

This is a repository copy of *RiskStructures:A Design Algebra for Risk-Aware Machines*.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/172137/>

Version: Published Version

---

**Article:**

Gleirscher, Mario [orcid.org/0000-0002-9445-6863](https://orcid.org/0000-0002-9445-6863), Calinescu, Radu [orcid.org/0000-0002-2678-9260](https://orcid.org/0000-0002-2678-9260) and Woodcock, Jim [orcid.org/0000-0001-7955-2702](https://orcid.org/0000-0001-7955-2702) (2021) *RiskStructures:A Design Algebra for Risk-Aware Machines*. *Formal Aspects of Computing*. pp. 763-802. ISSN 1433-299X

<https://doi.org/10.1007/s00165-021-00545-4>

---

**Reuse**

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.



# RISKSTRUCTURES: A design algebra for risk-aware machines

Mario Gleirscher <sup>1,2</sup>, Radu Calinescu <sup>2,3</sup> and Jim Woodcock 

Computer Science Department, University of York, York, UK

**Abstract.** Machines, such as mobile robots and delivery drones, incorporate controllers responsible for a task *while handling risk* (e.g. anticipating and mitigating hazards; preventing and alleviating accidents). We refer to machines with this capability as risk-aware machines. *Risk awareness* includes robustness and resilience and complicates monitoring (i.e., introspection, sensing, prediction), decision making, and control. From an engineering perspective, risk awareness adds a range of dependability requirements to *system assurance*. Such assurance mandates a correct-by-construction approach to controller design, based on mathematical theory. We introduce RISKSTRUCTURES, an algebraic framework for risk modelling intended to support the design of *safety controllers* for risk-aware machines. Using the concept of a *risk factor* as a modelling primitive, this framework provides facilities to construct, examine, and assure these controllers. We prove desirable algebraic properties of these facilities, and demonstrate their applicability by using them to specify key aspects of safety controllers for risk-aware automated driving and collaborative robots.

**Keywords:** Correct construction; Formal development; Risk awareness; Run-time mitigation; Safety controllers; Robots and autonomous systems

## 1. Introduction

Humans identify and predict dangers and take actions to avoid or recover from dangerous situations. Likewise, autonomous machines can be expected to handle *risk* (e.g. hazardous environmental situations, erroneous control inputs, machine faults) on their own. Such machines should continuously judge risk in a situation, using introspection, estimating the current situation, and predicting future situations. While some risks can be avoided, for most of them only the *likelihood of their occurrence* or the *likelihood and severity of their consequences* can be reduced. This circumstance requires a careful analysis of risk causes and consequences in terms of *risk factors*, and the investigation of temporal and causal relationships between these. Risk can be handled by various measures [HM99, AASB<sup>+</sup>06], ranging from passive measures to *active control schemes*. The former include, for example, flexible surfaces or soft materials used on robot arms and predetermined breaking points in machine chassis [HASH09]. The latter include low-level continuous control mechanisms [SC88], high-level discrete-event

---

Correspondence to: Mario Gleirscher, e-mail: mario.gleirscher@york.ac.uk

<sup>1</sup>Mario Gleirscher was supported in part by the German Research Foundation (DFG) under the Fellowship Grant no. 381212925.

<sup>2</sup>Work by Radu Calinescu and Mario Gleirscher was partially supported by the Lloyd's Register Foundation under the Autonomy Assurance International Programme (AAIP) Grant CSI:Cobot.

<sup>3</sup>Radu Calinescu was additionally supported by the UKRI Project EP/V026747/1 "Trustworthy Autonomous Systems Node in Resilience".

controllers [MGW<sup>+</sup>18, Gle17], and remote or supervisory control schemes [She03]. Overall, active control schemes detect violations of safety properties and try to reestablish a safe state.

**Objective and hypothesis** Decades ago, it was suggested that machine autonomy requires extensive sensor fusion, qualitative reasoning (e.g. [HM99]), and the internalisation of models of consciousness [HG03]. However, to this date, domains such as human-robot collaboration and autonomous driving are far from full autonomy in unrestricted environments, to a good extent because of difficulties in *assurance* [KW17, AZI<sup>+</sup>17]. As similar engineering challenges have been solved through formal design techniques (see, e.g. the survey papers [GM20, GFW19]), we hypothesise that these difficulties can be reduced by *bespoke rigorous design techniques*<sup>4</sup> beyond what can be achieved by employing general-purpose formal design methods alone. Such techniques could guide the correct construction of *safety controllers* that, inspired by [HM99, HG03], integrate a simple form of consciousness of risk, hereafter called *risk awareness*, into a machine.

**Contributions** Based on this hypothesis, we introduce RISKSTRUCTURES, an algebraic framework for designing correct-by-construction safety controllers for risk-aware machines. In line with conventional risk analysis, this framework uses the notion of a *risk factor* as a modelling primitive and, going beyond preliminary results from our previous work [Gle17, Gle18], enhances its formalisation and proves key algebraic properties of a factor-based design calculus for risk-aware controllers. Our work focuses on qualitative risk modelling, enabling the use of different approaches to dealing with uncertainty. A risk structure can be used to formally specify *acceptance criteria*<sup>5</sup> and from these, one can synthesise safety controllers [GC20] and reason about their correctness. The presented theory is supported by the research tool YAP [Gle20], which interprets risk structures encoded in a domain-specific language and automates analysis and synthesis steps. To the best of our knowledge, this work is the first to give an algebraic account of risk modelling while building a bridge to safety controller design for risk-aware machines, adding design guidance to state-of-the-art dependability and risk assessment techniques, monitoring approaches, and optimal control models (Sect. 8).

The remainder of this article is structured as follows. After summarising background material and the used formalism (Sect. 2), Sects. 3 to 7 present our framework, using examples (Sect. 3 and sects. 4.3, 6.2 and 7.4) and discussing the rationale for its components (Sects. 5.5 and 6.3). Sect. 7.6 summarises RISKSTRUCTURES as an approach to model uncertain negative outcomes of a machine’s actions. Based on [BFK13], Table 5 highlights how low-level uncertainties can be dealt with in RISKSTRUCTURES. After a comparison to previous research in Sect. 8, we add concluding remarks in Sect. 9. Proof details are listed in Appendix A.

## 2. Preliminaries

Beyond the avoidance and removal of development mistakes [ALRL04, Sec. 3.3.2] by rigorous design [GFW19] and dependability techniques [ALRL04, Sec. 5.3.1], there are *immanent operational risks* to be dealt with by bespoke machine features. This section introduces terminology, gives an overview of risk analysis and handling for robots and autonomous systems, and explains the formalism used in this work.

### 2.1. Risk analysis, assessment, and handling

*Risk* is the possibility of *undesired outcomes* (e.g. loss, injury, damage) of an action (whether nominal, faulty, or malicious) with *uncertain alternative outcomes* [KG81].<sup>6</sup> Particularly relevant are the actions (e.g. failures, incidents, accidents) in hazardous states (e.g. faults). Risk can be assessed by the *likelihood of a hazard*, the *likelihood of exposure* of an asset to this hazard, the *likelihood of an undesired outcome* if this hazard occurs, and the *severity* of this outcome [Lev95, Kum07]. As such, it is useful to identify causal relationships between events and to quantify risk using stochastic models (cf. reliability analysis of repairable systems [Kum07]).

<sup>4</sup>That is, techniques dedicated to capturing risk factors.

<sup>5</sup>For readers familiar with refinement-based development: a risk structure  $\mathfrak{R}$  refined by a process  $P$ , after hiding  $P$ ’s irrelevant events, can be used to express risk acceptance in  $P$  and verify this criterion of  $P$ .

<sup>6</sup>For the sake of simplicity, we neglect for whom the outcome is undesired and who performs the action. However, such aspects are useful and can be added.

Knowledge about risk often stems from accidents [SR02]. Methods such as *hazard and operability studies* (HazOp), event tree analysis, and layer of protection analysis help engineers to disclose hazardous event chains and design countermeasures (e.g. interventions) in response to and for the prevention of accidents [Lev95]. *Fault tree analysis* (FTA) and *failure mode and effects analysis* (FMEA) specialise in the identification, assessment, and reduction of faults or failures. There are many variants and combinations of these techniques, enriched with intuitive models and visual languages (e.g. AcciMaps [SR02]; STAMP [Lev04] for systemtheoretic process analysis [Lev12]; UML for HazOp [GMGP10]; CORAS diagrams [LSS11]), tailored for specific stages in the system life cycle [Eri15], and used to shape assurance arguments [McD94, GC17].

Using FTA, one can produce a *fault tree* (FT), a versatile causal model of fault propagation in a system, relating an *undesired event*  $e_u$  with a set  $B$  of *basic events* connected by various *gates* (e.g. AND, OR, NOT). A *minimum cut set* (MCS) is a minimal set of events causing  $e_u$ . *Dynamic FTs* are an extension modelling the ordering of such events, leading to *minimal cut sequences*. Such sequences describe minimal sets of events to occur in a particular order to activate  $e_u$ . FTs can be expressed by connecting sets  $minCS \subseteq 2^B$  in the disjunctive normal form in the antecedent and with  $e_u$  in the consequent of the implication  $(\bigvee_{S \in minCS} \bigwedge_{e_b \in S} e_b) \Rightarrow e_u$ , where  $minCS$  denotes all MCSs for  $e_u$ , and  $e_b$  stands for a *basic event*.

For *quantitative analysis*, FTA and FMEA have been further formalised using probabilistic models [Kum07, ORS06], particularly, input/output Markov chains [BCS07] for checking probabilistic temporal properties, or Markov automata [VJK16] for efficient calculation of failure rates. FTs can be synthesised from failure annotations of stochastic Petri nets [UPM18] or from counterexamples generated by continuous stochastic logic model checking of continuous time Markov chains [LFL13].

**Risk handling.** Many operational risks cannot be avoided or mitigated by passive measures (e.g. physical segregation, safe materials). Mechanical or electronic *equipment failures* are usually addressed by fault-tolerant, robust, or diverse designs [Bir17], *human errors* by ergonomic interfaces, controllers robust to unsafe inputs [Lev12], and *adverse environmental events* by bespoke safety functions (e.g. safety modes [AZI+17]). An elegant transition from formal FTs into the specification of such functions is described in [HRS98].

*Run-time verification* (RV) checks event traces recorded during operation for violation (safety) or acceptance (co-safety) of a property [LS09]. RV can range from simple watchdogs, the checking of constraints on complex algorithm outputs [BBH+17], to Bayesian prediction to warn about violated health parameters such as low battery time [IMM18] and violated probability thresholds for rear-end collisions [CLC+19]. Huang et al. [HEZ+14] show how *monitoring-oriented programming* (MOP) [MJG+11] allows the checking of safety properties in message communications of the *robot operating system*<sup>7</sup> (ROS). While these works focus on the synthesis of *property monitors*, there are also approaches to the design of *property enforcers* (e.g. [GPBB08]). Already in the 1980s, Sobek and Chatila [SC88] proposed the interruption of robot plan execution by *safety monitors* (e.g. for obstacle detection) responding with *corrective actions* (e.g. obstacle avoidance) and *mitigation monitors* resuming the planner after action success. Simmons [Sim94] speaks of deliberative components for normal situations and reactive behaviours for exceptional situations. Sorin et al. [SLJS16] describe the generation of safety monitors for ROS-based autonomous robots using rules with corrective actions.

*Risk-aware control* offers a versatile approach to property enforcement, for example, to minimise collision risk of *autonomous vehicles* (AVs). Althoff et al. [AKWB11] work with a probabilistic variant of inevitable collision states [FA04] to approximate collision probability and cost beyond the planning horizon by Monte Carlo sampling of trajectories. These metrics allow the ranking of simulated trajectories in navigation decisions. Pereira et al. [PBHS13] experiment with a grid-based minimal collision risk planner and a risk-aware *Markov decision process* (MDP) for navigating an underwater vehicle. Sanger [San14] formalises risk-aware movement by a risk estimate based on state uncertainty and the cost of control errors, implemented in a neuronal network-based AV controller. Feyzabadi and Carpin [FC14] propose a risk-aware planning algorithm using constrained MDPs, illustrated for autonomous indoor navigation. Müller and Sukhatme [MS14] encode collision risk by a Gamma distribution of the state and uncertain distance to a nearest obstacle. Shalev-Shwartz et al. [SSSS18] investigate collision-free navigation based on driving rules following the Duty of Care approach from Tort law [Har00]. This approach assumes proper response of other traffic participants in typical driving scenarios. The vehicle action space is discretised to simplify the control problem. Inspired by the notion and versatility of FTs, the present work seeks to identify the commonalities of works in risk-aware control and property enforcement and provide a unified risk handling framework.

<sup>7</sup>See <https://www.ros.org>

## 2.2. Formalism and notation

Our framework uses the *communicating sequential processes* (CSP, [Hoa85, Ros10]) approach to algebraic specification and *labelled transition system* (LTSs, [BK08]) to reason about its operational semantics. Let  $\Sigma$  be the set of all *concrete events* of a process, called its *alphabet*. We distinguish two special events:  $\checkmark$  for the *termination* of a process and  $\tau$  for the *invisible event*, resulting from *hiding* observations. We require  $\checkmark, \tau \notin \Sigma$  and define  $\Sigma^\checkmark = \Sigma \cup \{\checkmark\}$ ,  $\Sigma^\tau = \Sigma \cup \{\tau\}$ , and  $\Sigma^{\checkmark, \tau} = \Sigma^\checkmark \cup \Sigma^\tau$ .

**Definition 1 (Process)** A *process* is an expression of the form

$$P ::= a \rightarrow P \mid ?x : A \rightarrow P \mid P \sqcap P \mid P \square P \mid P \parallel_A P \mid P ; P \mid P \setminus A \mid SKIP \mid STOP \quad (1)$$

with event  $a \in \Sigma$ , shared alphabet  $A \subseteq \Sigma$ , event prefix ( $\rightarrow$ ), prefix choice ( $?x : \rightarrow$ ), non-deterministic choice ( $\sqcap$ ), external choice ( $\square$ ), generalised parallel composition ( $\parallel$ ), sequential composition ( $;$ ), hiding ( $\setminus$ ), termination (*SKIP*), and deadlock (*STOP*). Let  $\mathcal{P}$  be the set of all *processes*.

Each expression  $P \in \mathcal{P}$  corresponds to observable behaviours represented in a *trace model*. Traces are finite sequences over  $\Sigma^\checkmark$  abstracting from details (e.g. internal states) of  $P$ 's executions. In the trace model,  $traces(P)$  yields the *trace semantics* of  $P$ . For example,  $traces(STOP) = \{\langle \rangle\}$  and  $traces(SKIP) = \{\langle \rangle, \langle \checkmark \rangle\}$ .  $initials(P) = \{a \mid \langle a \rangle \in traces(P)\}$  denotes the set of all *initial events* of  $P$ . For processes  $P$  and  $Q$ , we write  $P \sqsubseteq_s Q$  to say that  $Q$  *s-refines*  $P$  (or  $P$  is *s-refined* by  $Q$ ) and that  $traces(P) \supseteq traces(Q)$ , with the usual abbreviation  $P =_s Q \iff P \sqsubseteq_s Q \wedge Q \sqsubseteq_s P$ .<sup>8</sup> We use the term *control state* for a process and write  $\parallel$  if  $A = \Sigma$ . Below, we consider  $P = Sy \parallel En$ , with a system  $Sy$  interacting with its environment  $En$ , or  $P = En(Sy)$ , with the CSP expression  $En$  using  $Sy$ .

**Definition 2 (Labelled transition system)** A LTS is a tuple  $T = (S, \Sigma^T, \rightarrow, S_0)$  with a set  $S$  of *states*, a set  $\Sigma^T \subseteq 2^{\Sigma^\tau} \setminus \{\emptyset\}$  of *abstract events* (i.e., sets of concrete events), a relation  $\rightarrow \subseteq S \times \Sigma^T \times S$  of state transitions when engaging in these events, and a set  $S_0 \subseteq S$  of *initial states*. Omission of  $S_0$  results in an *uninitialised* LTS  $(S, \Sigma^T, \rightarrow)$ , further omission of  $\Sigma^T$  in a directed graph  $(S, \rightarrow)$  with  $\rightarrow \subseteq S \times S$ , called *transition system* (TS). Let  $\mathcal{T}$  be the set of all LTSs.

We write  $s \xrightarrow{e} s'$  if  $(s, e, s') \in \rightarrow$ ,  $s \xrightarrow{e} s'$  if  $\exists s' \in S : (s, e, s') \in \rightarrow$ ,  $s \rightarrow$  if  $\exists e \in \Sigma^T, s' \in S : (s, e, s') \in \rightarrow$ ,  $s \xrightarrow{e}$  if  $\exists s' \in S : (s, e, s') \in \rightarrow$ , and  $s \rightarrow$  if  $\exists e \in \Sigma^T, s' \in S : (s, e, s') \in \rightarrow$ . A  $\square$  indicates discharged cases in proofs.

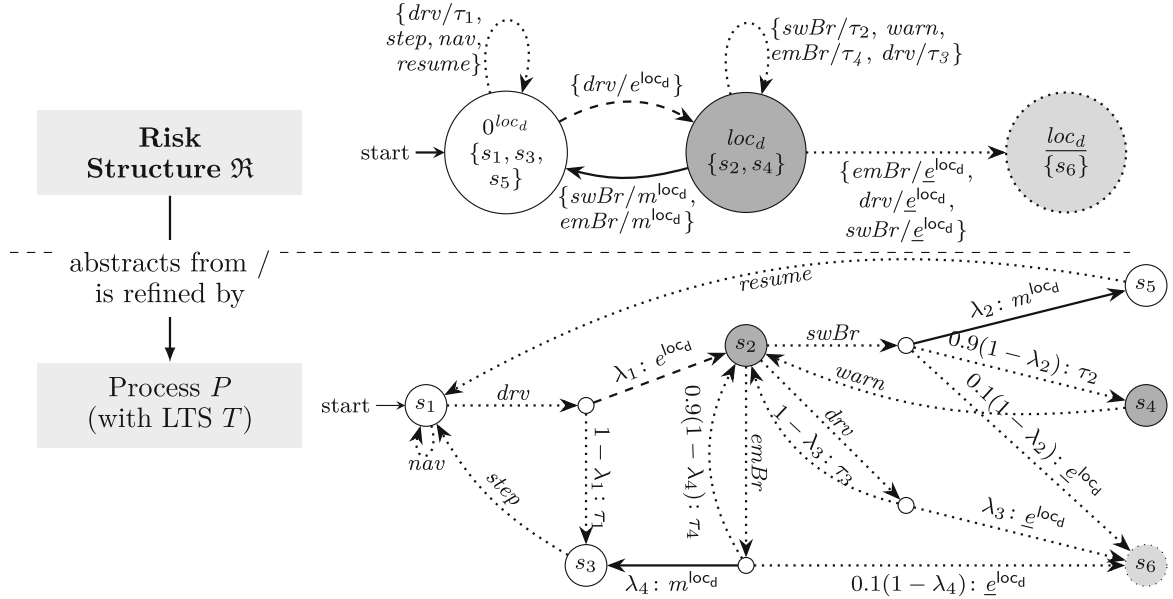
## 3. RISKSTRUCTURES: overview

In this section, we give an overview of the components of the proposed framework. We describe the abstraction facilitated by RISKSTRUCTURES, building a bridge between (probabilistic) process modelling, risk analysis, and controller design. Moreover, we illustrate the main concepts, terminology, and work steps using an example from the road vehicle domain, which will be revisited in later sections.

Figure 1 (bottom right) shows an example of a process  $P$  (Definition 1) associated with an LTS  $T$  (Definition 2).  $P$  describes actions of one or several agents (e.g. a machine, an operator) in a physical environment, for example, a manually driven vehicle with some kind of driving assistance on an urban road.  $P$ 's actions include data processing, control stimuli via actuators, and the monitoring of process outcomes. In this example, data processing occurs with *nav*, where the machine performs calculations for navigation (stored internally and displayed to the driver), *step* of the driver interpreting data for a next control action, *warn* of the vehicle sending information to the driver, and *resume* of the vehicle returning from a risk handling procedure. Control stimuli include the action *drv*, where the driver performs a default driving manoeuvre, *swBr* where the vehicle executes a swerve and brake manoeuvre, and *emBr* representing an emergency braking action conducted by the vehicle. Monitored outcomes include, for example, the occurrence ( $e^{loc_d}$ ) of the risk factor *loss of driving control* ( $loc_d$ ), a potential *accident* ( $e^{loc_d}$ ) following  $e^{loc_d}$ , and the performance and success of a mitigation ( $m^{loc_d}$ ) of this factor, recognised by the machine, and idling or no output ( $\tau_i$ ). Underspecification and uncertainty in  $P$  can be described by non-determinism or quantified by probabilities ( $\lambda_i$ ) on action outcomes.  $traces(P)$  denotes possible chains of events and  $T$  the reachable states ( $s_i$ ).<sup>9</sup>

<sup>8</sup>We use  $\sqsubseteq_T$  for *trace* refinement and  $\sqsubseteq_{FD}$  for *failures-divergences* refinement and refer the inclined reader to [Ros10].

<sup>9</sup> $traces(P)$  contains finite traces reaching such states such that causes can be represented by well-founded sets.



**Fig. 1.** Risk structure  $\mathfrak{R}$  partitioning the states of the process  $P$ , that is, labelling  $P$ 's states with the phases *inactive* ( $0^{loc}$ , white) and *active* ( $loc$ , dark gray) of the risk factor  $loc_d$ ;  $P$ 's transitions are classified into endangerments (dashed arcs), mitigations (solid arcs), and other transitions (dotted arcs); incomplete abstractions may contain unclassified states (light-grey)

Risk, as introduced in Sect. 2.1, can then be qualitatively (or quantitatively) characterised in  $P$  as the possibility (probability) of an undesired state (e.g.  $s_4$ ) being finitely reachable from a particular state (e.g.  $s_1$ ) in  $T$ . Next, Fig. 1 shows the abstraction from  $P$  into an exemplary *risk structure*  $\mathfrak{R}$  comprising the single factor  $loc_d$ . This abstraction allows one to focus on the events relevant from the perspective of this factor, that is, its activation ( $e^{loc_d}$ ) and mitigation ( $m^{loc_d}$ ). Accordingly, this abstraction step includes the labelling of regions of  $P$ 's state space with predicates (e.g. for hazards and, more generally, situations) and the construction of  $\mathfrak{R}$  over the resulting *abstract state space*. Moreover, this representation helps one to focus on the structure of risk of a certain application by using information about factors encoded in predicates (e.g. active, inactive) and real-valued functions (e.g. [AKWB11, PBHS13]). These predicates and functions can be derived from hazard analysis (Sect. 2.1) and be used to allow a risk-aware machine to estimate risk levels of reachable states and to plan and execute actions reaching or avoiding the corresponding states.

For example, the factor  $loc_d$  partitions  $T$ 's state space (cf. nodes in the bottom right of Fig. 1) into regions where  $loc_d$  is *inactive* ( $0^{loc_d}$ , white), *active* ( $loc_d$ , dark grey), or where an *accident* ( $loc_d$ , light-grey) has happened. To focus on the practical handling of  $loc_d$  in  $P$ , outcomes (e.g.  $e^{loc_d}, m^{loc_d}, \tau_i$ ) are paired with their control actions (e.g.  $drv, swBr, emBr$ ; dotted arcs) and merged into sets of transition labels called events (e.g.  $\{swBr/m^{loc_d}, \dots\}, \{swBr/\tau_2, \dots\}$ ). To keep track of mishap states, such as  $s_6$ , we assign the label  $loc_d$  to states where undesired consequences after an unmitigated  $e^{loc_d}$ -event will be materialised. Such states might be reached, for example, after the execution of  $emBr, drv$ , and  $swBr$  with a *collision*  $e^{loc_d}$  as the possible outcome.  $loc_d$  enables risk assessment in terms of likelihood and severity. To design mitigations for  $loc_d$ , we introduce an extra risk factor (e.g.  $col$ ) handling  $e^{loc_d}$  in analogy to  $e^{loc_d}$ . This technique will be discussed below. Overall, for  $loc_d$ ,  $\mathfrak{R}$  reduces to just two abstract states ( $0^{loc_d}, loc_d$ ) and two events, the endangerment  $e^{loc_d}$  (dashed arc) and the mitigation  $m^{loc_d}$  (solid arcs). We assume there to be observational *refinement* or some simulation relation between  $T$  and  $\mathfrak{R}$ .

Overall, risk awareness amounts to the machine incorporating  $\mathfrak{R}$  in its decisions such that expected risk from its own actions does not exceed an acceptable level. *So, how can we engineer  $\mathfrak{R}$ ?*

From an engineering perspective, our approach requires neither  $P$  nor  $\mathfrak{R}$  to exist in a complete form to facilitate the aforementioned abstraction or refinement. One can rather think of these two artefacts as being developed concurrently.

**Table 1.** A generic situation/action/factor table (a) and per-cell risk analysis checklists (b, c)

	<b>Action</b>	$a_1$	$a_2$	...		
<b>Situation</b>	<b>Factor</b>	$f_{1,1}$	$f_{1,2} \dots$	$f_{2,1}$	$f_{2,2} \dots$	...

(a) *Generic situation/action/factor table*

$\sigma_1$	Each cell filled with results from a cause/consequence and mitigation analysis along the lines of Tables 1b and 1c.	...
$\sigma_2$		...
$\vdots$	$\vdots$	$\vdots$

(b) *Factor activation analysis (causes/consequences)*

Why would  $f_i$  occur with/during action  $a_j$  in situation  $\sigma_k$ ?  
What if  $f_i$  occurs with/during action  $a_j$  in situation  $\sigma_k$ ?

---

1. What is the likelihood of occurrence of  $f_i$ ?  
Quantified by  $Pr[f_i | a_j, \sigma_k]$ .
2. What is the risk of an accident  $\underline{f}_i$  following  $f_i$ ?  
Quantified by the likelihood  $Pr[\underline{f}_i | f_i, a_j, \sigma_k]$  and the severity  $sev(\underline{f}_i; f_i, a_j, \sigma_k)$ .
3. Which other factors are related to the activation of  $f_i$ ?  
Qualified by relationships such as  $f_i$  **requires**  $f_j$ ,  $f_i$  **causes**  $f_k$ ; and quantified by  $Pr[f_k | f_i, a_j, \sigma_k]$ .

