

- Hasuo+, IEEE Trans. Intell. Vehicles, early access
<https://doi.org/10.1109/TIV.2022.3169762>
- <https://arxiv.org/abs/2207.02387>

S O K E N D A I

NII



Goal-Aware RSS for Complex Scenarios via Program Logic

Ichiro Hasuo^{1,7,*}, Clovis Eberhart^{1,8,*}, James Haydon^{1,*}, Jérémy Dubut^{1,8}, Rose Bohrer²,
Tsutomu Kobayashi¹, Sasinee Pruekprasert¹, Xiao-Yi Zhang¹, Erik André Pallas³,
Akihisa Yamada^{4,1}, Kohei Suenaga^{5,1}, Fuyuki Ishikawa¹,
Kenji Kamijo⁶, Yoshiyuki Shinya⁶, and Takamasa Suetomi⁶

1: National Institute of Informatics, Tokyo, Japan

2: Worcester Polytechnique Institute, USA (work done at NII)

3: University of Augsburg, Germany (work done at NII)

4: AIST, Japan

5: Kyoto University, Japan

6: Mazda Motor Corporation, Japan

7: SOKENDAI (The Graduate University for Advanced Studies), Japan

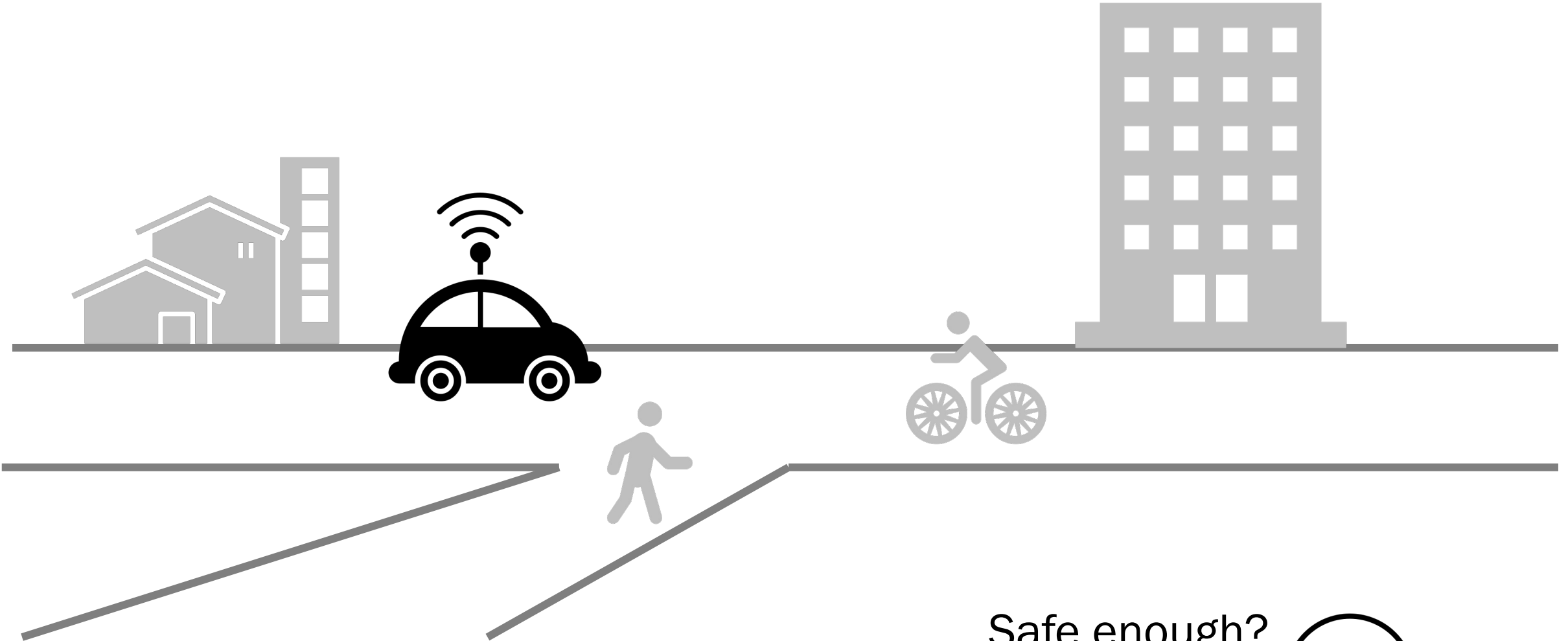
8: Japanese-French Laboratory for Informatics (IRL 3527), Tokyo, Japan

Supported by ERATO HASUO Metamathematics for Systems Design Project (No. JPMJER1603) and ACT-I (No. JPMJPR17UA), JST;
and Grants-in-aid No. 19K20215 & 19K20249, JSPS.

*: equal contribution

Outline

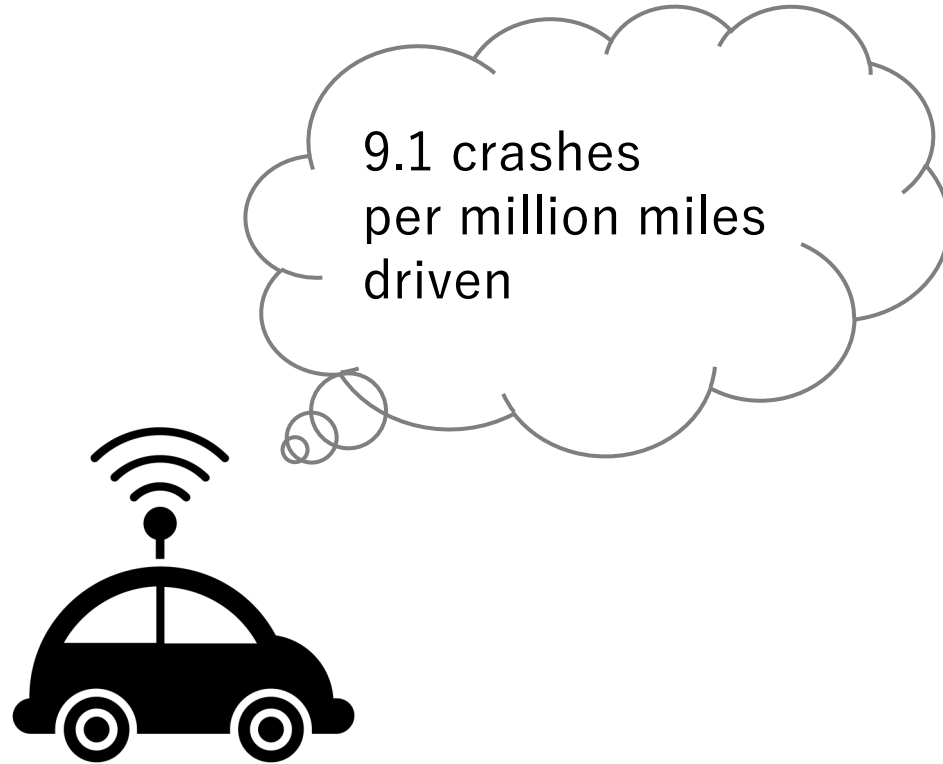
- A non-technical overview
- The modeling problem
- The RSS answer to the modeling problem
- Technical contributions: the logic dFHL
- Perspectives, practical & theoretical



Safe enough?



Guarantee by statistical data



Guarantee by testing and simulation



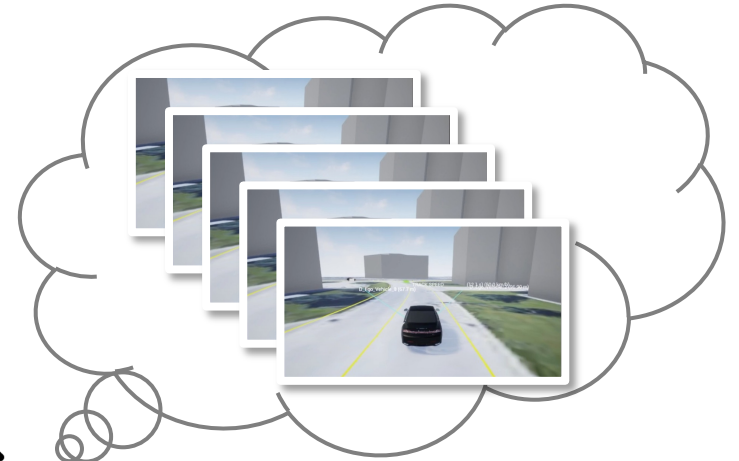
Guarantee strong enough?

Guarantee
by statistical data

9.1 crashes
per million miles
driven



Guarantee
by testing and simulation



Explainability?

Proof.

We prove the first statement. The rest is shown symmetrically.

Let $S \subseteq L$ be an arbitrary subset. We let S^\downarrow be the downward closure of S , that is,

$$S^\downarrow := \{y \in L \mid y \sqsubseteq s \text{ for each } s \in S\}$$

Since $S^\downarrow \subseteq L$ is a subset of L , it has its supremum in the semilattice (L, \sqsubseteq) . We claim that $\bigsqcup S^\downarrow$ is the infimum of S .

To prove the claim, it suffices to show the two-way characterization in (2.1), that is, we need to show

$$\frac{y \sqsubseteq s \text{ for each } s \in S}{y \sqsubseteq \bigsqcup S^\downarrow}.$$

For the downward implication in ??,

$$\begin{aligned} & y \sqsubseteq s \text{ for each } s \in S \\ \implies & y \in S^\downarrow && \text{by def. of } S^\downarrow \\ \implies & y \sqsubseteq \bigsqcup S^\downarrow && \text{since } \bigsqcup S^\downarrow \text{ is an upper bound} \end{aligned}$$

For the upward implication in ??, we first observe

$$\bigsqcup S^\downarrow \sqsubseteq s \text{ for each } s \in S.$$



Responsibility-Sensitive Safety (RSS)

[Shalev-Shwartz et al., arXiv preprint, 2017]

Lemma. (Conditional safety)

If all cars comply with
RSS rules, then
there is no collision

mathematically
proved

+

Assumption. (Rule compliance)

All cars comply with
RSS rules

the manufacturer's
responsibility to
ensure this



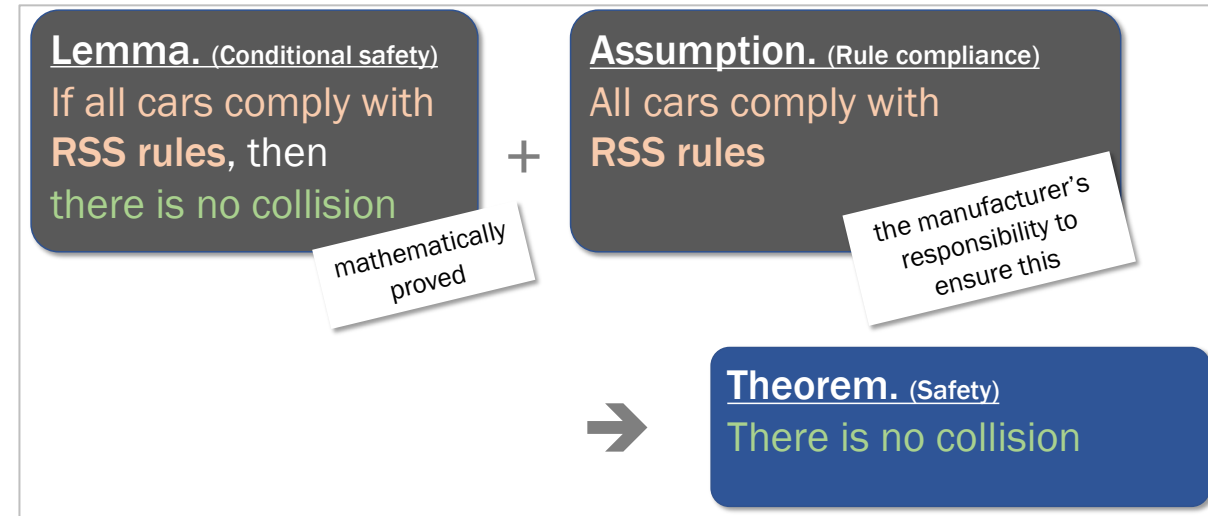
Theorem. (Safety)

There is no collision

Responsibility-Sensitive Safety (RSS)

[Shalev-Shwartz et al., arXiv preprint, 2017]

- “Let’s put all dirty details in an assumption” ...
Isn’t this cheating?
Isn’t the assumption too big?
- → No!
 - RSS rules are rigorous,
their compliance is verifiable by the third party
 - RSS rules can be enforced by the safety architecture (later)
 - Overall, RSS rules have the right granularity to impose as **social contracts**
- (Fresh view on proofs for us logicians...)



RSS Rule, an Example

[Shalev-Shwartz et al., arXiv preprint, 2017]

- An RSS rule is a pair (A, α) of an *RSS condition* A and a *proper response* α



RSS condition A :

Maintain an inter-vehicle distance at least

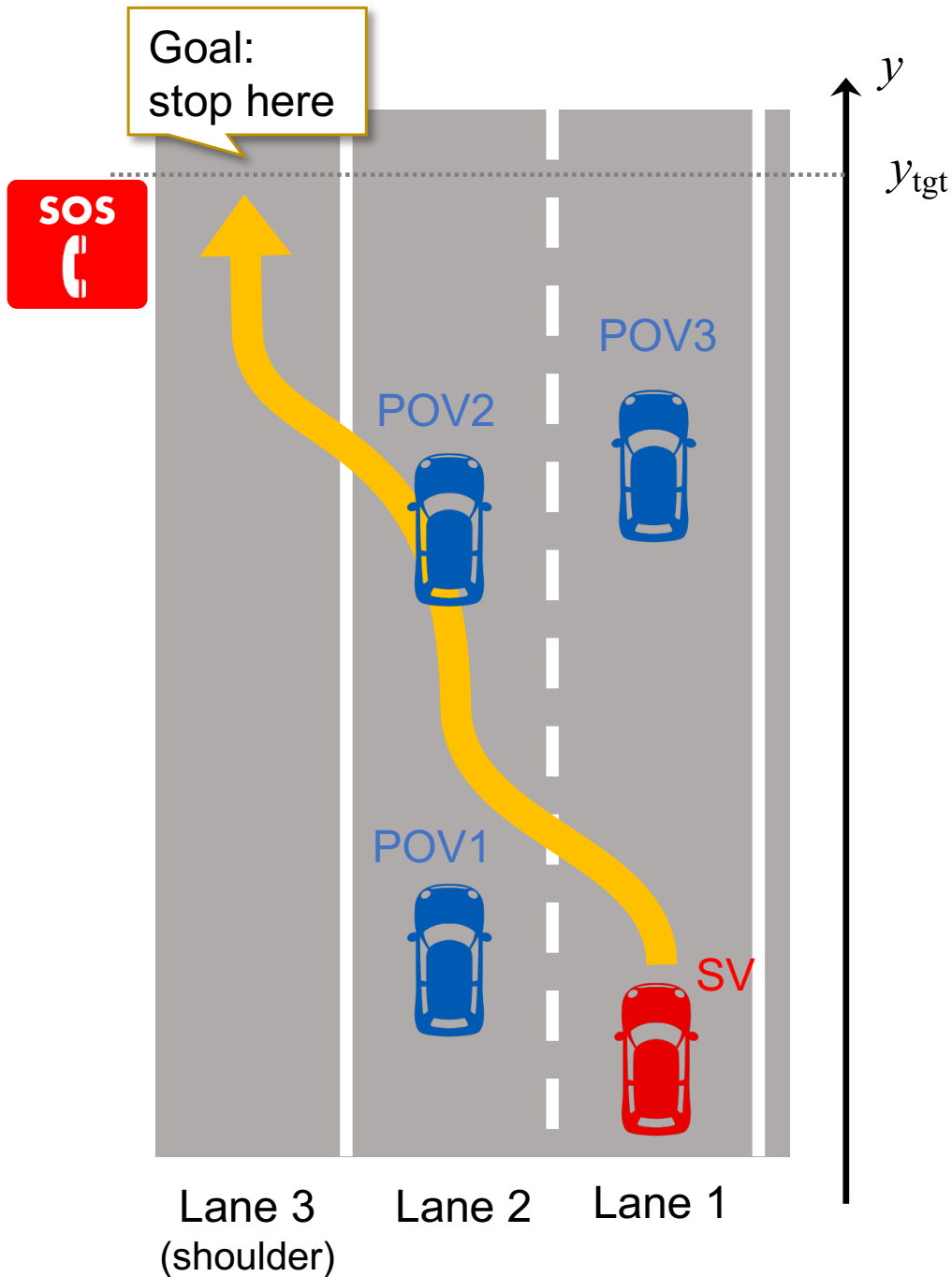
$$d_{\min} = \left[v_r \rho + \frac{1}{2} a_{\max, \text{accel}} \rho^2 + \frac{(v_r + \rho a_{\max, \text{accel}})^2}{2a_{\min, \text{brake}}} - \frac{v_f^2}{2a_{\max, \text{brake}}} \right]_+$$

