

# Responsibility-Sensitive Safety for Automated Driving

Clovis Eberhart

Software Construction Laboratory, Research Institute of Electric Communications,  
Tohoku University

“AI-Physical Systems” ASPIRE Kick-off Meeting, March 3-4,  
2026

# Challenges in Automated Driving

## Pros of automated driving

- can reduce travel time
- can reduce fuel consumption
- may become safer than human driving

## Cons of automated driving

- little public trust
- no clear regulations
- no clear delimitation of liability
- **safety problem** (twice as many crashes per mile as human driving)

# Safety of Automated Driving

Two types of methods for safety of automated driving.

## Empirical methods

- statistical methods, simulations, testing. . .
- **great scalability**
- **works on black-box systems**
- **no formal guarantees**

## Formal methods

- logical methods, automata-based methods, **semantic methods**
- **not scalable**
- **need formal definitions**
- **formal guarantee**

# Responsibility-Sensitive Safety<sup>1</sup>

## Problem with logical methods

Need formal definitions of systems, which are too complex.

Abstract away part of the systems to keep them simple and prove something about them.

$$\begin{cases} \dot{x} = v \\ \dot{v} = a \end{cases} \quad a \in [-b_{\max}, a_{\max}]$$

## Responsibility-Sensitive Safety

- RSS condition: condition under which the system is safe
- RSS response: how to maintain safety
- based on responsibility

---

<sup>1</sup>Shalev-Shwartz, Shammah, and Shashua, *On a formal model of safe and scalable self-driving cars*, arXiv, 2017.

# Responsibility-Sensitive Safety

Driving situation:



Who is responsible in case of a crash?

Strategy (RSS response)

Brake enough to leave at least a distance of

$$\max\left(0, v_r \rho + \frac{a_{\max} \rho^2}{2} + \frac{(v_r + a_{\max})^2}{2b_{\min}} - \frac{v_f^2}{2b_{\max}}\right)$$

between the two vehicles.

# Goal of the research<sup>1</sup>

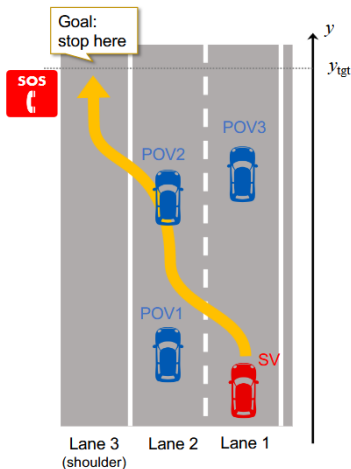
## Limitation of RSS

Pen-and-paper proofs:  
prone to errors, do not  
scale.

## Goal

Given a “strategy pattern”,  
get

- an RSS condition (a safety zone)
- an RSS response (a strategy to stay safe)



<sup>1</sup>Hasuo, Eberhart, Haydon, Dubut, Bohrer, Kobayashi, Pruekprasert, Zhang, Pallas, Yamada, Suenaga, Ishikawa, Kamijo, Shinya, and Suetomi. *Goal-Aware RSS for Complex Scenarios via Program Logic*, IEEE T-IV, 2023

## Ingredients

- **hybrid programs**: programs with continuous behaviours
- **Hoare quadruples**: properties of hybrid programs
- **dFHL**: logic to prove Hoare quadruples

## Proving safety of driving situations

- Turn a “strategy pattern” into a hybrid program
- Turn strategy safety into a Hoare quadruple
- Prove the quadruple with dFHL
- Infer pattern to get a concrete strategy

# Scenarios as hybrid programs

Hybrid programs:

$$\alpha, \beta ::= \text{skip} \mid \alpha; \beta \mid x := e \mid \text{if } (A) \alpha \text{ else } \beta \mid \\ \text{while } (A) \alpha \mid \text{dwhile } (A) \{ \dot{\mathbf{x}} = \mathbf{f} \}.$$

## Example of strategy pattern

$$\alpha = l := 0; \text{dwhile } (P_1) \{ \dot{x} = v, \dot{v} = a_{\max} \}; \\ l := 0.5; \text{dwhile } (P_2) \{ \dot{x} = v, \dot{v} = 0 \}; \\ l := 1; \dots$$

## Hoare quadruple

A **Hoare quadruple**  $\{P\} \alpha \{Q\} : S$  is valid if for all stores  $\rho \models P$ :

- there exists  $\rho'$  s.t.  $\langle \alpha, \rho \rangle$  converges to  $\rho' \models Q$ ,
- for all  $\langle \alpha, \rho \rangle \rightarrow^* \langle \alpha', \rho' \rangle$ ,  $\rho' \models S$ .

