

Interim report: Optimizing traffic flow using learning and the shout-ahead agent architecture

Christian Roatis
Department of Computer Science
University of Calgary
Calgary, Canada
christian.roatis@ucalgary.ca

Abstract—We present an attempt to improve traffic throughput in the traffic simulator SUMO by implementing the shout-ahead agent architecture and a hybrid behavior learning method for it into traffic signal controller agents. Problems with current traffic signal controllers, namely their inefficient, rigid phase changing algorithms, served as motivation for the project. Our proposed solution to this problem involves implementing a shout-ahead agent architecture and a hybrid learning behavior method into traffic signal controller agents to better optimize their performance. The performance of the shout-ahead enabled traffic signal control agents will be determined over a set of experimental evaluations, using simulation time and traffic throughput as performance measures, and will compare the shout-ahead traffic signal controller's performance relative to that of the standard signal control algorithm utilized by SUMO to determine the effectiveness of our implementation.

Keywords—traffic throughput, cooperative systems, optimization, shout ahead

I. INTRODUCTION/MOTIVATION

For modern suburb dwellers in large metropolitan cities across the globe, the drive to and from work, often featuring long commute times, has become a part of life [1]. A 2016 Statistics Canada report found that 7% of all car commuters in Canada spend at least 60 minutes travelling to work, with that number jumping to 10% for car commuters in metropolitan areas [2]. Long commute times are generally a result of traffic jams caused by bottlenecks, rather than long distances, and traffic signal timing at intersections is a major factor in creating said bottlenecks [3]. Traffic signal timing can also hinder traffic flow at times without congestion, such as at night, where for instance vehicles may be held up at an intersection that is giving way to a direction with no oncoming vehicles.

Problems such as these string from traffic signal algorithms that do not effectively take into account traffic volume and flow, and lack the flexibility to adapt to the ever-changing traffic state. Specifically, traffic signal algorithms rely on pre-set timer cycles that can be adjusted based on demand, but do so in predetermined ways that do not optimally accommodate for the current state of an intersection [4]. These pre-set timing cycles are also why vehicles may be held up at an empty intersection in the middle of the night. Adjusting traffic signal behavior based on current traffic states could theoretically improve traffic flow, by jettisoning through heavier flows for longer during rush hour, and responding to demand immediately at empty intersections,

to name a just a few examples. Theoretically, traffic flows could be optimized further if traffic signal algorithms were able to account for expected traffic flows in the immediate future and adjust their behaviors accordingly. In this paper, we propose a solution embodying such characteristics by implementing the shout-ahead architecture and hybrid behavior learning method for it into traffic signal controller agents.

II. BASIC DEFINITIONS

In this section, we provide a high-level definition of an agent that we will apply the shout-ahead architecture to in the next section, and a description of the Sarsa algorithm for reinforcement learning (see [26]) which the reinforcement learning component of the hybrid learning method will be based on.

As is described in [5], an agent Ag is represented as a quadruple $Ag = (Sit, Act, Dat, f_{Ag})$, where Sit is the set of situations Ag can distinguish, Act is the set of actions Ag can perform, Dat is the set of possible value combinations of Ag 's internal data areas (all internal "states" Ag can have) and $f_{Ag} : Sit \times Dat \rightarrow Act$ is Ag 's decision function, taking in the current situation and value combination of the internal data areas, and deciding on the next action to take based on them.

Reinforcement learning is a major component for the learning of cooperative behavior of agents. We define behavior as it is in [5]; that is, as a sequence of actions, meaning over time, reinforcement methods will attempt to create sequences of actions that are optimal. To achieve this, they reward actions that yield optimal outcomes, eventually forming whole behaviors composed of the most rewarded actions. Agents employing reinforcement learning generally base their decisions on a situation-action matrix contained within a data area in Dat . Each combination pair of situations and action contain a Q-value ("quality"-value) which is updated by rewards allocated to the agent based on its performance in the environment. This value represents the "quality" of any one combination pair relative to its cohorts, with higher Q-values reflecting successful combination pairs. The decision function f_{Ag} of an agent Ag generally selects the action with the highest Q-value for a given situation, but random factors may be involved in the decision process as well to facilitate the exploration of possibilities. There are numerous ways to update an evaluation $Q(s, a)$ of a situation-action pair (s, a) ($s \in Sit, a$

$\in Act$). The method we chose for this purpose is based on the same variant of the temporal difference learning method called Sarsa (see [26]) used in [5].

Specifically, if at time t an agent is in a situation s and takes the action a , leading to the situation s' at time $t + 1$, a reward r and if in situation $s' f_{Ag}$ selects $a' \in Act$ as the action to take, we would update $Q(s, a)$ as follows:

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)]$$

where γ , $0 \leq \gamma \leq 1$, is the discount rate and α , $0 \leq \alpha \leq 1$, is the learning factor. Both the discount rate γ and learning factor α are set by the user of the method, where γ signifies the emphasis on the importance of future evaluations, while α dictates the rate at which the value of $Q(s, a)$ converges against the optimal evaluation value of the situation-action pair.

