

چارچوبی جامع مبتنی بر یادگیری تقویتی برای کنترل سیگنال‌های ترافیکی در محیط‌های V2X
نویسندگان :

Maonan Wang *, Xi Xiong †, Yuheng Kan, Chengcheng Xu ‡, Man-On

Wang, Xi Xiong, Yuheng Kan, Chengcheng Xu, Man On Pun

: IEEE Transactions on Vehicular Technology, 2024. DOI / (IEEE Xplore via DOI):
<https://doi.org/10.1109/TVT.2024.3403879>.) IEEE Transactions on Vehicular Technology

6.1-8.6 : impact factor

چکیده

ترافیک متراکم یکی از چالش‌های مزمن در نواحی شهری است که نیازمند توسعه سامانه‌های مؤثر کنترل سیگنال ترافیکی¹ (TSC) می‌باشد. در حالی که روش‌های مبتنی بر یادگیری تقویتی² (RL) نشان داده‌اند می‌توانند عملکرد TSC را بهبود بخشند، تعمیم‌پذیری این روش‌ها میان تقاطع‌هایی با ساختارهای متفاوت همچنان چالش برانگیز است. در این کار یک چارچوب عمومی مبتنی بر RL برای محیط‌های خودرو به همه چیز³ (V2X) ارائه می‌شود. چارچوب پیشنهادی طراحی نوآورانه‌ای برای عامل (agent) معرفی می‌کند که ماتریس تقاطع (junction matrix) را برای توصیف وضعیت تقاطع وارد می‌سازد، و بدین ترتیب مدل پیشنهادی برای تقاطع‌های متنوع قابل اعمال می‌شود. به منظور توانمندسازی بیشتر چارچوب RL پیشنهادی برای مقابله با ساختارهای مختلف تقاطع، روش‌های افزایش داده (augmentation) ویژه‌ای برای حالت ترافیکی مناسب سیستم‌های کنترل چراغ راهنمایی طراحی شده‌اند. در نهایت، نتایج تجربی گسترده حاصل از پیکربندی‌های مختلف تقاطع، مؤثر بودن چارچوب پیشنهادی را تأیید می‌کنند. کد منبع این کار در دسترس است:

https://github.com/wmn7/Universal_Light

واژگان کلیدی: کنترل سیگنال ترافیکی • مدل‌های عمومی • یادگیری تقویتی • افزایش حالت ترافیکی

¹ Traffic Signal Control

² Reinforcement learning

³ Vehicle to everything

ترافیک متراکم به عنوان چالشی حیاتی در نواحی شهری مطرح است که منجر به اتلاف زمان شهروندان، مصرف سوخت بیش از حد و افزایش انتشار گازهای گلخانه‌ای می‌شود [۱]. برای کاهش تراکم، روش‌های سنتی کنترل سیگنال ترافیکی مانند کنترل با چرخه ثابت [۲]، روش وبستر [۳] و کنترل خودسازمان‌دهنده چراغ‌ها (SOTL) [۴] توسعه یافته‌اند. با این حال، با رشد شهرها، این رویکردهای سنتی اغلب در مواجهه با افزایش حجم ترافیک و شرایط جاده‌ای پویا ناکافی ظاهر می‌شوند [۵]. ظهور فناوری‌های V2X در سامانه‌های حمل‌ونقل نویدبخش راه‌حلی است که امکان برقراری ارتباط و تبادل داده میان وسایل نقلیه، زیرساخت‌ها و سایر کاربران جاده‌ای را فراهم می‌آورد [۶]. با استفاده از داده‌های بلادرنگ استخراج‌شده از وسایل نقلیه، زیرساخت‌های مدیریت ترافیک می‌توانند حرکت وسایل نقلیه و عابران را در تقاطع‌ها به صورت پویا و بر اساس شرایط جاری ترافیک تنظیم کنند [۷]. برای کنترل چراغ‌ها بر اساس شرایط ترافیکی لحظه‌ای، روش‌های متعددی مبتنی بر RL پیشنهاد شده‌اند. این روش‌ها معمولاً از سه رویکرد اصلی برای تنظیم چراغ‌ها استفاده می‌کنند: «انتخاب فاز بعدی»، «نگه‌داشتن یا تغییر فاز» و «تنظیم مدت زمان فاز». به طور مشخص، در رویکرد «انتخاب فاز بعدی»، عامل RL فاز بعدی را تعیین می‌کند و اجازه می‌دهد توالی فازها به جای ثابت بودن، انعطاف‌پذیر باشد [۸-۲۲]. هرچند این رویکرد انعطاف‌پذیری می‌آورد، ممکن است برای رانندگان نامأنوس یا گیج‌کننده به نظر برسد زیرا انتخاب فاز می‌تواند ظاهراً تصادفی شود که منجر به افزایش خطرات حوادث می‌گردد. در رویکرد «نگه‌داشتن یا تغییر فاز»، عامل RL تصمیم می‌گیرد که آیا فاز کنونی را حفظ کند یا تغییر دهد [۲۳-۲۵]. نهایتاً، در رویکرد «تنظیم مدت زمان فاز»، عامل طول مدت فاز جاری را از میان مجموعه‌ای از گزینه‌های از پیش تعیین‌شده انتخاب می‌کند [۲۶-۲۸]. از طریق تعامل مستقیم با محیط، عامل RL یاد می‌گیرد که بر مبنای تجربیات واقعی زمان، خود را با تغییرات شرایط ترافیک وفق دهد.