Proper response α :

If A is about to be violated, brake at rate $a_{\min, \text{brake}}$ within ρ seconds

Conditional safety lemma:

Any execution of α , from a state that satisfies A , is collision-free.



- Now what about this pull over scenario?
- Essential for eyes-off ADVs to hand the control over to human drivers
- Requires complex decision making
 - Merge before POV1, or after?
 - Accelerate to pass POV1...
 - ➔ Risk of overrun?



Our Contribution: Logical Formalization of RSS → More Scenarios

RSS

Responsibility-Sensitive Safety, Shalev-Shwartz et al., 2017

- Basic methodology of logical safety rules
- Standardization (IEEE 2846)
- Lack of formal implementation → appl. to complex scenarios is hard
- Guarantees only collision-freedom so far

↓ Software science research

differential program logic dFHL (our contribution)

$$\begin{array}{l} \text{inv: } A \Rightarrow e_{\text{inv}} \sim 0 \quad e_{\text{var}} \geq 0 \wedge e_{\text{inv}} \sim 0 \Rightarrow \mathcal{L}\dot{x} = f \ e_{\text{inv}} \geq 0 \\ \text{var: } A \Rightarrow e_{\text{var}} \geq 0 \quad e_{\text{var}} \geq 0 \wedge e_{\text{inv}} \sim 0 \Rightarrow \mathcal{L}\dot{x} = f \ e_{\text{var}} \leq e_{\text{ter}} \\ \text{ter: } A \Rightarrow e_{\text{ter}} < 0 \quad e_{\text{var}} \geq 0 \wedge e_{\text{inv}} \sim 0 \Rightarrow \mathcal{L}\dot{x} = f \ e_{\text{ter}} \leq 0 \end{array} \quad (\text{DWH})$$

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

- A logical system for deriving and proving safety rules

Compositional rule derivation workflow by dFHL (our contribution)



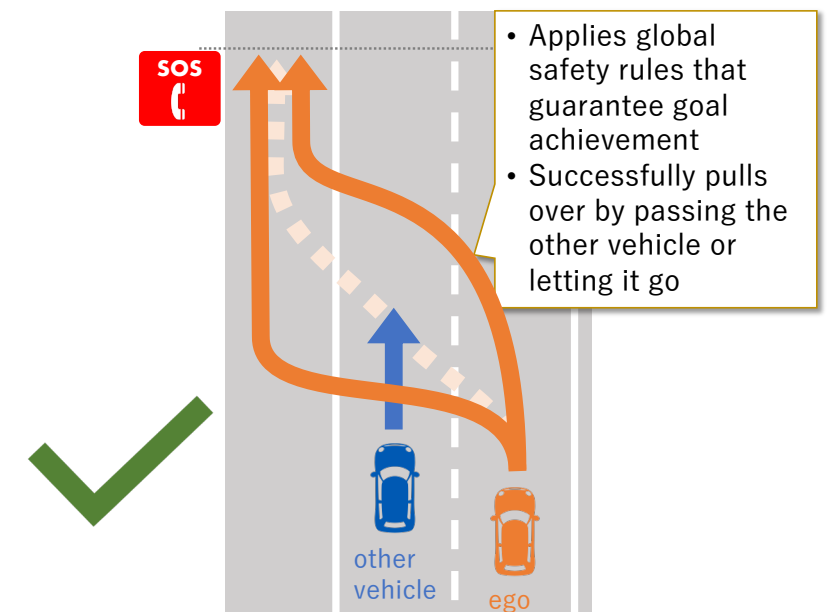
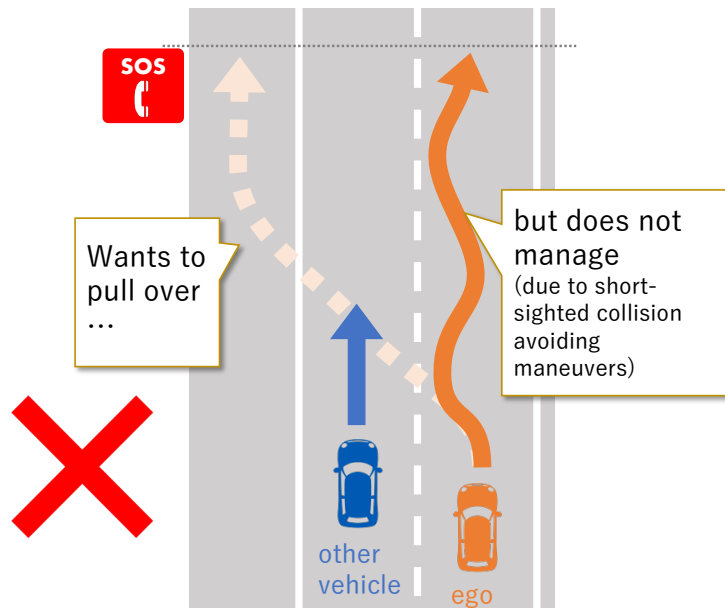
- "Divide and Conquer" complex driving scenarios
- Tool support by autom. reasoning

GA-RSS (our contribution)

Goal-Aware

Responsibility-Sensitive Safety [Hasuo+, IEEE T-IV, to appear]

- Guarantees goal achievement (e.g. successful pull over) and collision-freedom
- Global safety rules that combine mult. maneuvers
- Necessary for real-world complex driving scenarios



What is Formalization?

Informal pen-and-paper proofs



- Error-prone
- Poor traceability



Formal software-assisted proofs

$$\left\{ \begin{array}{l} \text{Env } \wedge I = 3 \wedge 0 \leq v \leq v_{\max} \\ \wedge y_{\text{tgt}} - y - \frac{v^2}{2a_{\text{brake}}} \geq 0 \end{array} \right\} \text{dwhile} \left(\frac{v^2}{2a_{\text{brake}}} < y_{\text{tgt}} - y \right) \left\{ \begin{array}{l} \text{Env } \wedge I = 3 \wedge 0 \leq v \leq v_{\max} \\ \wedge y_{\text{tgt}} - y - \frac{v^2}{2a_{\text{brake}}} = 0 \end{array} \right\}; \text{Env } \wedge I = 3 \wedge 0 \leq v \leq v_{\max} \\ \wedge -h_{\text{min}} \leq a \leq a_{\text{max}} \\ \wedge y_{\text{tgt}} - y - \frac{v^2}{2a_{\text{brake}}} \geq 0 \quad (18)$$

by (DWH) with $(\text{Env} \sim 0) = (v > 0), a_{\text{br}} = a_{\text{brake}}, y - y = \frac{v^2}{2a_{\text{brake}}}, a_{\text{br}} = -v$

$$\left\{ \begin{array}{l} \text{Env } \wedge I = 3 \wedge 0 \leq v \leq v_{\max} \\ \wedge \frac{v^2}{2a_{\text{brake}}} \leq y_{\text{tgt}} - y \end{array} \right\} \text{dwhile} \left(\frac{v^2}{2a_{\text{brake}}} < y_{\text{tgt}} - y \right) \left\{ \begin{array}{l} \text{Env } \wedge I = 3 \wedge 0 \leq v \leq v_{\max} \\ \wedge \frac{v^2}{2a_{\text{brake}}} = y_{\text{tgt}} - y \end{array} \right\}; \text{Env } \wedge I = 3 \wedge 0 \leq v \leq v_{\max} \\ \wedge -h_{\text{min}} \leq a \leq a_{\text{max}} \\ \wedge y \leq y_{\text{tgt}} \quad (19)$$

by (LIM) and (18)

$$\left\{ \begin{array}{l} \text{Env } \wedge I = 3 \wedge 0 \leq v \leq v_{\max} \\ \wedge \frac{v^2}{2a_{\text{brake}}} = y_{\text{tgt}} - y \end{array} \right\} \text{dwhile} (v > 0) \left\{ \begin{array}{l} \text{Env } \wedge I = 3 \wedge 0 \leq v \leq v_{\max} \\ \wedge \frac{v^2}{2a_{\text{brake}}} < y_{\text{tgt}} - y \end{array} \right\}; \text{Env } \wedge I = 3 \wedge 0 \leq v \leq v_{\max} \\ \wedge -h_{\text{min}} \leq a \leq a_{\text{max}} \\ \wedge \frac{v^2}{2a_{\text{brake}}} = y_{\text{tgt}} - y \quad (20)$$

$$\text{In}(\cdot) := \text{vSVBrake} = \text{vSVCruise} - \text{tBrake} * \text{aBrakeMin}$$

$$\left\{ \begin{array}{l} \text{Env } \wedge \\ \wedge \frac{v^2}{2a_{\text{brake}}} \end{array} \right\} \text{Out}(\cdot) := -\text{aBrakeMin} * \text{tBrake} + \text{vSVInit}$$

$$\text{In}(\cdot) := \text{xSVBrake} = \text{xSVCruise} + \text{Integrate}[\text{vSVCruise} - \text{t} * \text{aBrakeMin}, \{t, \theta, \text{tBrake}\}]$$

$$\left\{ \begin{array}{l} \text{Env } \wedge \\ \wedge \frac{v^2}{2a_{\text{brake}}} \end{array} \right\} \text{Out}(\cdot) := \frac{\text{aBrakeMin} * \text{tBrake}^2}{2} + \text{tBrake} * \text{vSVInit} + (-\text{tBrake} + \text{timeLaneChg}) * \text{vSVInit} + \text{xSVInit}$$

Fig. 13: replacing not explicitly present

$$\text{In}(\cdot) := \text{xSVFinal} = \text{xSVBrake}$$

$$\text{vSVFinal} = \text{vSVBrake}$$

$$\text{Out}(\cdot) := \frac{\text{aBrakeMin} * \text{tBrake}^2}{2} + \text{tBrake} * \text{vSVInit} + (-\text{tBrake} + \text{timeLaneChg}) * \text{vSVInit} + \text{xSVInit}$$

$$\text{Out}(\cdot) := -\text{aBrakeMin} * \text{tBrake} + \text{vSVInit}$$

$$\text{In}(\cdot) := \text{Equal} @ \text{postcond}$$

$$\text{Out}(\cdot) := \frac{\text{aBrakeMin} * \text{tBrake}^2}{2} - \text{tBrake} * \text{vSVInit} - (-\text{tBrake} + \text{timeLaneChg}) * \text{vSVInit} - \text{xSVInit} + \text{xTgt} = \frac{-\text{aBrakeMin} * \text{tBrake} + \text{vSVInit}}{2 * \text{aBrakeMin}}$$

- Symbolic proofs in our formal logical system dFHL
- Software tool checking the validity of each logical step of reasoning

Outline

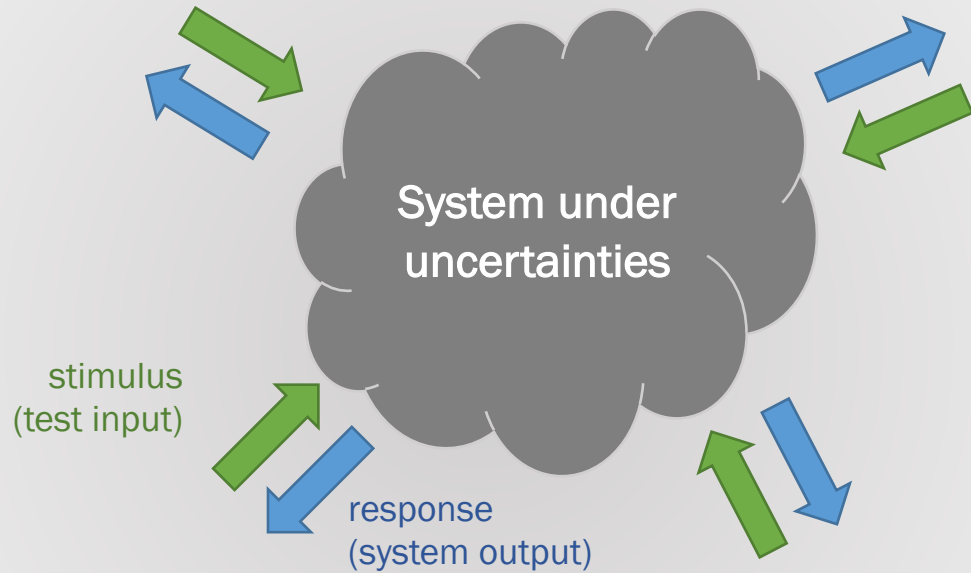
- A non-technical overview
- **The modeling problem**
- The RSS answer to the modeling problem
- Technical contributions: the logic dFHL
- Perspectives, practical & theoretical

