# Research Announcement TRI-01: Joint TRI-University Projects

## REVISION HISTORY

| Version | Date | Author | Description |
|---------|------|--------|-------------|
| 1.0 | 12/16/2019 | Eric Krotkov | Created |
| 2.0 | 4/22/2020 | Eric Krotkov | Updated Program Schedule (Section 3) and Awards (Section 4) to show start date of 4/1/2021. Updated Schedule (Section 3) and Awards (Section 4) with information about the annual review cycle. Updated Section 5.4 (proposal format) with multiple changes. |
| 2.1 | 4/23/2020 | Eric Krotkov | Updated Awards (Section 4) to discuss exercising options. Updated Cost Proposal (Section 5.4.4) to say that personnel make definite commitment. |
| 2.2 | 4/28/2020 | Eric Krotkov | Copy editing throughout |

## TABLE OF CONTENTS

# 1   INTRODUCTION

From its inception, TRI has engaged with university partners to conduct sponsored research in artificial intelligence.  From 2016 through 2020, TRI has supported work at Michigan, MIT, and Stanford.

Beginning in 2021, TRI will continue to support those three institutions, and in addition will support other institutions.  This document announces research opportunities starting in 2021 and solicits participation by an expanded university partner base.  The Frequently Asked Questions (FAQ) document on the program website provides more information about this research announcement.

A separate document, Research Announcement TRI-02, is available only by invitation and announces research opportunities for a separate solicitation targeting early stage researchers.

# 2   PROGRAM DESCRIPTION

The objectives of the Joint TRI-University program include the following:
- Add significant new knowledge and understanding to the field of artificial intelligence;
- Demonstrate the potential to radically advance areas beyond the state of the art;
- Promote the transfer of knowledge through the meaningful exchange of scientific and technical information between TRI researchers and university partners; and
- Create the potential for the creation and sharing of community infrastructure, including data and software, to further research, promote reproducibility, and support education.

Proposed topics should align with and contribute to the TRI mission, including automated driving, home robotics, and machine assisted cognition.  In exceptional cases, TRI may consider topics with high potential for revolutionary advance despite not aligning strongly with the TRI mission.

A TRI researcher should serve as an active member of the research team (not simply a sponsor). In exceptional cases, the program may consider topics with high potential for revolutionary advance despite the absence of a TRI researcher as an active member.

# 3   PROGRAM SCHEDULE

The following table presents the overall program schedule.

| Date | Description |
|---|---|
| 12/16/2019 | Research Announcement (this document) |
| 1/1/2020 – 1/31/2020 | Open Discussion with TRI Researchers |
| 3/2/2020 | Request for White Papers |
| 3/31/2020 | White Papers Due |
| 5/1/2020 | Request for Proposals |
| 5/29/2020 | Proposals Due |
| 6/30/2020 | Funding Decisions Announced |

| | |
|---|---|
| 7/1/2020 – 3/31/2021 | Principal Investigators hire students and staff (if needed) |
| 4/1/2021 | PI meeting for Year One at one location.  Funding begins |
| 9/1/2021 – 10/30/2021 | Site visits Year One (at university location) |
| 1/1/2021 – 1/15/2021 | Review of Year One Progress for every project at one location |
| 4/1/2022 | Kickoff meeting Year Two |

## 4   AWARDS

The award will take the form of an industry-sponsored research agreement.

The proposal will describe the time required to complete the proposed work.  TRI expects the typical duration to be 2-3 years, structured as follows:

- Base:  One year (for example, 4/1/2021 – 3/31/2022)
- Option Period 1:  One additional year (for example, 4/1/2022 – 3/31/2023)
- Option Period 2:  One more additional year (for example, 4/1/2023 – 3/31/2024)

Funding for the entire proposed duration of the project is not guaranteed.  Based on the Review of Year One Progress (Section 3), TRI will choose whether to exercise Option 1.  The Review of Year One Progress meeting will focus on achievement of proposed milestones.  Projects that do not meet their proposed Year One milestones will need to justify their continuation for Year Two.

Historically, TRI has acted to discontinue projects based on non-performance.  TRI expects that it is the nature of high-risk research with challenging objectives that some approaches will be found to be less fruitful than anticipated, and the project may need to change approaches or directions.

## 5   APPLICATION

The application process will occur in three stages: (i) discussion with TRI researchers, (ii) if invited, white paper submissions, and (iii) if invited, proposal submissions.

White papers will outline the proposed technical approach.  Proposals will describe how the team will conduct the proposed research.  Proposals should not re-state technical discussion provided in the white paper.

### 5.1   ELIGIBILITY

The proposer must be employed at a North American educational institution.

### 5.2   OPEN DISCUSSION WITH TRI RESEARCHERS

During the month of January 2020, proposers may reach out to TRI Researchers identified in Section 7 to conduct an open discussion of a research topic.

One of the reasons that TRI requires relatively shorter white papers and proposals is that TRI believes that verbal, face-to-face discussion may be more effective than relatively longer documents. As the simplest, most effective path to a mutual understanding of a research project, TRI encourages proposers to engage with TRI Researchers.

## 5.3 FORMAT OF WHITE PAPER

Length: Cover Page plus one (1) page

Format: Body font no less than 11 point, figure font no less than 8 point, margins > 0.75 inches.

Contents of White Paper

- Cover Page: Proposal Title, University Principal Investigator Name, University Name, TRI Researcher Name (or state "None"), Proposer Team, Topic Identifier (this is the number of the topic in Section 7, or state "None")

- Description – What are you trying to do, in general terms with no jargon?

- Technical Objective – What are the technical goals, in specific and measurable terms?

- Technical Approach – What is the proposed approach to meet the TRI research need?

- Novelty and Innovation - What is novel and innovative about your approach?

## 5.4 FORMAT OF PROPOSAL

Length: Cover Page plus three (3) pages

Format: Body font no less than 11 point, figure font no less than 8 point, margins >= 0.75 inches.

There is no need for new technical discussions; the white papers covered this already.

### 5.4.1 Cover Page

- Proposal Title

- University Principal Investigator Name and University Name

- Other University Investigator(s) and University Name(s)

- TRI Researcher Name(s)

- White Paper Identifier (this is the number listed as the ID in the first column of the Proposals to Request list)

- Period of Performance – The start and end dates of the project (entire duration, not just Year One)

- Intent to Subcontract – Yes or no

### 5.4.2 Technical Proposal

- Work Breakdown – Identify the high-level tasks to be performed. There is no need to elaborate on or justify the technical tasks, because the white paper already covered them.

- Organizational Structure – Describe how team members will coordinate their work. Explain which team member does each high-level task listed in the Work Breakdown.

- Performance Metrics – Describe in specific and measurable terms how progress will be evaluated. Well-formulated performance metrics have an unambiguous meaning; a naive observer will be able to evaluate whether they have been achieved. Examples of well-formulated metrics include the following: precision, recall, intersection over union, maximum operating range, power consumption, probability of detection.

- Technical Milestones – Using the performance metrics, specify milestones to be reached at 6-month intervals.

- Cost – Identify the cost to perform the proposed research for each year over the proposed period of performance. In addition, identify the total cost.

### 5.4.3 Supplementary Materials (Not Included in the Proposal Page Limits)

Experience with the white papers demonstrated that some authors felt more comfortable providing supplementary materials. Proposals may contain one page of supplementary materials. In general, these will not be reviewed.

### 5.4.4 Cost Proposal (Not Included in the Proposal Page Limits)

- The cost breakdown should be prepared by the sponsored projects office of the University Principal Investigator using their own format. This breakdown will allocate costs to university-standard categories such as labor, other direct costs, capital equipment, and overhead. The cost breakdown documents do not count against the page limits.
- One university may subcontract to another university, however the one university serving as the prime contractor may not charge overhead on subcontractors (double overhead is not permitted).
- Personnel listed on the cover page (Section 5.4.1) are expected to make a definite commitment to participate actively in the proposed project, including making all reasonable efforts to attend the PI meeting, Site Visit, and Review Meeting.
- For costing purposes, assume the following requirements:
  o A PI meeting in San Francisco, California on Tuesday – Thursday of the first week of April
  o A site visit meeting at the university in September or October with 10 Toyota personnel
  o A Review meeting in Atlanta, Georgia on Tuesday – Thursday of the second week of January
  o Note that this assumption is for costing purposes, and that the locations of the PI meeting and the Review meeting are subject to change

## 5.5 SUBMISSION

To submit a white paper or proposal, submit a PDF file to tri-ra01-submit@tri.global. The submitted file must include all materials, including the cover page, the technical sections, the supplementary material (if any), and the cost breakdown.

## 5.6 QUESTIONS

To submit questions about the application process, send a message to tri-ra01-questions@tri.global

To view answers to the submitted questions about the application process, visit the FAQ on the program website.

# 6 EVALUATION CRITERIA

Proposals will be evaluated using the following criteria listed in descending order of importance.

1. Overall Scientific and Technical Merit. The proposed technical approach is innovative, feasible, and achievable.
2. Proposer Team. The proposed technical team has the expertise and experience to accomplish the proposed tasks.
3. Alignment with TRI Mission. A TRI researcher serves as an active member of the research team. The proposed work directly contributes to an automated driving, home robotics, or machine assisted cognition project at TRI.
4. Potential for Revolutionary Advance. The proposed approach improves performance by an order of magnitude or develops a disruptive technology that could doom current approaches.

A TRI Steering Committee will evaluate proposals. The Steering Committee will not include TRI Researchers who have proposed to the program.

# 7 TOPIC AREAS

The three following sections list the topic areas of interest. Each topic area includes the name and email address of the TRI researcher, for use in the discussion period (see table in Section 3). The order in which the topic areas appear is not significant and does not imply any kind of ranking.

The topic writeups have been written by the interested TRI Researchers, in their own voice and with their own style. As a result, the writeups exhibit variability in structure and layout. We view this as a feature not a "bug."

The groupings into Automated Driving, Robotics, and Machine Assisted Cognition reflect the organizational home of the topic author. It is entirely possible for a topic listed in, for example Driving, to be valuable for and conceptually relevant to Robotics or Machine Assisted Cognition.

# 8 TOPICS: AUTOMATED DRIVING

## 8.1 GUARANTEED GUARDIAN: PLANNING WITH SAFE SETS FOR INTERACTIVE AUTONOMOUS DRIVING

TRI Researcher:     Avinash Balachandran
Email Address:      Avinash@tri.global
TRI Thrust:          Driving

With the concept of the "Toyota Guardian" system in mind, in this proposal we focus on developing an on-line planning framework to allow the human driver broad flexibility in controlling their vehicle, while providing a guaranteed safe intervention in the case that the

human gives a potentially unsafe control command. By taking into account human-robot interactions in real-world driving scenarios, we propose to develop algorithms to characterize the safe behavior of a vehicle in the presence of other agents in terms of safe sets. We will consider two approaches to constructing safe sets: (1) game-theoretic reachability analysis and (2) probabilistic analysis. With the former approach the objective is to analyze safety and performance in the worst-case (with respect to the other agents) while, with the latter approach, the objective is to reason probabilistically with respect to the behavior of the other agents. For both approaches, we will study methodologies to properly trade off safety and efficiency, possibly through hybridization schemes. In both approaches we will stress a focus on mathematically guaranteed performance, on-line computational speed, and experimental verification in simulation and on hardware systems.

**I. Algorithmic Aspects of Reachability Set-Based Planning**
In order to make same decisions, it is necessary to determine states that a vehicle can operate in while adhering to a set of prescribed rules and constraints. We propose to obtain such safe states through formal verification tools such as Hamiltonian-Jacobi (HJ) reachability analysis. analysis provides a set of "keep out states" that, if breached, can eventually lead to a violation of the constraints. It also provides an optimal policy to enact to avoid this set of "keep out states" when a vehicle encounters the boundary of this set. For example, we are interested in determining the set of states in the joint space of all vehicles in the neighborhood of the ego vehicle that can lead to a collision for the ego vehicle despite the best efforts of the driver. A guardian system would enact the optimal policy to avoid this set, overriding the driver's commands if the vehicle is found to be on the set boundary.

Since HJ analysis is based on dynamic programming, current approaches are generally not computationally efficient enough to provide a solution for high-dimensional problems such as driving in the presence of multiple other vehicles. Therefore, to compute solutions online in real-time, we plan to explore how to (1) split the problem into a hierarchy of easier-to-solve subproblems, (2) parallelize computations to take advantage of multi-core hardware, and (3) rely on offline pre-computation where possible. HJ reachability analysis also typically assumes perfect state knowledge, which is never the case in autonomous driving. Therefore, we will investigate how to incorporate perception uncertainty, including uncertainty in deep learning perception modules, into HJ reachability to maintain safety despite compounded uncertainty along the perception-planning pipeline.

**II. Algorithmic Aspects of Probabilistic Set-Based Planning**
In some cases, deterministic reachability can be too conservative, in the sense that a state may be unsafe only under a sequence of highly unlikely control actions by the other vehicles in an ego vehicle's neighborhood. Reachability tools typically do not distinguish between "formally unsafe but a collision is highly unlikely" versus "formally unsafe and a collision is highly likely." Therefore, in parallel with our research on reachability, we will also explore probabilistic models of safety, which will give a finer resolution on the degree of safety of a state or a trajectory. By taking into account the stochasticity of the system, we will develop a dynamic occupancy mapping technique based on approximate Gaussian processes that can be run in real-time. Gaussian process regression is used to represent a nonlinear pattern through an infinite number of functions with associated uncertainties. These uncertainties about the maneuverable areas allow a vehicle to trade off between safety and conservativeness, allowing a

guardian system to choose the best action while incorporating to a risk tolerance for collision. We treat safety in this case as a probabilistic concept, either as a chance constraint (the probability of violating the constraint must be below a threshold), or as risk sensitivity (we penalize variance or some other notion of uncertainty directly in an objective function).

To incorporate uncertainties into occupancy mapping, we need to reason about other vehicles' beliefs and other vehicles' perception models. For this, we propose to develop a driver model through (possibly risk-sensitive) inverse reinforcement learning in simulation, which can later be transferred to the real world using state-of-the-art 'sim-to-real transfer' techniques. Also, as for HJ reachability, probabilistic methods for computing safe sets also tend to be computationally intensive, and not well-suited to on-line implementation. We therefore will investigate approximate methods to reasoning probabilistically that maintain provable conservativeness, while given the speed required for online execution. We will consider spatially limited GPs computed with proximity graphs, as well as sparse GPs computed using variational inference techniques.

### III. Computational and Numerical Aspects of Set-Based Planning

Set-valued and game-theoretic planning problems are often naturally expressed as nested bi-level or multi-level optimization problems in which the feasible set of an outer problem is determined by the solution of an inner constrained optimization problem. Such problems, even when they are convex, are difficult to solve numerically with standard algorithms. We propose the development of a dedicated solver that can be tailored to specific problem instances that occur in autonomous driving, such as planning collision-free trajectories in the presence of other drivers. To ensure that the algorithm is scalable and fast enough for real-time applications, sparsity structure and parallelism (e.g. across timesteps in a trajectory and other vehicles in collision checking calculations) will be exploited to take maximum advantage of multi-core processors. Similarly, HJ reachability and probabilistic safe-set computations can both be expressed as partial differential equations (PDEs), or partial difference equations in discrete time. In this case, one of the core challenges to computing safe sets is in the numerical aspects of solving these PDEs. While these issues are typically not treated together with analytical and algorithmic aspects of reachability, we believe that to provide online speed, a comprehensive approach incorporating numerics together with analysis and algorithmics is essential. We will investigate methods for simplifying, parallelizing, and compressing the solution of these PDEs, e.g. through "concentration" techniques which represent the PDE as a collection of samples, neural network techniques that compress the value function that is the solution of the PDE, and model reduction techniques that represent the PDE through a lower order set of ODEs. In all cases we will stress mathematical guarantees of these solution techniques together with practical computational speed.

### IV. Experiments and Hardware Implementation

We propose to pursue these research thrusts in the context of specific traffic scenarios that are known to be difficult cases for autonomous driving. For example, we will consider merging on a freeway in traffic, and making an unprotected left turn at an intersection. Our research will include theoretical and algorithmic results, as well as verification on simulation platforms such as CARLA, and scale hardware platforms we have developed in our labs. Finally, we will test our algorithms on full scale driving platforms such as the X1 test vehicle and TRI test vehicles.

## 8.2   AUTONOMOUS DRIFTING

TRI Researcher:        Avinash Balachandran
Email Address:        [Avinash@tri.global](mailto:Avinash@tri.global)
TRI Thrust:            Driving

Our goal is to make a fully autonomous tandem drift scenario akin to Formula Drift. Our previous work has help us converge on a definition of vehicle stability internally and has shown how technologies required for drifting can be applied to aggressive lane changes.

A comprehensive statement of research needs is in preparation.

## 8.3   INCORPORATING CONTROL BARRIER FUNCTIONS WITHIN STOCHASTIC MODEL PREDICTIVE CONTROL

TRI Researcher:        Brian Goldfain
Email Address:        [Brian.Goldfain@tri.global](mailto:Brian.Goldfain@tri.global)
TRI Thrust:            Driving

A key enabling technology for capable self-driving vehicles is the ability to plan a collision free path on a road while considering all relevant information about the environment, vehicle, and goals. Human drivers perform this task seamlessly as they interact with other drivers, pedestrians, and cyclists, all while balancing the rules of the road and progressing toward their goal. Current approaches to solving the planning problem have proven capable in limited operating domains, but as that domain expands and more information must be considered to make safe driving decisions, they tend to become computationally intractable, may fail to return a feasible solution at all, or simply cannot incorporate the requisite amount or type of information. New approaches are needed to handle the problem complexity and reason over possible futures in real time during the dynamic driving task to enable the next generation of automated vehicle capabilities.

When deciding how to drive, the planning layer of an automated vehicle reasons over its own expected motion, predicted motions of dynamic objects in the environment, road rules, and navigation objectives. Proposed approaches should be able to handle the complex interactions among multiple agents required as the operational design domain expands. To benefit from the scale of data collection from a fleet of vehicles, considerations should be taken to utilize large amounts of training data.

In addition to a wide scope of valid input data and predictive capabilities, a planner should constrain its requested behavior to what is physically achievable by the vehicle. As the complexity of the planning problem increases, it becomes more difficult to find a plan that satisfies all the task constraints. A typical solution is to relax constraints or deploy simplified models throughout the computational pipeline. While these approximations may enable a solution to be found, they can suffer from the previously mentioned problem where the requested maneuvers may be physically impossible. Proposed approaches should be able to balance model complexity with accuracy and monitor and potentially correct for any infeasible results.

Particular areas of interest are machine learning, online adaptation, nonlinear dynamics, generic cost criteria, guarantees on feasibility/reachability, control barrier functions, parallel

computing, and real time operation without specialized hardware. Integration and testing in traffic scenarios and on real platforms are essential.

## 8.4 BAYESIAN ADAPTIVE CONTROL AND FAST RE-PLANNING FOR SAFETY CRITICAL SYSTEMS

TRI Researcher:      Brian Goldfain
Email Address:       Brian.Goldfain@tri.global
TRI Thrust:          Driving

"In an automated vehicle software stack, the bottom two layers are typically the planning and control modules. Planning is responsible for deciding how a vehicle should move through the world to successfully accomplish its task. Control is then responsible for turning that plan into a sequence of steering, throttle, and brake commands that are executed by the vehicle. Splitting the computation allows decisions to be optimized over different time scales using methods suitable for the types of information consumed at each stage in the pipeline. Typically, planning operates on task and environment information such as navigation goals, traffic signage and rules, and prediction of dynamic objects to generate a new plan for the vehicle to follow at a rate of a few times per second. Controls uses that plan, and knowledge of the capabilities and constraints of the vehicle, to compute the steering, throttle, and brake commands that will realize the desired driving behavior, typically at a much faster rate.

In order to maintain performance guarantees within planning and control, the modules should maintain contracts over the information they ingest. Example contract components include required data rates, time horizons, a shared understanding of uncertainty, and properties of the solution itself such as feasibility. Ideally, there would also be leeway for a module to refine the decisions made by upstream modules according to relevant information, which may not be available to upstream modules. For example, the trajectory that a planner requests a controller to follow must be within the capabilities of the vehicle and should maintain comfort and other objectives whenever possible. However, the planner considers relevant task information, but only a simplified vehicle dynamics model, whereas the controller consumes a simplified task description and has a much more detailed dynamics model. In an ideal world, contract violations would never occur, but unpredictable events encountered while driving, design considerations, computational restrictions, and access to information at each level can cause the requests passed between modules to occasionally violate the agreed-upon contracts.

Proposed solutions should have the ability to monitor its inputs for contract violations from modules above the controller and overcome some amount of violation with a "best-effort" solution to continue driving safely. Exploring approaches for the control module to adapt a plan and continue driving in the face of violations will have a direct impact on driving performance and will allow an exploration of how contracts can be handled throughout the entire software stack. Proposed solutions must operate in real time in concert with the entire self-driving software stack to demonstrate a clear benefit in terms of on road driving performance, which will also be verified in computer simulations and closed course testing. An emphasis will be put on real world driving performance.

