

CONCLUDING REMARKS

6.1 Thesis Contributions

This thesis covered the following two themes: i) evaluating perception models using metrics that are relevant to system-level specifications as well as the downstream planning logic, and ii) synthesis of reactive test environments and strategies.

Task-Relevant Evaluation of Perception: In this thesis, we considered the problem of evaluating the object detection and classification task of perception given system-level specifications. We identified confusion matrices as an appropriate model of sensor error for the object detection and classification task, and using principles of automata theory and probabilistic model checking, we formally defined probability functions to relate the confusion matrix to a probabilistic model of system state evolution. Confusion matrices are a popular choice for comparing and evaluating detection models in computer vision. This work is the first to formally establish a link between confusion matrices and system-level probability of satisfying a temporal logic specification. Qualitatively, our theoretical approach matches empirical observations in experimental work conducted in industry [13]. Furthermore, our approach lends a quantitative framework for designers to choose appropriate detection models based on their specifications. For instance, the precision-recall tradeoff which is well-known in detection tasks, is manifested in the system-level performance, and is quantified in the form of probabilistic guarantees. Due to this, we can compute desired lower bounds on detection performance (e.g., lower bounds on precision, false negative rate etc.) from desired quantitative system guarantee. We did this as a case study using the system design optimization tool, Pacti [68].

Based on these theoretical fundamentals, the second contribution in this direction is proposing new metrics for evaluating detection models that are more relevant to the system-level specification and the downstream controller. For this, we introduced proposition-labeled confusion matrices, in which traditional class labels are replaced by propositional formulas that are relevant to controller design. Furthermore, evaluations can be grouped at the same level of abstraction as the downstream controller that receives these detections as input. We evaluated a pre-trained

Pointpillars model that detects objects based on LiDaR data on the entire nuScenes dataset, and illustrate the result for a car-pedestrian example.

Reactive Test Synthesis: Automated test synthesis is a technical challenge motivated by the need for certification of safety-critical autonomous systems. These systems are expected to reason over both discrete and continuous states and inputs. This thesis studies synthesis of reactive test plans for high-level decision-making over discrete states and inputs. Specifications of the system are encoded in the system objective. In addition, user-defined specification of desired test behavior are encoded in the test objective, which is not revealed to the system under test. In this thesis, test objectives are manually specified. However, there is potential for automating this process as discussed in the future directions section later in the chapter.

In this thesis, we covered a test synthesis framework to restrict system actions reactively via a test harness. These restrictions can be implemented by the test environment using static and/or reactive obstacles, and dynamic test agents. First, we construct a product graph that tracks the system dynamics, and realization of the system and test objectives. Effectively, a path on the product graph represents a test execution. The routing problem is formulated as an optimization in which the test execution to realize the test objective without making it impossible for the system to realize the system objective. Via a reduction from 3-SAT, the computational complexity of this routing problem is shown to be NP-hard. This thesis covers two main approaches to solve the routing problem: Stackelberg game with coupled constraints and a mixed-integer linear program. The mixed-integer linear program can be solved more efficiently, and with guarantees that a feasible solution is a feasible test strategy. For different test environment types, the mixed-integer program can be easily modified by adding linear constraints to account for different environment types and to exclude any solutions. Static obstacles can be implemented as a special case of the reactive setting by adding constraints to enforce the non-reactivity of these restrictions. Furthermore, the dynamic agent strategy is synthesized to realize the reactive test strategy by matching the optimization solution.

Finally, we conducted hardware experiments using a pair of quadrupedal robots. The framework is agnostic to the specific controllers at the lower levels of the control stack, thus illustrating that the high-level tests synthesized by this framework can be effectively translated to hardware with test environments comprising of static and reactive obstacles, and dynamic test agents.

6.2 Future Directions

There are several exciting future directions for research on specification, testing, and verification of autonomous cyber-physical systems, guided by compelling demonstrations in hardware and simulation.

Layered, hierarchical test synthesis

Oftentimes, system-level failures in complex systems emerge from poor interfaces and interactions between subsystems. Current approaches to identifying failures with respect to specifications relies on black-box optimization methods, which are typically limited to identifying input signals in the continuous domain. While there is some work on identifying discrete-valued test inputs, it is often limited to variables that remain constant throughout the test (e.g., color of objects, the decision to place static barriers in the scene).

Figure 6.1 illustrates the vertical stack of the planning and control modules. The high-level planner, which operates at a slower timescale, is responsible for long-horizon decision-making which involves reasoning over fundamentally discrete variables. At the mid-level, trajectories with waypoints are planned for the robotic system. Finally the low-level controller, operating at faster speeds, executes the mid-level plan. In this thesis, we studied test environment synthesis for the high-level planner. The falsification approaches to identifying test cases are traditionally used to find failures at the mid-level and low-level. Furthermore, falsification algorithms often output open-loop trajectories, instead of reactive test strategies.

An open question is to identify falsifying instances resulting from a *combination* of poor *high-level decision making* together with continuous *nonlinear dynamics* at the low-level (e.g., incorrectly switching to a different dynamical mode, causing the system to violate safety or progress specifications). This is non-trivial even in simple hybrid system examples, especially when the system architecture and control design are black-box to the tester, and attempts to identify these fail-

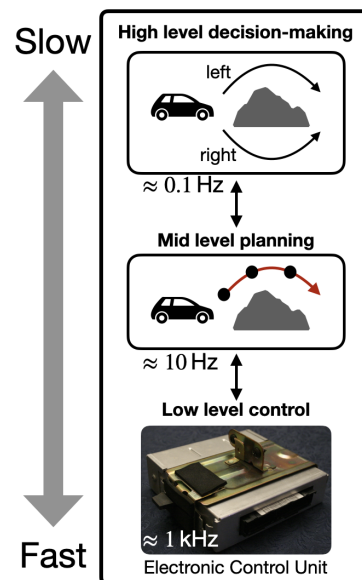


Figure 6.1: Overview of the planning and control software stack.