The Modeling Problem

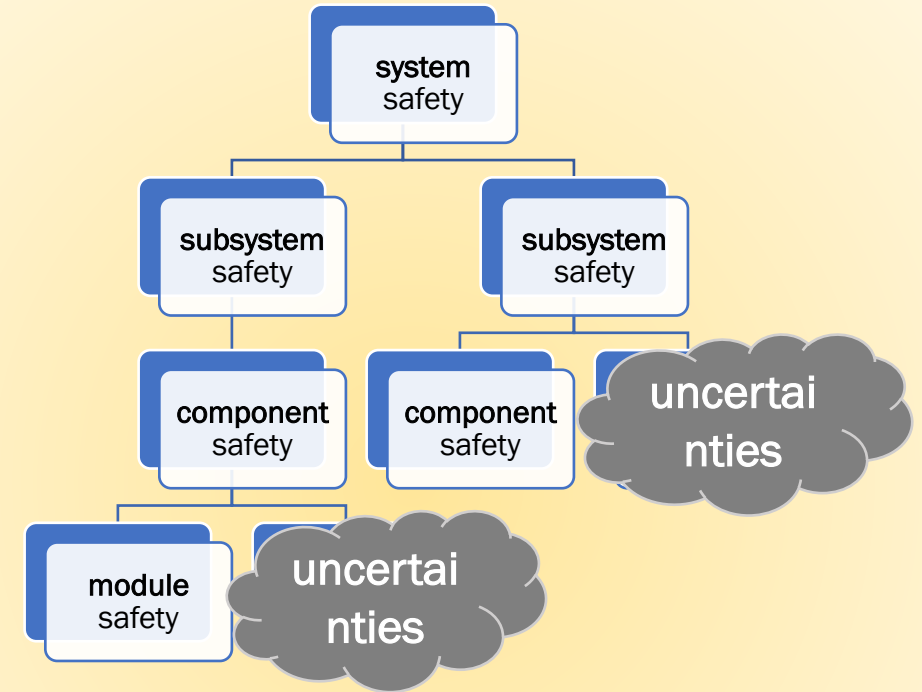
- Theorems need *definitions*;
formal verification needs *modeling*
- Automated driving systems (ADS) are **assively complex** system
 - Hundreds of chips, millions of LoC
 - Physical components. Internal combustion
 - ML components, especially for perception
 - Unpredictable road conditions
 - Other cars
 - Pedestrians
 - ...
- Modeling is hard (a grand challenge for us)



Logical Confinement of Uncertainties



- The whole system as a monolithic blackbox
- Analyzed by statistical and empirical means
- E.g. automated driving:
"1 fatality per XXX miles driven"
→ Doesn't exclude a scenario that is always fatal



- Logical argumentation of safety cases
- Impose rules/contracts on uncertain components
→ runtime monitoring, accountability, identifying causes of accident
- Finding a good "logical angle" is crucial, which takes theoretical insights and experience

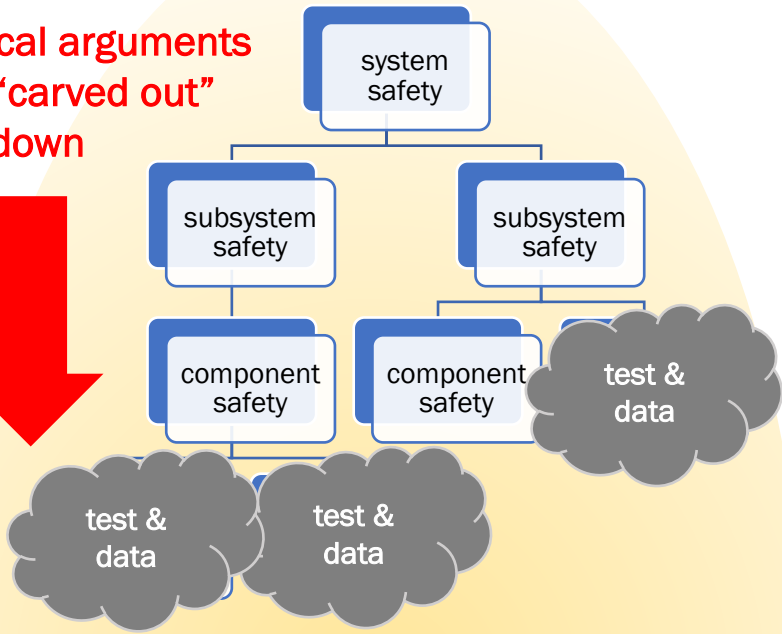


Purely data-driven approach to safety assurance

e.g. “one derailment every 10,000 miles” in automated driving

- ✓ Scalability, automation by efficient processing of big data
- ✗ Accountability. Hard to convince the customer/public of safety, or that duties of care have been fulfilled

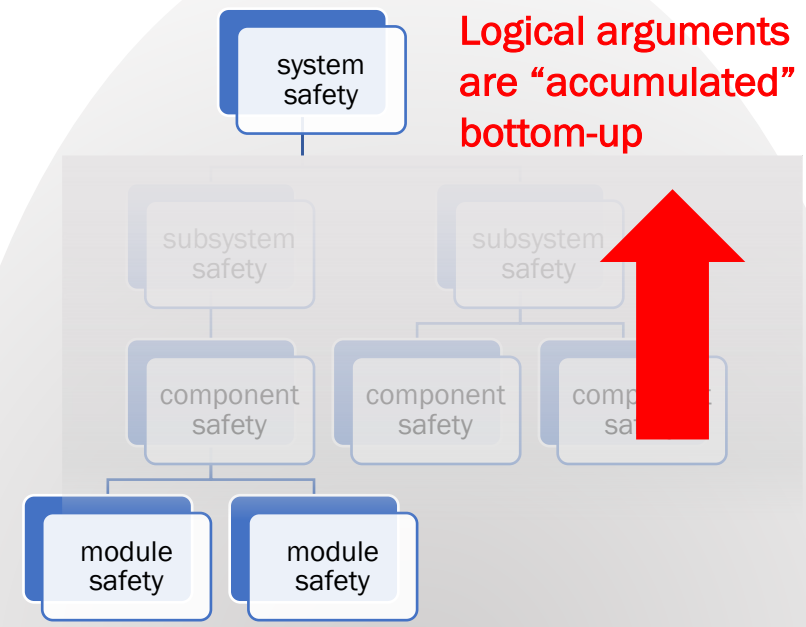
Logical arguments are “carved out” top-down



(Our approach) Logical confinement of uncertainties

- Start from the conclusion (system safety), and carve out logical arguments that lead to it
- Use test & data once the limit of logical arguments is reached
- ✓ Best-effort logical guarantee
Smaller resources/efforts yield non-zero assurance (if smaller)
- ✓ Explainability by logic.
Crucial for public acceptance of new ICT paradigms (such as automated driving)

Logical arguments are “accumulated” bottom-up



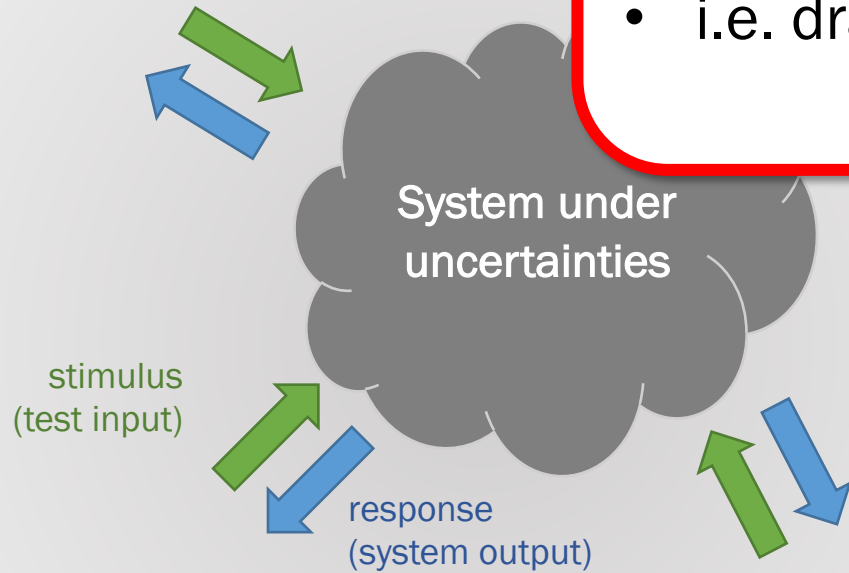
Purely logical approach to safety assurance

- Formal verification, a software science tradition
- Start with mathematical modeling of the target system, and build up logical consequences
- ✓ Traceability. Accountability. Trust. Every deduction step is explicit and rule-based.
- ✗ Complexity of modern ICT systems
→ Bottom-up efforts might never reach the final goal (namely the system safety)
- ✗ Moreover, an incomplete proof is totally useless. Huge cost until a non-zero value is produced

Logical Con

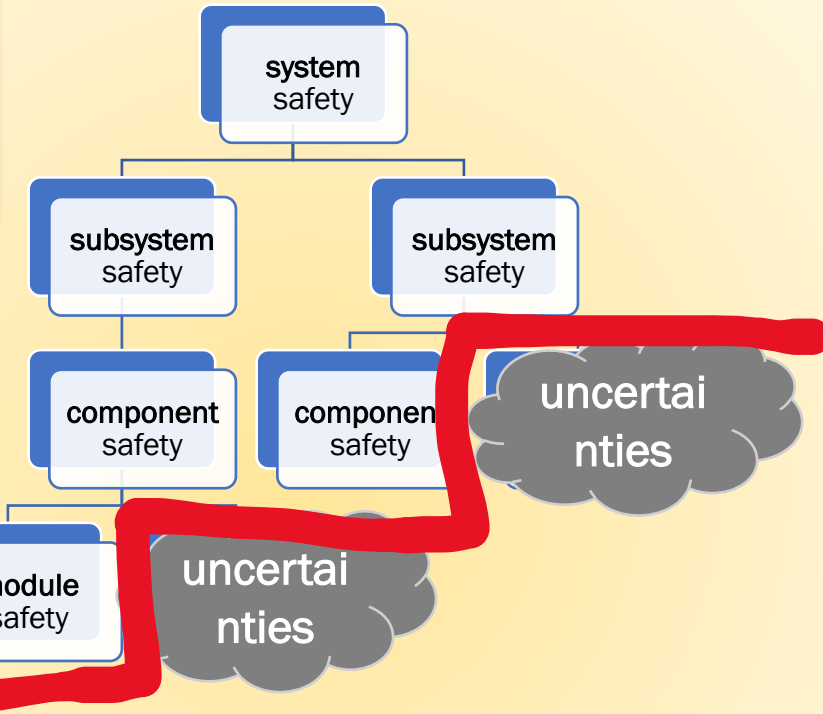
The modeling problem:

- Deciding **what to model** and **what not to model**
- i.e. drawing a good border



- The whole system as a monolithic blackbox
- Analyzed by statistical and empirical means
- E.g. automated driving:
"1 fatality per XXX miles driven"
→ Doesn't exclude a scenario that is always fatal

nties



- **Logical argumentation of safety cases**
- Impose **rules/contracts** on uncertain components
→ runtime monitoring, accountability, identifying causes of accident
- Finding a good "**logical angle**" is crucial, which takes theoretical insights and experience

Outline

- A non-technical overview
- The modeling problem
- The RSS answer to the modeling problem
- Technical contributions: the logic dFHL
- Perspectives, practical & theoretical

Logical Con

nties

The modeling problem:

- Deciding **what to model** and **what not to model**
- i.e. drawing a good border

RSS's answer

Responsibility-Sensitive Safety (RSS)

[Shalev-Shwartz et al., arXiv preprint, 2017]

Lemma.

If all cars comply with RSS rules, then there is no collision

mathematically proved

+

Lemma.

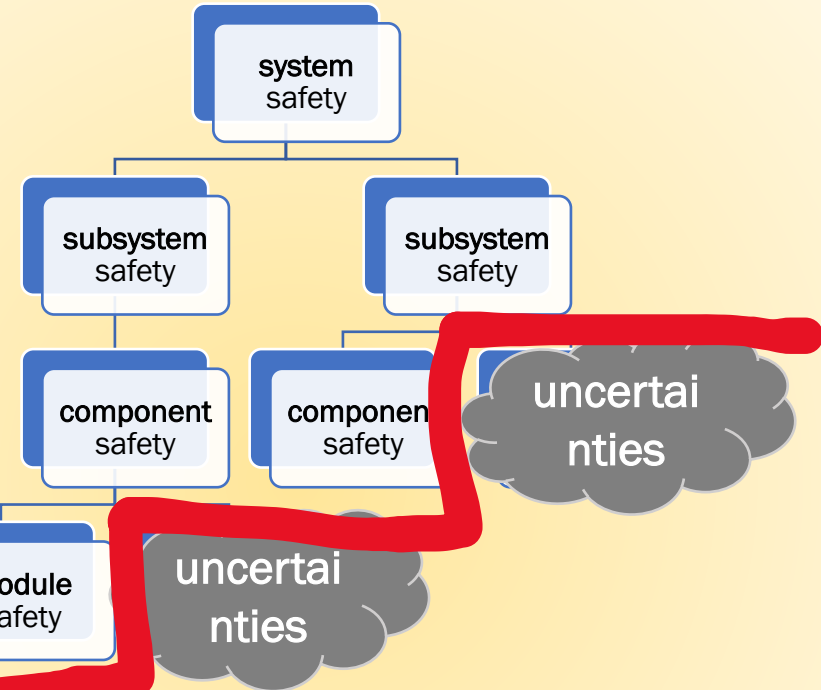
All cars comply with RSS rules

guaranteed by the manufacturer (testing, proofs, etc.)

→

Theorem.

There is no collision



- Logical argumentation of safety cases
- Impose **rules/contracts** on uncertain components
→ runtime monitoring, accountability, identifying causes of accident
- Finding a good “**logical angle**” is crucial, which takes theoretical insights and experience

Safety Guarantee for Automated Driving via Logical Safety Rules and Mathematical Proofs



“I’m safe since I respect the safety rules R_1, R_2, \dots ”

“I’m safe since I respect the safety rules R_1, R_2, \dots ”



“I’m safe since I respect the safety rules R_1, R_2, \dots ”

- Decompose safety (a complex goal) into logical safety rules (explicit, easy to check and enforce)
- “Ultimate assurance” in the form of mathematical proofs. Logical explanation by following their reasoning steps
- Safety rules are generic and reusable
→ regulation, standard → social acceptance
- Attribution of liabilities (collision → someone must have broken the rules)

Safety Rule R_1

In the *same-lane same-direction* driving scenario,

- Maintain the safety distance

$$d_{\min} = \left[v_r \rho + \frac{1}{2} a_{\max, \text{accel}} \rho^2 + \frac{(v_r + \rho a_{\max, \text{accel}})^2}{2a_{\min, \text{brake}}} - \frac{v_f^2}{2a_{\max, \text{brake}}} \right]_+$$

from the preceding car

- When that’s hard, brake at acceleration $a_{\max, \text{brake}}$

Theorem (Safety)

There is no collision attributed to the ego vehicle as long as the safety rule R_1 is respected

Proof (of the safety thm.)