---

(c) *Factor mitigation analysis*

How can we mitigate  $f_i$  during or after action  $a_j$  in the situation  $\sigma_k$ ?

---

4. Which mitigation options  $m$  are available to reach  $\overline{f}_i$ ?
5. What is the risk of an accident  $\underline{f}_i$  following  $\overline{f}_i$ ?  
Quantified by the likelihood  $Pr[\underline{f}_i | \overline{f}_i, m, a_j, \sigma_k]$  and the severity  $sev(\underline{f}_i; \overline{f}_i, m, a_j, \sigma_k)$ .
6. Which other factors are related to the mitigation of  $f_i$ ?  
Qualified by relationships such as  $f_i$  **prevents**  $f_j$ ,  $f_i$  **causes**  $f_k$ ; and quantified by  $Pr[f_k | \overline{f}_i, m, a_j, \sigma_k]$ .

For example, one might start with a part of  $P$ , then, after risk analysis, continue with a part of  $\mathfrak{R}$ , and in a second iteration, maybe even after a period of operation, proceed with another part of both  $P$  and  $\mathfrak{R}$ . In this sense,  $\mathfrak{R}$  can be incrementally developed. The notion of a situation will help us to structure such a development. A *situation*,  $\sigma$ , is an abstract state that describes a particular context (i.e., a scene or environment), an activity (i.e., a task, operating mode, a use case, or scenario), and a configuration of hazardous events occurring in  $P$ . Engineers and risk analysts are familiar with tables and graphs, for example, when applying HazOp or FMEA to assess risk factors of certain machine actions in particular situations. We will, therefore, show in the following how one can develop  $\mathfrak{R}$  by means of what we call a *situation/action/factor* table (Table 1a).

One may want to start with an initial set of application-specific *situations* and *actions* that the machine under consideration can perform, together with an initial set of *factors*, maybe obvious from accident experience and preliminary hazard analyses. Then, each resulting cell in Table 1a would involve a focused cause/consequence analysis (e.g. [ORS06]) and risk assessment (Sect. 2.1), for example, using risk graphs [Kum07, p. 53]. Tables 1b and 1c elaborate ideas in [Gle18] for such analyses, including the estimation of the conditional probability  $Pr[X | Y]$  and severity  $sev(X; Y)$  of an event  $X$  under the condition  $Y$ . We use  $0^f$ ,  $f$ ,  $\overline{f}$ , and  $\underline{f}$  to denote the *inactive*, *active*, *mitigated*, and *mishap* phases of a factor  $f$ , explained in Sect. 4.1. Incrementally, one would add further situations, actions, and factors to the table. At some point, one can split the table into several smaller ones. Please, note that such tables are not subject of this work. We presented them merely to hint to the data to be collected for building  $\mathfrak{R}$ . Moreover, when formalising situations, we will only consider the hazardous events part. In our future work, we will discuss situations in more detail.

**Table 2.** Example of a situation/action/factor table for road driving

Situation	Action Factor Descr.	drive ( <i>drv</i> )			releaseAirbag ( <i>rA</i> )	
		<i>loc<sub>d</sub></i>	<i>ncol</i>	<i>col</i>	<i>st<sub>rA</sub></i>	<i>fod<sub>rA</sub></i>
		loss of control	near-collision	collision	spurious trip	fail. on dem.
manual mode		3. <i>Pr</i> [ <i>ncol</i>   <i>loc<sub>d</sub></i> ]: h, requires <i>st<sub>rA</sub></i>	3. <i>Pr</i> [ <i>col</i>   <i>ncol</i> ]: h 4. <i>swerve</i> ; <i>brake</i>	2. <i>sev</i> ( <i>col</i> ): m 3. requires <i>ncol</i>	2. <i>sev</i> ( <i>shock</i> ): m	–
autonomous mode	–		3. <i>Pr</i> [ <i>col</i>   <i>ncol</i> ]: h 4. <i>swerve</i> ; <i>brake</i>	2. <i>sev</i> ( <i>col</i> ): m 3. requires <i>ncol</i>	3. causes <i>loc<sub>d</sub></i> 2. <i>sev</i> ( <i>shock</i> ): m	–
following vehicle		3. equal to manual/autonomous			2. <i>sev</i> ( <i>col</i> ): l	3. equal to manual/autonomous
leading vehicle		3. equal to manual/autonomous			2. <i>sev</i> ( <i>col</i> ): h, 4. <i>rA</i> , 5. <i>sev</i> ( <i>col</i> ): m	3. equal to manual/autonomous
oncoming vehicle		3. equal to manual/autonomous			2. <i>sev</i> ( <i>col</i> ): h, 4. <i>rA</i> , 5. <i>sev</i> ( <i>col</i> ): h	3. equal to manual/autonomous
collision ( <i>col</i> )	–	–	–	2. <i>sev</i> ( <i>col</i> ): m	–	3. equal to *
no airbag ( <i>fod<sub>rA</sub></i> )		1. <i>Pr</i> [ <i>loc<sub>d</sub></i> ]: l 3. eq. to manual	3. eq. to manual	* 2. <i>sev</i> ( <i>col</i> ): h	–	3. idempotent

*Legend:* Numbers indicate the analysis steps applied from the Table 1b and 1c; situation/action/factor parameters for *Pr* and *sev* are determined by the cell (hence omitted); values after the colon are exemplary; probabilities are abstracted to *h/m/l* for high/medium/low; relationship sources are determined by the column (hence omitted); \*...cross-reference to other cell; -irrelevant/not applicable

**Extending the Example from Fig. 1** In Fig. 1, we illustrated basic concepts by discussing a single risk factor (*loc<sub>d</sub>*). Now, we extend our discussion to a more realistic setting by investigating situations, system actions, and further risk factors. In Table 2, we consider actions such as *releaseAirbag* (*rA*), responsible for mitigating a collision *col*, that is,  $m^{col} \equiv rA$ . However, *rA* itself may be associated with two failure modes: *failure on demand* *fod<sub>rA</sub>* (i.e., *rA* not performed when requested) and *spurious trip* *st<sub>rA</sub>* (i.e., *rA* performed when not requested). In forward reasoning, one can ask: What if *fod<sub>rA</sub>* or *st<sub>rA</sub>*, both influencing or determining the behaviour of *rA*, are activated in manual or autonomous driving mode? *st<sub>rA</sub>* in manual mode would increase the probability of a crash because of a *loss of control* *loc<sub>d</sub>* from shock and distraction. *st<sub>rA</sub>* in autonomous mode may cause shock injuries but would unlikely cause a crash. *fod<sub>rA</sub>* will be irrelevant in normal situations but in a collision, risk will be quite similar to the case of no airbag.

Table 2 can store probabilistic relationships, for example, between *near-collision* *ncol* and *collision* *col*. Probabilities for *ncol* and *col* will be highly dynamic and need to be estimated and used for decision making at run-time. Probabilities for *fod<sub>rA</sub>* maybe more stable, known from experiments, and can be used in risk assessment at design-time. Requirements on the confidence in probability estimates are specific to each factor. Note the variety of relationships between situations, actions, and factors, for example, in manual mode *loc<sub>d</sub>* requires *st<sub>rA</sub>*, *fod<sub>rA</sub>* in collision is similar to *col* in no airbag, *st<sub>rA</sub>* in manual mode extends *st<sub>rA</sub>* in autonomous mode, and the airbag, subject of risk assessment in the two *rA*-columns, mitigates *col*.

Graphs (e.g. FTs) can be used to visualise and elicit the information stored in Table 2 if these tables get more complex. The tabular representation is, however, integrative, scales, and conveys the scheme recurring across the cells as well as the relationships between situations, actions, and factors. We will discuss and use some of these relationships in our framework.

The components of a design algebra for risk-aware machines will be described in the following sections: risk factors and risk spaces (Sect. 4), mitigation orders (Sect. 5), factor dependencies (Sect. 6), and, finally, RISKSTRUCTURES (Sect. 7), illustrated by examples and a discussion of applications.

## 4. Risk spaces

This section introduces *risk factors* and *risk spaces*, frequently abbreviated by “factors” and “spaces”, as the basic elements of the algebraic framework.

### 4.1. Risk factors

We first define the notion of a risk factor using a LTS and describe its properties and meaning.



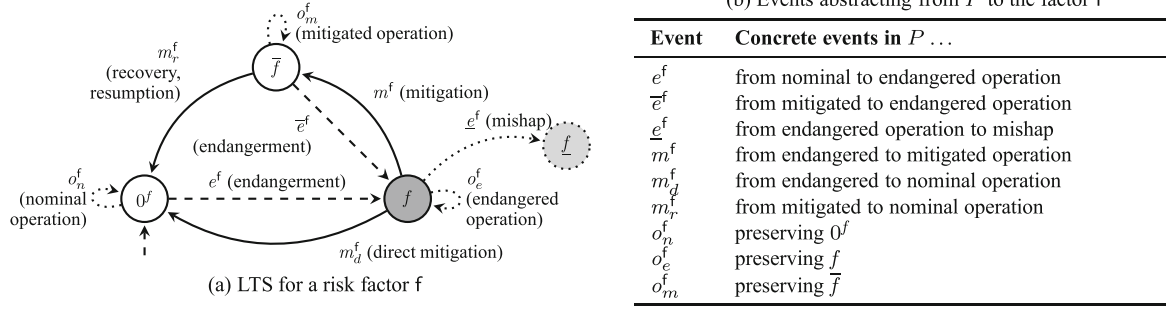


Fig. 2. Basic template of a risk factor  $f$  (a) with a description of  $f$ 's events (b)

**Definition 3 (Risk factor)** Let  $(Ph_f, \Sigma^f, \rightarrow_f)$  be an uninitialised LTS according to Definition 2. Extending this LTS, a *risk factor* is a tuple  $f = (Ph_f, \Sigma^f, \rightarrow_f, \leq_f, s_f)$  with

- A finite set  $Ph_f$  of *phases* of  $f$ ,
- A finite set  $\Sigma^f \subset 2^{\Sigma^f} \setminus \{\emptyset\}$  specifying *significant events* for  $f$ ,
- A labelled transition relation  $\rightarrow \subseteq Ph_f \times \Sigma^f \times Ph_f$ ,
- A partial order  $\leq_f \subseteq Ph_f \times Ph_f$  called *phase order* of  $f$ , and
- A pair  $s_f \in \mathbb{R}_+^2$  with  $s_f^{(1)} \leq s_f^{(2)}$  denoting the *severity of the least and worst expected impact* of  $f$ .<sup>10</sup>

Let  $\mathcal{F}$  be the set of all risk factors.

We consider factors  $f$  with  $\rightarrow_f$  according to Fig. 2a with  $Ph_f = \{0^f, f, \bar{f}\}$  for the phases *inactive* ( $0^f$ , typically, the initial phase), *active* ( $f$ ), and *mitigated* ( $\bar{f}$ ), where  $\leq_f \subseteq_f$  is at least<sup>11</sup> the reflexive transitive closure of  $\{(f, 0^f), (f, \bar{f})\}$ , and with  $\Sigma^f = \{e^f, \bar{e}^f, m^f, m_d^f, m_r^f, o_n^f, o_m^f, o_e^f\}$ .

The interval  $s_f$  represents the *severity* (also known as the detriment or negativity) of  $f$ 's expected materialised impact or consequences for any relevant assets if  $f$  gets and stays active. Impact can include, for example, damage of the environment, injury of humans, loss of a valuable, or any combination of these.

**Events and factor types.** By Definition 3, for any  $p, p' \in Ph_f$  if  $(p, e, p') \in \rightarrow_f$  then  $e \neq \emptyset$ .  $f$  is *deterministic* if and only if for every phase the events of all outgoing transitions are pairwise disjoint. We require  $\forall p \in Ph_f: (\bigcup_{(p, e, p') \in p \rightarrow e} e) \setminus \{\tau\} = \Sigma$ .<sup>12</sup> Figures 2a and 2b indicate the meaning of the events in  $\Sigma^f$  for factor modelling. The three phase-preserving events  $o_n^f, o_e^f$ , and  $o_m^f$  complement the endangerment and mitigation events. Moreover, we distinguish some special types of factors:

- If  $m^f \cup m_d^f = \emptyset$  then we call  $f$  *final*, otherwise *reducible*.
- For any reducible  $f$  with  $m^f \neq \emptyset$ , if  $\bar{e}^f \subset e^f$  then we call  $f$  *strongly reducible*.
- If  $m_d^f = \emptyset$  then we call  $f$  *indirectly reducible*.

**Modelling mishaps.** Observing from the example in Sect. 3, the *mishap* phase  $\underline{f}$ , modelling incidents or accidents from  $f$ , complements the severity information encoded in  $f.s$  and enables risk assessment solely based on a single risk factor. It turns out, however, that  $\underline{f}$  is non-essential for the basic framework we want to discuss in the following sections. While we keep  $\underline{f}$  for factor-specific risk assessment, to simplify algebraic reasoning about factors as discussed below, we collapse  $\underline{f}$  into  $f$  and  $\underline{e}^f$  into  $o_e^f$ . Furthermore, we introduce a final factor  $f'$  to model the corresponding mishap. This approach simplifies the basic factor model and allows the conversion of  $f'$  into a reducible factor when a mitigation for  $f'$  is available. For example, as illustrated in Table 2, the factor  $col$  provides accident handling for the factor  $ncol$ , whose likelihood of occurrence under the condition of  $loc_d$  is high. In [Gle18, GC20], however, we have started to work with a refined factor model including  $\underline{f}$  as well

<sup>10</sup>By usual convention, for an ordered  $n$ -tuple  $t$  and  $i \in [1..n]$ , we write  $t^{(i)}$  to refer to the value of the  $i$ -th element of  $t$ . Furthermore, if  $t$  has a uniquely named element  $e$ , we write  $t.e$  to refer to the value of  $e$  in  $t$ .

<sup>11</sup>Some applications might give rise to a *linear*  $\leq_f \subseteq_f$  by adding at most one out of  $\{(0^f, \bar{f}), (\bar{f}, 0^f)\}$ .

<sup>12</sup>In testing, this is known as input-enabledness [Tre08].

as further phases and events (e.g. for mishap alleviation). In order to preserve the main results presented here, such extensions require  $\preceq_f \preceq_f$  to have a unique maximal element (discussed in Sect. 5.3) and an adjustment of constraint definitions in Sect. 6.1. We conclude that  $Ph_f$  and  $\Sigma^f$  comprise a useful minimal set of elements for a *generic* risk factor  $f$ .

## 4.2. Risk states, spaces, and space composition

Risk factors give rise to *risk states* and *risk spaces*. Let  $F \subset \mathcal{F}$  be a finite factor set (Definition 3) where each  $f \in F$  has the form  $f = (Ph_f, \Sigma^f, \rightarrow_f, \preceq_f \preceq_f, s_f)$ .

**Definition 4 (Risk state)** Assume that risk factors are unique, that is,  $\forall f, g \in F: f \neq g \Rightarrow Ph_f \cap Ph_g = \emptyset$ . Then, a *risk state* is a faithful total injection  $\sigma: F \rightarrow \bigcup_{f \in F} Ph_f$ , that is,  $\forall f \in F: \sigma(f) \in Ph_f$ .

A risk state abstracts from states of  $P$  by focusing on risk-related information in form of state propositions associated with the factor phases. Such propositions will, however, not be formalised in this work.

**Definition 5 (Risk space)** For a factor set  $F$ , a *risk space*  $R(F)$  is the function space given by

$$R(F) = \{ \sigma \in F \rightarrow \bigcup_{f \in F} Ph_f \mid \sigma \text{ is total} \wedge \sigma \text{ is an injection} \wedge \forall f \in F: \sigma(f) \in Ph_f \}.$$

We omit the parameter  $F$  from  $R$  if it is clear from the context and denote the set of all risk spaces by  $\mathcal{R}$ .

By Definition 5,  $R$  is non-empty and finite if and only if  $F$  is non-empty and finite.  $R$  defines the set of all states an arbitrary<sup>13</sup> combination of factor phases might give rise to. Note that  $\sigma(f)$  can now be used to refer to the phase of factor  $f$  in a risk state  $\sigma$ .

**Definition 6 (Compatibility of risk states)** Given  $F_1, F_2 \subset \mathcal{F}$ ,  $\sigma \in R(F_1)$  and  $\sigma' \in R(F_2)$  are *compatible*, written

$$\sigma \approx \sigma' \iff \forall f \in F_1 \cap F_2: \sigma(f) = \sigma'(f). \quad (2)$$

Notice that any two risk states  $\sigma \in R(F_1)$  and  $\sigma' \in R(F_2)$  are compatible if  $F_1 \cap F_2 = \emptyset$  and, furthermore, that state equality implies  $F_1 = F_2$  and, therefore, state compatibility. This compatibility is a prerequisite for risk space composition as follows.

**Definition 7 (Risk space composition)** The *composition*  $\otimes: \mathcal{R} \times \mathcal{R} \rightarrow \mathcal{R}$  of the two risk spaces  $R(F_1)$  and  $R(F_2)$  is defined by

$$R(F_1) \otimes R(F_2) = \{ \sigma_1 \cup \sigma_2 \mid \sigma_1 \in R(F_1) \wedge \sigma_2 \in R(F_2) \wedge \sigma_1 \approx \sigma_2 \}. \quad (3)$$

Now, we can derive a basic law relating the union of risk factors and the composition of risk spaces. Furthermore, it will turn out that  $R$  is a homomorphism.

**Lemma 1 (Exchange of  $\cup$  and  $\otimes$ )**

$$R(F_1 \cup F_2) = R(F_1) \otimes R(F_2)$$

*Proof sketch.* The proof is by mutual existence and uniqueness: For each  $\sigma \in R(F_1 \cup F_2)$ , (i) there exists a  $\sigma_1 \cup \sigma_2 \in R(F_1) \otimes R(F_2)$  and (ii) this pair is unique, and (iii, iv) conversely. Details on the proof can be taken from Appendix A.  $\square$

**Lemma 2 ( $R$  is homomorphic)**  $R$  is a homomorphism in the context of  $(\mathcal{F}, \cup)$  and  $(\mathcal{R}, \otimes)$ .

$$\begin{array}{ccc} F_1, F_2 & \xrightarrow{\cup} & F_1 \cup F_2 \\ \downarrow R & & \downarrow R \\ R(F_1), R(F_2) & \xrightarrow{\otimes} & R(F_1 \cup F_2) \end{array}$$

<sup>13</sup>Below, we also view  $R$  as “the most general (risk) structure” for a specific factor set.

*Proof sketch.* We first make sure that we deal with semi-groups and then show by algebraic manipulation that  $\otimes$  is associative. Details on the proof can be taken from Appendix A.  $\square$

Two specific classes of risk states and a basic lemma close this section.  $R(\{f\}) = \{[f \mapsto 0^f], [f \mapsto f], [f \mapsto \bar{f}]\}$  forms the *trivial risk space* for  $f$ , and  $R(\emptyset) = \emptyset$  the *empty risk space*, used in the following equality.

**Corollary 1** For finite  $F \subset \mathcal{F}$ , Lemma 1 yields  $R(\emptyset)$  to be the *zero* element of composition with  $\otimes$ :

$$R(\emptyset) \otimes R(F) = R(F) . \quad (\otimes\text{-zero})$$

### 4.3. Example: risk factors on the road

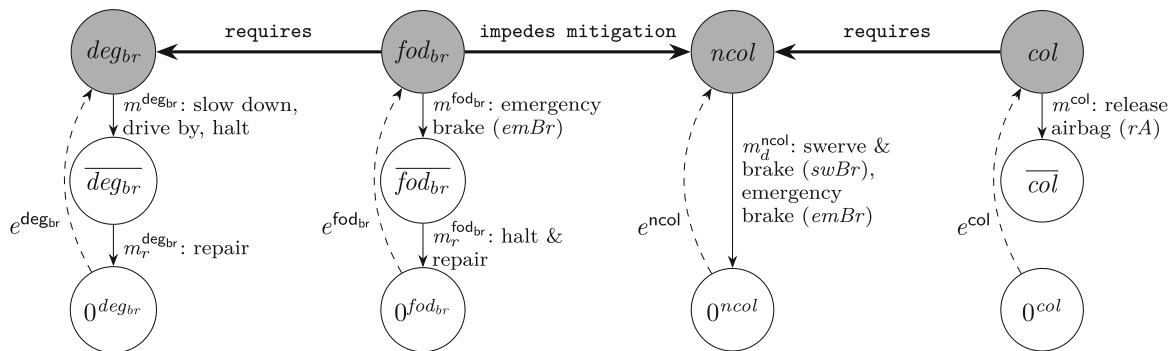
In the following, we discuss Table 2 more deeply and provide further examples for the event sets in Fig. 2b. Tracing back through the causal relation from  $col$ , we find that *collisions* require *near-collisions* to occur beforehand. Based on that observation, we identify a combination of *swerve & brake* ( $swBr$ ) to be a possible mitigation for  $ncol$ . However, the action  $swBr$  gives rise to the *failure on demand*  $fod_{br}$  of the car's braking action  $br$ . This extension of the analysis in Table 2 results in another dependency, namely,  $fod_{br}$  *impedes the mitigation of*  $ncol$ , technically, it disables the braking action  $br$ .

Tracing forward through the causal relation, we may stop at application-specific factors. Such factors can be modelled as *final*, that is, the machine will not mitigate them. This way, final factors define the *scope* of  $\mathfrak{R}$ . However, as explained in Sect. 4.1, a design increment of the machine can turn a final into a *reducible* factor. For example, designing an airbag into a car makes  $col$  reducible. This epistemic limit explains why safety always has to be considered *relative* to a known factor set  $F$  [Gle14].

Tracing back again from  $fod_{br}$  discloses the failure mode *degradation of brake*  $deg_{br}$ .  $deg_{br}$  causally relates to  $fod_{br}$  like  $ncol$  relates to  $col$ :  $deg_{br}$  will make  $br$  ineffectual but not ineffective, if one does not intervene,  $deg_{br}$  will actually cause  $fod_{br}$ . We want  $deg_{br}$  to be *strongly reducible*, based on a mitigation  $m^{deg_{br}}$  that establishes  $\overline{deg_{br}}$  from which only a strict subset  $\bar{e}^{deg_{br}} \subset e^{deg_{br}}$  can occur. In phase  $deg_{br}$ , *drive by and halt* would be a conservative option for  $m^{deg_{br}}$  to be implemented by a safety controller. These reasoning steps are summarised in Fig. 3.

While the state bi-partition for  $loc_d$  shown in Fig. 1 seems sufficient, some factors need a tri- or  $n$ -partition. For example,  $fod_{br}$  suggests a vehicle repair. Therefore, Fig. 2a introduces a third phase to separate the state where a brake failure ( $fod_{br}$ ) is mitigated ( $\overline{fod_{br}}$ ) from the state  $0^{fod_{br}}$  where the brake is fully operational. Because  $fod_{br}$  requires an off-line repair action  $m_r^{fod_{br}}$ , it is *indirectly reducible*. Two further examples for such a factor would be *leaking or damaged battery* and *run out of fuel*. Reaching the “safe state”  $\overline{fod_{br}}$  is possible via an intermediate stable state such as *halted at car repair shop*. From this phase, recovery ( $m_r^{fod_{br}}$ ) to the inactive phase  $0^{fod_{br}}$  should be feasible. In contrast,  $ncol$  or *too close to the front vehicle* are factors that can often be dealt with by braking or swerving correspondingly (i.e.,  $swBr$ ). Thus, reaching  $0^{ncol}$  should be possible. In this sense,  $ncol$  is *directly reducible*.

Risk awareness emerges from combining several factors (e.g.  $F = \{loc_d, ncol, col, st_{rA}, fod_{br}, deg_{br}\}$ ) into a risk space  $R(F)$ , detecting the current risk state, and estimating the likelihood of neighbouring, potentially worse, risk states. Our example suggests that risk awareness contains a notion of machine health. One might expect an autonomous vehicle to take full responsibility of emergency control in any of the states in  $R(F)$ .



**Fig. 3.** Factor instances ( $deg_{br}$ ,  $fod_{br}$ ,  $ncol$ ,  $col$ ) and their dependencies (thick arcs, e.g. requires)

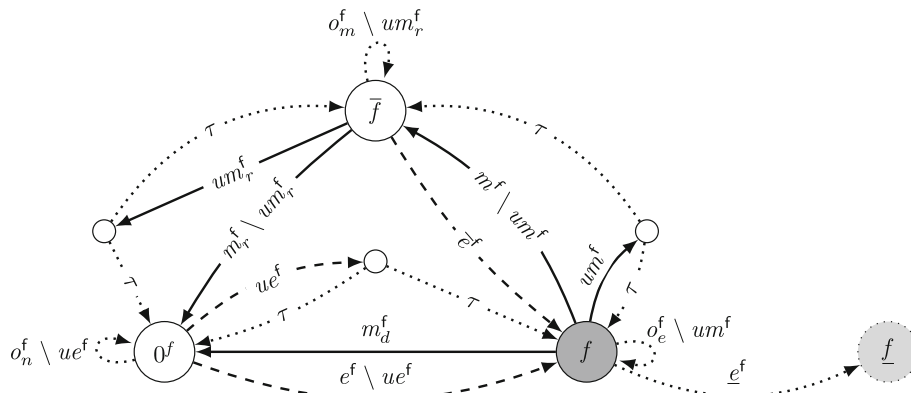


Fig. 4. Risk factor  $f$  with non-deterministic fractions  $ue^f$ ,  $um^f$ , and  $um_r^f$ . The label  $\tau$  is to be read as  $\{\tau\}$

This example illustrates how the proposed formalism captures and guides the way of thinking of safety engineers responsible for developing and assuring safety controllers of autonomous machines.

#### 4.4. Discussion: abstraction and types of risk factors

**Types of risk factors and risk states** Factors can be used to model faults, failure modes, hazards, incidents, and accidents. Final factors can model, for example, *permanent and off-line repairable faults*, and reducible factors, for example, *transient and on-line repairable faults*.