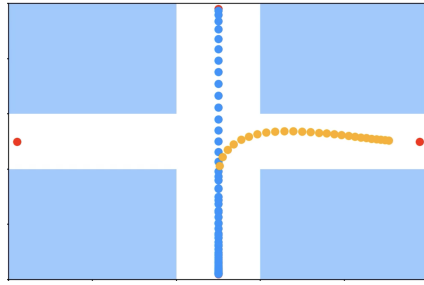
ures by jointly searching over discrete and continuous parameters will not scale. To address this challenge, we would need to i) infer how high-level commands affect continuous dynamics, and ii) infer when switches to different high-level modes result in dynamically unsafe trajectories. In addition to falsifying components at different levels of the control stack, we would also need to falsify the interfaces that map between them.

As an example, consider a simple switched system shown in Figure 6.2. The system in this example is a point-mass which must avoid the unsafe regions shaded in blue, and has two operating modes: a north-south mode, and an east-west mode. The tester has access to a switch command: when a switch command is issued, the system must eventually switch to the other operating mode. The system is a 2-dimensional double integrator, and has a simple model-predictive controller with a quadratic cost consisting of: i) position error with respect to the north-south and east-west axes, ii) control effort, and iii) terminal cost defined on the system state. The terminal cost function is a 2-norm of the distance to the goal (north/south/east/west) with zero velocity along either axis. The optimization includes constraints on control effort that can be exerted in each direction.

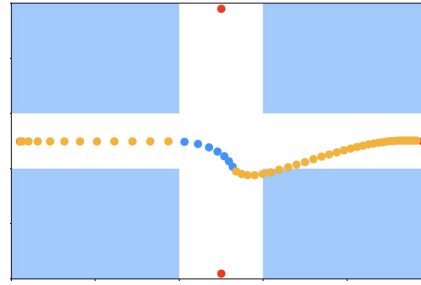
With no knowledge of the control design, and a couple of trials of commanding a single switch, the trajectories of this system would indicate a safe implementation. For example, see Figure 6.2a, in which the system safely avoids the unsafe region in executing the switch command. However, two switch commands in quick succession result in a falsifying trajectory of the system, highlighting the flaw in controller. In this example, the decision to switch twice is a fundamentally discrete choice, which current falsification algorithms typically cannot handle. The role of reactivity and layered architecture also becomes evident: i) the second switch is command in reaction to the system response to the first switch, and ii) the dynamical behavior (low-level behavior) of the system in response to the switch command combined with its decision to switch modes immediately without regard for safety (high-level decision) is what ultimately results in the failing trace.

Incorporating low-level dynamics

A challenge in hierarchical test and evaluation is in identifying the appropriate surrogate model of the system on which to test. As discussed previously, suppose the system-under-test is black-box to the test engineer, that is, the planning and control architecture and implementation is unknown to the tester. This would also imply



(a) Point mass (system) starting from the north position and in north-south mode (in blue), and commanded to switch to east-west mode (in orange).



(b) Point mass (system) starting from the west position and in east-west mode (in orange), and then commanded to switch twice in quick succession. The blue portion of the trajectory indicates the system switching to north-south mode before reverting to the east-west mode as shown in orange.

Figure 6.2: Simple switched system example, with position of the system shown at discrete time intervals.

that the system models used by the designers is also unknown to the test engineer. Therefore, the only entities known to the test engineer are the system specifications, the operational domain, and a black-box simulator or the physical system.

Though specifications are defined on the overall system, the subsystems responsible for the satisfaction of these specifications depends partly on the system implementation. For example, the requirement that the robot must remain safe likely requires multiple subsystems working together to satisfy safety. However, the requirement that the system exhibit certain motion primitives, e.g., quadruped must walk at a certain speed, might be the responsibility of a low-level controller. Depending on the scenario, certain specifications become prerequisites for other specifications. For example, the safety specification of the quadruped requires to evade an adversary might require it to consistently execute the walk motion primitive at a certain speed. For a broadly defined scenario such as evasion, if we are able to automatically order specifications according to this notion of pre-requisites, we can use the pre-requisite specifications to construct a model of high-level behavior of the system. In the literature on hybrid cyber-physical systems, this is related to work on constructing high-level system abstractions from low-level controllers using tools including reachability analysis. A concrete first step would be to review the literature on abstraction for control synthesis [130–133], and leveraging these methods to build a similar paradigm for abstraction, and discretization for testing.

The previous subsection is largely concerned with the case in which an abstraction of the system is assumed, and we need an efficient approach to testing that includes reasoning over both discrete and continuous variables. In this subsection, we are concerned with finding abstractions suitable for testing. There are two high-level directions for future work. First, how can we use prerequisite specifications to gather data from the system, and construct abstractions to model high-level behavior? These abstractions can include quantitative metrics of difficulty of executing lower level motion primitives. Second, what are the fundamental limitations of constructing these abstractions given the black-box nature of the system? Constructing these abstractions would potentially require a combination of model-based and data-driven methods — model-based in the sense that the physical dynamics, but not the control implementation, of the system might be known to the tester, and data-driven to get statistical data on the specific system implementation.

Criticality, coverage, and compositionality of test plans

Identifying critical test objectives is ill-posed when the operational design domain (ODD) is vast and difficult to characterize. While it might be hard to define a critical scenario in general, can we comparatively evaluate the criticality of two test scenarios? We can begin by studying the criticality of scenarios from a controls perspective (assuming perfect perception) before expanding to criticality from the system perspective (including all tasks pertaining to perception, reasoning, and control). There are two contending perspectives. On the one hand, criticality of a scenario will depend partly on the system controller — a scenario that is challenging (e.g., in terms of control effort, optimality, robustness) for one controller need not be equally challenging for another. Yet, at the same time, some scenarios seem universally less critical than others. For example, a safe unprotected left turn at a T-intersection amidst busy two-way traffic is more compelling than taking the same unprotected left turn without traffic. When are scenarios comparable? How do we quantify comparative criticality, and can we identify a class of controllers for which one scenario is more critical than the other?

