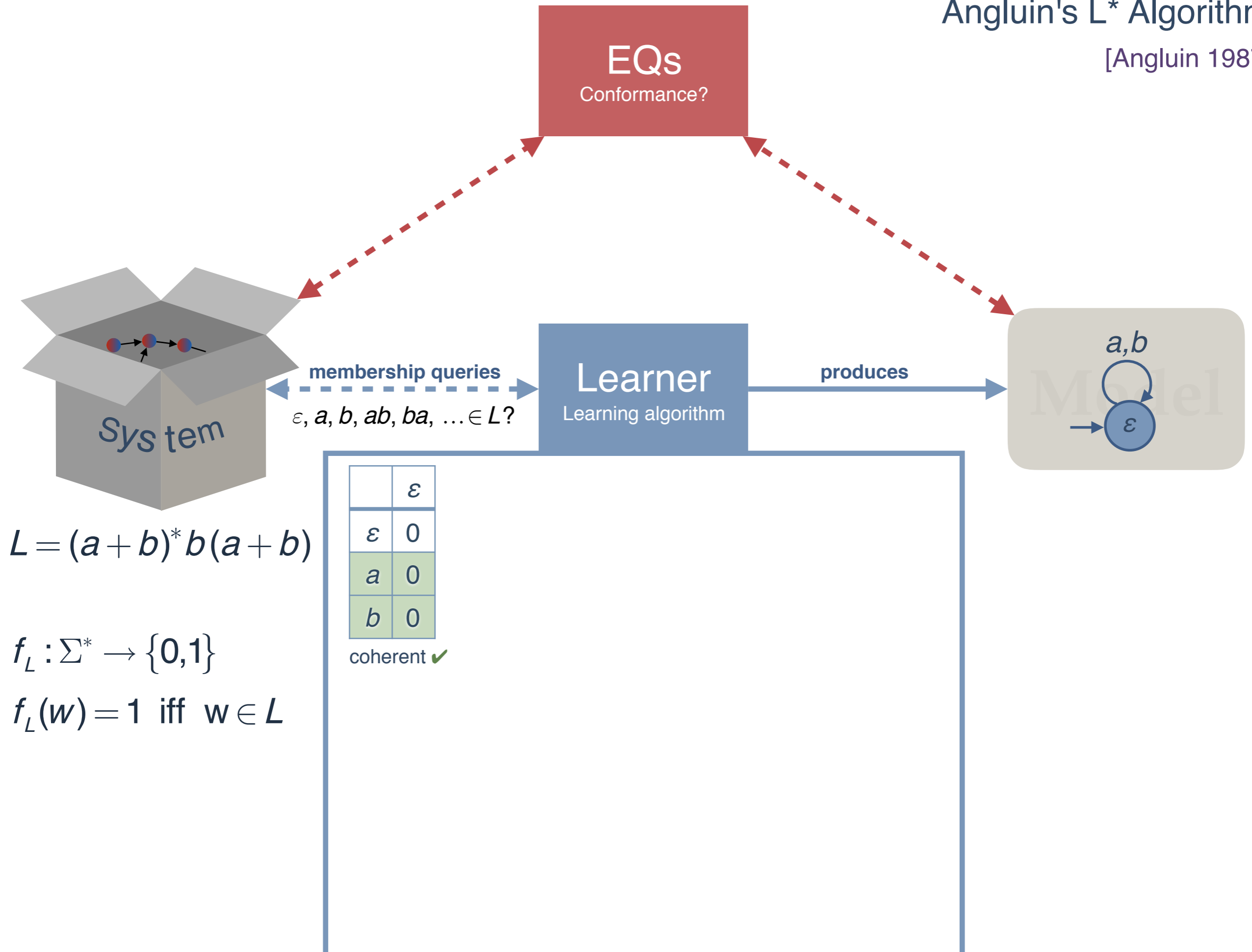
Angluin's L* Algorithm

[Angluin 1987]



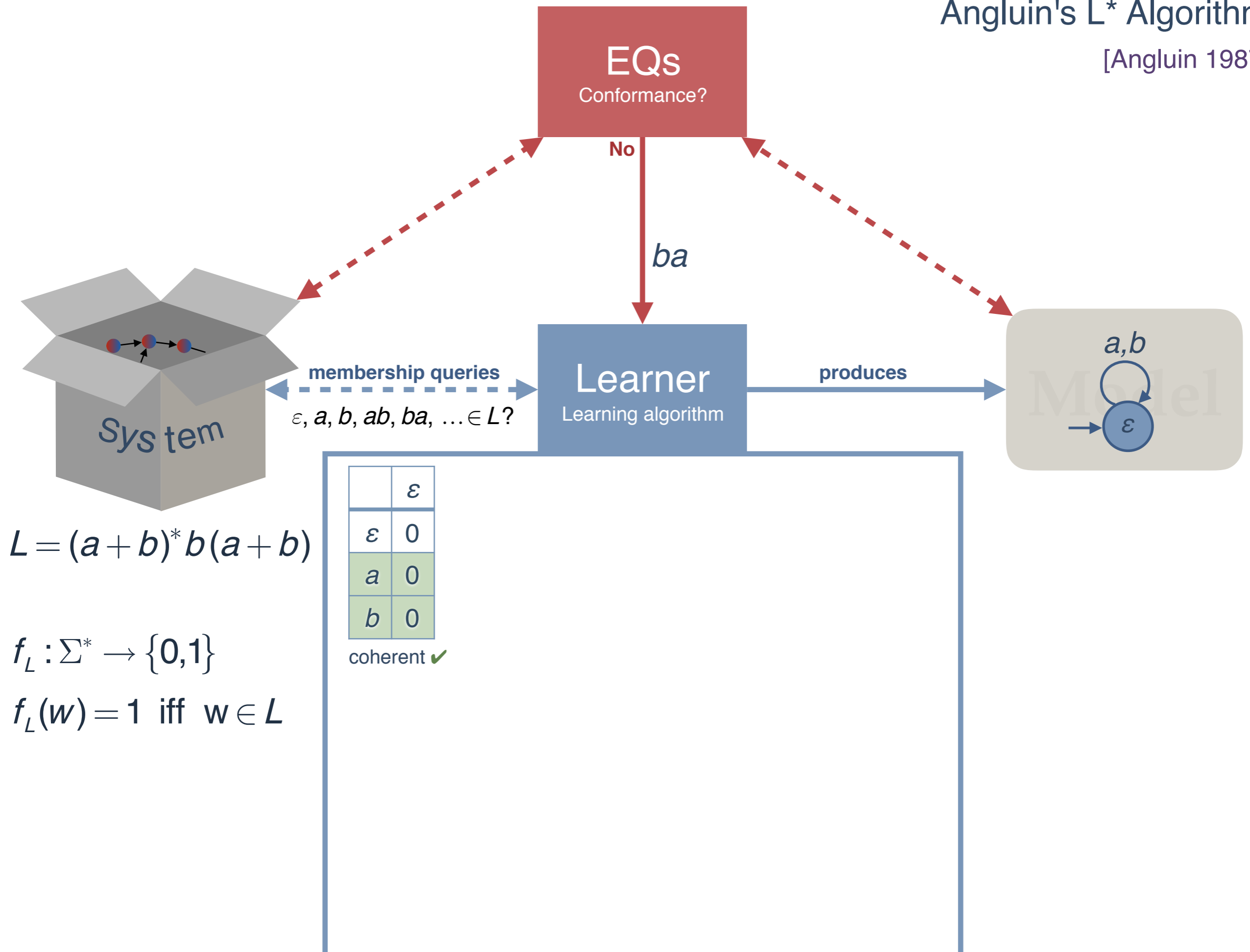
Angluin's L* Algorithm

[Angluin 1987]



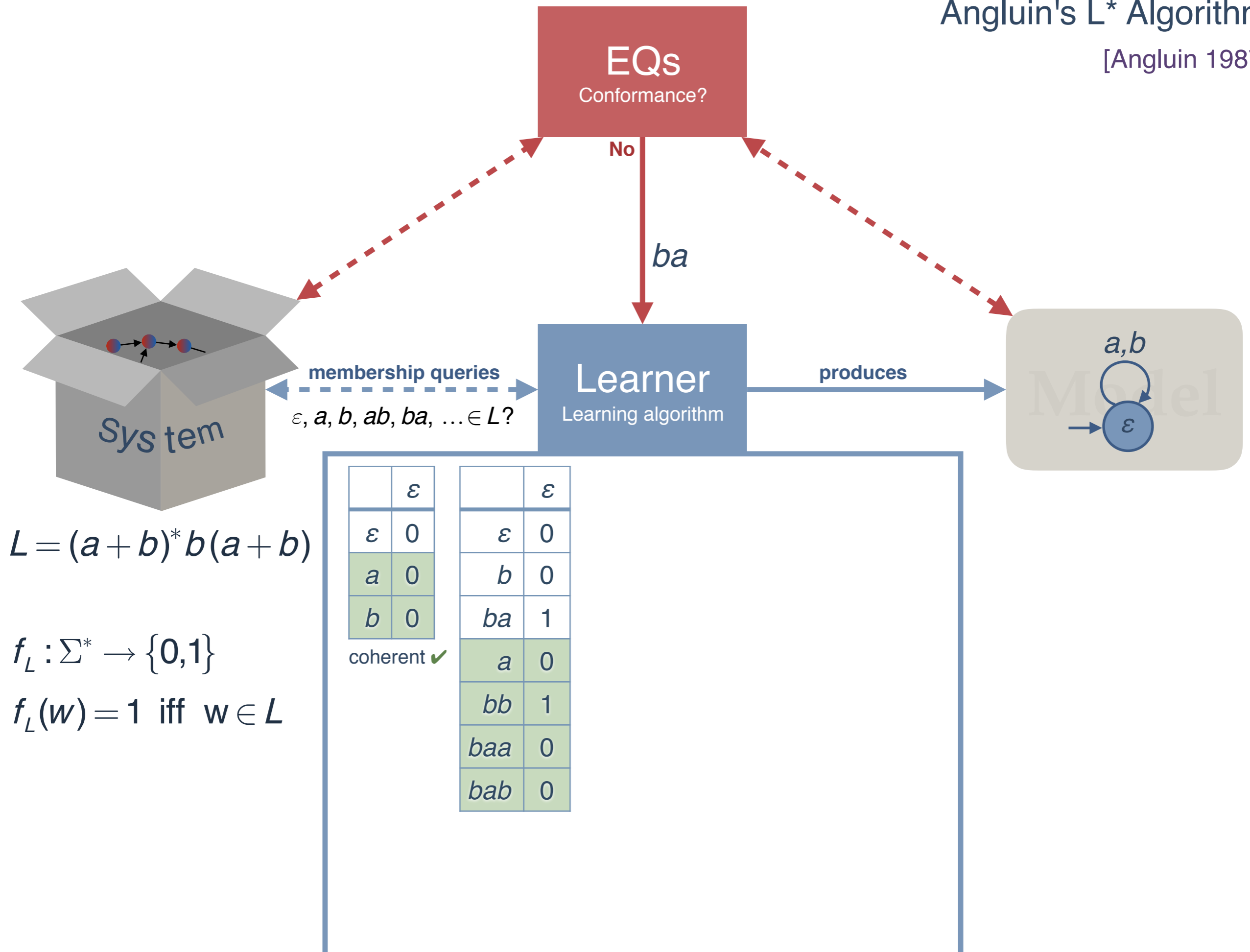
Angluin's L* Algorithm

[Angluin 1987]



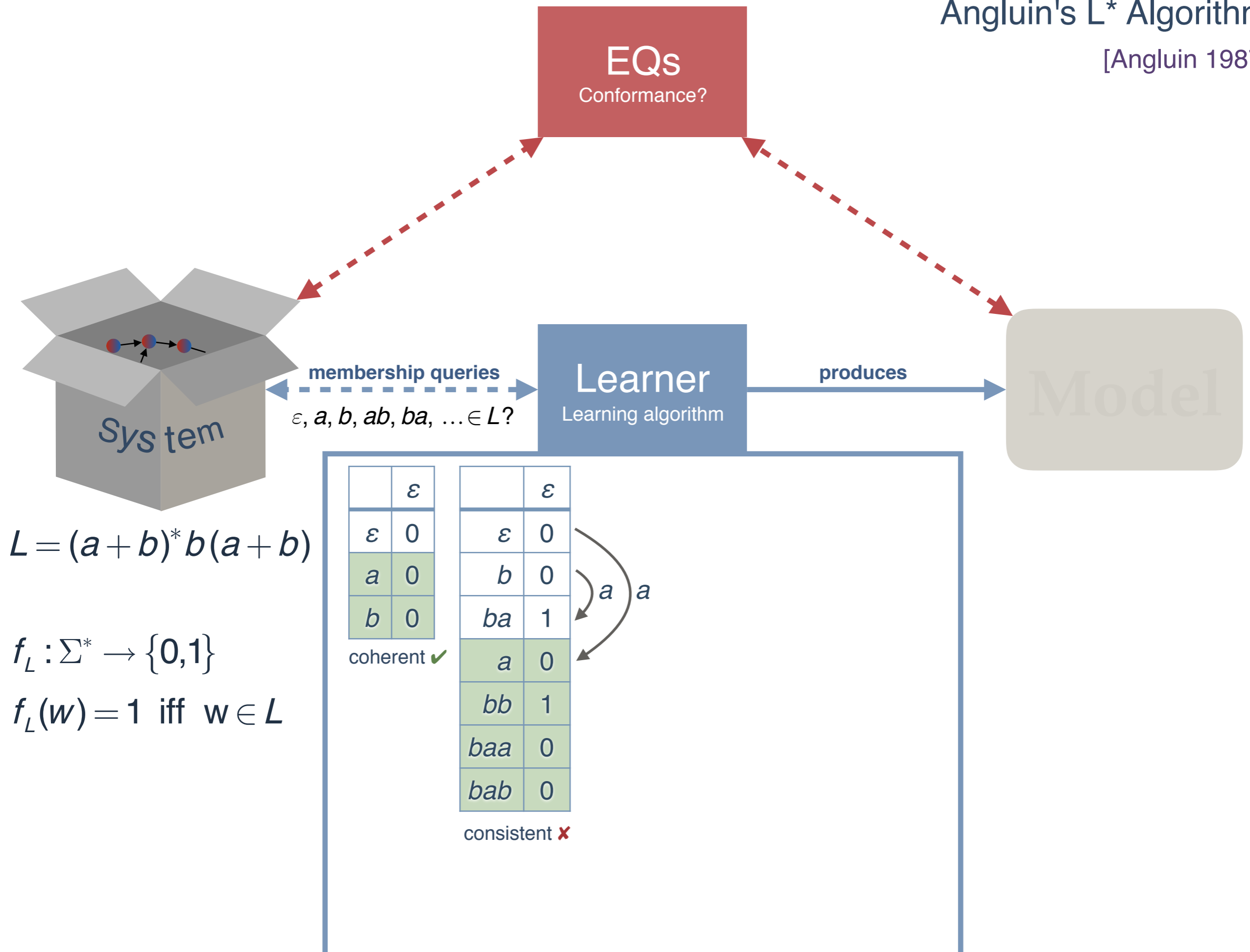
Angluin's L* Algorithm

[Angluin 1987]



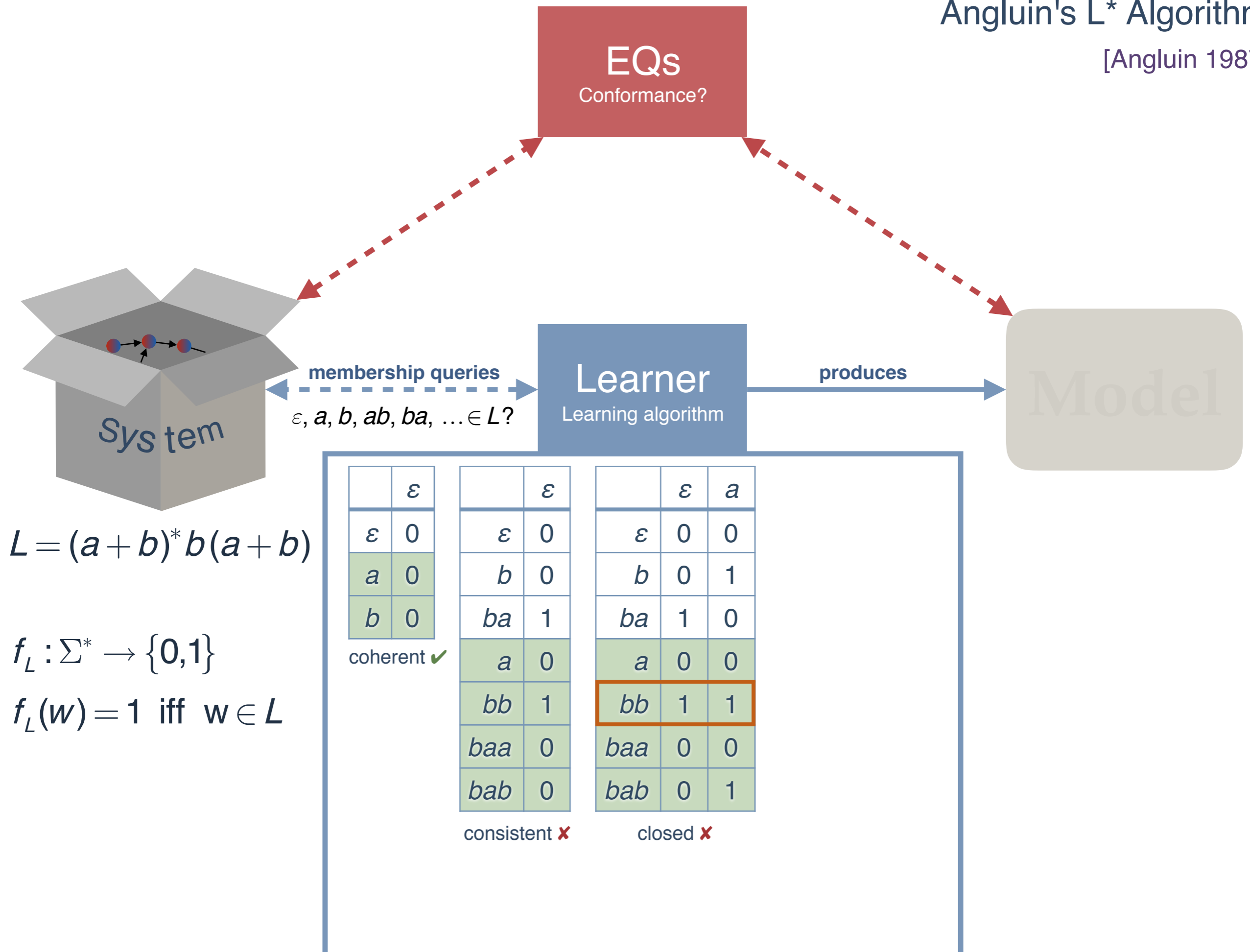
Angluin's L* Algorithm

[Angluin 1987]



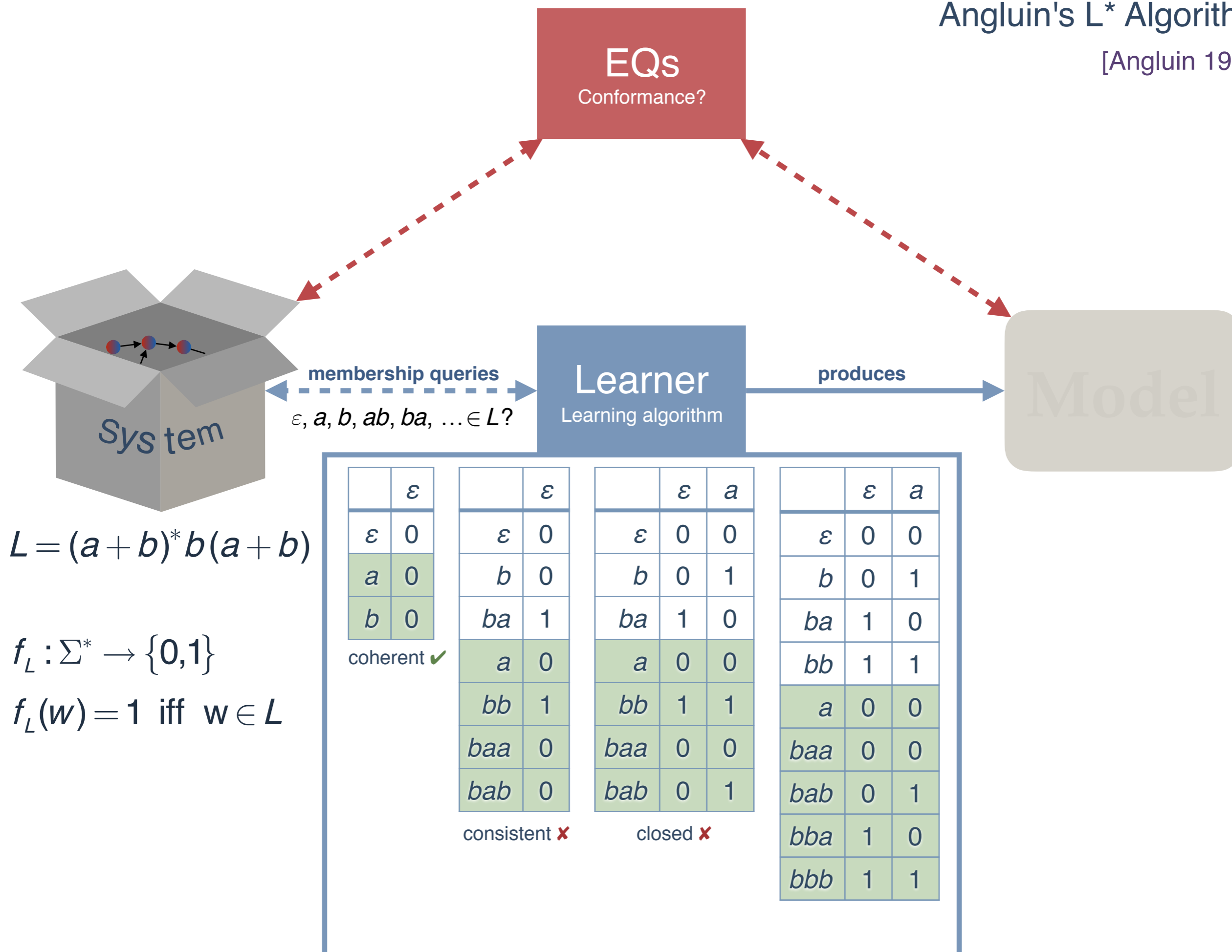
Angluin's L* Algorithm

[Angluin 1987]



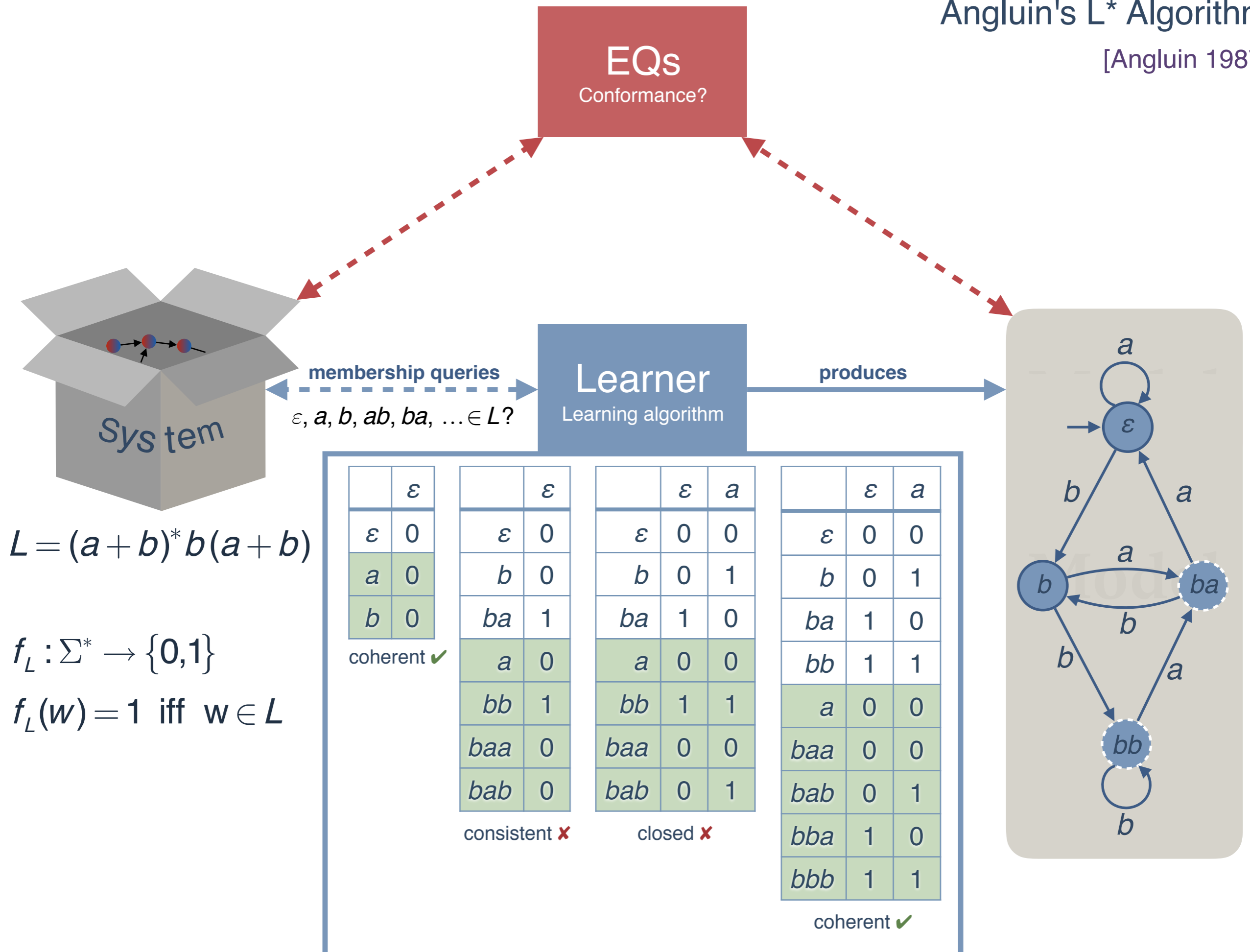
Angluin's L* Algorithm

[Angluin 1987]