States with an *active final factor*  $f$  (Sect. 4.1) expose the process  $P$  to *residual risk infinitely long and often*, thus, almost certainly materialising the consequences of  $f$ . By using  $\mathfrak{R}$ , a risk-aware machine in a process  $P$  should govern its (often probabilistic) choices to not enter such states.

**Uncertainty in risk factors** Figure 1 illustrates the abstraction from a possibly probabilistic process  $P$  into a deterministic risk structure  $\mathfrak{R}$  with a single factor, now call it  $f$ . This abstraction could preserve the probability of occurrence of  $f$ , if estimated, in  $\mathfrak{R}$ . However, the quantities of  $f$  (especially the probabilities) will often not be confidently known even if they are continuously estimated during operation. But, as indicated in Fig. 4, assumptions about the events of  $f$  might still be made, for example,

1. uncertainty about the actual occurrence of an endangerment  $e^f$  expressed as  $ue^f = o_n^f \cap e^f$ ,
2. uncertainties  $um^f = o_e^f \cap m^f$  and  $um_r^f = o_m^f \cap m_r^f$  in the success of mitigation  $m^f$  and recovery  $m_r^f$ , or
3. the likelihood of an accident  $\underline{f}$  from  $f$ .

Deterministic risk factors ( $ue^f = \emptyset$ , Definition 3) assume that  $\mathfrak{R}$  is certain about the preconditions that activate  $f$ , that is, from observing  $e^f$ ,  $\mathfrak{R}$  knows that the machine enters  $f$ . However,  $ue^f \neq \emptyset$  means that  $\mathfrak{R}$  can observe  $e^f$  but  $P$  does not necessarily move to a state fulfilling  $f$ .

An example for 1. could be an obstacle tracked within safe braking distance with low confidence, justifying both states  $0^f$  and  $f$ . This situation requires further state estimation (e.g. a homing algorithm [NSV03]) to disclose further information about the tracked obstacle. We also discuss this below in Sect. 7.6.

Moreover,  $e^f \setminus ue^f$  can describe known causes of  $f$  and  $ue^f$  *potential* causes of  $f$ .  $o_n^f \supseteq e^f$  means that  $\mathfrak{R}$  is uncertain about all of  $f$ 's causes. Following [Ros10, p. 116],  $ue^f$  is followed by non-deterministic choice over  $0^f$  and  $f$ , modelled by  $\tau$  from an anonymous state ( $\circ$ ) in the factor LTS. Analogously, with  $um^f$  for  $f$  and  $um_r^f$  for  $\bar{f}$ . The factor in Fig. 4 is equivalent to the one in Fig. 2a according to the following lemma.

**Lemma 3 (Isolating uncertainty preserves factor properties)**

$$f_{(4)} =_{FD} f_{(2a)}$$

*Proof sketch.* The proof is by showing that in both cases, considering uncertainty explicitly by isolating non-determinism (Fig. 4) and considering uncertainty implicitly by not isolating non-determinism (Fig. 2a), we deal with the same factor model. Details on the proof can be taken from Appendix A.  $\square$

For the sake of simplicity of the discussion, we neglect further non-determinism possible for the three other actions  $\bar{e}^f$ ,  $m_d^f$ , and  $e^f$ . However, our discussion suggests that in future work we can explore a factor theory addressing those factor types that are likely to occur frequently in practical applications.

## 5. Mitigation orders

This section investigates three basic orders over risk spaces, all intended to support the evaluation (e.g. comparison) of successors of a particular risk state reached after further endangerments or deployed mitigations. The *fully and partially comparable inclusive mitigation orders* (Sect. 5.1), both partial orders, are based on qualitative information, and the *strong mitigation order* (Sect. 5.2), a linear order, is based on quantitative information. These orders depend on the available information about risk and are related (Sect. 5.3).

### 5.1. Qualitative mitigation orders

Let  $R(F)$  be a risk space (Definition 5) for a factor set  $F \subseteq \mathcal{F}$ . Then, we define a partial order  $\preceq_m \subseteq R \times R$  as follows.

**Definition 8 (Fully comparable inclusive mitigation order)** For any pair of states  $\sigma, \sigma' \in R$ , define

$$\sigma \preceq_m \sigma' \iff \forall f \in F: \sigma(f) \preceq_f \sigma'(f).$$

By  $\sigma <_m \sigma' \iff \sigma \preceq_m \sigma' \wedge \sigma \neq_m \sigma'$ , we induce the corresponding strict order.  $\sigma$  and  $\sigma'$  are said to be *incomparable* if and only if  $\sigma \not\preceq_m \sigma' \wedge \sigma' \not\preceq_m \sigma$ . Intuitively,  $\sigma \preceq_m \sigma'$  signifies that “ $\sigma'$  is a better achievement in risk mitigation than  $\sigma$ .”<sup>14</sup> However,  $\preceq_m$  requires full comparability of two states. It might be cumbersome to require such comprehensive knowledge to determine which state is “more or less risky” than another. Allowing for both epistemic and aleatory uncertainty, we might instead account for partial knowledge in the phase orders (Definition 3) at the level of  $R$  by providing a relaxed partial order as follows.

**Definition 9 (Partially comparable inclusive mitigation order)** For states  $\sigma, \sigma' \in R$ , define

$$\sigma \lesssim_m \sigma' \iff \forall f \in F: \sigma(f) \preceq_f \sigma'(f) \vee ((\sigma(f), \sigma'(f)) \notin \preceq_f \wedge (\sigma'(f), \sigma(f)) \notin \preceq_f).$$

We use  $<_{\sim m}$  and  $=_{\sim m}$  to distinguish the corresponding strict order and equality for  $\lesssim_m$  from  $<_m$ . Intuitively, Definition 9 requires a “betterment in risk from  $\sigma$  to  $\sigma'$ ” based exactly on the comparable phases.

**Lemma 4** For any pair of risk states  $\sigma, \sigma' \in R$ , we have that

$$\sigma \preceq_m \sigma' \Rightarrow \sigma \lesssim_m \sigma'.$$

*Proof of Lemma 4.*  $\preceq_f$  is antisymmetric. By definition of  $\preceq_m$ , we may assume

$$\begin{array}{ll} \forall f \in F: \sigma(f) \preceq_f \sigma'(f) & (\forall\text{-elim}) \\ \vdash \sigma(f) \preceq_f \sigma'(f) & (\forall\text{-intro1}) \\ \vdash \sigma(f) \preceq_f \sigma'(f) \vee ((\sigma(f), \sigma'(f)) \notin \preceq_f \wedge (\sigma'(f), \sigma(f)) \notin \preceq_f) & (\forall\text{-intro, assumption for each } f) \\ \vdash \forall f \in F: \sigma(f) \preceq_f \sigma'(f) \vee ((\sigma(f), \sigma'(f)) \notin \preceq_f \wedge (\sigma'(f), \sigma(f)) \notin \preceq_f) & \square \end{array}$$

### Corollary 2

$$\sigma' \not\lesssim_m \sigma \Rightarrow \sigma' \neq_m \sigma$$

<sup>14</sup>This notation might feel unusual as risk reduction is about lowering risk. So, if  $\sigma'$  is a state with *lower risk* than state  $\sigma$  then we could write “ $\sigma \succeq_m \sigma'$ ”. But we might agree that being in  $\sigma'$  is better than residing in  $\sigma$ . Moreover, for risk mitigation, reaching *better* states from worse is in the foreground. Thus, it seems reasonable to use the inverted notation  $\sigma \preceq_m \sigma'$ .

*Proof of Corollary 2.*

$$\begin{aligned}
 \sigma' \not\sim_m \sigma &\Rightarrow \sigma' \neq_m \sigma && \text{(by definition)} \\
 \neg(\sigma' \lesssim_m \sigma \wedge \sigma \lesssim_m \sigma') &\Rightarrow \neg(\sigma' \leq_m \sigma \wedge \sigma \leq_m \sigma') && \text{(by conversion)} \\
 \sigma' \lesssim_m \sigma \wedge \sigma \lesssim_m \sigma' &\Leftarrow \sigma' \leq_m \sigma \wedge \sigma \leq_m \sigma' && \text{(by Lemma 4)} \square
 \end{aligned}$$

In Sect. 7.6, we revisit how features such as partial orders account for uncertainty in RISKSTRUCTURES.

## 5.2. Quantitative mitigation orders

So far, we have seen how partial orders account for a lack of knowledge and potential uncertainties about risk. Now, we will investigate the use of *impact or consequence* data in form of severity intervals, if available for specific factors, to interpolate knowledge gaps, model uncertainty, and derive a linear order over  $R$ .

We continue with definitions for dealing with intervals. Given two intervals  $[l_1, u_1], [l_2, u_2] \subset \mathbb{R}_+$ , the *convex hull* is a map  $\sqcup: \mathbb{R}_+^2 \times \mathbb{R}_+^2 \rightarrow \mathbb{R}_+^2$  given by

$$[l_1, u_1] \sqcup [l_2, u_2] = [\min\{l_1, l_2\}, \max\{u_1, u_2\}]. \quad (4)$$

For a family of  $n \in \mathbb{N}$  intervals  $I = ([l_i, u_i])_{i \in [1..n]}$ , we use the abbreviation  $\sqcup I = [l_1, u_1] \sqcup \dots \sqcup [l_n, u_n]$ . Furthermore, let *active*:  $R(F) \rightarrow 2^F$  with  $\text{active}(\sigma) = \{f \in F \mid \sigma(f) = f\}$  be the map returning the set of active factors of a risk state. Moreover, let  $S: R \rightarrow \mathbb{R}_+^2$  with

$$S(\sigma) = \bigsqcup (f.s)_{f \in \text{active}(\sigma)} \quad (5)$$

be a map for the construction of the severity interval of a risk state from the intervals of its factors.<sup>15</sup> Whereas the *minimal severity* of a factor  $f$  is given by  $f.s = [0, 0)$ , the minimal severity  $S(\sigma)$  of a state  $\sigma$  is the empty interval  $\square$ , which we equate with the empty set,  $\square = \emptyset$ . For any two real-valued intervals  $[a, b], [c, d] \in \mathbb{R}_+^2$ , Ishibuchi and Tanaka [IT90] define with  $[a, b] \leq [c, d] \iff a \leq c \wedge b \leq d$  a partial order over such intervals.

Moreover, we say that two risk states  $\sigma, \sigma' \in R$  are *severity-equivalent* if and only if their accumulated severity intervals are equal, that is,  $\sigma \sim_s \sigma' \iff S(\sigma) = S(\sigma')$ . We have that  $\sigma = \sigma' \Rightarrow \sigma \sim_s \sigma'$  because the factors that are in their active phases are identical. The relation  $\sim_s$  is an equivalence relation because it is reflexive, symmetric, and transitive (all by the usual equivalence over intervals). Furthermore,  $\sim_s$  induces equivalence classes  $[\sigma]_{\sim_s} = \{\sigma' \in R \mid \sigma' \sim_s \sigma\}$  over  $R$  for any  $\sigma \in R$  with the corresponding quotient class  $R/\sim_s$ . With the family  $(f.s)_{f \in F}$  of *severity intervals* of  $F$ , we now define an order over  $R/\sim_s$ .

**Definition 10 (Strong mitigation order)** For  $[\sigma]_{\sim_s}, [\sigma']_{\sim_s} \in R/\sim_s$ , define

$$[\sigma]_{\sim_s} \leq_m [\sigma']_{\sim_s} \iff \forall \check{\sigma} \in [\sigma]_{\sim_s}, \check{\sigma}' \in [\sigma']_{\sim_s}: S(\check{\sigma}) \geq S(\check{\sigma}') \vee S(\check{\sigma}') \subset S(\check{\sigma}).$$

$[\sigma]_{\sim_s} \leq_m [\sigma']_{\sim_s}$  can be dropped from  $R/\sim_s$ , yielding

$$\forall \check{\sigma} \in [\sigma]_{\sim_s}, \check{\sigma}' \in [\sigma']_{\sim_s}: \check{\sigma} \leq_m \check{\sigma}' \iff S(\check{\sigma}) \geq S(\check{\sigma}') \vee S(\check{\sigma}') \subset S(\check{\sigma}). \quad (6)$$

$\leq_m$  codifies that  $\sigma'$  (a) *reduces or* (b) *focuses risk* if the union of its severity intervals is

1. *Lower* in the ranking  $\leq$  of interval numbers or
2. *Strictly narrower* than the corresponding union for  $\sigma$ .

Condition (a) seems immediately intuitive. Condition (b) conveys the intuition that the interval carries less uncertainty about the consequences expected from  $\sigma'$  than from  $\sigma$ . Equivalence classes in  $R/\sim_s$  abstract from the factors from which the merged severity intervals originate. This abstraction has to be carefully taken into account when using  $\leq_m$  and, therefore, when specifying severity. Note that  $\leq_m$  is based on the convex hull of severity intervals from the active phases of a pair of risk states. Apart from the convex hull, interval addition and multiplication are relevant for alternative mitigation orders as we shall see below. However, a detailed investigation is left for future work. Let us now consider some properties of  $\leq_m$ .

<sup>15</sup>Note that  $S$  over-approximates (i.e., constructs the convex hull from) sparsely distributed severity intervals.

**Lemma 5**  $\leq_m$  is linear over  $R/\sim_s$ .

*Proof sketch.* We show by case analysis that any two risk states are comparable and  $\leq_m$  is antisymmetric. The complete proof is stated in Appendix A.  $\square$

**Corollary 3** After dropping Lemma 5 by Formula (6), we have that  $\leq_m$  is also linear over  $R$ .

**Lemma 6**  $(R, \leq_m)$  and  $(R/\sim_s, \leq_m)$  are well ordered.

*Proof sketch.* Lemma 6 follows from a finite  $R$  (by definition) and, thus, finite  $R/\sim_s$ , and linearity of  $\leq_m$  (by Lemma 5).  $\square$

**Definition 11** For  $\sigma \in R(F)$ , we also write  $\mathbf{0}^F \equiv \forall f \in F: \sigma(f) = 0^f$  and  $\mathbf{F} \equiv \forall f \in F: \sigma(f) = f$ . We denote by  $\top^F$  the set of *maximal elements* and by  $\perp^F$  the set of *minimal elements* of  $(R, \leq_m, \lesssim_m, \leq_m)$ . We characterise the minimal elements in  $(R(F), \leq_m)$  by

$$\perp^F \equiv \{\sigma \in R(F) \mid \forall \sigma' \in R(F): \sigma' \leq_m \sigma \Rightarrow \sigma' =_m \sigma\}$$

and analogously for  $(R, \lesssim_m)$  and  $(R, \leq_m)$  and the maximal elements.

**Corollary 4** If  $\forall f \in F: (Ph_f, \leq_f)$  is linear, then  $(R(F), \leq_m) = (R(F), \lesssim_m)$ . If  $R \neq R(\emptyset)$  then  $\perp^F$  and  $\top^F$  are non-empty and, therefore, have a proper manifestation. For  $(R/\sim_s, \leq_m)$ ,  $\perp^F$  and  $\top^F$  are singletons.

*Proof of Corollary 4.* The proof is by contradiction. For the sake of brevity, we only consider a sketch of this proof. Assume we have two state classes  $[\sigma]_{\sim_s}, [\sigma']_{\sim_s}$  in  $\perp^F$  with  $[\sigma]_{\sim_s} \neq_m [\sigma']_{\sim_s}$ . Because of our assumption, state classes are in linear order. Thus, by definition of  $\perp^F$ , one of these state classes causes a violation of the universal quantification in Definition 11 and, therefore, one of the classes cannot be in  $\perp^F$  which contradicts our assumption. The proof is analogous for  $\top^F$ .  $\square$

**Corollary 5**  $(R/\sim_s, \leq_m)$  forms a complete lattice.

*Proof of Corollary 5.* Linearity of  $\leq_m$  implies that every non-empty subset of  $R/\sim_s$  has a greatest lower bound and a least upper bound.  $\square$

### 5.3. Relating mitigation orders

The strong mitigation order characterised by the Lemmas 5 and 6 is driven by the number of active factors and their severity intervals (because of the definition of  $S$ ) but not by the equality of factor phases among the compared risk states. This offers the possibility of abstraction from individual factors and focusing on severity estimates. To avoid infeasible models (e.g. specifications that get too strong to be realisable), we *require* that the addition of severity intervals constitutes a *relational extension* of either  $\leq_m$  or  $\lesssim_m$ , formally,

$$(\forall \check{\sigma} \in [\sigma]_{\sim_s}, \check{\sigma}' \in [\sigma']_{\sim_s}: \check{\sigma} \leq_m \check{\sigma}' \vee \check{\sigma} \lesssim_m \check{\sigma}') \Rightarrow [\sigma]_{\sim_s} \leq_m [\sigma']_{\sim_s}.$$

Dropped to  $R$ , this implies  $\check{\sigma} \leq_m \check{\sigma}' \vee \check{\sigma} \lesssim_m \check{\sigma}' \Rightarrow \check{\sigma} \leq_m \check{\sigma}'$  for all pairs  $(\check{\sigma}, \check{\sigma}') \in [\sigma]_{\sim_s} \times [\sigma']_{\sim_s}$ , therefore,

$$\check{\sigma} \leq_m \check{\sigma}' \vee \check{\sigma} \lesssim_m \check{\sigma}' \Rightarrow S(\check{\sigma}) \geq S(\check{\sigma}') \vee S(\check{\sigma}') \subset S(\check{\sigma}) \quad (7)$$

for full and partial comparability, otherwise implying

$$(\check{\sigma}, \check{\sigma}') \notin \leq_m \wedge (\check{\sigma}, \check{\sigma}') \notin \lesssim_m \Rightarrow \top. \quad (8)$$

Intuitively, if  $\check{\sigma}$  is “worse” than  $\check{\sigma}'$  then its accumulated severity interval  $S(\check{\sigma})$  has to be greater than that of  $\check{\sigma}'$  and, therefore, must not be contained in that of  $\check{\sigma}'$ . Moreover, if  $\check{\sigma}$  and  $\check{\sigma}'$  are incomparable in  $\leq_m$  and  $\lesssim_m$  (i.e., some factors have inversely ordered or incomparable phases) then  $S(\check{\sigma})$  and  $S(\check{\sigma}')$  are allowed to form any relationship (signified by  $\top$  for “true”), for example, Formula (7).

What is the (necessary and) *sufficient condition* on  $F$  to satisfy the requirement expressed by Formula (7)? Risk spaces and risk state pairs are the interpretations and, therefore, potential models satisfying the relational extension imposed by Formula (7). Answering this question suggests the following lemma.

**Lemma 7** For  $\sigma, \sigma' \in R(F)$ ,

$$\sigma \leq_m \sigma' \vee \sigma \lesssim_m \sigma' \Rightarrow \text{active}(\sigma) \supseteq \text{active}(\sigma').$$

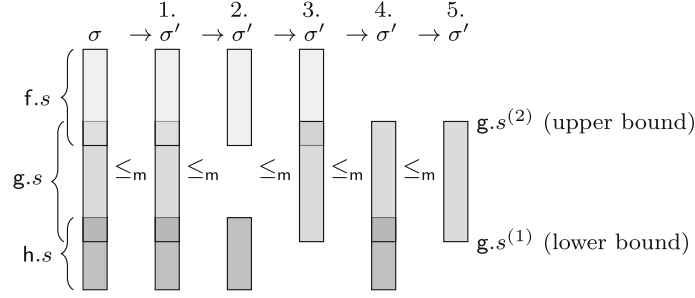


Fig. 5. Visualisation of the case distinction to establish Theorem 1, with  $f, g, h \in \mathcal{F}$ ,  $\text{TS}(R, \rightarrow)$ , and  $\sigma \rightarrow \sigma'$

*Proof sketch.* The proof is by induction over  $F$  and relies on the assumptions that, for any  $f \in F$ ,  $f$  is the *unique maximal element* in  $\leq_f$  and that the map *active* (Sect. 5.2) only returns such elements. The whole proof is stated in Appendix A.  $\square$

Again, fix a finite  $F$  and a pair  $\sigma, \sigma' \in R(F)$  and assume  $\sigma \leq_m \sigma' \vee \sigma \lesssim_m \sigma'$ . Then, by Lemma 7,  $\sigma'$  incorporates a *subset* of  $\sigma$ 's active factors. To show that  $S(\sigma')$  preserves  $S(\sigma') \subset S(\sigma)$ , the right-hand part of the disjunction in the consequent of Formula (7), it is sufficient to investigate  $S(\sigma')$  for any factor deactivation possible in a transition from  $\sigma$  to  $\sigma'$ . We can summarise all such possibilities in five cases:

1. all intervals remain ( $\sigma' = \sigma$ ),
2. only intervals in the convex hull of the remaining intervals are removed,
3. intervals only increasing the lower bound of this hull are removed,
4. intervals only decreasing the upper bound of this hull are removed, and
5. intervals increasing the lower bound and decreasing the upper bound of this hull are removed.

This case distinction is visualised in Fig. 5 and also works for  $\sigma' \prec_m \sigma$ . Formula (7) is also satisfied if all factors in  $F$  are assigned the same interval, call it  $s^F$ . In conclusion, the sufficient condition on  $F$  to satisfy Formula (7) is the “unique maximal element” precondition in the proof of Lemma 7. Apart from this precondition, Formula (7) holds of an arbitrary finite  $F \subseteq \mathcal{F}$ . Below, we shall call  $\leq_m$  and  $\lesssim_m$  *inclusive mitigation orders*, and  $\leq_m$  a *strong mitigation order*. We arrive at the following theorem.

**Theorem 1** The *strong mitigation order*  $\leq_m$  extends the *partially comparable inclusive mitigation order*  $\lesssim_m$  which, in turn, extends the *fully comparable inclusive mitigation order*  $\leq_m$ . Formally, for  $\sigma, \sigma' \in R$ :

$$\sigma \leq_m \sigma' \xrightarrow{\text{Lemma 4}} \sigma \lesssim_m \sigma' \xrightarrow{\text{Lemma 7}} \sigma \leq_m \sigma'.$$

#### 5.4. Application: evaluating local, regional, and global safety

The orders  $\leq_m$ ,  $\lesssim_m$ , and  $\leq_m$  allow a *local* evaluation of safety in the sense that their definitions only require the comparison of pairs of risk states. Two further qualitative notions of safety seem to be useful.

Let  $R$  be non-empty and finite and  $\text{reach}: R \times \mathcal{P} \rightarrow 2^R$ . Given a process  $P \in \mathcal{P}$  and a risk state  $\sigma \in R$ ,  $\text{reach}(\sigma, P) \subseteq R$  denotes the set of risk states reachable from  $\sigma$  by finite executions (prefixes) of  $P$  where  $\sigma$  itself is always reachable and, thus,  $\sigma \in \text{reach}(\sigma, P)$ . Figure 6 (light-grey area) exemplifies such a set. Then, we use the presented orders to determine non-empty sets of minimal and maximal elements in  $R$  reachable by  $P$ , such as  $\max_{\leq_m} \text{reach}(\sigma, P)$  or  $\min_{\lesssim_m} \text{reach}(\sigma, P)$ . In Fig. 6, the risk states 1 to 3 and 5 to 7, connected by double arcs, are the ones reachable from  $\sigma$  after at most two consecutive endangerments or mitigations. Reachability arcs are labelled with the weakest applicable order. Single arcs illustrate potential modifications or extensions of the reachability set, for example, by the states 4 and 8 to 10 after a reassessment of the operational environment by the machine.

In the situation described by the process  $P$  in the risk state  $\sigma$ , these two sets signify the *regionally safest* (max, white circles) and the *regionally most hazardous* (min, dark-grey circles) states, respectively. The smallest such set will only and exactly contain  $\sigma$ , meaning  $P$  cannot reduce risk in  $(R, \lesssim_m)$ .  $\leq_m$  and  $\lesssim_m$  enable a *regional* evaluation of safety inasmuch as once a maximal element in  $\max_{\lesssim_m} \text{reach}(\sigma, P)$  is reached,  $\lesssim_m$  limits reasoning about safer states that  $P$  could reach from  $\sigma$ .



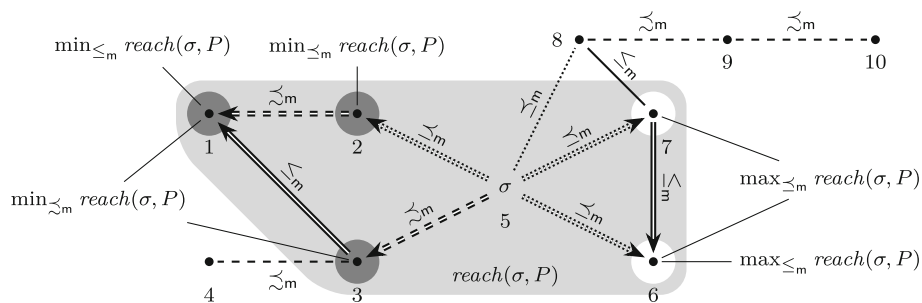


Fig. 6. Visualisation of local, regional, and global safety in a risk space  $R$

Besides local and regional safety,  $(R/\sim_s, \leq_m)$  supports *global* evaluation because of its linearity (Lemma 5). With  $\leq_m$ , there are always *unique* safest and riskiest states in  $\text{reach}(\sigma, P)/\sim_s$  (Corollaries 4 and 5). In contrast,  $\leq_m$  and  $\lesssim_m$  will not guarantee this uniqueness. Note that the use of equivalence classes leads to more abstract forms of safest and riskiest states. Overall, Lemma 6 and Corollary 5 provide necessary conditions for deriving *finite strategies* (i.e., policies, choice resolutions) that stabilise or terminate  $P$  in a safest state.

Quantitative mitigation orders (e.g.  $\leq_m$ ) will have to be calculated at run-time, according to  $\mathfrak{R}$  by estimating probabilities of factor occurrence and severity of factor consequences from situational data only available during operation. In  $P$ , *safety* can then be observed as the *gradual presence or absence of risk over traces* ( $P$ ). On a side note, the *reliability* as the *gradual presence or absence of defective behaviour* can be seen as a special case when considering only factors that model faults.