The only non-obvious point is that $e_{\text{inv},2}$ is preserved by the dynamics. We first observe

$$\mathcal{L}_{\delta_j, \delta_r} e_{\text{inv},2} = \begin{cases} 0 & \text{if } \text{dRSS}_{\pm}(v_f, v_r, \rho - t) \geq 0 \\ v_f - v_r & \text{otherwise,} \end{cases}$$

where $\text{dRSS}_{\pm}(v_f, v_r, \rho)$ is given by

$$\text{dRSS}_{\pm}(v_f, v_r, \rho) = v_r \rho + \frac{a_{\max} \rho^2}{2} + \frac{(v_r + a_{\max} \rho)^2}{2b_{\min}} - \frac{v_f^2}{2b_{\max}}$$

Therefore, we can infer as follows.

$$\begin{aligned} & \text{dRSS}_{\pm}(v_f, v_r, \rho - t) < 0 \\ \iff & v_r(\rho - t) + \frac{a_{\max}(\rho - t)^2}{2} + \frac{(v_r + a_{\max}(\rho - t))^2}{2b_{\min}} - \frac{v_f^2}{2b_{\max}} < 0 \end{aligned}$$

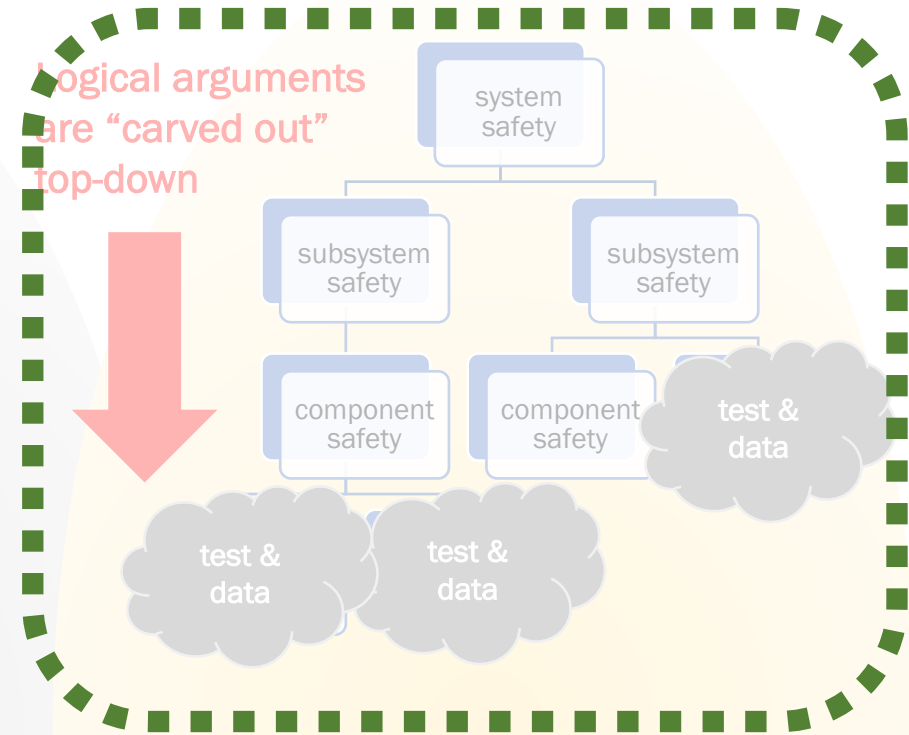
R_1
 R_2
 R_3
 \dots



Purely data-driven approach to safety assurance

e.g. "one derailment every 10,000 miles" in automated driving

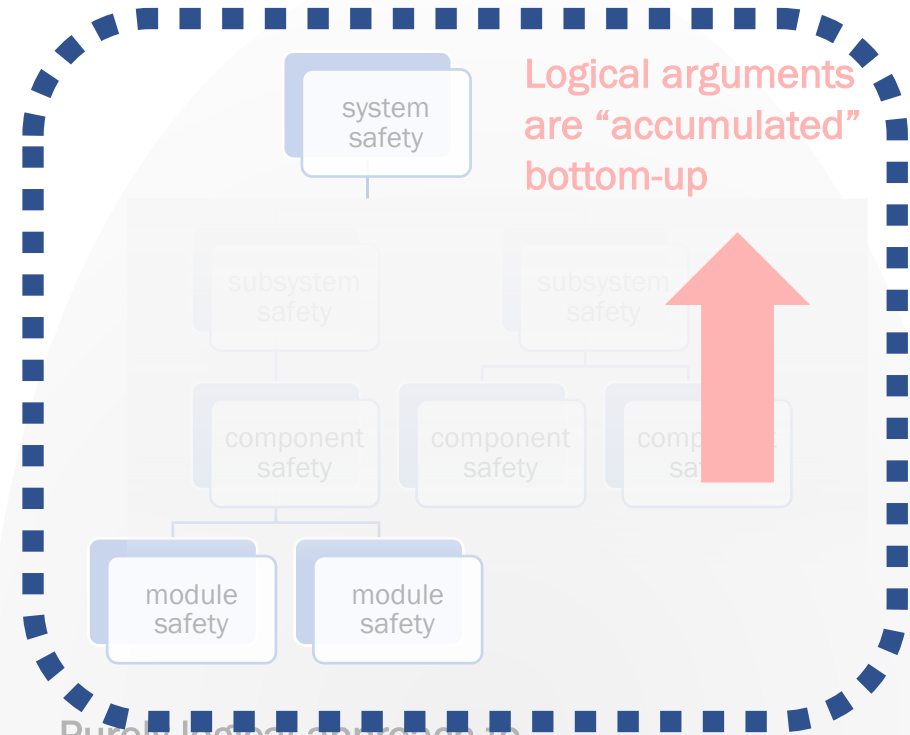
- ✓ Scalability, automation by efficient processing of big data
- ✗ Accountability. Hard to convince the customer/public of safety, or that duties of care have been fulfilled



(Our approach)

Logical safety rules for ADS

- Safety "theorems" are reduced to "axioms" (namely safety rules)
- The reduction is math. rigorous
- Rule compliance can be logically verified, but also be tested or monitored
- Not a full safety proof, but feasible. Enough for many usages



Purely logical approach to safety assurance

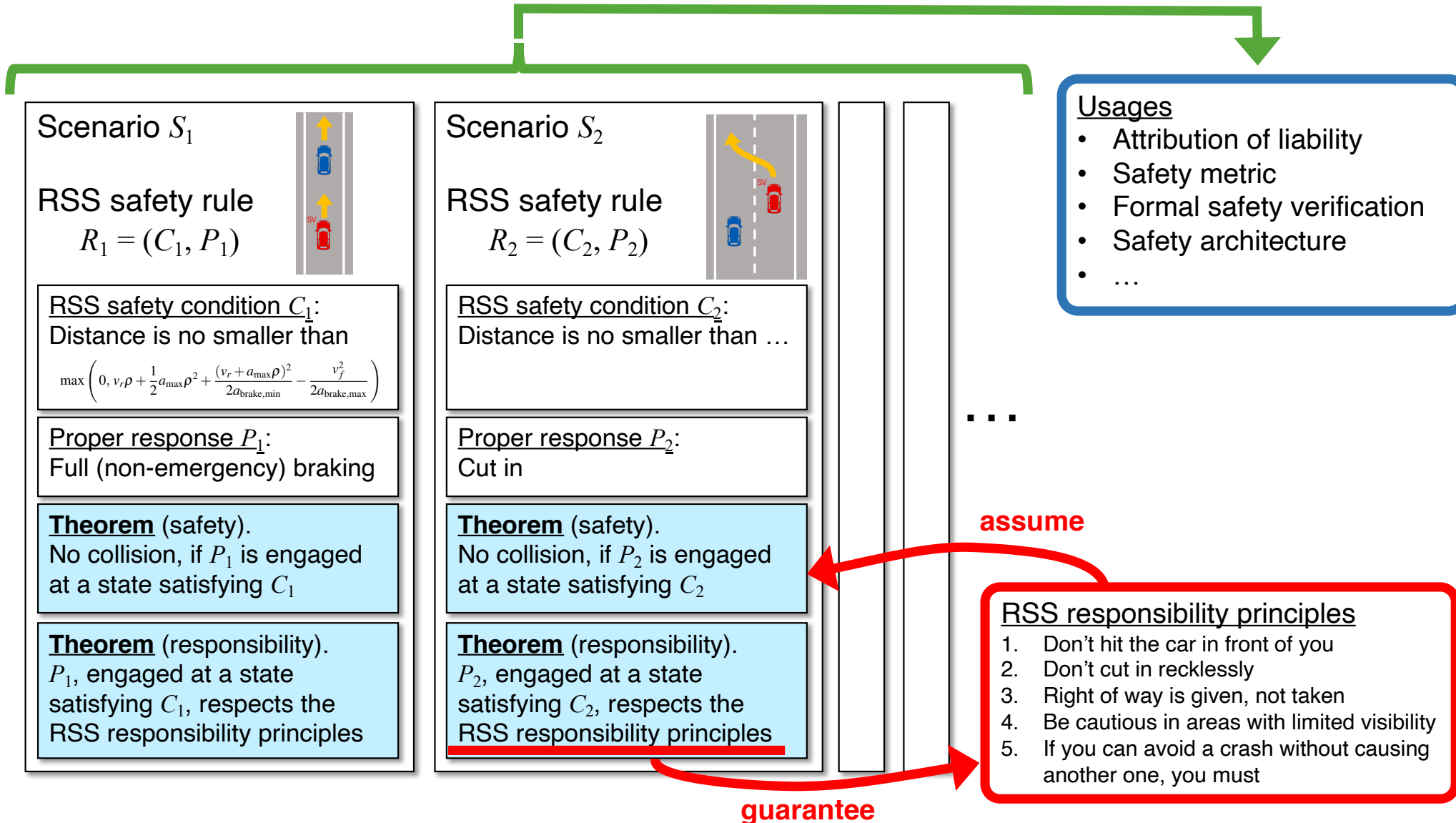
(Full) formal verif. of ADS safety

- Strong mathematical guarantee... if the proof is completed
- and completion is **very unlikely** (lack of models, budget limit, uncertainties in driving scenarios, etc.)
- (Full verif. is more viable in aerospace. No pedestrians, more budget)

RSS Framework

[Shalev-Shwartz et al., arXiv, 2017]
See also [Hasuo, arXiv 2206.03418]

Each rule consists of
a **condition** and a **proper response**



Outline

- A non-technical overview
- The modeling problem
- The RSS answer to the modeling problem
- Technical contributions: the logic dFHL
- Perspectives, practical & theoretical

Our Contribution: Formal Logic Foundations of RSS → More Scenarios

RSS

Responsibility-Sensitive Safety, Shalev-Shwartz et al., 2017

- Basic methodology of logical safety rules
- Standardization (IEEE 2846)
- Lack of formal implementation → **appl. to complex scenarios is hard**
- Guarantees only collision-freedom so far

↓ Software science research

differential program logic dFHL (our contribution)

$$\frac{\begin{array}{l} \text{inv: } A \Rightarrow e_{\text{inv}} \sim 0 \quad e_{\text{var}} \geq 0 \wedge e_{\text{inv}} \sim 0 \Rightarrow \mathcal{L}\dot{x} = f e_{\text{inv}} \geq 0 \\ \text{var: } A \Rightarrow e_{\text{var}} \geq 0 \quad e_{\text{var}} \geq 0 \wedge e_{\text{inv}} \sim 0 \Rightarrow \mathcal{L}\dot{x} = f e_{\text{var}} \leq e_{\text{ter}} \\ \text{ter: } A \Rightarrow e_{\text{ter}} < 0 \quad e_{\text{var}} \geq 0 \wedge e_{\text{inv}} \sim 0 \Rightarrow \mathcal{L}\dot{x} = f e_{\text{ter}} \leq 0 \end{array}}{\{A\} \text{dwhile}(e_{\text{var}} > 0) \dot{x} = f \{e_{\text{var}} = 0 \wedge e_{\text{inv}} \sim 0\} : e_{\text{inv}} \sim 0 \wedge e_{\text{var}} \geq 0} \text{ (DWH)}$$

- A logical system for deriving and proving safety rules

GA-RSS (our contribution)

Goal-Aware

Responsibility-Sensitive Safety

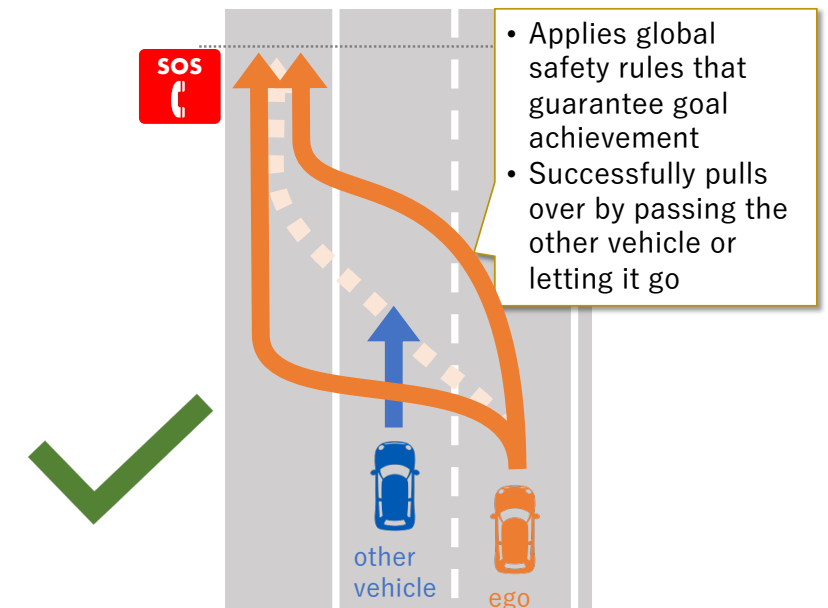
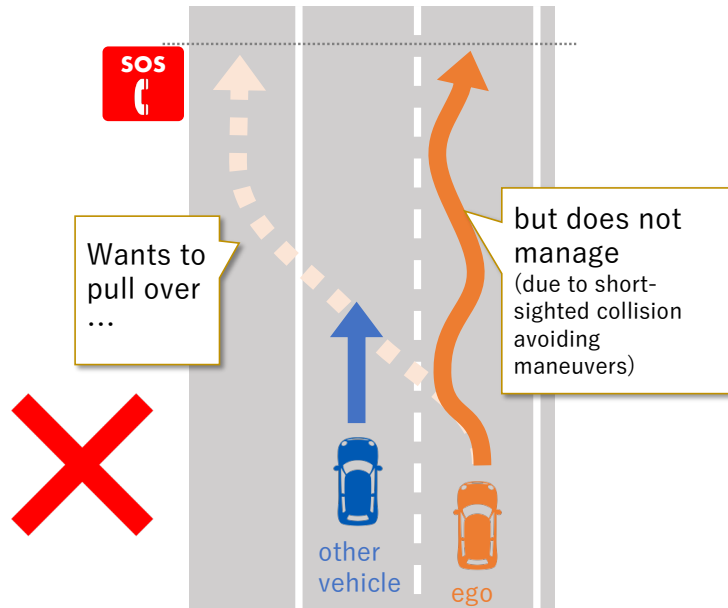
- Guarantees **goal achievement** (e.g. successful pull over) and collision-freedom
- Global safety rules that combine mult. maneuvers
- Necessary for real-world complex driving scenarios