Meaning of Hoare quadruple:

- if  $P$  holds before  $\alpha$  runs
- then  $Q$  holds after the run
- and  $S$  holds along the whole execution

For us:

- $P$ : RSS condition (to be inferred in our case)
- $Q$ : goal condition (progress)
- $S$ : safety (absence of collision)

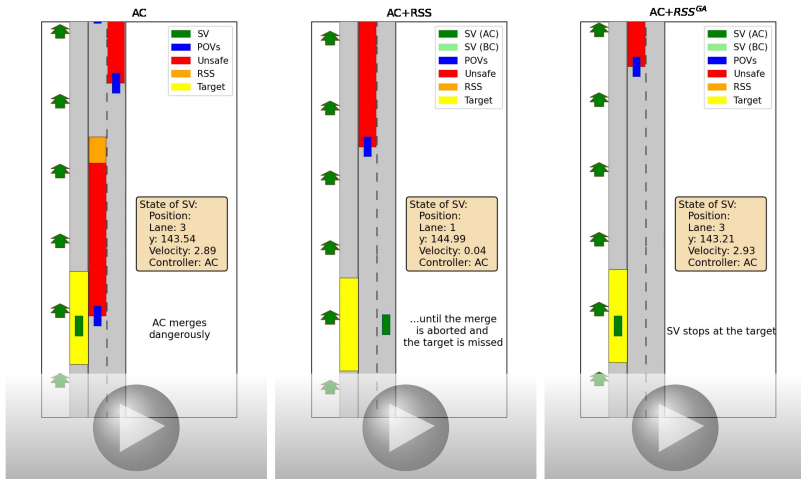
# The logic dFHL

$$\begin{array}{c}
 \frac{}{\{A\} \text{ skip } \{A\} : A} \text{ (SKIP)} \quad \frac{\{A\} \alpha \{B\} : S \quad \{B\} \beta \{C\} : S}{\{A\} \alpha; \beta \{C\} : S} \text{ (SEQ)} \quad \frac{}{\{A[e/x]\} x := e \{A\} : A \vee A[e/x]} \text{ (ASSIGN)} \\
 \\
 \frac{\{A \wedge B\} \alpha \{C\} : S \quad \{\neg A \wedge B\} \beta \{C\} : S}{\{B\} \text{ if } (A) \alpha \text{ else } \beta \{C\} : S} \text{ (IF)} \quad \frac{\{A \wedge B \wedge e_{\text{var}} \gtrsim 0 \wedge e_{\text{var}} = x\} \alpha \{B \wedge e_{\text{var}} \gtrsim 0 \wedge e_{\text{var}} \leq x - 1\} : S}{\{B \wedge e_{\text{var}} \gtrsim 0\} \text{ while } (A) \alpha \{\neg A \wedge B \wedge e_{\text{var}} \gtrsim 0\} : S} \text{ (WH)}^\dagger \\
 \\
 \begin{array}{l}
 \text{inv: } A \Rightarrow e_{\text{inv}} \sim 0 \quad e_{\text{var}} \geq 0 \wedge e_{\text{inv}} \sim 0 \Rightarrow \mathcal{L}_{\dot{x}=\mathbf{f}} e_{\text{inv}} \simeq 0 \\
 \text{var: } A \Rightarrow e_{\text{var}} \geq 0 \quad e_{\text{var}} \geq 0 \wedge e_{\text{inv}} \sim 0 \Rightarrow \mathcal{L}_{\dot{x}=\mathbf{f}} e_{\text{var}} \leq e_{\text{ter}} \\
 \text{ter: } A \Rightarrow e_{\text{ter}} < 0 \quad e_{\text{var}} \geq 0 \wedge e_{\text{inv}} \sim 0 \Rightarrow \mathcal{L}_{\dot{x}=\mathbf{f}} e_{\text{ter}} \leq 0
 \end{array} \\
 \frac{}{\{A\} \text{ dwhile } (e_{\text{var}} > 0) \dot{x} = \mathbf{f} \{e_{\text{var}} = 0 \wedge e_{\text{inv}} \sim 0\} : e_{\text{inv}} \sim 0 \wedge e_{\text{var}} \geq 0} \text{ (DWH)}^\dagger \quad \frac{\{A'\} \alpha \{B'\} : S' \quad A \Rightarrow A' \quad S' \wedge B' \Rightarrow B \quad S' \Rightarrow S}{\{A\} \alpha \{B\} : S} \text{ (LIMP)} \\
 \\
 \frac{\{A\} \alpha \{B\} : S \quad \{A\} \alpha \{B'\} : S'}{\{A\} \alpha \{B \wedge B'\} : S \wedge S'} \text{ (CONJ)} \quad \frac{A_0 \Rightarrow (\exists t \geq 0. C_{<t} \wedge \neg C_t \wedge B_t \wedge S_{\leq t})}{\{A\} \text{ dwhile } (C) \dot{x} = \mathbf{f} \{B\} : S} \text{ (DWH-SOL)}
 \end{array}$$