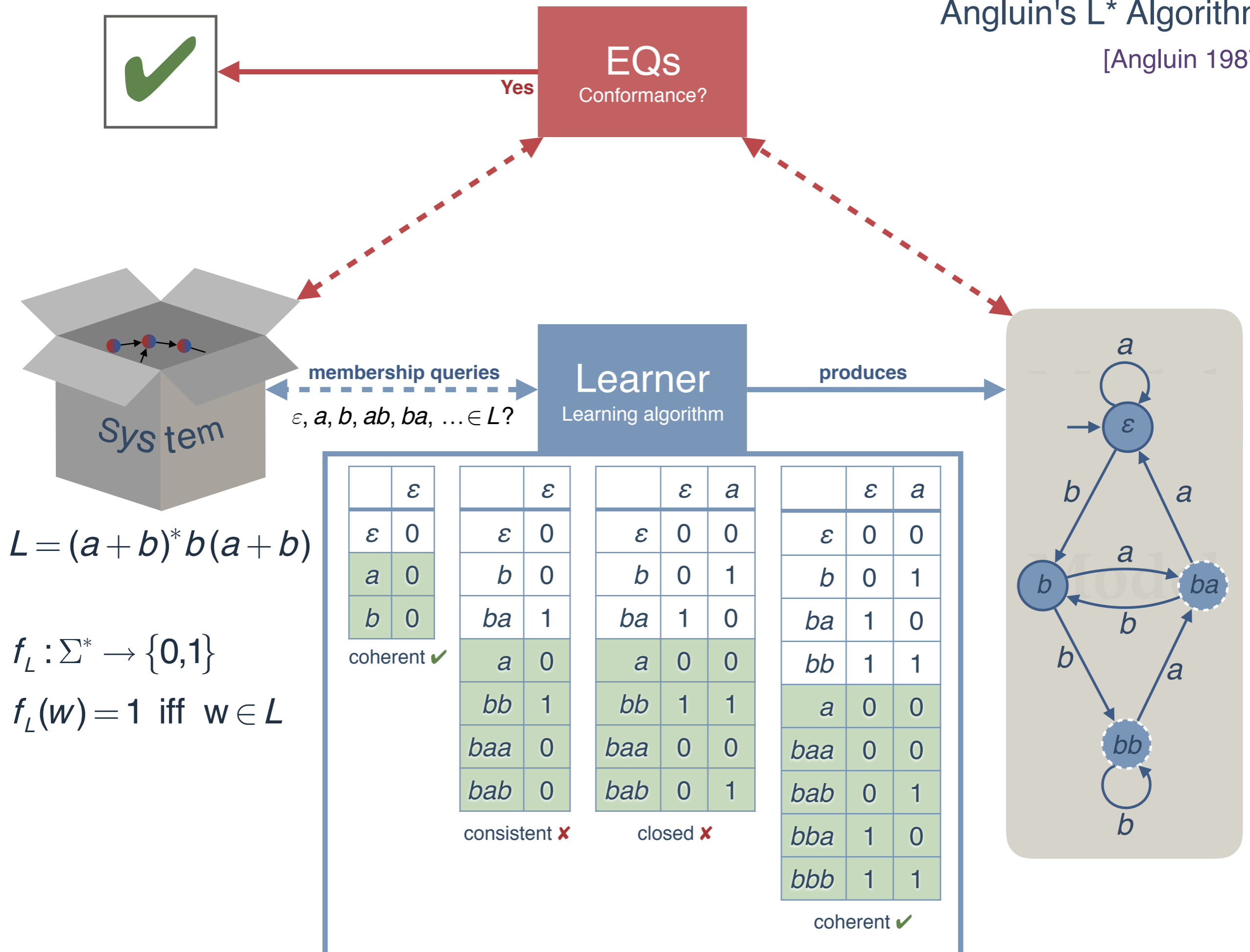
Angluin's L* Algorithm

[Angluin 1987]



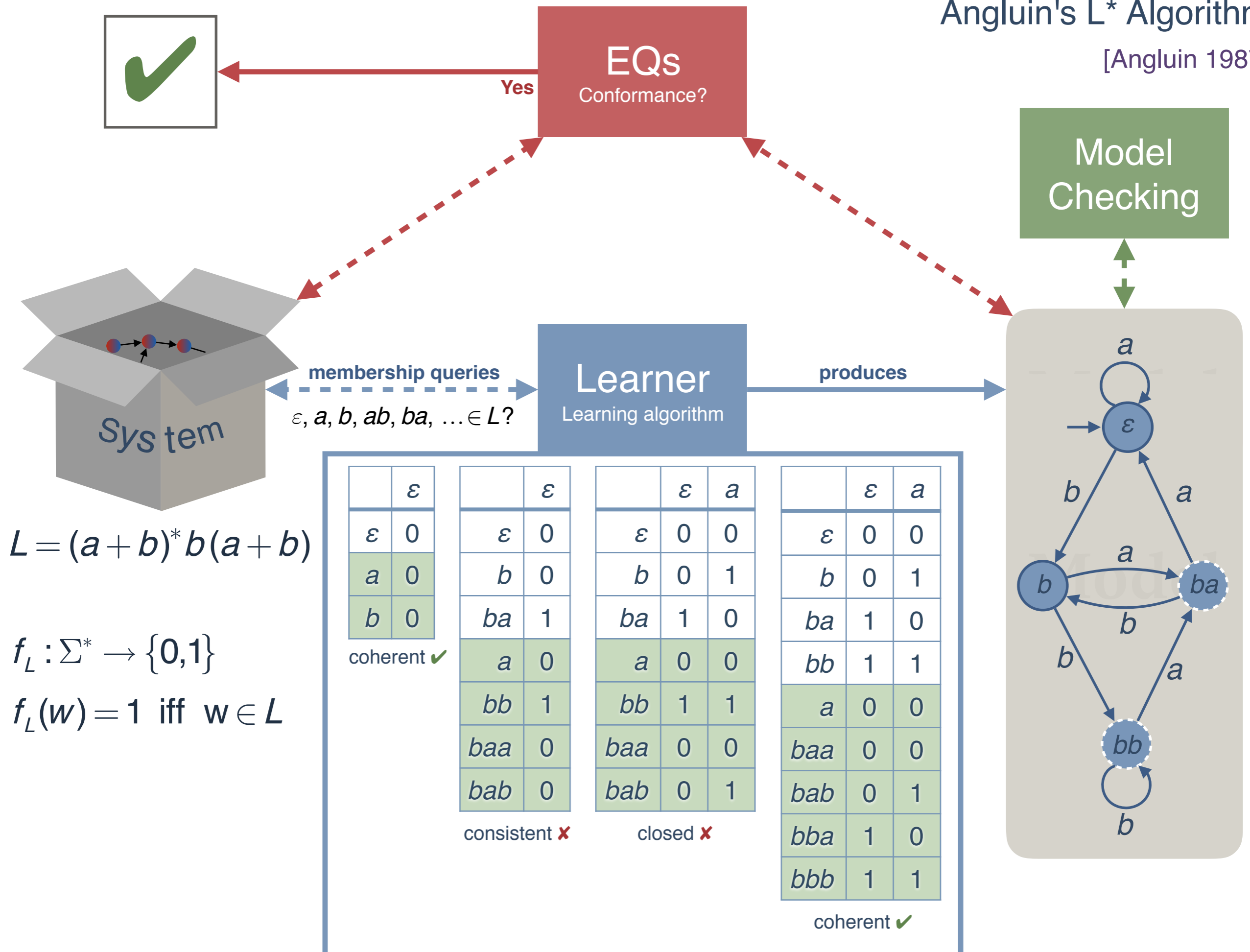
Angluin's L* Algorithm

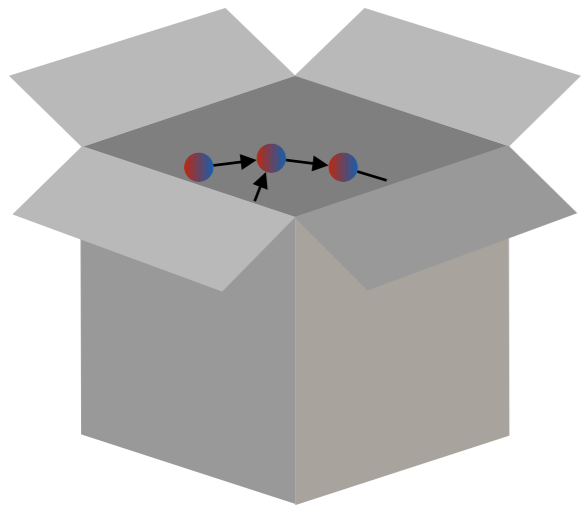
[Angluin 1987]



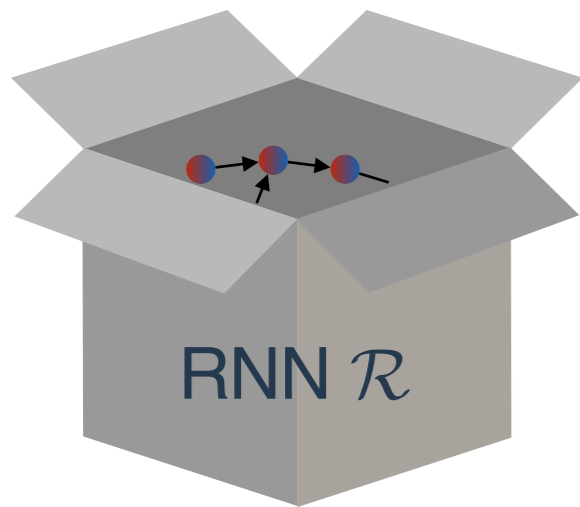
Angluin's L* Algorithm

[Angluin 1987]

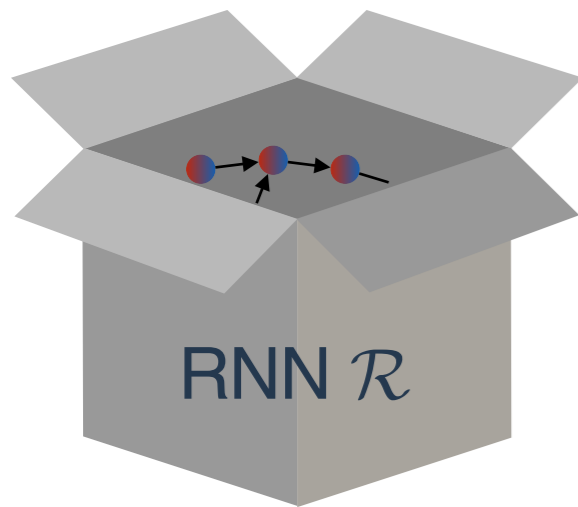




- legacy software
- code not open source
- third-party software
- embedded systems

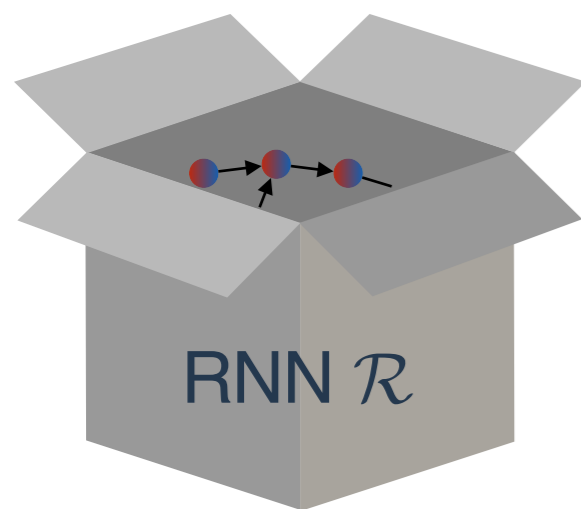


- legacy software
- code not open source
- third-party software
- embedded systems
- recurrent neural networks (RNNs)

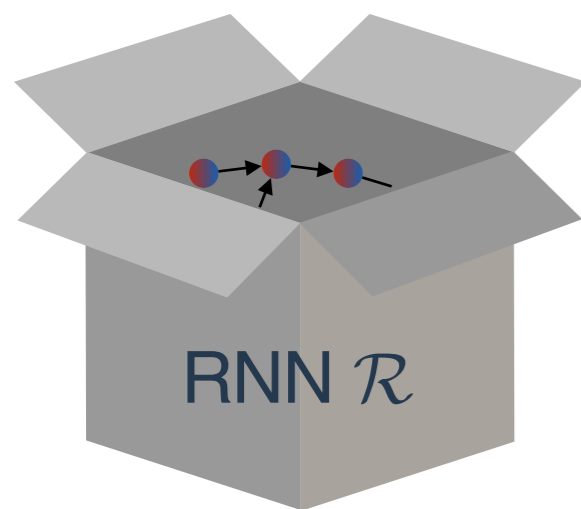


- legacy software
- code not open source
- third-party software
- embedded systems
- recurrent neural networks (RNNs)

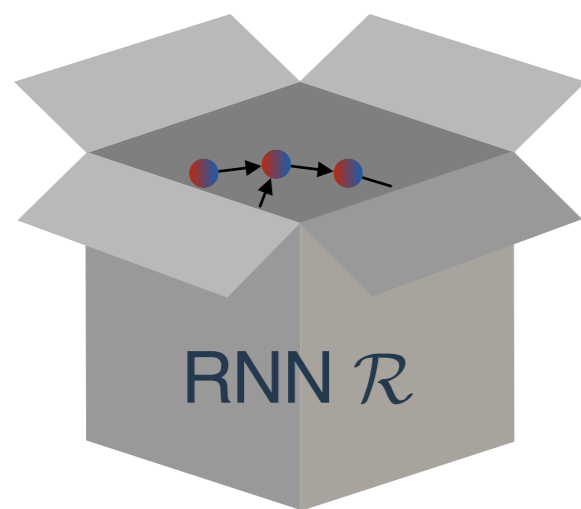
(approach works for binary classifiers)



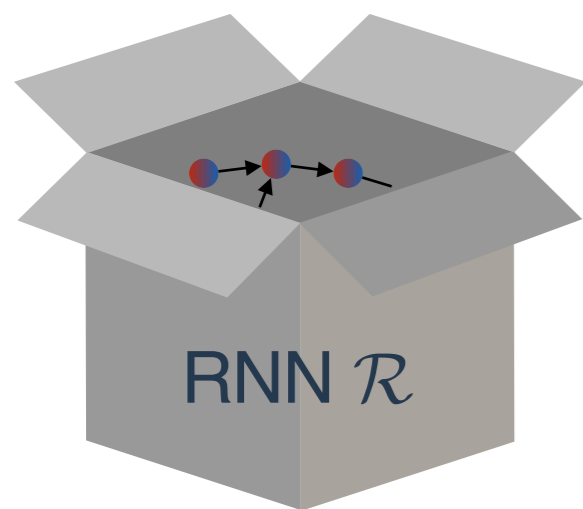
$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\mathbb{R}^\ell}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$



$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\Sigma}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$

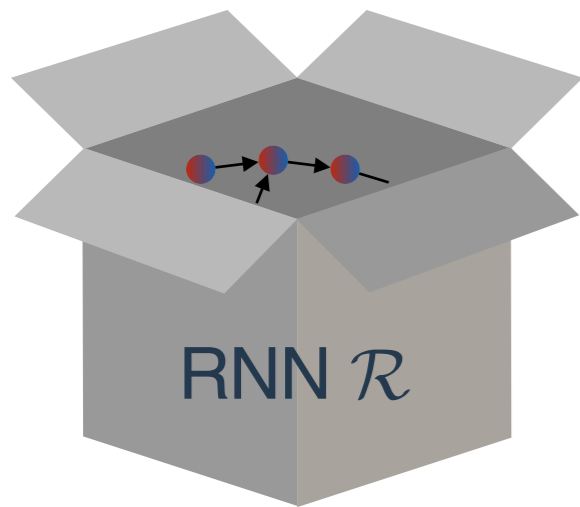


$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\Sigma^*}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$



$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\Sigma^*}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$

$$g: \underbrace{\mathbb{R}^k}_{\text{state}} \rightarrow \{0,1\}$$



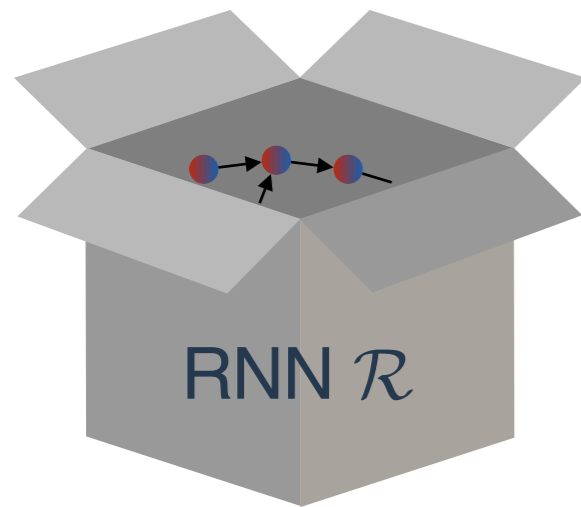
$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\Sigma^*}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$

$$g: \underbrace{\mathbb{R}^k}_{\text{state}} \rightarrow \{0,1\}$$

\rightsquigarrow language acceptor

$$L(\mathcal{R}) = \{w \in \Sigma^* \mid g(f(\text{init}, w)) = 1\}$$

What does correctness for RNNs mean?



$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\Sigma^*}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$

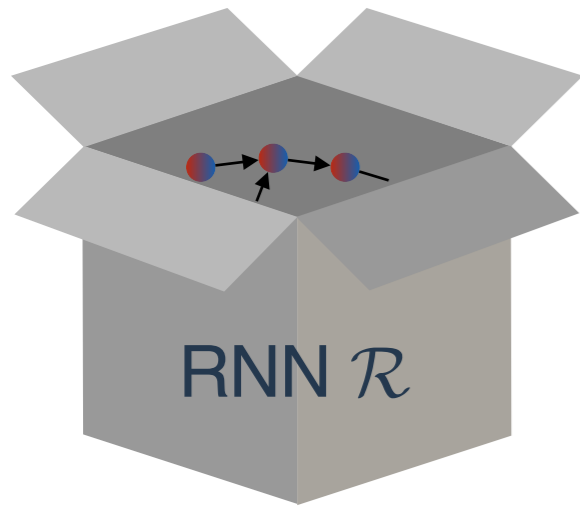
$$g: \underbrace{\mathbb{R}^k}_{\text{state}} \rightarrow \{0,1\}$$

\rightsquigarrow language acceptor

$$L(\mathcal{R}) = \{w \in \Sigma^* \mid g(f(\text{init}, w)) = 1\}$$

What does correctness for RNNs mean?

Specify correctness using
finite automata \mathcal{E} over Σ



$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\Sigma^*}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$

$$g: \underbrace{\mathbb{R}^k}_{\text{state}} \rightarrow \{0,1\}$$

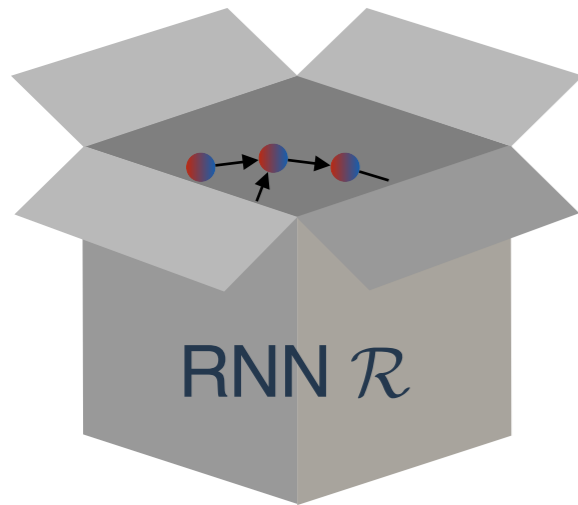
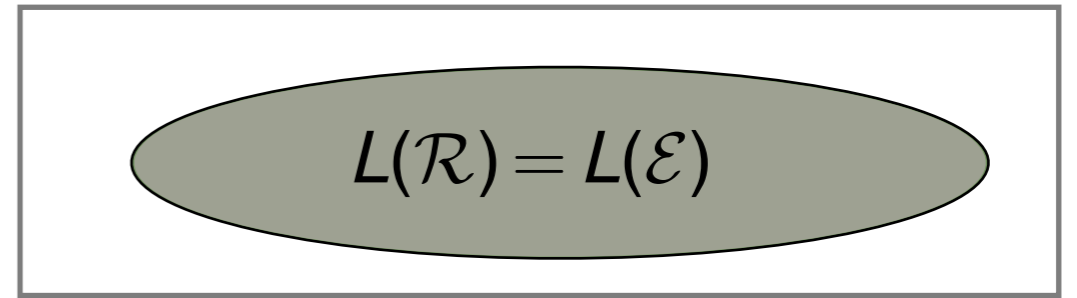
\rightsquigarrow language acceptor

$$L(\mathcal{R}) = \{w \in \Sigma^* \mid g(f(\text{init}, w)) = 1\}$$

What does correctness for RNNs mean?

Specify correctness using
finite automata \mathcal{E} over Σ

$$L(\mathcal{R}) = L(\mathcal{E})$$



$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\Sigma^*}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$

$$g: \underbrace{\mathbb{R}^k}_{\text{state}} \rightarrow \{0,1\}$$

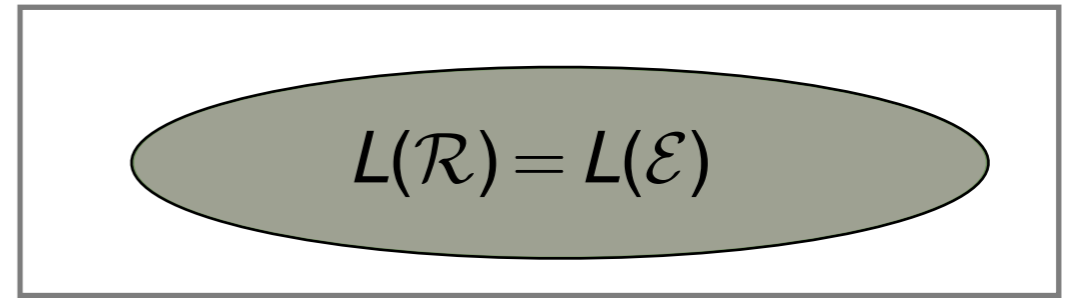
\rightsquigarrow language acceptor

$$L(\mathcal{R}) = \{w \in \Sigma^* \mid g(f(\text{init}, w)) = 1\}$$

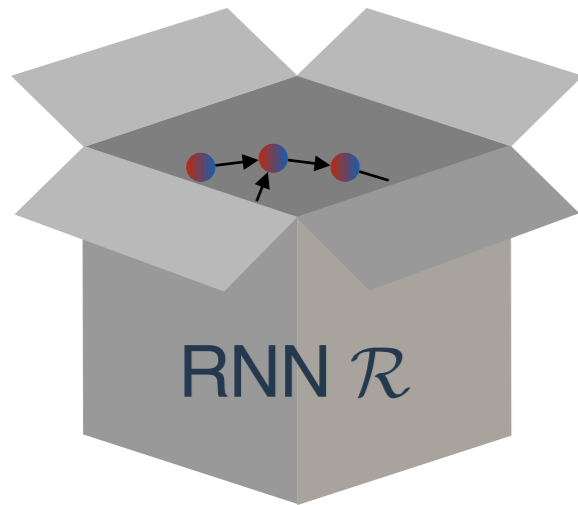
What does correctness for RNNs mean?

Specify correctness using
finite automata \mathcal{E} over Σ

$$L(\mathcal{R}) = L(\mathcal{E})$$



$\rightsquigarrow L(\mathcal{R})$ is a regular language



$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\Sigma^*}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$

$$g: \underbrace{\mathbb{R}^k}_{\text{state}} \rightarrow \{0,1\}$$

\rightsquigarrow language acceptor

$$L(\mathcal{R}) = \{w \in \Sigma^* \mid g(f(\text{init}, w)) = 1\}$$

What does correctness for RNNs mean?

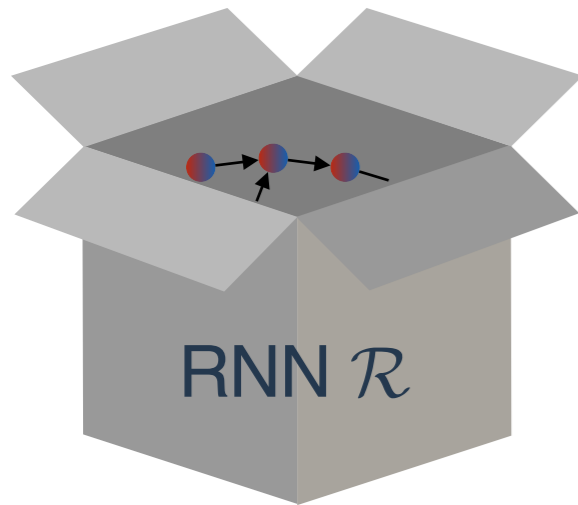
Specify correctness using
finite automata \mathcal{E} over Σ

$$L(\mathcal{R}) = L(\mathcal{E})$$

$$L(\mathcal{R}) = L(\mathcal{E})$$

$\rightsquigarrow L(\mathcal{R})$ is a regular language

restrictive: how about RNNs recognizing XML documents?