Compositional rule derivation workflow by dFHL (our contribution)

(our contribution)



- "Divide and Conquer" complex driving scenarios
- Tool support by autom. reasoning



Differential program logic dFHL



- Hoare logic
+ ODEs (dwhile)
+ “safety condition”

$$\{A\} \alpha \{B\} : S$$

postcondition \uparrow
(true at the end of α)

“safety condition” \uparrow
(true throughout α)

- Reasoning about ODEs via differential invariants (barrier cert.) and ranking/Lyapunov functions
- Theoretically not so much different from Platzer’s dL.
Simplified, aiding proof engineers

Def. (dFHL programs)

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

Def. (dFHL rules)

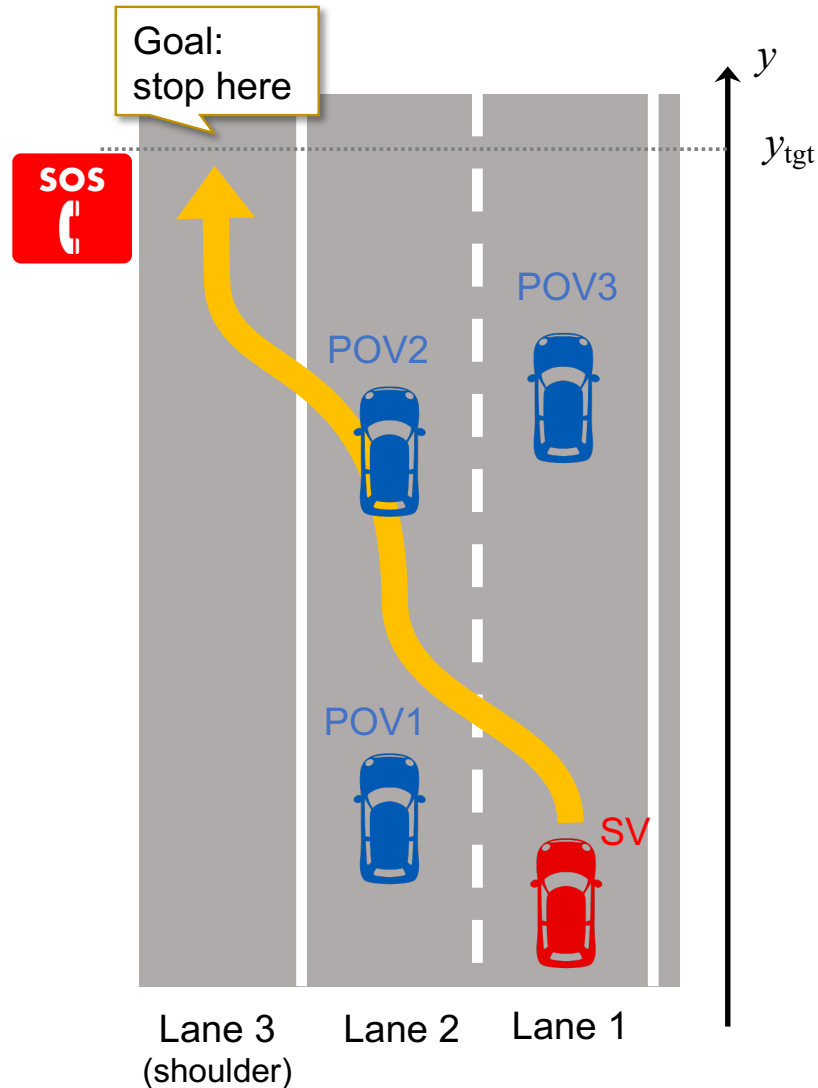
$$\frac{\{A\} \alpha \{B\} : S \quad \{B\} \beta \{C\} : S}{\{A\} \alpha; \beta \{C\} : S} \text{ (SEQ)}$$

$$\frac{\{A'\} \alpha \{B'\} : S' \quad \begin{array}{l} A \Rightarrow A' \\ S' \wedge B' \Rightarrow B \\ S' \Rightarrow S \end{array}}{\{A\} \alpha \{B\} : S} \text{ (LIMP)}$$

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

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

Compositional Rule Derivation



- We shall derive

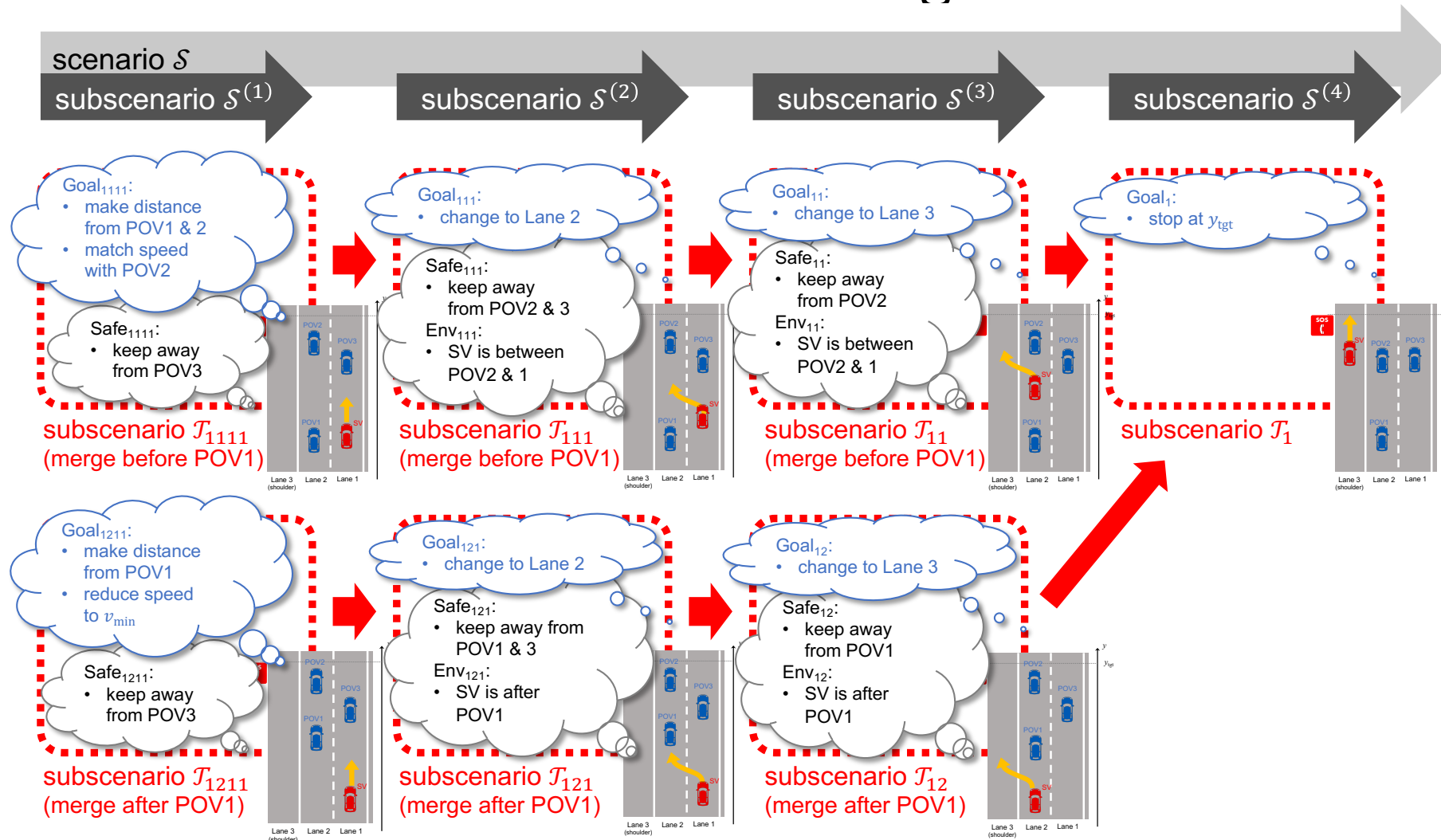
$$\{A\} \alpha \{B\} : S$$

for the following given data

- B is the **goal**: “stopping on the shoulder at y_{tgt} ”
- S is the **safety**: “no collision,” or better “securing RSS distance from every other car”
- We shall identify
 - α as an **RSS proper response**: “executing α will safely achieve the goal”
 - A as an **RSS condition**: “when A is true, B and S are guaranteed by executing α ”

Compositional Rule Derivation

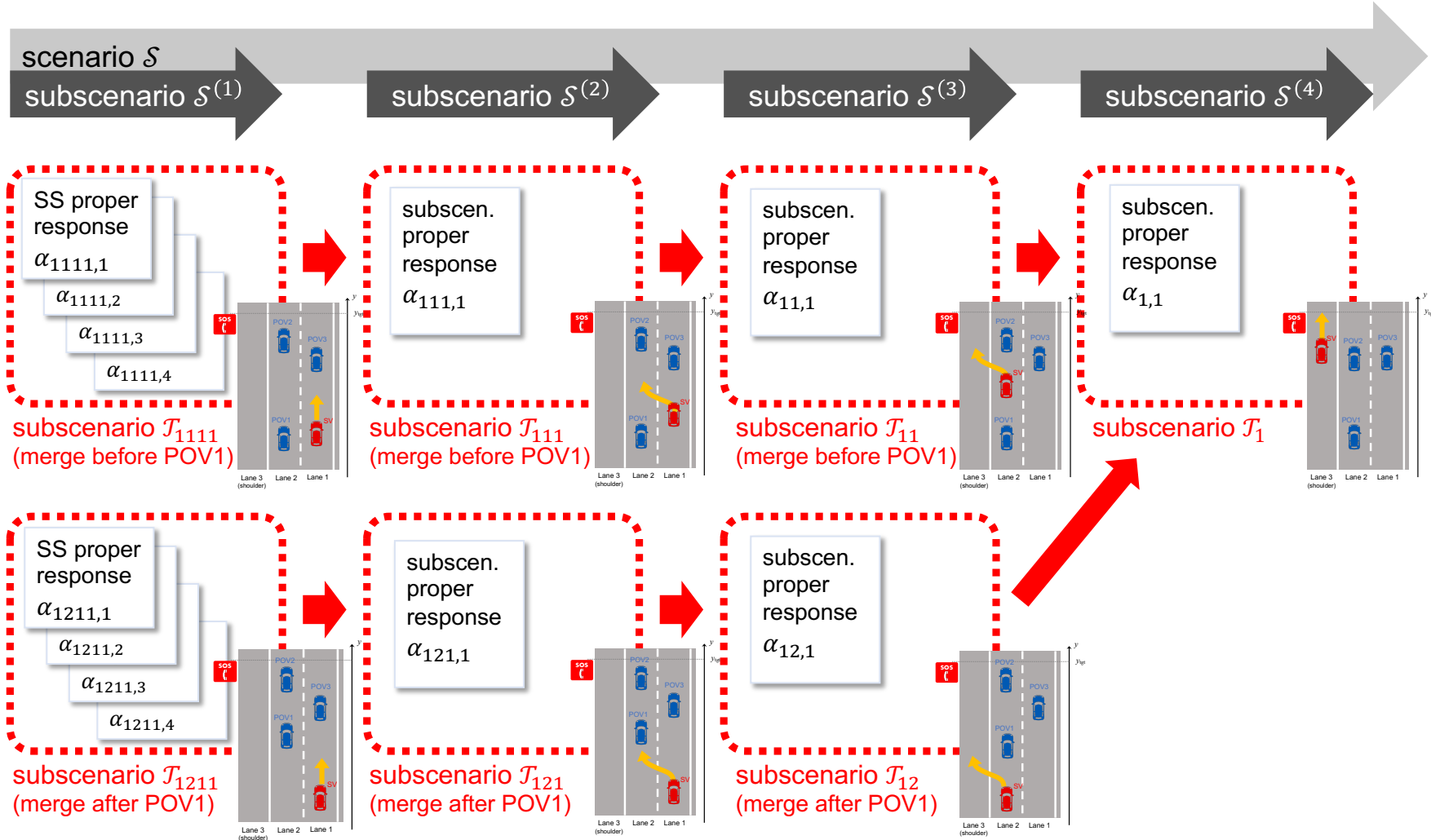
(1) Decompose the scenario into **subscenarios**, each of which has clearer focuses and goals