The aforementioned question also relates to coverage. Ideally, successfully passing a more critical test should imply high confidence that the system would pass the less critical test. We would need to define a coverage metric that is consistent with this notion of criticality while also capturing the diversity of possible scenarios that are not easily comparable. For example, would reactive test scenarios (e.g., test agents moving in an office space) cover “open-loop” tests in which the test environment is

completely static (e.g., static yet cluttered office environment)?

Another direction to tackle the coverage question is to decompose it to the subsystem level as opposed to the scenario level. Since the ODD is often vast and difficult to define, we can focus on coverage for inputs to various subsystems in the control stack, which can be relatively low-dimensional in comparison to multi-modal sensor data received by the perception system. Can we then *compose* coverage guarantees from testing the individual subsystems during development to infer coverage at the system-level? Can we rely on this analysis to identify operational tests that check multiple unit-level tests at once?

Task-relevant metrics for perception

While this thesis identifies task-relevant metrics for object detection and classification tasks, corresponding metrics for other perception tasks such as tracking and behavior prediction still remain to be studied. Secondly, it is not clear which metrics offer the tightest system-level guarantees. This thesis offers preliminary results — confusion matrices chosen based on system-level specifications and downstream control logic result in less conservative evaluations overall. However, further research on investigating the tightness of these guarantees needs to be studied, along with hardware validation of the derived system-level guarantees.

Thirdly, one can extend these principles to perception-planning-control architectures where the interfaces between modules are less distinct. The confusion matrix only accounts for the final layer of the model's neural network, which outputs a scalar value to classify the object. However, higher dimensional learned features of the model could contain rich information on the model's performance which can be exploited. In which system-level architectures is the confusion matrix a sufficient metric of detection error? Which learned features of the detection models are a better representation of model performance with respect to the system-level task? Finally, these task-relevant evaluation criteria are often meant for offline evaluations of perception models, and the resulting system-level guarantees are limited to scenarios from the distribution used for model evaluations. Future work should study the how these guarantees can be updated via runtime monitoring or how they can degrade if a scenario is out-of-distribution.

Bibliography

- [1] Zoox, “Putting Zoox to the test: preparing for the challenges of the road,” 2021. <https://zoox.com/journal/structured-testing/>, Last accessed on 2024-04-11.
- [2] Waymo, “A blueprint for av safety: Waymo’s toolkit for building a credible safety case,” 2020. <https://waymo.com/blog/2023/03/a-blueprint-for-av-safety-waymos/#:~:text=A%20safety%20case%20for%20fully,evidence%20to%20support%20that%20determination.>, Last accessed on 2024-05-05.
- [3] F. Favaro, L. Fraade-Blanar, S. Schnelle, T. Victor, M. Peña, J. Engstrom, J. Scanlon, K. Kusano, and D. Smith, “Building a credible case for safety: Waymo’s approach for the determination of absence of unreasonable risk,” 2023. www.waymo.com/safety.
- [4] N. Kalra and S. M. Paddock, “Driving to safety: How many miles of driving would it take to demonstrate autonomous vehicle reliability?,” *Transportation Research Part A: Policy and Practice*, vol. 94, pp. 182–193, 2016.
- [5] N. Webb, D. Smith, C. Ludwick, T. Victor, Q. Hommes, F. Favaro, G. Ivanov, and T. Daniel, “Waymo’s safety methodologies and safety readiness determinations,” 2020.
- [6] I. S. Organization, “Road vehicles: Safety of the intended functionality (ISO Standard No. 21448:2022),” 2022. <https://www.iso.org/standard/77490.html>, Last accessed on 2024-04-11.
- [7] L. Li, W.-L. Huang, Y. Liu, N.-N. Zheng, and F.-Y. Wang, “Intelligence testing for autonomous vehicles: A new approach,” *IEEE Transactions on Intelligent Vehicles*, vol. 1, no. 2, pp. 158–166, 2016.
- [8] H. Winner, K. Lemmer, T. Form, and J. Mazzega, “Pegasus—first steps for the safe introduction of automated driving,” in *Road Vehicle Automation 5*, pp. 185–195, Springer, 2019.
- [9] “DARPA Urban Challenge.” <https://www.darpa.mil/about-us/timeline/darpa-urban-challenge>.
- [10] “Technical Evaluation Criteria.” <https://archive.darpa.mil/grandchallenge/rules.html>.
- [11] P. Koopman and M. Wagner, “Challenges in autonomous vehicle testing and validation,” *SAE International Journal of Transportation Safety*, vol. 4, no. 1, pp. 15–24, 2016.