## 5.5. Discussion: ethical aspects of mitigation orders

Linear mitigation orders such as  $\leq_m$  promote machines with *negative utilitarian* decision ethics [War12, p. 51]. For example, severity intervals could be calculated at run-time based on sensor data about the possible operational situation of the machine. Expected outcomes of enabled mitigations, if any, will then be comparable according to  $\leq_m$ . This comparability allows the assessment of the actual reachability of states with strictly lower risk. Any resolution of a near-accident situation (Sect. 4.3) or a *tram*<sup>16</sup> *problem* [Foo78] would then consist in the choice of the mitigation leading to the state with the lowest risk or *the least severe of the expected negative outcomes*. This scheme characterises negative utilitarianism.

Linear mitigation orders globally resolve decisions based on explicit and, therefore, disputable criteria. Consequently, utilitarian ethics have been criticised to lead to oversimplified approaches to resolve indecision. Such critiques stress the difficulty of predicting the positive and negative effects (i.e., “double effects” [Foo78]) of certain actions, in our case, the estimation of the severity of an activated factor and the reached state [War12, pp. 48-49]. Appreciating the equality of any two groups of the human family [Uni48], one could conclude that tram problems should be solved by random decisions, independent of whether these are made by humans or machines. However, in analogy to risk-averse humans (Sect. 3), a risk-aware machine should make (self-inflicted) tram problems unlikely, that is, avoid such decisions rather than aiming at their random or utilitarian resolution. The likelihood of tram problems is an important subject to be investigated.

Accordingly, a structured risk model with dependencies between factors could be used to complement utilitarian decision ethics with Kantian ethics, that is, to use  $\leq_m$  to decide about conservative measures *before* high-severity factors get activated. This results in what is called “rule utilitarianism” [War12, p. 52]. For example, we can model the *necessary preconditions* of certain tram problems as risk factors and a machine based on this model could use these factors to constrain its behaviour. Although the presented model can be used with linear orders, the discussions below stay agnostic of the mitigation order.

<sup>16</sup>Also known as the “trolley problem”, a situation where any of the enabled actions leads to an unacceptable outcome.

## 6. Factor dependencies

Usually, only subsets of a risk space  $R$  and a process  $P$ 's traversals of  $R$  are relevant. Identifying these subsets is the task of risk analysis and can be difficult if the relevance of each risk state can only be determined at run-time based on its context in  $R$  and the state of  $P$ . As highlighted in Sect. 3, it is then useful to model (*causal*) *relationships* between (phases of) factors, their events (i.e., activation, mitigation), and, consequently, risk states. For example, we might want to specify that (i) the *activation* of a risk factor *causes* the *activation* of another factor, (ii) the *mitigation* of one factor *causes* the *activation* of several other factors, or (iii) the *activation* of a factor *requires* the *activation* of some other factors. We can take account of dependencies, such as illustrated in Table 2 (e.g. *requires*  $\text{ncol}$ , *causes*  $\text{loc}_d$ ), by imposing *constraints* on pairs of consecutive risk states and their comprising factors' phases. For this, we use binary relations over  $R$  to hypothesise *causality assumptions* about less known or controllable parts of  $P$  (i.e., the environment  $En$ ) and express *causality requirements* for known and controllable parts of  $P$  (i.e., the system  $Sy$ ).

### 6.1. Relations over risk spaces

Let  $F \subseteq \mathcal{F}$  and consider the space  $R(F)$  and two distinct factors  $f, g \in F$ . In the following, we employ relational specification to formalise factor dependencies as relations over  $R$ , that is, as subsets of  $R \times R$ . For example, a *causes*<sup>17</sup> constraint requires of a pair  $(\sigma, \sigma') \in R \times R$  that if  $f$  is active then, within at most one (logical) step,  $g$  must be active until  $f$  gets either inactive or mitigated.

Let  $\mathcal{C}$  be the set of all constraints. For the translation of constraints into a form useful for the discussion below, we use the map  $\llbracket \cdot \rrbracket_c^R: \mathcal{C} \rightarrow 2^{R \times R}$  to denote the relational semantics of constraints over  $R$ . To distinguish different semantic maps, as will be defined later, we use subscripts, for example,  $c$  in  $\llbracket \cdot \rrbracket_c^R$ . Recall that we use  $\sigma(f)$  to refer to the phase of factor  $f$  in state  $\sigma$  (Sect. 4.2).

For example, *causes* constraints can be encoded as relations over  $R$  by the definition

$$\llbracket f \text{ causes } g \rrbracket_c^R = \{(\sigma, \sigma') \in R \times R \mid \underbrace{\sigma(f) \neq f \wedge \sigma'(f) = f}_{\text{activation of } f \dots} \wedge \underbrace{\sigma(g) \neq \bar{g}}_{\text{when } g \text{ not mitigated}} \Rightarrow \underbrace{\sigma'(g) = g}_{\dots \text{activates } g}\} . \quad (9)$$

This constraint implies that factor activation takes a logical time step corresponding to a strictly positive real-time duration (i.e.,  $> 0$ ) in  $P$ . Note that this definition of *causes* encodes the assumption that  $g$ , once or as long as mitigated, cannot be caused (again) by activating  $f$ . While this behaviour is often to be expected of any mitigation designed for  $g$ , one can define a more general constraint, say *alwaysCauses*, to capture this behaviour. Moreover, *causes* is to be read as “is sufficient to cause” rather than as “is necessary to cause”, and it extends to factor sets. Given  $F, G \subseteq \mathcal{F}$  with  $F \cap G = \emptyset$ , we define

$$\llbracket F \text{ causes } G \rrbracket_c^R = \{(\sigma, \sigma') \in R \times R \mid \exists f \in F \forall g \in G: \sigma(f) \neq f \wedge \sigma'(f) = f \wedge \sigma(g) \neq \bar{g} \Rightarrow \sigma'(g) = g\} .$$

Note that all pairs  $(\sigma, \sigma')$  violating the antecedent of the conditional are in  $\llbracket f \text{ causes } g \rrbracket_c^R$  as well.

As another example, the *requires* constraint can be defined in relational form by

$$\llbracket f \text{ requires } g \rrbracket_c^R = \{(\sigma, \sigma') \in R \times R \mid \underbrace{\sigma'(f) = f}_{\text{an active } f \text{ requires} \dots} \Rightarrow \underbrace{\sigma(g) = g}_{\text{an active } g \text{ in the preceding state}}\} . \quad (10)$$

The lifting of *requires* to factor sets  $F, G \subseteq \mathcal{F}$  with  $F \cap G = \emptyset$  is described as

$$\llbracket F \text{ requires } G \rrbracket_c^R = \{(\sigma, \sigma') \in R \times R \mid \exists f \in F: \sigma'(f) = f \Rightarrow \forall g \in G: \sigma(g) = g\} .$$

This variant of the *requires* constraint refers to *all* factors specified on its right-hand side and, this way, resembles an AND-gate as used in FTA. The side condition  $F \cap G = \emptyset$  avoids that a factor requires or causes itself (i.e., it avoids the “chicken and egg” problem).

<sup>17</sup>Used in form of MCSs in FTA and, less frequently formally, in FMEA (see Sect. 2.1).

**Table 3.** A selection of constraints more commonly used in causal reasoning

Name	PC	D	P	C	M	FC	Y	Description
causes	$F \rightleftarrows G$	$\rightarrow$	+	w	1:n	=n	✓	The <i>activation</i> of a factor in $F$ <i>causes</i> the (successive) <i>activation</i> of the factor set $G$ , Formula (9).
causes <sup>-1</sup>	$\overline{F} \rightleftarrows G$	$\rightarrow$	+	w	m:n	=n	-	The <i>mitigation</i> of the factor set $F$ <i>causes</i> the <i>activation</i> of the factor set $G$ .
requires	$F \rightleftarrows G$	$\leftarrow$	+	s	1:n	=n	✓	The <i>activation</i> of a factor in $F$ <i>requires</i> the (prior) <i>activation</i> of all factors in $G$ . AND-gate in FTA, Formula (10).
requires <sub>1</sub>	$F \rightleftarrows G$	$\leftarrow$	+	s	1:n	$\geq 1$	✓	The <i>activation</i> of a factor in $F$ <i>requires</i> the (prior) <i>activation</i> of at least <i>one</i> factor in $G$ ; OR-gate in FTA.
prevents	$F \rightleftarrows G$	$\rightarrow$	-	s	1:n	=n	✓	The <i>activation</i> of a factor in $F$ <i>prevents</i> or <i>impedes</i> the <i>activation</i> of all factors in $G$ .
preventsMit	$F \rightleftarrows \overline{G}$	$\rightarrow$	-	s	1:n	=n	✓	The <i>activation</i> of a factor in $F$ <i>prevents</i> or <i>impedes</i> the <i>mitigation</i> of all factors in $G$ .
excludes	$F \rightleftarrows 0^G$	$\rightarrow$	+	s	1:n	=n	✓	The <i>activation</i> of a factor in $F$ <i>deactivates</i> (superposes or invalidates) all factors in $G$ .
direct	$F \rightleftarrows 0^G$	$\rightarrow$	+	s	1:0	=1	✓	The <i>deactivation</i> of a factor in $F$ <i>must</i> only happen directly from its <i>activated</i> phase.
offRepair	$F \rightleftarrows 0^G$	$\rightarrow$	-	s	1:0	=1	✓	The <i>deactivation</i> of a factor in $F$ <i>may not</i> happen directly from its <i>activated</i> phase.

Legend: See the *dimensions* discussed in the text of Sect. 6.1; ✓ ... implemented in YAP [Gle21]

**Remark 1** *causes* models weak causality and *requires* models a strong causality. *causes* and *requires* constitute basic dependency templates and exemplify how dependencies support the safety engineer or risk analyst in causal modelling in a state-based relational way.

Analogously, there are many possible *factor dependencies* over  $R$  that resemble causal reasoning of techniques such as FTA (Sect. 2.1). Such dependencies can be classified along several dimensions:

- *phase combination* (PC): active to active ( $F \rightleftarrows G$ ), active to mitigated ( $F \rightleftarrows \overline{G}$ ), active to inactive ( $F \rightleftarrows 0^G$ ), mitigated to mitigated ( $\overline{F} \rightleftarrows \overline{G}$ ), mitigated to active ( $\overline{F} \rightleftarrows G$ ), mitigated to inactive ( $\overline{F} \rightleftarrows 0^G$ ); combinations thereof are possible;
- *direction* (D) of cause-effect analysis: forward ( $\rightarrow$ , e.g. for modelling sufficient conditions or propagation), backward ( $\leftarrow$ , e.g. for modelling necessary conditions or explanation);
- *polarity* (P): obligation (+), permission ( $\circ$ ), inhibition (-);
- *causality* (C): strong (s, sequential events), weak (w, simultaneous events);
- *multiplicity* (M): one-to-one (1 : 1), one-to-many (1 : n), many-to-one (n : 1), many-to-many (m : n) where  $m, n > 0$ ; self-referential<sup>18</sup> (1 : 0);
- *factor combination* (FC): “ $\sim m$  out of  $n$ ” where  $\sim \in \{\leq, =, \geq\}$  and  $1 \leq m \leq n$ .

These dimensions<sup>19</sup> characterise factor dependencies commonly used in risk analysis. Most of the constraints listed in Table 3 are realised in our tool YAP [Gle21], which enables their use based on previous discussions in [Gle18, Gle17]. However, their comprehensive treatment would exceed the scope of this work.

We have now seen that constraints prune irrelevant state pairs. This mechanism is reflected by the following definition.

**Definition 12 (Relational semantics of constraints)** For a set of constraints  $C \subseteq \mathcal{C}$ ,

$$\llbracket C \rrbracket_c^R = \begin{cases} R \times R, & C = \emptyset \\ \bigcap_{k \in C} \llbracket k \rrbracket_c^R, & \text{otherwise} \end{cases} \subseteq R \times R. \quad (11)$$

Constraints are an instrument for specifying structures over risk spaces. Regarding the composition of spaces, this instrument requires the definition of the following well-formedness condition.

<sup>18</sup>In this case, we allow  $F = G = \{f\}$ .

<sup>19</sup>Naïve combination of all dimensions results in more than 1000 possibilities to define constraints. Many of these possibilities are not essentially different and some might not even be useful. Table 3 does, however, not aim to cover all useful constraints.

Listing 1: YAP script modelling the situation transferObj from Table 4

```

1 factormodel {
2   dmgo desc "object damage"
3   mishap
4   requiresNOF(1|hp,fodg)
5   sev=[5,10];
6   fodg desc "failure on demand of grabber"
7   requires(slp)
8   causes(dmgo);
9   slp desc "slippery grabber"
10  requires(wet)
11  direct
12  detectedBy(.senseSlip)
13  mitigatedBy(.incPr)
14  sev=[2,6];
15  wet desc "wet grabber surface"
16  causes(slp)
17  direct;
18  hp desc "high grabber pressure"
19  mitPreventsMit(slp)
20  direct
21  detectedBy(.sensePr)
22  mitigatedBy(.limPr)
23  sev=[3,8];
24  /* str desc "spurious release"; equal to fod */
25 }
26
27 controlloop collRob {
28  mode senseSlip desc "identify slip";
29  mode incPr desc "increase grabbing pressure"
30  update "pressure++";
31  mode sensePr desc "sense pressure"
32  guard "pressure > threshold(obj)";
33  mode limPr desc "limit grabbing pressure"
34  update "pressure = threshold(obj)";
35 }

```

Table 4. Situation/action/factor table for grabbing scenarios performed by collaborative robots

	Act.	work ( $w$ )	grab & hold ( $g$ )			release ( $r$ )	
Situ.	Fact.	dmgo object damage	fodg failure on demand	slp slippery grabber	wet wet surface	hp high pressure	str spurious trip
transferObj	3. requires <sub>1</sub> hp, fodg	3. requires <sub>1</sub> slp, str; causes dmgo	2. Pr[slp   slp] = h	3. requires wet 4. increase grab- bing pressure	3. causes slp	2. Pr[hp   hp] = m 3. mitPrevents- Mit slp 4. limit grabbing pressure	–

Legend:  $h$ ... high,  $m$ ... medium, –not analysed; see also Table 2 for guidance

**Definition 13 (Well-formedness of constraints)** Let  $R(F)$  be a space formed by a factor set  $F \subseteq \mathcal{F}$ . We say that a constraint  $k \in C$  is *well-formed for*  $R(F)$  if and only if it does not refer to factors other than the ones in  $F$ , formally, if and only if  $\llbracket k \rrbracket_c^R \subseteq R(F) \times R(F)$ . We say that  $C \subseteq \mathcal{C}$  is well-formed for  $R(F)$  if and only if each element in  $C$  is well-formed for  $R(F)$ .

## 6.2. Example: collaborative robots

To illustrate relationships for structuring risk perception, we apply our framework to a scenario in human-robot collaboration. Table 4 shows a risk analysis similar to Table 2 for a robotic arm with a grabber. Consider such a robot ( $w$ )orking by repetitively ( $g$ )rabbing, *holding*, and ( $r$ )eleasing work pieces. Water or oil on the robot’s grabber ( $wet$ ) can *cause* the grabber to be *slippery* ( $slp$ ) such that *holding* an object may increase the likelihood of *accidentally dropping* it ( $Pr[slp | slp]$ ). For our further analysis, we generalise the negative outcome from the phase  $slp$  to an additional *final* risk factor  $dmgo$ . Obviously, increased *grabbing pressure* can mitigate  $slp$ . However, forward reasoning leads to the conclusion that increased pressure could be *too high* ( $hp$ ) *causing* damage to the object ( $dmgo$ ). However, because we plan to install an *intervention* in our machine, we omit the dependency  $hp$  causes  $dmgo$ . Moreover, a FTA of the sensor software and the actuator hardware results in the factor *mistaken loosening* of the grabber ( $str$ ), which, for the sake of simplicity, we equate with a failure on demand of the “grab & hold” action ( $fodg$ ). Applying final backward causal reasoning, the object’s high *falling* onto a hard surface causes damage of the object ( $dmgo$ ) and requires<sub>1</sub> (at least one of) slippery grabber ( $slp$ ),  $fodg$ , or a spurious trip of the “release” action ( $str$ ).

Listing 1 encodes this analysis in a script processable by YAP, which calculates the risk space and the phase transition relation. The resulting risk graph in Fig. 7 can be useful as a controller design template for a risk-aware machine. Further details on how this graph can be constructed are discussed in Sect. 7.

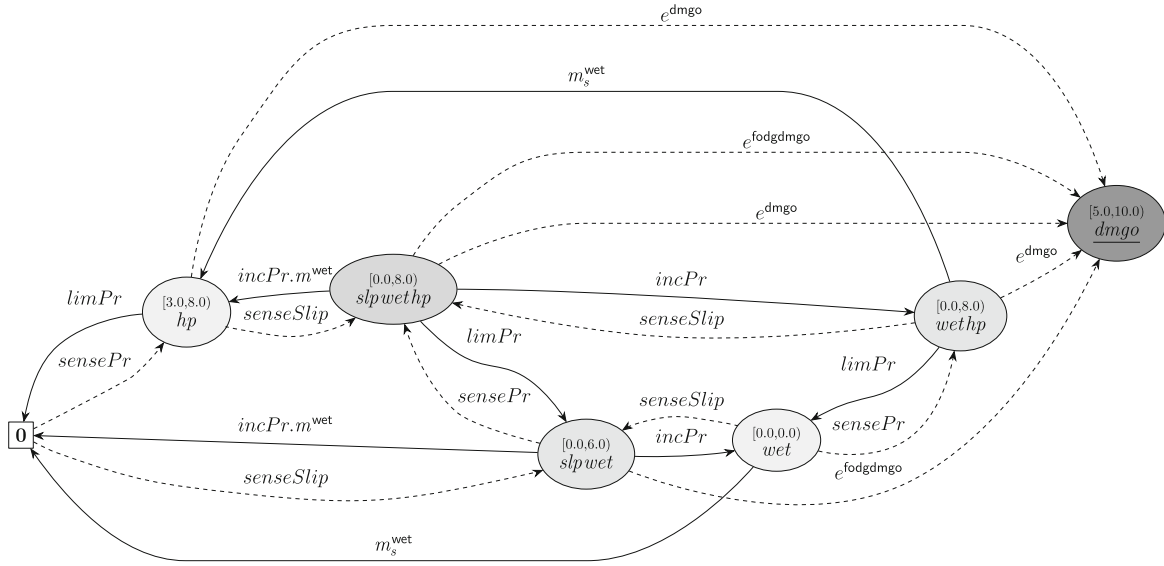


Fig. 7. Graph for risk structure  $\mathfrak{R}_{\text{transferObj}}$  with factor set  $F = \{\text{dmg}, \text{fod}, \text{slp}, \text{wet}, \text{hp}\}$  generated by YAP from Listing 1. To increase readability, inactive factors are not shown in the nodes

### 6.3. Discussion: abstraction, compatibility, and characteristics of risk factors

**Abstraction** Two structurally different risk states  $\sigma, \sigma' \in R$  (i.e.,  $\sigma \neq \sigma'$ ) can be severity-equivalent (i.e.,  $\sigma \sim_s \sigma'$ ). Moreover, the family  $(f.s)_{f \in F}$  of severity intervals of  $F$  forms a *cut of causal chains* and, thus, defines the *scope* of  $\mathfrak{R}$ . These intervals abstract from potential *consequences of active factors*. This abstraction is left to the modeller (i.e., the risk analyst or safety engineer) and can vary significantly, depending on the *assets* (e.g. humans, animals, environment) and *impact categories* (e.g. certain injuries, damage, loss) under discussion. For example, assume two independent consequences  $c_A$  and  $c_B$  (e.g. a class of injuries of an operator  $A$ , a class of damages of an object  $B$ ) associated with *possible*<sup>20</sup> intervals  $[l_A, u_A]$  and  $[l_B, u_B]$ . Consider a *single* factor  $f$  causing both  $c_A$  and  $c_B$ . This circumstance suggests the consistency condition  $f.s \subseteq [l_A + l_B, u_A + u_B]$ , requiring that all consequences of  $f$  (i.e.,  $c_A, c_B$ ) have to be measurable in their severity along the *same scale*. Our example stresses the difficulty of comparing injury with damage. However, if  $c_B$  was another class of injuries of that operator then  $c_B$  would fit into the same scale as  $c_A$ . Practitioners often formalise these concepts for quantitative analysis. However, we leave this for future work.

**Factor compatibility** Consequences of *several* factors should be *compatible* such that the convex hull of their intervals (recall  $\sqcup$  from Formula (4)) has a consistent meaning of severity in  $(R, \leq_m)$ . Consider the two factors  $f$  and  $g$ . Both can describe three cases with their *individually specified* intervals  $f.s$  and  $g.s$ :

1.  $f$  and  $g$  share all their consequences (e.g. both cause  $c_A$ ).  $f, g \rightarrow \textcircled{c_A}$
2.  $f$  and  $g$  share some of their consequences (e.g.  $f$  causes  $c_A$  and  $c_B$ ,  $g$  causes  $c_B$ ).  $f \rightarrow \textcircled{c_A} \textcircled{c_B} \leftarrow g$
3.  $f$  and  $g$  do not share any consequences (e.g.  $f$  causes  $c_A$ ,  $g$  causes  $c_B$ ).  $f \rightarrow \textcircled{c_A} \quad g \rightarrow \textcircled{c_B}$

If both factors get active in case 1, the convex hull can be backed by the condition  $f.s \sqcup g.s \subseteq [l_A, u_A]$ . Case 2 can be split into the  $\leq_m$ -compatible case 1 for  $c_B$  and case 3 for  $c_A$ . In case 3, generally, the convex hull may extend both the range of consequences and the severity intervals. However, the consistency condition  $f.s \sqcup g.s \subseteq [l_A, u_A] \sqcup [l_B, u_B]$  implied by  $\leq_m$  seems inappropriate. For example, if  $c_A$  and  $c_B$  signify damage of two independent objects with the same interval, the hull would not account for this because of idempotency of interval union.

<sup>20</sup>For  $[l_A, u_A]$  and  $[l_B, u_B]$ , we assume sufficient knowledge about possible consequences and their accurate evaluation.

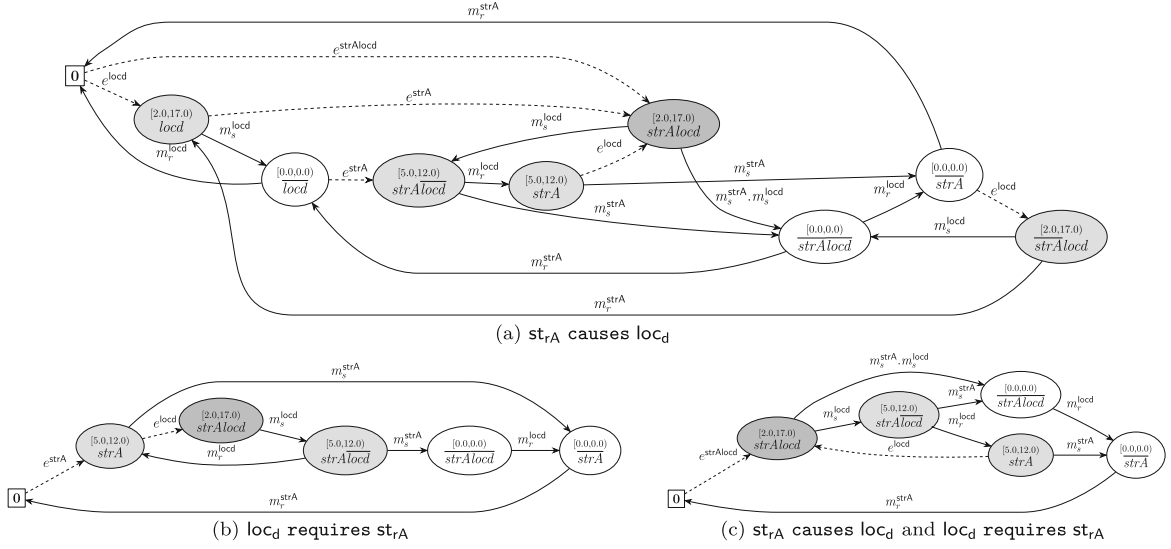


Fig. 8. Risk graphs for  $F = \{st_{rA}, loc_d\}$  after applying causes (a), requires (b), and both (c), using YAP

In summary, while  $\leq_m$  can deal with consequences shared by all risk factors (case 1), partially shared (case 2) or independent consequences (case 3) require a further investigation out of scope.

**Condition 1 (Factor convergence)** The use of  $\leq_m$  with several risk factors requires severity specifications to be based on a single compound consequence and a single severity scale across all factors.

The problematic case 3 could be modelled by an additional factor  $h$  that is caused by  $f$  and  $g$  and carries an up-shifted severity interval, for example,  $h.s = f.s + g.s$ , and  $f$ ,  $g$ , and  $h$  are all defined with respect to a single new consequence  $c_{A+B}$ . In Sect. 6.1 and Table 3, we have seen how such a dependency can be specified by constraints on the risk space, that is, by  $\{f, g\}$  causes  $h$  and  $h$  excludes  $\{f, g\}$ .  $excludes$  assures that severity calculation for states in phase  $h$  is only based on the interval  $h.s$ .

Based on this analysis, we call a factor set  $F \leq_m$ -compatible if the way of combining the severity intervals for each subset  $F' \subseteq F$  fulfils Condition 1. For support in achieving and maintaining compatibility, the next paragraphs exemplify how factor severity and dependencies play together with mitigation orders. We also explore conditions for the well-formedness of severity intervals.

**Factor characteristics and dependencies** Constraints over risk spaces have implications on the characteristics of risk factors such as their severity intervals. Conversely, severity intervals govern  $\leq_m$  and can impose well-formedness conditions on the definition of constraints. In the following, we explore the impact of two different well-formedness conditions ( $f.s \subseteq g.s$  and  $f.s \supseteq g.s$ ) on  $\leq_m$  and the influence of choosing either one on the monotonicity of mitigation strategies. For example, for a constraint

$$f \text{ causes } g, \quad \text{the condition } f.s \subseteq g.s, \quad \text{and for} \quad (12)$$

$$F \text{ causes } G, \quad \text{the condition } \bigcup_{f \in F} f.s \subseteq \bigcup_{g \in G} g.s \quad (13)$$

seems coherent with the (probabilistic) relationship between the activation of causal and consequential factors. Particularly, the condition in Formula (12) expresses that if  $f$  causes  $g$  then  $f$  should have at most the potential impacts of  $g$ .