Compositional Rule Derivation

(2) Devise **subscenario proper responses** for each subscenario

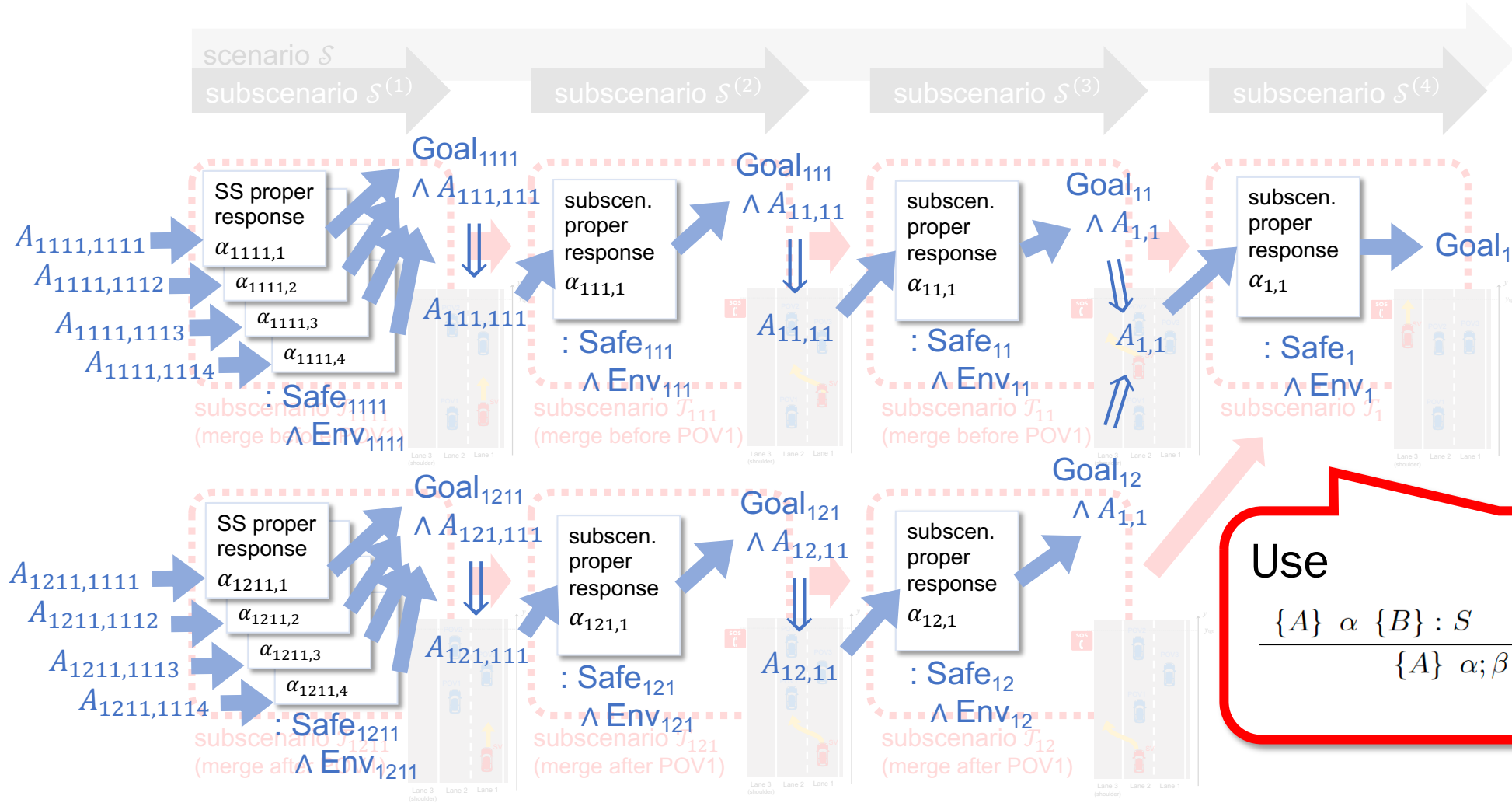
$$\{A\} \alpha \{B\} : S$$



Compositional Rule Derivation

(3) Backpropagate pre/postconditions, leading to the scenario-wide precondition

$$\{A\} \alpha \{B\} : S$$



Use

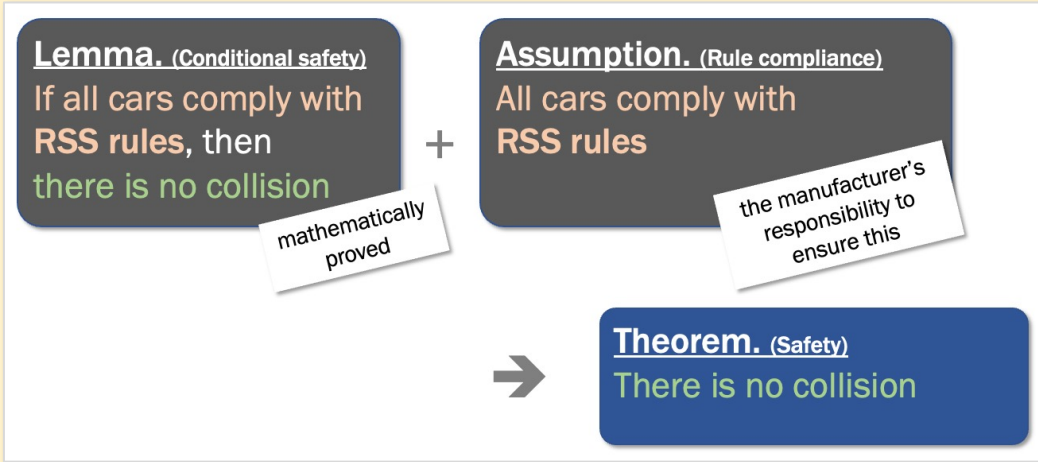
$$\frac{\{A\} \alpha \{B\} : S \quad \{B\} \beta \{C\} : S}{\{A\} \alpha; \beta \{C\} : S} \text{ (SEQ)}$$

Outline

- A non-technical overview
- The modeling problem
- The RSS answer to the modeling problem
- Technical contributions: the logic dFHL
- Perspectives, practical & theoretical

Logical Formalization of RSS

Covering More Scenarios → Real-World Deployment



- RSS as in [Shalev-Shwartz et al., arXiv, 2017] is a **methodology** – it is sensible and promising, but came with no proof technologies
- thus application was limited to simple driving scenarios



What is Formalization?

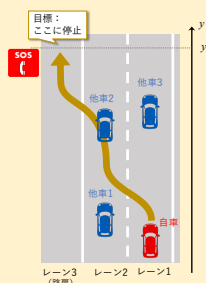
Informal
pen-and-paper proofs

- Error-prone
- Poor traceability

Formal
software-assisted proofs

- Symbolic proofs in our formal logical system dFHL
- Software tool checking the validity of each logical step of reasoning

- Our contribution [Hasuo+, IEEE T-IV, to appear]: **Logical technologies** to prove *conditional safety lemmas* for complex scenarios
- Compositional proofs, ensuring goal achievements, ...
- Much more scenarios proved safety by RSS → RSS at work → social acceptance of ADV



Safety Guarantee for Automated Driving via Logical Safety Rules and Mathematical Proofs



“I’m safe since I respect the safety rules R_1, R_2, \dots ”

“I’m safe since I respect the safety rules R_1, R_2, \dots ”



“I’m safe since I respect the safety rules R_1, R_2, \dots ”

- Decompose safety (a complex goal) into logical safety rules (explicit, easy to check and enforce)
- “Ultimate assurance” in the form of mathematical proofs. Logical explanation by following their reasoning steps
- Safety rules are generic and reusable
→ regulation, standard → social acceptance
- Attribution of liabilities (collision → someone must have broken the rules)

Safety Rule R_1

In the *same-lane same-direction* driving scenario,

- Maintain the safety distance

$$d_{\min} = \left[v_r \rho + \frac{1}{2} a_{\max, \text{accel}} \rho^2 + \frac{(v_r + \rho a_{\max, \text{accel}})^2}{2a_{\min, \text{brake}}} - \frac{v_f^2}{2a_{\max, \text{brake}}} \right]_+$$

from the preceding car

- When that’s hard, brake at acceleration $a_{\max, \text{brake}}$

Theorem (Safety)

There is no collision attributed to the ego vehicle as long as the safety rule R_1 is respected

Proof (of the safety thm.)

The only non-obvious point is that $e_{\text{inv},2}$ is preserved by the dynamics. We first observe

$$\mathcal{L}_{\delta_j, \delta_r} e_{\text{inv},2} = \begin{cases} 0 & \text{if } d\text{RSS}_{\pm}(v_f, v_r, \rho - t) \geq 0 \\ v_f - v_r & \text{otherwise,} \end{cases}$$

where $d\text{RSS}_{\pm}(v_f, v_r, \rho)$ is given by

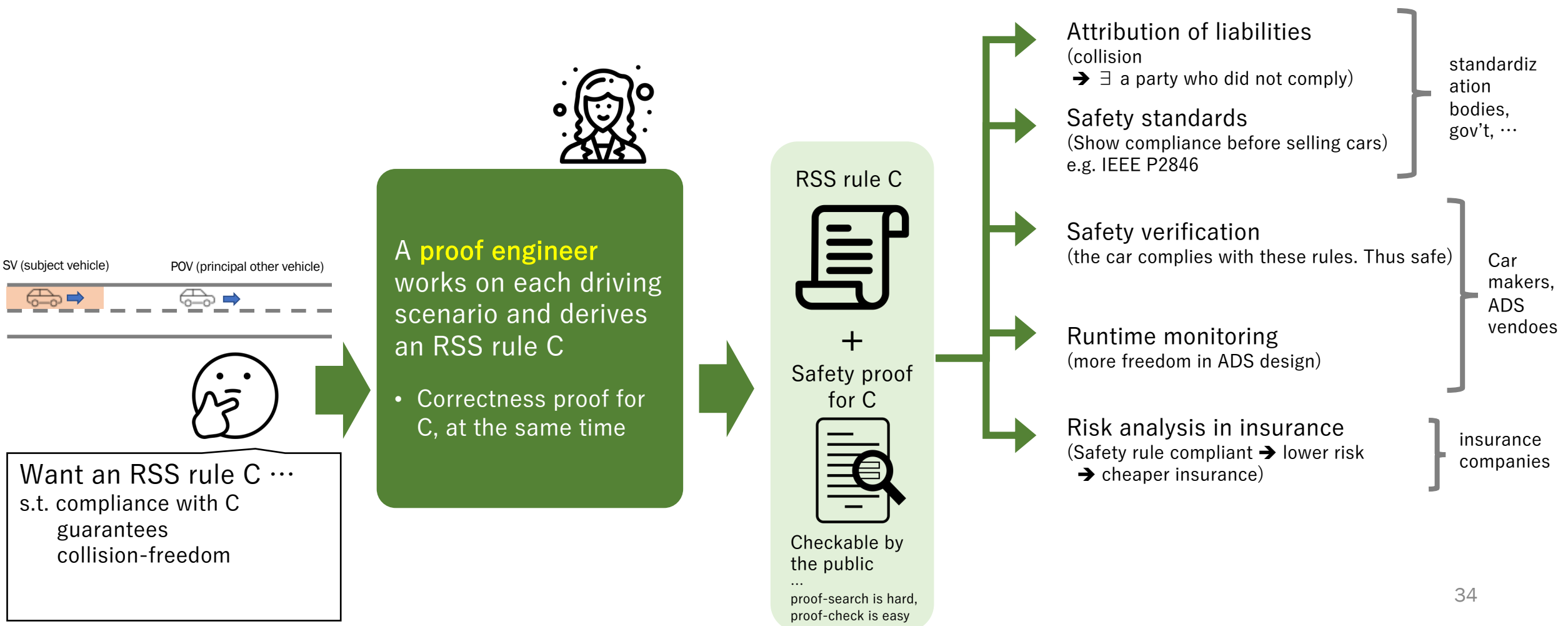
$$d\text{RSS}_{\pm}(v_f, v_r, \rho) = v_r \rho + \frac{a_{\max} \rho^2}{2} + \frac{(v_r + a_{\max} \rho)^2}{2b_{\min}} - \frac{v_f^2}{2b_{\max}}$$

Therefore, we can infer as follows.

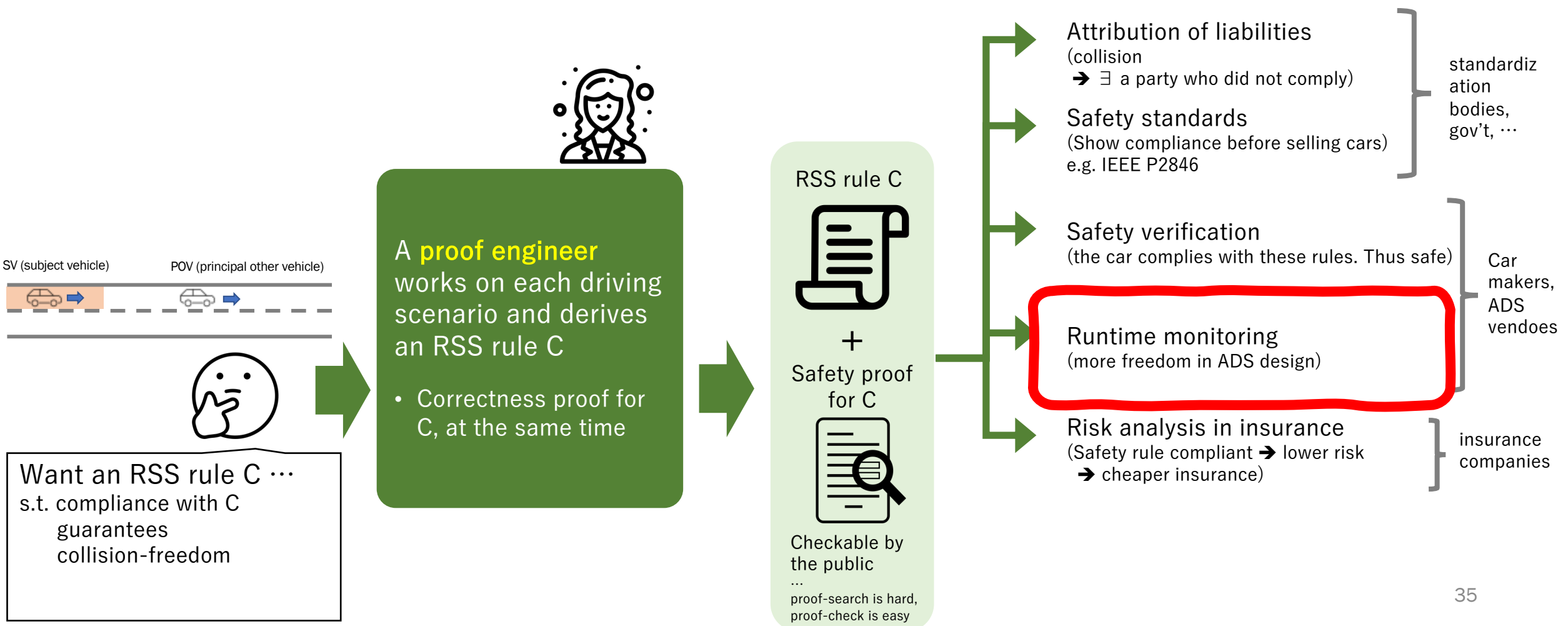
$$\begin{aligned} d\text{RSS}_{\pm}(v_f, v_r, \rho - t) &< 0 \\ \iff v_r(\rho - t) + \frac{a_{\max}(\rho - t)^2}{2} + \frac{(v_r + a_{\max}(\rho - t))^2}{2b_{\min}} - \frac{v_f^2}{2b_{\max}} &< 0 \end{aligned}$$

R_1
 R_2
 R_3
 \dots

RSS Rules as Social Contracts Impacts Everywhere in the ADV Ecosystem



RSS Rules as Social Contracts Impacts Everywhere in the ADV Ecosystem



Safety Envelope by RSS Rules

Can Be Retrofit to Any ADV Controller Monitor & Intervene → Runtime Safety Guarantee

RSS Rule, an Example

[Shalev-Shwartz et al., arXiv preprint, 2017]



- An RSS rule is a pair (A, α) of an **RSS condition** A and a **proper response** α

RSS condition A :

Maintain an inter-vehicle distance at least

$$d_{\min} = \left[v_r \rho + \frac{1}{2} a_{\max, \text{accel}} \rho^2 + \frac{(v_r + \rho a_{\max, \text{accel}})^2}{2a_{\min, \text{brake}}} - \frac{v_f^2}{2a_{\max, \text{brake}}} \right]_+$$

Proper response α :

If A is about to be violated, brake at rate $a_{\min, \text{brake}}$ within ρ seconds

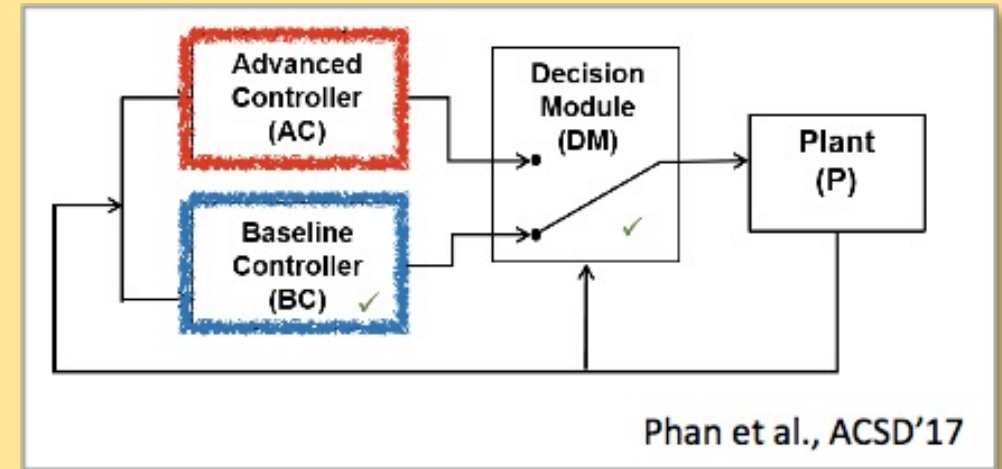
Conditional safety lemma:

Any execution of α , from a state that satisfies A , is collision-free.