Particular areas of interest include adaptive control, machine learning, tube model predictive control, reachability, and quadratic programming.

## 8.5 OPTIMAL POLICY BEHAVIOR PLANNING

TRI Researcher:        Constantin Hubmann
Email Address:         Constantin.Hubmann@tri.global
TRI Thrust:            Driving

Most state-of-the-art Motion Planning architectures for autonomous driving are split up over different layers:

- A kinodynamic trajectory planner
- A Behavior planner which sets the constraints of trajectory planning problem(s)
- A state estimation and prediction module which infers latent states of the environment and predicts the trajectories of the other drivers in a probabilistic fashion.

This separation into different modules can be robustly implemented but does not allow for global optimal behavior, correct modeling of interaction (the influence of ego actions on future trajectories of others) and makes correct, probabilistic handling of uncertainties difficult. This is the case as splitting up the motion planning problem under uncertainty into different, independent smaller problems is a simplification of the problem. It may ultimately lead to overly conservative behavior compared to global problem formulations, that may provide solutions in the policy space. Additionally, such an approach requires to tailor behavior generation to specific cases, which may become difficult for the magnitude of possible scenarios.

Non-deterministic, global problem formulations such as MDPs or even its extension for probabilistic states, POMDPs, provide a policy instead of a trajectory which allows for less conservative, globally optimal behavior. The global, optimal formulation allows for a more generic framework, which adapts easier to new scenarios due to its generic nature and formulation.

Possible topics in this area are:

- Online solving of (PO)MDPs by use of Monte Carlo Tree Search
  - Parallelization of sampling
  - Learning of the optimal Value function and action selection (especially interesting with the vast amount of upcoming Guardian data) for speeding up MCTS
  - Fundamental research: transfer from Monte Carlo Tree Search to Monte Carlo Graph Search, balancing of (Q-)Value estimation, Action selection
- Deep Imitation learning for policy generation
  - Very interesting because of the massive amount of upcoming Guardian data
  - Fundamental research in input state representation, safety guarantees, etc
- Deep Reinforcement Learning for policy generation
  - May be very promising in the long run
  - TRI is one of the few companies that has the possibility to do this: Strong ML Team and strong financial backing to train on huge simulation data
  - Results can also be used as heuristics for MCTS

## 8.6   SEMI-SUPERVISED LEARNING FOR UNDERSTANDING OBJECTS IN DRIVING DATA

TRI Researcher:        Dennis Park, Sudeep Pillai
Email Address:         Dennis.Park@tri.global, Sudeep.Pillai@tri.global
TRI Thrust:            Driving

**What do you want to do?**
The goal of this project is to train robust 3D perception models with the least amount of labels annotated by humans. The ideal outcome is a set of machine learning training strategies that can operate with partially labeled sequences of sensory data, for example, when only 1 out of 10,000 frames are labeled.

We will focus on the task of detecting 3D objects in driving scenes, using multi-modal sensory input. Our approach will be lower-bounded by the performance of current state-of-the-art fully-supervised learning algorithms, which leverage only the small portion of labeled data. The key problem is how to surpass these algorithms by leveraging the large amount of unlabeled data that will become available to Toyota starting from 2020.

We plan to exploit a set of inductive biases that has already proven to be effective in the machine learning community. This includes, but not limited to, object permanence (i.e. the same object is found in the neighboring frames) and consistency in confidence (i.e. the confidence of detector is invariant to a set of transformations that preserves the object identity).

In addition, we also want to explore how to exploit the structure of the sensor data as supervisory signals for object detection. Tentative questions we plan to explore includes how to impose a geometry-based consistency as constraints of inference (e.g. ground-plane assumption), or using detectors operating on different modalities (e.g. lidar vs. rgb) to derive geometric constraints.

**Why do we care?**
Robust perception system is a critical component of autonomous driving system. Although the research community has seen large improvements using data-driven machine learning models, training such models requires a large set of labeled data.

A unique advantage of Toyota as #1 automobile manufacturer is that we will soon have access to "Toyota-scale" sensory data that is order-of-magnitude larger than what our competitors have. This, however, imposes a unique challenge as well: manually labeling such data will quickly become infeasible. Therefore, developing a successful semi-supervised learning for perception system will bring the biggest advantage to Toyota.

**Technical problem statement**
Given a dataset $D = \{X_i, Y_i\} \cup \{X_j\}$, where are multi-modal sensory $D = \{X_i, Y_i\} \cup \{X_j\}$ $X_k$ data (RGB + Lidar) and $Y_k$ are corresponding 3D bounding box labels, train a discriminative model $f(X) = Y$ such that it minimizes a semi-supervised loss on $D$. The loss is derived from principles of object permanence, invariance of confidence over a set of transformations, and 3D geometric constraints.

## 8.7 IMITATION LEARNING FOR HIGH-UNCERTAINTY DRIVING SCENARIOS

TRI Researcher:       Dennis Park, Adrien Gaidon
Email Address:        Dennis.Park@tri.global, Adrian.Gaidon@tri.global
TRI Thrust:           Driving

**What do you want to do?**

The goal of this project is to train an imitation-learning agent that can learn to manipulate unmanned vehicles in challenging driving scenarios with high uncertainty (e.g. making unprotected left-turn, complex nudge-around scenarios). The agent learns the behavior by mimicking large scale expert demonstrations, without an explicitly designed cost function. The ideal outcome is a set of machine learning training strategies that can learn prototype controllers in challenging driving scenarios, at least in simulation environment.

While imitation learning (IL) has been widely explored in the RL community, the unique challenge of TRI is to extend the state-of-the-art IL algorithms from toy-like environments to complex driving environments. This requires designing algorithms that are highly intertwined with the underlying perception system. To reduce the dimensionality of the state space and remain consistent with existing state-of-the-art perception algorithms, we constrain the perceptual space to be in Bird's Eye View (BEV), where traffic agents are represented on an orthographic planar surface viewed from the top. This also allows for adapting the action space to simplified orientation control, plus a few more low-dimensional variables (e.g. brake or accelerator).

The challenges lie in multiple layers. First, the agent is in a "PU learning" setting (Positive Unlabeled): it needs to learn mostly from positive instances only (expert demonstrations), while being robust to noisy demonstrations and implicitly learning to avoid negative scenarios. Second, the BEV state space represents the (non-linear) propagation of uncertainty through a typical modular driving architecture (coming from mapping, localization, detection, tracking, prediction). Hence, the agent needs to learn to reason under that complex uncertainty using only demonstrations to predict a distribution of optimal actions. Third, real-world driving scenarios are more complex than typically explored in RL research.

We propose to first define a set of high-uncertainty driving scenarios where analytical controllers face challenges and collect expert's behavior from our driving records. Then we will explore principles of state-of-the-art IL algorithms while adapting them to high-dimensional probabilistic BEV state spaces and realistic actuation space in driving to enable end-to-end training.

**Why do we care?**

The decision process in driving is complicated enough, so that manual design becomes quickly infeasible. Although deep reinforcement learning enjoyed great success recently, directly applying them to driving scenarios is limited, since defining cost function over complicated state space is non-trivial.

Toyota has accumulated high-quality driving logs that record expert's driving examples in complex urban scenarios. We aim to use this raw behavioral experience to learn optimal control policy. A unique advantage of Toyota is that we will soon have access to large scale driving record, which will potentially bring unbounded improvement on our autonomous

controller.

A key part of the problem is defining the state and action space that are compatible with existing perception system, at the same time amenable to probabilistic reasoning involved in IL. The form of demonstration is flexible, from the full trajectory of atomic action and sensor input to a pair of first / last BEV representations.

## 8.8   SCALING UP PREDICTION AND PLANNING

TRI Researcher:        Guy Rosman, Stephen McGill, Jonathan DeCastro
Email Address:        Guy.Rosman@tri.global,  Stephen.McGill@tri.global,
                      Jonathan.Decastro@tri.global
TRI Thrust:          Driving

### 1 Background
Given the requirements of a Guardian system, and the need to understand risky behaviors on the road for both the driver and ado-vehicles, we are looking for prediction that will scale to multiple agents, to longer horizons, and rare events. Impact to TRI: We must predict road agent behavior robustly; modeling failure likelihoods, handling diverse scenarios, reasoning about multiple agents, and identifying rare events. Addressing these aspects across the large scale of Toyota fleet data will provide major safety gains in new applications.

### 2 Aims
In this large-scale effort, we want to address several questions that arise in planning and modeling of road agents:
- Scaling up prediction - Mining the prediction data and learning in heavily biased datasets, which metrics, which costs hierarchical, multi-scenario reuse of learned rules.
- Prediction and planning - confidence and impact on planning, how to quantify the impact of prediction on a plan execution?
- Explainable prediction - Can we mine explainable behaviors, including their use for verification and online use with confidence, from observed data at a large scale? In TRI we are interested to see how, as we get more data, we can get a finer granularity of understanding and prediction of road user behaviors.
- Prediction conditional on the ego car - what is a good representation to learn from data the conditional distribution of the ado-vehicles given ego-car actions that can be either human, or planner-based?

### 3 Technical Problem Definition
With these different aspects of prediction and planning, we are looking to extend current results in several directions:
- Approaches for prediction with lower error at longer prediction horizons – beyond 5 seconds, towards events full intersection negotiation length, > 10 seconds. Proposals should include milestones of accuracy and horizon, possibly with new metrics for prediction given its multicriteria nature (horizon, error, coverage of cases, etc).
- Efficient ways to integrate prediction of multiple agents into shared autonomy plans and autonomous planners. Demonstrations include reasoning about plans in complex interactions such as multi-agent intersection navigation on a long horizon. This includes

confidence measures for predictions as part of a way to mitigate the effects of errors in prediction.

- Extracting road agents prediction that cover what humans interpret as maneuvers and multi-agent interactions on the road. Detecting the patterns in the data in terms of multi-agent interactions, and reasoning about them efficiently.
- Efficiently represent other road agents' future prediction given the ego-vehicle's behavior, in a way that can be learned from arbitrary datasets (e.g. not a full autonomous driving stack), but is useful for planning.
- Approaches of covering rare events in the predictions and gauging their risk. Proposals should provide milestones for measuring the effectiveness of risk estimation and level of uncertainty in prediction.

**4 Data, Platforms, Deliverables**
In order to handle both the control, reactive part of the problem, and the large-scale data aspect, we will test multiple data sources and platforms. Data sources include driving data from both external and TRI sources, DataFromSky annotated data or similar large-scale data. Deliverables with specific APIs and benchmarks can include prediction benchmarks, or planning. Overall, milestones should include both existing and future datasources, and demonstration of robustness, sample efficiency, and other properties in a reproducible manner. For the reactive part, CARLA (or other simulators) and RC cars (or similar robots) can be used for risky or multi-agent events.

## 8.9   PLANNING AND LANGUAGE

| | |
|---|---|
| TRI Researcher: | Guy Rosman |
| Email Address: | Guy.Rosman@tri.global |
| TRI Thrust: | Driving |

**1 Background**
In both the robotics and autonomous driving domains, we are looking for planning and prediction algorithms that integrate with human understanding of tasks and plans. This is relevant both when we want to understand when people do on the road, at the maneuver and interactions level, and when we want robotics at the home to scale up in their set of objects and tasks. We are looking for approach to learn the association of language and tasks, to understand when our vocabulary is limited or inappropriate to describe the tasks.

**2 Impact to TRI**
For automated driving, we want to use the language as a tool to reason about more complex interactions, as planning / RL is assumed to allow better understanding of long-term, complex, interactions. For home robotics, having approaches that naturally can scale robots' vocabulary of tasks and objects can be crucial as we scale to new deployment in arbitrary settings.

**3 Aims**
The project would explore approaches that connect language and planning., We are looking for approaches to connect language, planning, and prediction, in both the driving and robotics domain including but not limited to methods that –

- Explain plans with language, express planner rollouts via language or query them.
- Plan according to common language directions.

- Explain failure cases and critical events, near-accidents and accidents.
- Demonstrate how planning with language scales automatically to more diverse actions and objects, possibly with mixed-supervision learning.

## 4 Technical Problem Definition

Measurable capabilities of these systems should explore these directions, extending the state of the art in several directions –

- Demonstrating planning and prediction based on language commands and descriptions, showing these approaches scale to large datasets with many tasks and state, and that learning can be done in a sample-efficient way.
- Demonstrate how planners and predictions based on a language can obtain good coverage of multi-agent interactions.
- Demonstrate how planners can describe failure modes for plans and/or execution; Demonstrate predictors that can identify possible risky events as trained with a handful of such cases in the dataset.
- Demonstrate extension of plans beyond current limitation of language-based planning, for example to include full paragraphs, reasoning about different modifiers to the plans, or extension to more complete, in-the-wild scenarios.

## 5 Data, Platforms, Deliverables

In order to handle both the control, reactive part of the problem, and the large-scale data aspect, we encourage to test on multiple data sources and platforms. Data sources include driving data from both external and TRI sources, DataFromSky annotated data or similar large-scale data. For critical events, proposals can include mining of different sources for accidents data. For robotics demonstration, both simulation, passive datasets, and demonstrations on real robots are encouraged. For the reactive part in driving, CARLA (or other simulators) and RC cars (or similar robots) can be used for risky or multi-agent events. As part of the projects, a deliverable with a clear API is encouraged to showcase some of these approaches' capabilities.

## 8.10 EVENT-BASED VISUAL PERCEPTION FOR NAVIGATION

TRI Researcher:      Hanme Kim
Email Address:      Hanme.Kim@tri.global
TRI Thrust:      Driving

Rapid status update of nearby dynamic agents in our autonomous driving scenarios is highly critical to plan and control our autonomous vehicles (AVs) fast and reliable enough not to cause/get involved with any accident, especially in our Guardian mode. The current update rate is however limited by the measurement rate of our sensors (e.g. LiDAR and camera, typically 10-30Hz). We therefore rely on motion models (e.g. constant velocity model) to predict how tracked objects could move from one sensor measurement to the next during so-called the blind time interval (e.g. 33~100ms). Such motion prediction cannot be always correct, and it gets worse at a higher speed (e.g. while we predict an agent moves straight, it could actually make a sudden turn).

To overcome this challenge, we propose to utilize an event camera, a paradigm shift in visual sensing. Unlike a standard camera, which generates video by regularly and synchronously opening its shutter to expose all pixels and capture frames, the event camera has no shutter. Its

pixels are independent elements which continuously monitor the intensity of light reaching them, and send out asynchronous reports ("events") only when the intensity changes by a threshold amount. By encoding only brightness changes, it offers the potential to transmit the information in a standard video but at vastly reduced bitrate, and with huge added advantages of very high dynamic range and temporal resolution. Moreover, since Samsung has started mass production of their Dynamic Vision Sensor, the price of the cameras has gone down 10-fold.

The main idea is therefore that we track and match dynamic objects in autonomous driving scenarios using the almost continuous stream of events from the event camera (possibly fuse with other sensor measurements). We expect that our perception to be more robust and efficient by reducing the motion uncertainty between measurements (e.g. increase the accuracy of the motion prediction, require a smaller search region, etc) and updated at a much higher rate (e.g. 1000Hz), and eventually enable us to plan and control our AVs in a way that avoids any potential accident. But it has proven very challenging to use this novel sensor in most computer vision problems, because it is not possible to apply well established computer vision techniques, which require synchronous intensity information, to its fundamentally different visual measurements.

## 8.11  LOW COST STATE AND FRICTION ESTIMATION FOR VARIOUS DYNAMIC REGIMES

TRI Researcher: Jeff Walls, Carrie Bobier-Tiu, Manuel Ahumada
Email Address: Jeff.Walls@tri.global, Carrie.Bobier-Tiu@tri.global, Manuel.Ahumada@tri.global
TRI Thrust: Driving

**What do you want to do?**
We are looking to develop new capabilities for low cost friction and vehicle state estimation. We are interested in exploring methods for estimation in various driving regimes, from low speed (or stopped) and low dynamic excitation to limit handling maneuvers.

**Why do we care?**
TRI's approach to autonomous driving lies in two systems: Chauffeur, a level 4 or 5 autonomous vehicle; and Guardian, a vehicle with a driver in the loop, but full or partial usage of the autonomous sensor suite and algorithm stack to provide enhanced safety features. One of the underlying concepts behind both Guardian and Chauffeur is that these systems must be at least as good as the best human drivers. This means being able to handle the vehicle in a large set of maneuvers including up to the vehicle's handling limits (e.g., aggressive lane-changes, driving on low friction surfaces, etc). Accurate methods for real-time state and friction estimation is critically important for this task.

In order to push the capabilities of advanced ADAS systems, and to allow the Guardian and Chauffeur systems to function in a variety of operating conditions, it is crucial to develop low cost and effective means of estimating the vehicle's capabilities under its present, and potentially future, operating regime.

**Technical problem statement**

If we consider a simple vehicle dynamics model used in ADAS systems, the bike model, we can capture a large portion of the vehicle's nominal behavior using 3 states: longitudinal velocity, lateral velocity (or sideslip), and yaw rate. With the addition of a nonlinear tire model (Fiala or Pacejka, for example), the limit behavior of the vehicle is also well-captured through the definition of unstable equilibrium points.

In order to use these models effectively in a system like Guardian, some quantities that are not directly measurable with low-cost sensors need to be estimated for peak performance. For example, estimation of lateral velocity or sideslip is key to understanding the lateral stability limits of the vehicle. Similarly, estimation of tire-road friction defines the limitations of the tire capabilities through the friction circle, but the quantity is historically difficult to estimate accurately until the limits are approached through high levels of state excitation. On the opposite spectrum, it is often difficult to estimate even longitudinal velocity at very low speeds due to wheel speed encoder limitations.

We are looking for novel and cost-sensitive approaches to solve these type of estimation problems in regimes that are typically difficult to achieve good results. Expanding our estimation performance over a high range of operating conditions will allow Toyota's Guardian and Chauffeur systems to provide higher quality safety.

## 8.12  3D REGISTRATION OF HETEROGENEOUS DATA

TRI Researcher:        Jeff Walls
Email Address:         [Jeff.Walls@tri.global](mailto:Jeff.Walls@tri.global)
TRI Thrust:            Driving

**What do you want to do?**
Consider the age-old problem of registering two sets of corresponding data–computing the rigid body transformation between sensor coordinate frames associated with each data set. The data sets may take many different forms including low-level points (e.g., 2D pixel location or 3D lidar point clouds) or even higher-level geometric primitives (e.g., lines produced as the output of a feature detection process). We would like to consider generalizable methods to register corresponding heterogeneous data sets, i.e., each data set is produced from a different sensing modality (e.g., register 2D image data with 3D point cloud).

**Why do we care?**
Localization (estimating pose with respect to some map) is a core competency for TRI's self driving vehicle effort. The localization problem boils down to computing a rigid body transformation between online perception data and a map. Maps may be represented in a variety of ways: point cloud, sparse features (e.g., lane lines), etc. Moreover, vehicles may be instrumented in different ways, e.g., camera, LIDAR, etc. To enable all vehicles to reliably localize within a map, we require the ability to register sets of heterogeneous data.

**Technical problem statement**
Core technical challenges include:
- Developing robust registration techniques that can generalize to heterogeneous data representations.
- Represent uncertainty of computed registration.
- Meeting real-time requirements.

## 8.13  INTRODUCING ASSISTIVE TECHNOLOGIES TO DRIVERS: ESTABLISHING TRUST WITH EMOTIONAL UNDERSTANDING

TRI Researcher:      John Gideon
Email Address:       John.Gideon@tri.global
TRI Thrust:          Driving

**Background & motivation**

As more assistive technologies are incorporated into our vehicles, it becomes increasingly necessary to establish trust between the car and driver. One way to accomplish this is by creating a conversation between the vehicle and the driver, using dialogue systems (e.g. Amazon's Alexa). To build trust, this conversation should continuously make vehicle intentions clear, as well as ensure that driver preferences and comfort are accommodated[1].

This is especially important for drivers using assistive technologies for the first time, such as lane keeping and automatic cruise control. A bad initial impression can lead to the driver permanently deactivating a feature that could otherwise greatly improve safety. Furthermore, without fully understanding a feature's limits, this can result in unnecessary risk-taking by drivers seeking to find these limits themselves. However, different training strategies are needed to allow individuals to learn at their own pace. Furthermore, as people become more familiar with these features, the vehicle should adjust to their preferences.

In order to best facilitate this process and foster trust, an understanding of driver emotion is key[2]. For example, frustration could indicate that the training process should be altered or the feature parameters (e.g. braking distance) should be changed. Alternatively, happiness could indicate satisfaction with the feature, suggesting that similar techniques could be used for other interactions. By improving these initial impressions, we can increase the perceived value of assistive technologies and increase adoption[3].