Meanwhile, evolutionary approaches for the learning of cooperative behaviors differ from reinforcement methods in that instead of individual step-wise refinement, they create full behaviors that are evaluated as a whole. As such, the evolutionary process creates new behaviors out of old ones, with the hope better whole behaviors are created among the children. This process is as follows: the evolutionary methods work on sets of “individuals”, which represent one possible solution to a problem, and evaluate them using a fitness function *fit*. These sets are called “generations”, and the next generation set is created by copying a number of the best individuals from the previous generation into the new one, and filling the remaining space with newly generated individuals. Said new individuals are created by selecting a set of parents from the previous generation, and combining their most successful attributes (as determined by some fitness measure), often with some random new pieces as well, to construct a whole new individual. This process is done by “genetic operators”. After the creation of a new generation, all individuals are evaluated by *fit* again, and a number of the weakest rated individuals are culled off. Following this, the cycle continues. The first generation is generally created randomly, and the process stops after a certain number of iterations, or when a predetermined milestone is reached.

III. PROPOSED SOLUTION TO THE PROBLEM

In this section, we present my proposed solution to the traffic signal controller problem via the rule-based shout-ahead agent architecture and accompanied hybrid behavior learning method detailed in [5]. In the following, every set of individual traffic signals in an intersection is controlled by an agent, which is responsible for changing their states. As discussed previously, the current publicly employed algorithms for these agents are relatively static and inefficient, running predetermined cycles with little room for adjustment based on traffic volumes and flows [4]. The Shout-ahead agent architecture allows the use of communicated intentions of other agents, which in concert with information about an agent’s current state and environment, allows the agent the opportunity to take a more informed action at any given state than it would otherwise. The hybrid behavior learning method for the Shout-ahead architecture then works to

refine the behavior of the intersection controller agent, aiming to improve the agent’s decision making such that the action being taken at any one time step, is the optimal action at said time step, or very near to it. The application of this approach was shown to be effective when applied to units in a computer game, thus the motivation for extending it to intersection controllers [5]. Unlike the shout-ahead enabled agents in [5], traffic signal controller agents deal with completely different game and action states, so adjustments and improvisations had to be made to fit the architecture to this application area.

Within the Shout-ahead architecture, an agent holds two rule sets, each containing a number of rules with corresponding actions that the agent can apply at a given time-step. Each rule has a weight, which factors into both the rule selection process later on, and the hybrid behavior learning method. One rule set makes decisions without any communicated intentions of other agents, and the other makes them using mainly these intentions. This way, the agent can incorporate both information about itself, its environment and other agents into its overall decision-making process. This makes the agent’s behavior more dynamic and responsive to its environment at a given state. The agent will first choose a rule from the first set using rule weights, as well as some randomness for exploration, to determine its intended action, and communicates it to the other agents. It then chooses a rule from the second set, containing communicated intentions of other agents, using the same process involving rule weights and randomness. Finally, based on a mixture of looking at rule weights and using probabilities, the agent will select one of the rules and apply its corresponding action. We did however make one adjustment to this process. An additional rule set was added to each of the two previously described rule sets with user defined rules. Rules within this set have maximum priority, such that if a rule within this set is applicable at a time step, it will always be chosen. If there are multiple applicable rules within this set at a given time, the one with the highest weight will be chosen. If no rules within this set are applicable, then the selection process will continue as described previously. This allows for prioritization of both rare, but critical circumstances such as giving way to an emergency vehicle, and rules we may deem necessary in certain circumstances. While the pertinence of the latter will become apparent one way or another in testing, but the former is clearly relevant and alone justifies the extra rule set.

To implement the Shout-ahead architecture for our problem case, we developed a set of 81 predicates, which together describe the agent’s environment. 57 of these are range predicates, meaning they pertain to a certain range of values. This is done because information about our environment, an intersection, includes data on waiting times and vehicle counts, which can best be represented in this manner. Thus, these 57 predicates actually serve as sub-predicates to 5 general ones: (1) LongestTimeWaitedToProceedStraight, which ranges from 0 to 300 seconds over 14 sub-predicates, (2) LongestTimeWaitedToTurnLeft, which ranges from 0 to 300 seconds over 14 sub-predicates, (3) NumOfCarsWaitingToProceedStraight, which ranges from 0 to 45 vehicles over 8 sub-predicates, (4)

NumOfCarsWaitingToTurnLeft, which ranges from 0 to 15 vehicles over 7 sub-predicates and (5) TimeSpentInCurrentPhase, which ranges from 0 to 300 seconds over 14 sub-predicates. The other predicates are concerned with the state of the traffic light signals in the intersection. A traffic light state is comprised of three components: (1) the direction traffic is arriving from, including “vertical” and “horizontal” for traffic flowing bilaterally, or “north-south”, “south-north”, “east-west” or “west-east” for traffic flowing in only one direction, with the destination (2) the phase, which describes whether traffic will flow “straight” (and naturally, “right” as well), such as in Fig. 1, “left”, such as in Fig. 2, or “straight left” (indicating both straight and left flow is permitted) during it, such as in Fig. 3, and (3) the state of the traffic signal, which is either “green” or “yellow”.