- [12] J. Eskenazi and W. Jarett, “Explore: See the 55 reports — so far — of robot cars interfering with SF fire dept.,” 2023. <https://missionlocal.org/2023/08/cruise-waymo-autonomous-vehicle-robot-taxi-driverless-car-reports-san-francisco/>, Last accessed on 2024-04-11.
- [13] H. Zhao, S. K. Sastry Hari, T. Tsai, M. B. Sullivan, S. W. Keckler, and J. Zhao, “Suraksha: A framework to analyze the safety implications of perception design choices in avs,” in *2021 IEEE 32nd International Symposium on Software Reliability Engineering (ISSRE)*, pp. 434–445, 2021.
- [14] H. Kress-Gazit, G. E. Fainekos, and G. J. Pappas, “Temporal-logic-based reactive mission and motion planning,” *IEEE Transactions on Robotics*, vol. 25, no. 6, pp. 1370–1381, 2009.
- [15] M. Kloetzer and C. Belta, “A fully automated framework for control of linear systems from temporal logic specifications,” *IEEE Transactions on Automatic Control*, vol. 53, no. 1, pp. 287–297, 2008.
- [16] M. Lahijanian, S. B. Andersson, and C. Belta, “A probabilistic approach for control of a stochastic system from LTL specifications,” in *Proceedings of the 48th IEEE Conference on Decision and Control (CDC) held jointly with 2009 28th Chinese Control Conference*, pp. 2236–2241, IEEE, 2009.
- [17] V. Raman, A. Donzé, M. Maasoumy, R. M. Murray, A. Sangiovanni-Vincentelli, and S. A. Seshia, “Model predictive control with signal temporal logic specifications,” in *53rd IEEE Conference on Decision and Control*, pp. 81–87, IEEE, 2014.
- [18] T. Wongpiromsarn, U. Topcu, and R. M. Murray, “Receding horizon temporal logic planning,” *IEEE Transactions on Automatic Control*, vol. 57, no. 11, pp. 2817–2830, 2012.
- [19] G. Katz, C. Barrett, D. L. Dill, K. Julian, and M. J. Kochenderfer, “Reluplex: An efficient SMT solver for verifying deep neural networks,” in *International Conference on Computer Aided Verification*, pp. 97–117, Springer, 2017.
- [20] M. Fazlyab, M. Morari, and G. J. Pappas, “Probabilistic verification and reachability analysis of neural networks via semidefinite programming,” in *2019 IEEE 58th Conference on Decision and Control (CDC)*, pp. 2726–2731, IEEE, 2019.
- [21] M. Fazlyab, M. Morari, and G. J. Pappas, “Safety verification and robustness analysis of neural networks via quadratic constraints and semidefinite programming,” *IEEE Transactions on Automatic Control*, 2020.
- [22] H.-D. Tran, X. Yang, D. M. Lopez, P. Musau, L. V. Nguyen, W. Xiang, S. Bak, and T. T. Johnson, “NNV: The neural network verification tool for

- deep neural networks and learning-enabled cyber-physical systems,” in *International Conference on Computer Aided Verification*, pp. 3–17, Springer, 2020.
- [23] T. Dreossi, S. Jha, and S. A. Seshia, “Semantic adversarial deep learning,” in *International Conference on Computer Aided Verification*, pp. 3–26, Springer, 2018.
- [24] S. A. Seshia, A. Desai, T. Dreossi, D. J. Fremont, S. Ghosh, E. Kim, S. Shivakumar, M. Vazquez-Chanlatte, and X. Yue, “Formal specification for deep neural networks,” in *International Symposium on Automated Technology for Verification and Analysis*, pp. 20–34, Springer, 2018.
- [25] T. Dreossi, A. Donzé, and S. A. Seshia, “Compositional falsification of cyber-physical systems with machine learning components,” *Journal of Automated Reasoning*, vol. 63, no. 4, pp. 1031–1053, 2019.
- [26] S. Topan, K. Leung, Y. Chen, P. Tupekar, E. Schmerling, J. Nilsson, M. Cox, and M. Pavone, “Interaction-dynamics-aware perception zones for obstacle detection safety evaluation,” in *2022 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1201–1210, IEEE, 2022.
- [27] K. Chakraborty and S. Bansal, “Discovering closed-loop failures of vision-based controllers via reachability analysis,” *IEEE Robotics and Automation Letters*, vol. 8, no. 5, pp. 2692–2699, 2023.
- [28] A. Dokhanchi, H. B. Amor, J. V. Deshmukh, and G. Fainekos, “Evaluating perception systems for autonomous vehicles using quality temporal logic,” in *International Conference on Runtime Verification*, pp. 409–416, Springer, 2018.
- [29] A. Balakrishnan, A. G. Puranic, X. Qin, A. Dokhanchi, J. V. Deshmukh, H. B. Amor, and G. Fainekos, “Specifying and evaluating quality metrics for vision-based perception systems,” in *2019 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pp. 1433–1438, IEEE, 2019.
- [30] B. Bauchwitz and M. Cummings, “Evaluating the reliability of Tesla model 3 driver assist functions,” 2020.
- [31] H. Kress-Gazit, D. C. Conner, H. Choset, A. A. Rizzi, and G. J. Pappas, “Courteous cars,” *IEEE Robotics & Automation Magazine*, vol. 15, no. 1, pp. 30–38, 2008.
- [32] H. Kress-Gazit and G. J. Pappas, “Automatically synthesizing a planning and control subsystem for the DARPA Urban Challenge,” in *2008 IEEE International Conference on Automation Science and Engineering*, pp. 766–771, IEEE, 2008.