$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\Sigma^*}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$

$$g: \underbrace{\mathbb{R}^k}_{\text{state}} \rightarrow \{0,1\}$$

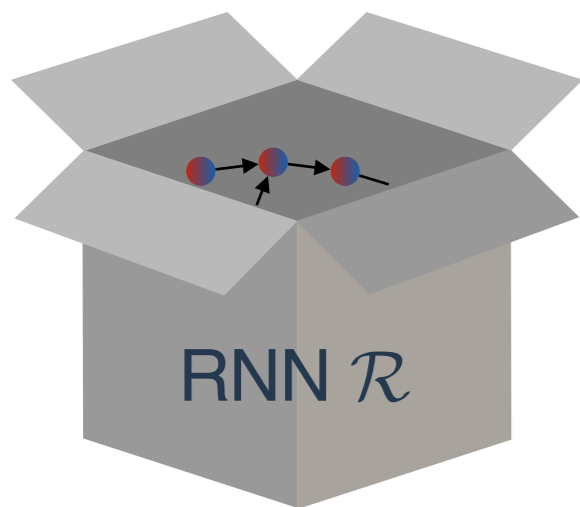
\rightsquigarrow language acceptor

$$L(\mathcal{R}) = \{w \in \Sigma^* \mid g(f(\text{init}, w)) = 1\}$$

What does correctness for RNNs mean?

Specify correctness using
finite automata \mathcal{E}, \mathcal{N} over Σ

$$L(\mathcal{R}) = L(\mathcal{E})$$

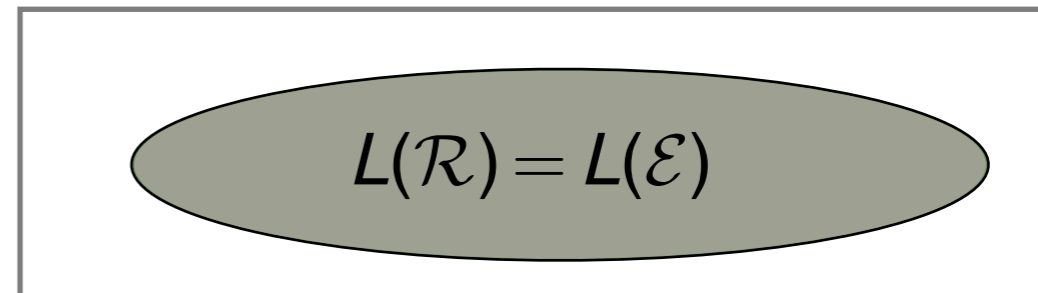


$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\Sigma^*}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$

$$g: \underbrace{\mathbb{R}^k}_{\text{state}} \rightarrow \{0,1\}$$

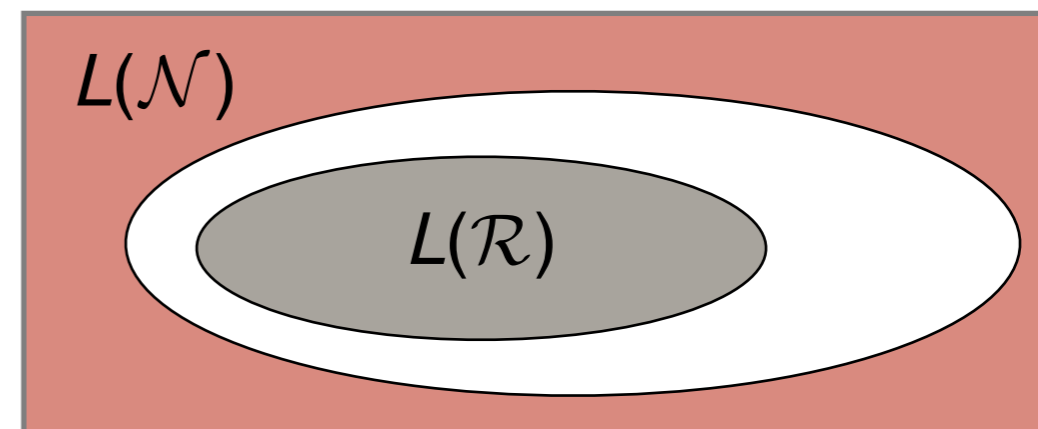
\rightsquigarrow language acceptor

$$L(\mathcal{R}) = \{w \in \Sigma^* \mid g(f(\text{init}, w)) = 1\}$$



$\rightsquigarrow L(\mathcal{R})$ is a regular language

restrictive: how about RNNs recognizing XML documents?



each $w \in L(\mathcal{N})$ must be classified as negative
 \mathcal{R} does not produce false positives

What does correctness for RNNs mean?

Specify correctness using
finite automata \mathcal{E}, \mathcal{N} over Σ

$$L(\mathcal{R}) = L(\mathcal{E})$$

$$L(\mathcal{R}) = L(\mathcal{E})$$

$\rightsquigarrow L(\mathcal{R})$ is a regular language

restrictive: how about RNNs recognizing XML documents?

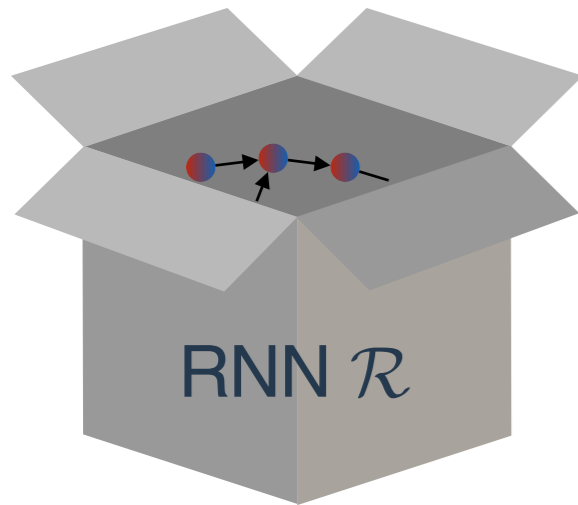
complement

$$L(\mathcal{R}) \subseteq L(\overline{\mathcal{N}})$$

$$L(\mathcal{N})$$

$$L(\mathcal{R})$$

each $w \in L(\mathcal{N})$ must be classified as negative
 \mathcal{R} does not produce false positives



$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\Sigma^*}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$

$$g: \underbrace{\mathbb{R}^k}_{\text{state}} \rightarrow \{0,1\}$$

\rightsquigarrow language acceptor

$$L(\mathcal{R}) = \{w \in \Sigma^* \mid g(f(\text{init}, w)) = 1\}$$

What does correctness for RNNs mean?

Specify correctness using
finite automata \mathcal{E} , \mathcal{N} , \mathcal{P} over Σ

$$L(\mathcal{R}) = L(\mathcal{E})$$

$$L(\mathcal{R}) = L(\mathcal{E})$$

$\rightsquigarrow L(\mathcal{R})$ is a regular language

restrictive: how about RNNs recognizing XML documents?

complement

$$L(\mathcal{R}) \subseteq L(\overline{\mathcal{N}})$$

$$L(\mathcal{N})$$

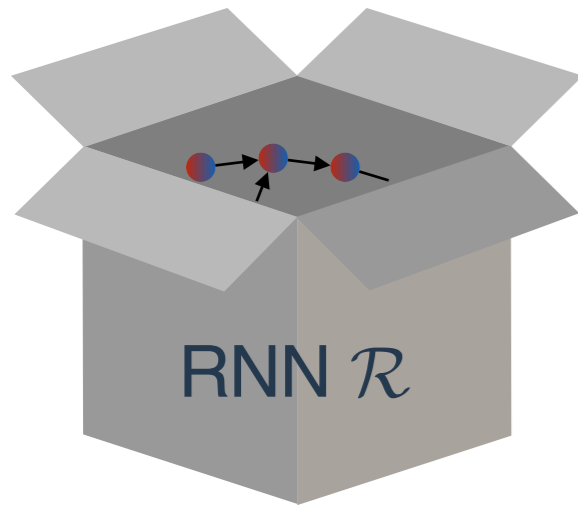
$$L(\mathcal{R})$$

each $w \in L(\mathcal{N})$ must be classified as negative
 \mathcal{R} does not produce false positives

$$L(\mathcal{R})$$

$$L(\mathcal{P})$$

each $w \in L(\mathcal{P})$ must be classified as positive
 \mathcal{R} does not produce false negatives



$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\Sigma^*}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$

$$g: \underbrace{\mathbb{R}^k}_{\text{state}} \rightarrow \{0,1\}$$

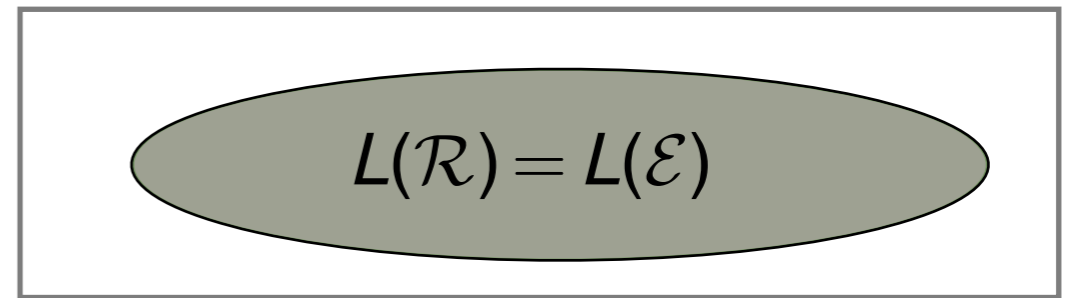
\rightsquigarrow language acceptor

$$L(\mathcal{R}) = \{w \in \Sigma^* \mid g(f(\text{init}, w)) = 1\}$$

What does correctness for RNNs mean?

Specify correctness using finite automata \mathcal{E} , \mathcal{N} , \mathcal{P} over Σ

$$L(\mathcal{R}) = L(\mathcal{E})$$

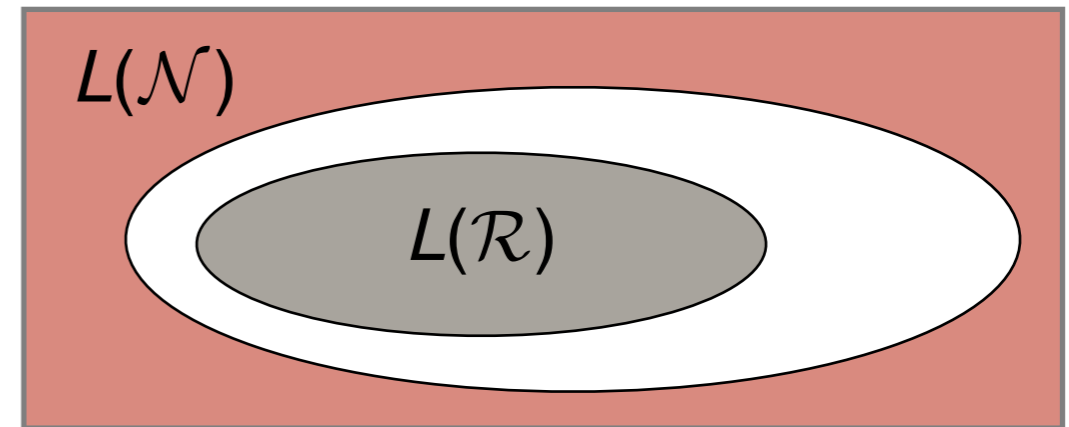


$\rightsquigarrow L(\mathcal{R})$ is a regular language

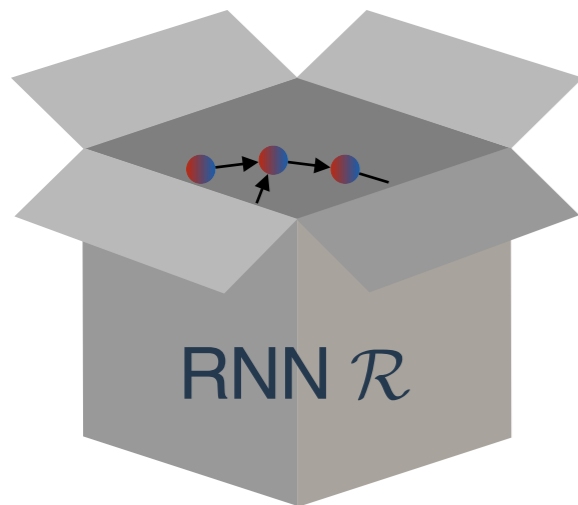
restrictive: how about RNNs recognizing XML documents?

complement

$$L(\mathcal{R}) \subseteq L(\overline{\mathcal{N}})$$



each $w \in L(\mathcal{N})$ must be classified as negative
 \mathcal{R} does not produce false positives



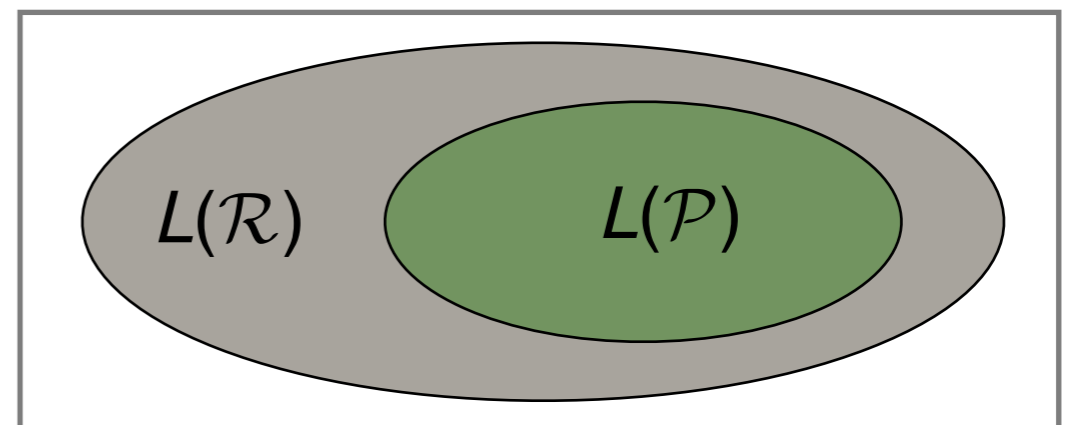
$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\Sigma^*}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$

$$g: \underbrace{\mathbb{R}^k}_{\text{state}} \rightarrow \{0,1\}$$

\rightsquigarrow language acceptor

$$L(\mathcal{R}) = \{w \in \Sigma^* \mid g(f(\text{init}, w)) = 1\}$$

$$L(\overline{\mathcal{R}}) \subseteq L(\overline{\mathcal{P}})$$



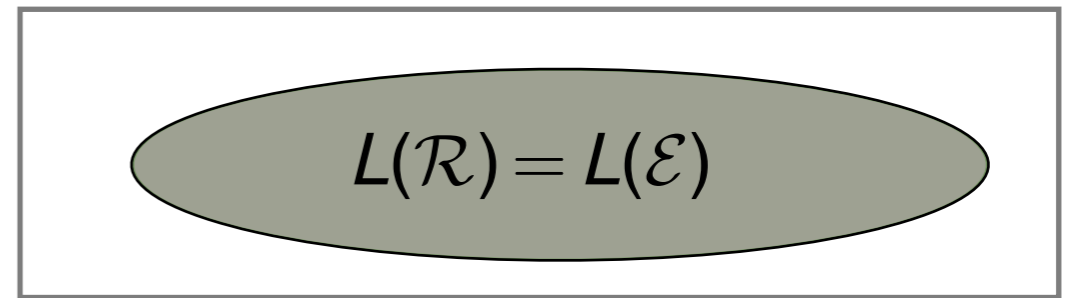
each $w \in L(\mathcal{P})$ must be classified as positive
 \mathcal{R} does not produce false negatives

complement

What does correctness for RNNs mean?

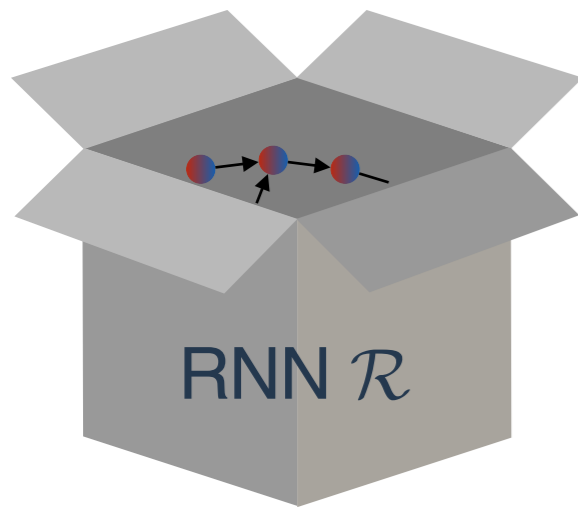
Specify correctness using
finite automata $\mathcal{E}, \mathcal{N}, \mathcal{P}$ over Σ

$$L(\mathcal{R}) = L(\mathcal{E})$$



$\rightsquigarrow L(\mathcal{R})$ is a regular language

restrictive: how about RNNs recognizing XML documents?



$$f: \underbrace{\mathbb{R}^k}_{\text{state}} \times \underbrace{\Sigma^*}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^k}_{\text{state}}$$

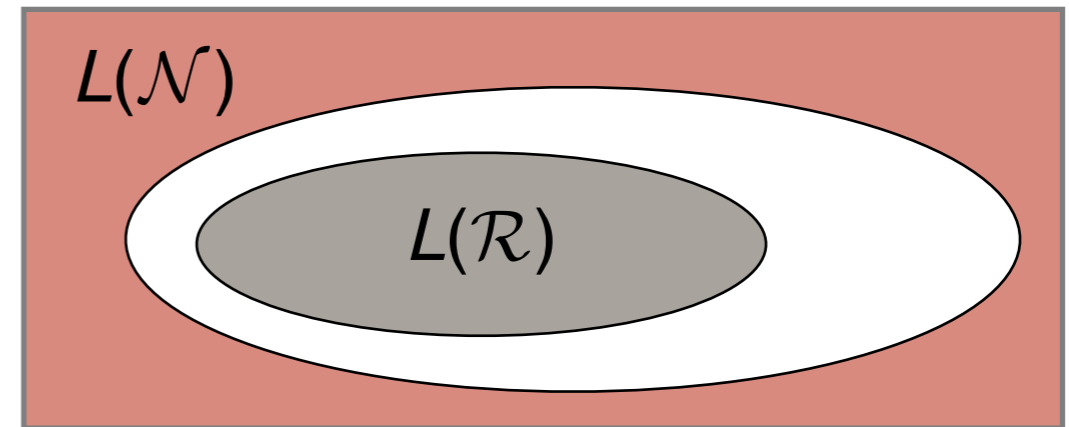
$$g: \underbrace{\mathbb{R}^k}_{\text{state}} \rightarrow \{0,1\}$$

\rightsquigarrow language acceptor

$$L(\mathcal{R}) = \{w \in \Sigma^* \mid g(f(\text{init}, w)) = 1\}$$

complement

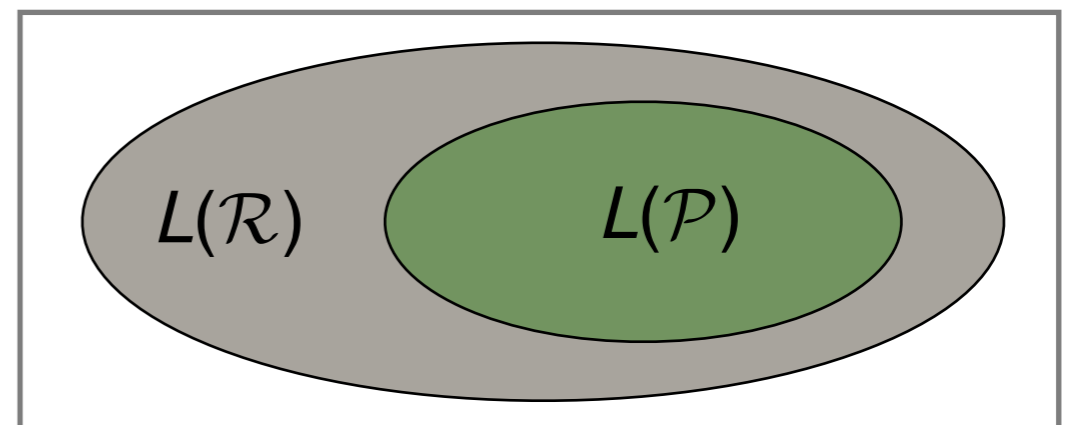
$$L(\mathcal{R}) \subseteq L(\overline{\mathcal{N}})$$



each $w \in L(\mathcal{N})$ must be classified as negative
 \mathcal{R} does not produce false positives

**inclusion
checking**

$$L(\overline{\mathcal{R}}) \subseteq L(\overline{\mathcal{P}})$$



each $w \in L(\mathcal{P})$ must be classified as positive
 \mathcal{R} does not produce false negatives

complement

Assume \mathcal{R} is supposed to recognize XML documents over

$$\Sigma = \{ \text{<list>, </list>, <item>, </item> } \}.$$

```
<list>
  <item>
    <list>
      <item>
      </item>
      <item>
      </item>
    </list>
  </item>
  <item>
  </item>
</list>
```

Assume \mathcal{R} is supposed to recognize XML documents over

$$\Sigma = \{ \langle \text{list} \rangle, \langle / \text{list} \rangle, \langle \text{item} \rangle, \langle / \text{item} \rangle \}.$$

```
<list>
  <item>
    <list>
      <item>
      </item>
      <item>
      </item>
    </list>
  </item>
  <item>
  </item>
</list>
```

Regular negative specification \mathcal{N} :

Assume \mathcal{R} is supposed to recognize XML documents over

$$\Sigma = \{ \langle \text{list} \rangle, \langle / \text{list} \rangle, \langle \text{item} \rangle, \langle / \text{item} \rangle \}.$$

```
<list>
  <item>
    <list>
      <item>
      </item>
      <item>
      </item>
    </list>
  </item>
  <item>
  </item>
</list>
```

Regular negative specification \mathcal{N} :

"there is an opening tag that is not
eventually followed by a corresponding closing tag"

Assume \mathcal{R} is supposed to recognize XML documents over

$$\Sigma = \{ \langle \text{list} \rangle, \langle / \text{list} \rangle, \langle \text{item} \rangle, \langle / \text{item} \rangle \}.$$

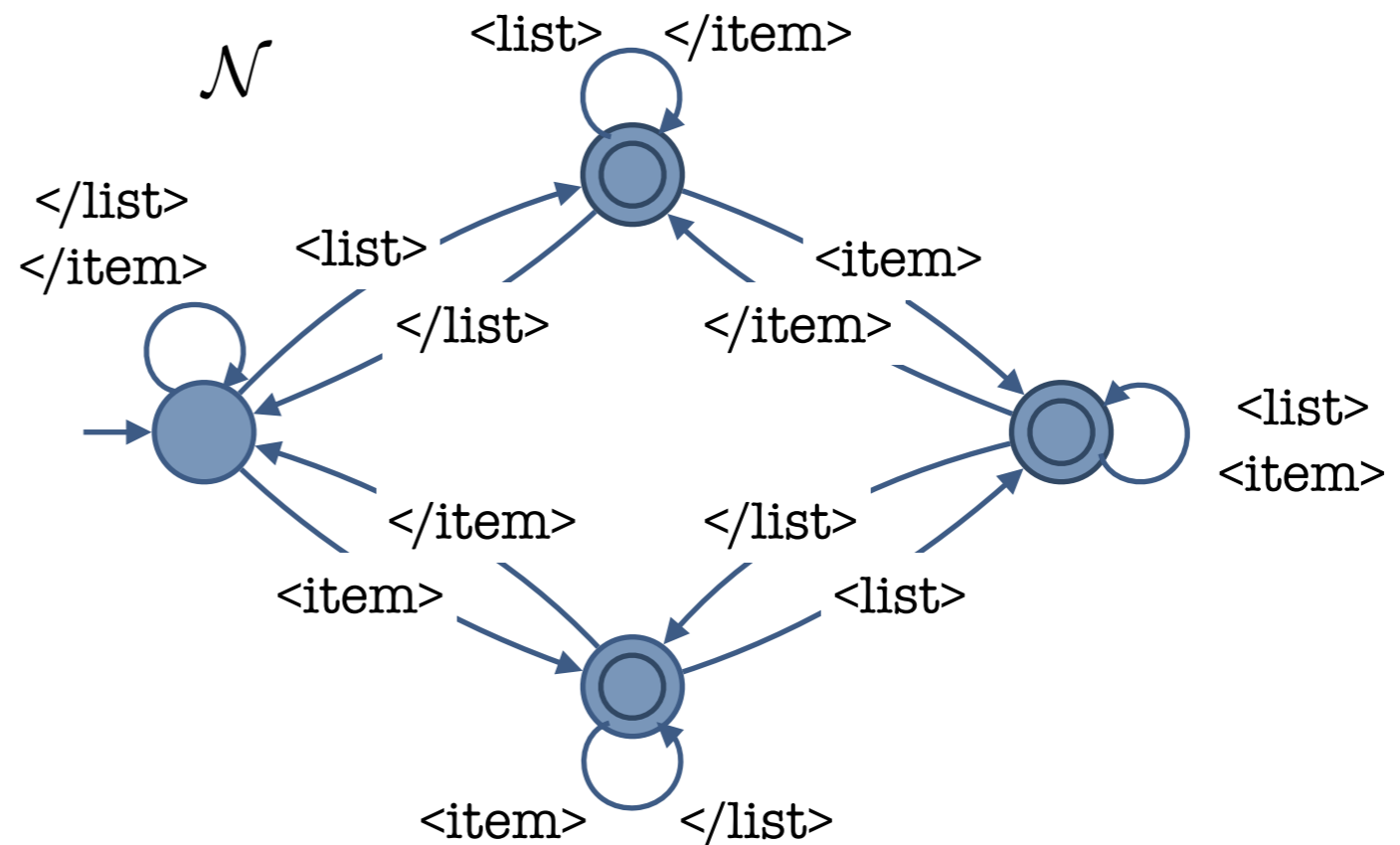
Regular negative specification \mathcal{N} :

"there is an opening tag that is not eventually followed by a corresponding closing tag"

```

<list>
  <item>
    <list>
      <item>
        </item>
      <item>
        </item>
    </list>
  </item>
  <item>
    </item>
  </list>

```



Assume \mathcal{R} is supposed to recognize XML documents over

$$\Sigma = \{ \langle \text{list} \rangle, \langle / \text{list} \rangle, \langle \text{item} \rangle, \langle / \text{item} \rangle \}.$$