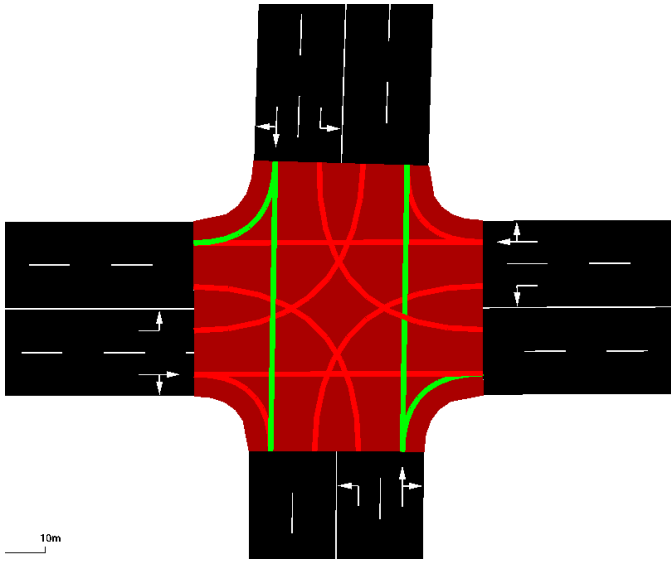


Fig. 1. Example of a “straight phase” in SUMO.

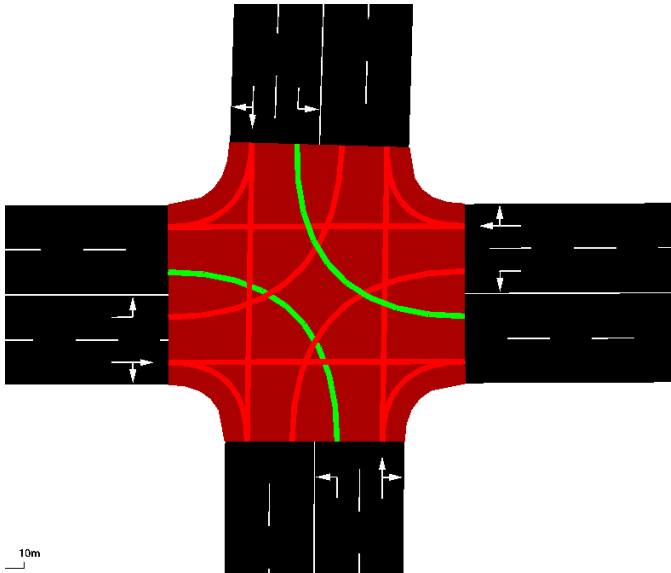


Fig. 2. Example of a “left phase” in SUMO.

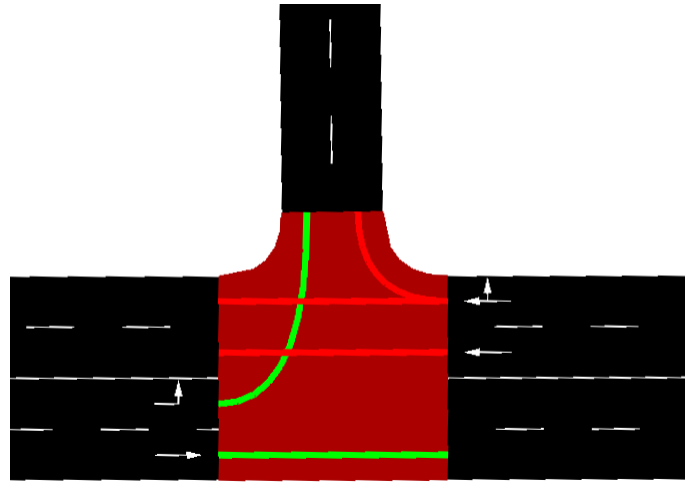


Fig. 3. Example of a “straight left” phase in SUMO.