با وجود پیشرفت‌های قابل توجه روش‌های RL مطرح‌شده، یک محدودیت اصلی این است که بسیاری از آن‌ها برای ساختارهای تقاطع مشخص طراحی شده‌اند. به عبارت دیگر، این مدل‌های RL باید هنگام مواجهه با تقاطع‌هایی که از حیث تعداد راه‌های ورودی، خطوط و فازها متفاوت‌اند، از نو طراحی و دوباره آموزش داده شوند که این امر مستلزم صرف منابع قابل توجهی برای جمع‌آوری داده، توسعه مدل و آزمایش است [۲۹]. بنابراین توسعه مدل‌های عمومی که بتوان آن‌ها را به آسانی در طیف وسیعی از تقاطع‌ها به کار گرفت، اهمیت دارد تا پیاده‌سازی V2X با نیاز کمتر به سفارشی‌سازی و توسعه مجدد در هر تقاطع، به طور مؤثر مقیاس‌پذیر شود [۳۰].

در ادبیات موضوع، چندین مدل تعمیم‌یافته برای انواع مختلف تقاطع پیشنهاد شده‌اند [۱۰-۱۲، ۱۸]. با وجود عملکرد خوب آن‌ها، این مدل‌های تعمیم‌یافته تنها برای پیکربندی‌های مشخصی که در طراحی آن‌ها مورد نظر قرار گرفته‌اند قابل استفاده‌اند و هنگام مواجهه با پیکربندی‌های ناآشنا دچار افت عملکرد می‌شوند. با توجه به این انگیزه‌ها، ما چارچوبی مبتنی بر یادگیری تقویتی با تقویت حالت ترافیکی عمومی^۱ (UniTSA) ارائه می‌کنیم. چارچوب پیشنهادی امکان آموزش یک عامل عمومی با استفاده از داده‌های افزوده‌شده را برای کنترل سیگنال ترافیکی فراهم می‌سازد. برای رسیدگی به تقاطع‌هایی با ساختارهای متنوع، ماتریس تقاطع توسعه داده شده است تا وضعیت تقاطع را توصیف کند. با استفاده از این ماتریس نوآورانه، تقاطع‌های دارای ساختارهای مختلف را می‌توان با ماتریس‌هایی از همان اندازه نمایش داد. علاوه بر این، در طراحی عمل (action) در این کار، رویکرد «نگه‌داشتن یا تغییر فاز فعلی» به کار گرفته شده است تا ساختار مدل در سراسر تقاطع‌های مختلف حفظ شود. برای مقابله با تقاطع‌های نامشخص (unseen)، پنج روش افزایشی حالت ترافیکی توسعه یافته‌اند تا جمع‌آوری داده‌های عامل را در طول فرایند آموزش غنی سازند. این روش‌های افزایشی، آموزش جامع‌تری فراهم می‌آورند که در نهایت به عملکرد بهتر در تقاطع‌هایی منجر می‌شود که پیکربندی آن‌ها در مجموعه آموزشی وجود نداشته است. علاوه بر این، سازگار سازی کم‌رتبه (Low-Rank Adaptation — LoRA) [۳۱] برای انجام بهینه سازی بیشتر مدل در تقاطع‌های حیاتی به کار گرفته می‌شود. در نهایت، آزمایش‌های گسترده‌ای با استفاده از پلتفرم SUMO^۲ بر روی تقاطع‌هایی با شمار مختلف راه‌های نزدیک‌شونده، خطوط و فازها انجام گرفت. نتایج تجربی نشان می‌دهد مدل UniTSA پیشنهادی عملکرد بسیار خوبی حتی در تقاطع‌های ناآشنا ارائه می‌دهد.