**Goals**

The goals of this project are:
1. To collect a dataset of drivers being introduced to new assistive features using a variety of (sometimes frustrating) training strategies. Recordings of teenagers learning to drive, as well as older drivers would be of particular interest.
2. To improve a dialogue system to better recognize and respond to driver emotions when using assistive technologies. We aim to demonstrate that such a system increases interest in activating safety features versus those introduced with minimal prior training.
3. To facilitate future development of driver understanding and assistive technologies. Emotional events could trigger naturalistic data collection and improve the diversity of our models.

**Technical problem statement**

This project will need to address the following:
- Which modalities (video, audio, speech content) are most useful to estimate emotion and the effectiveness of the current training method?
- How can we isolate driver audio from other passengers and background noise? Can this separation and automatic speech recognition be improved with lip reading?

- How should the system incorporate information about the context of emotional events, including unusual braking patterns and distractions in the environment?
- How does the expression of emotion under simulated driving conditions differ from the real-world? What is the best method of combining these different sources of data?
- What is the importance of diversity in the initial (relatively small) training set? What methods of subject personalization over time can be used to improve model performance?
- How can we mitigate alarm fatigue and ensure that the dialogue system itself is not deactivated?

[1]Bobbie D. Seppelt, Int J Hum Comput Stud., 2019.
[2]Jeamin Koo et al., IJIDeM, 2015.
[3]Min Kyung Lee, Big Data & Society, 2018.

## 8.14  SPECIFICATION LANGUAGES FOR DATA-GROUNDED SCENARIO SYNTHESIS

TRI Researcher:      Jonathan DeCastro
Email Address:       Jonathan.DeCastro@tri.global
TRI Thrust:          Driving

### What are we trying to do, and why?

Improving the testing process for autonomous vehicles (AV) requires new techniques for automatically constructing scenarios aligned with reality.  Approaches that involve manual scenario generation from a basic set of requirements or directly from log data are limited by the subjective biases, and moreover do not scale.  Furthermore, it is difficult to achieve the right abstraction of a scenario containing what is essential for testing from a corpus of data without over-fitting to data or generalizing away important features.

A virtuous cycle of simulation-enabled development will entail a data-driven approach, where the user has only to worry about the testing requirements where it is easy to introduce variations on scenarios with minimal cognitive overhead, while being able to take cues from a large corpus of data.  Such approaches will allow users to construct scenarios that have introspection on the realism or other metrics that data can provide, will allow creation of data-grounded testing strategies, and will provide an implicit certificate or benchmark that remains invariant to the state of the AV system development.

### Technical problem statement

TRI is interested in ideas and solutions that enable modelling and creation of realistic scenarios in a way that is intelligently informed by large corpuses of naturalistic driving data.  Specifically, we seek:

- Approaches that strengthen the connection between 'data' and 'data-driven' agent models and scenario instantiations with an aim to improve characterization, coverage and robustness of such models.  This will allow for closing the simulation-to-reality gap using reasonable scenario abstractions over data and allow introspection of causal factors for simulation runs that violate some specification.

- Use of formal languages or domain-specific languages with the expressive power to synthesize scenarios, requirements, and evaluators declaratively (i.e. able to specify the "what" without the "how").

Such data-driven scenario descriptions will offer: 1) a concrete set of tests and probabilistic analyses to apply to an AV system and 2) implicit certificates for the AV system that are backed up by real-world data, with human-interpretable abstractions.

Proposals are sought that address the temporal aspects of parameterizable, data-driven agent models. The awardee will interact closely with TRI researchers to develop the software tools, use appropriate data sources, and integrate within TRI's testing infrastructure.

## 8.15  MAP COMPRESSION: PRUNING FACTOR GRAPHS WHILE MAINTAINING LOCALIZATION PERFORMANCE AND ANNOTATIONS

TRI Researcher:      Karl Rosaen
Email Address:       Karl.Rosaen@tri.global
TRI Thrust:          Driving

The SLAM team maintains maps in the form of factor graphs of vehicle poses and landmarks, built from sensor data collected by driving logs. Landmarks facilitate localization. Both the vehicle poses and landmarks serve as the basis for human annotation of the rules of the road which are published to the planner while driving to make safe and feasible driving decisions.

We are working on scaling up our map system to cover the world and be able to benefit from the sensor data of a growing fleet, allowing us to cover a wider area and update it frequently. One key challenge introduced by the increased volume of data is continuously improving the map while avoiding having its size grow with every new driving episode. We wish to preserve localization performance for any trajectory within the map and maintain consistent annotations.

While some heuristics come to mind such as automatically removing graph nodes tied to sensor data older than a fixed window and transferring annotations to newer graph nodes, we'd like to explore more rigorous approaches to graph reduction through removal or other means without sacrificing localization performance. What if sensor data from a year ago is still the best coverage of lane lines in a particular region? What if newer data has nothing to add? We would like to explore these questions and based on the results deploy the insights into a large-scale mapping system.

## 8.16  SPECIFICATION-DRIVEN INTELLIGENT TESTING

TRI Researcher:      Nikos Arechiga
Email Address:       Nikos.Arechiga@tri.global
TRI Thrust:          Driving

Design processes for safety-critical systems have traditionally relied on manually generated tests. Human engineers with insight into the functionality of the system craft test cases that will exercise behaviors that are believed to maximally stress the design. However, this type of insight-driven testing is unlikely to scale to advanced autonomous systems. The sheer scale and complexity of these systems means that it is difficult to gain good coverage of the system design. Furthermore, test cases that seem difficult to a human may be trivial to an autonomous

system and vice versa. The goal of this initiative is to develop techniques to automatically generate test cases that stress the autonomy stack and provide good confidence of coverage.

## 8.17 COMPOSITIONAL AND DECOMPOSITIONAL REASONING

TRI Researcher:        Nikos Arechiga
Email Address:         [Nikos.Arechiga@tri.global](mailto:Nikos.Arechiga@tri.global)
TRI Thrust:            Driving

**What do you want to do?**

This initiative seeks to develop technologies capable of (1) automatically reasoning about system components and specifications to find contradictions, (2) checking conformance with higher-level specifications, and (3) suggesting modifications to interface specifications that would repair inconsistencies.

**Why do we care?**

Development of advanced autonomous systems requires separation of responsibilities across multiple engineering teams. At integration time, there is a risk that teams make inconsistent assumptions, or that the integrated system does not satisfy high-level safety and correctness specifications.

**Technical Problem Statement**

Specifications for system components may be given in a variety of different forms, including formal logic, UML diagrams, and domain-specific languages. The goal of this project is to develop a reasoning engine that can parse a set of such specifications and reason compositionally and de-compositionally in the following sense.

1.   Are the interface components between specifications consistent, or do they incur contradictions?

2.   Do the interface specifications ensure that high-level system requirements are met?

3.   If high-level requirements are not met, can the system suggest modifications of the interface requirements?

## 8.18 ASSURANCE FOR PEREPTION COMPONENTS

TRI Researcher:        Nikos Arechiga, Sudeep Pillai, Wadim Kehl
Email Address:         [Nikos.Arechiga@tri.global](mailto:Nikos.Arechiga@tri.global), [Sudeep.Pillai@tri.global](mailto:Sudeep.Pillai@tri.global),
                       [Wadim.Kehl@tri.global](mailto:Wadim.Kehl@tri.global)
TRI Thrust:            Driving

Perception components are a key part of autonomous functionality. However, ML techniques are inherently unpredictable because they rely on implicit, statistical specifications, and lack the mathematical framework to provide guarantees on safety and correctness. Commercial autonomous systems, however, require predictable behavior and strong assurance. The goal of this initiative is to develop an assurance strategy based for perception systems which is based on partial specifications for perception systems. These partial specifications include things like model assertions, logical scaffolds, uncertainty modeling, and model calibration. These partial specifications will enable (1) reasoning about the responsibilities that the perception system

must satisfy for assurance purposes, (2) guide the testing strategy of these systems and help prioritize directions for improvement, and (3) serve as runtime checks to detect abnormal conditions and violations of critical assumptions.

## 8.19 USE THE FORCE: MULTI-AXIAL HAPTIC SYNTHESIS FOR DRIVER-VEHICLE STEERING AUTOMATION

TRI Researcher:      Selina Pan
Email Address:      Selina.Pan@tri.global
TRI Thrust:          Driving

**What do you want to do?**
Develop a grip force sensing and shape-changing steering wheel to serve as an intuitive means of communication between a human driver and an automated driving system.

**Why do we care?**
Clear and intuitive communication must be established to support driver/automation interactions, especially during control transitions when authority is shifted from one agent to the other. Each agent should be aware of the other's control actions and the current delegation of control authority (which agent has more control at a given time). Visual, audio, and limited haptic modalities have been explored in current ADAS systems; however, given the physical contact between a human driver's hands and the steering wheel, the third modality has huge potential for serving as the most intuitive means of communication in a system like Guardian.

**Technical problem statement**
This project would explore the use of a finite and possibly varying automation impedance to facilitate control sharing and to facilitate transitions of control authority, as appropriate to the driving situation. In part, this human/automation control sharing paradigm is inspired by human/human haptic communication. Two humans manipulating an object cooperatively can read the other's control actions by comparing haptically monitored force or motion responses to expected responses even while pushing and pulling on an object. They can also read the other's control authority by monitoring mechanical impedance (the relationship between force and motion), even through the object. However, two cooperating humans will typically supplement their pushing and pulling with other communication channels, often leading to significant improvements in performance on the shared task. In this spirit, this project posits that a grip-force sensing and shape-changing steering wheel can set up an additional communication channel between the human driver and automation system that will significantly improve shared driving performance.

This project will use haptic communication in the axis of grip to communicate control authority and thereby to make the transitions of control authority between the human driver and the automation system smoother and more intuitive. Thus this modality would communicate information that is redundant with the impedance in the steering axis. The idea is rooted in the observation that human drivers increase their grip on the steering wheel either when they want to take over control of the vehicle or when they are surprised by something that they encounter on the road. To ensure safe and comfortable driving, a well-designed automation system would 1) sense the grip force applied by the driver, 2) understand the traffic situation and decide whether to override or acquiesce to the human control action, and 3) inform the driver of its

decision. The communication aspects can be addressed by a shape-changing steering wheel with pressure or grip force sensing. A driver squeezing the steering wheel to request greater control authority who feels the steering wheel expand in response will immediately know that their takeover request was denied. On the other hand, if the steering wheel deflates under an increase in grip, the driver will immediately know that the takeover request was granted and that they are now in control. Likewise, transitions can be initiated and communicated by the automation system with shape changes. Further, at the same time shape change is produced by driving fluid mass in or out of bladders on the steering wheel, pressure changes can be used to measure the grip force applied by the driver. Advantageously, this two-way communication channel in the axis of grip is orthogonal to the two-way communication taking place in the axis of control (steering). The existence of an independent communication channel in an orthogonal axis supports simultaneous negotiation of control authority and execution of control action.

## 8.20  NEURO-ADAPTIVE OBSERVERS FOR SIMULTANEOUS PARAMETER AND STATE ESTIMATION IN TIRE AND VEHICLE MODELS

TRI Researcher:       Selina Pan
Email Address:        Selina.Pan@tri.global
TRI Thrust:           Driving

**What do you want to try to do?**
The overall objectives of this proposal are to develop a vehicle and tire model consisting of a combination of physically meaningful differential equations and adaptive machine-learning-based neural networks, to develop associated hierarchical algorithms for estimation of both slip/force variables as well as tire model parameters, and to validate the complete model and estimation system using data from CARSIM and from limited real vehicle experimental measurements.

**Why do we care?**
Tracking of snow and ice profiles on roads will be required in order to roll out Guardian or Chauffeur systems. This project will also have applications in autonomous driving on winter roads, since lane markers can be completely invisible on snow covered roads.

**Technical problem statement:**
This project will use a modeling approach consisting of a combination of physically meaningful differential equations and adaptive deep-learning-based neural networks to represent vehicle dynamics and tire models. In particular, well-understood phenomena such as force balances, mechanical motion per Newton's laws, aerodynamic drag, rolling resistance, and combined acceleration terms for lateral and roll accelerations will be modeled using analytical differential equations. Tire models for both lateral and longitudinal forces, and the friction circle will be modeled using neural networks whose weights can be initially obtained using training via backpropagation. In addition to initial training, model parameters for the neural networks and a subset of parameters for the physically meaningful differential equations will also be updated automatically online during regular vehicle use, based on a hierarchical estimation system and based on the type of vehicle maneuver being executed by the vehicle. The project will include development of a rigorous neuro-adaptive observer that enables estimation of states and model parameters. The algorithm will be reconfigurable, with the available measurements determining the set of parameters that can be updated online. A hierarchical architecture will enable slip,

force and friction coefficients to be updated quickly in real-time. Tire model parameters will be updated on a slower time scale using significantly larger sets of data which allows assumptions on average friction coefficient values and use of repetitive windows of similar data sets. The overall activities in the project will include development of the architecture for the combined modeling approach, development of the rigorous neuro-adaptive estimation algorithms for both parameter and state estimation, and multi-stage validation of the complete model and estimation system using data from CARSIM and from limited real vehicle experimental data available in the PI's lab.

A number of tire-road friction coefficient estimation algorithms have been developed by several researchers in literature, in addition to algorithms for estimation of slip angle, slip ratio and tire forces. However, these algorithms are largely based on analytical tire models and their known parameter values. The work proposed in this project is potentially ground-breaking in that it utilizes data-based tire models and further uses hierarchical estimation which updates not only state variables and tire- road friction coefficients, but also tire model parameters over longer time periods. Further, the tire-road friction coefficient estimates with the new approach will work more reliably over a wider range of operating slip conditions. This project will enable further refinement of the developed theory and application to a commercially useful real-world intelligent vehicle application.

## 8.21  PROTECTING BICYCLISTS WITH MODERN TECHNOLOGY

TRI Researcher:       Selina Pan
Email Address:        Selina.Pan@tri.global
TRI Thrust:           Driving

**What do you want to try to do?**
Leverage sensor systems developed for smart bicycles into combined bicycle-vehicle testing for mixed-traffic scenarios, with wider applications to heterogeneous road users.

**Why do we care?**
Over 48,000 bicyclist injuries and approximately 700 bicyclist fatalities due to crashes with cars are reported annually. With emerging autonomous vehicle technology, we have the opportunity to actively focus on bicyclist-vehicle interactions and technology. These interactions also have wider use cases extended to motorcycles, scooters, hoverboards, and future mixed-use road spaces in evolving city centers.

**Technical problem statement:**
Protecting a bicycle from potential car-bicycle crashes requires addressing a number of difficult challenges. Only inexpensive sensors can be utilized on the bicycle. Sensors on autonomous cars typically cost many thousands, or even tens of thousands of dollars, and utilize more electrical and computational power. The system proposed for use in this project, on the other hand, is aimed at a market price below $500 and is suitable for a cost-sensitive bicyclist. The technology needs to operate in complex traffic scenarios as well. The sensing and estimation algorithms on the bicycle need to track trajectories of vehicles on local urban roads where traffic scenarios are more complex than on highways. Riding on local roads involves traffic intersections with many left-turning, right-turning, and cross-traffic vehicles, all of which need to be tracked. Finally, the collision prevention system needs to be useful on today's roads. It cannot rely on all cars being

autonomous, all cars having wireless connectivity, or other futuristic assumptions which may take many years/decades to bear fruit.

There exist prototypes that have been developed for smart bicycles which address all of the above challenges and functions effectively on today's roads. The instrumentation on these bicycles include a single-beam laser sensor, a custom triad sonar unit and a low-density LiDAR sensor. Unique solutions include active rotational control of the single-beam laser sensor to continuously track moving cars behind the bicycle, algorithms to reliably identify vehicles from other objects, and novel estimation algorithms to accurately track lateral and longitudinal positions of cars. On-board speakers are used to create a custom audio "horn" to alert car drivers to the presence of the bicycle. The audio alert presentation is designed to minimize the trade-offs between low reaction time and unnecessarily intrusive disturbances to the bicycle rider and to nearby motorists.

There are many different directions that TRI can go with this project, depending on its needs and the university PI. Primary objectives can be, but are not limited to:
a) Field testing using smart bicycle technologies that rely on user study data gleaned from bicyclist volunteers with significant daily urban commutes. Extensive analysis of bicycle data recorded during these commutes can be used in snippet testing for our driving stack.
b) Focused analysis of traffic intersections (or other such road cases) for specific cities targeted for Toyota technology rollout, or extrapolated to cost-sensitive applications in developing countries, can serve as test cases for our driving stack.
c) Results from smart bicycle prototypes can be used for other agents that share the road, including motorcyclists, scooters, e-bikes, etc. Cities and campuses are starting to see increased popularity of these multiple forms of travel. Automated vehicles or vehicles with highly advanced driver assistance systems must be able to react appropriately to these agents sharing space with them.
d) Bicycle-vehicle communication can be explored using these sensor technologies.
e) Some of the smart bicycle technology uses laser sensors on rotationally controlled stepper motor platforms. This can be used for finding lateral position in the lane on snow covered roads. Current camera technology cannot find lateral position when the road is snow-covered. This could be in conjunction with another friction estimation project. This proposal is more open-ended due to the forward-thinking nature of the technology involved; however, I believe this is an important area and a platform of which TRI could take great advantage in many different directions.

## 8.22  DATA-DRIVEN MODELING OF DRIVER ATTENTION AND SITUATIONAL AWARENESS

TRI Researcher:     Simon Stent, Guy Rosman, Luke Fletcher
Email Address:      Simon.Stent@tri.global , Guy.Rosman@tri.global, Luke.Fletcher@tri.global
TRI Thrust:         Driving

**Background & motivation**

In a survey of 5,471 light vehicle crashes between 2005-7, the National Highway Traffic Safety Administration (National Motor Vehicle Crash Causation Survey, 2008) found that inadequate

driver situational awareness (driver recognition error, caused primarily by insufficient surveillance, internal and external distraction and inattention) was a critical reason in over a third of cases. The statistics are dated, but considering the growth in smart-phone usage since 2007, they are unlikely to be any less significant today.

For Advanced Driver Assistance Systems (ADAS) to help alleviate this class of problem, it will be important to be able to build and maintain models of driver situational awareness. Such models will allow vehicles to respond in a more timely manner in certain scenarios (e.g. by being more assertive if drivers are observed to be unaware of perceived risks in the environment). Beyond ADAS, in Level 2 and 3 automated driving, modeling a driver's awareness will be equally important to safely negotiate exchanges in control.

This project is concerned with fusing and advancing recent developments in two separate research fields: mechanistic models of driver situational awareness from human vision and psychology (e.g. [1]), and data-driven models of task-conditioned human attention from computer vision (e.g. [2]). We will take advantage of the relatively recent availability of high precision eye trackers and advances in scene representation provided by deep learning to develop and evaluate data-driven models of driver attention and awareness.

### Goals

The anticipated goals of this project are to:

1. create two or more large-scale benchmark datasets of gaze behavior in driving scenes (in-sim and on-road, potentially including closed course testing), to extend existing TRI efforts

2. develop data-driven, mechanistically-grounded predictive models of overt attention and awareness, for both idealized drivers and individual drivers

3. analyze model performance on specific scenarios of driver recognition error gathered from naturalistic driving data and closed course testing.

### Technical problem statement

Technical problems of relevance to this project are:

- can we learn a joint model of driving scene, task and minimal set of what a safe driver should attend to or be aware of? How much variation in inattentive behavior exists across individuals, particularly with respect to high risk scenarios?

- given a driver's noisy registered gaze in a scene, how can we estimate what they are aware of? Can overt attention be used acausally as a proxy for awareness for high-risk objects?

- how can we validate an awareness model fairly? What percentage of accidents caused by driver recognition error might be mitigated with better models of attention and awareness?

### References

[1] B. A. Wolfe, B. Sawyer, and R. Rosenholtz. Toward a mechanistic understanding of situation awareness in driving. In revision, Human Factors. 1

[2] Y. Huang, M. Cai, Z. Li, and Y. Sato. Predicting gaze in egocentric video by learning task-dependent attention transition. In Proceedings of the European Conference on Computer Vision (ECCV), pages 754–769, 2018.

## 8.23 DATASETS AND HIGH-FIDELITY MODELS FOR IN-CABIN BEHAVIOR UNDERSTANDING

TRI Researcher: Simon Stent
Email Address: Simon.Stent@tri.global
TRI Thrust: Driving

**Background & motivation**

Over the coming decade, driver and cabin-facing cameras will become increasingly commonplace in new privately-owned vehicles. The trend will be driven partly by our desire to communicate by video while on the move, fueled by faster cellular connections, lower-cost hardware and increasing vehicle automation. More importantly, it will also be driven by the many applications in safety (e.g. Euro NCAP[1]), security and user experience that can be unlocked by applying modern computer vision to in-cabin video. Such applications include: detecting driver inattention or drowsiness; estimating driver takeover readiness; left-child detection; seatbelt detection and adaptive airbag deployment; identity verification; emotional state estimation to match cabin settings/vehicle behavior to mood; audio-visual speech separation and recognition; and unconstrained gesture recognition.