A predicate exists for every combination of direction, phase and signal state, so to account for any type of intersection. Combinations of these predicates make up rules in the set RS , and dictate which rule may be applicable by an intersection controller agent at a given time. The set RS_{int} meanwhile has its own set of 39 predicates it uses to create rules. Unlike the predicates used in RS , these describe communicated intentions from other agents and the time since a given communication was received. The predicates describing the time since a communication was received take on similar range descriptions as those mentioned above. The intention of this is to provide an agent the opportunity to choose an action based on the intended actions of its neighbors. This principle lies at the heart of the shout-ahead architecture. The design of these predicates greatly affects how the agent behaves, and they will likely be adjusted as the evaluation process evolves to account for circumstances unforeseen at this point in time. The user-defined rule set contains 4 predicates, two of which specify which direction an emergency vehicle is approaching from. The other two check if an intersection has spent the maximum time in a green or yellow state. These two predicates are intended to keep intersections from getting stuck in a particular phase.

We also defined a set of actions to be applied when a rule is selected by an intersection controller agent. Traffic lights in our study will adhere to the same signal phase changing conventions Canadian traffic lights do, and so possible actions are bounded by said conventions, too. As such, actions when a given rule is selected to be applied by an agent include: (1) transitioning from a red phase to a green phase, (2) transitioning from a green phase to yellow phase, or (3) do-nothing. It may appear odd to omit the action of transitioning from yellow to red, but SUMO ignores explicitly describing this action, since it is implicit when a different direction is granted a green phase. Actual action sets vary from agent to agent, depending on their intersection configuration, which SUMO defines automatically, and the system cycles through it as it sees appropriate.

The hybrid behavior learning method for the Shout-ahead architecture is then responsible for optimizing the agent’s

behavior. Two learning algorithms are employed by this method. The main one is an evolutionary algorithm, with each individual consisting of the two rule sets for each agent, and the two user-defined rule sets. New individuals are created using crossover and mutation operators, with respect paid to rule weights in selecting individuals to be used in these operators, and some additional random factors. Naturally, we need a fitness measure for an individual, which would be determined by doing simulation (training) runs. An individual's fitness is calculated by the following equation:

$$\tau * (1/\text{simulationTime}) + \delta * (1 - 1/\text{sumOfAllRuleWeights})$$

where τ , $0 \leq \tau \leq 1$, is the simulation time weight factor and δ , $0 \leq \delta \leq 1$, is the rule weight factor. This equation rewards the individual with the fastest simulation time and highest aggregate rule weights. The weights of these two factors relative to each other can be adjusted by the user, and optimal values will be sought in testing. The agents perform a SARSA variant of Reinforcement Learning on their rule sets, to refine the weights of their existing rules, giving a higher priority to rules that have yielded positive results in the past and lower priority to the opposite [7]. This variant was shown to be effective with the Shout-ahead architecture when applied on unit agents in a video game [5]. These two learning phases each play crucial roles in improving the performance of the intersection controller agents, with the evolutionary algorithm providing new individuals that could potentially perform better than existing ones, and the reinforcement algorithm painting a better picture of what rules could be used to create better individuals. Altogether, if the objective is to increase traffic throughput, over time it is expected that individuals will be created that operate more efficiently than current intersection controller algorithms.

The Shout-ahead architecture and hybrid behavior learning method were implemented within the traffic simulator SUMO (Simulation of Urban MObility) [6]. SUMO provides a robust environment for both the implementation of the agent architecture for multiple intersection types, as well as for simulating a comprehensive set of test cases and training simulations to determine the effectiveness of the Shout-ahead architecture in this context. The extendibility of the SUMO platform and the fact it was designed for research of this very nature make it an ideal tool for the purpose of this project.

Within SUMO, we instantiate three types of intersection types, and thus three different classes of intersection controller agents: (1) *Standard 4-arm Intersection*: Fig. 4 for example, is an obvious choice, given its prevalence in road networks globally, (2) *T- Intersection*: Fig. 5 for example, is another prevalent intersection in road networks. Much like the Standard 4-arm Intersection, it is simple in design but presents many interesting test cases given the forced turning scenario presented by one of its straight phases ending, (3) *4-arm Incoming Layout Intersection*: Fig. 6 for example. Unlike the other two, this intersection type is more complex, and potentially presents more bottleneck, or inefficient signal timing, opportunities.



Fig. 4. Example of a Standard 4-arm Intersection in SUMO.

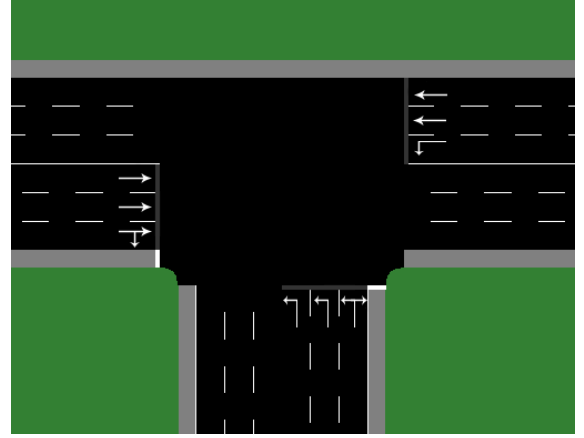


Fig. 5. Example of a T- Intersection in SUMO.