- [33] T. Wongpiromsarn, S. Karaman, and E. Frazzoli, “Synthesis of provably correct controllers for autonomous vehicles in urban environments,” in *2011 14th International IEEE Conference on Intelligent Transportation Systems (ITSC)*, pp. 1168–1173, IEEE, 2011.
- [34] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, “CARLA: An open urban driving simulator,” in *Conference on Robot Learning*, pp. 1–16, PMLR, 2017.
- [35] D. J. Fremont, T. Dreossi, S. Ghosh, X. Yue, A. L. Sangiovanni-Vincentelli, and S. A. Seshia, “Scenic: a language for scenario specification and scene generation,” in *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pp. 63–78, 2019.
- [36] Y. Annpureddy, C. Liu, G. Fainekos, and S. Sankaranarayanan, “S-taliro: A tool for temporal logic falsification for hybrid systems,” in *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pp. 254–257, Springer, 2011.
- [37] G. E. Fainekos and G. J. Pappas, “Robustness of temporal logic specifications for continuous-time signals,” *Theoretical Computer Science*, vol. 410, no. 42, pp. 4262–4291, 2009.
- [38] G. E. Fainekos, S. Sankaranarayanan, K. Ueda, and H. Yazarel, “Verification of automotive control applications using s-taliro,” in *2012 American Control Conference (ACC)*, pp. 3567–3572, IEEE, 2012.
- [39] S. Sankaranarayanan and G. Fainekos, “Falsification of temporal properties of hybrid systems using the cross-entropy method,” in *Proceedings of the 15th ACM international conference on Hybrid Systems: Computation and Control*, pp. 125–134, 2012.
- [40] S. Bak, S. Bogomolov, A. Hekal, N. Kochdumper, E. Lew, A. Mata, and A. Rahmati, “Falsification using reachability of surrogate koopman models,” in *Proceedings of the 27th ACM International Conference on Hybrid Systems: Computation and Control, HSCC ’24*, (New York, NY, USA), Association for Computing Machinery, 2024.
- [41] A. Donzé, “Breach, a toolbox for verification and parameter synthesis of hybrid systems,” in *International Conference on Computer Aided Verification*, pp. 167–170, Springer, 2010.
- [42] C. E. Tuncali, G. Fainekos, H. Ito, and J. Kapinski, “Simulation-based adversarial test generation for autonomous vehicles with machine learning components,” in *2018 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1555–1562, IEEE, 2018.

- [43] C. Menghi, P. Arcaini, W. Baptista, G. Ernst, G. Fainekos, F. Formica, S. Gon, T. Khandait, A. Kundu, G. Pedrielli, *et al.*, “Arch-comp 2023 category report: Falsification,” in *10th International Workshop on Applied Verification of Continuous and Hybrid Systems. ARCH23*, vol. 96, pp. 151–169, 2023.
- [44] T. Dreossi, D. J. Fremont, S. Ghosh, E. Kim, H. Ravanbakhsh, M. Vazquez-Chanlatte, and S. A. Seshia, “Verifai: A toolkit for the formal design and analysis of artificial intelligence-based systems,” in *International Conference on Computer Aided Verification*, pp. 432–442, Springer, 2019.
- [45] A. Corso, P. Du, K. Driggs-Campbell, and M. J. Kochenderfer, “Adaptive stress testing with reward augmentation for autonomous vehicle validation,” in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, pp. 163–168, IEEE, 2019.
- [46] S. Feng, H. Sun, X. Yan, H. Zhu, Z. Zou, S. Shen, and H. X. Liu, “Dense reinforcement learning for safety validation of autonomous vehicles,” *Nature*, vol. 615, no. 7953, pp. 620–627, 2023.
- [47] X. Qin, N. Arechiga, J. Deshmukh, and A. Best, “Robust testing for cyber-physical systems using reinforcement learning,” in *Proceedings of the 21st ACM-IEEE International Conference on Formal Methods and Models for System Design, MEMOCODE ’23*, (New York, NY, USA), p. 36–46, Association for Computing Machinery, 2023.
- [48] S. A. Seshia, D. Sadigh, and S. S. Sastry, “Toward verified artificial intelligence,” *Commun. ACM*, vol. 65, p. 46–55, jun 2022.
- [49] B. Johnson and H. Kress-Gazit, “Probabilistic analysis of correctness of high-level robot behavior with sensor error,” 2011.
- [50] A. Géron, *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems*. O’Reilly Media, 2019.
- [51] X. Wang, R. Li, B. Yan, and O. Koyejo, “Consistent classification with generalized metrics,” 2019.
- [52] P. Antonante, H. Nilsen, and L. Carlone, “Monitoring of perception systems: Deterministic, probabilistic, and learning-based fault detection and identification,” *arXiv preprint arXiv:2205.10906*, 2022.
- [53] M. Hekmatnejad, S. Yaghoubi, A. Dokhanchi, H. B. Amor, A. Shrivastava, L. Karam, and G. Fainekos, “Encoding and monitoring responsibility sensitive safety rules for automated vehicles in signal temporal logic,” in *Proceedings of the 17th ACM-IEEE International Conference on Formal Methods and Models for System Design*, pp. 1–11, 2019.

- [54] T. Wongpiromsarn and E. Frazzoli, “Control of probabilistic systems under dynamic, partially known environments with temporal logic specifications,” in *2012 IEEE 51st IEEE Conference on Decision and Control (CDC)*, pp. 7644–7651, 2012.
- [55] A. Badithela, T. Wongpiromsarn, and R. M. Murray, “Leveraging classification metrics for quantitative system-level analysis with temporal logic specifications,” in *2021 60th IEEE Conference on Decision and Control (CDC)*, (Austin, TX, USA (virtual)), pp. 564–571, IEEE, 2021.
- [56] C. S. Pasareanu, R. Mangal, D. Gopinath, S. G. Yaman, C. Imrie, R. Calinescu, and H. Yu, “Closed-loop analysis of vision-based autonomous systems: A case study,” *arXiv preprint arXiv:2302.04634*, 2023.
- [57] S. Beland, I. Chang, A. Chen, M. Moser, J. Paunicka, D. Stuart, J. Vian, C. Westover, and H. Yu, “Towards assurance evaluation of autonomous systems,” in *Proceedings of the 39th International Conference on Computer-Aided Design*, pp. 1–6, 2020.
- [58] Y. V. Pant, H. Abbas, K. Mohta, R. A. Quaye, T. X. Nghiem, J. Devietti, and R. Mangharam, “Anytime computation and control for autonomous systems,” *IEEE Transactions on Control Systems Technology*, vol. 29, no. 2, pp. 768–779, 2021.
- [59] P. Karkus, B. Ivanovic, S. Mannor, and M. Pavone, “Diffstack: A differentiable and modular control stack for autonomous vehicles,” in *Proceedings of The 6th Conference on Robot Learning* (K. Liu, D. Kulic, and J. Ichnowski, eds.), vol. 205 of *Proceedings of Machine Learning Research*, pp. 2170–2180, PMLR, 14–18 Dec 2023.
- [60] C. Baier and J.-P. Katoen, *Principles of model checking*. MIT press, 2008.
- [61] O. Koyejo, N. Natarajan, P. Ravikumar, and I. S. Dhillon, “Consistent multilabel classification,” in *NeurIPS*, vol. 29, (Palais des Congrès de Montréal, Montréal CANADA), pp. 3321–3329, Advances in Neural Information Processing Systems, 2015.
- [62] M. Kwiatkowska, G. Norman, and D. Parker, “Prism 4.0: Verification of probabilistic real-time systems,” in *International conference on computer aided verification*, pp. 585–591, Springer, 2011.
- [63] C. Dehnert, S. Junges, J.-P. Katoen, and M. Volk, “A Storm is coming: A modern probabilistic model checker,” in *International Conference on Computer Aided Verification*, pp. 592–600, Springer, 2017.
- [64] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liang, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, “nusenes: A multimodal dataset for autonomous driving,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11621–11631, 2020.