Consider the airbag spurious trip ( $st_{rA}$  with  $s = [5, 12]$ ) causing loss of driving control ( $loc_d$  with  $s = [2, 17]$ ) from the example in Table 2. For  $st_{rA}$  causes  $loc_d$ , the two intervals would fulfil the aforementioned condition. The structure fulfilling this constraint is shown in Fig. 8a, it includes the states  $\mathbf{0}$  with  $S(\mathbf{0}) = [0, 0]$ ,  $loc_d$  with  $S(loc_d) = [2, 17]$ , and  $loc_d st_{rA}$  with  $S(loc_d st_{rA}) = [2, 17]$ , with three equivalence classes in  $R(\{st_{rA}, loc_d\})/\equiv_m$  ordered according to

$$loc_d \overline{st_{rA}}, loc_d st_{rA}, loc_d <_m \overline{loc_d st_{rA}}, st_{rA} <_m \overline{loc_d}, \overline{st_{rA}}, \mathbf{0}, \overline{loc_d st_{rA}}.$$

On the contrary, assume a risk structure where  $st_{rA}.s = [2, 17)$  and  $loc_d.s = [5, 12)$  in order to be compliant with the converse condition  $f.s \supseteq g.s$ . From this structure, we obtain the order

$$\overline{loc_d st_{rA}}, loc_d st_{rA}, st_{rA} <_m \overline{loc_d st_{rA}}, loc_d <_m \overline{loc_d}, st_{rA}, \mathbf{0}, \overline{loc_d st_{rA}}.$$

In both orders, mitigation paths (i.e., sequences of solid arcs) in Fig. 8a, for example,

$$m^{loc_d} \rightarrow m^{st_{rA}} \rightarrow m_r^{st_{rA}} \rightarrow m_r^{loc_d} \rightarrow SKIP,$$

remain in their equivalence class or lead to a state in a better class according to  $\leq_m$ . Mitigations (solid arcs) deactivate at least one factor and, thus, lead to a better state according to  $<_{\tilde{m}}$ . We call such paths through  $\mathfrak{R}$  *mitigation monotonous* [GK17, Def. 8]. However, if intervals are specified in coherence with Formula (12), seven instead of five out of the twelve mitigations *strictly* reduce risk.

Analogously, for a constraint  $f$  requires  $g$ , the condition  $f.s \supseteq g.s$  seems useful. The two factors in our example fulfil this condition for  $loc_d$  requires  $st_{rA}$ . The structure fulfilling this constraint is shown in Fig. 8b, with the states ordered according to

$$loc_d st_{rA} <_m \overline{loc_d st_{rA}}, st_{rA} <_m \overline{st_{rA}}, \mathbf{0}, \overline{loc_d st_{rA}},$$

whereas intervals specified the opposite way (following  $f.s \subseteq g.s$ ) would lead to a less distinctive order

$$\overline{loc_d st_{rA}}, loc_d st_{rA}, st_{rA} <_m \overline{st_{rA}}, \mathbf{0}, \overline{loc_d st_{rA}}.$$

Again, mitigation paths are monotonous and coherent with the above condition; four instead of three out of seven mitigations *strictly* reduce risk according to  $\leq_m$ . Overall, a violation of the mentioned conditions does not influence mitigation monotonicity of this risk structure with respect to  $\leq_m$  but it reduces non-determinism in a strategy over  $\mathfrak{R}$ .

Figure 8c shows a risk graph for the discussed causes/requires constraint pair fulfilling both conditions.

Note that the states are equivalent to Fig. 8b but the arc  $\mathbf{0} \xrightarrow{e^{st_{rA}}} st_{rA}$  in Fig. 8b is replaced by  $\mathbf{0} \xrightarrow{e^{loc_d st_{rA}}} st_{rA} loc_d$  and the edge  $st_{rA} loc_d \xrightarrow{m_s^{st_{rA}}, m_s^{loc_d}} \overline{st_{rA} loc_d}$  is added. In summary, the causes constraint streamlines the transition relation of the risk structure as opposed to the requires constraint. This can be useful for merging factors during reduction of a risk structure. However, Fig. 8a shows that causes is less restrictive than requires, particularly, allowing  $e^{loc_d}$  to occur without a previous  $e^{st_{rA}}$ , for example, due to other causes not made explicit as factors in the model.

An in-depth analysis of techniques for making  $F$  compatible, a discussion of a set of rules relating factor characteristics and factor dependencies, as well as an alternative to  $\leq_m$  are subject of future work.

## 7. RISKSTRUCTURES: the language

In the previous sections, we have explored how a risk space  $R$  can be used for assessing *risk mitigation capabilities* of a process  $P$  (cf. Fig. 1) and for equipping  $P$  with a form of *risk awareness*. Risk in specific situations has a *specific causal structure*. We have seen how *dependencies* in a factor set  $F$  focus on this structure. Starting from the full extension of  $R$ , we *select subsets* of  $R$ , making assumptions if the structure is not entirely known in advance. This focus and these assumptions can be expressed by a constrained *composition*. Specifically, based on *constraints* on the combinations of factor phases, on the phase transitions, and on the *synchronisation* of events corresponding to the CSP model of concurrency, we can specify

- which region of  $R$  we want to pay attention to (i.e., the scope of safety guarantees) and
- which region of  $R$  we consider to be safe for a process  $P$  (i.e., conventional safety).

We use LTSs for investigating the operational semantics of RISKSTRUCTURES, beginning with factors and their constrained composition, and accompanied by a discussion of the *consistency, well-formedness, and validity* of their algebraic semantics.

**Definition 14 (Risk structure)** Given a factor set  $F \subseteq \mathcal{F}$ , a *risk structure*  $\mathfrak{R}$  is an expression of the form

$$\mathfrak{R} ::= p \mid \mathfrak{R} \parallel \mathfrak{R} \mid [\mathfrak{R}]_C$$

where, for a factor  $f \in F$  (Definition 3), phase  $p \in Ph_f$  (e.g.  $0^f, f, \bar{f}$ ) is a risk structure and  $C \subseteq \mathcal{C}$  is a set of *constraints* (Definition 12). Let  $\mathcal{S}$  be the set of all risk structures, including all phases  $\bigcup_{f \in \mathcal{F}} Ph_f \subset \mathcal{S}$ .

The binary operator  $\parallel$  signifies the parallel composition of two structures, and the binary operator  $[\cdot]_C$  applies all constraints in the set  $C$  to a risk structure.

**Operational semantics** Let the map  $\llbracket \cdot \rrbracket_r : \mathcal{S} \rightarrow \mathcal{T}$  uniquely assign a LTS to a CSP process corresponding to a risk structure. With  $\mathfrak{R} \in \mathcal{S}$  we associate a LTS  $\llbracket \mathfrak{R} \rrbracket_r = (R, \Sigma^u, \rightarrow, R_0)$  with the risk space  $R$  (Definition 5), the abstracted used alphabet  $\Sigma^u \subset 2^{\Sigma^r} \setminus \{\emptyset\}$  (Definition 1, Sects. 7.1, 7.2, 7.3), the transition relation  $\rightarrow \subseteq R \times \Sigma^u \times R$  (Definitions 15 and 16, Sect. 7.3), and a set of initial states  $R_0 \subseteq R$ . The operational semantics and algebraic properties of the operators and constructs of the language in Definition 14 are provided below. The construction of the used alphabet  $\Sigma^u$  will be explained along with the operators.

## 7.1. Risk factors as sequential processes

A risk factor  $f \in F$  (Figs. 2 and 3) can be represented as a sequential, mutually recursive process  $\mathfrak{R}_f$  according to Definition 1. We express  $f$  as

$$\begin{aligned} \mathfrak{R}_f &= 0^f && \text{(init)} \\ 0^f &= ?x : e^f \rightarrow f \square ?x : o_n^f \rightarrow 0^f && \text{(inactive)} \\ f &= ?x : m_d^f \rightarrow 0^f \square ?x : o_e^f \rightarrow f \square ?x : m^f \rightarrow \bar{f} && \text{(active)} \\ \bar{f} &= ?x : \bar{e}^f \rightarrow f \square ?x : m_r^f \rightarrow 0^f \square ?x : o_m^f \rightarrow \bar{f}. && \text{(mitigated)} \end{aligned}$$

Note the use of  $f$  as a symbol for the risk factor as a transition system (Fig. 2a) and the use of  $f$  for the CSP process that models this factor in its active phase. Different fonts signify the semantic difference. Also notice that  $p \xrightarrow{\tau}$  for any  $p \in Ph_f$  of a non-deterministic factor (Fig. 4) corresponds to the  $\square$ -fraction in  $\square$ .

This mapping enables an algebraic treatment of factors in CSP.

The semantics of a single factor  $f$  in phase  $p \in Ph_f$  is given by  $\llbracket p \rrbracket_r = (R(\{f\}), \Sigma^f, \rightarrow_f, \{\{f \mapsto p\}\})$  and the elements of this tuple are given by Definition 3 and described in Fig. 2. For example,  $\llbracket 0^f \rrbracket_r = (R(\{f\}), \Sigma^f, \rightarrow_f, \{\{f \mapsto 0^f\}\})$  provides the semantics of  $f$  initialised with the phase  $0^f$ . For each factor, we assume a phase order  $\leq_f$  as given in Sect. 4.1.

**Factor alphabets for global input-enabledness** Based on the definitions in Sect. 4.1, the condition

$$\forall f, g \in F : \left( \bigcup_{e \in \Sigma^f} e \right) \setminus \{\tau\} = \left( \bigcup_{e' \in \Sigma^g} e' \right) \setminus \{\tau\} = \Sigma \quad (14)$$

requires all factors to support the alphabet  $\Sigma$  of the process  $P$ . This choice does not affect risk space composition (Sect. 4.2) but allows factors to be composed in parallel synchronously and be ready to agree with some event  $P$  offers. This attempt to avoid interleaving requires input-enabledness of  $\mathfrak{R}$  (Sect. 4.1). Given  $F$  in terms of a factor family  $(f_i)_{i \in [1..n]}$  with  $n \in \mathbb{N}$ , one way to satisfy Formula (14) is a structured alphabet with compound (i.e., multi-part) events

$$r.\alpha_1.\alpha_2 \dots \alpha_n$$

with  $\alpha_1.\alpha_2 \dots \alpha_n \in \Sigma^{f_1}.\Sigma^{f_2} \dots \Sigma^{f_n}$  and a central channel  $r$  of type  $\Sigma^{f_1}.\Sigma^{f_2} \dots \Sigma^{f_n}$  through which  $\mathfrak{R}$  can observe and influence  $P$ . For  $i, j \in [1..n]$  with  $i \neq j$ , two (partially defined) factor event sets, say an endangerment  $e^{f_i} \in \Sigma^{f_i}$  activating factor  $f_i$  and a mitigation  $m^{f_j} \in \Sigma^{f_j}$  of factor  $f_j$ , then correspond to the set of (fully defined) compound events given by

$$events_{\mathfrak{R}}(r?e^{f_i}!m^{f_j}) = \{r.\alpha_1 \dots ?e^{f_i} \dots !m^{f_j} \dots \alpha_n \mid \alpha_1 \in \Sigma^{f_1}, \dots, e^{f_i} \in \Sigma^{f_i}, \dots, m^{f_j} \in \Sigma^{f_j}, \dots, \alpha_n \in \Sigma^{f_n}\}.$$

At the level of  $P$ , these compound events are equivalently given by

$$events_P(r?e^{f_i}!m^{f_j}) = \{r.x_1 \dots x_i \dots x_j \dots x_n \mid x_1 \in \bigcup_{e \in \Sigma^{f_1}} e, \dots, x_i \in e^{f_i}, \dots, x_j \in m^{f_j}, \dots, x_n \in \bigcup_{e \in \Sigma^{f_n}} e\}.$$

This channel type is inspired by suggestions in [Ros10, Sec. 1.2.4] and by [BS01, Sec. 3.4]. Thus, each factor  $f$  synchronises with the projection from  $\Sigma$  into  $\Sigma^f$  and is always ready to agree with some event offered by



$P$ .<sup>21</sup> This construction models partial observability of  $P$  through factors implemented by sensors and actuators and, for a deterministic factor set  $F$ , guarantees deadlock freedom of  $P \parallel \mathfrak{R}$ , modelling that a corresponding implementation of  $F$  can always monitor and influence  $P$ .

## 7.2. Composition

Let  $A, B, F \subseteq \mathcal{F}$  be factor sets with  $A, B \subseteq F$  and, for  $i \in \{1, 2, 3\}$ ,  $\mathfrak{R}_i \in \mathcal{S}$  be structures according to Definition 14 with  $\llbracket \mathfrak{R}_i \rrbracket_r = (R_i, \Sigma_i^u, \rightarrow_i, R_{i,0})$ . Now, we may want to combine two structures, say  $\mathfrak{R}_1$  and  $\mathfrak{R}_2$ , with possibly intersecting factor sets, for example, if two safety engineers are trusted with two HazOps of the same machine. Let the map  $scope: \mathcal{S} \rightarrow 2^{\mathcal{F}}$  identify the factors referred to by a structure. Then,  $scope(\mathfrak{R}_1) \cap scope(\mathfrak{R}_2)$  indicates all shared factors.

Because a single factor cannot be in two different phases at the same time, we apply Definition 6 for  $F_1 = scope(\mathfrak{R}_1)$  and  $F_2 = scope(\mathfrak{R}_2)$  to uniquely define the transition relation resulting from parallel composition: *only those states can be combined whose factors in the shared scope are in their identical phases.*

**Definition 15 (Parallel composition)** We define  $\llbracket \mathfrak{R}_1 \parallel \mathfrak{R}_2 \rrbracket_r = (R, \Sigma^u, \rightarrow, R_0)$  with  $R = R_1 \otimes R_2$  (see Definition 7 and Lemma 1) and  $R_0 = R_{1,0} \otimes R_{2,0}$ . For states  $\sigma_1, \sigma'_1 \in R_1, \sigma_2, \sigma'_2 \in R_2$ , the composed transition relation  $\rightarrow$  is given by the *left, right, and synchronous* step rules

$$\frac{\sigma_1 \xrightarrow{e} \sigma'_1 \quad \sigma_2 \xrightarrow{f} \sigma'_2}{\sigma_1 \cup \sigma_2 \xrightarrow{e \setminus f} \sigma'_1 \cup \sigma'_2} [\sigma_1 \approx \sigma_2, \sigma'_1 \approx \sigma'_2, e \neq f], \quad (\parallel\text{-l-step})$$

$$\frac{\sigma_1 \xrightarrow{e} \sigma'_1 \quad \sigma_2 \xrightarrow{f} \sigma'_2}{\sigma_1 \cup \sigma_2 \xrightarrow{f \setminus e} \sigma_1 \cup \sigma'_2} [\sigma_1 \approx \sigma_2, \sigma_1 \approx \sigma'_2, e \neq f], \text{ and} \quad (\parallel\text{-r-step})$$

$$\frac{\sigma_1 \xrightarrow{e} \sigma'_1 \quad \sigma_2 \xrightarrow{f} \sigma'_2}{\sigma_1 \cup \sigma_2 \xrightarrow{e \cap f} \sigma'_1 \cup \sigma'_2} [\sigma_1 \approx \sigma_2, \sigma'_1 \approx \sigma'_2, e \cap f \neq \emptyset]. \quad (\parallel\text{-s-step})$$

Based on  $\rightarrow$ , the *used alphabet* is  $\Sigma^u = \{e \in 2^{\Sigma^r} \mid \exists \sigma, \sigma' \in R: \sigma \xrightarrow{e} \sigma'\}$ .

Note that  $\parallel\text{-l-step}$  covers the special case  $f = \emptyset$  corresponding to  $\sigma_2 \rightarrow_2$ .  $\parallel\text{-r-step}$  works analogously. These step rules together resemble the step law of generalised parallel composition in CSP [Ros10, Sec. 3.4].

**Algebraic properties of composition** In the following, we show some desirable properties of the composition of two or more risk structures.

**Lemma 8 (Idempotency of  $\parallel$  under determinism)** For any *deterministic*  $\mathfrak{R}_1 \in \mathcal{S}$ , we have

$$\mathfrak{R}_1 \parallel \mathfrak{R}_1 = \mathfrak{R}_1. \quad (\parallel\text{-idem})$$

*Proof of Lemma 8.* From Definition 5, by Lemma 1, we obtain  $R_1 \otimes R_1 = R_1$ . Based on Sect. 4.2 and the alphabet of  $\mathfrak{R}_r$  defined in Sect. 7.1, we can apply the rules from Definition 15 as follows.

Because of the uniqueness of  $\mathfrak{R}_1$  and the compatibility requirement (Formula (2)), the composition takes two consistent views of  $\mathfrak{R}_1$  offering identical source and target states ( $\sigma_1 = \sigma_2$ ) at each step. By definition, all factors share the alphabet  $\Sigma$  and  $\mathfrak{R}_1$  is deterministic. Then, any two event sets  $e, f$  in  $initials(\mathfrak{R}_1)$  (i.e., in a state  $\sigma$ ) are either disjoint or equal. Hence, there is no way for a concrete event in  $P$  to be a member of more than one such event set. In line with Definition 3,  $\parallel\text{-l-step}$  and  $\parallel\text{-r-step}$  do not apply.  $\parallel\text{-s-step}$  turns into a tautology and both views exhibit an identical transition:

$$\frac{\sigma_1 \xrightarrow{e} \sigma'_1}{\sigma_1 \cup \sigma_1 \xrightarrow{e} \sigma'_1 \cup \sigma'_1}.$$

<sup>21</sup>According to the CSP notation for channel events, we have  $\Sigma = \llbracket r \rrbracket = events(r)$  for the channel  $r$ .

For  $\sigma_1, \sigma'_1 \in R_1$  and  $e \in \Sigma_1^u$ , this tautology preserves  $\Sigma_1^u, \rightarrow_1$ , and  $R_{1,0}$ .  $\square$   
 The composition of a deterministic  $\mathfrak{R}_1$  with itself turns into a synchronous (i.e., non-interleaved) composition, essentially, a sequence of external choices where both  $\mathfrak{R}_1$ s synchronously agree with the environment.

**Lemma 9 (Commutativity of  $\parallel$ )** For any  $\mathfrak{R}_1, \mathfrak{R}_2 \in \mathcal{S}$ , we have

$$\mathfrak{R}_1 \parallel \mathfrak{R}_2 = \mathfrak{R}_2 \parallel \mathfrak{R}_1 . \quad (\parallel\text{-comm})$$

*Proof of Lemma 9.* Lemma 1 in Sect. 4.2 makes it easy to show  $R_1 \otimes R_2 = R(\text{scope}(\mathfrak{R}_1) \cup \text{scope}(\mathfrak{R}_2)) = R_2 \otimes R_1$ . Furthermore, the symmetric duals (i.e.,  $\parallel\text{-r-step}$  is the dual of  $\parallel\text{-l-step}$ ) of the rules in Definition 15 yield the same  $\rightarrow$  and, hence, the same  $\Sigma^u$ .  $\square$

**Lemma 10 (Associativity of  $\parallel$ )** For any  $\mathfrak{R}_1, \mathfrak{R}_2, \mathfrak{R}_3 \in \mathcal{S}$  with *disjoint* scopes, that is,

$$\text{scope}(\mathfrak{R}_1) \cap \text{scope}(\mathfrak{R}_2) \cap \text{scope}(\mathfrak{R}_3) = \emptyset,$$

we have

$$(\mathfrak{R}_1 \parallel \mathfrak{R}_2) \parallel \mathfrak{R}_3 = \mathfrak{R}_1 \parallel (\mathfrak{R}_2 \parallel \mathfrak{R}_3) . \quad (\parallel\text{-assoc-1})$$

*Proof sketch for Lemma 10.* Because of empty shared scopes, the side conditions always hold on both sides and states are merged by disjoint union. Set operations on states and events are commutative and associative. By definition,<sup>22</sup> all factors (Definition 3) share the alphabet  $\Sigma$ , that is, they synchronise on all events (Sect. 4.2). Thus, interleaving is avoided, the alphabetised parallel operator takes  $\Sigma$  on both sides and, this way, reduces to synchronous parallel composition for which a general associative law is available [Ros10, p. 60].  $\square$   
 Lemma 10 allows the use of  $\parallel^{f \in F} 0^f$  as a shortcut for  $((\dots (0^{f_1} \parallel 0^{f_2}) \parallel \dots) \parallel 0^{f_n})$  with  $f_i \in F$ .

**Lemma 11 (Associativity of  $\parallel$ )** For any  $\mathfrak{R}_1, \mathfrak{R}_2, \mathfrak{R}_3 \in \mathcal{S}$  with *equal* scopes, that is,

$$\text{scope}(\mathfrak{R}_1) = \text{scope}(\mathfrak{R}_2) = \text{scope}(\mathfrak{R}_3) ,$$

we have

$$(\mathfrak{R}_1 \parallel \mathfrak{R}_2) \parallel \mathfrak{R}_3 = \mathfrak{R}_1 \parallel (\mathfrak{R}_2 \parallel \mathfrak{R}_3) . \quad (\parallel\text{-assoc-2})$$

*Proof sketch of Lemma 11.* The proof is analogous to the proof of Lemma 10.  $\square$

### 7.3. Constraints

In Sect. 6, we specified relations over risk spaces to determine  $\rightarrow$  by pruning  $R \times R$ . We now use constraints as a construct in the language of RISKSTRUCTURES (Definition 14). The redundancy coming with this construct can be used to identify inconsistencies between  $\mathfrak{R}$  and the real world, potentially helpful in *model refinement*, *completion*, and *validation* [Gle14]. Particularly, these inconsistencies allow choices for their resolution. We can make such inconsistencies explicit as follows.

For factors  $F \subseteq \mathcal{F}$ , constraints  $C \subseteq \mathcal{C}$ , and a risk state  $\sigma \in R(F)$ , we call

$$\mathfrak{R}_\sigma^C = [\parallel^{f \in F} \sigma(f)]_C$$

with  $R_0 = \{\sigma\}$  the *characteristic risk structure* of  $\sigma$  under constraints  $C$ . *But which definition of  $\sigma \rightarrow$  is most useful?* To investigate this, we split the determination of  $\sigma \rightarrow$  for  $\mathfrak{R}_\sigma^C$  into three fractions  $\sigma \rightarrow_{c,f,p}$  with

$$\begin{aligned} \sigma \rightarrow_c & \text{ derived from } C, \text{ with} & \sigma \xrightarrow{\tau_c} \sigma' \in \sigma \rightarrow_c & \iff (\sigma, \sigma') \in \llbracket C \rrbracket_c^R, \\ \sigma \rightarrow_f & \text{ from Definitions 3 and 15, with} & \sigma \xrightarrow{e} \sigma' \in \sigma \rightarrow_f & \iff \sigma \xrightarrow{e} \sigma', \text{ and} \\ \sigma \rightarrow_p & \text{ from process } P, \text{ with} & \sigma \xrightarrow{e} \sigma' \in \sigma \rightarrow_p & \iff e \in \text{initials}(P) \wedge \sigma' \in \text{reach}(\sigma, P) . \end{aligned}$$

These fractions give rise to the following *inconsistencies*:

1. The relation  $\sigma \rightarrow_c \setminus \sigma \rightarrow_f$  describes *sensible* transitions with *invisible events* signified by  $\tau_c$ . We choose to prune transitions from  $\rightarrow$  if they deviate from what is provided in  $\rightarrow_f$ .

<sup>22</sup>We provide the rationale for this definition in Sect. 7.5.

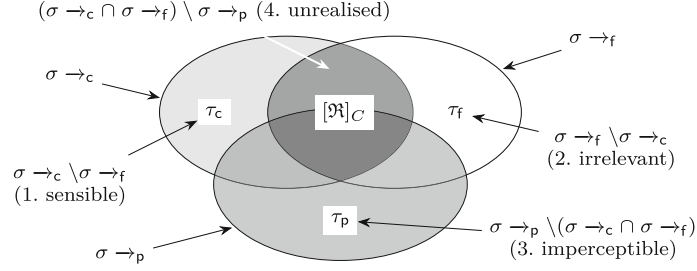


Fig. 9. Semantics of  $\sigma \rightarrow$  for the constrained form  $[\mathfrak{R}]_C$  according to Definition 16

2. The relation  $\sigma \rightarrow_f \setminus \sigma \rightarrow_c$  describes *irrelevant* transitions labelled with a  $\tau_f$ . We choose to prune transitions from  $\rightarrow$  if they violate constraints and, hence, lead to inconsistencies in  $\mathfrak{R}$ .
3. The relation  $\sigma \rightarrow_p \setminus (\sigma \rightarrow_c \cup \sigma \rightarrow_f)$  describes *imperceptible* transitions. In  $\rightarrow$ , we could label such transitions with a  $\tau_p$ , making them subject of *process-driven disclosure* of  $\mathfrak{R}$ .
4. The relation  $(\sigma \rightarrow_c \cap \sigma \rightarrow_f) \setminus \sigma \rightarrow_p$  describes *unrealised* transitions. We could prune such transitions from  $\rightarrow$ , because they are not realised in  $P$  and, hence, would only add little value to  $\mathfrak{R}$ .

This case analysis is visualised in Fig. 9 and suggests several possibilities to design a semantics for constraints. As indicated in the case list, one would desire the following possibility for  $\rightarrow$ :

$$\begin{aligned} \sigma \rightarrow &= (\sigma \rightarrow_c \cup \sigma \rightarrow_f \cup \sigma \rightarrow_p) \setminus \left( \underbrace{(\sigma \rightarrow_c \setminus \sigma \rightarrow_f)}_{1. \text{ reduce to } F} \cup \underbrace{(\sigma \rightarrow_f \setminus \sigma \rightarrow_c)}_{2. \text{ reduce to } C} \cup \underbrace{(\sigma \rightarrow_c \cup \sigma \rightarrow_f \setminus \sigma \rightarrow_p)}_{4. \text{ reduce to } P} \right) \\ &= \sigma \rightarrow_c \cap \sigma \rightarrow_f \cap \sigma \rightarrow_p. \end{aligned} \quad (\text{by set-algebraic manipulation})$$

**Corollary 6** For any characteristic structure  $\mathfrak{R}_\sigma^C$ , we have  $\sigma \rightarrow_c \setminus \sigma \rightarrow_f = \emptyset$  because of the input-enabledness [Tre08] of factors and, thus,  $\sigma \rightarrow = \sigma \rightarrow_c \cap \sigma \rightarrow_p$ .