Fig. 6. Example of a 4-arm Incoming Layout Intersection in SUMO.

The design of this intersection is more irregular than the previous two selected, with a two-way road having two separate one-way roads converging in on it orthogonally, requiring at a

minimum three separate green phases. This is unlike the Standard 4-arm and T-Intersections, whom theoretically only require a minimum of two green phases, should they choose to exclude dedicated left-turning phases. In real world applications, this intersection type is far less prevalent than the other two intersections, so its inclusion in our study allows us to see the effectiveness of the Shout-ahead architecture when applied to potentially more complex, irregular scenarios.

IV. INTENDED EVALUATION OF PROPOSED SOLUTION

In this section, we will lay out my proposed experimental evaluation of the solution outlined in Section III. Given the purpose of this project is to demonstrate that a traffic signal controller integrated with learning and the shout-ahead architecture is superior to currently available traffic signal control algorithms, we must naturally pit the two against each other in evaluations.

SUMO provides a few default signal-timing algorithms that are similar to the preset ones used in intersections today [4][8]. One such algorithm asserts that all traffic lights work on a fixed cycle, with a default cycle time of 90 seconds. That is, all possible phases in an intersection are completed within the 90 second interval before beginning again. The cycle time can be adjusted by the user. Another signal-timing algorithm built into SUMO involves “actuated traffic lights”. This algorithm supports gap-based actuated traffic control, which prolongs traffic phases when a continuous vehicle stream is detected in the intersection. The purpose of this is to more effectively distribute green phases based on demand, affecting cycle durations based on traffic conditions as a result. Such a signal control scheme is common in Germany, and is similar to more modern systems used around the world, making it an attractive candidate to evaluate the learning and shout-ahead equipped agents in testing [4]. Additionally, SUMO allows for the importing of other signal-timing algorithms, however they are likely to require some retrofitting to prepare them for the SUMO interface, an effort proposition that greatly lessens its viability for our current purposes [8].

In addition to testing the learning and shout-ahead integrated traffic signal control agents against a “control” algorithm discussed above, it would be equally interesting to compare its performance against agents with the shout-ahead component turned off, leaving only learning enabled. Precisely, this would mean making use of only one rule set RS instead of two, RS and RS_{int} . Doing so would provide us insight into the benefits and overall effectiveness of the shout-ahead architecture in this application area. Experimental evaluation in [5], in the context of the Battle for Wesnoth, showed agents using shout-ahead were 2.5 times more successful than agents not using shout-ahead. By comparing the performance of shout-ahead and shout-ahead-less agents in a different context, a traffic network for instance, would lend more insight on shout-ahead’s general effectiveness as an agent architecture. Should it improve agent performance by a similarly significant margin, we can begin to confidently hypothesize that learning cooperative behavior with shout-ahead holds merit in a wide array of applications, for at the moment it only has one proven application space.

Before experimental evaluation can take place, the agent pools must learn. Given the dynamic nature of traffic networks, individuals (or, solutions to the problem) will learn on the same networks they are tested on, with varying traffic flows adjusted in various manners (as will be discussed later in this section). The exact number of networks to be used for learning and experimental evaluation is not set at this time, but they quantity will likely range between 3 and 5. The complexities of road networks in SUMO can vary widely, just as they do in the real world. For the purpose of this project, the most simplistic road network must at minimum incorporate one of each type of intersection described above: a standard 4-arm intersection, a T-intersection and a 4-arm Incoming Layout intersection, such as the one in Fig. 7.

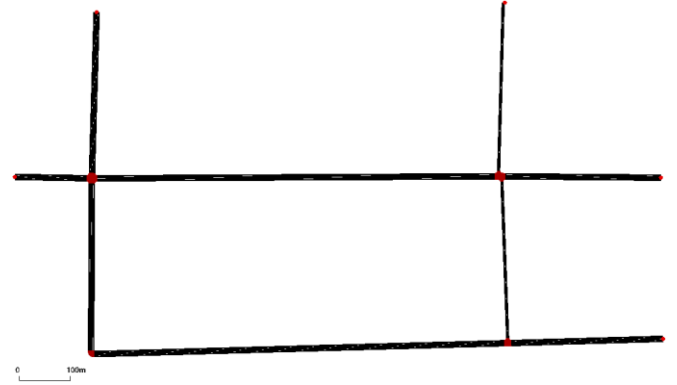


Fig. 7. Example of a simple road network in SUMO.

In SUMO, intersections are represented by the relatively large, red hexagons or octagons where different roads meet, while the red dots represent boundary points of the network, where vehicles originate and exit the network from. Any other red markers serve as SUMO’s representation of bends in a road. The maximum complexity of a SUMO road network for the purposes of this project is a difficult thing to define, given there a number of constraints to consider when doing so.