- [65] A. H. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang, and O. Beijbom, “Pointpillars: Fast encoders for object detection from point clouds,” in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, (Los Alamitos, CA, USA), pp. 12689–12697, IEEE Computer Society, jun 2019.
- [66] M. Contributors, “MMDetection3D: OpenMMLab next-generation platform for general 3D object detection.” <https://github.com/open-mmlab/mmdetection3d>, 2020.
- [67] S. Gupta, J. Kanjani, M. Li, F. Ferroni, J. Hays, D. Ramanan, and S. Kong, “Far3det: Towards far-field 3d detection,” in *2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, (Los Alamitos, CA, USA), pp. 692–701, IEEE Computer Society, jan 2023.
- [68] I. Incer, A. Badithela, J. Graebener, P. Mallozzi, A. Pandey, S.-J. Yu, A. Benveniste, B. Caillaud, R. M. Murray, A. Sangiovanni-Vincentelli, *et al.*, “Pacti: Scaling assume-guarantee reasoning for system analysis and design,” *arXiv preprint arXiv:2303.17751*, 2023.
- [69] A. Badithela, T. Wongpiromsarn, and R. M. Murray, “Evaluation metrics of object detection for quantitative system-level analysis of safety-critical autonomous systems,” in *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, (Detroit, MI, USA), p. To Appear., IEEE, 2023.
- [70] A. Donzé and O. Maler, “Robust satisfaction of temporal logic over real-valued signals,” in *International Conference on Formal Modeling and Analysis of Timed Systems*, pp. 92–106, Springer, 2010.
- [71] E. Plaku, L. E. Kavradi, and M. Y. Vardi, “Falsification of ltl safety properties in hybrid systems,” *International Journal on Software Tools for Technology Transfer*, vol. 15, no. 4, pp. 305–320, 2013.
- [72] G. Chou, Y. E. Sahin, L. Yang, K. J. Rutledge, P. Nilsson, and N. Ozay, “Using control synthesis to generate corner cases: A case study on autonomous driving,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 37, no. 11, pp. 2906–2917, 2018.
- [73] T. Wongpiromsarn, M. Ghasemi, M. Cubuktepe, G. Bakirtzis, S. Carr, M. O. Karabag, C. Neary, P. Gohari, and U. Topcu, “Formal methods for autonomous systems,” *arXiv preprint arXiv:2311.01258*, 2023.
- [74] G. Fainekos, H. Kress-Gazit, and G. Pappas, “Hybrid controllers for path planning: A temporal logic approach,” in *Proceedings of the 44th IEEE Conference on Decision and Control*, pp. 4885–4890, 2005.

- [75] R. Majumdar, A. Mathur, M. Pirron, L. Stegner, and D. Zufferey, “Paracosm: A language and tool for testing autonomous driving systems,” *arXiv preprint arXiv:1902.01084*, 2019.
- [76] L. Tan, O. Sokolsky, and I. Lee, “Specification-based testing with linear temporal logic,” in *Proceedings of the 2004 IEEE International Conference on Information Reuse and Integration, 2004. IRI 2004.*, pp. 493–498, IEEE, 2004.
- [77] G. Fraser and F. Wotawa, “Using LTL rewriting to improve the performance of model-checker based test-case generation,” in *Proceedings of the 3rd International Workshop on Advances in Model-Based Testing*, pp. 64–74, 2007.
- [78] G. Fraser and P. Ammann, “Reachability and propagation for LTL requirements testing,” in *2008 The Eighth International Conference on Quality Software*, pp. 189–198, IEEE, 2008.
- [79] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to algorithms*. MIT press, 2009.
- [80] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to algorithms*. MIT press, 2022.
- [81] C. Menghi, C. Tsigkanos, P. Pelliccione, C. Ghezzi, and T. Berger, “Specification patterns for robotic missions,” *IEEE Transactions on Software Engineering*, vol. 47, no. 10, pp. 2208–2224, 2019.
- [82] R. Bloem, G. Fey, F. Greif, R. Könighofer, I. Pill, H. Riener, and F. Röck, “Synthesizing adaptive test strategies from temporal logic specifications,” *Formal methods in system design*, vol. 55, no. 2, pp. 103–135, 2019.
- [83] J. Tretmans, “Conformance testing with labelled transition systems: Implementation relations and test generation,” *Computer Networks and ISDN Systems*, vol. 29, no. 1, pp. 49–79, 1996.
- [84] B. K. Aichernig, H. Brandl, E. Jöbstl, W. Krenn, R. Schlick, and S. Tiran, “Killing strategies for model-based mutation testing,” *Software Testing, Verification and Reliability*, vol. 25, no. 8, pp. 716–748, 2015.
- [85] R. Hierons, “Applying adaptive test cases to nondeterministic implementations,” *Information Processing Letters*, vol. 98, no. 2, pp. 56–60, 2006.
- [86] A. Petrenko and N. Yevtushenko, “Adaptive testing of nondeterministic systems with FSM,” in *2014 IEEE 15th International Symposium on High-Assurance Systems Engineering*, pp. 224–228, IEEE, 2014.
- [87] A. Pnueli and R. Rosner, “On the synthesis of a reactive module,” in *Proceedings of the 16th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pp. 179–190, 1989.