Structure of an RSS rule

- RSS Condition A :
“You can still **escape** if A is true”
- Proper response α :
“control strategy to **escape**”

escape =
MRM
(minimum risk maneuver)



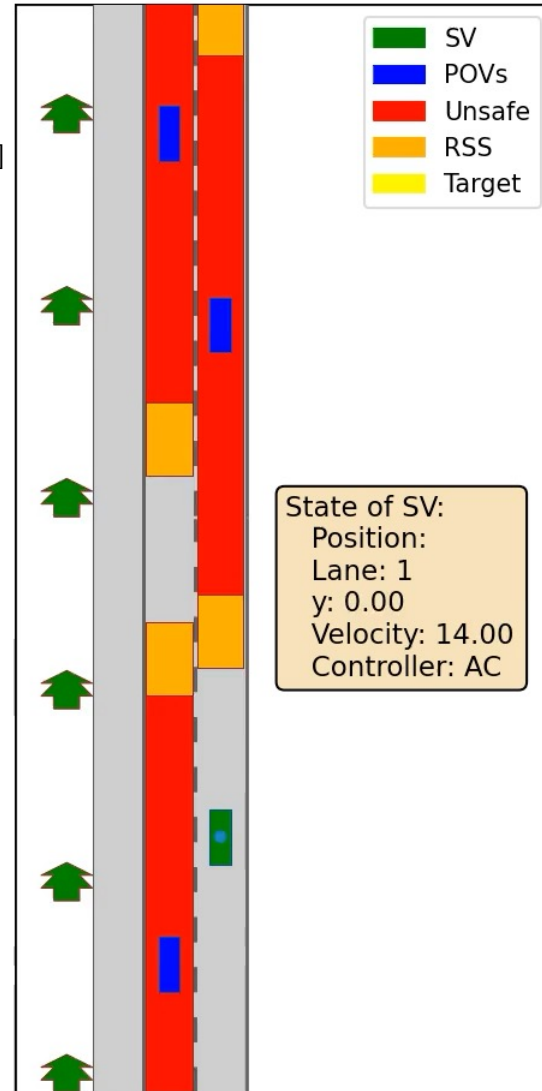
Phan et al., ACSD'17

Simplex architecture

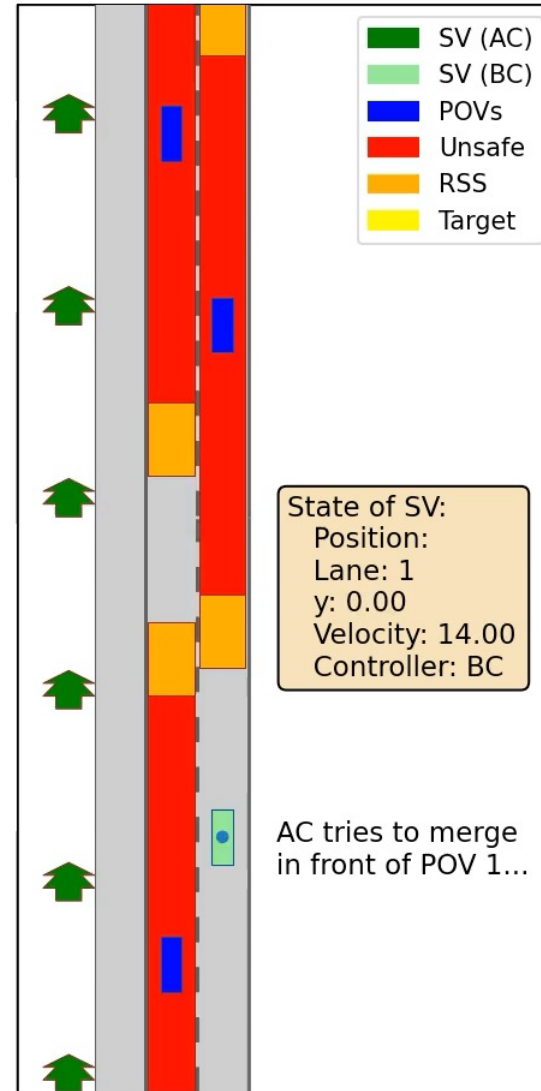
- AC pursues performance and safety
 - BC pursues safety (only)
 - DM (decision module) switches between them—
“use BC to escape”
- RSS rules fit perfectly!
- AC: existing controller (optimization-based, ML, ...)
 - BC: executes a proper response
 - DM: monitors an RSS condition. Violation foreseen → switch to BC

RSS Safety Envelopes in Action, Scenario I

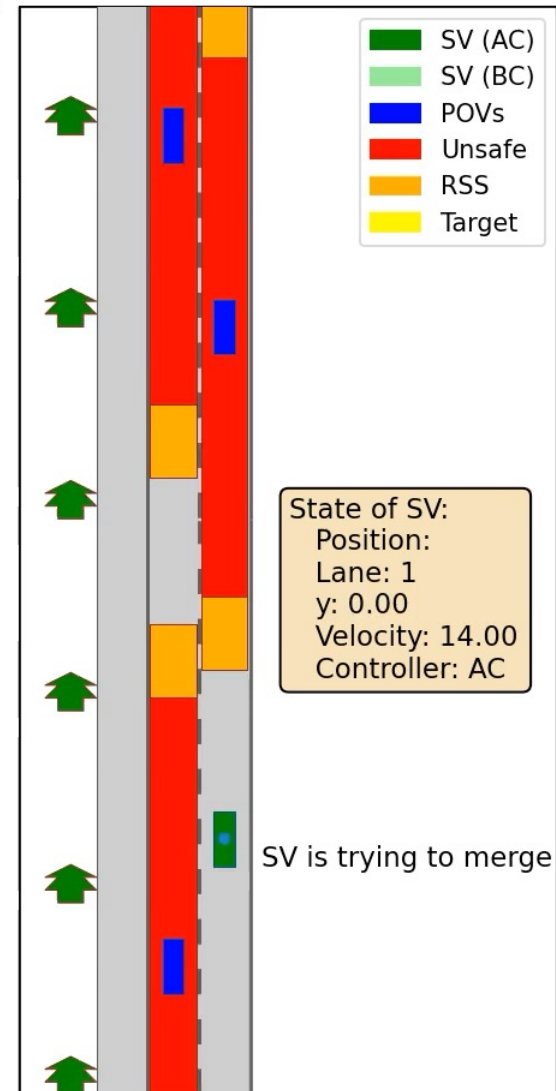
- **AC**: no safety envelope
 - **AC+RSS**:
Original RSS rule [Shalev-Shwartz et al., arXiv, 2017]
as a safety envelope
("short-sighted" collision avoidance)
 - **AC+RSS^{GA}**:
Our RSS rule [Hasuo+, IEEE T-IV]
as a safety envelope
(goal achievement too with longer-term
planning)
- **AC** is not safe (hazardous cut-in)
 - **AC+RSS** does not reach the shoulder
 - **AC+RSS^{GA}** successfully deployed the long term strategy of (brake → merge behind).
Achieved both safety and the goal



AC



AC+RSS



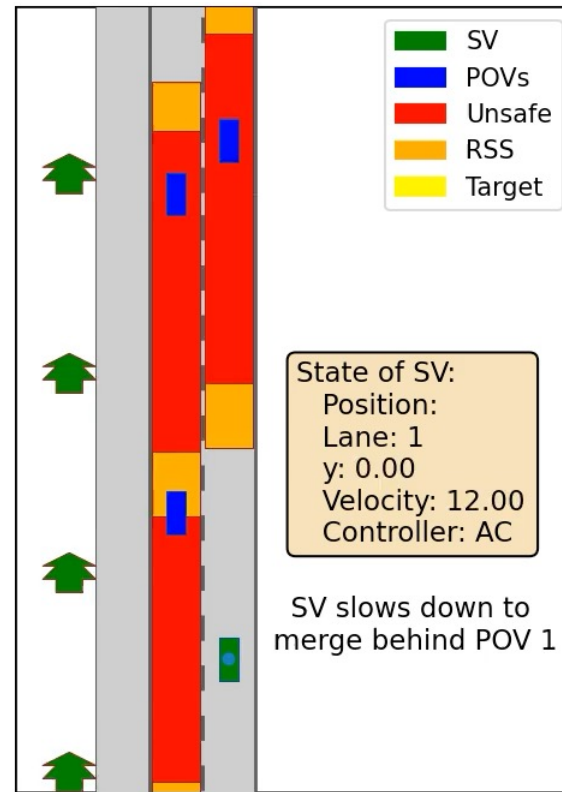
AC+RSS^{GA}

RSS Safety Envelopes in Action, Scenario II

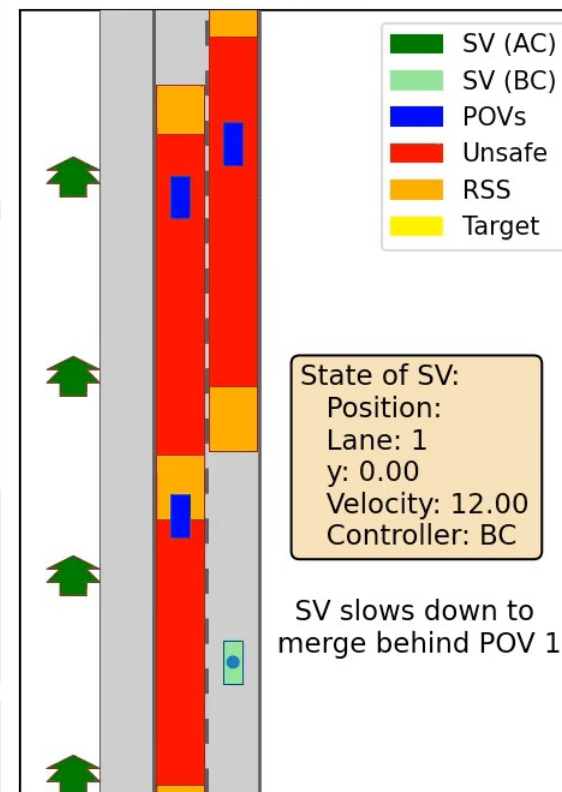
- **AC**: no safety envelope
- **AC+RSS**:
Original RSS rule
[Shalev-Shwartz et al., arXiv, 2017]
as a safety envelope
("short-sighted" collision avoidance)
- **AC+RSS^{GA}** :
Our RSS rule [Hasuo+, IEEE T-IV]
as a safety envelope
(goal achievement too
with longer-term planning)

- **AC** & **AC+RSS** safety achieve the goal, but are **slow**

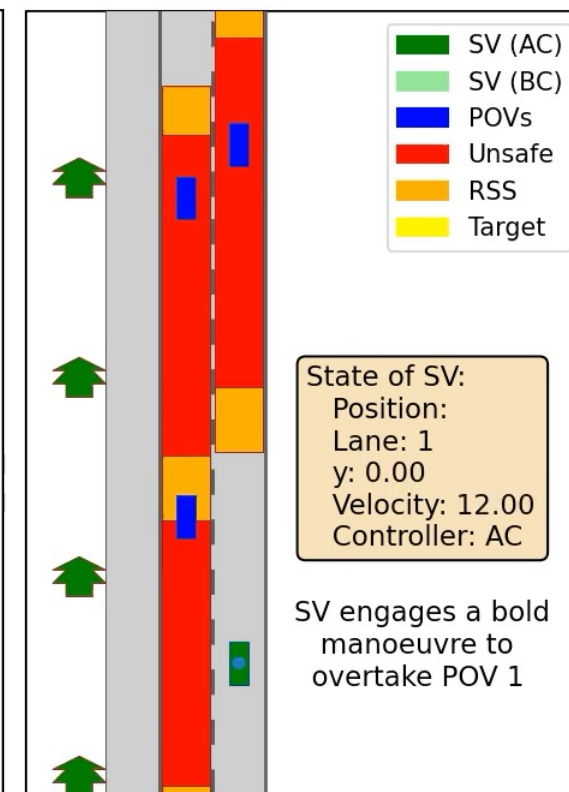
- **AC+RSS^{GA}**,
under mathematical safety guarantee,
boldly accelerates and merge in front