To weaken the dependency of  $\mathfrak{R}$  on a possibly partially known process, we apply the following definition.

**Definition 16 (Constraint)** Let  $\mathfrak{R}$  be a risk structure (Definition 14) with  $[\mathfrak{R}]_r = (R, \Sigma^u, \rightarrow, R_0)$  and  $C$  a set of constraints well-formed for  $R$  (Sect. 6, Definition 13). The *constrained form*  $[\mathfrak{R}]_C$  is defined by  $[[[\mathfrak{R}]_C]]_r = (R, \Sigma_C^u, \rightarrow_C, R_0)$ .  $\rightarrow_C$  is determined by the following two rules. For  $\sigma, \sigma' \in R$

$$\frac{\sigma \xrightarrow{e} \sigma' \quad (\sigma, \sigma') \in [[C]]_c^R}{\sigma \xrightarrow{e}_C \sigma'} \quad \text{and} \quad ([\cdot]_C\text{-step})$$

$$\frac{\sigma \xrightarrow{e} \sigma' \quad (\sigma, \sigma') \notin [[C]]_c^R}{\sigma \xrightarrow{\tau}_C \sigma} \quad ([\cdot]_C\text{-cancel})$$

Based on  $\rightarrow_C$ ,  $\Sigma_C^u = \{e \in 2^{\Sigma^r} \mid \exists \sigma, \sigma' \in R: \sigma \xrightarrow{e}_C \sigma'\}$ .

For arbitrary structures, the  $[\cdot]_C\text{-step}$  rule implements  $\rightarrow = \rightarrow_c \cap \rightarrow_f$  (i.e., cases 1 and 2). The reduction of unrealised transitions and the coverage of imperceptible transitions (i.e., cases 3 and 4) could be addressed by an incremental construction of  $\mathfrak{R}$ .  $[\cdot]_C\text{-cancel}$  produces  $\tau_f$  self-loops to preserve input-enabledness.<sup>23</sup> Notice that  $\tau_f$ -loops remain whether or not they are contained in further constraints applied to  $\mathfrak{R}$ , and also that Corollary 6 and Definition 16 prevent constraints such as causes from producing  $\tau_c$ . Overall, the rules  $\|\text{-l-step}$ ,  $\|\text{-r-step}$ ,  $\|\text{-s-step}$ ,  $[\cdot]_C\text{-step}$ , and  $[\cdot]_C\text{-cancel}$  determine the operational semantics of a risk structure as an LTS and, together with *reach* (Sect. 5.4), form a basis for automated reasoning about RISKSTRUCTURES. We demonstrate such automation with our tool prototype YAP [Gle21].

**Remark 2** The meaning of *f* causes *g*, specified as a relation in Sect. 6.1, can now be investigated in the LTS semantics of  $\mathfrak{R}$ . On a particular *path* through  $R$  according to  $\rightarrow_C$  from a state in  $R_0$ , whenever *f* gets active, *g* must

<sup>23</sup>Operationally,  $[\cdot]_C\text{-cancel}$  acts like a *SKIP*, non-blocking *STOP*, or a hiding of trace suffixes. Hence, we define  $\tau_f = \surd$ .

get active in the same or the following state. This way, *causes* constraints allow *immediate or weak causation* and requires constraints *delayed or strong causation* at most one logical step apart in  $\mathfrak{R}$ . The weak form of *causes* enforces the *simultaneous* occurrence of  $f$  and  $g$ , it *synchronises* the activation of  $g$  with the prior activation of  $f$ . Moreover,  $g$  can stay active forever and may already be active before the activation of  $f$ . Consider that  $g$  might have been activated by a factor or a causal relationship not (yet) captured in  $\mathfrak{R}$ . Such activation is by default possible in  $\mathfrak{R}$  once  $g$  is in the factor set.

**Algebraic properties of constraints**  $([\cdot]_C)$ . Let  $C_i \subseteq C$  for  $i \in \{1, 2, 3\}$ . We have that  $[\mathfrak{R}]_\emptyset = \mathfrak{R}$  by Definition 12. The following lemma establishes that the order in which single constraints are applied to  $\mathfrak{R}$  does not matter.

**Lemma 12 (Exchange of  $[\cdot]_C$  and  $\cup$ )**

$$[[\mathfrak{R}]_{C_1}]_{C_2} = [\mathfrak{R}]_{C_1 \cup C_2} \quad (15)$$

*Proof sketch.* The proof is by induction over  $C_1$ , supported by an additional lemma for the induction step, and takes advantage of associativity of set intersection as used in Definition 12. The detailed proof is stated in Appendix A.  $\square$

**Lemma 13 (Idempotency)**

$$[[\mathfrak{R}]_C]_C = [\mathfrak{R}]_C \quad ([\cdot]_C\text{-idem})$$

*Proof sketch.* This lemma follows from Lemma 12 and idempotency of set union.  $\square$

**Lemma 14 (Commutativity)**

$$[[\mathfrak{R}]_{C_1}]_{C_2} = [[\mathfrak{R}]_{C_2}]_{C_1} \quad ([\cdot]_C\text{-comm})$$

*Proof sketch.* This lemma follows from Lemma 12 and commutativity of set union.  $\square$

**Lemma 15 (Associativity)**

$$[[\mathfrak{R}]_{C_1 \cup C_2}]_{C_3} = [[\mathfrak{R}]_{C_1}]_{C_2 \cup C_3} \quad ([\cdot]_C\text{-assoc})$$

*Proof sketch.* This lemma follows from Lemma 12 and associativity of set union.  $\square$

For an arbitrary  $C \subseteq C$ , we relax Definition 16 by establishing following equivalence

$$[\mathfrak{R}]_C = [\mathfrak{R}]_{C'} \quad (16)$$

for the largest  $C' \subseteq C$  well-formed for  $R$  (Definition 13). Then,  $[\mathfrak{R}]_C$  denotes the structure resulting from applying only and exactly the constraints in  $C'$ . Moreover, Definition 16 complements Definition 15 by guarding the  $\|\text{-l-step}$ ,  $\|\text{-r-step}$ , and  $\|\text{-s-step}$  rules. In this way, the application of a constraint  $c \in C$  to a risk structure  $\mathfrak{R}$  restricts its transition relation  $\rightarrow$ . Definition 16, together with the relational pruning for three out of the four mentioned cases (cf. Fig. 9), yields the following refinement law.

**Lemma 16 (Constraints trace-refine)**

$$\mathfrak{R} \sqsubseteq_T [\mathfrak{R}]_C$$

*Proof sketch.* The proof is by showing that  $\text{traces}(\mathfrak{R}) \supseteq \text{traces}([\mathfrak{R}]_C)$ . Fix  $t \in \text{traces}([\mathfrak{R}]_C)$  with  $t = f \wedge l$  for induction over the split of  $t$  into  $f$  and  $l$ . Induction step: Assume that  $f \in \text{traces}(\mathfrak{R})$  as *induction hypothesis* (IH). With  $l = \langle e \rangle \wedge l'$  and  $t = f \wedge \langle e \rangle \wedge l'$ , there exist  $\sigma, \sigma' \in R$  such that  $\sigma \xrightarrow{e}_C \sigma'$  and, according to rule  $[\cdot]_C\text{-step}$ , such that  $e \in \Sigma^u$ ,  $\sigma \xrightarrow{e} \sigma'$ , and  $(\sigma, \sigma') \in \llbracket C \rrbracket_C^R$ . Because of IH and  $\sigma \xrightarrow{e} \sigma'$ , there must be a trace  $f \wedge \langle e \rangle \wedge l'' \in \text{traces}(\mathfrak{R})$  and, because of downward closure, a trace  $f \wedge \langle e \rangle$ .  $e$  is part of  $l$  and, therefore,  $t$  such that we complete the induction step and establish the new IH  $f \wedge \langle e \rangle$ . (We do not need to prove the equivalence  $l' = l''$ .) The case  $f = \langle \rangle$  provides the induction start and  $l = \langle e \rangle$  terminates the induction.  $\square$

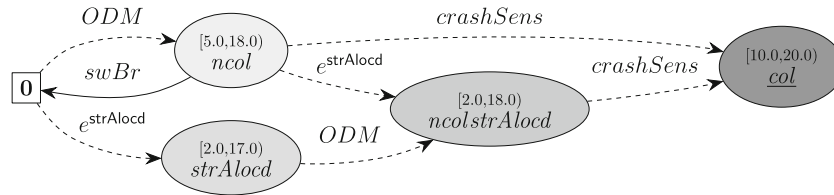
Following the pruning semantics described in Fig. 9, Lemma 16 establishes the desirable effect of constraints, the pruning of the transition relation of  $\mathfrak{R}$  according to the relational specification of the added constraints. However, an investigation of the conditions under which constraints are compositional, that is, when they can be exchanged with, for example,  $\|\text{-}$ , is left for future work.

Listing 2: YAP script modelling the situation manual from Table 2

```

1 factormodel for manual {
2   // single compound consequence on a single asset: impact on
   human driver
3   col desc "collision"
4   accident
5   detectedBy(.crashSens)
6   mitigatedBy(.rA)
7   requires(ncol)
8   sev=[10,20];
9
10  strA desc "airbag spurious trip"
11  causes(locd)
12  sev=[5,12];
13  locd desc "loss of driving control"
14  preventsMit(ncol,locd) // mitigation lock at loc
15  requires(strA)
16  sev=[2,17];
17  ncol desc "near collision"
18  detectedBy(.ODM)
19  mitigatedBy(PREVENT_CRASH.swBr)
20  direct
21  sev=[5,18);
22 }
23
24 controlloop vehOnRd {
25   mode ODM desc "object detection and distance measurement
   system"
26   guard "distToNearestObject < SBD" // distance To Nearest Object on
   Trajectory Is Less Than Safe Braking Distance
27   embodiedBy lidar;
28
29   mode SwBr desc "swerve and brake"
30   type PREVENT_CRASH
31   target (act=SwBr); // perform low-level procedure
32
33   mode crashSens desc "crash sensor"
34   update "shock=true"; // forward validated shock signal to
   airbag
35
36   mode rA desc "release driver and passenger airbags"
37   guard "shock=true"
38   update "release=DrvAndPsgrr"
39   embodiedBy airbag;
40
41   lidar;
42   airbag;
43 }

```

Fig. 10. Risk structure  $\mathfrak{R}_{\text{manual}}$  with the factor set  $F = \{\text{loc}_d, \text{ncol}, \text{str}_A, \text{col}\}$  generated by YAP from Listing 2

Finally, we discuss the special case that constraints are applied to a phase  $p \in Ph_f$  of a factor  $f$ . Observe that  $scope(p) = \{f\}$  and, by Definition 13,  $C$  is well-formed for  $R(scope(p))$  if  $C$  contains only constraints that refer to  $f$ . Moreover, if  $C$  does not contain self-referential constraints (Sect. 6.1), Formula (16) entails  $C' = \emptyset$  for the largest well-formed  $C' \subseteq C$ . From this observation, we obtain

$$[p]_C \quad \stackrel{\text{by Formula (16)}}{=} \underbrace{[p]_{C'}}_{\text{by Definition 13}} \quad \stackrel{\text{by Definition 13}}{=} \underbrace{[p]_{\emptyset}}_{\text{by Definition 12}} \quad \stackrel{\text{by Definition 12}}{=} p .$$

The effect of this last result is that single factors are not affected by constraints other than self-referential ones. Self-referential constraints such as `f direct` can, however, restrict  $p$  or convert  $p$  to `SKIP`. Often, such constraints can be avoided because they do not offer much expressiveness over bespoke factor models.

#### 7.4. Example: risk factors on the road (cont'd)

We develop the example in Sect. 3 and sects. 4.3 and 6.3, modelling instances of the factor set  $F = \{\text{loc}_d, \text{ncol}, \text{str}_A, \text{col}\}$  with corresponding events and compose these factors into the structure

$$\mathfrak{R}_{\text{manual}} = [((0^{\text{str}_A} \parallel 0^{\text{loc}_d}) \parallel 0^{\text{ncol}}) \parallel 0^{\text{col}}]_C$$

for the situation manual.  $C$  contains dependencies such as `locd preventsMit {ncol, locd}}` and `strA causes locd`. The full dependency list can be derived from Listing 2. Figure 10 shows the transition graph for  $\mathfrak{R}_{\text{manual}}$  after applying composition and constraints. The edges are labelled with factor events and the nodes with active factors. Listing 2 is written in YAP script, the input language for our tool YAP [Gle21], which allows one to encode a situation/action/factor table (Table 1a) and elaborate the example in Sect. 3 in several steps.

*Collision* (`col`) is the first phenomenon we model as a factor. With the `detectedBy` directive, we specify that a collision is detected by a crash sensor (`crashSens`) that signals the release of an airbag (`rA`). We associate an interval `col.s = [10, 20)` of 10 to 20 “negativity units” to a collision.

To illustrate *scoping* of a risk structure, we use the keyword `accident` to declare `col` as final. Mitigation of final factors is *out of scope*, so vehicle control would not take into account `crashSens` and `rA` in this case. However, in practice, `col` is not final because of being mitigated by an airbag. Capturing a step of accident analysis—technically, backwards causal reasoning—the dependency `col` requires `ncol` specifies that a near-collision is assumed to occur before any collision.

A *near-collision* (`ncol`) is detected by an object distance measurement system (ODM). The ODM is responsible for monitoring the predicate  $distToNOoT < SBD$ , which is true if the distance measured to the nearest object on the road within the planned vehicle trajectory ( $distToNOoT$ ) is smaller than the safe braking distance  $SBD$ . We assume an interval  $ncol.s = [5, 18)$  with a smaller upper bound than the one for `col` because the indirection to the accident reduces the possibility of bad consequences. Concretely, the likelihood of consequences from `col` is reduced by the action `swBr`. Based on the discussion in Sect. 4.4, we plan in future work to quantify this likelihood through probabilistic factor transitions.

The mitigatedBy(`PREVENT_CRASH.swBr`) directive specifies as a mitigation of `ncol` the *swerve-and-break* manoeuvre (`swBr`). The constraint `direct` specifies that `swBr` to be the event  $m_d^{ncol}$  (cf. Fig. 2a). Shown in the `controlloop` block, `swBr` is a low-level control mode to be elaborated in an external model (e.g. in MATLAB or its free alternative GNU/Octave<sup>24</sup>). `PREVENT_CRASH` is a classifier from which generic parameters for low-level controllers of this kind could be inherited. The modelling of modes is akin to the specification of guarded commands (i.e., using the guard and update attributes). We demonstrate in [GC20] how YAP modes can be used for controller synthesis, and in [FGC20] how factor LTSs can be modelled as hybrid automata, simulated in GNU/Octave, and formally verified in Isabelle/UTP.

By FMEA of the airbag, we identify its *spurious trip* (`strA`) and assume it to be *one of possibly several sufficient causes* of a *loss of driving control* (`locd`). This finding is captured by the two dependencies `strA causes locd` and `locd requires strA`. `strA` gets assigned the interval  $st_{rA}.s = [5, 12)$ .<sup>25</sup> We specify the interval  $loc_d.s = [2, 17)$ . Moreover, `locd` preventsMit {`ncol`, `locd`} captures the fact that after a loss of control, manual driving mode prevents the mitigation of `ncol`. Additionally, `locd` locks itself using a circularity so that, in  $\mathfrak{R}_{\text{manual}}$ , it cannot be mitigated.

From the severity intervals calculated for each risk state according to Formula (5) (see the state labels in Fig. 10), we establish

$$col <_m ncol =_m loc_d st_{rA} ncol <_m loc_d st_{rA} <_m \mathbf{0}.$$

Notice that, according to  $\leq_m$ , the state `ncol` is equivalent to `locdstrAncol`. This equivalence results from the equality of the upper bounds of both intervals [5, 18) and [2, 18). Although one might think that the state with the smaller lower bound (i.e., `locdstrAncol`) is better, consider that `ncol` has two active factors less. This circumstance gives rise to investigate refinements of  $\leq_m$ . In fact, our tool YAP [Gle21] implements a special case of  $\leq_m$  distinguishing pairs  $(\sigma, \sigma')$  of risk states that normally are  $\sigma =_m \sigma'$  because they only differ in the lower bounds of their severity intervals.

The state `F` with all factors activated (cf. Definition 11) is not reachable from `0`. Moreover, only one of the five relevant states in  $\mathfrak{R}_{\text{manual}}$  is covered by a mitigation `swBr` for `ncol`. In practice, one desires mitigations that cover many states or mitigate several factors at once. Although it is useful to statically assign severity intervals for analysis purposes, it is obvious that these intervals will in most applications of risk-aware machines have to be estimated online. From a particular state, online estimation allows one to trigger a mitigation based on a severity threshold that codifies an unacceptable level of risk. Moreover, a risk structure makes this estimation explicit and potentially simpler. Overall, the  $\mathfrak{R}_{\text{manual}}$  is a simple risk structure, yet showing many features of RISKSTRUCTURES. Figures 11a and 11b in Appendix B show a more complex structure for a set of six factors, generated by YAP in 66 milliseconds. The use of constraints leads to a focus on 24 out of the 4096 ( $4^6$  from using four-phase factors) risk states.

## 7.5. Discussion: compositional abstraction

**Event abstraction** Events in the CSP interpretation are atomic observations [Ros10, Ch. 1.5]. The initiation and termination of events representing complex enduring real-world phenomena are to be viewed as non-separable aspects of these events. Consequently, care is necessary when making assumptions about the atomicity of such

<sup>24</sup>GNU/Octave, <https://octave.sourceforge.io>.

<sup>25</sup>To keep the model simple for the sake of illustration, the upper bound 12 takes into account the possibility that the driver might be able to regain control after the airbag deflates.

events interfering with other events. Nevertheless, we consider the event sets in Fig. 2b and their embedding into compound events (Sect. 7.1) useful to model the *initiation, termination, or other significant events* of the corresponding sub-processes (i.e., endangerment and mitigation processes) of  $P$ .

Any implementation of a risk structure as a safety controller needs to make an assumption about how *shared resources*, such as actuators, are used. In CSP, two processes (e.g. two factors) cannot engage in two different events on the same channel at the same time. Hence, resource sharing has to be implemented outside a risk structure, creating two cases: (a) mitigation actions of two factors mapping onto the same resource simultaneously occur and (b) a mitigation action competes with an action outside the safety controller. In case (a), signal *priorities* or more sophisticated signal *merge functions* may determine the control inputs to be forwarded. Case (b) suggests the use of a *safety override*, that is, giving priority to the safety controller.

**Shared factor alphabet for compositional abstraction** The side conditions in the rules of  $\parallel$  (Definition 15) prohibit any behaviour leading to inconsistent states. The composition constrains the behaviour of two structures with an overlapping scope. Corollary 6 summarises our decisions on defining the constraint operator in such a way that Lemma 16 confirms this operator to form a refinement of a risk structure.

Were we to use arbitrary CSP processes for factors and were we to allow different interfaces for each use of the parallel composition operator, several structures, when composed, would not have guaranteed freedom from interference, exhibit different event and state traces and, consequently, the order of their composition would have led to different risk models. Associative laws for generalised parallel composition in CSP are not universally applicable. Discussions by Oliveira et al. [OCW09] and Roscoe [Ros10, p. 60] highlight how differences in the alphabets shared between each pair in a set of risk structures would entail meaning to the order in which these pairs are composed. For example, in CSP, the equality  $(P \parallel_X Q) \parallel_Y R = P \parallel_X (Q \parallel_Y R)$  holds generally only if  $X = Y$ .

In case of  $X \neq Y$ , one has to compare the traces of the composed processes to prove specific guarantees to be preserved by their composition. This can be computationally complex.

Overall, we obtain two advantages in RISKSTRUCTURES when using factors with an alphabet as described in Sect. 7.1. First, all factors have to always agree on their view of the same process  $P$ . Second, overlapping scopes in CSP terms then mean copies of factors that will be reduced by idempotency if deterministic (see Lemma 8). The transformation of factors into CSP (Sect. 7.1) encodes all the information from the risk space  $R$  into the processes  $0^f$ ,  $f$ , and  $\bar{f}$ . This way, we restrict the use of generalised parallel composition to the use of synchronous parallel composition [Ros10, p. 45].

**Risk graphs and Yap** The risk graphs in Figs. 7, 8 and 10 omit the transition self-loops for nominal, endangered, and mitigated operation because these are non-essential for endangerment and mitigation. Currently, YAP explores risk by changing (i.e., activating, mitigating) one factor at a time unless constraints such as *causes* are used. *causes* enforces the synchronisation of phase changes of involved factors. Because of input-enabledness from the compound event construction (Sect. 7.1), interleaving is avoided, but non-determinism can arise from non-deterministic factors. However, composition according to Definition 15 is more flexible because it allows several factors to change their phases simultaneously. This ability is exploited for the synthesis of safety controllers from YAP models in [GC20]. Furthermore, Definition 15 does not require input-enabledness but it relies on the alphabet defined in Sect. 7.1. Anyway, the use of a risk structure as a concurrent monitor or, more actively, as a safety controller, makes input-enabledness a useful requirement.

## 7.6. Discussion: uncertainty in RISKSTRUCTURES

Recall the notion of risk as an action with an *uncertain* and *undesired* outcome. As summarised in Fig. 1, a risk structure focuses on such outcomes.<sup>26</sup> Consequently, we propose risk structures as one way of formalising the irreducible (or aleatory) uncertainty associated with the actions of a risk-aware machine operated in a particular domain. Table 5 highlights *low-level* epistemic and aleatory uncertainties about a risk structure itself as a formal model of aleatory uncertainty in that domain, and as a model of a safety controller of that machine. Table 5 classifies these low-level uncertainties based on the discussion in [BFK13] and indicates how RISKSTRUCTURES incorporate and support the handling of these uncertainties.

<sup>26</sup>Desired or certain outcomes are not typically subject of operational risk analysis of autonomous machines.

**Table 5.** Summary of uncertainties inherent to and controllable in RISKSTRUCTURES

RISKSTRUCTURES can be uncertain in their ...	at design-time <i>epistemic</i> <sup>1</sup> (imprec.)	at run-time <i>aleatory</i> <sup>2</sup> (variab.)	<i>epistemic</i> <sup>3</sup> (meas.)	Mathematical approach	Can be implemented by ...
<i>Factor uncertainties</i>					
factor observation	yes	–	yes	subsetting	$e^f$ only covers domain subset
impact quantification	yes	yes	yes	intervals	severity bounds
phase changes	yes	–	yes	non-determinism	overlapping factor event sets
risk qualification	yes	–	–	partial order	partial phase order $\preceq_f$
<i>Structural uncertainties</i>					
analysis scope	yes	–	–	relations	factor dependencies/omission
risk qualification	yes	–	–	partial order	partial mitigation order $\lesssim_m$
risk quantification	yes	yes	yes	interval arithmetic	strong mitigation order $\leq_m$
<i>To be investigated in future work</i>					
phase changes	yes	yes	yes	MDPs	transition probabilities

Legend: <sup>1</sup>analytic imprecision, <sup>2</sup>physical parameter variability, <sup>3</sup>parameter measurement, – not applicable

**Non-deterministic and partial risk factors** With *non-deterministic factors* (Sect. 4.4) in  $\mathfrak{R}$ , from an observed phase/event,  $P$  can reach several distinct phases and, thus, risk states. Without further information,  $\mathfrak{R}$  ceases to know which of these states was actually reached.  $\mathfrak{R}$  would be in a *risk state region*, ordered and bounded by  $\min_{\leq_m} / \max_{\leq_m}$  (Sect. 5.4), suggesting to express  $\mathfrak{R}$ 's uncertainty by a *risk interval*. If factor phases carry information about  $P$ 's state in terms of disjoint state invariants, we can build state estimators into  $\mathfrak{R}$ , able to gradually restore lost state information, narrow the interval, and again uniquely identify  $P$ 's actual risk state. Following up on Sect. 4.4, if we can show that a risk structure is an Input/Output LTS and it can be transformed to a non-deterministic finite Mealy machine, the state estimation problem can be solved, for example, by homing algorithms used for passive testing [NSV03].

*Partially specified risk factors* (e.g. because of partial activation or mitigation analysis according to Tables 1b and 1c) may lead to *sensible transitions*  $\tau_c$  (cf. Sect. 7.3). Occurrences of  $\tau_c$  could be used to incrementally *learn a property enforcer* guided by unrealised events from risk factors.

## 8. Related work

Following up Sect. 2.1, we summarise research on dependability modelling for repairable systems, algebraic methods for risk assessment, run-time monitoring, and risk-aware planning and control. We discuss the role of RISKSTRUCTURES in the design of safety controllers.

**From dependability analysis to RISKSTRUCTURES** Leveson and Stolzy [LS87] investigate reachability graphs of timed Petri nets with failure places to assess *fault-tolerance* of safety-critical real-time systems. RISKSTRUCTURES form a projection and compositional generalisation of such graphs when risk factors are used to model faults. Based on the annotation of a component architecture with a fault model, Unanue et al. [UPM18] show the syntheseses of a failure automaton, a temporal fault/repair tree, and MCSs from this tree. From the failure automaton, they construct an extended Petri net to calculate failure probabilities to identify critical MCSs. Unanue et al.'s use of graphical models is intuitive and practical. A failure automaton corresponds to a risk structure with directly reducible risk factors modelling failure and repair. The commonalities of our framework with theirs suggests that RISKSTRUCTURES can use input from state-of-the-art dependability assessments of *repairable systems* for the design of safety controllers. At the moment, however, RISKSTRUCTURES provide no direct support for modelling probabilistic actions.