مشارکت‌های ما را می‌توان به صورت زیر خلاصه کرد:

(الف) یک چارچوب تطبیقی کنترل سیگنال ترافیکی به نام UniTSA برای محیط‌های V2X پیشنهاد شده است که بر پایه یک مدل RL عمومی ساخته شده و طراحی عامل نوینی دارد که می‌تواند ساختارهای متنوع تقاطع را مدیریت کند. علاوه بر این، یک مکانیزم فاین تیون برای ارتقای عملکرد در تقاطع‌های کلیدی طراحی شده است؛

(ب) روش‌های افزایش حالت ترافیکی برای چارچوب TSC پیشنهادی توسعه یافته‌اند تا درک عامل از تقاطع‌های متنوع را افزایش داده و عملکرد را در هر دو مجموعه آموزش و آزمون بهبود بخشند؛

(ج) آزمایش‌های گسترده روی ۱۲ تقاطع با ساختارهای گوناگون نشان می‌دهد که مدل UniTSA پیشنهادی به‌طور قابل توجهی از مدل‌های عمومی مرسوم بهتر عمل می‌کند. افزون بر این، برای یک تقاطع جدید، UniTSA

^۱ Universal Traffic State Augmentation

^۲ Simulation of Urban MObility

می‌تواند مدل پیش‌آموزش‌یافته خود را با فاین‌تیون کردن به‌سرعت تطبیق دهد و با زمان آموزش به‌مراتب کمتر نسبت به آموزش یک مدل نو از ابتدا، عملکردی قابل‌مقایسه یا بهتر به‌دست آورد.

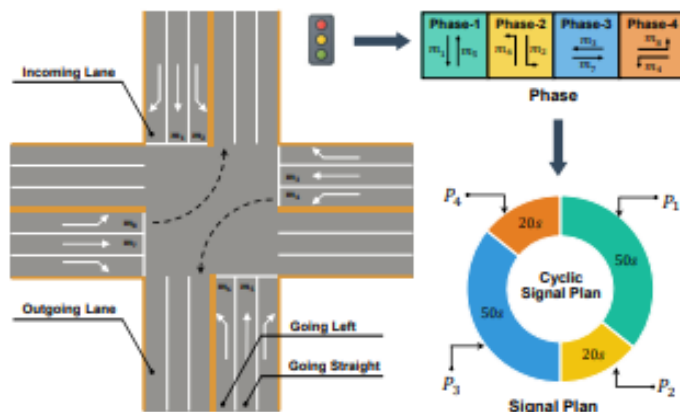
باقی ساختار مقاله به‌صورت زیر است: بخش ۲ مروری بر کارهای مرتبط با TSC ارائه می‌دهد، در حالی که بخش ۳ اصطلاحات مرتبط با راه و چراغ‌های ترافیکی را معرفی می‌کند. سپس بخش ۴ چارچوب UniTSA پیشنهادی و پنج روش افزایش حالت ترافیکی را مطرح می‌سازد و بخش ۵ تنظیمات تجربی و نتایج UniTSA را تشریح می‌کند. در نهایت بخش ۶ مقاله را جمع‌بندی می‌کند.