Unlike outward-facing vehicular sensing, which has seen an explosion in large-scale academic datasets and simulators over the past decade that have helped to foster progress in perception, prediction and planning (e.g. KITTI, Cityscapes, Argoverse, nuScenes, CARLA), as of today, there is much less open research being conducted for training and evaluating driver and passenger behavioral understanding (with a few exceptions[2]). This is partly due to the various difficulties of collecting diverse, large-scale and publishable driver-facing data. Since driver error is estimated to be the critical reason in over 90% of road accidents[3], building such datasets to enable better understanding of drivers and their passengers is of high social importance.

**Goals**

The anticipated goals of this project are to:

1. create a hardware setup for efficiently acquiring large-scale, multi-view video data of humans

2. use (1) to generate diverse, domain-specific datasets of drivers and passengers in vehicles, with implicit annotation provided by synchronized and calibrated RGB and NIR video streams

3. use (2) to study both domain-specific and more generalized models for human behavioral coding and understanding, particularly using more constrained sensor setups (e.g. dense 3D pose estimation or action recognition under heavy occlusion from one or two cameras)

4. demonstrate the generalization of (2-3) by application to TRI-owned large-scale naturalistic driving datasets, enabling behavioral analysis.

To help maximize impact, members of the TRI-AD (Advanced Development) Driver and Passenger Monitoring team will also be involved in the project.

**Technical problem statement**

This project will face several technical challenges:

- how to adapt pre-existing multi-camera setups such as Panoptic Studio[4], to target the idiosyncrasies of the in-cabin environment? (e.g. high occlusions, constrained camera positions, wider-angle cameras making calibration and registration challenging, extreme lighting and shadows)

- how to develop models which efficiently acquire invariance to viewpoint? How to best combine information from multiple viewpoints? How to create efficient multi-task models to feed the wide range of downstream applications?

- how important is diversity in the subject pool for generalization? What inductive biases might allow better generalization with fewer training subjects? How does synthetic data compare?

- for a given task, what is the best trade off between camera placement and performance when it comes to deployment in commercial vehicles?

[1] Euro NCAP 2025 Roadmap

[2] For example, driveandact.com, ICCV 2019

[3] National Motor Vehicle Crash Causation Survey, 2008

[4] http://domedb.perception.cs.cmu.edu/

## 8.24 HIGH-SPEED GAZE CONTROLLERS FOR EFFICIENT HIGH-RESOLUTION SCENE UNDERSTANDING

TRI Researcher: Simon Stent, Velin Dimitrov
Email Address: Simon.Stent@tri.global, Velin.Dimitrov@tri.global
TRI Thrust: Driving

**Background & motivation**

Considering that humans are only capable of "20/20" vision within a ~1.5° cone of our visual fields (in photopic illumination conditions, i.e. daylight), our visual systems are incredibly efficient at resolving the relevant details of our surroundings. With the benefit of peripheral vision, powerful attention mechanisms, working memory and finely tuned ocular motor control, we use our vision to solve all sorts of tasks which require wide-angle, fine spatial understanding, from safely changing lanes while driving fast along a highway to social interactions around a dinner table.

In robotics, elegant practical approaches for using one or two cameras to achieve a similar degree of spatial coverage are rare. Brute force solutions are the norm: use a high-resolution camera and move the robot to move the camera, or don't move the robot and add more cameras. The former is often slow or infeasible: cars should not have to turn left to see left. For this reason, Tesla's "full self-driving" setup consists of eight external cameras, with the highest focal camera (i.e. longest range) having a horizontal field of view of ~35° over 1280 pixels -

slightly lower acuity than 20/20 human vision. The MobilEye EyeQ5 chip, to be launched in 2020, will enable processing of "more than sixteen multi-mega-pixel cameras." While clever engineering makes such solutions viable, fixed multi-camera setups will always be constrained by the trade-off between resolution and coverage: the higher the resolution needed in any particular area, the lower the overall spatial coverage of the system (without adding more cameras). Tesla's current setup is unlikely to be able to resolve the gaze and facial expressions of all three other drivers at a 4-way intersection.

Recently, attempts have been made towards more elegant solutions, using systems of mirrors and optics to steer otherwise static cameras (e.g. [1, 2]). Such systems, if made to work more practically, could vastly reduce the number of cameras required to acquire superhuman visual acuity and angular range.

**Goals**

The expected goals of this project are:

1. design and fabricate a working mirror-based saccading camera system, taking into account the pros and cons of existing approaches

2. develop methodologies to optimize camera control for specific tasks, using a wide-angle camera to inform glance strategies and bind together saccading camera inputs. Motivating tasks include: monitoring fixed points in space under robot motion (e.g. for long-range obstacle avoidance on curved highways), improving trajectory prediction during social driving interactions, or finding a book on a large bookshelf

3. investigate the practical viability of such a system for applications in household robotics and driving.

**Technical problem statement**

Many interesting technical challenges are anticipated for this project. Hardware and design challenges include: how to make a practical controller which meets desirable constraints of form factor, saccade speed and precision, angular range, and supported camera focal length? How can the system rapidly accommodate varying intensity levels or scene depths in different regions of the visual field? Software challenges include: how can saccade strategies be learned through self-supervision for a particular task? How can information from the saccading camera be efficiently bound to information from the wide-angle camera?

**References**

[1] K. Okumura, H. Oku, and M. Ishikawa. High-speed gaze controller for millisecond-order pan/tilt camera. In 2011 IEEE International Conference on Robotics and Automation, pages 6186–6191. IEEE, 2011. 1

[2] K. Iida and H. Oku. Saccade mirror 3: High-speed gaze controller with ultra-wide gaze control range using triple rotational mirrors. In 2016 IEEE International Conference on Robotics and Automation (ICRA), pages 624–629. IEEE, 2016. 1

## 8.25 INFERRING AWARENESS, MOTIVATION, AND INTENT FROM IMAGES AND VIDEO

TRI Researcher:     Simon Stent
Email Address:     Simon.Stent@tri.global

**Background & motivation**

Humans are remarkably adept at inferring the awareness, motivation and intent (Awareness: is X aware of Y ? Intent: what does X want to do next? Motivation: why does X want to do it?) of other agents, even when observing highly abstract stimuli such as the Heider-Simmel Illusion. This cognitive capability is present early in infant development: from a very young age we understand that others around us have mental states like goals and beliefs, and this understanding strongly constrains our own learning and predictions.

The ability to reason from the perspective of other agents is important for any intelligent machine, such as a household robot or an autonomous vehicle, to predict – and therefore safely interact with – the physical, human world. However, there is no straightforward method to instill such ability in a general sense. while modern-day object trackers are increasingly able to capture and predict complex kinematics (i.e. motion), they cannot yet adequately leverage visual and temporal context to infer the awareness, motivation or intent of agents (i.e. causes for motion), which limits their capacity to truly generalize. No machine can draw the same rich narrative from observing the Heider-Simmel Illusion as a toddler. Kinematics alone may be sufficient for prediction in the vast majority of autonomous driving, but in the long tail of driving, there are likely to be many scenarios which are better predicted by models equipped with stronger psychological intuition [1, 2].

**Goals**

This project aims to study methods for interpreting the awareness, motivation and intent of agents from observations of their behavior and other context in images and video. While the ultimate goal is to develop autonomous "theory-of-mind" frameworks, the initial objective is more pragmatic: to investigate supervised methods through which machines can better leverage visual cues that improve behavioral prediction. This may touch on work from numerous sub-fields, such as gaze prediction, emotion and gesture recognition, and motion analysis. Expected project goals include:

1. generate several large-scale datasets of images and short video clips, where subjects in the scene are annotated with expectations of awareness, motivation and intent. One dataset should be abstract, synthetic and highly specifiable, to allow controlled, toy problem analyses. Further dataset(s) may progress to more general and diverse scenarios (e.g. movie sequences or specific driving or robotics applications such as pedestrian and vehicle interactions)
2. develop supervised methods for prediction in (1), incorporating spatial and temporal context and evaluating the importance of different visual cues
3. demonstrate the benefit of (2) in a real application, such as predicting the likelihood of pedestrians to cross or vehicles to go at an intersection
4. for toy data at least, compare supervised methods from (2) to existing "theory-of-mind" frameworks, particularly with respect to generalization.

**Technical problem statement**

Key technical questions to address are:

- the annotation task is very open-ended. How can we create a sensible ontology which is both rich enough to be interesting and useful, abstract enough to be generalizable, and simple enough to be feasible to annotate?

- how can estimates of awareness, intent and motivation for various agents in a scene be optimally used to predict future outcomes?
- is it possible to create supervised models or "theory-of-mind" models which can generalize better across datasets, to the point of anthropomorphizing abstract unseen video stimuli?

References
[1] R. Brooks. The big problem with self-driving cars is people. IEEE Spectrum, 2017. 1
[2] B.M. Lake, T.D. Ullman, J.B. Tenenbaum, and S.J. Gershman. Building machines that learn and think like people. Behavioral and brain sciences, 40, 2017. 1

## 8.26 CERTIFIED SAFETY FOR SELF-DRIVING CARS

TRI Researcher:      Soonho Kong
Email Address:       Soonho.Kong@tri.global
TRI Thrust:          Driving

**What do you want to try to do?**
The goal of this project is to explore and develop software components in autonomous driving that generate "evidence for safety (Certificates)" in addition to the main computational output. The generated evidence is used to check if the computational output is consistent and safe with respect to its input. When a certificate check fails, the system should have a mitigation plan to handle the anomaly.

The checking process should be fast and straightforward to implement. This exploits the computational asymmetry between finding a solution, which usually involves a search procedure, and checking the solution (usually in $O(n)$). As a result, the size of the trust computing base remains small, which provides high assurance.

The requirement for a component to provide consistent evidence can make the system more robust and secure. To have a compromised component output a malicious result, an attacker needs to do extra work to forge corresponding evidence for the result.

**Why do we care?**
Safety is in our mission statement. We are developing the technology to build a car that is incapable of causing a crash, regardless of the skill or condition of the driver. Toyota Guardian is our approach to this goal. This system oversees the driving environment to mitigate and avoid predicted crashes. However, a natural question follows: "Who is guarding the guardian?". We need an answer for this.

Neither formal verification nor testing is good enough to solve the problem. Formally verifying all components is not realistic: 1) there are known scalability issues in formal methods techniques, 2) for machine learning components, we do not have formal specifications. Testing is the strategy that we rely on in our development. However, this is not a complete technique. It can be used to show the presence of bugs, but never to show their absence.

**Technical problem statement**
Given a software component, we formalize it as a mathematical function, $f(X) = Y$ where $X$ is an input to this component and $Y$ is the output of the component.
- Extend $f$ to $f_c$ such that $f_c(X) = (Y, C)$ where $C$ is a certificate which is a proof establishing the consistency between $f$, $X$, and $Y$.

- Design and implement a proof system and *Checker(C, X, Y )* . Formally verify that the checker meets a requirement.
- Demonstrate an autonomous-driving system with certified components in perception, planning, and control stacks.

## 8.27  AI-BASED PLANNING AND CONTROL OF FUTURE MOBILITY SYSTEMS FROM THEORY TO DEPLOYMENTS

TRI Researcher:  Stephen Hughes, Masanori Yamato
Email Address:  Soonho.Kong@tri.global, Masori.Yamato@tri.global
TRI Thrust:  Driving

**What do you want to do?**
We want to harness research advances in the field of AI to design tools that enable optimal planning and control of future mobility systems. By future mobility systems we mean the growing ecosystem of mobility options of shared, specific-purpose, autonomously capable vehicles to provide on-demand mobility services. Specifically, we are interested in AI-based tools to (1) co-optimize the design of a vehicle (e.g., its capacity, range, level of autonomy, etc.) with the mobility service that such a vehicle will deliver (e.g., last-mile, point-to-point, shared, etc.), (2) intelligently plan the operations for such systems (e.g., in terms of charging infrastructure needs), (3) devise real-time control algorithms to manage such systems in the broader context of an optimized transportation network (e.g., accounting for public transport), (4) implement part of this research by partnering with Toyota-managed mobility systems (e.g., Ha:Mo) and business partners such as Grab.

**Why do we care?**
TRI and Toyota are increasingly investing in AV-related technologies and new forms of mobility (such as micro-mobility vehicles, autonomous shuttles, etc.), thus this research is particularly timely. Specifically, this research will be critical to (1) inform stakeholders on how to best deploy future mobility modes within the transportation network, (2) inform the trajectory for technology development for AVs (e.g., if we knew that the "killer app" for AVs are slow-moving autonomous shuttles, this would constrain the requirements for the development of an AV autonomy stack), and (3) provide state-of-the-art tools to manage large-scale on-demand vehicles fleets.

**Technical problem statement**
The problem of planning and controlling large-scale on-demand vehicle fleets combines features of networked control, optimization, transportation options, and decision making under uncertainty. We plan to leverage the wealth of optimization-based techniques we have developed in the past three years with our Stanford collaborators, along with new advances in AI techniques such as RL, neural network models for forecasting, and data-driven control to tackle such a challenging problem with a fresh perspective. Specifically, this project will entail novel advances in terms of (1) RL techniques for large-scale fleet optimization; (2) meta-learning models for accurate forecasting; (3) data-driven modeling along with optimization-based techniques to reason about interactions among mobility operators and with other infrastructures such as the power network, and (4) co-design techniques for the optimal planning of future mobility systems. We anticipate that these research efforts will provide key advancements to the field and will continue to produce award-winning scientific publications

as in the past three years. Also, along with the Stanford collaborators, we will infuse some of these techniques into real systems. In the past, we deployed a first-of-its-kind optimization-based controller for vehicle rebalancing within the Toyota Ha:Mo system. Currently, we are working with Grab on deploying real-time routing algorithms to best integrate a fleet of fully-controlled vehicles within the current fleet of ride-hailing vehicles. In the future, we plan to seek additional opportunities for technology infusion by leveraging the growing set of partnerships that Toyota is establishing in the mobility space.

## 8.28 COORDINATED ACTIVE SENSING AMONG VEHICLES AT JUNCTIONS

TRI Researcher: Stephen McGill
Email Address: Stephen.McGill@tri.global
TRI Thrust: Driving

### Goal
Vehicles that enter a junction inherently share information asymmetries due to their varying sensor viewpoints, and vehicles may not have enough information to make safe decisions when navigating junctions. Recently, Vehicle-to-Vehicle/Infrastructure techniques have attempted to improve vehicle safety by sharing information, leveraging research on optimal communication techniques. However, these communication techniques have not yet considered the ability for vehicles/infrastructure to change their behavior (as in autonomous vehicles) in order to exchange more useful information. This proposal seeks efforts to characterize scenarios that will benefit from active sensing approaches among vehicles and infrastructure and asks for algorithms that implement active sensing in these scenarios.

### Motivation
Toyota sells millions of vehicles a year, and many of them approach junctions (intersections, merges, roundabouts, etc.) at the same time. However, there is no way to leverage this ad-hoc multiagent system to improve safety. TRI, through its Guardian lens, seeks to make the uncrashable car. Along the way, TRI must ensure that a Toyota does not cause another Toyota to crash. This proposal will move TRI and TMC closer to its Guardian objectives by identifying, in a principled manner, scenarios in which V2V/V2X methods can reduce risk and embodiments that perform this risk reduction.

### Context
Research during University 1.0 led to development of a risk assessment algorithm that operates at junctions, without sharing information. This work is being transferred to TRI and TMC vehicles, highlighting the value of risk measures at junctions. There is still work to do in supporting semi-autonomous path planning, vehicular communication and enhanced risk metrics at junctions. At present, TRI fleet vehicles communicate with infrastructure in California and Ann Arbor, but TRI vehicles do not communicate amongst each other; therefore, this proposal highlights a huge opportunity for research, development and technology transfer to TMC.

### Technical Problem Statement
Real world deployments of V2V/V2X methods must provide clear improvements against baseline systems. This proposal seeks technical accomplishments that consider:

- Formal characterization of scenarios that do and do not benefit from active sensing methods, and by what metric.
- Theoretical metrics and bounds on information gain and risk reduction
- Closed course evaluation among multiple TRI vehicles, with performance metrics (Working with TRI PoC)
- Semi-autonomous action sets to improve safety and information exchange
- Adversarial behavior and approaches to mitigate them

## 8.29  TOYOTA MINIATURE VEHICLE PROVING GROUND

TRI Researcher:      Stephen McGill
Email Address:       Stephen.McGill@tri.global
TRI Thrust:          Driving

**Goal**
We propose a pipeline to characterize challenging driving scenarios and to prove out novel algorithms for modeling and acting in noisy environments and unpredictable agents. University research enters this pipeline, strongly impacting TRI/TMC fleet deployment. Pipeline stages include simulation, miniature vehicle demonstrations and dataset generation, aiding in transfer to TRI and TRI-AD.

**Motivation**
TRI improves road safety by bringing novel and proven technologies to Toyota vehicles. To prove new algorithms, TRI leverages its Test Engineering group to emulate challenging interactions safely in closed course environments. However, safely testing challenging interactions at scale in a closed course remains difficult due to requirements of manpower, space and coordination of multiple vehicles. Novel core autonomous vehicle technologies are driven, increasingly, by datasets, benchmarks and shared reference code. However, principled experimentation of the integration of core technologies continues to lag due to the aforementioned platform and environment constraints when testing field robotics systems. This proposal addresses issues in experimentation and evaluation by enabling a lower barrier to entry via standardized multi-robot closed-loop benchmarks.

**Context**
Research on small scale cars during University 1.0 led directly to technology transfer of a risk assessment algorithm that is on the path to market deployment within TRI-AD. Additionally, TRI sponsored the CARLA challenge, showcasing University research in a common simulated environment. TRI shall build on these results to facilitate additional successful deployments. We will rapidly assess University research on a common and safe platform, via a set of physically embodied vehicles and environments.
Feedback from partners during University 1.0 highlighted the need for real world datasets from the TRI fleet in order to validate and motivate research. After selecting promising algorithms from those showcased in simulation and on small scale cars, TRI will work with partners to collect and to share TRI vehicle data relevant to those algorithms.

**Technical Problem Statement**

When making safe decisions in the real world, autonomous driving algorithms require robustness in the presence of noisy measurements and unpredictable agents.

- We seek proposals on designing challenging urban environments and defining objectives and performance metrics for an autonomous vehicle in the environment. Environment design includes road topology, object placements, road agent behavior and task definitions. Designs must provide fidelity with respect to the real world in order to ensure reliable domain transfer from simulation and miniaturization to full scale environments. Environments should be implemented first in CARLA and then with mini-vehicles. Performance metrics must separate immature and promising approaches, thereby providing value for TRI and TMC. For instance, the recent CARLA challenge required autonomous vehicles to navigate a roundabout and provided success criteria that included the time within the roundabout and accident fault assignment.
- We seek proposals on autonomous vehicle technologies that aid in navigating urban environments, tackling the concepts of "stay on the road," "don't hit things" and "don't get hit." Technologies should operate on both perfect information/actuation, as well as onboard perception/actuation. Algorithms should be implemented first in CARLA and then with mini-vehicles.
- We seek proposals that provide both environments and autonomy algorithms to tackle these environments.
- We seek competition among partners, where teams attempt to outperform each other on common metrics in the challenging environments. TRI seeks to amalgamate complementary technologies from many of its partners, leveraging unique strengths and performance of each partner.

## 8.30  CONTEXTUAL ADAPTATION WITH RISK MEASURES FOR ENHANCED NAVIGATION UNDER INFORMATION DEGRADATION

TRI Researcher:        Velin Dimitrov
Email Address:         [Velin.Dimitrov@tri.global](mailto:Velin.Dimitrov@tri.global)
TRI Thrust:            Driving

**What do you want to try to do?**

We want to advance the capabilities of TRI's Guardian technology through the development of novel methods and tools to achieve safe, reliable, and adaptive navigation and decision-making under information degradation, both intrinsic and extrinsic to the vehicle. The future of Guardian depends heavily on a robust method to quantify and act on risk associated with decisions. The key barrier to successful deployment of Guardian on production vehicles is the development of adaptive behaviors based on contextual information. Navigation and control of autonomous vehicle can be enhanced in real-world environments with imperfect information by allowing adaptation based on quantifiable risk metrics. These metrics can be learned over time and are based on the contextual information from the operating environment around the vehicle, in addition to the underlying vehicle capabilities. The problem breaks down into two key areas which should be the focus of solicited projects. How do we quantify risk as a metric in a robust and repeatable manner and how do we utilize those risk metrics to their fullest extent to make faster, better, and safer decisions in the autonomous or semi-autonomous operation of vehicles?

**Why do we care?**

Long term, Guardian technology is positioned to enable two modalities of vehicle operation: a shared human-autonomy mode and hierarchical full multi-autonomy mode. We anticipate the immediate focus will be on the former to bring enhanced safety behaviors to production vehicles following in the footsteps of capabilities similar to automatic emergency braking (AEB). In addition, Guardian is expected to evolve over the next decade, monitoring and intervening on behalf of other higher-level controllers. Whether Toyota or Toyota-partner deployed, these controllers will likely be developed in a multi-layered hierarchical structure to enable robust, intelligent, and safe vehicle operation. The basis of these interventions will ultimately be decided by Guardian's understanding of risks and the sources of those risks, irrespective of human or autonomous origin.