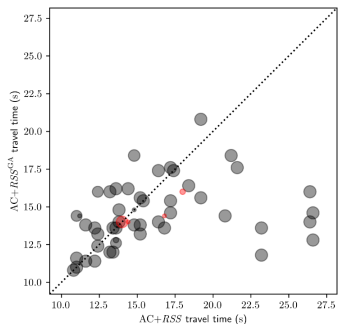
## Theorem

*Only valid Hoare quadruples can be proved in dHL.*

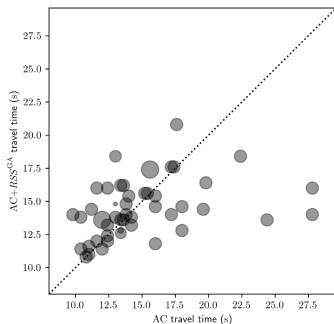
# Experimental results



# Experimental evaluation



(a)  $AC+RSS^{CA}$  vs  $AC+RSS^{GA}$ .

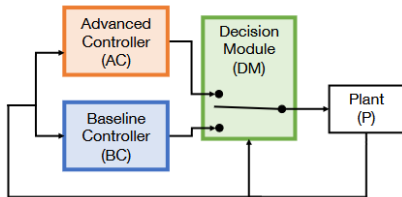


(b)  $AC$  vs  $AC+RSS^{GA}$ .

Fig. 16: Comparison for progress. A red disc indicates that  $AC+RSS^{CA}$  failed to achieve the goal in this scenario instance.

# Safety architectures<sup>1</sup>

Proving safety of “safety architectures”.



Compatible with RSS:

- DM: RSS precondition
- BC: RSS response
- AC: unknown (e.g., AI controller)

---

<sup>1</sup>Eberhart, Dubut, Haydon, and Hasuo. *Formal Verification of Safety Architectures for Automated Driving*, IV 2023

## Hoare quintuples

$A : \{P\} \alpha \{Q\} : S$  where  $A$  is an assumption.

We can then prove:

- system safe and makes progress under some strong assumption (e.g., other vehicles at constant speed)
- system safe under a weaker assumption (e.g., other vehicles may change speed)

where  $AC$  is an AI controller.

Our approach to RSS:

- is based on semantics of programming languages
- gives **formal** guarantees of safety for automated driving systems
- can **scale** beyond pen-and-paper proofs (needs a proof engineer)
- can be used to prove **safety of AI controllers** (in an architecture)