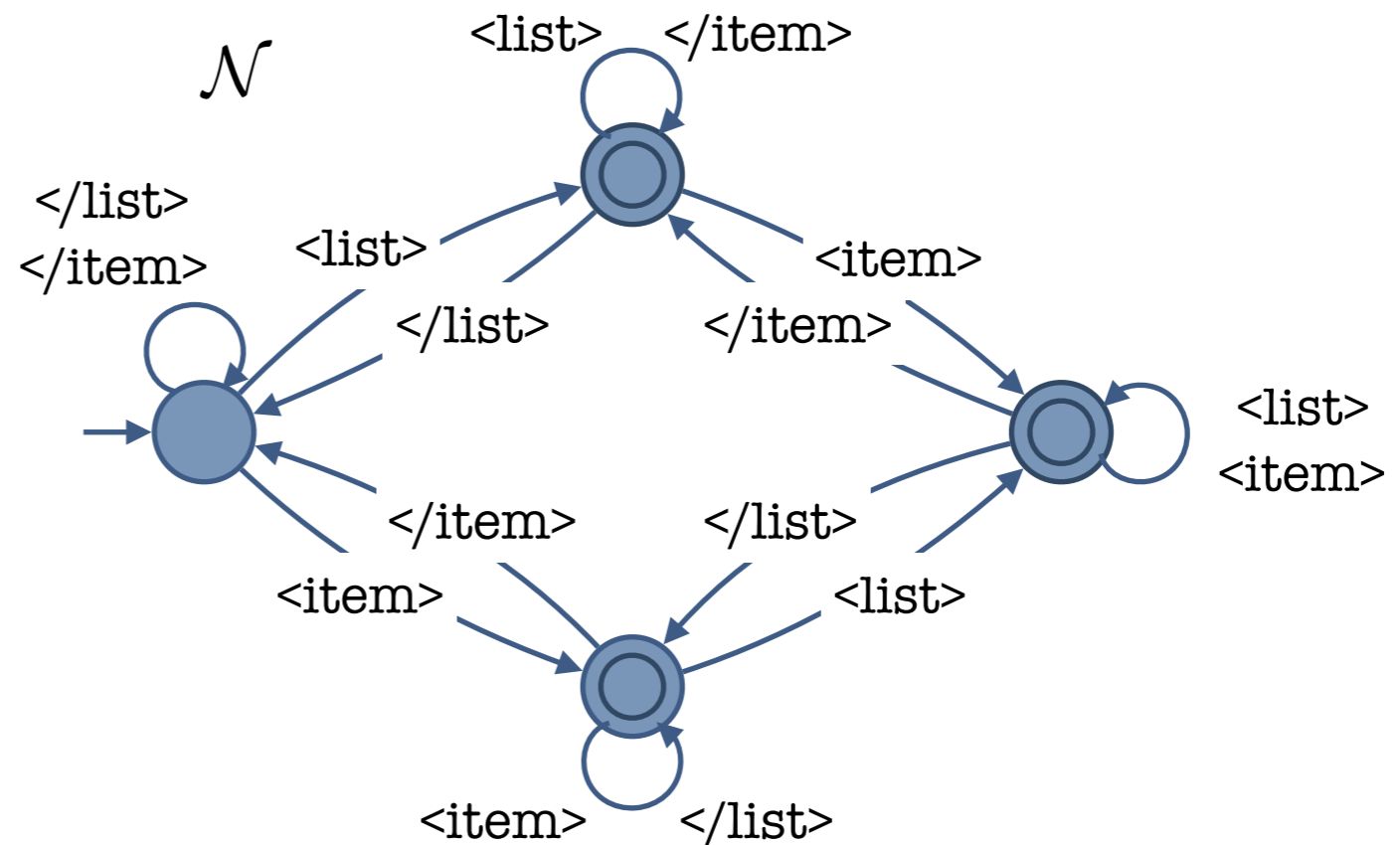
Regular negative specification \mathcal{N} :

"there is an opening tag that is not eventually followed by a corresponding closing tag"

```

<list>
  <item>
    <list>
      <item>
      </item>
      <item>
      </item>
    </list>
  </item>
  <item>
  </item>
</list>

```



$\langle \text{list} \rangle \langle \text{item} \rangle \langle / \text{item} \rangle \langle \text{item} \rangle \langle / \text{list} \rangle \in L(\mathcal{N})$

Assume \mathcal{R} is supposed to recognize XML documents over

$$\Sigma = \{ \langle \text{list} \rangle, \langle / \text{list} \rangle, \langle \text{item} \rangle, \langle / \text{item} \rangle \}.$$

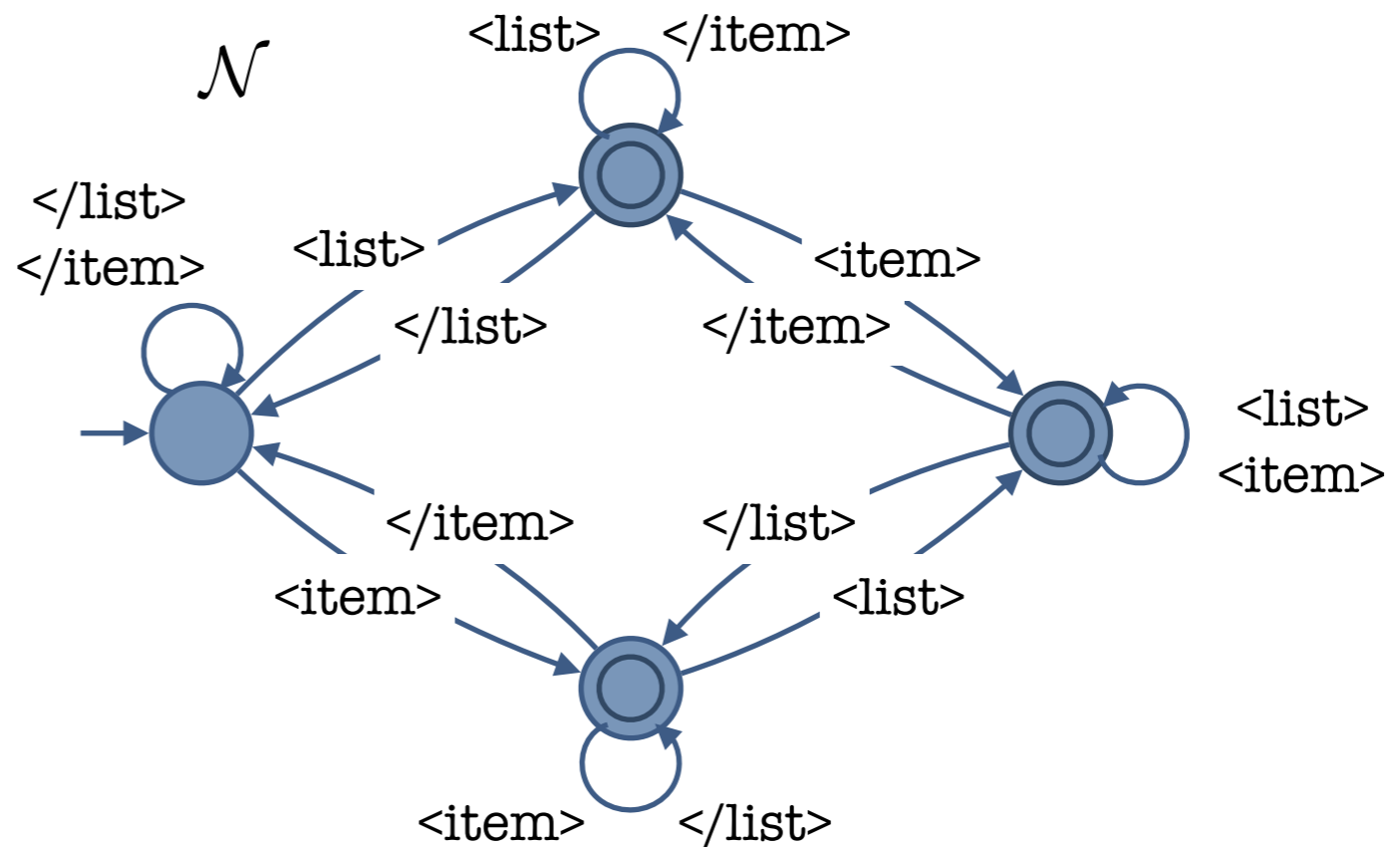
Regular negative specification \mathcal{N} :

"there is an opening tag that is not eventually followed by a corresponding closing tag"

```

<list>
  <item>
    <list>
      <item>
      </item>
      <item>
      </item>
    </list>
  </item>
  <item>
  </item>
</list>

```



$\langle \text{list} \rangle \langle \text{item} \rangle \langle / \text{item} \rangle \langle \text{item} \rangle \langle / \text{list} \rangle \in L(\mathcal{N})$

$\langle \text{list} \rangle \langle \text{item} \rangle \langle \text{item} \rangle \langle / \text{item} \rangle \langle / \text{list} \rangle \notin L(\mathcal{N})$

Assume \mathcal{R} is supposed to recognize XML documents over

$$\Sigma = \{ \langle \text{list} \rangle, \langle / \text{list} \rangle, \langle \text{item} \rangle, \langle / \text{item} \rangle \}.$$

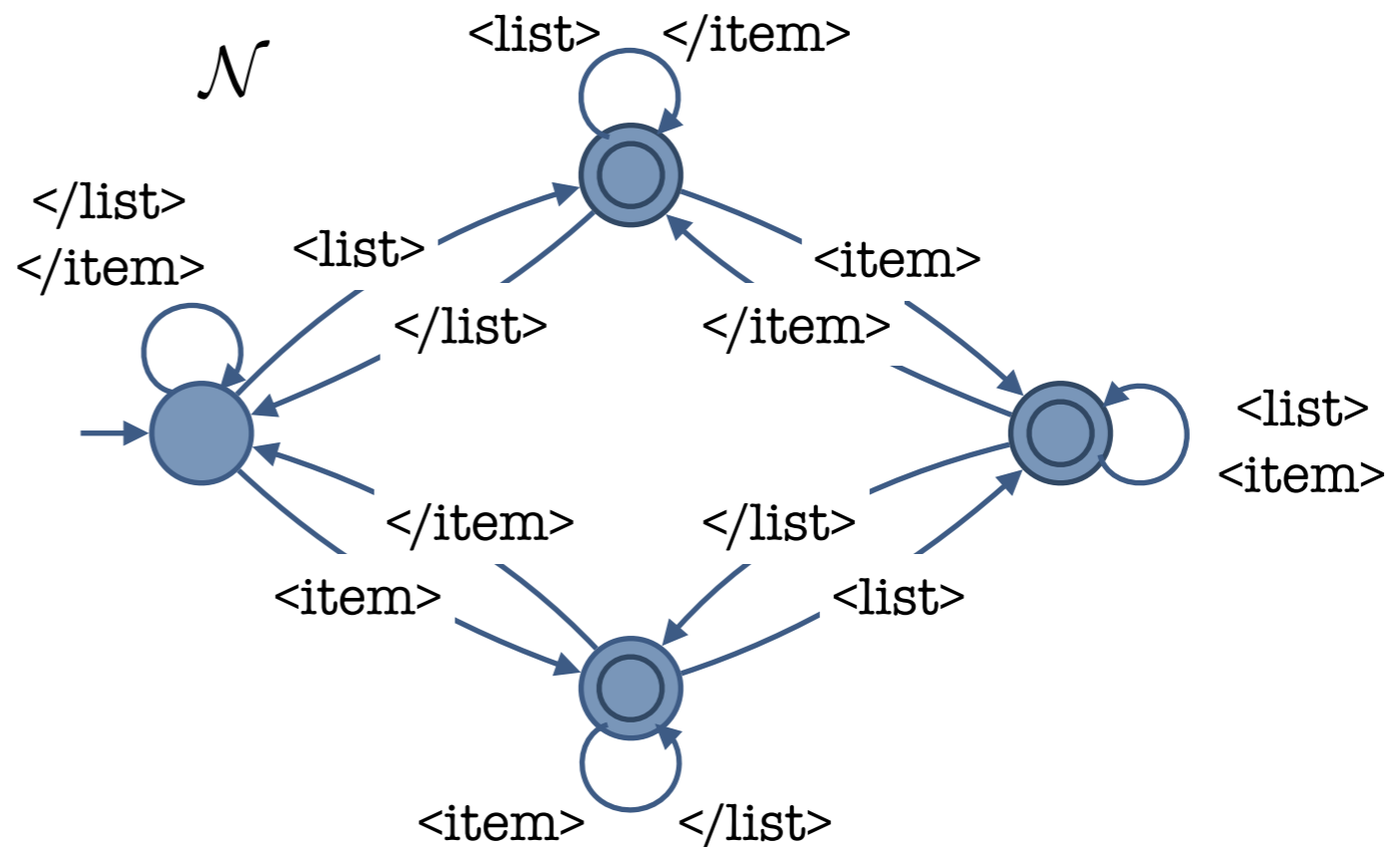
Regular negative specification \mathcal{N} :

"there is an opening tag that is not eventually followed by a corresponding closing tag"

```

<list>
  <item>
    <list>
      <item>
      </item>
      <item>
      </item>
    </list>
  </item>
  <item>
  </item>
</list>

```



We check $L(\mathcal{R}) \subseteq L(\overline{\mathcal{N}})$.

$\langle \text{list} \rangle \langle \text{item} \rangle \langle / \text{item} \rangle \langle \text{item} \rangle \langle / \text{list} \rangle \in L(\mathcal{N})$

$\langle \text{list} \rangle \langle \text{item} \rangle \langle \text{item} \rangle \langle / \text{item} \rangle \langle / \text{list} \rangle \notin L(\mathcal{N})$

Assume \mathcal{R} is supposed to recognize XML documents over

$$\Sigma = \{ \langle \text{list} \rangle, \langle / \text{list} \rangle, \langle \text{item} \rangle, \langle / \text{item} \rangle \}.$$

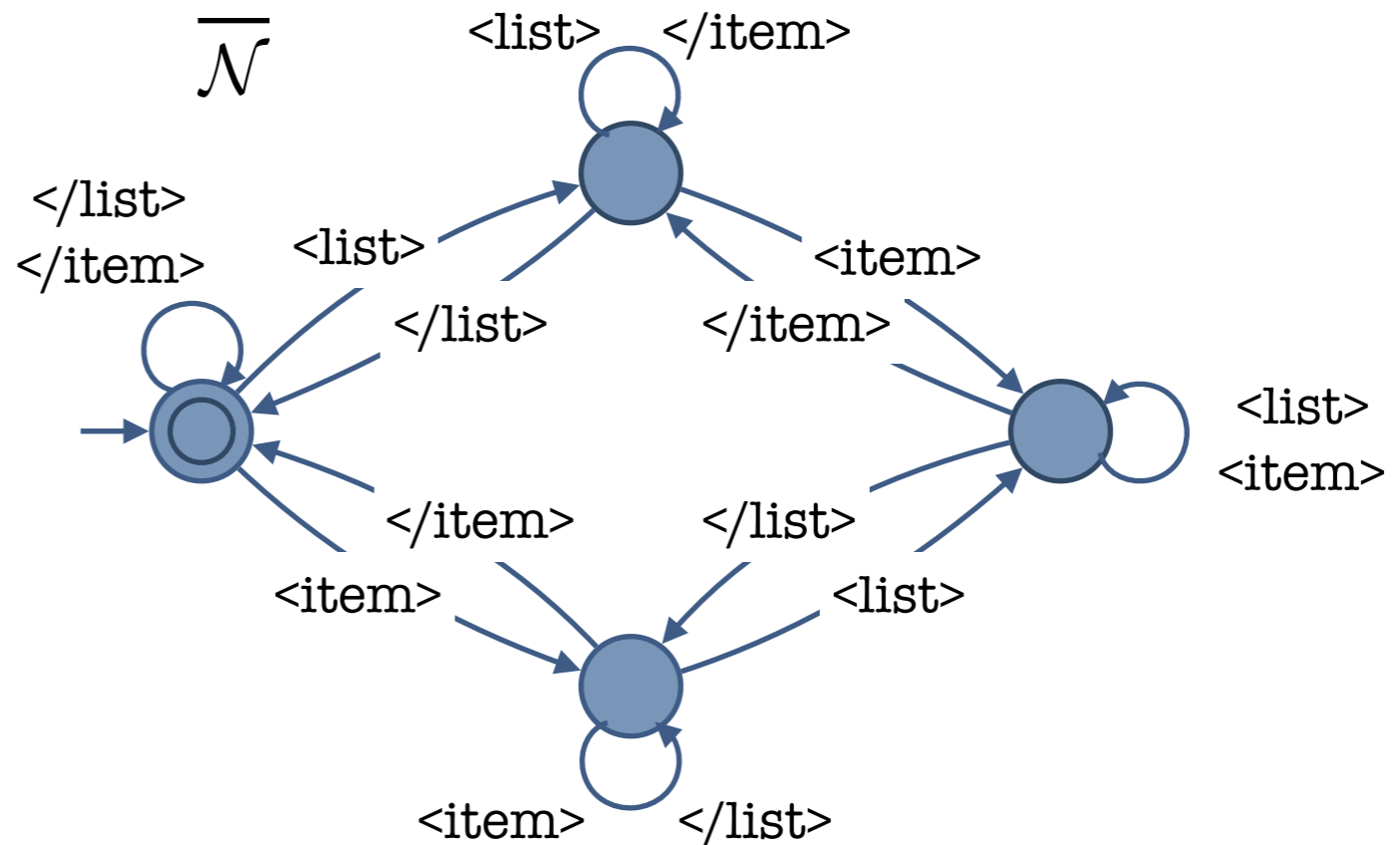
Regular negative specification \mathcal{N} :

"there is an opening tag that is not eventually followed by a corresponding closing tag"

```

<list>
  <item>
    <list>
      <item>
      </item>
      <item>
      </item>
    </list>
  </item>
  <item>
  </item>
</list>

```



We check $L(\mathcal{R}) \subseteq L(\overline{\mathcal{N}})$.

$\langle \text{list} \rangle \langle \text{item} \rangle \langle / \text{item} \rangle \langle \text{item} \rangle \langle / \text{list} \rangle \in L(\mathcal{N})$

$\langle \text{list} \rangle \langle \text{item} \rangle \langle \text{item} \rangle \langle / \text{item} \rangle \langle / \text{list} \rangle \notin L(\mathcal{N})$

Checking $L(\mathcal{R}) \subseteq L(\mathcal{A})$ for deterministic finite automaton \mathcal{A}

Checking $L(\mathcal{R}) \subseteq L(\mathcal{A})$ for deterministic finite automaton \mathcal{A}

- Statistical model checking (SMC)

Checking $L(\mathcal{R}) \subseteq L(\mathcal{A})$ for deterministic finite automaton \mathcal{A}

- Statistical model checking (SMC)
- Automaton abstraction & model checking (AAMC)
(model-learning approach)

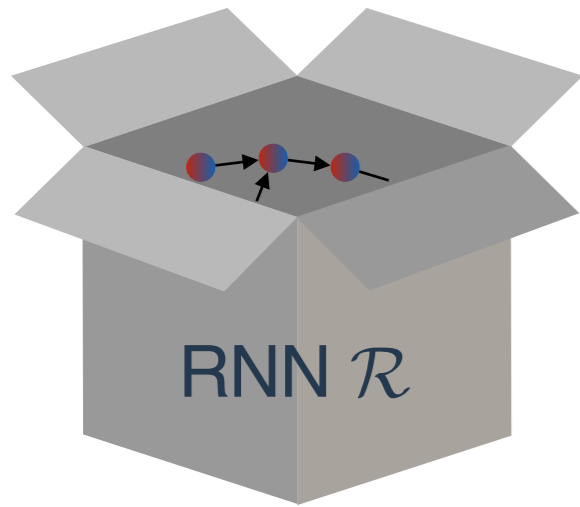
Checking $L(\mathcal{R}) \subseteq L(\mathcal{A})$ for deterministic finite automaton \mathcal{A}

- Statistical model checking (SMC)
- Automaton abstraction & model checking (AAMC)
(model-learning approach)
- Property-directed verification (PDV)

Statistical model checking (SMC)

Fix $\varepsilon, \gamma > 0$.

Sample $\log(2 / \varepsilon) / (2\gamma^2)$ words over Σ .



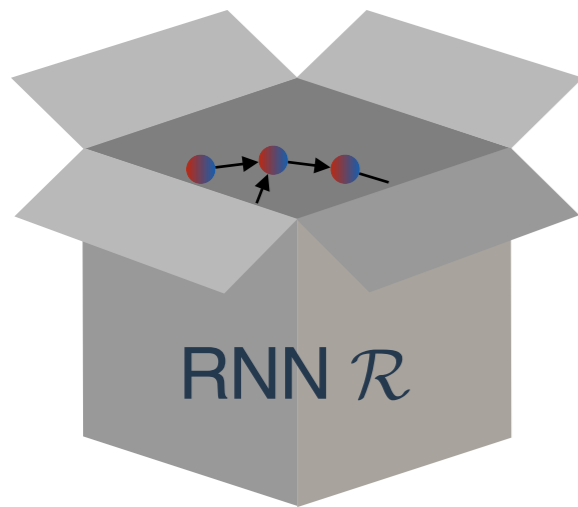
Statistical model checking (SMC)

Fix $\varepsilon, \gamma > 0$.

Sample $\log(2 / \varepsilon) / (2\gamma^2)$ words over Σ .

If, for some word w , we have $w \in L(\mathcal{R})$ and $w \notin L(\mathcal{A})$:

Property not satisfied.



Statistical model checking (SMC)

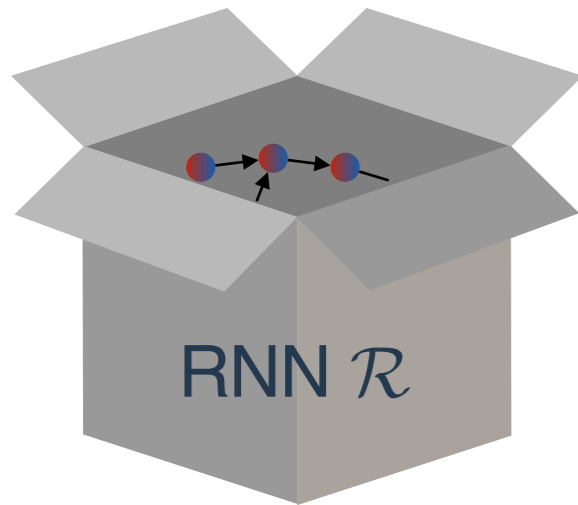
Fix $\varepsilon, \gamma > 0$.

Sample $\log(2 / \varepsilon) / (2\gamma^2)$ words over Σ .

If, for some word w , we have $w \in L(\mathcal{R})$ and $w \notin L(\mathcal{A})$:

Property not satisfied.

Else, \mathcal{R} is ε -approximately correct with probability at least $1 - \gamma$.



Statistical model checking (SMC)

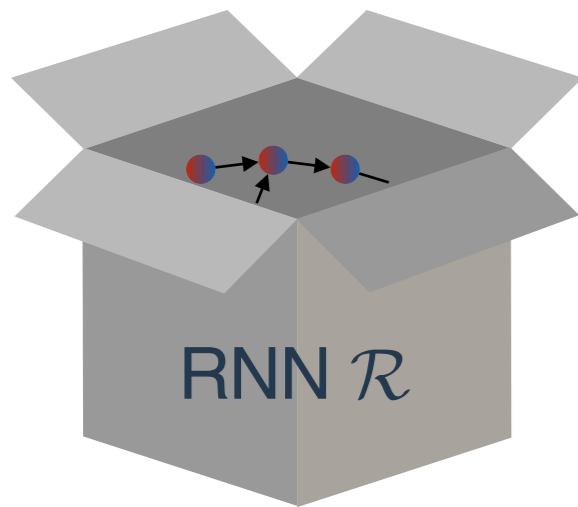
Fix $\varepsilon, \gamma > 0$.

Sample $\log(2 / \varepsilon) / (2\gamma^2)$ words over Σ .

If, for some word w , we have $w \in L(\mathcal{R})$ and $w \notin L(\mathcal{A})$:

Property not satisfied.

Else, \mathcal{R} is ε -approximately correct with probability at least $1 - \gamma$.



$$\Pr(L(\mathcal{R}) \setminus L(\mathcal{A})) < \varepsilon$$

Statistical model checking (SMC)

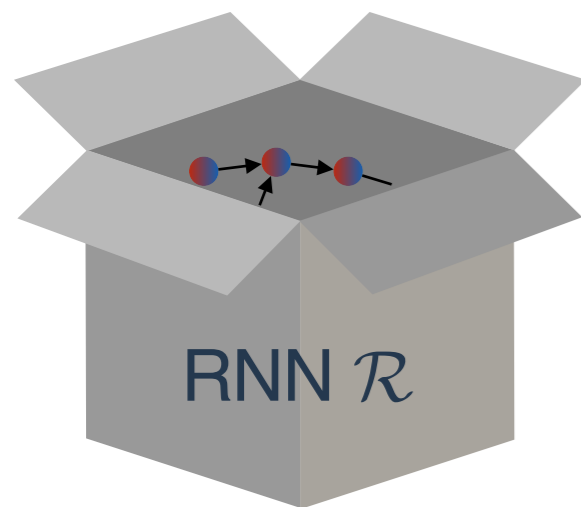
Fix $\varepsilon, \gamma > 0$.

Sample $\log(2 / \varepsilon) / (2\gamma^2)$ words over Σ .

If, for some word w , we have $w \in L(\mathcal{R})$ and $w \notin L(\mathcal{A})$:

Property not satisfied.