**Technical Problem Statement**

The problem is complex and requires a holistic approach utilizing concepts from artificial intelligence, model-based design, and traditional controls approaches. Experts in one singular domain are unlikely to come to scalable solutions that exhibit the robustness and reliability inherently needed for commercialization. Successful characterization of risk and risk-based autonomy has already been demonstrated in wide-ranging domains such as high frequency trading and robotic manipulation. A similar result in the field of automated driving could lead to comparable significant advancement of the state of the art. Proposals to address the identified problem will need to address the following technical problems:

- How are the sources of risk in autonomous and semi-autonomous systems quantified in a deterministic manner, when some elements of the system may be non-deterministic?
- How can the evolution of risk be modeled and learned to enable long-term autonomy and navigation of a constantly changing environment?
- Once quantified, how can the risk be controlled or at least mitigated so the vehicle operates in the desired manner, but also optimizes for safety and reliability?
- How does risk cascade in a system controlled by multiple levels and layers of controllers, of both autonomous and human nature?

## 8.31 DESIGN OF RESILIENT MOTION PLANNING FOR TOYOTA GUARDIAN SYSTEMS WITH APPLICATION TO COLLISION AVOIDANCE UNDER UNCERTAINTY

TRI Researcher:     Vishnu Desaraju
Email Address:     Vishnu.Desaraju@tri.global
TRI Thrust:     Driving

The goal of this project is to develop trajectory planning and/or control methodologies that will allow Guardian systems to improve collision avoidance capabilities in a way that is safe and reliable, even in the presence of various sources of uncertainty. Uncertainty may stem from a variety of sources including variations in vehicle dynamics, environmental conditions, and imperfect sensing. While uncertainty poses key challenges for Chauffeur systems as well, Guardian systems must be particularly robust

to in order to ensure a safe operating envelope for the human driver without overconstraining the driver's actions. This project will primarily focus on a subset of this safe operating envelope relating to challenging collision avoidance scenarios, e.g., involving highway speeds, evasive maneuvers that excite the nonlinearities of the vehicle dynamics, or requiring combinations of steering, braking, and accelerating to both avoid the collision and stabilize the vehicle afterwards.

This project will ideally achieve the following objectives:
- Analyze and identify the sources of uncertainty with the greatest impact on the vehicle's ability to safely avoid collisions
- Identify and develop techniques for estimating, learning, bounding, adapting to, or otherwise mitigating the effects of these uncertainties
- Develop a nonlinear trajectory planning/control algorithm that leverages these uncertainty mitigation techniques while accounting for the vehicle dynamics and stability requirements for safe collision avoidance
- Quantify the resiliency of the developed algorithm, e.g., in terms of convergence rates, bounds, operating ranges
- Translate the proposed techniques to computationally efficient implementations and empirically validate the techniques in different challenging collision avoidance scenarios, both at the university and in the TRI Planning and Control architecture

The main value to TRI stems from the fact that uncertainty mitigation will be one of the key factors in expanding the range of scenarios in which Guardian can reliably assist or intervene. So this project aims to provide an analysis of these uncertainty effects, as well as algorithms that are tailored to addressing these effects. The final objective also has a natural tech transfer aspect to it, which will facilitate TRI's in-house evaluation and potential deployment of these techniques.

## 8.32 ADAPTIVE SHARED AUTONOMY BASED ON OPERATOR INTENT FOR ENHANCED PERFORMANCE WITH TOYOTA GUARDIAN

TRI Researcher:     Vishnu Desaraju
Email Address:      Vishnu.Desaraju@tri.global
TRI Thrust:         Driving

The goal of this project is to develop an adaptive shared autonomy framework that leverages formalized models of operator intent to inform near-term motion planning and decision making, thereby improving the performance of Guardian systems.
There are two key components to this project that go beyond just formulating a human-in-the-loop architecture. The first is to develop models of operator intent that not only enable inference and prediction but also provide guarantees and strong correctness properties. This could include bound guarantees on uncertainty, measures of model accuracy relative to the underlying distributions, or model fidelity and convergence properties for different operators/levels of operator proficiency. The second component is to use these models to infer operator intent and integrate this information into a finite horizon decision making or motion planning algorithm as part of a shared autonomy framework. This integration aims to demonstrate the practical utility

of these models to enhance the human-in-the-loop system as measured, for example, in terms of improving performance and safety or reducing the magnitude of autonomy actions needed to keep the vehicle within a safe operating envelope.

This project will ideally achieve the following objectives:
- Develop a representation for near-term operator intent that enables inference of the underlying distribution and prediction of future operator actions
- Perform a rigorous analysis of the proposed intent model to establish its key properties and guarantees
- Develop finite horizon decision making/motion planning algorithms that leverage these models to improve shared control with a human operator
- Translate the proposed techniques to computationally efficient implementations and demonstrate performance with different operators and scenarios, both at the university and in TRI's Planning and Control architecture
- Quantify via empirical evaluation the key properties of these techniques, their adaptation performance, and their implications for statistically safe behavior

The main value to TRI stems from the development of operator intent models with formal guarantees, as leveraging these models will increase confidence that Guardian is strictly aiding and enhancing the human operator's performance. The last two objectives also have key tech transfer aspects that will facilitate TRI's in-house evaluation and potential deployment of these techniques.

## 8.33 MULTI-SENSORY AND MULTI-MODAL FUTURE PREDICTIONS FOR AUTONOMOUS DRIVING

TRI Researcher:     Wadim Kehl
Email Address:      Wadim.Kehl@tri.global
TRI Thrust:         Driving

Typical CNNs are over-confident in their predictions. Moreover, these networks tend to approximate the conditional averages of the target data resulting in over-smooth predictions. These undesired properties render the immediate outputs of those networks unsuitable for the quantification of calibrated uncertainty.

The literature contains approaches to this problem including Mixture Density Networks, tailored Winner Takes All (WTA), Relaxed-WTA, and Evolving-WTA (EWTA).  These approaches generally consider only low dimensional posterior with the assumption of a Euclidean space.  TRI seeks approaches with the multi-dimensional, non-Euclidean and highly non-convex posteriors that govern our real world.

## 8.34 AUGMENTING CONTINUOUS HIGH DEFINITION MAP UPDATES WITH HUMAN INTELLIGENCE

TRI Researcher:     Xipeng Wang
Email Address:      Xipeng.Wang@tri.global
TRI Thrust:         Driving

**One-sentence Summary**
Leverage crowdsourced human intelligence to facilitate continuous update of HD maps.

**Background**
HD maps are a key requirement for building Chauffeur and Guardian technologies. HD maps represent the world to centimeter resolution with sematic information. In our driving stack, HD maps contain two layers. One is a geometric map layer, which fuses data from sensors such as lidar, camera, GPS, IMU, etc. The other layer is a semantic map layer, which represents the road network topology and includes attributes like lane boundaries, crosswalks, and traffic lights.. Building HD maps is framed as an optimization problem shown in the following equation: find the map that is most consistent with the set of observations, where f is the sensor model, X is the world model, Z is the observation, and g is the cost function model.

$$X = argmax\ g(f(X),\ Z)\ X$$

Optimization is required here is because the perception is not perfect. The most challenging problem in mapping is data association (e.g. Is this stop sign the one you saw yesterday?). Bayes' theorem tells us the more information we get, the better map we can build. But this is under the assumption that we know the characteristics of sensor models, in other words, the noise level of your perception system. Because we have a better understanding of the sensor models, such as GPS, IMU, ranging, building SLAM Map is nearly an automated process.

**What do we want to do?**
We would like to build accurate continental scale maps with continuous updates. Two fundamental approaches to solve this problem: improving perception system and improving human intervention process. In this proposal, we mainly focus on the idea of improving human intervention. More specifically, the things we would like to achieve are:
- Leverage crowdsourcing to perform human intervention so that we can perform large scale mapping.
- Leverage human intelligence to validate and improve the quality of maps.
- Leverage online crowdsourcing to perform continuous online update of maps.

**Why do we care?**
Both Chauffeur and Guardian technology requires maps. Considering using maps as a sensor, the better map we can get, the safer technology we can provide. There are hundreds of million Toyota cars on the road every day across the world. If Toyota wants to provide Chauffeur and Guardian technology to all customers, we need to think about how to provide continental scale maps with high quality.

**What are the key technical challenges?**
- Leverage crowdsourcing to perform human intervention so that we can perform large scale mapping.
  - How can build tools for people with zero CS/SLAM/robotics knowledge?
- Leverage human intelligence to validate and improve the quality of maps.
  - How can perform minimal human interventions to achieve the best performance?
- Leverage online crowdsourcing to perform continuous online update of maps.

Comparing machines with humans for perception tasks (e.g. detecting map changes), machines can do a large amount of work but with lower quality. How can we combine the advantages of both humans and machines to provide accurate continuous map update to map customers?

## 8.35 LEARN DESCRIPTORS FOR OBJECT-LEVEL FEATURES

TRI Researcher:        Xipeng Wang
Email Address:        Xipeng.Wang@tri.global
TRI Thrust:            Driving

**One-sentence Summary**
Learn feature descriptors for high level structural features that encodes geometric relations in the sensor view to improve data association.

**Background**
HD maps are a key requirement for building Chauffeur and Guardian technologies. HD maps represent the world to centimeter resolution with sematic information. In our driving stack, HD maps contain two layers. One is a geometric map layer, which fuses data from sensors such as lidar, camera, GPS, IMU, etc. The other layer is a semantic map layer, which represents the road network topology and includes attributes like lane boundaries, crosswalks, and traffic lights.. Building HD maps is framed as an optimization problem shown in the following equation: find the map that is most consistent with the set of observations, where f is the sensor model, X is the world model, Z is the observation, and g is the cost function model.

$$X = argmax\ g(f(X), Z)\ X$$

Optimization is required here is because the perception is not perfect. The most challenging problem in mapping is data association (e.g. Is this stop sign the one you saw yesterday?). Bayes' theorem tells us the more information we get, the better map we can build. But this is under the assumption that we know the characteristics of sensor models, in other words, the noise level of your perception system. Because we have a better understanding of the sensor models, such as GPS, IMU, ranging, building SLAM Map is nearly an automated Process.

**What do we want to do?**
We would like to build accurate continental scale maps with continuous updates. Two fundamental approaches to solve this problem: improving perception system and improving human intervention process. In this proposal, we mainly focus on the idea of improving perception system. More specifically, the things we would like to achieve are:
- Build a perception system that can provide high precision/recall semantic labels using multimodal sensing, e.g., camera or lidar.
- Generate descriptors for high level semantic labels such as poles, signs, etc.

**Why do we care?**
Both Chauffeur and Guardian technology requires maps. Considering using maps as a sensor, the better map we can get, the safer technology we can provide. There are hundreds of millions of Toyota cars on the road every day across the world. If Toyota wants to provide Chauffeur and Guardian technology to all customers, we need to think about how to provide continental scale maps with high quality.

**What are the key technical challenges?**
- Build a perception system that can provide high precision/recall semantic labels using multimodal sensing, e.g., camera or lidar.
    - Push the state-of-the-art performance with real time performance.
- Generate descriptors for high level semantic labels such as poles, signs, etc.
    - Majority perception systems for SLAM learn descriptors only for keypoints. But these low-level keypoint features are not stable under different light conditions. Hence, we choose high level structural features such as pole, signs, etc. The key challenge here is to learn descriptors for these high-level features that encode spatial relations among them.

## 8.36 DATA-DRIVEN SCENARIO GENERATION WITH HETEROGENEOUS ROAD USERS

TRI Researcher:      Jonathan DeCastro
Email Address:       Jonathan.DeCastro@tri.global
TRI Thrust:          Driving

**What are we trying to do, and why?**
Understanding the risks of autonomous vehicle (AV) decisions will require leveraging large-scale, naturalistic data for all types of road users. Such data will be useful for constructing multi-agent prediction models, including those of vehicles, pedestrians, bicyclists, and motorcycles (i.e. vulnerable road users), as well as understanding critical scenarios for safety . The nature of the interactions and severity of interactions are complex and the heterogeneity of the risks involved in AV failure are especially important when analyzing AV safety.

Large-scale naturalistic driving data is important for several reasons. First, a viable safety case for AVs will require testing the AV in a manner guided by statistics of the population of road users (all possible configurations and parameterizations of agents within some statistical bound). Due to the vastness of the possible scenarios an AV could encounter, a sufficiently-rich data source will contain corner cases that can direct testing in a principled way. Second, a rich data source is useful for constructing predictive models for evaluating risk in online decision-making.

To date, studies have been made to better understand pedestrian-vehicle interactions, while extracting the critical factors involved in features and identifying representative corner-cases. Many of these rely on manual annotation of corner cases, which is limiting from the standpoint of both scalability and objectivity. Independent of corner cases, works have also focused on data-driven modeling of road user behaviors. There has also been much work in inference and prediction of road users, as have simulation-driven approaches, via Monte Carlo techniques including importance sampling and cross-entropy. In forming a viable testing strategy, it is required to understand risk from naturalistic datasets involving all possible road user types in order to make use of the techniques both for simulation and agent predictions for AV decisions.

**Technical problem statement**
We seek proposals that allow for leveraging statistically-significant data sources involving interactions between homogenous and heterogenous road users, including vulnerable road users, to allow for 1) automated analysis for creating an objective testing strategy that includes corner

cases and statistical coverage guarantees, and 2) training generative agent models for pedestrians, bicycles, and motorcycles. The focus of the effort will include the following elements:

- Learning metrics from large-scale data sources to allow for deciding a set of representative scenarios for achieving test coverage and profiling the space of possible scenarios from common events to corner cases, and everything in between.
- Working with TRI to construct probabilistic generative models for predicting agent trajectories for the dual purpose of agent prediction and building a data- and requirements-driven testing strategy. The data source will be used to characterize and construct statistical guarantees for these models.

## 8.37 CAUSAL MULTI-AGENT PREDICTION

TRI Researcher:        Kuan Lee, Adrien Gaidon
Email Address:        Kuan.Lee@tri.global, Adrien.Gaidon@tri.global
TRI Thrust:        Driving

**What do you want to try to do?**
Our goal in this project is to develop a multi-agent framework for ado vehicle trajectory prediction that can dynamically model and forecast vehicle interactions in a probabilistic, causal, and scene-aware fashion. We would like to explore a key missing component in the current state of the art: how to leverage observed and forecasted interactions to reduce uncertainty by reasoning about cause and effect via feedback loops across perception, prediction, and planning models.

**Why do we care?**
Prediction is an open research problem on the critical path to autonomy. Autonomous driving indeed requires reasoning about the uncertain future behavior of agents in a variety of driving situations. In multi-agent settings, each agent's behavior affects the behavior of others. Motivated by people's ability to reason in these settings, a multi-agent trajectory prediction is needed to forecast multi-agent interactions from observations extracted from cameras and LIDAR. Furthermore, because intents are not observable in general, this modeling and reasoning is probabilistic in nature. Careful reasoning is required to explore dynamically the best trade-offs between risk seeking (e.g., willing to pay a high cost in probability of collision) vs risk averse (often resulting in drastically reduced driving abilities compared to good human drivers).

**Technical problem statement**
Trajectory prediction is one of the major research areas in autonomous driving. Many existing works have proposed multi-agent prediction approaches/frameworks [1-3] that learn a probabilistic forecasting model for ado agents' trajectory prediction, where path history and environment cues (either dynamic or static) are usually leveraged as a prior in the framework. Recently, goal-conditioned methods [1,2] get more and more attention since they not only encode past information but also refer to the goal/destination in the future.

On the other hand, to precisely predict trajectory in the future, probabilistically modeling interactions between multi-agents (including ego and ago agents) is essential. To this end, social cues are typically used to capture agents' trajectories towards their inferred goals while avoiding collisions. Many studies applied RNN to model trajectory prediction with social cues

[4-7]. Several works [5-7] based on adversarial learning were recently proposed to model agent-to-agent interactions in a GAN framework. Moreover, Graph Networks [1] and attention mechanisms are also a promising direction to address such agent-to-agent relationship.

We would like to investigate a compositional multi-agent scene-aware framework for trajectory prediction based on causal inference and uncertainty modeling. Moreover, we would like to explore the capability of causal inference in goal-conditioned prediction, and how it benefits path planning and decision making in autonomous driving applications.

[1] B. Ivanovic and M. Pavone. The Trajectron: Probabilistic Multi-Agent Trajectory Modeling with Dynamic Spatiotemporal Graphs. ICCV 2019.
[2] N. Rhinehart, R. McAllister, K. Kitani and S. Levine. PRECOG: PREdiction Conditioned on Goals in Visual Multi-Agent Settings. ICCV 2019.
[3] Y. Tang and R. Salakhutdinov. Multiple Futures Prediction. NeurIPS 2019.
[4] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, F.-F. Li, S. Savarese. Social LSTM: Human Trajectory Prediction in Crowded Spaces. CVPR 2016.
[5] A. Gupta, J. Johnson, F.-F. Li, S. Savarese, A. Alahi. Social GAN: Socially Acceptable Trajectories with Generative Adversarial Networks. CVPR 2018.
[6] A. Sadeghian, V. Kosaraju, A. Sadeghian, N. Hirose, S. H. Rezatofighi, S. Savarese. SoPhie: An Attentive GAN for Predicting Paths Compliant to Social and Physical Constraints. CVPR 2019.
[7] J. Ho and S. Ermon. Generative Adversarial Imitation Learning. NeurIPS 2016.

## 8.38  SEMI-SUPERVISED LEARNING TOWARDS ROBUST 3D DETECTION AND TRACKING

TRI Researcher:      Jie Li, Kuan Lee, Rares Ambrus
Email Address:       Jie.Li@tri.global, Kuan.Lee@tri.global, Rares.Ambrus@tri.global
TRI Thrust:          Driving

**What do you want to do?**

Our goal in this project is to explore and develop object detection algorithms that work better with downstream tracking modules through better uncertainty estimation. In addition, we would like to deploy a better semi-supervised object detection learning framework by extending supervision from labeled frames to neighboring un-labeled frames with tracking in the loop.

Desired output of this project would include but are not limited to the following:

- A 3-D based object detection and tracking framework that generate better tracking results with high efficiency, ideally in real-time.
- A semi-supervised pipeline that could make use of unlabeled frames in a dataset to learn a better object detection and/or tracking algorithm beyond supervision.

**Why do we care?**

Object detection and tracking are key modules in both autonomous driving and indoor robotic platforms under development in TRI. While baseline solutions exists, heuristic approaches are currently employed to provide object detection uncertainty in location measurements, e.g. a diagonal constant matrix. A positional uncertainty aware object detection and tracking system will provide better uncertainty propagation to the tracker resulting in more reliable tracking and uncertainty estimation, and therefore, will increase the robustness of our systems of both autonomous driving and robotics.

On the other hand, the resulting better tracking result of the above mentioned framework can help generate more training data in unlabeled frames to the object detection, resulting in a more scalable solution to leverage Toyota's data advantage.

**Technical problem statement**

Current state-of-the-art object detectors can provide a proxy of epistemic and semantic uncertainty of a detection, c.f. object detection score and semantic label distribution, but no uncertainty is provided on the spatial aspect or object localization. On the other hand, a lot of widely used tracking system builds its state and uncertainty estimation on physical locations. Thus, when a tracking system receive a detection measurement from detector, an uncertainty over the location is essential. Most conventional systems use a fixed empirical matrix to model the detection uncertainty, resulting sub-optimal tracking results. This gap between two connected modules has recently draw more attention. More researchers are looking into this connection from the aspect of metrics, object detection, and tracking. One of our ongoing projects with Stanford is also looking at this problem, from the tracking side, and propose a filtering-based algorithm that won the 1st place of nuScenes Tracking Challenge in NeurIPs 2020.

In the area of data bootstrapping and learning beyond supervision, more recent work has also been proposed to improve the object detection algorithm, leveraging tracking algorithms in propagating key frame labels or building up temporal constraints between labeled and unlabeled frames.

Thus, proposing a better solution to fill the gap between the object detection output and the tracking algorithm input is a two-bird-with-one-stone effort. We are looking to improve the overall performance of a detection-tracking system while also making use of more data than the labeled ones.

## 8.39 PEDESTRIAN AND CROWD BEHAVIOR UNDERSTANDING

TRI Researcher:     Kuan Lee, Adrien Gaidon
Email Address:      Kuan.Lee@tri.global, Adrien.Gaidon@tri.global
TRI Thrust:         Driving


**What do you want to do?**

The bottleneck in deep learning is the dependence on a large amount of manual annotations. Labeling at scale is prohibitively costly for behavior understanding, especially in crowded scenes. Annotators indeed need to look at the entire data and all pedestrians to detect and label

events of interest without any false negatives or false positives. Our goal in this project is to develop a pedestrian and crowd action recognition system that can maximally leverage unlabeled data (self-supervised learning) or with as few manual labels as possible (semi-supervised learning).