- [88] R. Bloem, B. Jobstmann, N. Piterman, A. Pnueli, and Y. Sa’ar, “Synthesis of reactive (1) designs,” *Journal of Computer and System Sciences*, vol. 78, no. 3, pp. 911–938, 2012.
- [89] M. Yannakakis, “Testing, optimization, and games,” in *Proceedings of the 19th Annual IEEE Symposium on Logic in Computer Science, 2004.*, pp. 78–88, IEEE, 2004.
- [90] L. Nachmanson, M. Veanes, W. Schulte, N. Tillmann, and W. Grieskamp, “Optimal strategies for testing nondeterministic systems,” *ACM SIGSOFT Software Engineering Notes*, vol. 29, no. 4, pp. 55–64, 2004.
- [91] A. David, K. G. Larsen, S. Li, and B. Nielsen, “Cooperative testing of timed systems,” *Electronic Notes in Theoretical Computer Science*, vol. 220, no. 1, pp. 79–92, 2008.
- [92] E. Bartocci, R. Bloem, B. Maderbacher, N. Manjunath, and D. Ničković, “Adaptive testing for specification coverage in CPS models,” *IFAC-PapersOnLine*, vol. 54, no. 5, pp. 229–234, 2021.
- [93] T. Marcucci, J. Umenberger, P. Parrilo, and R. Tedrake, “Shortest paths in graphs of convex sets,” *SIAM Journal on Optimization*, vol. 34, no. 1, pp. 507–532, 2024.
- [94] T. Marcucci, M. Petersen, D. von Wrangel, and R. Tedrake, “Motion planning around obstacles with convex optimization,” *Science Robotics*, vol. 8, no. 84, p. eadf7843, 2023.
- [95] H. Zhang, M. Fontaine, A. Hoover, J. Togelius, B. Dilkina, and S. Nikolaidis, “Video game level repair via mixed integer linear programming,” in *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, vol. 16, pp. 151–158, 2020.
- [96] M. Fontaine, Y.-C. Hsu, Y. Zhang, B. Tjanaka, and S. Nikolaidis, “On the Importance of Environments in Human-Robot Coordination,” in *Proceedings of Robotics: Science and Systems*, (Virtual), July 2021.
- [97] J. R. Büchi, *On a Decision Method in Restricted Second Order Arithmetic*, pp. 425–435. New York, NY: Springer New York, 1990.
- [98] A. Duret-Lutz, A. Lewkowicz, A. Fauchille, T. Michaud, É. Renault, and L. Xu, “Spot 2.0 — a framework for ltl and omega-automata manipulation,” in *Automated Technology for Verification and Analysis* (C. Artho, A. Legay, and D. Peled, eds.), (Cham), pp. 122–129, Springer International Publishing, 2016.
- [99] F. Fuggitti, “Ltl2dfa,” June 2020.

- [100] S. Bansal, Y. Li, L. Tabajara, and M. Vardi, “Hybrid compositional reasoning for reactive synthesis from finite-horizon specifications,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 9766–9774, Apr. 2020.
- [101] N. Klarlund and A. Møller, *MONA Version 1.4 User Manual*. BRICS, Department of Computer Science, University of Aarhus, January 2001. Notes Series NS-01-1. Available from <http://www.brics.dk/mona/>.
- [102] D. Goktas and A. Greenwald, “Convex-concave min-max Stackelberg games,” *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [103] I. Tsaknakis, M. Hong, and S. Zhang, “Minimax problems with coupled linear constraints: computational complexity, duality and solution methods,” *arXiv preprint arXiv:2110.11210*, 2021.
- [104] M. L. Bynum, G. A. Hackebeil, W. E. Hart, C. D. Laird, B. L. Nicholson, J. D. Sirola, J.-P. Watson, and D. L. Woodruff, *Pyomo—optimization modeling in python*, vol. 67. Springer Science & Business Media, third ed., 2021.
- [105] V. V. Vazirani, *Approximation algorithms*, vol. 1. Springer, 2001.
- [106] M. Fischetti and M. Monaci, “A branch-and-cut algorithm for mixed-integer bilinear programming,” *European Journal of Operational Research*, vol. 282, no. 2, pp. 506–514, 2020.
- [107] J. B. Graebener, A. S. Badithela, D. Goktas, W. Ubellacker, E. V. Mazumdar, A. D. Ames, and R. M. Murray, “Flow-based synthesis of reactive tests for discrete decision-making systems with temporal logic specifications,” *arXiv preprint arXiv:2404.09888*, 2024.
- [108] T. Wongpiromsarn, U. Topcu, N. Ozay, H. Xu, and R. M. Murray, “Tulip: a software toolbox for receding horizon temporal logic planning,” in *Proceedings of the 14th international conference on Hybrid systems: computation and control*, pp. 313–314, 2011.
- [109] I. Filippidis, S. Dathathri, S. C. Livingston, N. Ozay, and R. M. Murray, “Control design for hybrid systems with tulip: The temporal logic planning toolbox,” in *2016 IEEE Conference on Control Applications (CCA)*, pp. 1030–1041, IEEE, 2016.
- [110] S. Maoz and J. O. Ringert, “Gr (1) synthesis for ltl specification patterns,” in *Proceedings of the 2015 10th joint meeting on foundations of software engineering*, pp. 96–106, 2015.
- [111] S. A. Cook, “The complexity of theorem-proving procedures,” in *Logic, Automata, and Computational Complexity: The Works of Stephen A. Cook*, pp. 143–152, 2023.