AC



AC+RSS

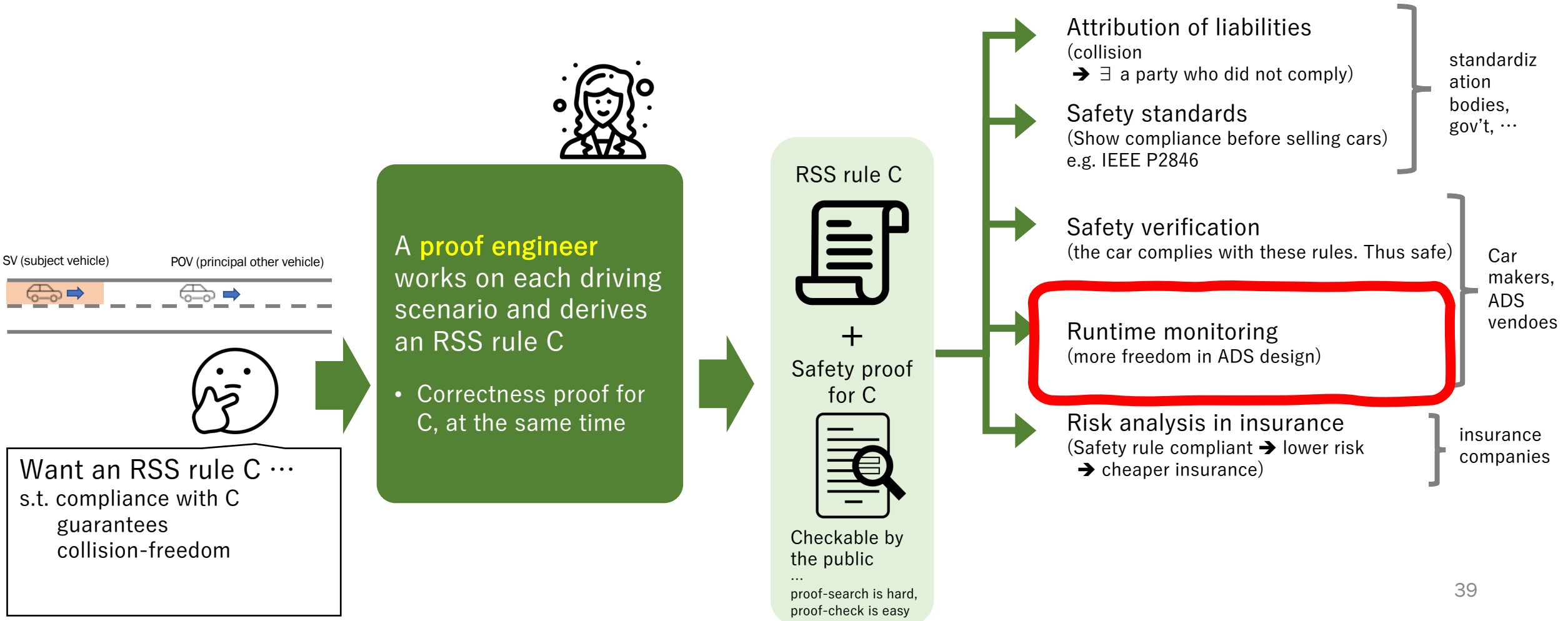


AC+RSS^{GA}

- ... who says safe ADVs are conservative and boring? 😊



RSS Rules as Social Contracts Impacts Everywhere in the ADV Ecosystem



Two Different Approaches, with Different Business Models



Fixed-route bus, taxi, delivery service



Consumer ADV

remote	human intervention	on-site (human driver)
offers fixed-route mobility and delivery service	business model	sells consumer vehicles with ADV functionality
yes (the route is known)	geofencing	no (should drive on all public roads)
full ODD (automated driving in the entire route)	ODD operational design domain "Under which condition can the ADV take responsibility?"	partial ODD (automated driving only in prescribed situations, e.g. highway)

Two Different Approaches, with Different Business Models



Fixed-route bus, taxi, delivery



Consumer ADV

Either way, to be responsible, we need to know driving scenarios in advance

→ We derive and verify RSS rules for those driving scenarios, and mathematically guarantee safety

remote

human intervention

on-site (human driver)

offers fixed-route mobility and delivery service

business model

sells consumer vehicles with ADV functionality

yes
(the route is known)

geofencing

no
(should drive on all public roads)

full ODD

ODD

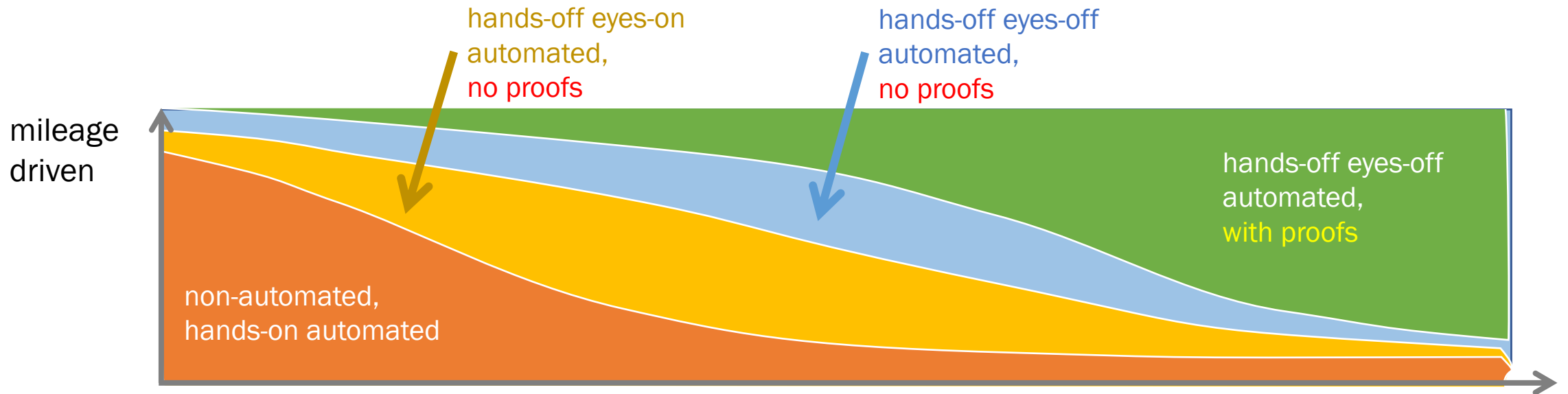
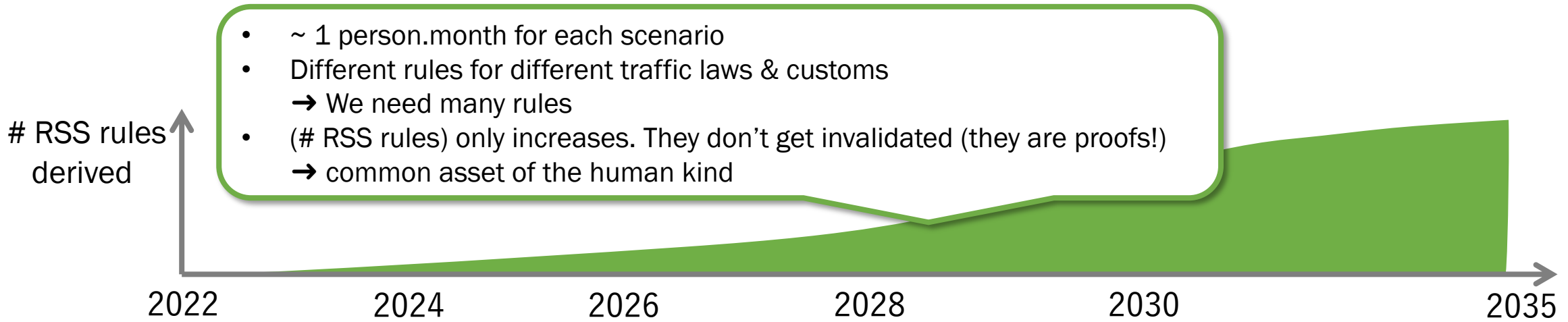
partial ODD

(automated driving in the entire route)

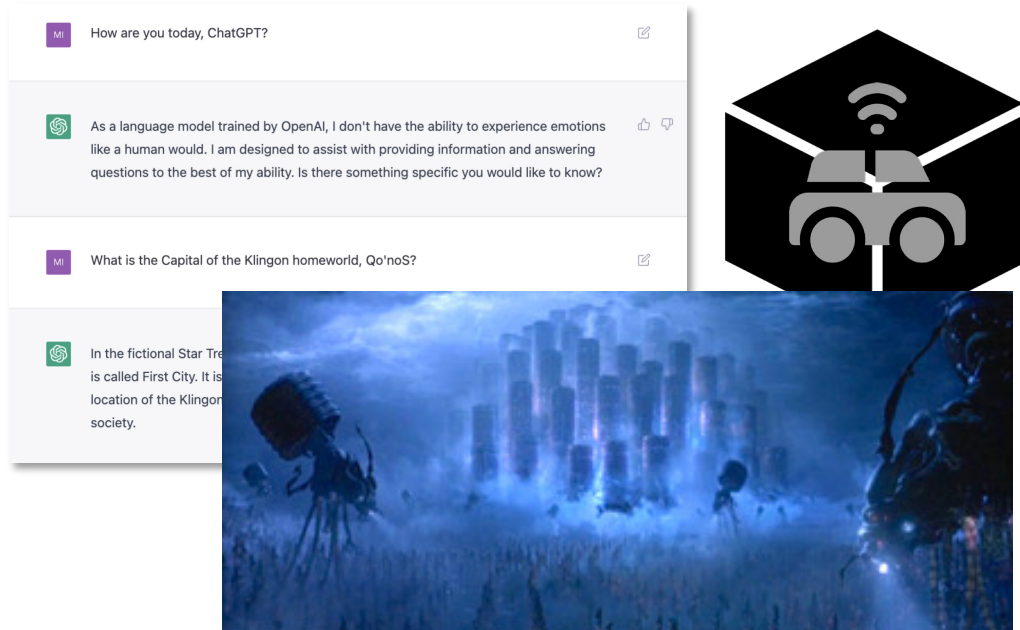
operational design domain
"Under which condition can the ADV take responsibility?"

(automated driving only in prescribed situations, e.g. highway)

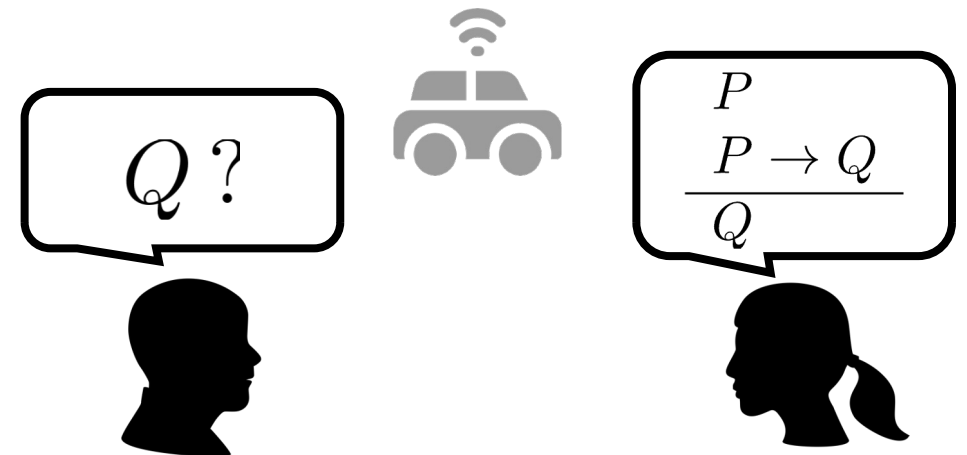
Incremental Accumulation of RSS Rules, Incremental ODD Expansion of “ADVs with Proofs”



Safety-Critical Systems Should Never be Blackbox Proofs Explicate Assumptions, Contracts, ODDs, and Responsibilities



- Many emerging technologies are statistical and blackbox
- We shouldn't let them operate in safety-critical domains
- (... fight against the “lawyer up” approach towards safety!)

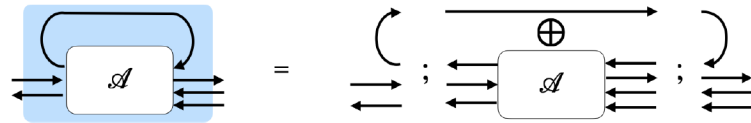


- Conventionally:
Proofs are for establishing absolute truths
- New: proofs are **communication media** for
 - explicating assumptions and contracts,
 - showing who's responsible for what, and
 - writing and assessing safety cases
- Logiic as a social infrastructure for trust in ICT

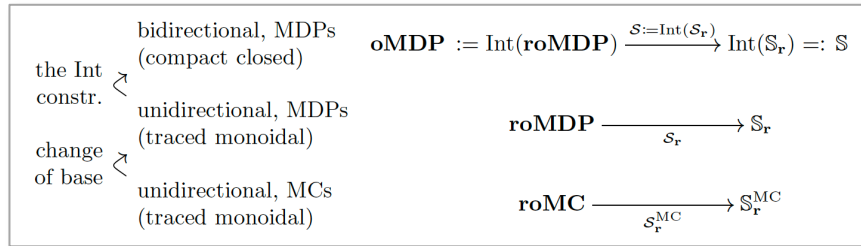
Compositional MDP model checking by string diagrams

[Watanabe, Eberhart, Asada & Hasuo, CAV'23]

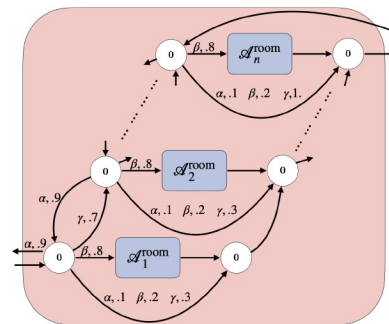
- MDP model checking can be **compositional** over **string diagrams** of MDPs



- Algorithm derived from the structural theory of **monoidal categories**



- ... which can be **arbitrary faster** than existing (non-compositional) algorithms



From mathematical abstraction to programming abstraction

[Kori, Urabe, Katsumata, Suenaga & Hasuo, CAV'22]

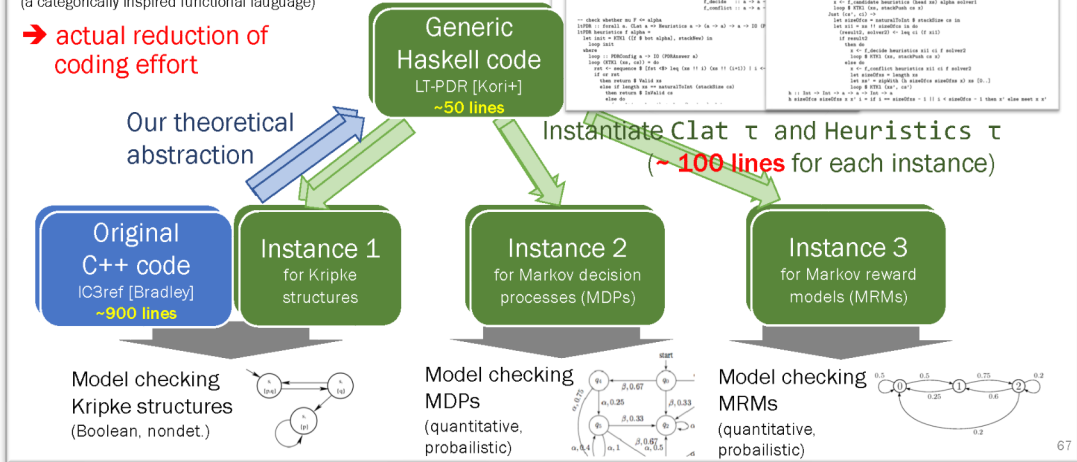
[Kori, Ascari, Bonchi, Bruni, Gori & Hasuo, CAV'23]

Programming Abstraction in LT-PDR

Mayuko Kori, Natsuki Urabe, Shin-ya Katsumata, Kohhei Suenaga, Ichiro Hasuo, The Lattice-Theoretic Essence of Property Directed Reachability Analysis. Proc. CAV 2022

Exploiting the power of Haskell (a categorically inspired functional language)

→ actual reduction of coding effort



- We can **literally code the abstract theory** thanks to Haskell
- Appl. to IC3/PDR (Bradley, Een, ...): 50 LOC (general) + ~100 LOC each (instant.)
 - vs. original IC3 impl., ~900 LOC in C++