**Why do we care?**

Recognizing and predicting actions of traffic participants, especially Vulnerable Road Users (VRUs), is a fundamental capability in Automated Driving and ADAS systems. According to the 2017 crash report by NHTSA, pedestrian fatalities in US urban areas increased by 46 percent since 2008. According to the NTSB report on the 2018 Uber ATG fatal crash, one of the causes for the crash is the inability of the Automated Driving system to model the behavior (appearance, motion, pose, intent, and context) of Elaine Herzberg, leading to the tragic accident. This evidence and the current progress in the research community highlights that even state-of-the-art Computer Vision and ML systems still fall short from basic safety requirements. As these techniques improve with data, the best strategy to solve this issue is by scaling up training sets (e.g., in order to see more jaywalkers), but this is limited by the need for manual inspection and labeling of all the data, hence the current insufficient amount of data used to develop this technology. Furthermore, the difficulty in modeling dynamic causal interactions between agents and their environments (as in the case of Elaine Herzberg) requires new models and learning algorithms.

**Technical problem statement**

In order to scale up training data, one needs to identify the right set of inductive biases that can replace large scale manual labeling. TRI is interested in two key problems in this space:

1) semi-supervised video captioning of driving scenes, i.e. associating rich textual descriptions to snippets of driving logs, by leveraging their inherent structure as inductive bias (e.g., rules of the road, driver demonstrations, attention mechanisms);

2) self-supervised action recognition using spatio-temporal graph networks to dynamically model (and forecast) rich interactions in traffic scenes (potentially crowded).

## 8.40 LEARNING TRANSPARENT HRI

TRI Researcher:      Guy Rosman
Email Address:       Guy.Rosman@tri.global
TRI Thrust:          Driving


**Background**

In both the Guardian system and more general applications in robotics, there's a need for systems that interact with humans in a self-explainable way. While a standard approach to UI design involves hard-coding of the interface and deployment after testing, we expect the capabilities of these systems to adapt on the go, and the pattern of interaction with the user to be more complex, meriting novel approaches that allow these systems to adapt their interface as a computational problem. Impact to TRI: As we extend the capabilities of our vehicle stack, we

need a way for machines to better align with the human driver, and make themselves as easily understood as possible. Proposals should explore computational approaches towards this goal, resulting in more natural paradigms for car interface upgrades.

**Aims**

The aim of this project is to explore novel approaches for UX design that define the interaction problem in terms of planning goals, or a joint computational model, and demonstrate ability to adapt to the users and create a low-cognitive-load, easily usable, HAI. Projects under this proposal should demonstrate the approach for either or both the driving and robotics domains.

**Technical Problem Definition**

In the driving domain, projects should define the interaction between the car and the human and reason about risk and human comfort/satisfaction in a quantitative, computable, way. They should allow for multiple ways of shared autonomy /dialogue systems, and optimize the approach taken by the car UI and actions on the road. The system should demonstrate ability to perform multiple types of shared autonomous interactions, and reason about their joint utility under the model. The project should demonstrate in experiments the system's ability to choose interaction modes, define metrics of optimality for the problem, and demonstrate the system's capability to adapt to different users and/or temporal changes. Standard approaches during the experiments should be used to validate the user reaction to the system and the computational model. In the robotics domain, experiments should validate system's capability to adapt the mode of interaction between a human and an assistive robot in a variety of scenarios, including both real demonstrations and simulated interactions.

**Data, Platforms, Deliverables**

The proposal should define specific deliverables, including APIs and a basic implementation in one of the domains. In simulated experiments, there should be a repeatable benchmark to compare results against reasonable baselines. For driving, CARLA or other simulators can be used as a benchmark to demonstrate shared autonomy and warnings systems. For the robotics domains, simulation environments should include standard environments and simulators, to be specified in the proposals.

## 8.41 LEARNED REPRESENTATIONS FOR 3D OBJECT DETECTION

TRI Researcher:     Rares Ambrus
Email Address:      [Rares.Ambrus@tri.global](mailto:Rares.Ambrus@tri.global)
TRI Thrust:         Driving

**What do you want to do?**
The aim of this project is to push the boundaries of 3D object detection in the case of imperfect data, and learn representations that can model and identify the underlying surface with increased robustness. Specifically we would like to develop representations that are robust to missing data (e.g. black cars or rain), sparse(r) data (e.g. a 4-beam lidar vs a 64-beam lidar), different sensing modalities (e.g. Velodyne vs Luminar vs depth-from-images). Additionally, we would like to close the loop between sensing and detection, i.e. our representations should not only be robust to imperfect data but ideally also be able to feedback and correct the raw data itself (e.g. depth

completion). Finally, by better modeling the underlying structure, we will also enable knowledge transfer between domains (indoor, outdoor, night, day, different cities) or across sensors (lidar, depth-from-images, radar).

**Why do we care?**
3D object detection from Lidar data has seen significant progress in recent years both at the fundamental level: points as unordered sets, uneven sampling, part decomposition, finding correspondences; as well as at the implementation level (real-time lidar detector). Our aim in this project is to further improve upon the state-of-the-art and tackle situations (i) where our sensors perform poorly (e.g. when dealing with reflective surfaces or during rain), (ii) when only sparse data is available and (iii) when different sensing modalities are available. The methods developed will also expand our understanding of what is needed for successful 3D object detection in challenging situations, as well as how to leverage experience from other domains (e.g. indoor) or sensors (e.g. sparse to dense lidar) to improve performance.

**Technical Problem Statement**
In order to increase the robustness and applicability of current 3D object detectors, TRI is interested in:
- Increasing the robustness of current detectors in challenging situations (e.g. rain), when the noise characteristics perform differently that under nominal conditions
- Developing approaches for dealing with sparse data (e.g. a 4-beam lidar vs a 64-beam lidar)
- Leveraging cross-domain and cross-sensor data to bootstrap learning and improve performance.

## 8.42 PRACTICAL DEEP LEARNING THEORY FOR GENERALIZATION IN THE LONG TAIL

TRI Researcher:      Adrien Gaidon, Nikos Arechiga
Email Address:       Adrien.Gaidon@tri.global, Nikos.Arechiga@tri.global
TRI Thrust:          Driving

**What do you want to do?**
- Develop theoretical understanding of Deep Learning generalization under real-world data distributions, in particular in the long tail regime (i.e. few frequent cases, many rare cases).
- Derive practical algorithms from the theory to improve empirical generalization performance in real-world datasets (imbalanced, biased, noisy, etc).

**Why do we care?**
Understanding generalization is currently one of the main pursuits in Machine Learning, because previous mathematical tools are not explaining observed behavior of modern deep neural networks (cf. Ben Recht's work, best paper at ICLR'17, and NeurIPS'19 best paper award on the inadequacy of uniform convergence as a tool to understand generalization). This is a big problem in safety critical applications like Automated Driving or Robotics, where currently there is no guarantee that a deep net will perform as expected after deployment, even in the same environment as the training one, especially on rare events. Getting a more solid understanding

(under real-world assumptions) will enable to lessen the better on physical and statistical testing, which are the biggest bottleneck to deploying ML models at scale for automated driving.

Furthermore, improving theoretical understanding typically results in new practical algorithms and breakthroughs (cf. our NeurIPS'19 paper with Tengyu Ma at Stanford https://papers.nips.cc/paper/8435-learning-imbalanced-datasets-with-label-distribution-aware-margin-loss). In addition, theoretical understanding of generalization performance under real-world data distributions will also enable better uncertainty modeling, especially to characterize when a model is "out of domain" and hence should not be trusted. Finally, understanding generalization is key to developing data-efficient ML algorithms that can reduce the cost of collecting human labels and improve the prediction accuracy of TRI products. Vision models on autonomous cars have to generalize quickly to rare events (e.g., animal crossing or children running on the road to catch a ball) and new environments (e.g., a small town in rural areas) without much human supervision. Understanding the rare events is critical for safety features on the cars, especially for scalable self- or semi-supervised models.

**Technical problem statement**
We aim to develop principled large-scale ML algorithms that can a) train faster and generalize better, b) generalize better in the heavy tail of the data including in noisy or partially labeled scenarios, and c) transfer to new domains with fewer labeled examples.

## 8.43  SCALABLE SELF-SUPERVISED LEARNING FOR 3D SCENE UNDERSTANDING

TRI Researcher:      Vitor Guizilini,Adrien Gaidon, Rares Ambrus
Email Address:       Vitor.Guizilini@tri.global, Adrien.Gaidon@tri.global,
                     Rares.Ambrus@tri.global
TRI Thrust:          Driving

**Technical problem statement**
Self-supervised learning is key for scalability, since it allows models to be trained with raw input data, without the need for ground-truth information of any kind. This is achieved by introducing priors that are inherent to the input data itself, and by learning to produce models capable of operating under these constraints the task at hand is solved as a by-product. This is the case with monocular depth estimation, where a single image is used as input and a depth map is generated as output. Even though this is an inherently ambiguous and ill-posed problem (the same image could have infinite possible depth maps), recent progresses in deep learning have shown that it is possible to recover accurate depth information from single images with unprecedented levels of details.

Furthermore, new developments are rapidly surfacing, and cameras are becoming a viable alternative to more expensive and power-consuming LiDAR sensors. Decreasing this gap, by leveraging the unique and attractive properties of monocular imagery (i.e. ubiquity, portability, ease of access, high frame-rate, high resolution, etc.) will be key in bridging the gap between camera and LiDAR information. One of these unique and attractive properties, ubiquity, enables training at massive scale, using information from any visual sensor, however most models currently available in the literature are developed with a single camera sensor in mind, and trained using information from that single sensor. Once training is done, the learned features are transferred via fine-tuning on different camera configurations, but this comes at the expense of

"forgetting" the old configuration. Recent works have started to introspect what is actually being learned inside depth networks, and there is a strong geometric prior embedded in the network (which is to be expected, since the self-supervision uses strong geometric cues in the training stage). This geometric prior includes parameters such as camera intrinsics/extrinsics, which cannot be easily transferred between different datasets. A few works, e.g. Angelova et al., aim to learn from video sequences without such information, by also learning the camera intrinsics, however we believe this is not the case in most real applications, where the intrinsics is actually known, but it changes between frames.

**What do you want to do?**
Our goal with this project is to increase the scalability of TRI's in-house self-supervised learning algorithms, particularly for monocular depth estimation. This will be done with the development of camera-agnostic algorithms, that are capable of producing accurate depth estimates regardless of which camera configuration is being used. The generation of such algorithms will have the following effects:

- Take scalability to the next level, allowing not only training on massive amounts of data from a single camera configuration (i.e. intrinsics + extrinsics) but on data from any vehicle in any condition.
- Enable 360 monocular depth estimation with the same model. The alternative right now is to learn individual models for each camera, which is very costly memory-wise, since all models have to be stored at the same time, and decreases the amount of data available for each model to learn from.
- Leverage our large amounts of unlabeled data, which will play a crucial role in the first stages of the project, since in a virtual environment it is possible to easily modify and benchmark the performance of various algorithms for different camera configurations.
- Camera-agnostic algorithms can also be seen as a form of domain adaptation (a geometric one), and overcoming that challenge could be one more step towards training in one scenario and testing in another (i.e. virtual vs real or United States vs Japan).

**Why do we care?**
Toyota is currently investing heavily in 3D object detection algorithms for autonomous driving. While these will be initially LiDAR-based, there is interest in extending these algorithms to also be vision-based, due to economic factors, inherent limitations in LiDAR (i.e. poor environment conditions, such as rain), Toyota's current data advantage (short and mid-term) lying in massive amounts of monocular video data (e.g., from tens of millions of cars with TSS3), and to enable robust sensor fusion. Depth estimation is the key component required to lift 2D images to 3D, and therefore any improvements in depth estimation will translate to better 3D representations of the environment, that will in turn lead to better vision-based 3D object detection results. Within depth estimation, TRI has chosen self-supervised monocular learning because it can leverage all the data it collects at training time, and has shown in recent published work that increasing the amount of training data indeed leads to better models, to the point where self-supervision is on-par with fully supervised methods, even though it operates on a much more challenging scenario. In order to continue in this direction and further increase scalability, camera-agnostic algorithms are a natural choice, since they would:

- Truly allow the use of all available video sequences (TRI or otherwise), regardless of vehicle and relative position.

- Enable 360 training and inference, which would make the resulting vision-based pointclouds closer to the LiDAR configuration.
- Lead to high-quality publications in top scientific venues, since this is a very sought-after topic with competitive leaderboards that are used to communicate scientific and technical leadership.

## 8.44  LEARNING SCENE GRAPHS FOR COMPOSITIONAL VISUAL REASONING

TRI Researcher:      Adrien Gaidon, Rares Ambrus, Kuan Lee
Email Address:       Adrien.Gaidon@tri.global, Rares.Ambrus@tri.global, Kuan.Lee@tri.global
TRI Thrust:          Driving

**What do you want to try to do?**
The aim of this project is to combine scene graph representations with deep learning methods to improve visual reasoning. Scene Graphs represent scenes as directed graphs, where nodes are objects and edges give relationships between objects. Graph convolutional networks represent a new class of operations that pass information along graph edges. We need methods that can predict and reason over scene graph representations over large scale (potentially novel and unmapped) environments for fast adaptation. Using this symbolic representation, we can model state transitions (e.g. planning prediction or planning), reason about realism (e.g. does it make sense that we have detected a car floating above the ground?), as well as generate synthetic data given a particular symbolic representation of the environment (e.g. generate realistic images containing two cars at an intersection).

**Why do we care?**
For safe and smooth navigation while driving, humans take important and intuitive decisions. These decisions are the result of sequences of actions and interactions with others in the scene. Human drivers perceive the scene and anticipate the driving conditions so that they can make decisions for next steps. However, machine vision is only a single (faulty) component and it requires higher complexity of visual reasoning for proper scene-level analysis. For autonomous driving applications, accurate scene analysis benefits not only perception, but also all further downstream tasks such as planning and control. Finally, a structured representation enables causal reasoning, principled uncertainty propagation, and faster adaptation to new scenarios and environments by leveraging higher order perceptual abstractions and their relations.

**Technical problem statement**
Given the large amounts of data collected by TRI, we are interested in learning symbolic representations that are structured and probabilistic, allowing us to reason under uncertainty about scenes at a high level that is interpretable. Leveraging these representations, we aim to (i) validate current models about the environment online; (ii) decompose dense scene representations (e.g. images or point clouds) into a key set of actors and interactions which would facilitate prediction and planning; (iii) generate synthetic scenes that can be used to validate our models.

## 8.45 SELF-SUPERVISED COMPOSITIONAL REPRESENTATION LEARNING OF VIDEOS

TRI Researcher:       Adrien Gaidon
Email Address:        Adrien.Gaidon@tri.global
TRI Thrust:           Driving

### What do you want to try to do?
Despite the large progress in many areas of computer vision due to adoption of deep learning, applying these techniques in the wild remains a challenge. One of the main hurdles is the fact that deep learning approaches are data-hungry, and the distribution of concepts in the world is inherently long-tailed. That is, only a few categories are frequent, and the others are increasingly rare. This problem is especially severe in the domain of videos, where the data is more expensive to obtain, and the distributions are even more skewed. We propose to address this issue by learning compositional video representations that don't treat video as uniform spatio-temporal volumes of pixels, but instead represent them as in a highly-semantic space of objects and interactions.

### Why do we care?
The problem of recognizing rare scenarios in videos is central to the domain of self-driving vehicles. Indeed, most of the road scenarios can be easily handled by existing solutions, it's the 0.01% percent of rare events, such as never before seen object appearing on the road, or unexpected maneuvers of other agents, which prevent the broad adoption of these systems. And it's precisely the rare nature of these events that does not allow simply collecting large amounts of training data for them. It is thus crucial to develop novel representations that are capable of generalization to unseen examples in videos from a few or no examples.

### Technical problem statement
Existing video representations naively extend deep architectures for image recognition to the video domain by treating videos as spatio-temporal volumes of pixels. This approach naturally requires large amounts of training data for representation learning, since it has to discover all the important information for understanding videos, from notions of object and categories to the types of interactions between them, from scratch. We propose to instead directly model videos as combinations of objects and interactions, learned on existing images and video datasets. This would greatly reduce the hypothesis space for recognizing higher-level semantic concepts in videos, thus increasing generalization abilities of the model.

Some of the most relevant papers that are publicly available:
- http://openaccess.thecvf.com/content_ICCV_2019/papers/Tokmakov_Learning_Compositional_Representations_for_Few-Shot_Recognition_ICCV_2019_paper.pdf
- http://openaccess.thecvf.com/content_CVPR_2019/papers/Zhang_A_Structured_Model_for_Action_Detection_CVPR_2019_paper.pdf
- http://openaccess.thecvf.com/content_ICCVW_2019/papers/HVU/Dave_Towards_Segmenting_Anything_That_Moves_ICCVW_2019_paper.pdf
- http://openaccess.thecvf.com/content_ICCV_2019/papers/Wang_Meta-Learning_to_Detect_Rare_Objects_ICCV_2019_paper.pdf

## 8.46  INTERACTION-AWARE CONTROL FOR CARS THAT ANTICIPATE AND EXPLAIN COMPLEX ENVIRONMENTS

TRI Researcher:     Adrien Gaidon, Guy Rosman, Allan Raventos
Email Address:      Adrien.Gaidon@tri.global, Guy.Rosman@tri.global, Allan.Raventos@tri.global
TRI Thrust:         Driving

**What do you want to do?**
Estimate, from large scale demonstrations, human-robot interaction models that can robustly predict and explain human actions and reactions to automated cars in challenging driving scenarios.

**Why do we care?**
An essential component for interaction with autonomous cars is learning a predictive model of human behavior through collected data. However, today's techniques usually learn human models that are trained on a small collected dataset hence are not guaranteed to be accurate in all scenarios. Ideally, the decision-making capabilities of autonomous agents must be robust to potential flaws in the learned human model in order to guarantee successful objective completion.

**Technical problem statement**
Typically, the problem of modeling human driving behavior is addressed as an inverse reinforcement learning (IRL) problem. However, these techniques are developed based on the assumption of collecting expert trajectories and are prone to challenging scenarios not present in the training data. In this proposal, we first plan to study the robustness of these learned models, generate other structures and models that explain the novel scenarios, and then automatically generate test cases that find these novel scenarios. There are three key research directions:
- robustness analysis of learned human models;
- explainable modeling of human driving behavior in complex scenarios;
- generating scenarios with risky learned human models.

## 8.47  PHYSICAL AND FUNCTIONAL INDUCTIVE BIASES FOR VISUAL REPRESENTATION LEARNING

TRI Researcher:     Adrien Gaidon, Wadim Kehl
Email Address:      Adrien.Gaidon@tri.global, Wadim.Kehl@tri.global
TRI Thrust:         Driving

**What do you want to do?**
Driving scenarios are rich and diverse, yet follow very rigorous assumptions about the general scene layout in terms of environment and physical actors. We want to leverage inductive biases rooted in physics and geometry to allow us to extrapolate structured prior knowledge to unseen scenarios. From this, our goal is to factorize the scene into separable, parametric entities such as cars, drivable surface, people, vegetation and man-made structures such as lane markings, poles and buildings.

**Why do we care?**

Supervised learning via manual labeling does not scale to large amounts of data. In order to leverage petabytes of raw sensory data, we must instead find ways to inject into the learning process information about the world that we know to hold true in generality such as gravity, object permanence, or the physics of light travel. This is required to perform automatic scene analysis/decomposition that would enable mining for rare events/scenarios, as well as promising paths towards resynthesizing novel scenes (e.g. American roads with Japanese cars) to create rich simulation from our data or check for robustness of our ML models. Furthermore, such decompositions could be used for HD map building, congestion and traffic flow analysis etc.

**Technical problem statement**
Provided with many driving logs that include sensory information, we would like to find the physical and functional inductive biases (i.e. priors or regularizers) that can help learning in a self- or semi-supervised setting.

## 8.48 EFFICIENT PERCEPTION FROM VIDEOS AT SCALE FOR AUTONOMOUS DRIVING

TRI Researcher:     Sudeep Pillai, Quincy Chen, Chao Fang Allan Raventos
Email Address:      sudeep.pillai@tri.global, quincy.chen@tri.global, chao.fang@tri.global, allan.raventos@tri.global
Thrust:        Driving

**What do you want to try to do?**
Automated vehicle fleets today typically collect PBs of data on a weekly basis, and are expected to significantly grow in the upcoming years. While typical autonomous vehicle platforms collect data from multiple sensors (LiDAR, Radar, Camera, CAN bus, GPS/INS etc), video data accounts for more than 50% of the total data collected in these systems. In order for Toyota to truly take advantage of its "Data Advantage," we need to be able to build and leverage efficient tools for video processing in order to query and learn from our autonomous vehicle fleets in a scalable manner. More specifically, we would like to build tools for efficient computation, indexing and querying of content in videos collected by autonomous vehicle fleets.