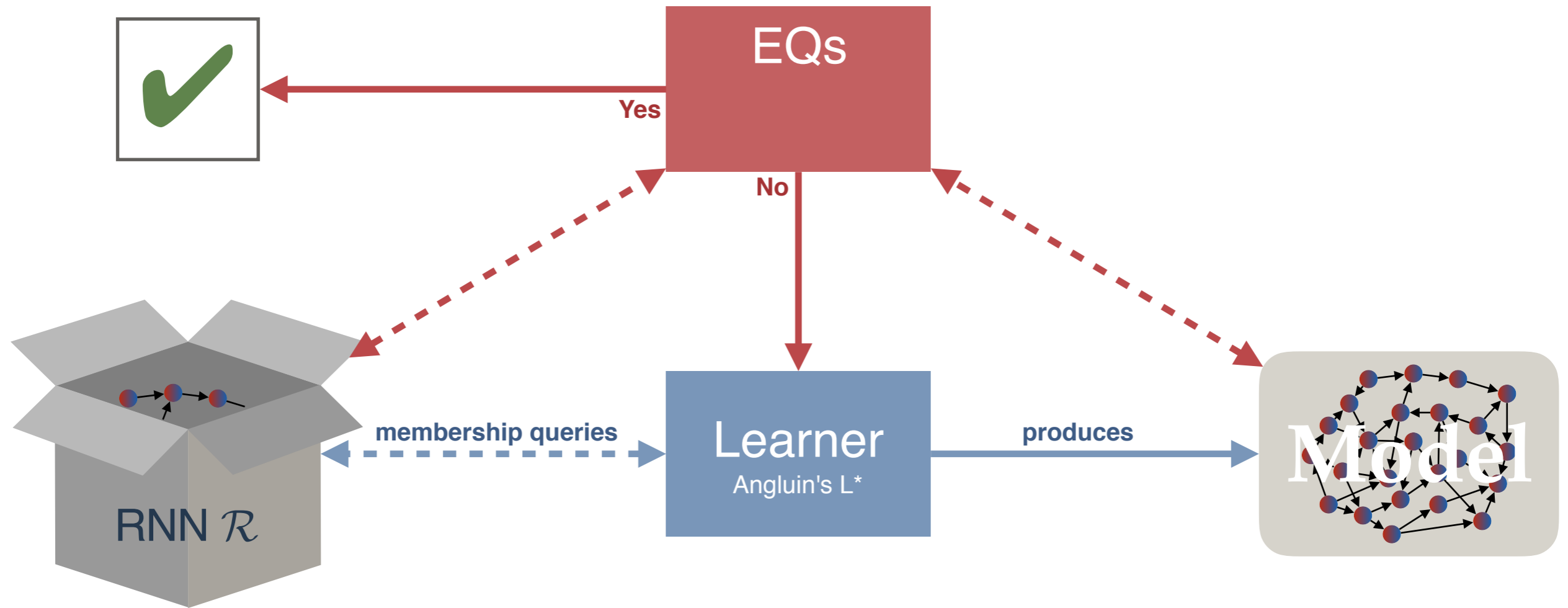
Else, \mathcal{R} is ε -approximately correct with probability at least $1 - \gamma$.



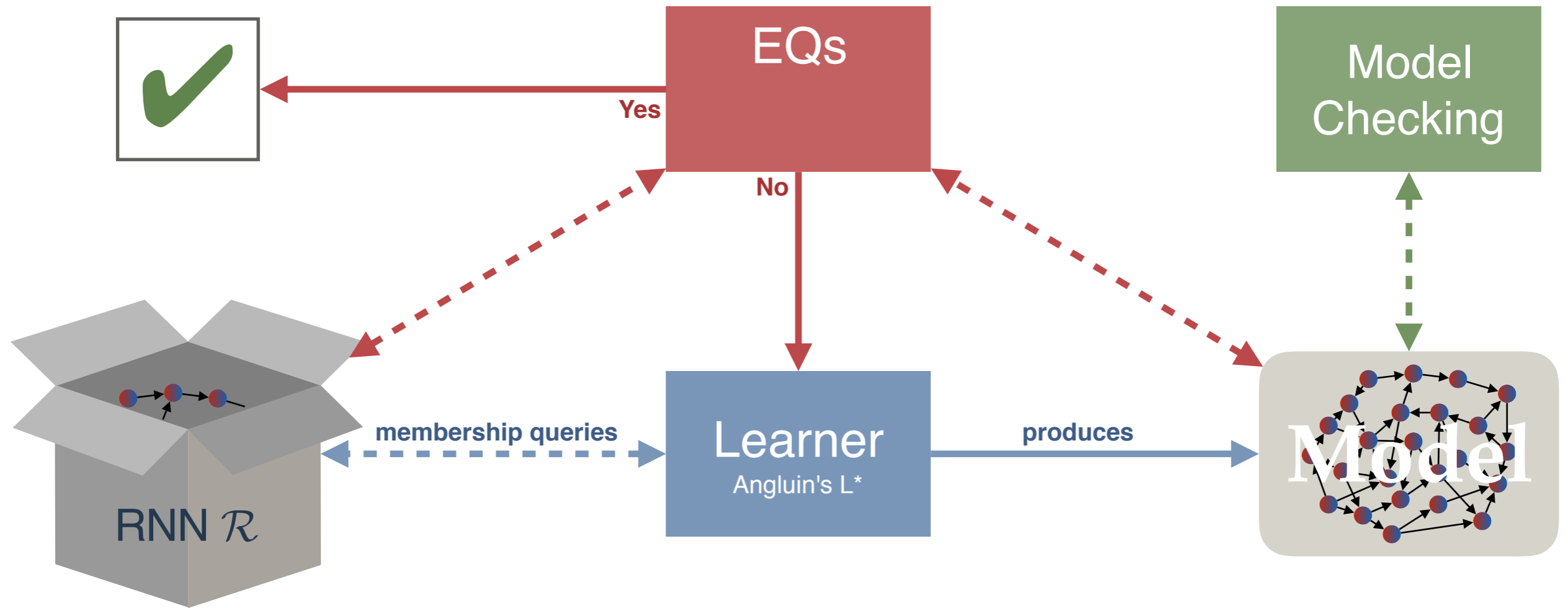
$$\Pr(L(\mathcal{R}) \setminus L(\mathcal{A})) < \varepsilon$$

- Relies on probability distribution.
- May require many queries.

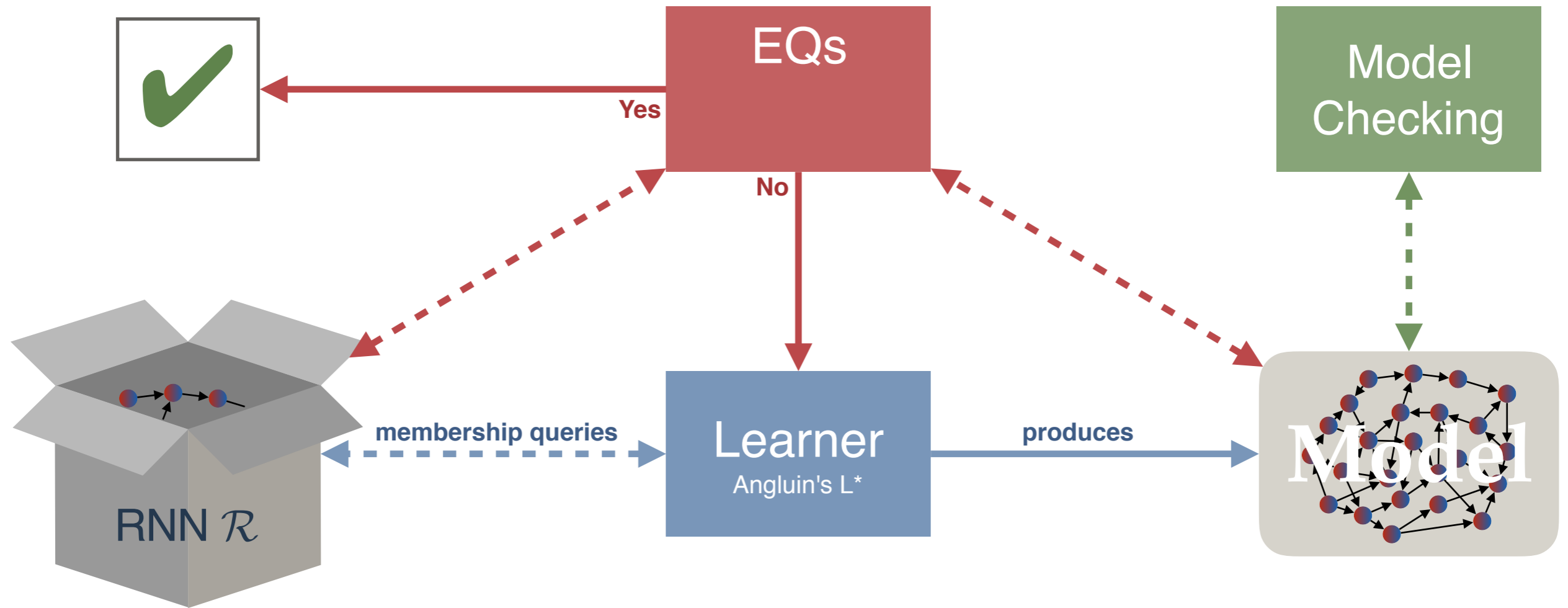
Automaton abstraction & model checking (AAMC)



Automaton abstraction & model checking (AAMC)

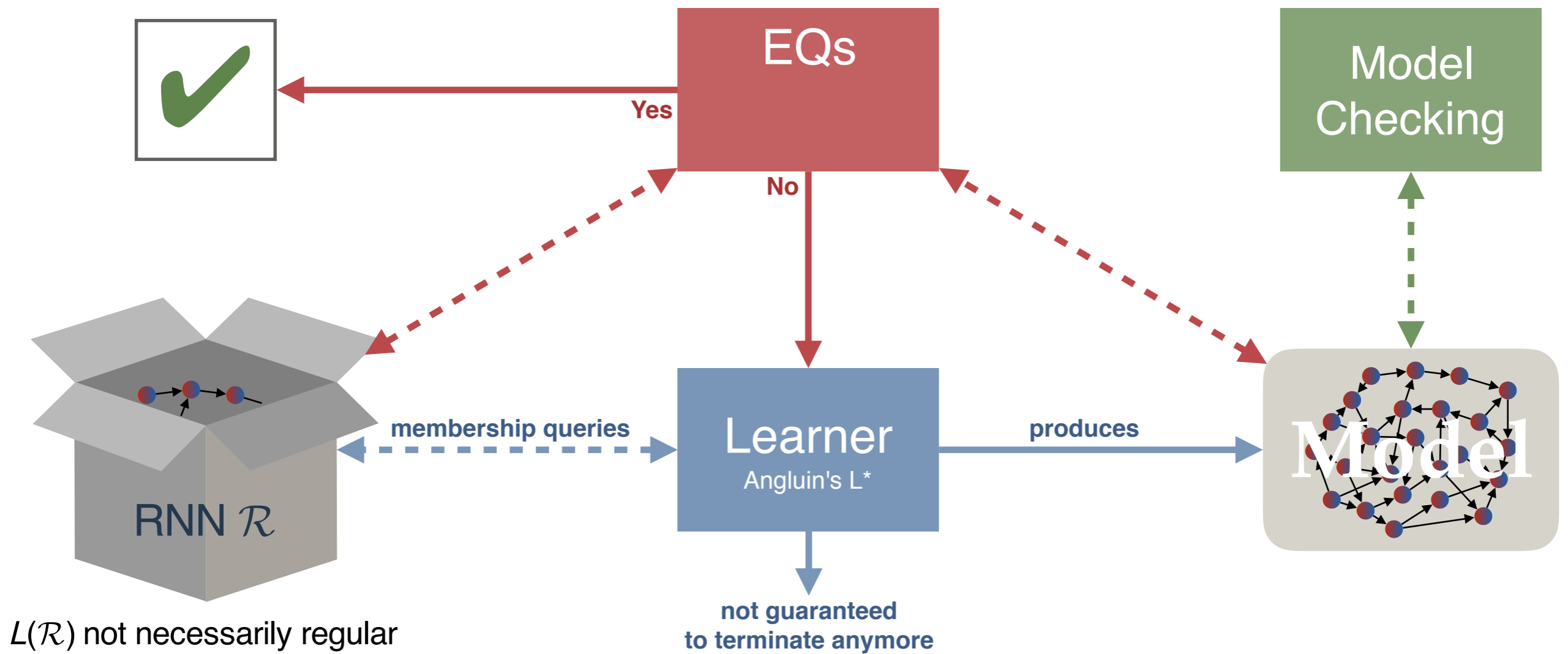


Automaton abstraction & model checking (AAMC)

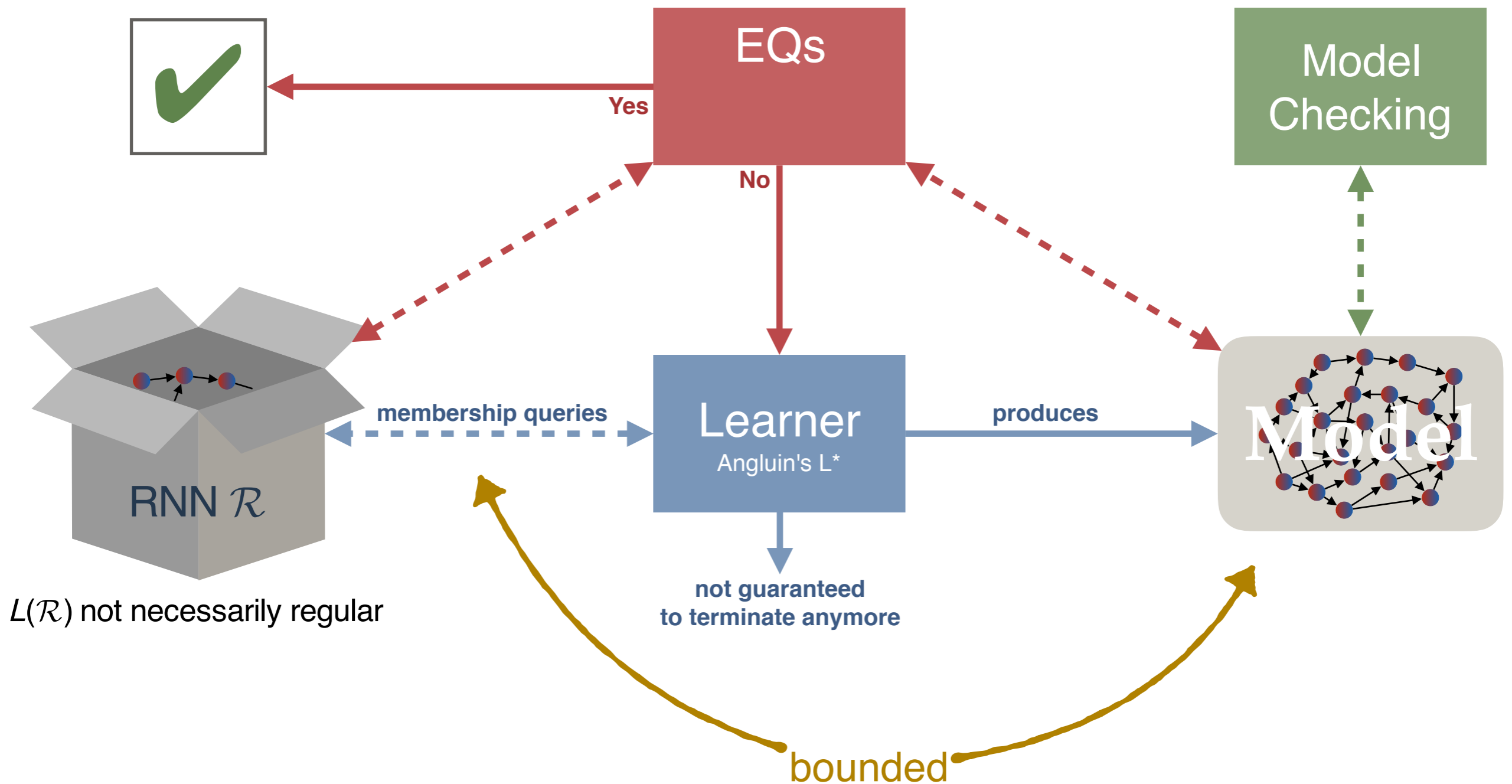


$L(\mathcal{R})$ not necessarily regular

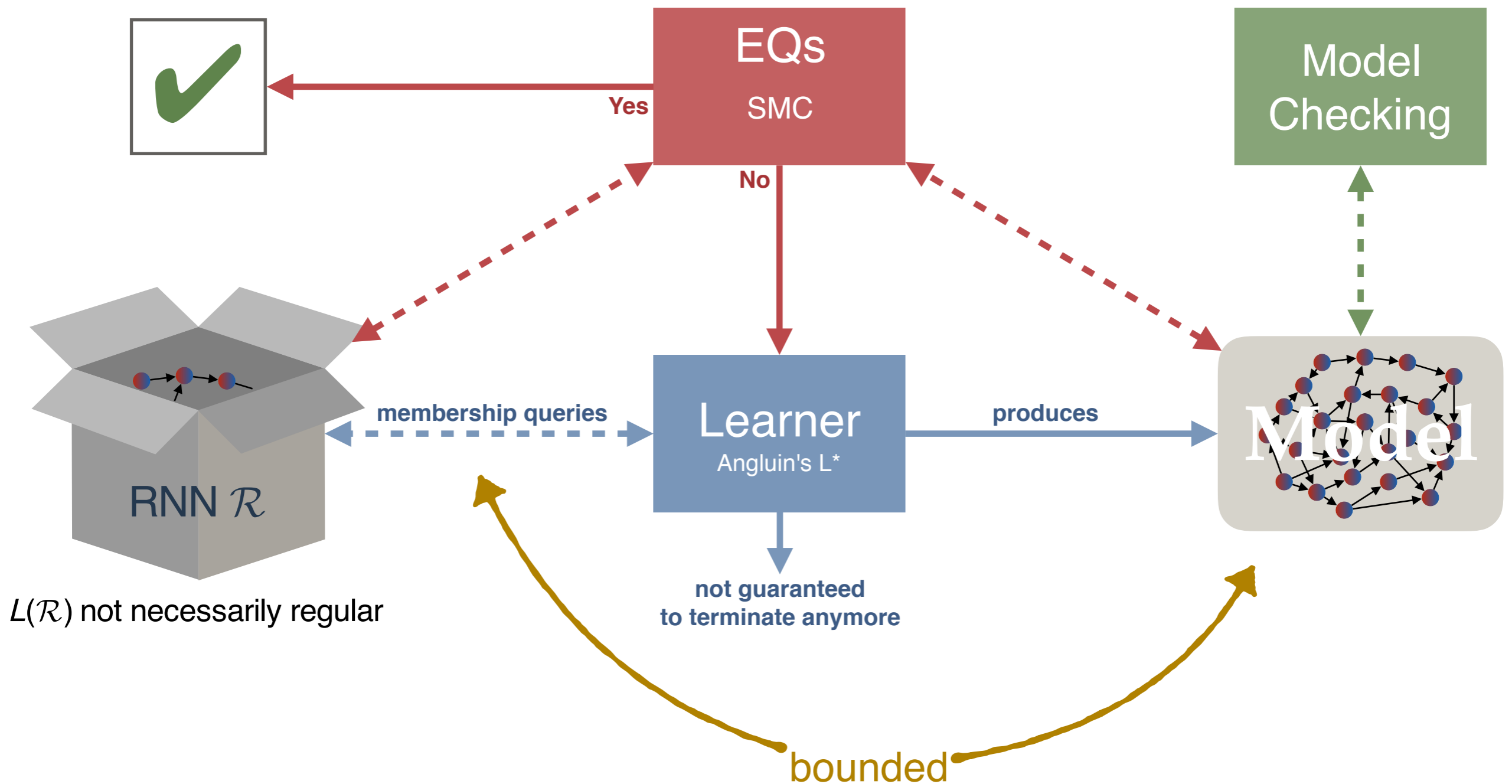
Automaton abstraction & model checking (AAMC)



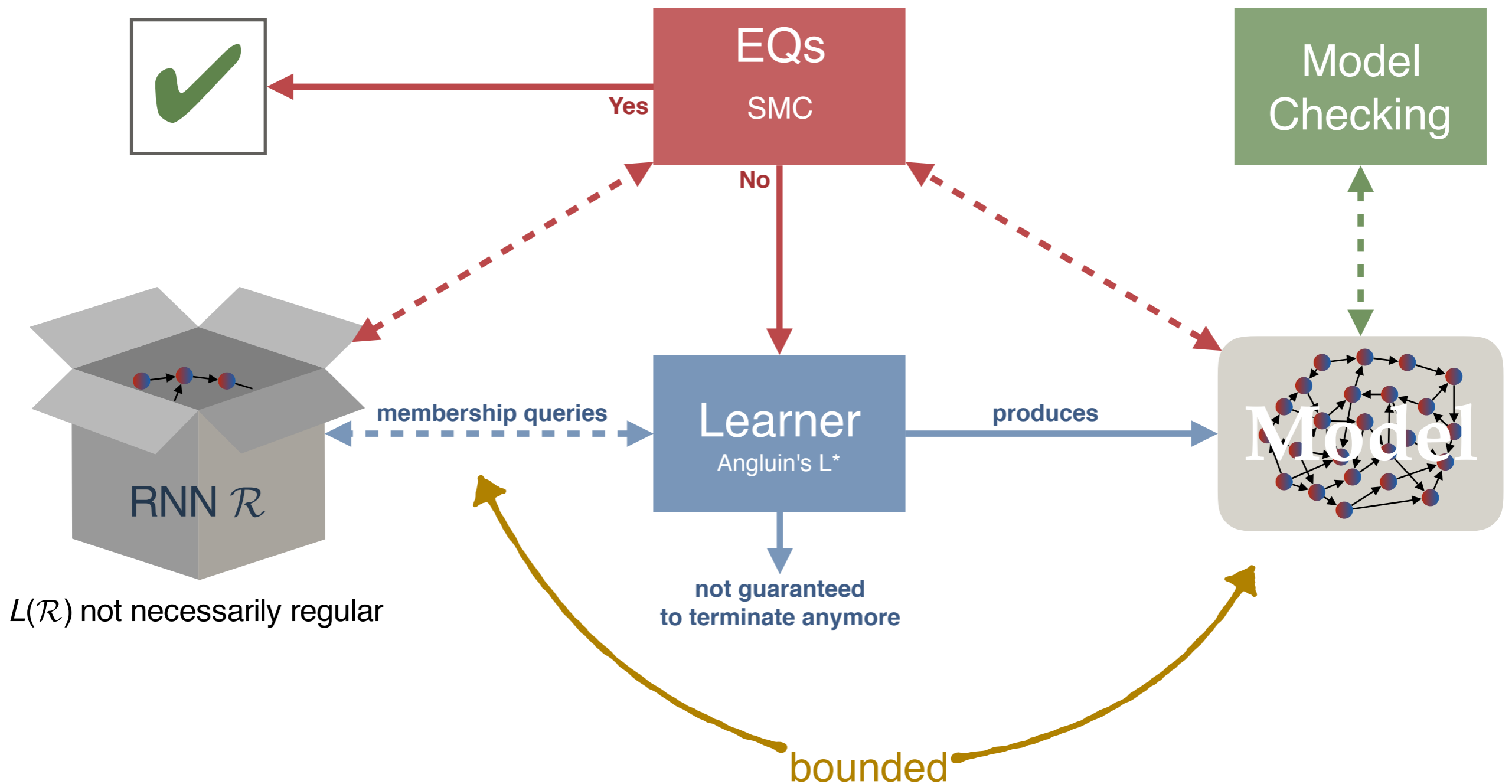
Automaton abstraction & model checking (AAMC)



Automaton abstraction & model checking (AAMC)

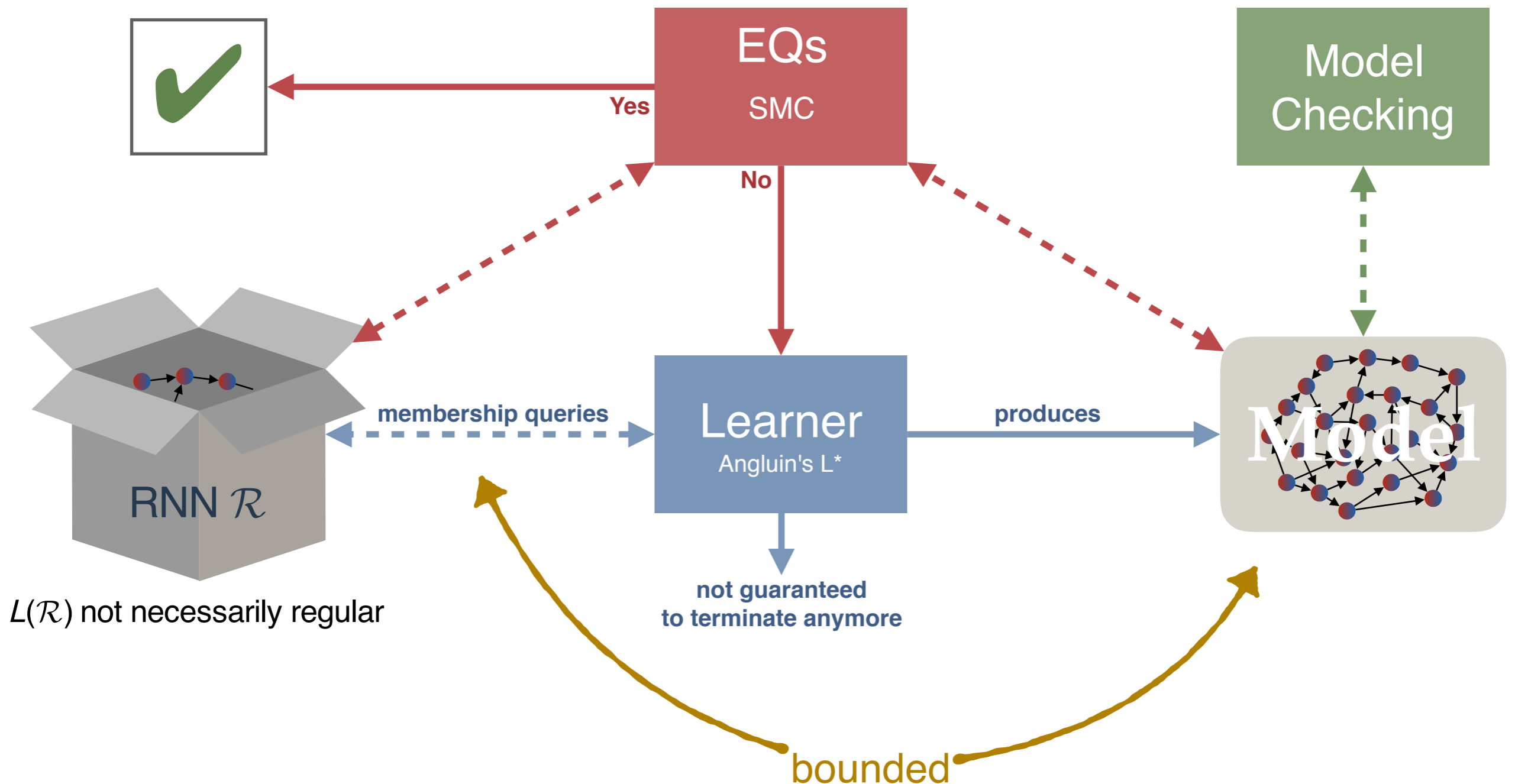


Automaton abstraction & model checking (AAMC)



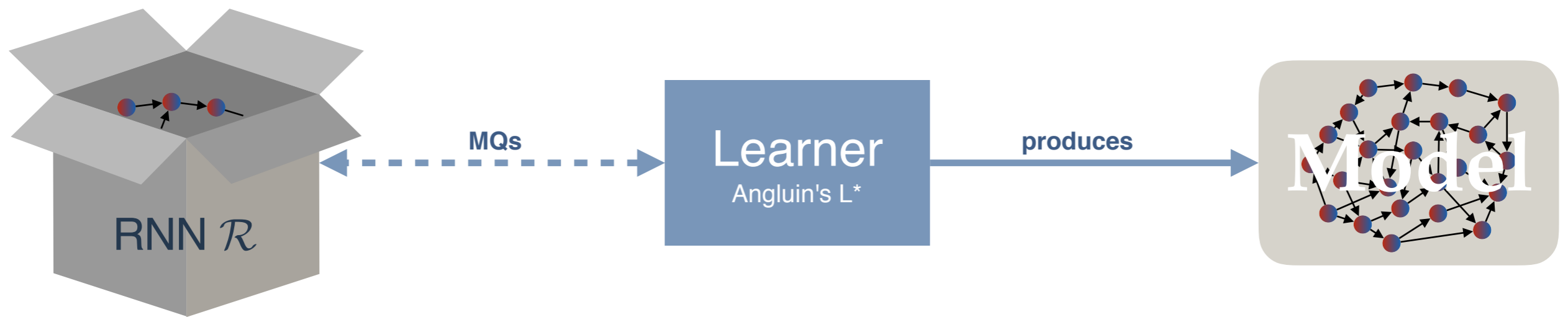
- Relies on probability distribution.