(۲) کارهای مرتبط

پژوهش‌های گسترده‌ای در حوزه حمل‌ونقل برای مطالعه کنترل سیگنال ترافیکی انجام شده است. به‌طور سنتی، کنترل زمان ثابت یکی از قدیمی‌ترین و گسترده‌ترین روش‌های TSC است [۳۳]. این روش‌ها بر زمان‌بندی‌های از پیش تعیین‌شده مبتنی هستند که براساس الگوهای تاریخی ترافیک یا دستورالعمل‌های مهندسی تعیین می‌شوند. تکنیک‌های بهینه‌سازی مختلفی برای تعیین برنامه‌های زمان ثابت بهینه برای پیکربندی‌های تقاطع خاص پیشنهاد شده‌اند که روش Webster یکی از موفق‌ترین روش‌ها برای مورد یک تقاطع است؛ این روش طول چرخه و تقسیم فاز را متناسب با حجم ترافیک در یک بازه معین (مثلاً ۱۵ یا ۳۰ دقیقه گذشته) محاسبه می‌کند. با این حال، چنین روش‌های زمان ثابتی اغلب به‌دلیل عدم انطباق با تغییرات پویا شرایط ترافیکی عملکرد زیرحد بهینه دارند. برای مقابله با این مشکل، روش‌های کنترلی عمل‌گرا مانند سیستم هم‌آهنگ‌شده و پویا سیدنی (SCATS) [۳۴]، کنترل فشار حداکثر [۳۵] و SOTL [۴] طراحی شده‌اند تا زمان‌بندی‌ها را براساس تقاضای ترافیکی بلادرنگ تنظیم کنند. با وجود مزایای متعدد، این روش‌های عمل‌گرا نیازمند تنظیمات تخصصی برای هر تقاطع هستند و در سناریوهای پیچیده عملکردشان کاهش می‌یابد.

در سال‌های اخیر، روش‌های مبتنی بر RL توجه زیادی را به‌خود جلب کرده‌اند به‌دلیل توانایی قابل‌توجه آن‌ها در سازگاری با شرایط ترافیک بلادرنگ و قابلیت یادگیری سیاست‌های کنترلی بهینه در سناریوهای پیچیده [۳۶]. به‌صورت کلی، این روش‌ها را می‌توان به سه دسته تقسیم کرد: روش‌های مبتنی بر ارزش [۸و۹و۱۰و۱۱و۱۲و۱۳و۲۳و۲۵] روش‌های مبتنی بر سیاست [۱۵و۱۶و۱۷و۱۸و۱۹] و روش‌های بازیگر-منتقد [۲۰و۲۱و۲۲و۲۸] باوجود عملکرد خوب، اکثر روش‌های RL موجود بر آموزش مدل‌هایی برای پیکربندی‌های تقاطع مشخص متمرکزند. تلاش‌هایی نیز برای ساخت مدل‌های تعمیم‌یافته TSC انجام شده است؛ برای نمونه

[۱۱] مدلی نسبتاً عمومی تر با بهره‌گیری از استراتژی متا-یادگیری [۳۷] آموزش می‌دهد. با این حال، مدل نام برده ۸ نیازمند بازآموزی پارامترهای مدل برای هر تقاطع جدید است. برای غلبه بر این نقطه ضعف، [۱۸، ۱۰، ۱۲] مدل‌های عمومی از طریق اشتراک‌گذاری پارامترها ایجاد کردند؛ اما این روش‌ها ساختار فاز اصلی چراغ راهنما را حفظ نمی‌کنند. در مقابل، روش پیشنهادی ما می‌تواند ساختار چراغ سیگنال موجود را حفظ کند و هم‌زمان با فاین‌تیون تقاطع‌های حیاتی عملکرد را به‌طور چشمگیری بهبود دهد



شکل ۱) نقشه یک چهارراه استاندارد

۳- مفاهیم اولیه

در این بخش، برخی از مفاهیمی که در این کار مورد استفاده قرار گرفته‌اند، با استفاده از یک تقاطع استاندارد چهارراه مطابق شکل ۱ تعریف می‌شوند. این مفاهیم به‌آسانی قابل تعمیم به تقاطع‌هایی با ساختارهای متفاوت هستند.

الف) خط عبور : خط عبور به بخشی از راه گفته می‌شود که مسیر تعریف‌شده‌ای را برای حرکت وسایل نقلیه در یک جهت خاص فراهم می‌کند. در یک تقاطع معمولی دو نوع خط وجود دارد:

- خطوط ورودی l_{in} که وسایل نقلیه از آن وارد می‌شوند
- خطوط خروجی l_{out} که وسایل نقلیه از آن خارج می‌شوند.