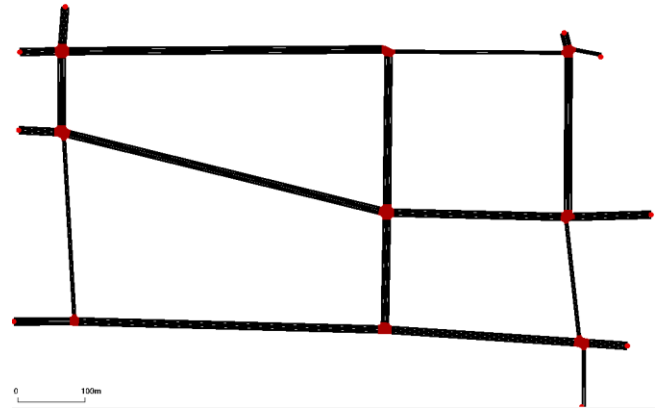


Fig. 8. Example of a complex road network in SUMO.

Firstly, we must consider the hardware constraints. The number of intersections of a road network has a direct relationship with the maximum number of vehicles that can be going through the network at any one time, meaning computational complexity increases significantly with the introduction of every

intersection. The second constraint to consider is time; the more complex a road network is, the longer it will take for the system to complete training and testing runs. Given the relatively limited timeframe this project has to be completed in (approximately 3 months from the time of this writing), we must be wary when designing a road network. For the time being, we have decided the most complex road network to be used will feature no more than three of each of our chosen intersection types. An example of such a network is in Fig. 8.

Finally, we must consider what traffic flows will be used for both learning and testing. This is not a trivial matter, and will require some thought, and trial and error to decide upon. Training will require multiple different traffic flows, however the differences between these must be carefully considered as well. Emulating volume at different times of day and the traffic that comes with it would be ideal, however using starkly different flows for the training of one agent pool, such as morning and afternoon rush hour, may just confuse the learner and be ultimately counterproductive. It is likely that the differences between flows will be minimal for learning, and the same ones will be utilized in testing. Additionally, new flows will be designed for testing, their variance from the learning ones remaining in question as well. We simply do not have enough concrete information about what constitutes a useful flow at this time to make any concrete assertions.

Once we have decided on our networks for learning and testing, and the accompanying traffic flows, we will evaluate the performances of the three types of traffic signal control agents (timing-cycle, shout-ahead enabled and shout-ahead disabled) to see which proves superior in which circumstances. As mentioned previously, the performance of an agent-type will be measured by how quickly they can facilitate a certain flow through the network. That is, the time it takes to get every vehicle in the network from its starting point to its destination. As the experimental evaluation goes on, there will likely be results that motivate new experiments to be carried out and adjustments to be made to our shout-ahead implementation.

V. RELATED WORK

Naturally, there has been a lot of work in the area of traffic signal optimization with artificial intelligence and machine learning. Most traffic light controllers rely on pre-timed cycle policies, but [11] showed that using an intelligent agent to generate traffic signal controller (TSC) policies was more effective, and [21] detailed further examples in which machine learning would be useful in the context of vehicular networks. Furthermore, [25] highlighted that many current solutions to the TSC problem are centralized, and argued that this centralized nature makes them infeasible for large scale adaptive traffic signal control due to the high dimension of the joint action space. As such, decentralized solutions must be explored in order to find an optimal solution to TSC problem, a characteristic shared by this project and its related works. While specific methods vary between the papers, their end-goals align: to optimize traffic throughput through the implementation of algorithms that utilize artificial intelligence and machine learning principles.

The use of multi-agent reinforcement learning (MARL) was the most popular method observed, solely utilized by [8], [9], [10], [12], [13], [16], [17], [18], [20], [23], [24] and [25]. Naturally, nuances exist between each of the studies, despite their shared solution to the TSC problem. [8] constructed the problem as a discounted cost Markov decision process, and applied multi-agent Q-learning to obtain dynamic traffic signal control policies, with Q-factors updated based on communication with neighboring TSC agents. [9] and [18] make use of similar neighbor communication approaches, except [9] utilizes the max-plus algorithm to achieve coordination between TSC agents, while [18] defines their own coordination algorithm. In both cases, agent protocols are improved via reinforcement learning. Unlike [8] and [18], [9] applies its solution directly to a large-scale problem to verify its efficacy in realistic settings, and found success with this approach.

Cooperation between individual TSC agents wasn't the only way MARL was applied to the TSC problem. [13] defined two individual agent classes which worked together to optimize traffic flow. They presented a detailed protocol making use of two different classes of agents, dubbed intersection managers and driver agents. The presence of the latter agent marked a departure from the singular focus applied on the TSC agent by most other papers. [23] meanwhile assigned multiple local controller agents to a single traffic junction, making use of local traffic state data, a learned probabilistic model of car behavior, and a learned value function which indicates how traffic light decisions affect long-term utility, in terms of the average waiting time of cars, to optimize the performance of the intersection.