Automaton abstraction & model checking (AAMC)

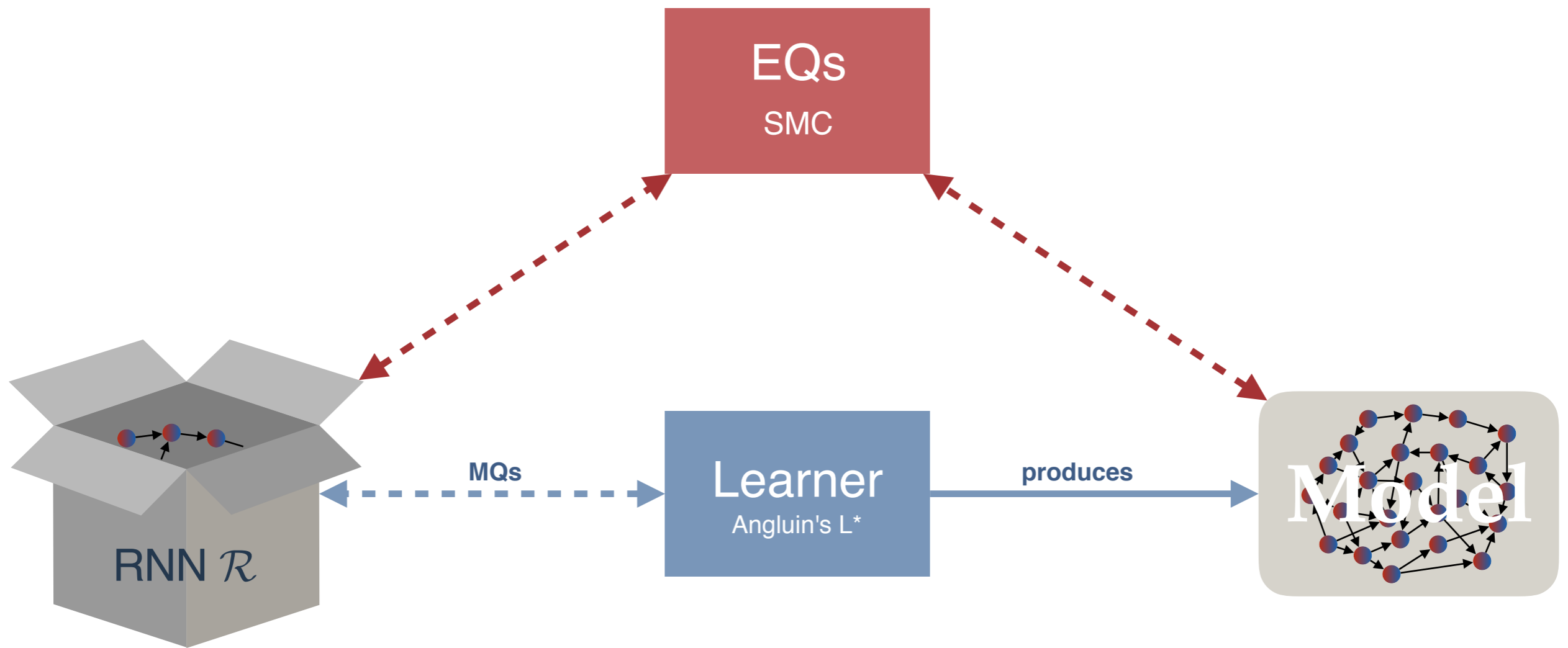


- Relies on probability distribution.
- Counterexample may be spurious.

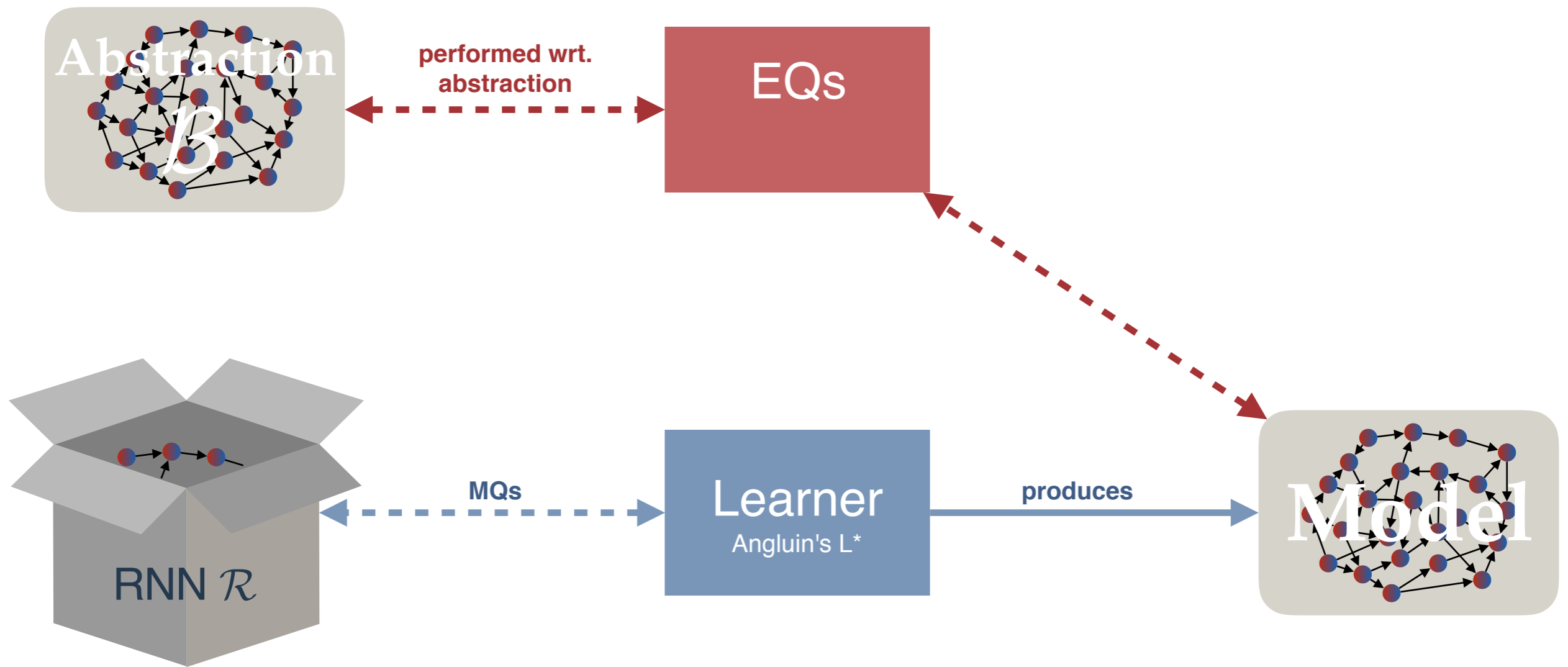
As EQs are implemented by SMC, a "counterexample" may be classified by RNN as negative.



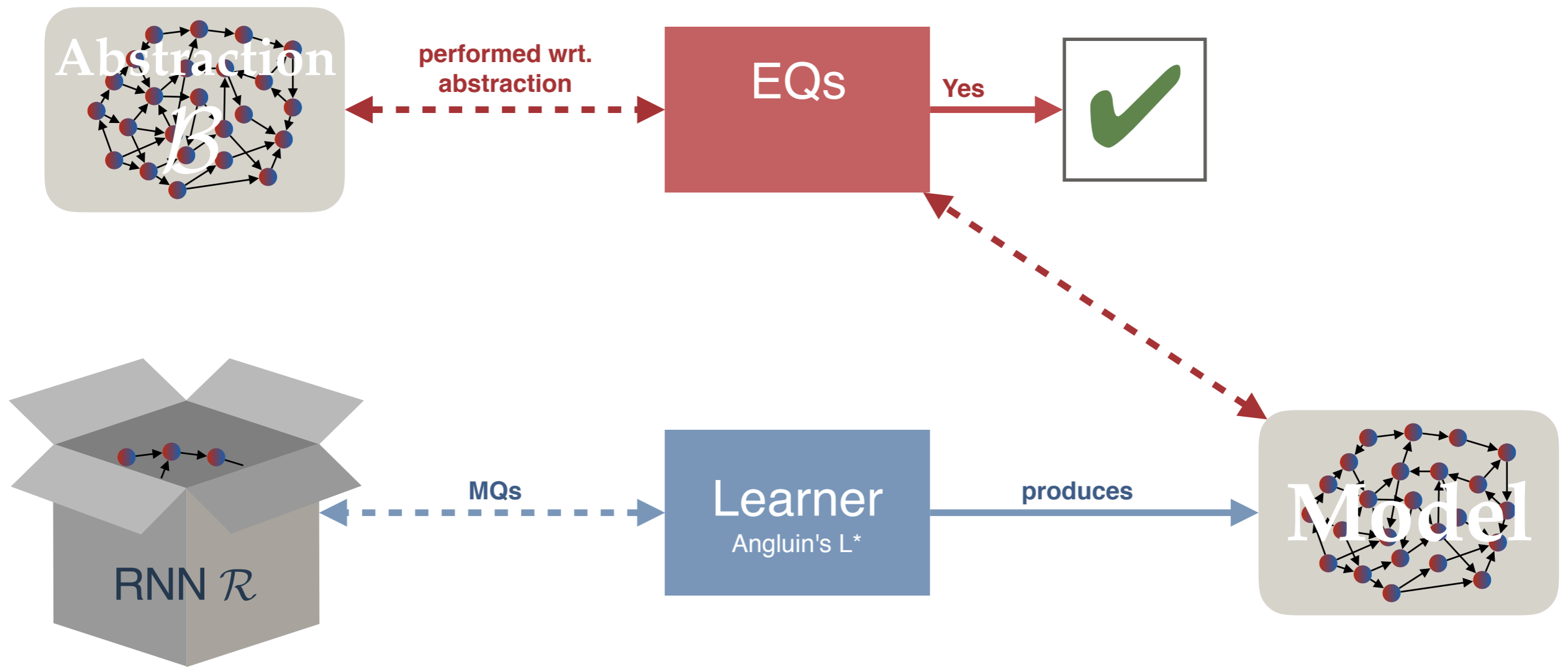
[Weiss, Goldberg, Yahav: *Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples*. ICML 2018]



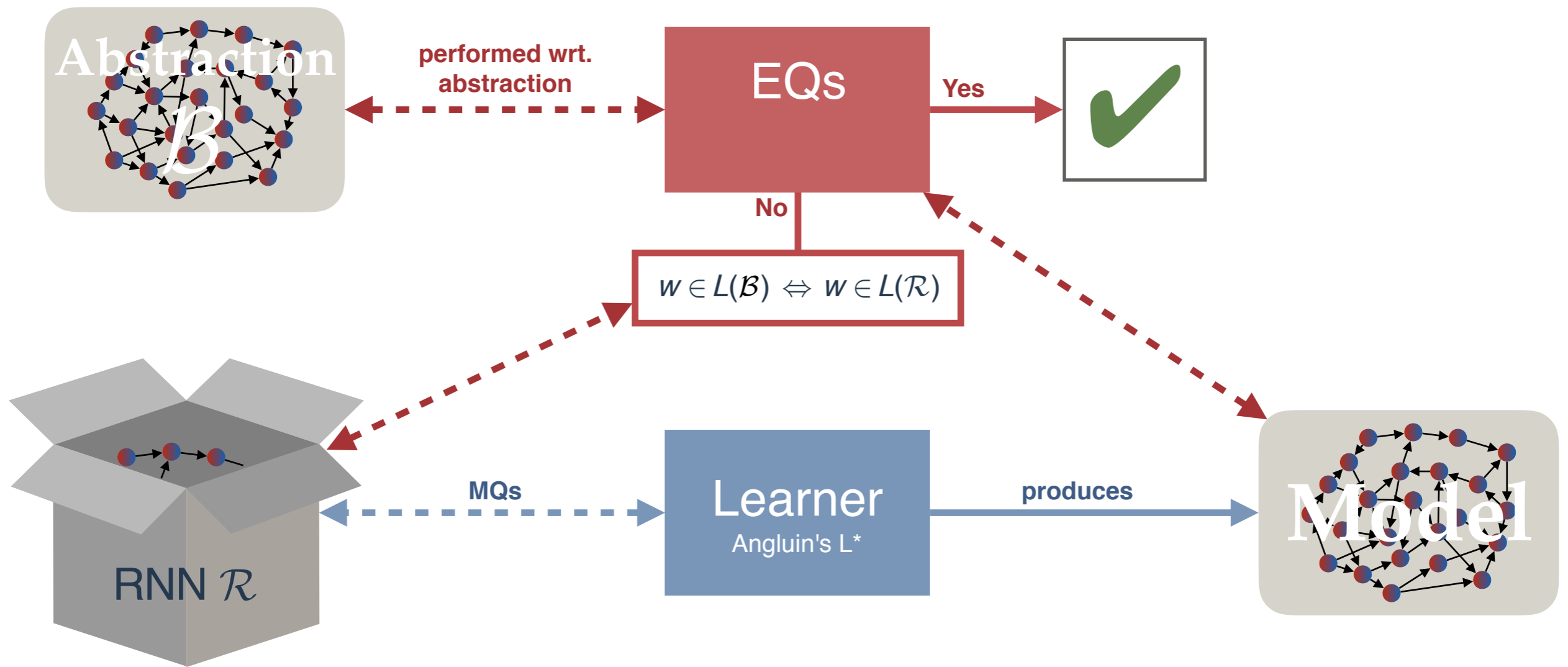
[Weiss, Goldberg, Yahav: *Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples*. ICML 2018]



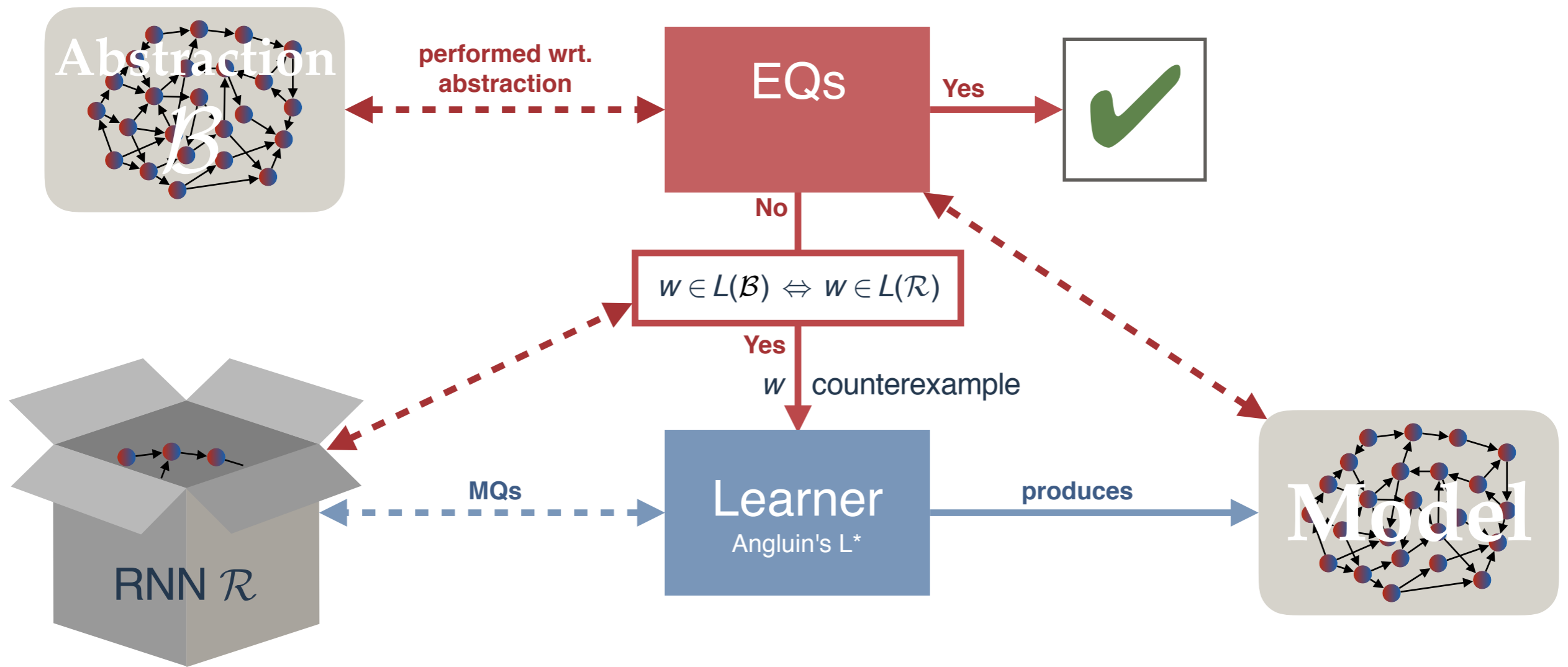
[Weiss, Goldberg, Yahav: *Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples*. ICML 2018]



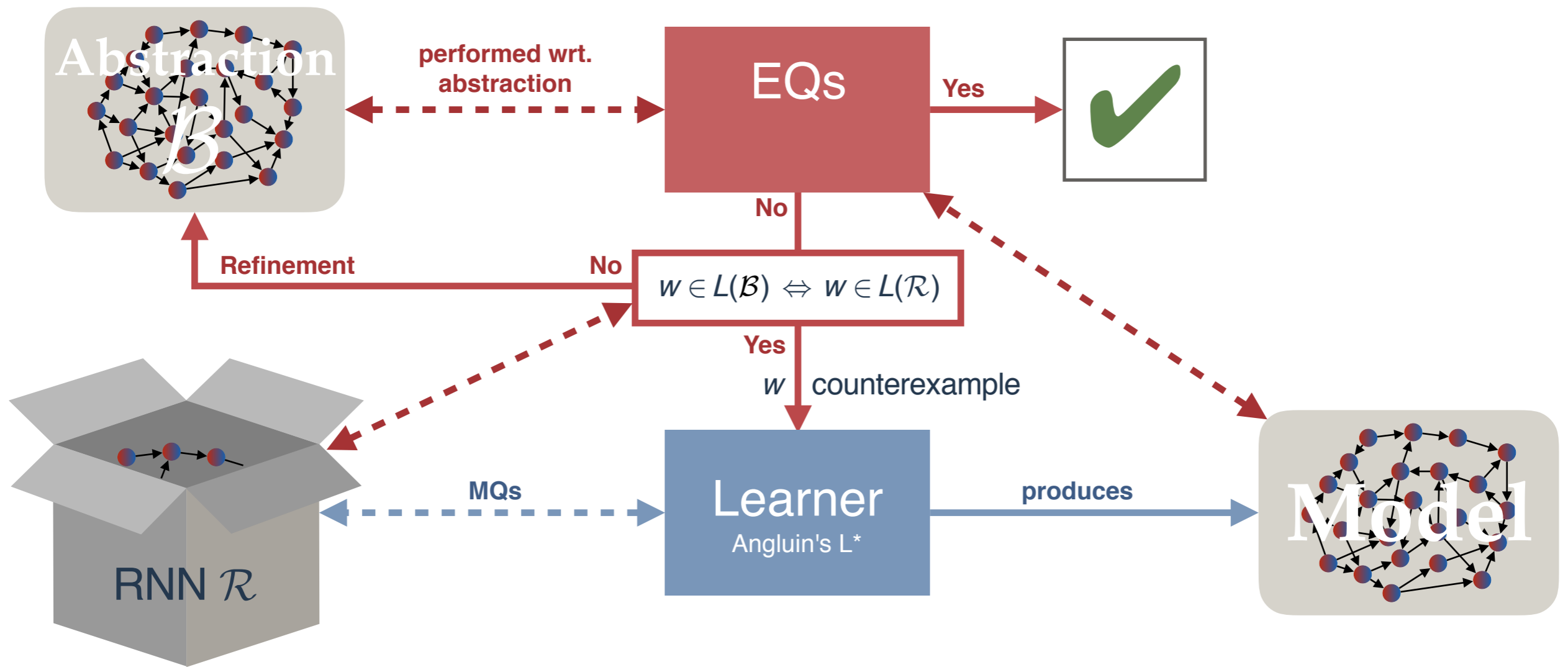
[Weiss, Goldberg, Yahav: *Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples*. ICML 2018]



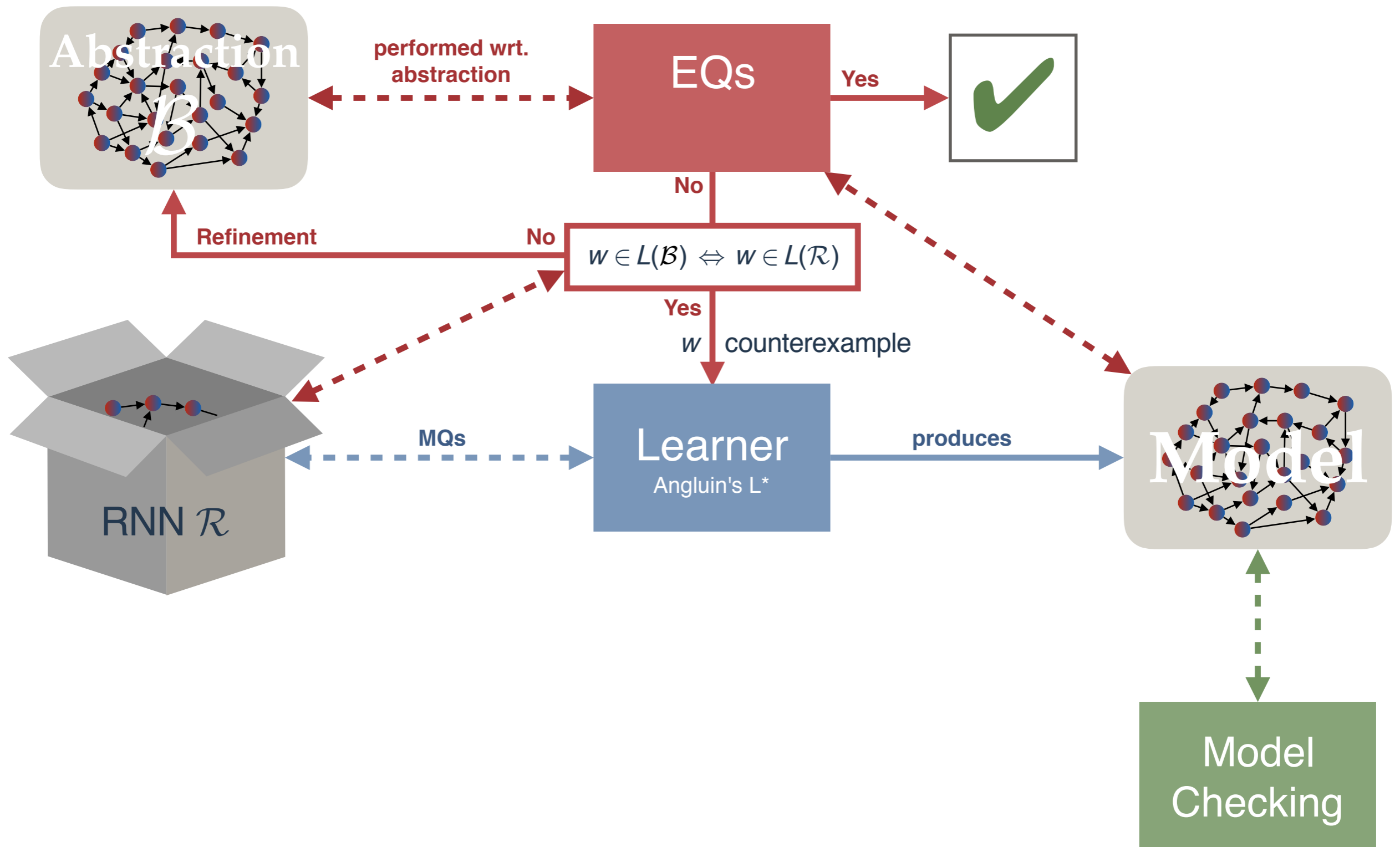
[Weiss, Goldberg, Yahav: *Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples*. ICML 2018]