**Why do we care?**
A unique advantage of Toyota as #1 automobile manufacturer is that we will soon have access to "Toyota-scale" sensory data that is order-of-magnitude larger than what our competitors have. This, however, imposes a unique challenge as well: As the fleet size grows, the volume of video data that needs to be processed will be challenging and especially expensive. In order to truly take advantage of the value of data collected by the Toyota fleet, we need to develop effective mechanisms to query video content in a large-scale fleet setting, and leverage this to efficiently learn from diverse scenarios/experience.

## 8.49 PROGRAMMATICALLY BUILDING AND MANAGING TRAINING DATA FOR AUTONOMOUS DRIVING

TRI Researcher:     Sudeep Pillai, Allan Raventos, Quincy Chen, Chao Fang, Adrien Gaidon, Dennis Park

Email Address:     sudeep.pillai@tri.global, allan.raventos@tri.global, quincy.chen@tri.global, chao.fang@tri.global, adrien.gaidon@tri.global, dennis.park@tri.global
Thrust:       Driving

**What do you want to try to do?**

Automated driving solutions have mostly resorted to supervised learning as their primary mode for training of ML models for robust perception. As a consequence, the autonomous vehicle industry has been focused on amassing large volumes of labeled data to have a competitive edge in the deep-learning era. At TRI, we realize the need to go beyond supervised learning for automated driving, especially in computer vision problems that are seeing great progress with strong supervision today.

The goal of this project is to identify and establish fundamentally novel methods to collecting, curating and managing training data for autonomous driving, specifically in order to motivate scientific efforts in semi-supervised, weakly-supervised and self-supervised learning. We believe that it is impractical for every sample collected by the fleet of Toyota vehicles to be labeled, especially due to the expensive labeling and growing compute costs. We envision that modern machine learning techniques will play an integral role in the programmatic generation and curation of autonomous driving training datasets that allow us to rapidly learn at scale from fleets of vehicles collecting experience.

**Why do we care?**

A unique advantage of Toyota as #1 automobile manufacturer is that we will soon have access to "Toyota-scale" sensory data that is order-of-magnitude larger than what our competitors have. This, however, imposes a unique challenge as well: manually labeling such data will quickly become infeasible. Therefore, developing new techniques in programmatically generating driving datasets to afford semi-supervised and weakly-supervised learning in a scalable manner will be critical for the success of Toyota's autonomous driving efforts.

## 8.50 BRIDGING PERCEPTION AND CONTROL WITH UNCERTAINTY MODELING

TRI Researcher:     Adrien Gaidon, Guy Rosman, Vitor Guizilini, Allan Raventos, Avinash Balachandran
Email Address:     Adrien.Gaidon@tri.global, Guy.Rosman@tri.global, Vitor.Guizilini@tri.global, Allan.Raventos@tri.global, Avinash.Balachandran@tri.global
TRI Thrust:       Driving

**What do you want to do?**

Our goal with this project is to make progress towards learning from demonstrations how to optimize end-to-end a typical modular automated driving stack following the "3 P decomposition": perception, prediction, and planning. In contrast to end-to-end deep sensorimotor policies (e.g., pixels to steering), a modular architecture is more interpretable, testable, and robust, but its implementation typically results in an engineered system that is hard to tune in a data-driven way. A particularly important challenge is making the whole architecture probabilistic, i.e. learning how to model and propagate uncertainty from the input

sensory information all the way to the output controls via passing through all intermediate modules (that alter the uncertainty in complex non-linear ways).

**Why do we care?**

Uncertainty estimation is a key aspect of decision-making. A model should be aware of its own limitations, be it due to the lack of training data, architectural shortcomings or unforeseen circumstances, and the ability to determine how much its own predictions should be trusted is highly valuable, as the basis for sensor fusion, planning and control. The introduction of uncertainty in deep learning, particularly in a fully differentiable way, is an open research problem. While classical methods (i.e. Gaussian Processes, Kalman Filters) have solid statistical foundations, large-scale learning via deep neural networks still lacks this capability, despite their overwhelming success in addressing a multitude of different tasks. The development of uncertainty-aware deep models would bridge the gap between perception and planning and control.

Furthermore, TRI has a suite of sensors capable of performing similar tasks (i.e. localization via mapping and perception, depth estimation via LiDAR and monocular imagery), and the introduction of uncertainty estimates would also allow principled fusion between different modalities, including redundancy and complementing each sensor's shortcomings. Additionally, such systems would enable improving the whole system with data, without having to optimize intermediate surrogate objectives that might be a poor proxy for the desired behavior of the system.

**Technical problem statement**

Our overarching goal here is to bridge the gap between perception and control with a learning algorithm that scales with data (in particular demonstrations). Handling constraints, modeling uncertainty end-to-end, causality, and generalization are key concerns that are hard to integrate in current methods.

## 8.51  INTERACTION-AWARE CONTROL FOR CARS THAT ANTICIPATE AND EXPLAIN COMPLEX ENVIRONMENTS

TRI Researcher:      Adrien Gaidon, Guy Rosman, Allan Raventos
Email Address:       Adrien.Gaidon@tri.global, Guy.Rosman@tri.global,
Allan.Raventos@tri.global
TRI Thrust:          Driving

**What do you want to do?**
Estimate, from large scale demonstrations, human-robot interaction models that can robustly predict and explain human actions and reactions to automated cars in challenging driving scenarios.

**Why do we care?**
An essential component for interaction with autonomous cars is learning a predictive model of human behavior through collected data. However, today's techniques usually learn human models that are trained on a small collected dataset hence are not guaranteed to be accurate in

all scenarios. Ideally, the decision-making capabilities of autonomous agents must be robust to potential flaws in the learned human model in order to guarantee successful objective completion.

**Technical problem statement**

Typically, the problem of modeling human driving behavior is addressed as an inverse reinforcement learning (IRL) problem. However, these techniques are developed based on the assumption of collecting expert trajectories and are prone to challenging scenarios not present in the training data. In this proposal, we first plan to study the robustness of these learned models, generate other structures and models that explain the novel scenarios, and then automatically generate test cases that find these novel scenarios. There are three key research directions:

- robustness analysis of learned human models;
- explainable modeling of human driving behavior in complex scenarios;
- generating scenarios with risky learned human models.

## 8.52 META-LEARNING AND INVERSE REINFORCEMENT LEARNING FROM LARGE SCALE DEMONSTRATIONS

TRI Researcher:        Adrien Gaidon, Dennis Park
Email Address:        Adrien.Gaidon@tri.global, Dennis.Park@tri.global
TRI Thrust:        Driving

**What do you want to try to do?**

Develop robust meta-learning algorithms (including meta inverse reinforcement learning) to leverage large scale driving demonstrations to quickly adapt models of a modular automated driving stack (perception, prediction, planning) to new environments or new driving tasks (e.g., changes in the ontology, new driving maneuvers).

**Why do we care?**

Learning good features or cost functions is key to the performance of ML models. This typically requires a lot of data for each task at hand. Meta-learning is a nascent set of techniques showing great promise towards learning general models that can be quickly adapted with little additional data, hence making adaptation cost effective.

**Technical problem statement**

Meta-learning for perception and prediction: learning to quickly learn to detect new object categories (few shots learning), predict new types of behaviors or specialize prediction models for finer-grained agent categories (i.e. conditioning on finer-grained semantics).
Meta-IRL for planning: learning general cost functions that can be quickly specialized for specialized maneuvers (e.g., zip-merges, unprotected left turns, etc).

## 8.53 TOWARDS LARGE SCALE EFFICIENT AND ROBUST MACHINE LEARNING SYSTEMS

TRI Researcher:        Chao Fang, Quincy Chen
Email Address:        Chao.Fang@tri.global, Quincy.Chen@tri.global

TRI Thrust:          Driving

**What do you want to do?**

Our goal of this project is to develop a large-scale ML training system which can leverage Peta size of data and scalable to thousands of GPUs. Typically model learning with large system trade off accuracy with speed. To this end, we propose two different directions to solve the problem: 1) design a good optimizer that could scale the model training without significantly losing accuracy, and 2) Optimized ML models that are more scalable to be trained with large training system. Desired deliverables of this project would include but are not limited to the following:

- New optimizers for large batch training.
- Auto learning-rate tuning.
- ML models optimized for large scale training.

**Why do we care?**

The amount of data to training data scale very fast for machine learning models particularly for Autonomous Vehicle applications. Training such models in a short time required a large ML system that could take advantage of the massive parallel computing.

Research in this domain provide a solid base for scaling up any machine learning models. It will fundamentally increase the capability of model to production for Toyota.

**Technical problem statement**

Neural networks have achieved great success in improving various computer vision tasks. However, these models are large in size, difficult to train, and slow in inference. Furthermore, it has been reported that these models are not robust to small perturbation – even though a near-perfect performance can be achieved on training data, it is often the case that the prediction will drop to chance level under small input perturbation. In this proposal, we aim to improve the training and inference speed of large machine learning models, and propose ways to make machine learning models more robust. We will apply the developed algorithms to several computer vision tasks, including image recognition, detection, 3D reconstruction and scene understanding models.

**1      Efficient Large-scale Training**

The past several years have seen tremendous growth in both the volume of data and the size of models. For example, a driving data-set in Toyota scale could soon reach multi-millions of images and hundreds of thousands of scenes. Progressed in self-supervised training also making large dataset more accessible to ML models. As a result, the long training time of Deep Neural Networks (DNNs) has become a bottleneck for Machine Learning (ML) researchers and developers. For example, it takes 29 hours to finish 90-epoch ImageNet/ResNet-50 training on eight P100 GPUs, and 81 hours to finish BERT pre-training on 16 v3 TPU chips. To speed up the training pipeline, we plan to investigate the following research topics:

- New optimizers for large batch training. Although companies are willing to speed up training using tens or hundreds of GPUs, the current SGD-based optimizers fail to fully utilize those computation resources since they usually converge to suboptimal solution under large batch size. Taking ResNet training on ImageNet as an example, when scaling the SGD batch size to 64K, the test accuracy of con-verged solution will drop from 76% to 66%. Thus a new optimizer is required for efficient training with multiple GPUs. Our previous works showed that a layer-wise step size is crucial for large batch train-ing. The proposed LARS algorithm can successfully train ResNet on ImageNet within 15 minutes [11], and the proposed LAMB algorithm (a variation of Adam) can train BERT in 76 minutes [12]. These algorithms have been used to achieve state-of-the-art models in large-scale. However, the performance of LARS and LAMB will still degrade when the batch size reaches certain level, so in the future we plan to develop better optimizers for large batch training. One idea is that the layer-wise LR for LARS and LABM are actually related to second order information, so in the future we will study how to formally utilize second order information for optimizing large-scale neural networks. Of course, full second or-der information will be too slow to use but we plan to exploit partial second order information, such as (layer-wise) Barzilai-Borwein learning rate or using sub-sampled Gauss-Newton or Fisher information matrix to approximate the second order information, in order to develop better optimization algorithms for large batch training.

  Furthermore, currently most of the large batch training experiments are conducted on ImageNet classification data. We will investigate large-batch training methods for other computer vision models, including detection models for images and videos. Furthermore, we will investigate large-batch training for model distillation and architecture search.

- Auto learning-rate tuning. In addition to the update rule, it has been observed that the performance of training highly depends on learning-rate scheduling. In addition to the standard decaying schedule, many other techniques have been designed including the warm-up scheduling and the cyclic scheduling [8]. In practice, a significant amount of engineering efforts have been spending on learning rate tuning. We propose to study a new family of automatic learning rate tuning procedure to remove human in the training loop. The main idea is to formulate learning rate scheduling as another learning problem and apply techniques such as reinforcement learning or Monte-Carlo Tree Search (MCTS) to automatically find a good scheduling that leads to superior generalization error.

## 2    Efficient inference and on-device computing

Currently, many complex and deep models cannot be deployed on real world systems due to their large model size and slow inference speed. Our second goal is to resolve these problems. More specifically, we consider the following directions:

- Fast inference with large output space. For models with large output space, such as multi-class prediction with tens of thousands of labels or other structural prediction tasks, the bottleneck of inference time is usually on the final layer of neural network. This mainly involves the operation of Maximum Inner Product Search (MIPS) – given a query vector

v, find the vector from a large database that has maximum inner product with v. To deal with this, our group has proposed state-of-the-art algorithms [13]. To further improve the speed of this part, we propose two novel future directions. First, notice that previous MIPS algorithms are focusing on improving the worst-case performance under any given query, while in practice queries are usually semantic features extracted from input images instead of arbitrary vectors. Based on this insight, we will develop faster algorithms by exploiting data distribution. For instance, our recent work [1] has shown that in NLP applications, the semantic features often follow clustering structure, and by exploiting this structure our algorithm achieves significant speedup over existing methods. We will thus study how to extend this to various of computer vision tasks.

- Architecture search for efficient structure. We also plan to develop efficient Neural Architecture Search (NAS) algorithms and apply them to get better network structure with smaller size and faster inference speed. Existing algorithms such as DARTS [5] constructs a differentiable search space and then optimizes it by gradient descent. However, DARTS is still slow as it updates an ensemble of all operations and keeps only one after convergence. Besides, DARTS can converge to inferior architectures due to the strong correlation among operations. We will investigate how to make NAS more efficient and convergence to better solution. To improve the efficiency of NAS, we propose to sparsify the operations such that for each edge in the computation graph, only one or few operations need to be computed in forward and backward computations. This will require novel formulations with sparse regularization and new solvers for solving this. Furthermore, to make NAS converges to a better architecture instead of stuck at local minimums, we propose a new framework to combine Bayesian optimization with existing iterative solvers (such as DARTS) to get improved solution without much more effort.

## 3      Improving the robustness of machine learning models

Although neural networks have achieved remarkable performance, it has been shown recently that a small adversarial or non-adversarial perturbation can easily lead to significantly degraded performance of state-of-the-art models. This finding indicates lacking of robustness of machine learning models and creates safety concerns in many real world applications, such as aircraft control systems and self-driving cars, and leads to the following questions: How to characterize the robustness of machine learning models? and how to make them more robust

To answer the first question, we plan to develop algorithms to formally verify the robustness of machine learning models. This has been identified as an important task for safety-critical systems, such as aircraft control systems [4, 3]. Mathematically, a verification algorithm aims to provably characterize the prediction function of a network within some specified region of the input space. For example, within a prescribed region of the input one may wish to provide simple (e.g., affine) upper and lower bounds of the network output. However, due to the complicated interactions of the nonlinearities in deep networks it is often NP-hard to compute these piecewise regions and corresponding linear/affine bounds exactly. To tackle this problem, several recent algorithms have been proposed [2, 7, 10] (including some seminal works from our group [9, 14]), obtain an approximate solution by relaxing the non-linearities in the network. Unfortunately, the bounds computed by current approaches are often not tight enough for real applications, and furthermore, current methods are exclusively focusing on feedforward ReLU

networks and have difficulties extending to more general network structures such as residual links and attention layers. We thus propose to tackle these problems to enable neural network verification for general network structures and perturbation models.

For the second question, making neural networks more robust has become one of the most important research topics in machine learning. Currently there are two successful ways for defense: adversarial training and randomization. Adversarial training aims to minimize the robust error of neural network instead of clean error, where the robust error can be computed by either adversarial attack or verification. Randomization improves the robustness through smoothing out the prediction in a neighborhood. Some of our previous works have explored both of these techniques [6, 15], and we will develop several algorithms to improve over existing approaches, including stratified adversarial training (learn better weighting for each sample in adversarial training), random smoothing with learned optimal distribution, and designing robust architectures. Furthermore, we will investigate how to distill the knowledge learned by non-robust networks into robust structures (such as decision trees) which will improve both robustness and inference speed.

References

1. P. H. Chen, S. Si, S. Kumar, Y. Li, and C.-J. Hsieh. Learning to screen for fast softmax inference on large vocabulary neural networks. In ICLR, 2019.
2. T. Gehr, M. Mirman, D. Drachsler-Cohen, P. Tsankov, S. Chaudhuri, and M. Vechev. Ai2: Safety and robustness certification of neural networks with abstract interpretation. In 2018 IEEE Symposium on Security and Privacy (SP), pages 3–18. IEEE, 2018.
3. K. D. Julian, S. Sharma, J.-B. Jeannin, and M. J. Kochenderfer. Verifying aircraft collision avoidance neural networks through linear approximations of safe regions. arXiv preprint arXiv:1903.00762, 2019.
4. 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, pages 97–117. Springer, 2017.
5. H. Liu, K. Simonyan, and Y. Yang. Darts: Differentiable architecture search. arXiv preprint arXiv:1806.09055, 2018.
6. X. Liu, M. Cheng, H. Zhang, and C.-J. Hsieh. Towards robust neural networks via random self-ensemble. In Proceedings of the European Conference on Computer Vision (ECCV), pages 369–385, 2018.
7. G. Singh, T. Gehr, M. Mirman, M. Puschel, and M. Vechev. Fast and effective robustness certification. In Advances in Neural Information Processing Systems, pages 10802–10813, 2018.
8. L. N. Smith. Cyclical learning rates for training neural networks. In 2017 IEEE Winter Conference on Applications of Computer Vision (WACV), pages 464–472. IEEE, 2017.
9. T.-W. Weng, H. Zhang, H. Chen, Z. Song, C.-J. Hsieh, L. Daniel, D. Boning, and I. Dhillon. Towards fast computation of certified robustness for relu networks. In International Conference on Machine Learning, pages 5273–5282, 2018.
10. E. Wong and Z. Kolter. Provable defenses against adversarial examples via the convex outer adversarial polytope. In International Conference on Machine Learning, pages 5283–5292, 2018.

11. Y. You, Z. Zhang, C.-J. Hsieh, J. Demmel, and K. Keutzer. Imagenet training in minutes. In Proceedings of the 47th International Conference on Parallel Processing, page 1. ACM, 2018.
12. Y. You, J. Li, S. Reddi, J. Hseu, S. Kumar, S. Bhojanapalli, X. Song, J. Demmel, and C.-J. Hsieh. Large batch optimization for deep learning: Training bert in 76 minutes. arXiv:1904.00962, 2019.
13. H.-F. Yu, C.-J. Hsieh, Q. Lei, and I. S. Dhillon. A greedy approach for budgeted maximum inner product search. In Advances in Neural Information Processing Systems, pages 5453–5462, 2017.
14. H. Zhang, T.-W. Weng, P.-Y. Chen, C.-J. Hsieh, and L. Daniel. Efficient neural network robustness certification with general activation functions. In Advances in neural information processing systems, pages 4939–4948, 2018.
15. H. Zhang, H. Chen, C. Xiao, B. Li, D. Boning, and C.-J. Hsieh. Towards stable and efficient training of verifiably robust neural networks. arXiv preprint arXiv:1906.06316, 2019.

# 9 TOPICS: ROBOTICS

## 9.1 DISTRIBUTIONAL ROBUSTNESS FOR OPEN-WORLD MANIPULATION

TRI Researcher:  Russ Tedrake, Hongkai Dai, Andres Valenzuela
Email Address:  Russ.Tedrake@tri.global, Hongkai.Dai@tri.global, Andres.Valenzuela@tri.global
TRI Thrust:  Robotics

TRI has been developing advanced manipulation capabilities for tasks requiring complex perception, planning, and control. One example of this that we have made public is a robot loading the dishwasher. We are interested in proposals that can address the fundamentally hard problems in making these systems robust. Let us say that we have run some number of experiments in the lab and we have run many more experiments in simulation… what can we say about the expected performance if we go to deploy the system in a new environment? How do we provide test coverage for every possible kitchen?

"Trustworthy AI" and "verification of machine learning components" are hot topics in research today. We are seeking proposals that address the particular challenges that become central in manipulation, and also the specific opportunities are enabled by the technology already under-development at TRI.

Some of the challenges central to manipulation include:
- Rich contact interactions with the environment. Unlike, e.g. the challenge of modeling other drivers for verifying autonomous driving, the rules of the game here are relatively known, but they are complicated (non-smooth mechanics, rich uncertainty models).
- Complex building blocks. Advanced manipulation systems today have deep networks for perception, sophisticated task-level planners (discrete) as well as sample-based or optimization-based motion planners (continuous), and low-level feedback controllers.
- Distributions over environments. How do we produce meaningful distributions? How accurate do they need to be? Is generalization theory and/or naive domain randomization sufficient to instill confidence? What are the critical sources of

randomness/variability that must be included to establish robust performance in reality?

There are also many opportunities. TRI has developed a mature manipulation toolchain and has made many of the components available as open-source code. We have a state-of-the-art dynamics engine that is capable of accurately simulating contact-rich interaction, and it is written in a framework that was built to support rigorous design and analysis (all state and all random variables are declared specifically; most components support autodiff and even symbolic analysis; our mathematical programming interface already supports interfacing with advanced optimization and SMT solvers). We also have infrastructure for extensive hardware testing.

We are particularly interested in proposals that can leverage and/or extend these open-source tools. In addition to novel research ideas and publications, we encourage research that can be published as open-source code; we believe that this helps to build the community and also improves TRI's ability to leverage your ideas and results.