Some algorithms expanded upon the basic MARL strategies implemented by papers mentioned above. [17] for instance, detailed a Q-learning algorithm with a feedforward neural network for value function approximation. This approach was shown to be a superior method in this context when compared against a longest-queue first algorithm. Examples of other solutions that make a point of their decentralized nature are [15], [20], [22] and [25]. [25] details a fully scalable and decentralized MARL algorithm for an advantage actor critic agent within the context of adaptive traffic signal control, as a proposed evolution to past solutions to the TSC problem using MARL. [20] meanwhile proposed a TSC system enabled by a hierarchical multi-agent modeling framework. In a similar vein, [14] proposed a new method called Reinforcement Learning with Context Detection (RL-CD), which it claimed would be superior to popular methods such as Q-learning due to its adaptability in dynamic environments, such as intersections. RL-CD employs multiple partial models of the environment to allow the learning system to partition its knowledge in a way that each model becomes responsible for understanding one kind of flow pattern. Empirical results supported the authors claim.

While communication and cooperation between agents seems intuitively pragmatic, it can be difficult to discern what information is useful to share and what the consequences of such cooperation are. [24] investigates this very topic in the

context of a distributed urban traffic control, where each junction is controlled by an independent agent. The authors analyze both the benefits of sharing information, and the consequences of doing so.

Naturally, it must be acknowledged the entire premise of this project is rooted in the work presented in [5], which lays out the shout-ahead architecture and hybrid behavior learning method in abstracted detail, creating a blueprint for future implementations to follow. Though the application area of [5] does not align with this project, the success of its experimental evaluations merit the implementation of the architecture in a different application area.

VI. CONCLUSION

Ineffective traffic signal timing algorithms can exacerbate poor traffic conditions, creating bottlenecks leading to traffic jams. Traffic signals within individual intersections are managed by a controller agent, which is responsible for changing signal states based on an algorithm that relies on a predetermined signal timing cycle. Though the cycle can be altered based on heavier traffic flow, applicable adjustments are also predetermined and thus don't take into consideration the full state of the intersection, information that could lead to a different action being taken. By implementing the Shout-ahead agent architecture and a hybrid behavior learning method for it into intersection controller agents, these agents gain a level of dynamic flexibility in their decision making, allowing them to take the most optimal action for the current state of their environment, instead of being bound to rigid timing-cycle. Theoretically, the ability for a traffic signal control agent to learn which actions are proper in which situations, coupled with communication with its neighboring cohorts, should result in more optimal actions being taken, and more effective intersection control than a near-static algorithm would provide.