$$\begin{array}{c}
 \frac{}{\langle \text{skip}; \beta, \rho \rangle \rightarrow \langle \beta, \rho \rangle} \quad \frac{\langle \alpha, \rho \rangle \rightarrow \langle \alpha', \rho' \rangle}{\langle \alpha; \beta, \rho \rangle \rightarrow \langle \alpha'; \beta, \rho' \rangle} \quad \frac{}{\langle x := e, \rho \rangle \rightarrow \langle \text{skip}, \rho[x \rightarrow \llbracket e \rrbracket_\rho] \rangle} \quad \frac{\rho \models A}{\langle \text{if } (A) \alpha \text{ else } \beta, \rho \rangle \rightarrow \langle \alpha, \rho \rangle} \\
 \\
 \frac{\rho \not\models A}{\langle \text{if } (A) \alpha \text{ else } \beta, \rho \rangle \rightarrow \langle \beta, \rho \rangle} \quad \frac{\rho \not\models A}{\langle \text{while } (A) \alpha, \rho \rangle \rightarrow \langle \text{skip}, \rho \rangle} \quad \frac{\rho \models A}{\langle \text{while } (A) \alpha, \rho \rangle \rightarrow \langle \alpha; \text{while } (A) \alpha, \rho \rangle} \\
 \\
 \frac{t \geq 0 \quad \hat{x}(0) = \rho \quad \frac{d\hat{x}}{dt}(t) = \llbracket \mathbf{f} \rrbracket_{\hat{x}(t)} \quad \rho' = \hat{x}(t) \quad \forall t' \leq t. \hat{x}(t') \models A}{\langle \text{dwhile } (A) \{ \dot{x} = \mathbf{f} \}, \rho \rangle \rightarrow \langle \text{dwhile } (A) \{ \dot{x} = \mathbf{f} \}, \rho' \rangle} (*) \\
 \\
 \frac{t \geq 0 \quad \hat{x}(0) = \rho \quad \frac{d\hat{x}}{dt}(t) = \llbracket \mathbf{f} \rrbracket_{\hat{x}(t)} \quad \rho' = \hat{x}(t) \quad \forall t' < t. \hat{x}(t') \models A \quad \hat{x}(t) \not\models A}{\langle \text{dwhile } (A) \{ \dot{x} = \mathbf{f} \}, \rho \rangle \rightarrow \langle \text{skip}, \rho' \rangle} (*)
 \end{array}$$

## Lemma

*Reduction is deterministic: if  $\langle \alpha, \rho \rangle \rightarrow^* \langle \alpha_1, \rho_1 \rangle$  and  $\langle \alpha, \rho \rangle \rightarrow^* \langle \alpha_2, \rho_2 \rangle$ , then  $\langle \alpha_1, \rho_1 \rangle \rightarrow^* \langle \alpha_2, \rho_2 \rangle$  or  $\langle \alpha_2, \rho_2 \rangle \rightarrow^* \langle \alpha_1, \rho_1 \rangle$ .*

# The dwhile rule

$$\begin{array}{lll} \text{inv:} & A \Rightarrow e_{\text{inv}} \sim 0 & e_{\text{var}} \geq 0 \wedge e_{\text{inv}} \sim 0 \Rightarrow \mathcal{L}_{\dot{\mathbf{x}}=\mathbf{f}} e_{\text{inv}} \simeq 0 \\ \text{var:} & A \Rightarrow e_{\text{var}} \geq 0 & e_{\text{var}} \geq 0 \wedge e_{\text{inv}} \sim 0 \Rightarrow \mathcal{L}_{\dot{\mathbf{x}}=\mathbf{f}} e_{\text{var}} \leq e_{\text{ter}} \\ \text{ter:} & A \Rightarrow e_{\text{ter}} < 0 & e_{\text{var}} \geq 0 \wedge e_{\text{inv}} \sim 0 \Rightarrow \mathcal{L}_{\dot{\mathbf{x}}=\mathbf{f}} e_{\text{ter}} \leq 0 \end{array}$$

$$\frac{\{A\} \text{ dwhile } (e_{\text{var}} > 0) \dot{\mathbf{x}} = \mathbf{f} \{e_{\text{var}} = 0 \wedge e_{\text{inv}} \sim 0\} : e_{\text{inv}} \sim 0 \wedge e_{\text{var}} \geq 0}{\text{(DWH)}^\dagger}$$

- makes use of: an invariant (for the safety condition), a variant (to leave condition), and a terminator (to prove termination)
- $\mathcal{L}_{\dot{\mathbf{x}}=\mathbf{f}} e$ : Lie derivative, how much  $e$  varies by following  $\dot{\mathbf{x}} = \mathbf{f}$

# Advantages of this approach

- diversity of driving situations can be treated by different rules
- explainable safety  $\leadsto$  better public trust
- clearly defined rules  $\leadsto$  can be investigated in case of crash
- clearly defined rules  $\leadsto$  car manufacturers can avoid liability