- [112] C. H. Papadimitriou, *Computational complexity*, p. 260–265. GBR: John Wiley and Sons Ltd., 2003.
- [113] W. Ubellacker and A. D. Ames, “Robust locomotion on legged robots through planning on motion primitive graphs,” in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 12142–12148, 2023.
- [114] Gurobi Optimization, LLC, “Gurobi Optimizer Reference Manual,” 2023.
- [115] E. W. Dijkstra, “Guarded commands, nondeterminacy and formal derivation of programs,” *Communications of the ACM*, vol. 18, no. 8, pp. 453–457, 1975.
- [116] L. Lamport, “win and sin: Predicate transformers for concurrency,” *ACM Transactions on Programming Languages and Systems (TOPLAS)*, vol. 12, no. 3, pp. 396–428, 1990.
- [117] B. Meyer, “Applying ‘design by contract’,” *Computer*, vol. 25, no. 10, pp. 40–51, 1992.
- [118] A. Benveniste, B. Caillaud, A. Ferrari, L. Mangeruca, R. Passerone, and C. Sofronis, “Multiple viewpoint contract-based specification and design,” in *Formal Methods for Components and Objects: 6th International Symposium, FMCO 2007, Amsterdam, The Netherlands, October 24-26, 2007, Revised Lectures* (F. S. de Boer, M. M. Bonsangue, S. Graf, and W.-P. de Roever, eds.), (Berlin, Heidelberg), pp. 200–225, Springer Berlin Heidelberg, 2008.
- [119] A. L. Sangiovanni-Vincentelli, W. Damm, and R. Passerone, “Taming Dr. Frankenstein: Contract-based design for cyber-physical systems,” *Eur. J. Control*, vol. 18, no. 3, pp. 217–238, 2012.
- [120] P. Nuzzo, A. L. Sangiovanni-Vincentelli, D. Bresolin, L. Geretti, and T. Villa, “A platform-based design methodology with contracts and related tools for the design of cyber-physical systems,” *Proceedings of the IEEE*, vol. 103, no. 11, pp. 2104–2132, 2015.
- [121] I. Incer, *The Algebra of Contracts*. PhD thesis, EECS Department, University of California, Berkeley, May 2022.
- [122] A. Benveniste, B. Caillaud, D. Nickovic, R. Passerone, J.-B. Racllet, P. Reinkemeier, A. L. Sangiovanni-Vincentelli, W. Damm, T. A. Henzinger, K. G. Larsen, *et al.*, “Contracts for system design,” *Foundations and Trends in Electronic Design Automation*, vol. 12, no. 2-3, pp. 124–400, 2018.
- [123] I. Incer, A. L. Sangiovanni-Vincentelli, C.-W. Lin, and E. Kang, “Quotient for assume-guarantee contracts,” in *16th ACM-IEEE International Conference on Formal Methods and Models for System Design, MEMOCODE’18*, pp. 67–77, October 2018.

- [124] R. Passerone, Í. Íncer Romeo, and A. L. Sangiovanni-Vincentelli, “Coherent extension, composition, and merging operators in contract models for system design,” *ACM Transactions on Embedded Computing Systems (TECS)*, vol. 18, no. 5s, pp. 1–23, 2019.
- [125] R. Negulescu, “Process spaces,” in *CONCUR 2000 — Concurrency Theory* (C. Palamidessi, ed.), (Berlin, Heidelberg), pp. 199–213, Springer Berlin Heidelberg, 2000.
- [126] J. B. Graebener^{*}, A. Badithela^{*}, and R. M. Murray, “Towards better test coverage: Merging unit tests for autonomous systems,” in *NASA Formal Methods* (J. V. Deshmukh, K. Havelund, and I. Perez, eds.), (Cham), pp. 133–155, Springer International Publishing, 2022. A. Badithela and J.B. Graebener contributed equally to this work.
- [127] R. Bloem, B. Könighofer, R. Könighofer, and C. Wang, “Shield synthesis,” in *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pp. 533–548, Springer, 2015.
- [128] L. Kocsis and C. Szepesvári, “Bandit based monte-carlo planning,” in *European conference on machine learning*, pp. 282–293, Springer, 2006.
- [129] I. Incer, L. Mangeruca, T. Villa, and A. Sangiovanni-Vincentelli, “The quotient in preorder theories,” *arXiv:2009.10886*, 2020.
- [130] O. Hussien, A. Ames, and P. Tabuada, “Abstracting partially feedback linearizable systems compositionally,” *IEEE Control Systems Letters*, vol. 1, no. 2, pp. 227–232, 2017.
- [131] P. Tabuada, G. J. Pappas, and P. Lima, “Composing abstractions of hybrid systems,” in *International Workshop on Hybrid Systems: Computation and Control*, pp. 436–450, Springer, 2002.
- [132] S. Coogan and M. Arcaç, “Efficient finite abstraction of mixed monotone systems,” in *Proceedings of the 18th International Conference on Hybrid Systems: Computation and Control*, HSCC ’15, (New York, NY, USA), p. 58–67, Association for Computing Machinery, 2015.
- [133] J. Liu and N. Ozay, “Abstraction, discretization, and robustness in temporal logic control of dynamical systems,” in *Proceedings of the 17th international conference on Hybrid systems: computation and control*, pp. 293–302, 2014.