As discussed in Sect. 6, FTs (Sect. 2.1) can be translated into RISKSTRUCTURES using factor dependencies. Because of the step semantics of constraints (cf. Definition 12), such a translation is also possible for PAND and POR gates<sup>27</sup> used in dynamic FTs. This way, RISKSTRUCTURES can include a translation of FTs generated from

<sup>27</sup>These are AND and OR gates that also express priority in a FT by taking into account the order of event occurrences.



system architectures (see, e.g. [UPM18]). The expressiveness of risk factors enables the combination of Boolean factors modelling failures internal to an autonomous robot with factors modelling hazards in physical processes in this robot’s environment.

Hansen et al. [HRS98] use the Duration Calculus to derive safety requirements from system FTs for the incremental design of safety-critical real-time controllers. We have shown an integration of FTA and incremental design in [GC17], however, from a less rigorous perspective. Because FTs can be modelled by RISKSTRUCTURES, an adoption of Hansen et al.’s approach suggests a refinement of the risk factors defined in Figs. 2a and 4 for deriving safety controllers for real-time systems.

**Algebraic methods for risk assessment** The only algebraic approach to design-oriented risk assessment, we could find, is from Hamdi and Boudriga [HB03]. They formalise IT security risk management as an algebraic datatype specification to check the consistency of security risk analyses viewed as algebras. Probability of occurrence and severity of consequence are modelled as metrics over attack actions to select optimal countermeasures using multi-objective optimisation. While focusing on security management, their algebraic datatype is an inspiration for a further formalisation of consequences and assets in Sect. 6.3. However, our action abstraction and the use of CSP offers a design method for safety controllers as opposed to the more abstract attack model proposed in [HB03], which focuses on risk assessment.

**From RISKSTRUCTURES to safety monitors and controllers** As mentioned for the comparison with [UPM18], each factor  $f$  (Fig. 2) can be seen as a *sequential monitor automaton* for a process  $P$ , with  $f$ ’s events being triggered (i.e., sensors, variable checkers) by  $P$ ’s actions. A risk structure  $\mathfrak{R}$  then forms a *concurrent monitor* of  $P$ ’s risk space, each risk state being concurrently monitored for the detection of endangerments (i.e., safety violations) and mitigation successes (i.e., mitigation acceptances).

Guiochet et al. [GPBB08] encode a risk model into a monitor automaton with partially ordered safety modes, each with a constraint limiting the parameters of robot actions to a mode-specific risk level. Refining this idea, Mekki-Mokhtar et al. [MMBG<sup>+</sup>12] distinguish between safe, warning, and catastrophic states with safety trigger conditions. Based on this framework, Machin et al. [MGW<sup>+</sup>18] apply model checking to determine whether catastrophic states can be reached despite such triggers. The authors present an algorithm for the synthesis of a mitigation policy, that is, minimal sequences of interventions to reach the nearest safe state from any warning state fulfilling validity, permissiveness, and safety. Our framework could enhance their refined hazard model with an algebraic method and the specification of dependencies between hazards.

Based on the MOP approach to RV [MJG<sup>+</sup>11], Huang et al. [HEZ<sup>+</sup>14] present a monitoring infrastructure for (i) checking trace properties (specified in different formalisms) observable from communications between ROS modules and (ii) the property-triggered execution of mitigations. Bogdiukiewicz et al. [BBH<sup>+</sup>17] describe how monitors can be formally developed in a step-wise manner using Event-B, an abstract state machine language supporting refinement-based development. RISKSTRUCTURES provide a bridge between formal risk modelling and the design of monitors in Event-B as shown in [BBH<sup>+</sup>17], in the ROS-specific MOP framework [HEZ<sup>+</sup>14], or the ROS-based framework proposed in [SLJS16].

**From RISKSTRUCTURES to risk-aware control** Works in (stochastic) optimal control [PBHS13, San14, FC14, MS14, SSSS18], summarised in Sect. 2.1, incorporate risk either as a (chance) constraint not to be violated or as a minimisation criterion for determining an optimal plan. The underlying stochastic models allow the per-state estimation of expected risk. Severity intervals would allow a per-state estimation of the *expected risk range*. Complementary to, for example, Sanger’s neuronal network-based controller [San14], RISKSTRUCTURES provide an algebraic method for structuring and composing complex risk models (i.e., state spaces, value functions, action spaces) in a verifiable manner. Once a risk model is found and validated, optimal control provides the low-level mathematical framework for detailed controller design.

## 9. Conclusion

The certification of autonomous systems requires the rigorous verification of their safety controllers. From the application of formal methods in other domains, we know that the emerging complexity of such controllers can be tackled by compositional verification. With RISKSTRUCTURES, we proposed a novel algebraic framework for the modelling and analysis of the risk profile of an automated or autonomous machine in its multi-scenario operational environment. This framework also supports the construction of discrete-event safety controllers that,

when integrated into the machine, improve this risk profile. These controllers incorporate an explicit and, thus, verifiable structure of the risk awareness of such a machine. The presented algebraic laws support the design of and reasoning about these controllers. The proposed theory was applied in form of situation/action/factor tables used to craft risk models with the YAP tool and resulted in the generation of risk graphs useful for controller design. To the best of our knowledge, this work is the first to provide an algebraic account of risk modelling and systematic safety controller design for risk-aware machines.

**Future work** Our work on this framework has led to a range of new research directions we want to pursue in the near future. The further development of our framework will include its extension to more general and more refined forms of risk factors; the elaboration of probabilistic factor semantics; the investigation of more specific mitigation orders, the construction of risk lattices from risk spaces, mitigation orders, and event structures; the development of analysers for risk estimation and synthesizers for safety controllers; the derivation of RISKSTRUCTURES from dynamical models; the further discussion of factor characteristics such as compatibility; the investigation of distributivity of composition and constraints; and a mechanisation and extension of further proofs using a specialised proof assistant such as Isabelle/UTP [FBC<sup>+</sup>20].

## Acknowledgements

Mario Gleirscher was supported in part by the German Research Foundation (DFG) under the Fellowship Grant no. 381212925. Work by Radu Calinescu and Mario Gleirscher was partially supported by the Lloyd's Register Foundation under the Autonomy Assurance International Programme (AAIP) Grant CSI:Cobot. Radu Calinescu was additionally supported by the UKRI Project EP/V026747/1 "Trustworthy Autonomous Systems Node in Resilience". We would like to thank Simon Foster for inspiring discussions on the use of relational specification; Ana Cavalcanti and Cliff Jones for insightful questions about the abstraction, composition, and methodology underlying RISKSTRUCTURES; James Baxter, Alvaro Miyazawa, and Pedro Ribeiro for enlightening conversations about CSP. We are also thankful to Sam Clark for helpful feedback on an early version of the introductory and closing sections.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Author Contributions.** • **Mario Gleirscher:** Conceptualization, Methodology, Formal analysis and investigation, Writing - original draft preparation, Writing - review and editing, Funding acquisition. • **Radu Calinescu:** Writing - review and editing, Funding acquisition, Resources, Supervision. • **Jim Woodcock:** Formal analysis and investigation, Writing - review and editing, Resources, Supervision.

## Abbreviations

AV	Autonomous vehicle
CSL	Continuous stochastic logic
CSP	Communicating sequential processes
CTMC	Continuous time Markov chain
ETA	Event tree analysis
FMEA	Failure mode and effects analysis
FT	Fault tree

<b>FTA</b>	Fault tree analysis
<b>HazOp</b>	Hazard and operability studies
<b>IH</b>	Induction hypothesis
<b>IOLTS</b>	Input/Output LTS
<b>IS</b>	Induction step
<b>LOPA</b>	Layer of protection analysis
<b>LTS</b>	Labelled transition system
<b>MCS</b>	Minimum cut set
<b>MDP</b>	Markov decision process
<b>MOP</b>	Monitoring-oriented programming
<b>ROS</b>	Robot operating system
<b>RV</b>	Run-time verification
<b>STAMP</b>	System-theoretic accident model and process
<b>STPA</b>	System-theoretic process analysis
<b>TS</b>	Transition system
<b>UML</b>	Unified modelling language

## A. Proof details

This section provides further proof details. The symbol  $\square$  indicates a discharged case distinction.

*Proof of Lemma 1.* The proof is by mutual existence and uniqueness: For each  $\sigma \in R(F_1 \cup F_2)$  (i) there exists a  $\sigma_1 \cup \sigma_2 \in R(F_1) \otimes R(F_2)$  and (ii) this pair is unique, and (iii, iv) conversely.

Below, we will need the infix binary operator scheme  $\cdot|_{F'} : R(F) \rightarrow R(F')$  that describes the (phase-invariant) projection from  $R(F)$  to  $R(F')$  with  $F' \subseteq F$ .

We show (i): Let  $\sigma \in R(F_1 \cup F_2)$ , then, by definition,  $\sigma$  is a total injection and, hence, every restriction of  $\sigma$  is a total injection, particularly, the restrictions  $\sigma|_{F_1}$  and  $\sigma|_{F_2}$ . Obviously, we have  $\sigma|_{F_1}(f) = \sigma|_{F_2}(f)$  for all  $f \in F_1 \cap F_2$ . Furthermore, by definition,  $\sigma(f)$  is faithful to  $Ph_f$  for all  $f \in F_1 \cup F_2$  and, thus, so are both these restrictions. These two results lead to  $\sigma_1 = \sigma|_{F_1} \in R(F_1)$  and  $\sigma_2 = \sigma|_{F_2} \in R(F_2)$  and, finally, to the existence of the wanted pair.  $\square$

We show (ii): Suppose there are two pairs  $(\sigma_1, \sigma_2) \neq (\sigma'_1, \sigma'_2) \in R(F_1) \otimes R(F_2)$ . Then there exists  $f \in F_1 \setminus F_2$  where  $\sigma_1(f) \neq \sigma'_1(f)$ , thus, also  $\sigma_1 \cup \sigma_2 \neq \sigma'_1 \cup \sigma_2$ . However, there can be no faithful total injection  $\sigma \in R(F_1 \cup F_2)$  such that  $\sigma = \sigma_1 \cup \sigma_2 \wedge \sigma = \sigma'_1 \cup \sigma_2$ .  $\square$

We show (iii): For each  $\sigma_1 \cup \sigma_2 \in R(F_1) \otimes R(F_2)$  there exists a  $\sigma \in R(F_1 \cup F_2) = \{\sigma \in (F_1 \cup F_2) \rightarrow \bigcup_{f \in F_1 \cup F_2} Ph_f \mid \sigma \text{ is a total injection} \wedge \forall f \in F_1 \cup F_2: \sigma(f) \in Ph_f\}$ : By definition, both  $\sigma_1$  and  $\sigma_2$  are faithful total injections (i.e., matching results for all  $f \in F_1 \cap F_2$ ). Then, we can construct a faithful total injection  $\sigma$  by applying set union to the domains and co-domains of these two. Thus,  $\sigma \in R(F_1 \cup F_2)$ .  $\square$

We show (iv): Suppose from  $\sigma_1 \cup \sigma_2$  we can construct  $\sigma \neq \sigma' \in R(F_1 \cup F_2)$  then there exists  $f \in F_1 \cup F_2$  such that  $\sigma(f) \neq \sigma'(f)$ . However, then either  $\sigma_1$  or  $\sigma_2$  must have violated injectivity which in turn would have violated the definition of  $R(F_1) \otimes R(F_2)$ .  $\square$

*Proof of Lemma 2.*  $(\mathcal{F}, \cup)$  is a semi-group because  $\cup$  is an associative binary operation on  $\mathcal{F}$ . We have that

$$\sigma_1 \approx (\sigma_2 \cup \sigma_3) \Leftrightarrow \sigma_1 \approx \sigma_2 \wedge \sigma_1 \approx \sigma_3 . \quad (17)$$

(The sub-proof based on the definition of  $\approx$  is omitted here.) Based on Formula (17), we show by algebraic manipulation that the binary operation  $\otimes$  on  $\mathcal{R}$  is associative:

$$\begin{aligned}
 & R(F_1) \otimes (R(F_2) \otimes R(F_3)) \\
 &= \{\sigma_1 \cup \sigma \mid \sigma_1 \in R(F_1) \wedge \sigma \in R(F_2) \otimes R(F_3) \wedge \sigma_1 \approx \sigma\} \\
 &= \{\sigma_1 \cup (\sigma_2 \cup \sigma_3) \mid \sigma_1 \in R(F_1) \wedge (\sigma_2 \cup \sigma_3) \in R(F_2) \otimes R(F_3) \wedge \sigma_1 \approx (\sigma_2 \cup \sigma_3)\} \\
 &= \{\sigma_1 \cup (\sigma_2 \cup \sigma_3) \mid \sigma_1 \in R(F_1) \wedge (\sigma_2 \in R(F_2) \wedge \sigma_3 \in R(F_3) \wedge \sigma_2 \approx \sigma_3) \wedge \sigma_1 \approx (\sigma_2 \cup \sigma_3)\} \\
 &= \{(\sigma_1 \cup \sigma_2) \cup \sigma_3 \mid (\sigma_1 \in R(F_1) \wedge \sigma_2 \in R(F_2)) \wedge \sigma_3 \in R(F_3) \wedge \sigma_2 \approx \sigma_3 \wedge \sigma_1 \approx (\sigma_2 \cup \sigma_3)\} \quad (\text{by Formula (17)}) \\
 &= \{(\sigma_1 \cup \sigma_2) \cup \sigma_3 \mid (\sigma_1 \in R(F_1) \wedge \sigma_2 \in R(F_2)) \wedge \sigma_3 \in R(F_3) \wedge \sigma_2 \approx \sigma_3 \wedge \sigma_1 \approx \sigma_2 \wedge \sigma_1 \approx \sigma_3\} \\
 &= \{(\sigma_1 \cup \sigma_2) \cup \sigma_3 \mid (\sigma_1 \in R(F_1) \wedge \sigma_2 \in R(F_2) \wedge \sigma_1 \approx \sigma_2) \wedge \sigma_3 \in R(F_3) \wedge (\sigma_1 \cup \sigma_2) \approx \sigma_3\} \\
 &= \{(\sigma_1 \cup \sigma_2) \cup \sigma_3 \mid (\sigma_1 \cup \sigma_2) \in R(F_1) \otimes R(F_2) \wedge \sigma_3 \in R(F_3) \wedge (\sigma_1 \cup \sigma_2) \approx \sigma_3\} \\
 &= (R(F_1) \otimes R(F_2)) \otimes R(F_3)
 \end{aligned}$$

Hence,  $(\mathcal{R}, \otimes)$  is a semi-group, too. Lemma 1 completes the proof.  $\square$

*Proof of Lemma 3.*

We express the basic factor template from Fig. 2a in CSP as

$$\begin{aligned}
 f_{(2a)} &= 0^f \\
 0^f &= ?x : e^f \rightarrow f \square ?x : o_n^f \rightarrow 0^f \\
 f &= ?x : m_d^f \rightarrow 0^f \square ?x : o_e^f \rightarrow f \square ?x : m^f \rightarrow \bar{f} \\
 \bar{f} &= ?x : \bar{e}^f \rightarrow f \square ?x : o_m^f \rightarrow \bar{f} \square ?x : m_r^f \rightarrow 0^f .
 \end{aligned}$$

(As indicated in the Sect. 4.1 and 4.4, we consistently omit  $e^f$  for the sake of simplicity.) Analogously, we express the factor template with isolated non-determinism from Fig. 4 in CSP as

$$\begin{aligned}
 f_{(4)} &= 0^f \\
 0^f &= ?x : e^f \setminus ue^f \rightarrow f \square ?x : o_n^f \setminus ue^f \rightarrow 0^f \square (?x : ue^f \rightarrow (0^f \square f)) \\
 f &= ?x : m_d^f \rightarrow 0^f \square ?x : o_e^f \setminus um^f \rightarrow f \square ?x : m^f \setminus um^f \rightarrow \bar{f} \square (?x : um^f \rightarrow (f \square \bar{f})) \\
 \bar{f} &= ?x : \bar{e}^f \rightarrow f \square ?x : o_m^f \setminus um_r^f \rightarrow \bar{f} \square ?x : m_r^f \setminus um_r^f \rightarrow 0^f \square (?x : um_r^f \rightarrow (0^f \square \bar{f})) .
 \end{aligned}$$

The inherent difference of the actions  $m_d^f, \bar{e}^f$ , and  $e^f$  from their phase competitors allows us to make the simplifying assumption  $\bar{e}^f \cap o_m^f = \bar{e}^f \cap m_r^f = e^f \cap o_e^f = e^f \cap m^f = e^f \cap m_d^f = m^f \cap m_d^f = \emptyset$ . We further define  $ue^f = o_n^f \cap e^f$ ,  $um^f = o_e^f \cap m^f$ ,  $um_r^f = o_m^f \cap m_r^f$ . For the equivalence proof, we have to distinguish whether or not the events in the three sets  $\{e^f, o_n^f\}$ ,  $\{m_d^f, o_e^f, m^f\}$ , and  $\{\bar{e}^f, o_m^f, m_r^f\}$  are pairwise disjoint.

*Case 1: Events are pairwise disjoint.* Then,

1. from the emptiness of  $ue^f$ ,  $um^f$ , and  $um_r^f$ ,
2. from applying the CSP laws  $?x : \emptyset \rightarrow P =_{FD} STOP$  and  $P \square STOP =_{FD} P$  [Sch99, p. 93] to  $f_{(4)}$ , and
3. from the resulting structure of  $f_{(4)}$  being equivalent to  $f_{(2a)}$ ,

we immediately derive  $f_{(4)} =_{FD} f_{(2a)}$ .  $\square$

*Case 2: Events are not pairwise disjoint.* Then,

1. from the potential non-emptiness of  $ue^f$ ,  $um^f$ , and  $um_r^f$ ,
2. from applying the step law of external choice [Sch99, p. 93] to  $f_{(2a)}$  in order to isolate, for example,  $?x : o_n^f \cap e^f \rightarrow (0^f \square f)$ , and
3. from our definitions (e.g.  $ue^f = o_n^f \cap e^f$ ),

we again immediately have that  $f_{(2a)} =_{FD} f_{(4)}$ . The cases 1 and 2 complete the proof.  $\square$

The following corollary and its proof rule out potential errors in the proof of Lemma 4.

**Corollary 7 (Converse of Lemma 4)**

$$\begin{aligned} \neg(\sigma \leq_m \sigma') &\Leftarrow \neg(\sigma \lesssim_m \sigma') && \text{(by definition of } \not\leq_m, \lesssim_m) \\ \sigma \not\leq_m \sigma' &\Leftarrow \sigma \not\lesssim_m \sigma' && \text{(by negation and definition of } \leq_m, \lesssim_m) \\ \sigma >_m \sigma' \vee ((\sigma, \sigma') \notin \leq_m \wedge (\sigma', \sigma) \notin \leq_m) &\Leftarrow \sigma >_{\sim} \sigma' \vee ((\sigma, \sigma') \notin \lesssim_m \wedge (\sigma', \sigma) \notin \lesssim_m) && (18) \end{aligned}$$

*Proof of Corollary 7. Case 1:* If  $(\sigma, \sigma') \notin \lesssim_m \wedge (\sigma', \sigma) \notin \lesssim_m$  then (by Definitions 8 and 9) also  $(\sigma, \sigma') \notin \leq_m \wedge (\sigma', \sigma) \notin \leq_m$  and, therefore, Formula (18).  $\square$

*Case 2:* If  $\sigma >_{\sim} \sigma' = \sigma' \lesssim_m \sigma \wedge \sigma' \not\sim_m \sigma$ , then we have either  $(\sigma, \sigma') \notin \leq_m \wedge (\sigma', \sigma) \notin \leq_m$  (because of some incomparable phases) which fulfils Formula (18). Alternatively, we have  $(\sigma', \sigma) \in \leq_m$  which means  $\sigma$  and  $\sigma'$  are fully comparable and, because of Corollary 2, we have  $\sigma >_m \sigma'$ .  $\square$

*Proof of Lemma 5.* For this we only need to show that any two risk states  $\sigma, \sigma' \in R$  are (i) comparable and (ii) antisymmetric:  $[\sigma]_{\sim_s} \leq_m [\sigma']_{\sim_s} \wedge [\sigma']_{\sim_s} \leq_m [\sigma]_{\sim_s} \Rightarrow [\sigma]_{\sim_s} =_m [\sigma']_{\sim_s}$ .

We show (i) by showing that the conditions (a) and (b) guarantee the comparability of any two risk states in  $R$  based on the interval order  $\leq$  as defined above and use the fact that comparability can be dropped from  $R/\sim_s$  to  $R$  using  $[\sigma]_{\sim_s} \leq_m [\sigma']_{\sim_s} \iff \forall \check{\sigma} \in [\sigma]_{\sim_s}, \check{\sigma}' \in [\sigma']_{\sim_s}: \check{\sigma} \leq_m \check{\sigma}'$  (Formula (6)). Recall that we equate the empty interval with the empty set,  $\perp = \emptyset$ . Now, we have to consider the following three cases to complete the sub-proof of (i):

*Case “no risk factors activated”:*  $\perp \geq \perp \vee \emptyset \subset \emptyset$  validates (b) and lack of comparability invalidates (a), yet we have  $[\sigma]_{\sim_s} \leq_m [\sigma']_{\sim_s}$ .  $\square$

*Case “at least one risk factor activated only in  $\check{\sigma}$ ”:* Let  $[a, b] = S(\check{\sigma})$ .  $[a, b] \geq \perp \vee [a, b] \subset \emptyset$  invalidates (a) and (b). However, the observation  $\perp \geq [a, b] \vee \emptyset \subset [a, b]$  yields  $[\sigma']_{\sim_s} \leq_m [\sigma]_{\sim_s}$ . The dual of this case works analogously.  $\square$

*Case “at least one risk factor activated in both  $\check{\sigma}, \check{\sigma}'$ ”:* Let  $[a, b] = S(\check{\sigma})$  and  $[c, d] = S(\check{\sigma}')$ . We need to show  $[a, b] \geq [c, d] \vee [a, b] \subset [c, d]$  for all  $a, b, c, d \in \mathbb{R}$ : We can assume  $a \leq b \wedge c \leq d$  by definition of intervals. Then, we face the case that

- $c \leq a \wedge d \leq b$  validates (a) by definition of  $\leq$  over intervals,
- $c > a \wedge d \leq b$  validates (b),
- $c \leq a \wedge d > b$  implies (b) for the dual case  $[\sigma']_{\sim_s} \leq_m [\sigma]_{\sim_s}$ , or
- $c > a \wedge d > b$  implies (a) for the dual case.

Hence, from each of these four cases, either  $[\sigma]_{\sim_s} \leq_m [\sigma']_{\sim_s}$  or  $[\sigma']_{\sim_s} \leq_m [\sigma]_{\sim_s}$  follows.

Then, we show (ii) by contradiction: Assuming  $[\sigma]_{\sim_s} \leq_m [\sigma']_{\sim_s} \wedge [\sigma']_{\sim_s} \leq_m [\sigma]_{\sim_s}$ , we have  $\forall \check{\sigma} \in [\sigma]_{\sim_s}, \check{\sigma}' \in [\sigma']_{\sim_s}: \check{\sigma} \leq_m \check{\sigma}' \wedge \check{\sigma}' \geq_m \check{\sigma}$  by dropping from  $R/\sim_s$ . Now, we claim that  $[\sigma]_{\sim_s} \neq_m [\sigma']_{\sim_s}$ . Hence,  $\exists \check{\sigma} \in [\sigma]_{\sim_s}, \check{\sigma}' \in [\sigma']_{\sim_s}: \check{\sigma} \not\sim_m \check{\sigma}'$  and, consequently,  $\exists \check{\sigma} \in [\sigma]_{\sim_s}, \check{\sigma}' \in [\sigma']_{\sim_s}: \check{\sigma} \not\leq_m \check{\sigma}' \vee \check{\sigma}' \not\leq_m \check{\sigma}$ . The latter contradicts our assumption.  $\square$

*Proof of Lemma 7.* Fix a finite  $F \subseteq \mathcal{F}$  and a pair  $\sigma, \sigma' \in R(F)$ . The proof is by induction over  $F$  and relies on the assumption that, for any  $f \in F$ , the phase  $f$  is the *unique maximal element* of  $\leq_f$  and that *active* only returns such elements:

Induction start  $F_0 = \emptyset$ :  $\sigma \upharpoonright_{\emptyset} \leq_m \sigma' \upharpoonright_{\emptyset}$  holds trivially and so does  $\text{active}(\sigma' \upharpoonright_{\emptyset}) \subseteq \text{active}(\sigma \upharpoonright_{\emptyset})$ .

*Induction step (IS)*  $F_{n+1} = F_n \cup \{f\}$  where  $n \geq 0$  and  $f \in F \setminus F_n$ : For the IH, assume  $\sigma \upharpoonright_{F_n} \leq_m \sigma' \upharpoonright_{F_n} \Rightarrow \text{active}(\sigma' \upharpoonright_{F_n}) \subseteq \text{active}(\sigma \upharpoonright_{F_n})$ . Then, we show that  $\sigma \upharpoonright_{F_{n+1}} \leq_m \sigma' \upharpoonright_{F_{n+1}} \Rightarrow \text{active}(\sigma' \upharpoonright_{F_{n+1}}) \subseteq \text{active}(\sigma \upharpoonright_{F_{n+1}})$  (i.e., the IS). For this, we prove

- the case of *incomparable phases*  $(\sigma(f), \sigma'(f)) \notin \leq_f$ : In this case, the state pair gets incomparable and the implication is trivially fulfilled,
- the *inverse case*  $\sigma(f) >_f \sigma'(f)$ : In this case, the state pair gets incomparable and the implication is again trivially fulfilled, and
- the *aligned case*  $\sigma(f) \leq_f \sigma'(f)$ : In this case, the state pair stays comparable. However,  $\sigma'(f)$  can either be  $f$  (hence,  $\sigma(f) = f$  and maintaining  $\subseteq$ ),  $\bar{f}$  (hence,  $\sigma(f) \in \{f, \bar{f}, 0^f\}$  and maintaining  $\subseteq$ ), or  $0^f$  (see former sub-case).

Having proved these cases completes the induction step by establishing the IH.