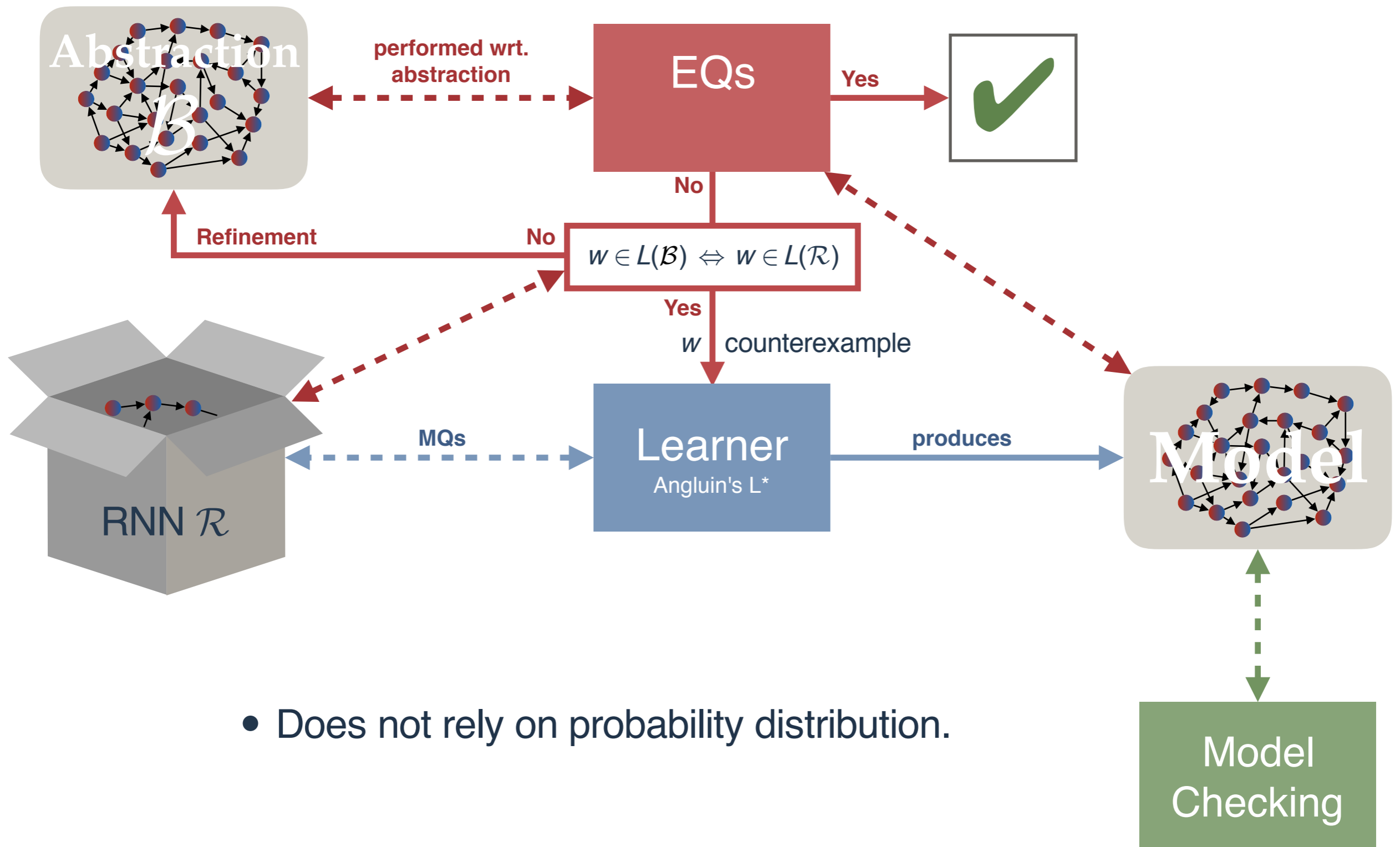
[Weiss, Goldberg, Yahav: *Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples*. ICML 2018]



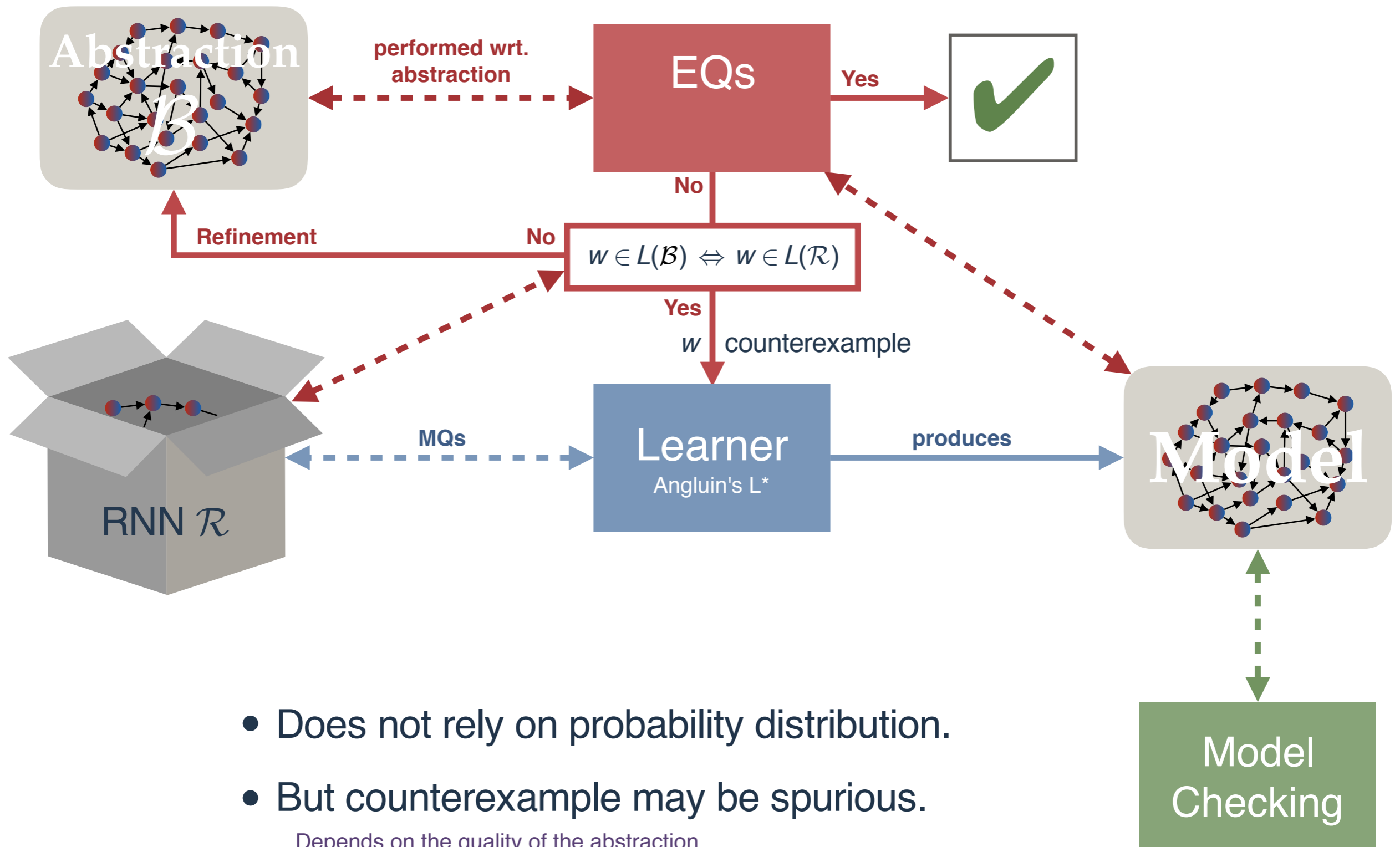
[Weiss, Goldberg, Yahav: *Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples*. ICML 2018]



[Weiss, Goldberg, Yahav: *Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples*. ICML 2018]



[Weiss, Goldberg, Yahav: *Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples*. ICML 2018]



- Does not rely on probability distribution.
- But counterexample may be spurious.
Depends on the quality of the abstraction.

[Weiss, Goldberg, Yahav: *Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples*. ICML 2018]

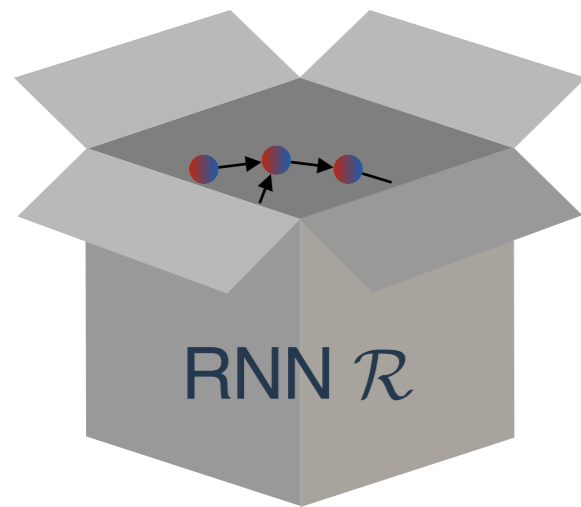
Checked on 30 DFAs / RNNs and 138 specifications:

| Type | <i>Avg time</i> (s) | <i>Avg len</i> | <i># Mistakes</i> | <i>Avg MQs</i> |
|------|---------------------|----------------|-------------------|----------------|
| SMC | 92 | 111 | 122 | 286063 |
| AAMC | 444 | 7 | 30 | 3701916 |

Checked on 30 DFAs / RNNs and 138 specifications:

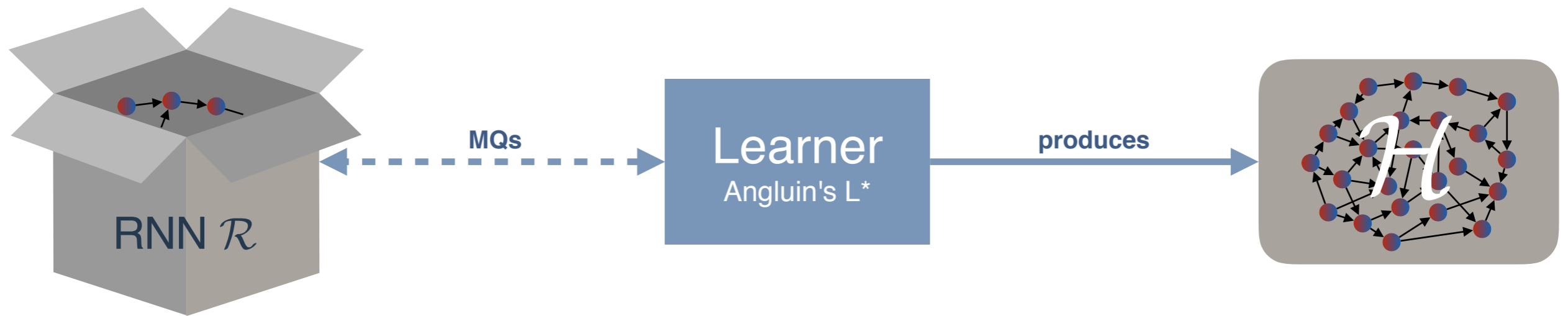
| Type | <i>Avg time</i> (s) | <i>Avg len</i> | <i># Mistakes</i> | <i>Avg MQs</i> |
|------|---------------------|----------------|-------------------|----------------|
| SMC | 92 | 111 | 122 | 286063 |
| AAMC | 444 | 7 | 30 | 3701916 |

 Property-directed verification (PDV)



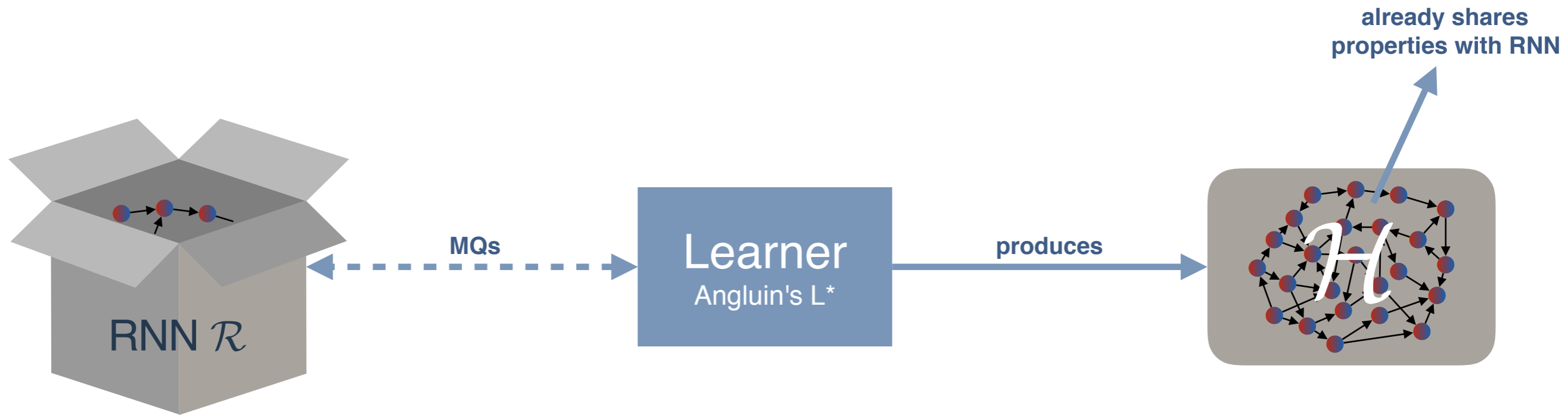
inspired by black-box checking: **interweave learning and model checking**

[Peled, Vardi, Yannakakis 1999]



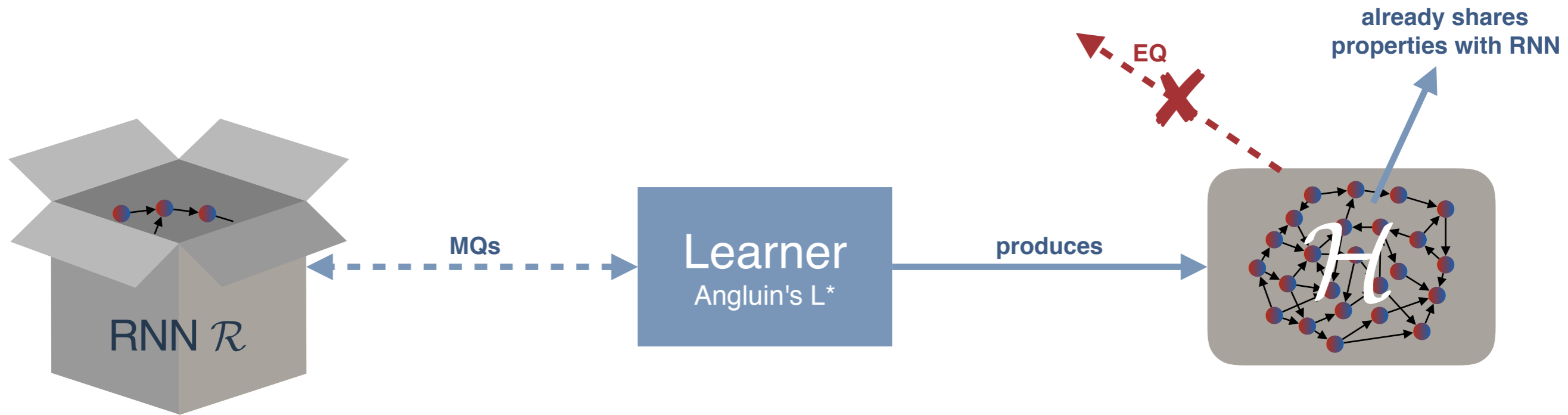
inspired by black-box checking: **interweave learning and model checking**

[Peled, Vardi, Yannakakis 1999]



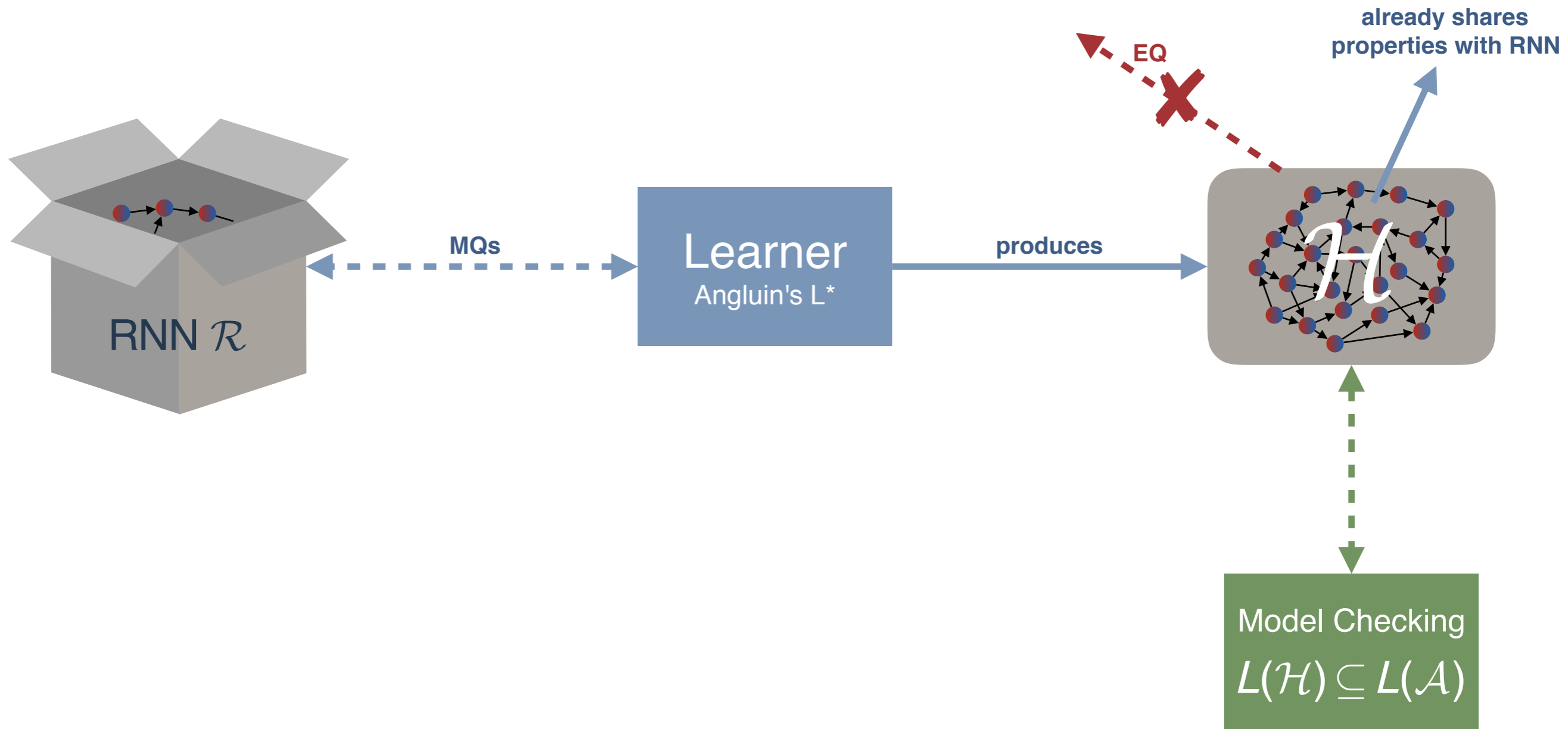
inspired by black-box checking: **interweave learning and model checking**

[Peled, Vardi, Yannakakis 1999]



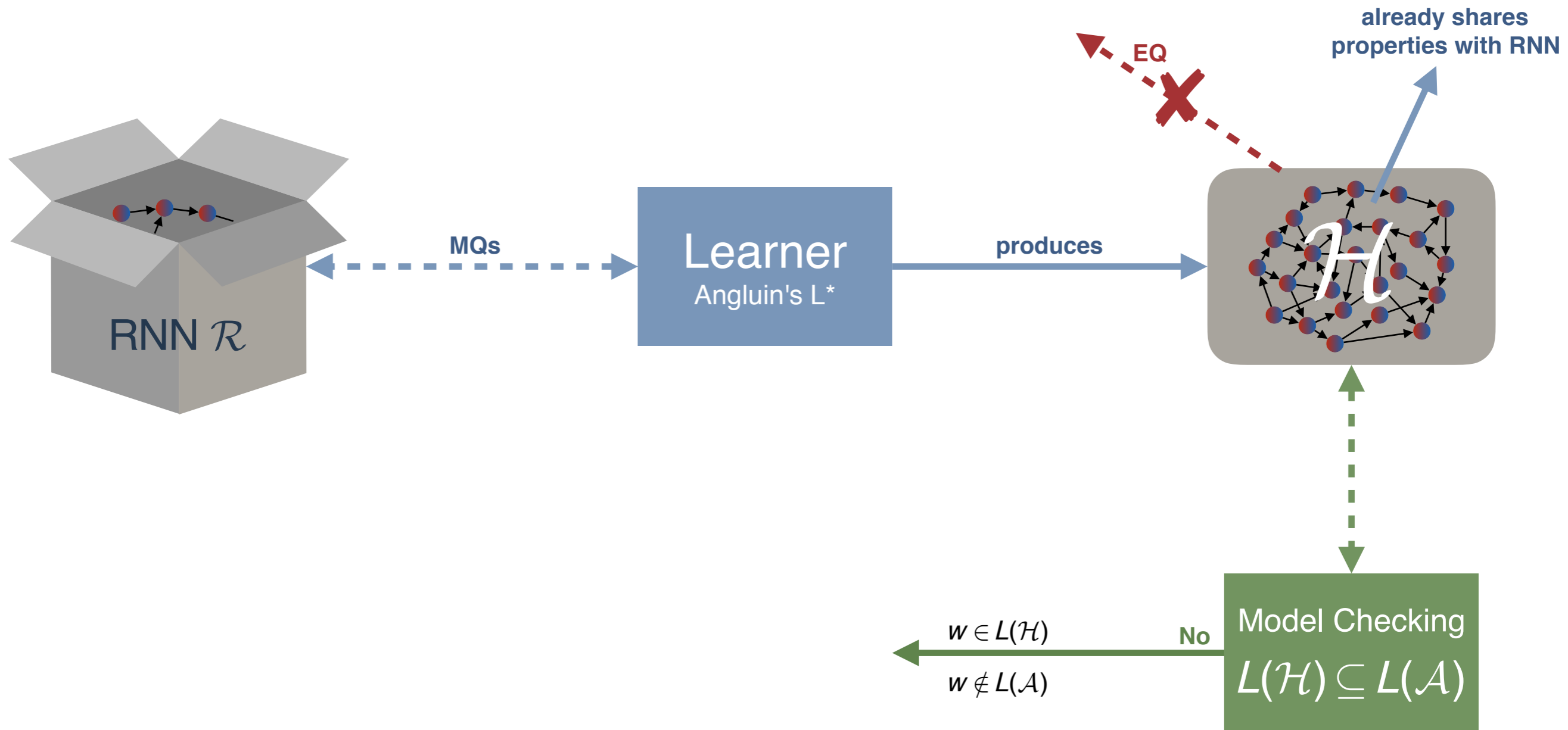
inspired by black-box checking: **interweave learning and model checking**

[Peled, Vardi, Yannakakis 1999]



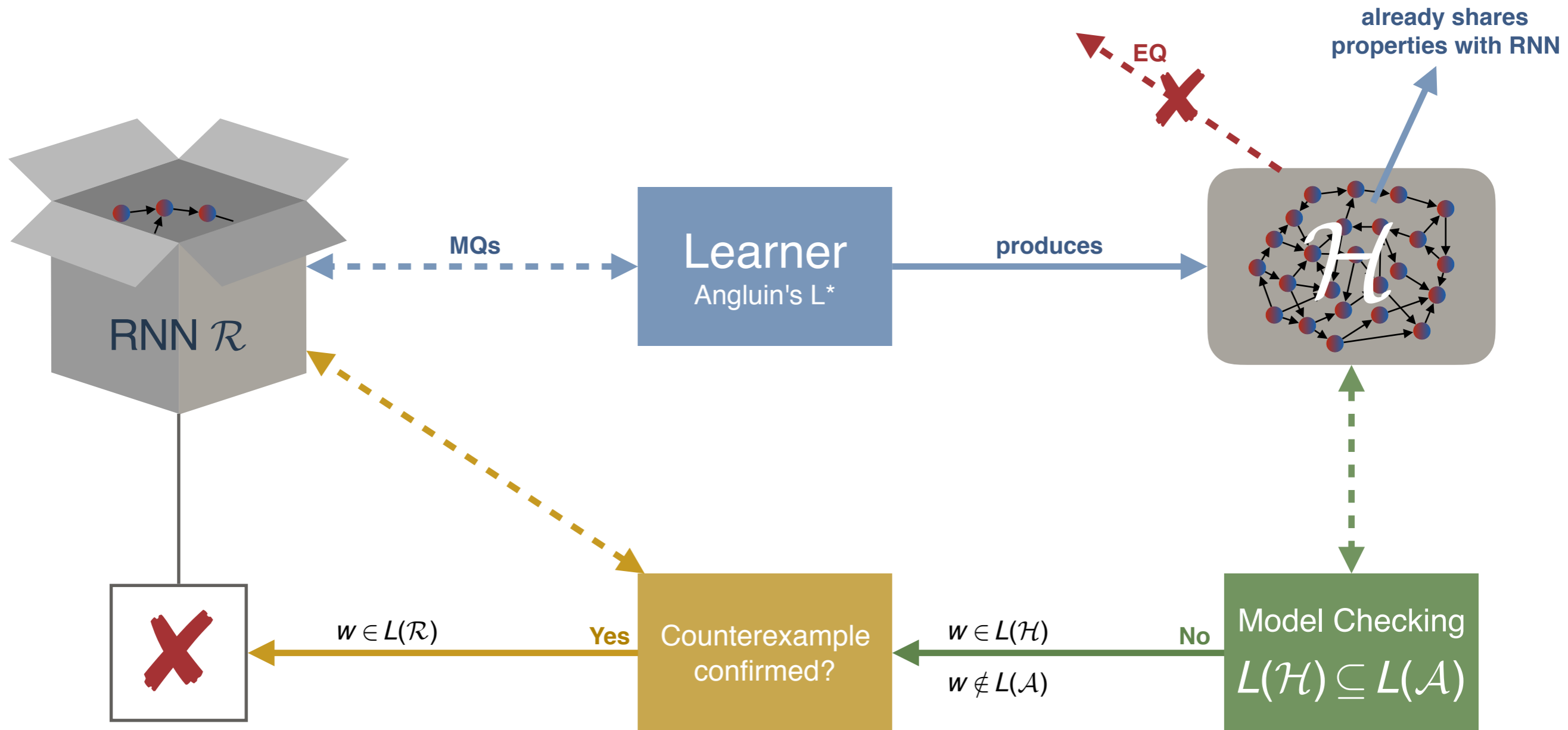
inspired by black-box checking: **interweave learning and model checking**