REFERENCES

- [1] M. Bennardo, "StatsCan study shows Canadian commute times are getting longer — and it's costing us," *CBC*, March 2019. [Online]. Available: <https://www.cbc.ca/news/business/statistics-canada-commute-times-study-1.5038796>. [Accessed: Sept. 21, 2019].
- [2] T. Yaropud, J. Gilmore, S. LaRoche-Côté, "Results from the 2016 Census: Long commutes to work by car," *Statistics Canada*, February 2019. [Online]. Available: <https://www150.statcan.gc.ca/n1/en/pub/75-006-x/2019001/article/00002-eng.pdf?st=wRNvAal2>. [Accessed: Sept. 21, 2019].
- [3] Y. Shaoxin, X. Zhao, Y. An, "Identification and optimization of traffic bottleneck with signal timing," *Journal of Traffic and Transportation Engineering (English Edition)*, vol. 2, pp. 353-361, October 2014.
- [4] National Academies of Sciences, Engineering, and Medicine, *Signal Timing Manual – Second Edition*, Washington, DC: The National Academies Press, 2015.
- [5] S. Paskaradevan, J. Denzinger, D. Wehr, "Learning cooperative behaviour for the shout-ahead architecture," *Web Intelligence and Agent Systems: An International Journal*, vol. 12, pp. 309-324, October 2014.
- [6] M. Behrisch, L. Bieker, J. Erdmann, and D. Krajzewicz, "SUMO—simulation of urban mobility: an overview," in *Proceedings of the Third International Conference on Advances in System Simulation (SIMUL 2011)*, 2011.
- [7] "Automatically Generated TLS-Programs," *Simulation/Traffic Lights – SUMO Documentation*. [Online]. Available: https://sumo.dlr.de/docs/Simulation/Traffic_Lights.html#actuated_traffic_lights. [Accessed: 17-Jan-2020].
- [8] R.S. Sutton and A.G. Barto, *Reinforcement Learning: An Introduction*, The MIT Press, 1998.
- [9] K. J. Prabuchandran, A. N. Hemanth Kumar and S. Bhatnagar, "Multi-agent reinforcement learning for traffic signal control," *17th International IEEE Conference on Intelligent Transportation Systems (ITSC)*, Qingdao, 2014, pp. 2529-2534.
- [10] L. Kuyer, S. Whiteson, B. Bakker, and N. Vlassis, "Multiagent Reinforcement Learning for Urban Traffic Control Using Coordination Graphs," *Machine Learning and Knowledge Discovery in Databases Lecture Notes in Computer Science*, pp. 656–671, 2008.
- [11] M. Abdoos, N. Mozayani and A. L. C. Bazzan, "Traffic light control in non-stationary environments based on multi agent Q-learning," *2011 14th International IEEE Conference on Intelligent Transportation Systems (ITSC)*, Washington, DC, 2011, pp. 1580-1585.
- [12] P. Gregoire, C. Desjardins, J. Laumonier and B. Chaib-draa, "Urban Traffic Control Based on Learning Agents," *2007 IEEE Intelligent Transportation Systems Conference*, Seattle, WA, 2007, pp. 916-921.
- [13] M. Wiering, "Multi-Agent Reinforcement Learning for Traffic Light Control," *Proceedings of the Seventeenth International Conference on Machine Learning (ICML 2000)*, Stanford, CA, 2000.
- [14] K. Dresner and P. Stone, "Multiagent Traffic Management: Opportunities for Multiagent Learning," *Learning and Adaption in Multi-Agent Systems Lecture Notes in Computer Science*, pp. 129–138, 2006.
- [15] D. de Oliveira Boschetti, A. Bazzan, B. da Silva, E. Basso and L. Nunes, "Reinforcement Learning based Control of Traffic Lights in Non-stationary Environments: A Case Study in a Microscopic Simulator," *Proceedings of the 4th European Workshop on Multi-Agent Systems EUMAS'06*, Lisbon, Portugal, 2006.
- [16] E. Camponogara and W. Kraus, "Distributed Learning Agents in Urban Traffic Control," *Progress in Artificial Intelligence Lecture Notes in Computer Science*, pp. 324–335, 2003.
- [17] M. A. Khamis and W. Gomaa, "Enhanced multiagent multi-objective reinforcement learning for urban traffic light control," *2012 11th International Conference on Machine Learning and Applications*, Boca Raton, FL, 2012, pp. 586-591.
- [18] I. Arel, C. Liu, T. Urbanik and A. G. Kohls, "Reinforcement learning-based multi-agent system for network traffic signal control," *IET Intelligent Transport Systems*, vol. 4, no. 2, pp. 128-135, 2010.
- [19] E. Van der Pol and F.A. Oliehoek, "Coordinated Deep Reinforcement Learners for Traffic Light Control," *NIPS'16 Workshop on Learning, Inference and Control of Multi-Agent Systems*, December 2016.
- [20] D. Houli, L. Zhiheng, and Z. Yi, "Multiobjective Reinforcement Learning for Traffic Signal Control Using Vehicular Ad Hoc Network," *EURASIP Journal on Advances in Signal Processing*, vol. 2010, no. 1, 2010.
- [21] J. Jin and X. Ma, "Hierarchical multi-agent control of traffic lights based on collective learning," *Engineering Applications of Artificial Intelligence*, vol. 68, pp. 236–248, 2018.
- [22] H. Ye, L. Liang, G. Y. Li, J. Kim, L. Lu and M. Wu, "Machine Learning for Vehicular Networks: Recent Advances and Application Examples," in *IEEE Vehicular Technology Magazine*, vol. 13, no. 2, pp. 94-101, June 2018.
- [23] S. El-Tantawy and B. Abdulhai, "Multi-Agent Reinforcement Learning for Integrated Network of Adaptive Traffic Signal Controllers (MARLIN-ATSC)," *2012 15th International IEEE Conference on Intelligent Transportation Systems*, 2012.
- [24] B. Bakker, S. Whiteson, L. Kester, and F. C. A. Groen, "Traffic Light Control by Multiagent Reinforcement Learning Systems," *Interactive Collaborative Information Systems Studies in Computational Intelligence*, pp. 475–510, 2010.
- [25] D. D. Oliveira and A. L. Bazzan, "Multiagent Learning on Traffic Lights Control: Effects of Using Shared Information," *Multi-Agent Systems for Traffic and Transportation*, pp. 307–321, 2009.
- [26] T. Chu, J. Wang, L. Codeca and Z. Li, "Multi-Agent Deep Reinforcement learning for large-scale Traffic Signal Control," *IEEE Transactions on Intelligent Transportation Systems*, 2019.
- [27] R.S. Sutton and A.G. Barto, *Reinforcement Learning: An Introduction*, The MIT Press, 1998.