## 9.2   REAL-WORLD MANIPULATION

| | |
|---|---|
| TRI Researcher: | Calder Phillips-Grafflin, Naveen Kuppuswamy, Russ Tedrake |
| Email Address: | Calder.Phillips-Grafflin@tri.global, naveen.kuppuswamy@tri.global, Russ.Tedrake@tri.global |
| TRI Thrust: | Robotics |

The robotic systems developed at TRI incorporate state-of-the-art techniques in perception, planning, sensing, and simulation to perform limited domestic tasks with some level of robustness; our robot loading the dishwasher is one example. In each of these areas, applying our systems to new real-world tasks will require capabilities that are either insufficiently robust, underexplored, or as-yet unavailable.

TRI would like to support research that tackles some of the most pressing challenges in real-world manipulation for the home. Topics of interest include:

- Manipulating deformable items, such as clothing, food, and other household objects

- Strategies that address high item variety and variance, where objects to be manipulated are not only too varied to be modelled a priori, but also variance between items of a given class or type is high enough to thwart existing approaches (ex. scan + model cleanup)

- Highly cluttered environments, where clutter must be perceived and manipulated to accomplish the intended task

- Uncertain or never-seen-before environments, where an existing model of the environment is unavailable or useless and the robot must develop its own model (both metric and semantic) of the environment and reason over its uncertain or unknown surroundings

- Highly dexterous behaviors, such as in-hand manipulation or manipulating objects with complex dynamics

- Highly complex tasks, where task complexity and the need for robustness means that behaviors cannot be authored or taught by expert users but must be planned or synthesized automatically and must handle unexpected behavior during execution

- Fault-tolerance and robustness, which requires that the limitations and corner cases of complex manipulation systems can be discovered automatically and improved

- Fleet-wide lifelong learning, where multiple robots can learn from each others' mistakes and grow in capability and robustness over time

- Safe and efficient physical interaction with humans or cooperative human-robot tasks

TRI is interested in proposals for new collaborations in which we can collaborate closely, sharing code and data and domain experience -- TRI's manipulation effort can provide code maturity and a scale of hardware and software testing that is hard to match at a university. We would like to find ways to test your ideas in our hardware benchmarks/experiments.

Specifically, the tasks we seek to handle involve the following challenges:

- Deformable items, such as clothing, food, and other household objects

- Difficult-to-perceive items (such as glassware or liquids) and environments (highly reflective surfaces, variable lighting, or outdoors)

- High item variety and variance, where objects to be manipulated are not only too varied to be modelled a priori, but also variance between items of a given class or type is high enough to thwart existing approaches (ex. scan + model cleanup)

- Highly cluttered environments, where clutter must be perceived and manipulated to accomplish the intended task

- Uncertain or never-seen-before environments, where an existing model of the environment is unavailable or useless and the robot must develop its own model of the environment and reason over its uncertain or unknown surroundings

- Highly dexterous behaviors, such as in-hand manipulation or manipulating objects with complex dynamics

- Highly complex tasks, where task complexity and the need for robustness means that behaviors cannot be authored or taught by expert users but must be planned or synthesized automatically and must handle unexpected behavior during execution

- Fault-tolerance and robustness, which requires that the limitations and corner cases of complex manipulation systems can be discovered automatically and improved

- Fleet-wide lifelong learning, where multiple robots can learn from each other's mistakes and grow in capability and robustness over time

- Safe and efficient physical interaction with humans or cooperative human-robot tasks

## 9.3 ADVANCED PERCEPTION FOR HOME ROBOTS

TRI Researcher: Duy-Nguyen Ta, Kunimatsu Hashimoto, Siyuan Feng
Email Address: Duy@tri.global, Kunimatsu.Hashimoto@tri.global, Siyuan.Feng@tri.global
TRI Thrust: Robotics

TRI is interested in developing advanced manipulation capabilities for tasks requiring complex perception, planning, and control; our [robot loading the dishwasher](#) is one example. Despite recent progress in object-based perception for manipulation, both in known-model and unknown-model (category-level) cases, current state-of-the-art robotic systems fall short of dealing with the complexity of perception in the home environment. TRI is interested in advanced research on scalable approaches to robot perception that can support robot manipulation in the home.

What information is required from a perception system in order to support manipulation in a home environment? Tasks like loading the dishwasher require enough semantic information to understand which rack to use in the washer, enough geometric information to know how to grasp the object, avoid collisions with the sink, and place it in the rack, but also enough understanding of the object to get the orientation in the rack correct. How can we efficiently acquire these types of understanding for novel objects (perhaps with very minimal supervision from a human)? Can we develop strategies that work for difficult-to-perceive items (such as glassware or liquids) and environments (highly reflective surfaces, variable lighting, or outdoors)? Can we help robots understand the hierarchical semantics of object parts to enable more intelligent planning? Can we transfer structured knowledge from one class of objects to a novel class, instead of retraining from scratch? Can the perception system reliably communicate its confidence or uncertainty about its current predictions in a format that can be consumed by downstream components in the manipulation stack? Can we reliably track and predict the future of things that dynamically change over time under different control policies (such as deformable objects) to find the optimal policy that achieves the goal? How do we efficiently test / establish confidence in the system to the level of maturity that we would be willing to deploy it in a consumer's home?

TRI is interested in proposals for new collaborations with academic partners who can explore fundamentally new approaches to this difficult problem. We are interested in models where we can collaborate closely, sharing code and data and domain experience -- TRI's manipulation effort can provide code maturity and a scale of hardware and software testing that is hard to match at a university. We would like to find ways to test your ideas in our hardware benchmarks/experiments.

## 9.4 DEXTEROUS ROBOT HANDS

TRI Researcher:       Alex Alspach, Naveen Kuppuswamy, Avinash Uttamchandani
Email Address:        Alex.Alspach@tri.global, Naveen.Kuppuswamy@tri.global,
                      Avinash.Uttamchandani@tri.global
TRI Thrust:           Robotics

At TRI, manipulation researchers have so far mostly stuck with parallel, off-the-shelf industrial grippers due to their reliability and toughness. We've been working around their limitations thus far, but as tasks and grasps become more complex, we realize the potential of having more sensing and dexterity at the end effector. We welcome proposals for research and development of a robust and dexterous robot hand with integrated high-fidelity sensing, designed for real-world, everyday tasks.

Form and function - A parallel jaw gripper can be cleverly employed to perform dexterous tasks like rooting through a cluttered sink. Most researchers here at TRI have performed the entertaining experiment of wielding a parallel gripper for a day to see what it takes to complete everyday tasks with our robot's limited end effector. One can get quite a bit done, but tasks take longer and stable grasps take effort.

With dexterity a combination of hardware and software capabilities, we imagine a cohesive codesign of gripper topology, sensors and control, focusing on reliability/uptime, robust ability to manipulate human-made objects, features and tools (e.g. spatula, mug handle, scissors), ability to fit into tight spaces (e.g. cluttered sink, back of a well-stocked cabinet), dexterity (e.g. hand positioning via wrist, in-hand manipulation), and demonstrated capability for in-home manipulation tasks.

Capabilities of interest include and are not limited to:

- Robust manipulation of household objects in household situations
- Motion/force transparent actuation
- Torque/force controllability
- Integrated tactile sensing
- Normal and shear force estimation
- Vibration and slip sensing
- High-fidelity proprioception
- In-hand pose estimation
- Passive and active compliance and jamming
- In-hand manipulation of smaller objects
- Tying and other precision manipulation tasks
- 5000+ hour uptime
- Reliability and repairability
- 10+kg payload
- Ability to use human tools and spaces
- Actuated, flexible wrists
- A standard and extensible data bus (e.g. Ethernet, Ethercat) that allows for the hand to be used with a variety of robot arms.
- Modularity between end effectors/fingers to allow for different configurations for different tasks

## 9.5 ADVANCED SIMULATION CAPABILITIES FOR HOME ROBOTS

TRI Researcher: Russ Tedrake, Michael Sherman, Alejandro Castro
Email Address: Russ.Tedrake@tri.global, Sherm@tri.global, Alejandro.Castro@tri.global
TRI Thrust: Robotics

TRI has been developing advanced manipulation capabilities for tasks requiring complex perception, planning, and control. One example of this that we have made public is a robot loading the dishwasher. The development of these results has made extensive use of simulation; we have invested heavily in robust simulation of the complex mechanics of rigid and nearly rigid contact (c.f. [1]), on simulated perception for training our computer vision components, and on advanced algorithms using simulations to quantify and improve robustness. But there are many aspects of the problem of robot manipulation for the home that we cannot yet capture in simulation.

We welcome proposals for improving the scope of simulation-based design and analysis for home robots. Relevant ideas include:

- Content generation. Even for (nearly) rigid objects, how do we capture the diversity of objects/scenes that we might find in the home? What fidelity do these models need to have (public datasets typically do not include sufficient information to use models for both physical simulation and rendering)? Often we use artists to "touch up" our art assets -- is there a way to automate this pipeline to dramatically improve its throughput? Must we solve the inverse rendering problem, or do you have specific proposals for avoiding it? Can a robot equipped for manipulation tasks in the home autonomously generate new art assets for simulation (e.g. without bespoke scanning equipment)? Do we need accurate material properties (e.g., BRDFs and friction coefficients), or can we overcome this with sufficient domain randomization?
- Simulating non-rigid objects in the home, including soft, potentially thin deformable objects, (cloth/laundry, lettuce leaves, ..), liquids, noodles, plush toys, etc. These models need not run at real-time to have value, but there is significant value in fast/approximate models which are sufficiently faithful to physics for robot software design and are also robust enough to withstand rigorous testing (e.g. with Monte-Carlo methods).
- Simulation models that are more amenable to optimization / analysis. Examples include a differentiable renderer that can capture the fidelity which we now better understand is required for sim2real transfer of perception, or even a symbolic rendering pipeline for formal analysis.

We are specifically interested in integrating your best ideas for advanced simulation capabilities into our open-source dynamics engine, Drake (http://drake.mit.edu), and in building the open-source community around Drake with expert users.

[1] Ryan Elandt, Evan Drumwright, Michael Sherman, Andy Ruina. A pressure field model for fast, robust approximation of net contact force and moment between nominally rigid objects. https://arxiv.org/abs/1904.11433

## 9.6 COMBINING MECHANICS AND CONTROL WITH MACHINE LEARNING FOR MOBILE MANIPULATION

TRI Researcher:      Krishna Shankar
Email Address:      Krishna.Shankar@tri.global
TRI Thrust:      Robotics

**Technical Problem Statement**
How can we leverage machine-learning to combine classical techniques in mobile-manipulation to perform complex tasks in real environments?

**Why do we care**
We have a thorough understanding of mobile-manipulation in a classical sense. On the other hand, there is a small but growing field of work on using machine-learning for one-off mobile-manipulation tasks end-to-end tasks. This newer work shows promise, but lacks the rigor, thoroughness or generality of prior work. We are looking for university partners to investigate the 'glue' between classical theory and modern learning approaches, with the belief that it is development in this interface that will allow us to build safe, useful and reliable mobile-manipulation systems capable of helping us do everyday tasks in human environments.

**What do we want to do**
Combine learning with existing theory in the following areas:
- Grasp/Cage Generation
- Optimal task/trajectory planning for complex, contact rich manipulation
- Whole-body task sequencing and planning
- Task Planning with guarantees
- Perception-based manipulation and control
- Safe, efficient reinforcement-learning and adaptive control
- Mechanism/Actuator design to enable novel mobile manipulation capabilities

## 9.7  SPEECH, LANGUAGE, AND DIALOG WITH ROBOTS

TRI Researcher:        Thomas Kollar
Email Address:         Thomas.Kolar@tri.global
TRI Thrust:            Robotics


The fleet learning team aims to significantly advance the state of the art in mobile manipulation technology by demonstrating the technical feasibility of a general purpose mobile manipulation robot that dramatically improves the quality of life when performing household tasks. For robots to be able to work collaboratively with humans, they must be able to be controlled and commanded in a flexible and natural way and be able to respond to events in their environment.

There are two primary goals of the fleet learning team that enable flexible and natural human-robot interaction. The first goal is to develop the capability to interact with the robot via speech. A second goal is to develop the capability to detect events in audio. In addition, we would like robots to be able to learn how to perform tasks given natural language instructions, to respond to physically grounded questions, and to be able to generate natural language descriptions of their state. In support of these goals, we are looking for University 2.0 partners in the following areas:
- Grounded speech and language understanding, such as connecting natural language commands and questions to complex task graphs, environment states, robot actions, and robot plans.

- Dialog systems for physically grounded interactions, including anaphora resolution, semantic parsing, named entity resolution and intent recognition. End-to-end learning for grounded dialog systems.
- Multi-modal teaching of objects, attributes, relations and tasks via speech interaction, natural language, gestures and other modalities.
- Visual question answering - Enabling robots to answer questions about their current state or the state of the environment. In addition, enabling a robot to answer longitudinal questions about their environment (e.g., "Where are my keys?", "Who came into the house today?").
- Enabling robots to collaboratively work with people in a household to perform a task, such as when executing a recipe.
- Speech and audio event recognition in adverse environments, such as when fan noise, joint noise, robot motion, background speech are present.
- Speaker recognition as a part of robotic tasks. For example, a task could involve only working with a specific person to achieve a goal.
- Enabling robots to ask for input or help when they are uncertain of what action to perform. For example, a robot might confirm which object to manipulate when there are multiple objects of the same type present, or there are other environmental ambiguities.
- Ontologies and representations for grounded natural language commands; semantic parsing of these representations.
- The ability to describe robot and environment states with natural language.
- End-to-end models that couple speech, natural language processing, dialog, and vision.

## 9.8 LEVERAGING SIMULATION TO LEARN MANIPULATION BEHAVIORS IN HOME ENVIRONMENTS

Toyota Research Institute (TRI) is focused on enabling mobile robots to provide in assistive care for the elderly. In order to do this, a robot must be capable of performing robust manipulation behaviors in unstructured home environments. One promising avenue to achieving this is through the use of simulation, in which a robot can be exposed to a large number of scenarios to train on.

Recent results suggest that simulation can lead to quite robust manipulation performance in unseen settings. However, more work needs to be done to make simulation useful in the home environment. This research call is for addressing three challenges currently facing learning from simulation: State Representation, Simulation Quality and Error Detection.

**Challenge 1: State Representation**
Traditionally, most robotic manipulation behaviors operate on a 3D representation of the world, such as a depth images, point cloud or voxel map. While it has been shown that one can learn transferable policies using these representation in simulations, home environments present a novel challenge since everyday objects are reflective, shiny or transparent. Thus, representations to transfer learned perception form simulation, which require dense 3D information, may no longer be viable. We are interested in research to focus on new ways to transfer manipulation from simulation that require only spares 3D information.

**Challenge 2: Simulation Quality**
Given a reasonable state representation it's still unclear how accurate a simulator needs to be for learned manipulation behaviors to transfer. Since, there is an inherent trade-off between between having high quality simulation and the level of data diversity that can be generated. We are interested in better understanding and improving how well techniques like noise injection, or domain randomization, can enable robustness with low quality simulation.

**Challenge 3: Error Detection**
When learned behaviors are deployed in home setting it's important to understand when they are likely to fail, so that damage to the environment can be prevented. If a behavior is a learned policy it can be difficult to understand when the estimated parameters are invalid. New algorithms for efficiently detecting when the robot encounters an unknown state of the world and communicating this uncertainty to a human operator is of high interest to TRI

## 9.9    TOWARDS REAL-TIME ROBOT AUTONOMY WITH HARDWARE-ACCELERATED PARALLEL OPTIMIZATION

TRI Researcher:        Huihua Zhao
Email Address:        Huihua.Zhao@tri.global
TRI Thrust:        Robotics

Motivated by the increasing needs in autonomous robots, the emerging research focus is to significantly advance the computational performance of real-time motion planning of complex robotic systems, such as home service robots and autonomous driving. Specifically, this project will seek to develop massively parallel optimization algorithms and computing hardware to tackle the current computational bottleneck in the real-time trajectory optimization of high degree-of-freedom and highly dynamic robots, particularly those whose dynamics have dominant effects on planning responsive motions.

Dynamics based motion planning methods can realize more agile behaviors, but the planning process requires an excessive amount of time and may not be able to converge reliably, and therefore, are only suitable for off-line planning. This has been a significant impediment toward applying these approaches to achieve dynamic maneuvers under complex environments or emergent scenarios. In a broader view, this is where the full-order dynamics-based approaches have had to pay the piper for its admissibility of highly dynamic behaviors. New approaches and platforms are needed to be explored and developed to solve these challenges in complex problems such as highly articulated legged and home service robotics. Despite recent achievements in the parallel computation, the motion planning problems are still too complex to fit on an efficient embedded platform to achieve real-time performance due to the large size and irregularities of the parallel threads in these optimization problems.

To facilitate hardware parallelization, a rigorous mathematical framework of parallel computing structure for robot trajectory optimization problems will be formulated. By exploiting the structures of the optimization algorithms and the patterns of the numerical computation, we will create a new motion planning processor (MPP) architecture for this class of problems. We will design MPP hardware prototypes to demonstrate higher performance and better efficiency than what is available with commodity hardware. An MPP compiler will be

developed in conjunction with the MPP hardware to enable the seamless translation from a motion planning algorithm to an MPP hardware. The MPP hardware prototype will be capable of performing real-time dynamic motion planning for complex robotics systems, such as home service robots, and speeding up the whole-body dynamics-based motion planning by two orders of magnitude compared to the current state-of-the-art approaches, while taking a fraction of the size, weight, and power. Moreover, the proposed optimization algorithms and computing architectures can be broadly applicable to other robotic platforms.

# 10  TOPICS:  MACHINE ASSISTED COGNITION (MAC)

TRI Researcher:      Fran Bell
Email Address:       [Fran.Bell@tri.global](mailto:Fran.Bell@tri.global)
TRI Thrust:          MAC

TRI recently created a new research thrust on Machine Assisted Cognition (MAC), adding to the research thrusts on Automated Driving, Home Robotics, and Accelerated Materials Design and Discovery.  TRI is in the process of hiring a group of researchers and establishing a new research agenda.  The specific goals for the year ahead will not be settled until the leadership of the group has been hired.

The goal of the MAC program is to develop and demonstrate computational aids to amplify and augment human cognitive abilities.  By computational aids we mean AI-powered software systems that directly assist users with tasks requiring careful, deliberate thought.  Cognitive abilities of interest include prediction, judgment, and decision making and exclude memory and search.  By "amplify and augment" we mean not to imitate or replace human abilities; instead, the goal is to leverage and extend them.

Technical objectives of the MAC program include the following:
1. Increase quality of predictions, judgments, and decisions by (a) Neutralizing biases, (b) Ensuring fitness, and (c) Ensuring ethics
2. Speed up predictions, judgments, and decisions
3. Scale up to more complex cognitive tasks and larger groups while countering the tendency of groups to be risk averse

The MAC approach is interdisciplinary, including both analytical disciplines (computer science, statistics, data science) and behavioral disciplines (psychology, cognitive science, behavioral economics).

To achieve Technical Objective 1, the MAC program seeks to develop new algorithms for human-aware AI, and to demonstrate AI systems that augment human mental capabilities. Topics of interest include the following:
- Model and predict human behavior
- Detect when the user is vulnerable to a cognitive bias, and "nudge" them toward a better outcome
- Assist the user in recalling past experiences applicable to the user's current situation
- Model how people make tradeoffs

- Model and predict consumer behavior; for example, in purchasing cars versus trucks or in attitudes toward electric vehicles
- Model how people perform selected cognitive tasks; for example, learning generalized value functions
- Assess user fitness for making predictions, judgments, and decisions; for example, infer user state through physiological measurements
- Given a low state of fitness for optimal performance on cognitive tasks, provide feedback to the user on how to improve fitness
- Reason about the consistency of predictions, judgments, or decisions with ethical principles and core values, and provide feedback to the user in cases of inferred ethical violations
- Achieve effective interactions between humans and AI systems by ensuring that the AI system design suffers from neither under-trust nor over-trust
    - Identify, quantify and mitigate biases in AI systems to develop trust
    - Develop systems that continuously evolve with evolving environments
    - Demonstrate consistent behavior across a range of augmentation tasks to develop predictability

For Technical Objective 2, the MAC program strives to introduce AI representations and methods into prediction, judgment, and decision-making. Topics of interest include the following:
- Model-based design that uses data science to accelerate search through a parameter space, or that uses learned models of the quality of the outputs of a simulation
- Accelerate design using generative models that both expand the palette of concepts and save time compared to human generation of concepts
- Accelerate design using predictive models that eliminate unpromising design concepts early in the process
- Retrieve information when the user needs it, even when they have not prompted the system explicitly for that information
- Increase human-machine collaboration speed by developing and surfacing human understandable reasoning

For Technical Objective 3, the MAC program aims to enhance problem-solving within mixed teams of people and AI systems. This work will enable a combination of humans and AI systems to work together to achieve application goals, where the collaborative interaction takes advantage of the complementary nature of humans and AI systems. An additional research goal is to explore issues associated with scaling up the size of teams comprising multiple machines and multiple humans.

For all the technical objectives, proposers should feel free to nominate and explore other goals during the open discussion period (January 2020) for this solicitation.