[Peled, Vardi, Yannakakis 1999]



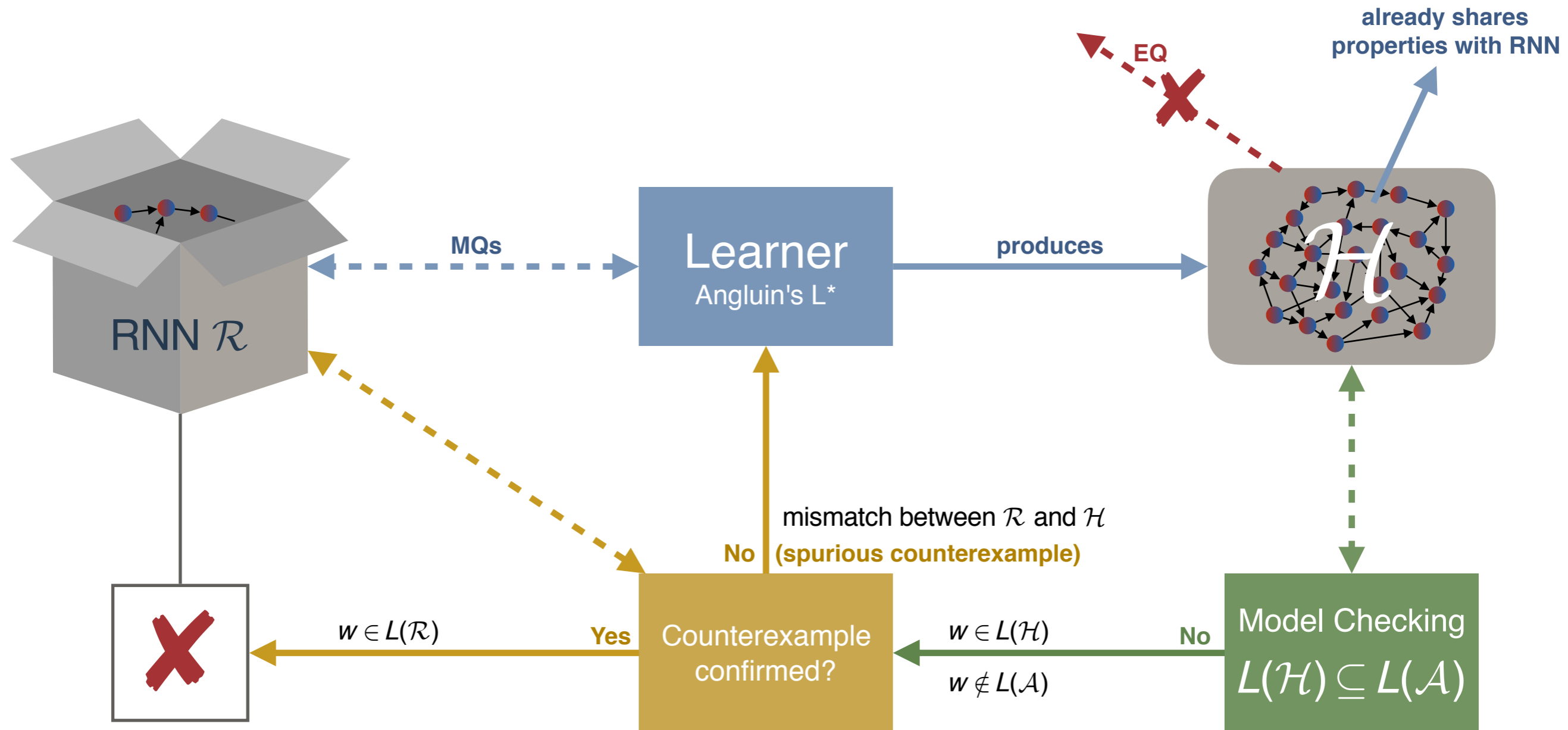
inspired by black-box checking: **interweave learning and model checking**

[Peled, Vardi, Yannakakis 1999]



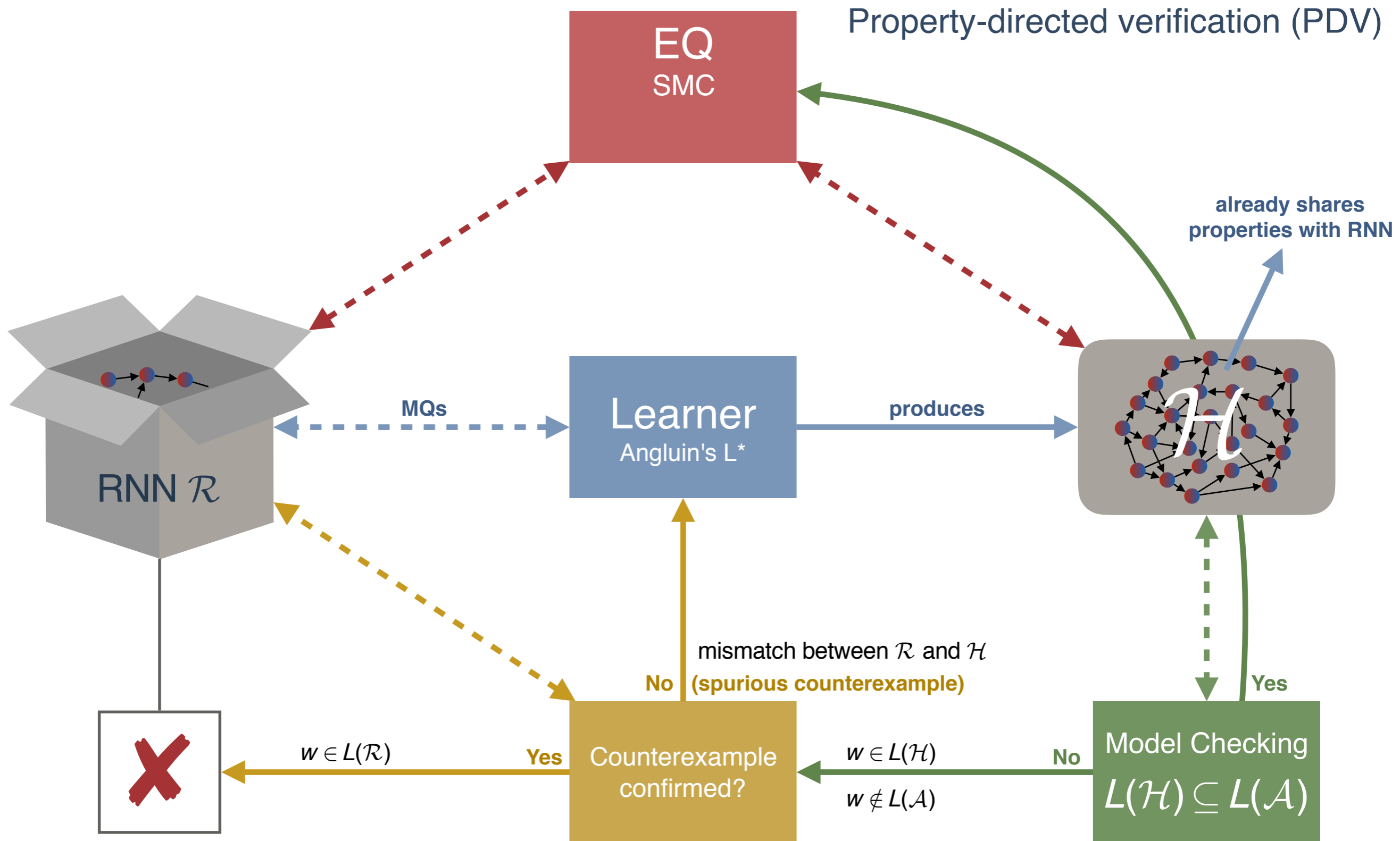
inspired by black-box checking: **interweave learning and model checking**

[Peled, Vardi, Yannakakis 1999]



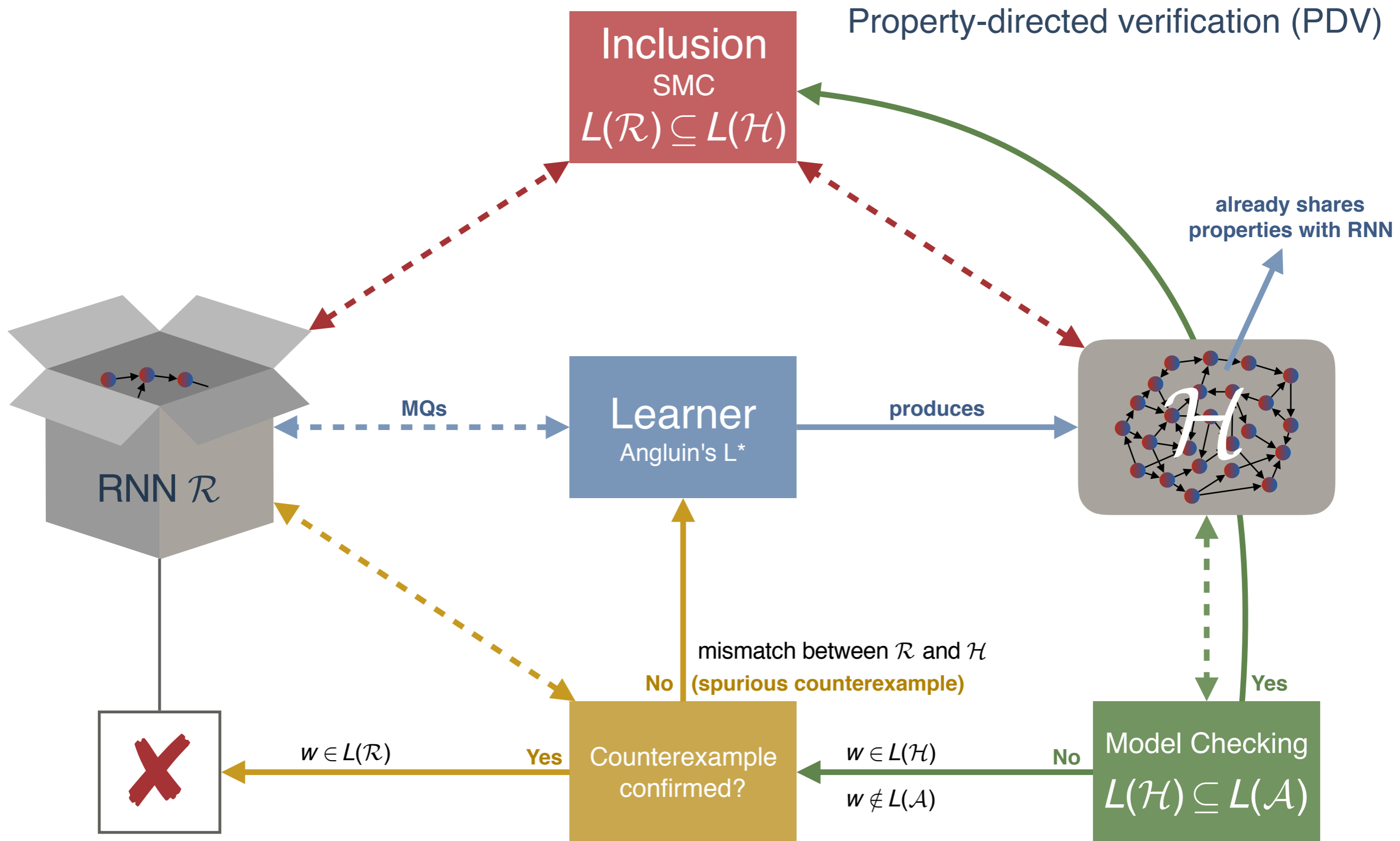
inspired by black-box checking: **interweave learning and model checking**

[Peled, Vardi, Yannakakis 1999]



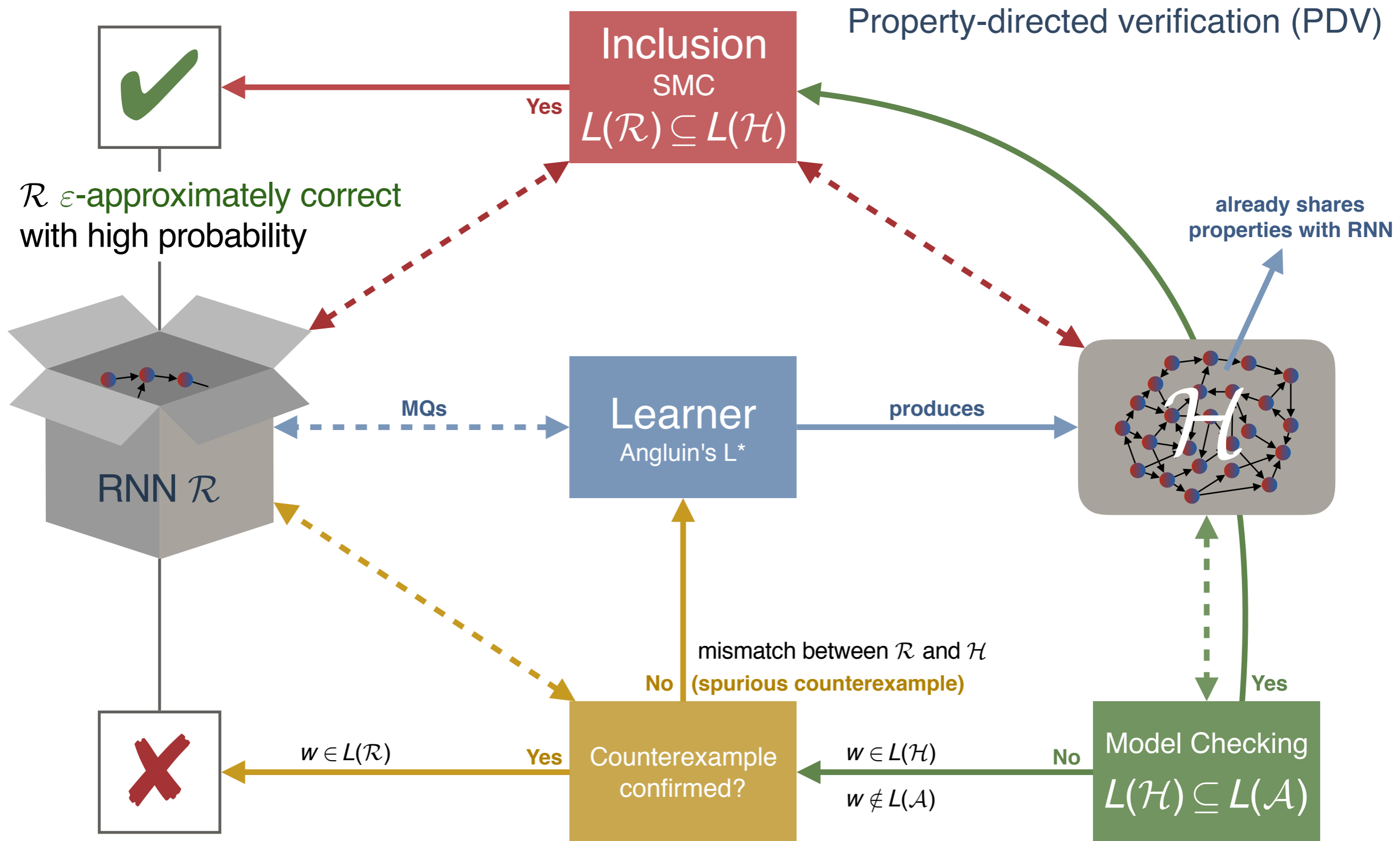
inspired by black-box checking: **interweave learning and model checking**

[Peled, Vardi, Yannakakis 1999]



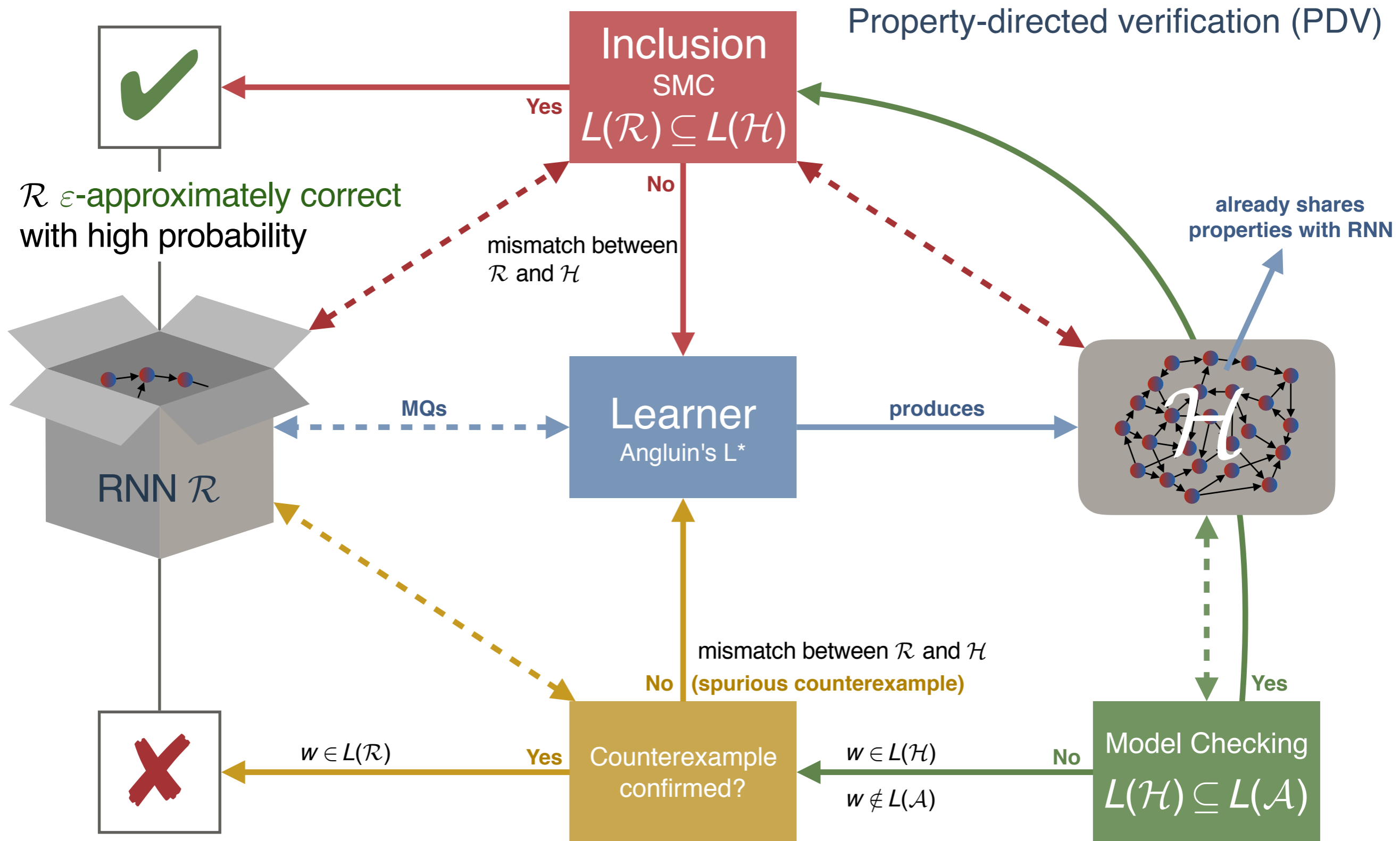
inspired by black-box checking: **interweave learning and model checking**

[Peled, Vardi, Yannakakis 1999]



inspired by black-box checking: **interweave learning and model checking**

[Peled, Vardi, Yannakakis 1999]



inspired by black-box checking: **interweave learning and model checking**

[Peled, Vardi, Yannakakis 1999]

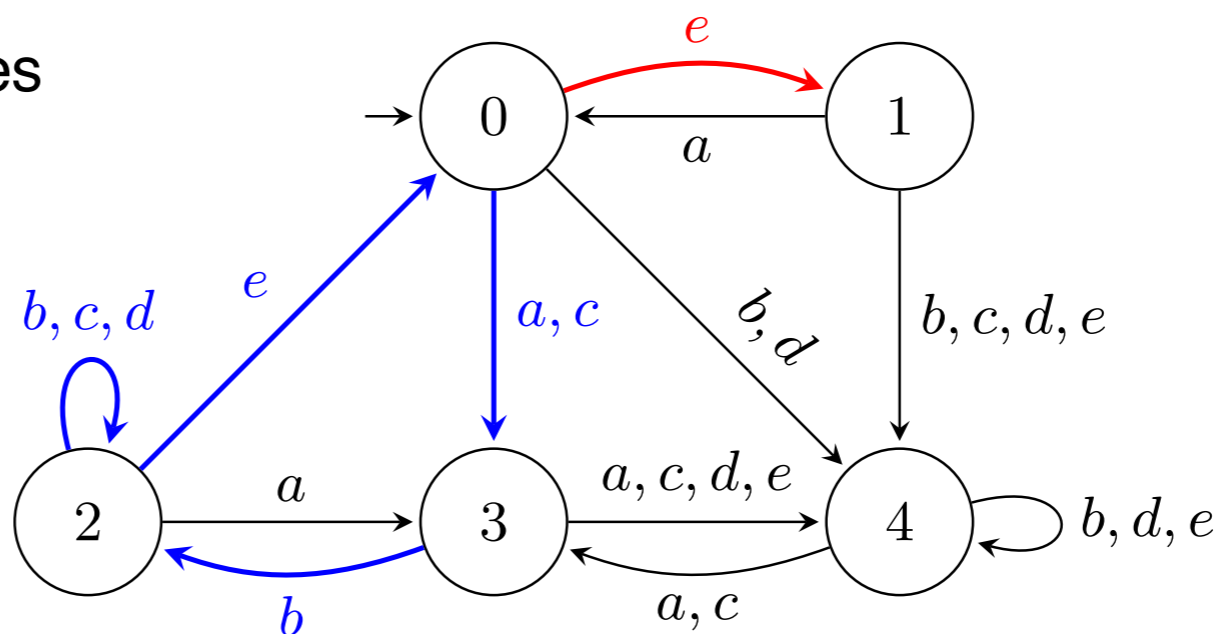
Checked on 30 DFAs / RNNs and 138 specifications:

| Type | <i>Avg time</i> (s) | <i>Avg len</i> | <i># Mistakes</i> | <i>Avg MQs</i> |
|------|---------------------|----------------|-------------------|----------------|
| SMC | 92 | 111 | 122 | 286063 |
| AAMC | 444 | 7 | 30 | 3701916 |
| PDV | 21 | 11 | 109 | 28318 |

Checked on 30 DFAs / RNNs and 138 specifications:

| Type | <i>Avg time</i> (s) | <i>Avg len</i> | <i># Mistakes</i> | <i>Avg MQs</i> |
|------|---------------------|----------------|-------------------|----------------|
| SMC | 92 | 111 | 122 | 286063 |
| AAMC | 444 | 7 | 30 | 3701916 |
| PDV | 21 | 11 | 109 | 28318 |

PDV detected "faulty flows"
in 81/ 109 of the counterexamples



Conclusion

Conclusion

1. Extraction of finite-state information from RNNs contributes to their understanding and verification.

Conclusion

1. Extraction of finite-state information from RNNs contributes to their understanding and verification.
2. Property-directed verification (PDV) is able to find short counterexamples fast.

Conclusion

1. Extraction of finite-state information from RNNs contributes to their understanding and verification.
2. Property-directed verification (PDV) is able to find short counterexamples fast.

Conclusion and Outlook

1. Practical applications of PDV?

Conclusion

1. Extraction of finite-state information from RNNs contributes to their understanding and verification.
2. Property-directed verification (PDV) is able to find short counterexamples fast.

Conclusion and Outlook

1. Practical applications of PDV?
2. Verification of RNN-based agent environment systems?

Conclusion

1. Extraction of finite-state information from RNNs contributes to their understanding and verification.
2. Property-directed verification (PDV) is able to find short counterexamples fast.

Conclusion and Outlook

1. Practical applications of PDV?
2. Verification of RNN-based agent environment systems?
3. Going beyond regular specifications: Inference of context-free languages.

Conclusion

1. Extraction of finite-state information from RNNs contributes to their understanding and verification.
2. Property-directed verification (PDV) is able to find short counterexamples fast.

Conclusion and Outlook

1. Practical applications of PDV?
2. Verification of RNN-based agent environment systems?
3. Going beyond regular specifications: Inference of context-free languages.

Thank you!