For  $\lesssim_m$ , we only have to substitute case 1 of the IS: For  $f$ , the smallest  $\leq_f$  after Definition 3 implies incomparability at most for  $(0^f, f)$  and  $(f, 0^f)$ . These two phase pairs of  $f$  would not alter the *active* maps of both states and hence maintain  $\subseteq$  as well.  $\square$

*Proof of Lemma 12.* From associativity of set intersection in Definition 12, we know that the order in which constraints are applied to  $R \times R$  does not matter. Consequently, we also have  $\llbracket k \rrbracket_c^R \supseteq \llbracket \{k, k'\} \rrbracket_c^R$  for any  $k, k' \in \mathcal{C}$ , that is, constraints only prune  $R \times R$ . First, we prove the lemma

$$\llbracket [\mathfrak{A}]_{\{k\}} \rrbracket_C = [\mathfrak{A}]_{\{k\} \cup C} . \quad (19)$$

Fix  $\sigma \xrightarrow{e} \sigma'$ .  
 $\Rightarrow$ :

1. Case  $(\sigma, \sigma') \in \llbracket k \rrbracket_c^R$ :  $[\cdot]_C$ -step applies to  $[\mathfrak{A}]_{\{k\}}$ .

(a) Case  $(\sigma, \sigma') \in \llbracket C \rrbracket_c^R$ :  $[\cdot]_C$ -step applies also to  $\llbracket [\mathfrak{A}]_{\{k\}} \rrbracket_C$ . Then,  $(\sigma, \sigma')$  must be in  $\llbracket k \rrbracket_c^R \cap \llbracket C \rrbracket_c^R$ , which by Definition 12 is  $\llbracket \{k\} \cup C \rrbracket_c^R$ . Because of  $(\sigma, \sigma') \in \llbracket \{k\} \cup C \rrbracket_c^R$ ,  $[\cdot]_C$ -step applies to  $[\mathfrak{A}]_{\{k\} \cup C}$  in the same way as it does to  $\llbracket [\mathfrak{A}]_{\{k\}} \rrbracket_C$ .  $[\cdot]_C$ -cancel does not apply. Consequently, if  $\sigma \xrightarrow{e} \sigma'$  is permitted (and added) by the left-hand side (LHS) then it is permitted by the right-hand side (RHS) of Formula (19).

(b) Case  $(\sigma, \sigma') \notin \llbracket C \rrbracket_c^R$ :  $[\cdot]_C$ -cancel applies to  $\llbracket [\mathfrak{A}]_{\{k\}} \rrbracket_C$  instead of  $[\cdot]_C$ -step, resulting in the production of  $\sigma \xrightarrow{\tau_f}_C \sigma$  rather than  $\sigma \xrightarrow{e}_C \sigma'$  on the LHS. Because  $(\sigma, \sigma') \notin \llbracket C \rrbracket_c^R$ , we have  $(\sigma, \sigma') \notin \llbracket \{k\} \rrbracket_c^R \cap \llbracket C \rrbracket_c^R$  and, by Definition 12,  $(\sigma, \sigma') \notin \llbracket \{k\} \cup C \rrbracket_c^R$ . Hence,  $[\cdot]_C$ -cancel applies to RHS and produces  $\sigma \xrightarrow{\tau_f}_C \sigma$ .

2. Case  $(\sigma, \sigma') \notin \llbracket k \rrbracket_c^R$ :  $[\cdot]_C$ -cancel applies to  $[\mathfrak{A}]_{\{k\}}$  instead of  $[\cdot]_C$ -step, producing  $\sigma \xrightarrow{\tau_f}_C \sigma$  instead of  $\sigma \xrightarrow{e}_C \sigma'$  on the inner LHS.

(a) Case  $(\sigma, \sigma) \in \llbracket C \rrbracket_c^R$ :  $[\cdot]_C$ -step is applicable to  $\llbracket [\mathfrak{A}]_{\{k\}} \rrbracket_C$ , producing  $\sigma \xrightarrow{\tau_f}_C \sigma$  for the outer LHS.

(b) Case  $(\sigma, \sigma) \notin \llbracket C \rrbracket_c^R$ :  $[\cdot]_C$ -cancel is applicable to  $\llbracket [\mathfrak{A}]_{\{k\}} \rrbracket_C$ , producing  $\sigma \xrightarrow{\tau_f}_C \sigma$  for the outer LHS.

Because of  $\llbracket \{k\} \cup C \rrbracket_c^R = \llbracket k \rrbracket_c^R \cap \llbracket C \rrbracket_c^R$ , in both cases (a) and (b), the  $[\cdot]_C$ -cancel rule applies to  $[\mathfrak{A}]_{\{k\} \cup C}$ , producing a  $\sigma \xrightarrow{\tau_f}_C \sigma$  on the RHS of Formula (19).  $\square$

$\Leftarrow$ :

1. If  $(\sigma, \sigma') \in \llbracket \{k\} \cup C \rrbracket_c^R$  then  $[\cdot]_C$ -step is applicable to  $[\mathfrak{A}]_{\{k\} \cup C}$  and, because of  $\llbracket k \rrbracket_c^R \cap \llbracket C \rrbracket_c^R$ , to the inner and outer constraint in  $\llbracket [\mathfrak{A}]_{\{k\}} \rrbracket_C$ . Thus, if  $\sigma \xrightarrow{e}_C \sigma'$  is produced on the RHS then so on the LHS.

2. If  $(\sigma, \sigma') \notin \llbracket \{k\} \cup C \rrbracket_c^R$  then  $[\cdot]_C$ -cancel applies to  $[\mathfrak{A}]_{\{k\} \cup C}$  and produces  $\sigma \xrightarrow{\tau_f}_C \sigma$  on the RHS. However, because of  $(\sigma, \sigma') \notin \llbracket k \rrbracket_c^R \cap \llbracket C \rrbracket_c^R$ , we have at least one of the two cases for the LHS:

(a) Case  $(\sigma, \sigma') \notin \llbracket k \rrbracket_c^R$ : This matches case (2a), that is,  $[\cdot]_C$ -cancel and  $[\cdot]_C$ -step produce  $\sigma \xrightarrow{\tau_f}_C \sigma$  for  $\llbracket [\mathfrak{A}]_{\{k\}} \rrbracket_C$ .

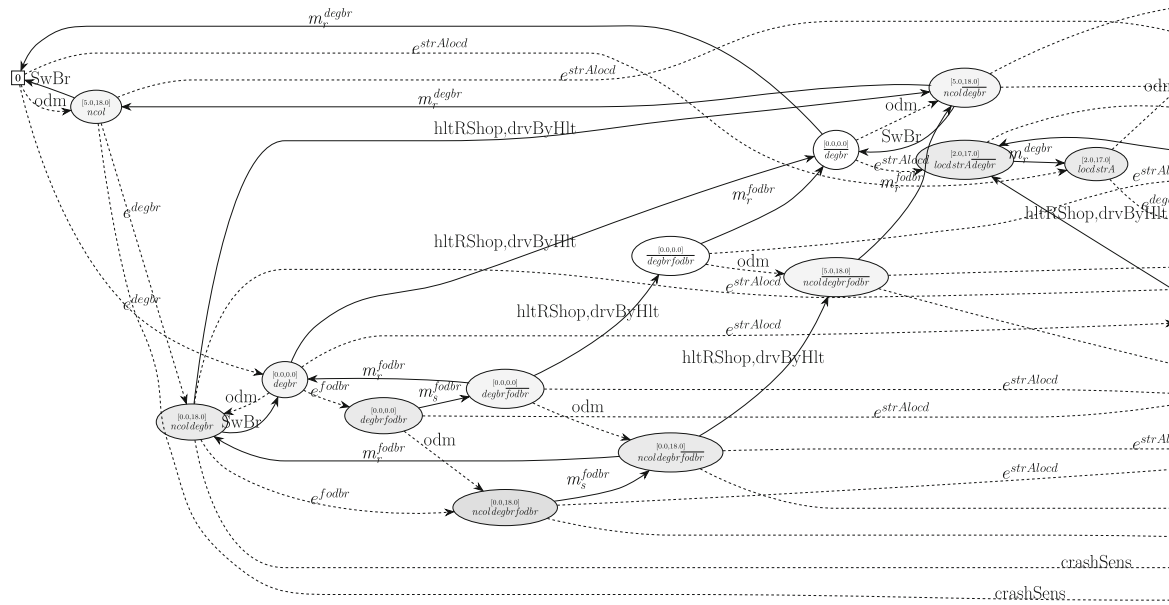
(b) Case  $(\sigma, \sigma') \notin \llbracket C \rrbracket_c^R$ : This matches either case (1b) or case (2b), that is, in both cases  $[\cdot]_C$ -cancel and  $[\cdot]_C$ -step produce  $\sigma \xrightarrow{\tau_f}_C \sigma$  for  $\llbracket [\mathfrak{A}]_{\{k\}} \rrbracket_C$ .  $\square$

Second, we show by induction over  $C_1$  that the  $[\cdot]_C$ -step and  $[\cdot]_C$ -cancel rules apply in the same way to both sides of Formula (15).

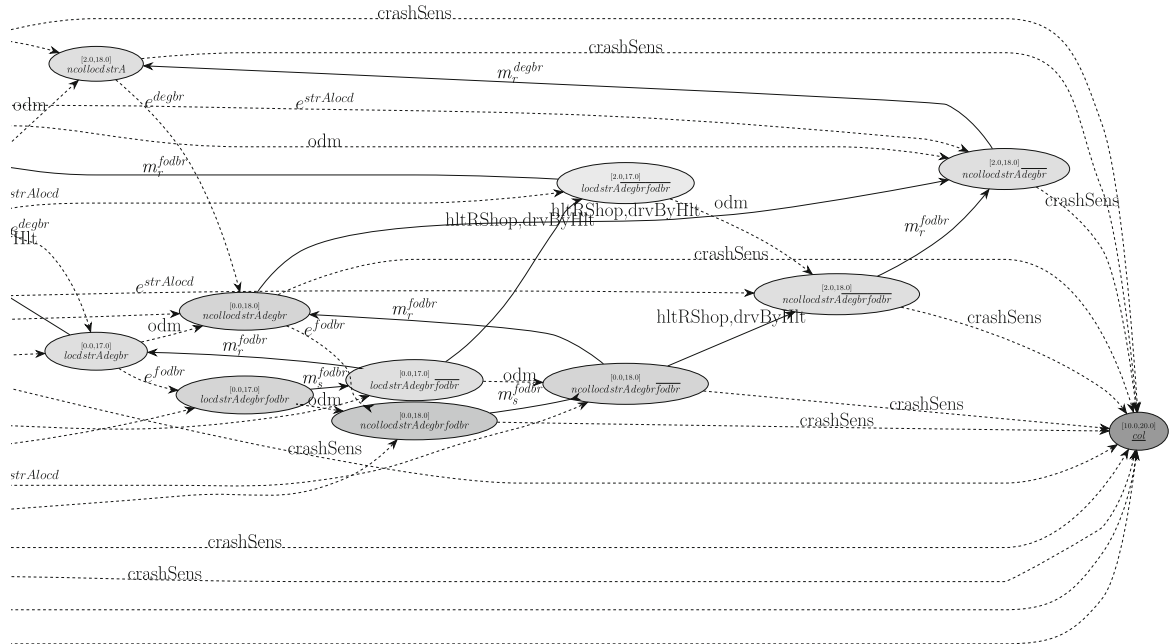
Induction starts with  $C_1^0 = \emptyset$ : In this case, we have  $\llbracket [\mathfrak{A}]_{\emptyset} \rrbracket_{C_2} = [\mathfrak{A}]_{\emptyset \cup C_2}$  and by applying  $[\cdot]_C$ -step or Definition 12 on the left and  $\emptyset \cup C_2 = C_2$  on the right, we get the tautology  $[\mathfrak{A}]_{C_2} = [\mathfrak{A}]_{C_2}$ .  $\square$

IS with  $C_1^{n+1} = C_1^n \cup \{k\}$  where  $n \geq 0$  and  $k \in C_1 \setminus C_1^n$ : By assuming  $\llbracket [\mathfrak{A}]_{C_1^n} \rrbracket_{C_2} = [\mathfrak{A}]_{C_1^n \cup C_2}$  (IH), we show

$$\begin{aligned} & \llbracket [\mathfrak{A}]_{C_1^{n+1}} \rrbracket_{C_2} && \text{(by def of IS)} \\ & = \llbracket [\mathfrak{A}]_{C_1^n \cup \{k\}} \rrbracket_{C_2} && \text{(by } \cup\text{-comm and Formula (19))} \\ & = \llbracket \llbracket [\mathfrak{A}]_{\{k\}} \rrbracket_{C_1^n} \rrbracket_{C_2} && \text{(by IH)} \\ & = \llbracket [\mathfrak{A}]_{\{k\}} \rrbracket_{C_1^n \cup C_2} && \text{(by Formula (19))} \\ & = [\mathfrak{A}]_{\{k\} \cup C_1^n \cup C_2} && \text{(by def)} \\ & = [\mathfrak{A}]_{C_1^{n+1} \cup C_2} && \square \end{aligned}$$



(a) Part 1



(b) Part 2

Fig. 11. Risk structure for all factors mentioned in Sects. 4.3, 7.4

## B. Supplemental material for the road vehicle example

Figures 11a and 11b contain two parts of a larger risk structure for the road vehicle example illustrated and discussed in the Sects. 7.4, 4.3, 3.

## References

- [AASB<sup>+</sup>06] Alami R, Albu-Schäffer A, Bicchi A, Bischoff R, Chatila R, De Luca A, De Santis A, Giralt G, Guiochet J, Hirzinger G, Ingrand F, Lippiello V, Mattone R, Powell D, Sen S, Siciliano B, Tonietti G, Villani L (2006) Safe and dependable physical human-robot interaction in anthropic domains: State of the art and challenges. In: Intelligent robots and systems (IROS), IEEE/RSJ international conference
- [AKWB11] Althoff D, Kuffner JJ, Wollherr D, Buss M (2011) Safety assessment of robot trajectories for navigation in uncertain and dynamic environments. *Auton Robots* 32(3):285–302
- [ALRL04] Avizienis A, Laprie J-C, Randell B, Landwehr C (2004) Basic concepts and taxonomy of dependable and secure computing. *Dependable Secure Comput IEEE Trans* 1(1):11–33
- [AZI<sup>+</sup>17] Ajoudani A, Zanchettin AM, Ivaldi S, Albu-Schäffer A, Kosuge K, Khatib O (2017) Progress and prospects of the human-robot collaboration. *Auton Robots*, 42(5):957–975
- [BBH<sup>+</sup>17] Bogdiukiewicz C, Butler M, Hoang TS, Paxton M, Snook J, Waldron X, Wilkinson T (2017) Formal development of policing functions for intelligent systems. In: Software Reliability Engineering (ISSRE), 28th IEEE international symposium
- [BCS07] Boudali H, Crouzen P, Stoelinga M (2007) Dynamic fault tree analysis using input/output interactive Markov chains. In: Dependable systems and networks (DSN), 37th annual IEEE/IFIP international conference, pp 708–717
- [BFK13] Beer M, Ferson S, Kreinovich V (2013) Imprecise probabilities in engineering analyses. *Mech Syst Signal Process*, 37(1-2):4–29
- [Bir17] Birolini A (2017) Reliability Engineering. Springer, 8th edition
- [BK08] Baier C, Katoen J-P (2008) Principles of Model Checking. MIT Press
- [BS01] Broy M, Stølen K (2001) Specification and Development of Interactive Systems. Springer
- [CLC<sup>+</sup>19] Chen C, Liu X, Chen H-H, Li M, Zhao L (2019) A rear-end collision risk evaluation and control scheme using a Bayesian network model. *IEEE Trans Intell Transp Syst*, 20(1):264–284
- [Eri15] Ericson CA (2015) Hazard analysis techniques for system safety. Wiley, 2nd edition
- [FA04] Fraichard T, Asama H (2004) Inevitable collision states – a step towards safer robots?. *Adv Robotics*, 18(10):1001–1024
- [FBC<sup>+</sup>20] Foster S, Baxter J, Cavalcanti A, Woodcock J, Zeyda F (2020) Unifying semantic foundations for automated verification tools in Isabelle/UTP. *Sci Comput Program*, 197:102510
- [FC14] Feyzabadi S, Carpin S (2014) Risk-aware path planning using hierarchical constrained Markov decision processes. In: Automation science and engineering (CASE), IEEE international conference IEEE
- [FGC20] Foster S, Gleirscher M, Calinescu R (2020) Towards deductive verification of control algorithms for autonomous marine vehicles. In: Engineering of complex computer systems (ICECCS), 25th international conference, pp 113–118, Singapore
- [Foo78] Foot P (1978) The problem of abortion and the doctrine of the double effect. *Virtues and Vices and Other Essays in Moral Philosophy*, 19. Originally published in 1967.
- [GC17] Gleirscher M, Carlan C (2017) Arguing from hazard analysis in safety cases: A modular argument pattern. In: High Assurance Systems Engineering (HASE), 18th international symposium, pp 53–60. IEEE
- [GC20] Gleirscher M, Calinescu R (2020) Safety controller synthesis for collaborative robots. In: Engineering of complex computer systems (ICECCS), 25th international conference, pp 83–92, Singapore
- [GFW19] Gleirscher M, Foster S, Woodcock J (2019) New opportunities for integrated formal methods. *ACM Comput Surv*, 52:117:1–117:36
- [GK17] Gleirscher M, Kugele S (2017) From hazard analysis to hazard mitigation planning: The automated driving case. In: Barrett C, et al., editor, NASA Formal Methods (NFM), 9th international symposium, volume 10227 of *LNCS*. Springer
- [Gle14] Gleirscher M (2014) Behavioral Safety of Technical Systems. Dissertation, Technische Universität München
- [Gle17] Gleirscher M (2017) Run-time risk mitigation in automated vehicles: A model for studying preparatory steps. In: Bulwahn L, Kamali M, Linker S, (eds), Formal Verification of Autonomous Vehicles (FVAV), 1st iFM Workshop, EPTCS
- [Gle18] Gleirscher M (2018) Strukturen für die Gefahrenerkennung und -behandlung in autonomen Maschinen. In: Jürgen B, Petra W (eds), Beiträge zu einer Systemtheorie Sicherheit, acatech DISKUSSION, Chapter 8.4, pp 154–167. Herbert Utz Verlag, München
- [Gle20] Gleirscher M (2020) YAP: Tool support for deriving safety controllers from hazard analysis and risk assessments. In: Luckuck M, Farrell M (eds), Formal Methods for Autonomous Systems (FMAS), 2nd Workshop, volume 329 of EPTCS, pp 31–47. Open Publishing Association
- [Gle21] Gleirscher M (2021) YAP Against Perils: Application Guide and User’s Manual. University of York and Technical University of Munich
- [GM20] Gleirscher M, Marmsoler D (2020) Formal methods in dependable systems engineering: A survey of professionals from Europe and North America. *Empir Softw Eng*, 25(6):4473–4546
- [GMGP10] Guiochet J, Martin-Guillerez D, Powell D (2010) Experience with model-based user-centered risk assessment for service robots. In: High Assurance Systems Engineering (HASE), 12th IEEE international symposium
- [GPBB08] Guiochet J, Powell D, Baudin É, Blanquart J-P(2008) Online safety monitoring using safety modes. In: Technical challenges for dependable robots in human environments, workshop, pp 1–13. Rapport LAAS No. 08339
- [Har00] Harpwood V (2000) Principles of Tort Law. Cavendish, 4th edn
- [HASH09] Haddadin S, Albu-Schäffer A, Hirzinger G (2009) Requirements for safe robots: Measurements, analysis and new insights. *Int J Robot Res*, 28(11-12):1507–1527
- [HB03] Hamdi M, Boudriga N (2003) Algebraic specification of network security risk management. In: ACM workshop on formal methods in security engineering (FMSE). ACM Press
- [HEZ<sup>+</sup>14] Huang J, Erdogan C, Zhang Y, Moore B, Luo Q, Sundareshan A, Rosu G (2014) ROSRV: Runtime verification for robots. In: Runtime verification, pp 247–254. Springer
- [HG03] Holland O, Goodman R (2003) Robots with internal models: A route to machine consciousness?. *J Conscious Stud*, 10(4-5):77–109
- [HM99] Howe RD, Matsuoka Y (1999) Robotics for surgery. *Ann Rev Biomed Eng*, 1(1):211–240



- [Hoa85] Hoare T (1985) Communicating sequential processes. International series in computer science, Prentice-Hall
- [HRS98] Hansen KM, Ravn AP, Stavridou V (1998) From safety analysis to software requirement. *IEEE Trans Softw Eng*, 24(7):573–84
- [IMM18] Iamsung C, Mosleh A, Modarres M (2018) Monitoring and learning algorithms for dynamic hybrid Bayesian network in on-line system health management applications. *Reliab Eng Syst Safety*, 178:118–129
- [IT90] Ishibuchi H, Tanaka H (1990) Multiobjective programming in optimization of the interval objective function. *Eur J Oper Res*, 48(2):219–225
- [KG81] Kaplan S, Garrick BJ (1981) On the quantitative definition of risk. *Risk Anal*, 1(1):11–27
- [Kum07] Kumamoto H (2007) Satisfying safety goals by probabilistic risk assessment. *Reliability engineering*. Springer
- [KW17] Koopman P, Wagner M (2017) Autonomous vehicle safety: An interdisciplinary challenge. *IEEE Int Transp Syst Mag*, 9(1):90–96
- [Lev95] Leveson NG (1995) *Safeware: system safety and computers*. Addison-Wesley
- [Lev04] Leveson NG (2004) A new accident model for engineering safer systems. *Safety Sci*, 42(4):237–70
- [Lev12] Leveson NG (2012) *Engineering a safer world: systems thinking applied to safety*. Engineering systems. MIT Press
- [LFL13] Leitner-Fischer F, Leue S (2013) Probabilistic fault tree synthesis using causality computation. *Int J Critical Comput-Based Syst*, 4(2):119–43
- [LS87] Leveson NG, Stolzy JL (1987) Safety analysis using Petri nets. *IEEE Trans Softw Eng*, 13(3):386–97
- [LS09] Leucker M, Schallhart C (2009) A brief account of runtime verification. *J Logic Algeb Program*, 78(5):293–303
- [LSS11] Lund MS, Solhaug B, Stølen K (2011) Model-driven risk analysis: The CORAS approach. Springer
- [McD94] McDermid John A (1994) Support for safety cases and safety arguments using SAM. *Reliability Engineering & System Safety*, 43(2):111–127
- [MGW<sup>+</sup>18] Machin M, Guiochet J, Waeselynck H, Blanquart J-P, Roy M, Masson L (2018) SMOF – a safety monitoring framework for autonomous systems. *IEEE Trans Syst Man Cybern: Syst*, 48(5):702–715
- [MJG<sup>+</sup>11] Meredith PO, Jin D, Griffith D, Chen F, Roşu G (2011) An overview of the MOP runtime verification framework. *Int J Softw Tools Technol Trans*, 14(3):249–289
- [MMBG<sup>+</sup>12] Mekki-Mokhtar A, Blanquart J-P, Guiochet J, Powell D, Roy M (2012) Safety trigger conditions for critical autonomous systems. In: *Dependable Computing (PRDC)*, 18th IEEE Pacific Rim International symposium. IEEE
- [MS14] Müller J, Sukhatme GS (2014) Risk-aware trajectory generation with application to safe quadrotor landing. In: *Intelligent Robots and Systems (IROS)*, IEEE/RSJ International conference
- [NSV03] Netravali AN, Sabnani KK, Viswanathan R (2003) Correct passive testing algorithms and complete fault coverage. In: *Formal techniques for networked and distributed systems (FORTE)*, pp 303–318. Springer
- [OCW09] Oliveira M, Cavalcanti A, Woodcock J (2009) A UTP semantics for *Circus*. *Formal Aspects Comput*, 21(1-2):3–32
- [ORS06] Ortmeier F, Reif W, Schellhorn G (2006) Deductive cause-consequence analysis (DCCA). *IFAC Proc Vol*, 38(1):62–67
- [PBHS13] Pereira AA, Binney J, Hollinger GA, Sukhatme GS (2013) Risk-aware path planning for autonomous underwater vehicles using predictive ocean models. *J Field Robot*, 30(5):741–762
- [Ros10] Roscoe AW (2010) *Understanding concurrent systems*. Springer
- [San14] Sanger TD (2014) Risk-aware control. *Neural Comput*, 26(12):2669–2691
- [SC88] Sobek RP, Chatila RG (1988) Integrated planning and execution control for an autonomous mobile robot. *Artif Intell Eng*, 3(2):103–113
- [Sch99] Schneider S (1999) *Concurrent and real-time systems: the CSP approach*. Wiley, New York
- [She03] Sheridan TB (2003) *Telerobotics, automation, and human supervisory control*. The MIT Press
- [Sim94] Simmons RG (1994) Structured control for autonomous robots. *IEEE Trans Robot Autom*, 10(1):34–43
- [SLJS16] Sorin A, Larsen M, Jensen K, Schultz UP (2016) Rule-based dynamic safety monitoring for mobile robots. *J Softw Eng Robot*, 7(1):120–141
- [SR02] Svedung I, Rasmussen J (2002) Graphic representation of accident scenarios: mapping system structure and the causation of accidents. *Safety Sci*, 40(5):397–417
- [SSSS18] Shalev-Shwartz S, Shammah S, Shashua A (2018) On a formal model of safe and scalable self-driving cars. Technical report, Mobileye
- [Tre08] Tretmans J (2008) Model based testing with labelled transition systems. In: *Formal methods and testing*, pp 1–38. Springer
- [Uni48] United Nations General Assembly (1948) Universal declaration of human rights. Technical report research 217 A (III), United Nations General Assembly. [hereinafter “UDHR”].
- [UPM18] Unanue JIA, Papadopoulos Y, Merle G (2018) Explicit modelling and treatment of repair in prediction of dependability. *IEEE Trans Depend Secure Comput*, pp 1–16
- [VJK16] Volk M, Junges S, Katoen J-P (2016) Advancing dynamic fault tree analysis – get succinct state spaces fast and synthesise failure rates. In: *Computer safety, reliability, and security (SAFECOMP)*, 35th international conference, pp 253–265
- [War12] Warburton N (2012) *Philosophy: the basics*. Taylor & Francis

Received 15 September 2020

Accepted in revised form 13 March 2021 by Cliff Jones