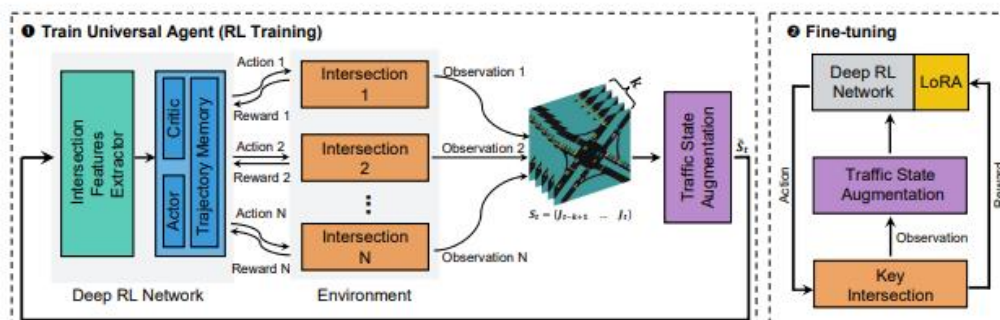
ب) **حرکت ترافیکی** : حرکت ترافیکی به اتصال میان یک خط ورودی l_{in} و یک خط خروجی l_{out} اطلاق می‌شود. در یک تقاطع چهارراه متداول در سمت چپ شکل ۱، مجموعاً ۱۲ حرکت وجود دارد که شامل گردش‌های راست، گردش‌های چپ و حرکات مستقیم در چهار جهت می‌باشند.

ج) **سیگنال حرکت**: سیگنال حرکت بر روی حرکت ترافیکی تعریف می‌شود. چراغ سبز یعنی حرکت مربوطه مجاز است و چراغ قرمز یعنی حرکت ممنوع است. از آنجا که در بسیاری از کشورها گردش راست مستقل از چراغ مجاز است، تنها ۸ سیگنال حرکت از ۱۲ حرکت ممکن در یک چهارراه مورد استفاده قرار می‌گیرند. این ۸ حرکت که با M_1, M_2, \dots, M_8 نمایش داده شده‌اند عبارت‌اند از: حرکت شمال سو (N)، شمال سو با گردش چپ (NL)، شرق سو (E)، شرق سو با گردش چپ (EL)، غرب سو (W)، غرب سو با گردش چپ (WL)، جنوب سو (S)، و جنوب سو با گردش چپ (SL). برای مثال، m_8 نشان می‌دهد که خودروها می‌توانند از سمت غرب به سمت شمال حرکت کنند.

د) فاز : فاز ترکیبی از چند سیگنال حرکت است. هر فاز مجموعه خاصی از حرکات ترافیکی را مجاز کرده و سایر حرکات را محدود می‌کند. قسمت بالا-راست شکل ۱ چهار فاز مربوط به یک چهارراه را نشان می‌دهد. مثلاً، فاز ۱ شامل حرکات‌های m_1 و m_5 است و اجازه می‌دهد خودروهای شمال-جنوب حرکت کنند، در حالی که سایر حرکات ممنوع می‌باشند.

ه) **طرح سیگنال**: طرح سیگنال دنباله‌ای از فازها و مدت زمان هر فاز است که برای کنترل چراغ‌های یک تقاطع به کار می‌رود. به صورت ریاضی به صورت مجموعه‌ای از $\{(p_1, t_1), (p_2, t_2), \dots\}$ نمایش داده می‌شود، که p_i فاز و t_i مدت زمان آن است. معمولاً توالی فازها چرخه‌ای است. قسمت پایین-راست شکل ۱ یک طرح چرخه‌ای سیگنال را نشان می‌دهد که در آن مدت زمان هر فاز برابر است با:

$$\begin{aligned} t_1 &= 50, \\ &= 20, \\ &= 20. \end{aligned}$$



شکل ۲) نمای کلی ساختار UniTSA

۴- روش‌شناسی

۴-۱ چارچوب

همان‌گونه که در شکل ۲ نشان داده شده است، چارچوب پیشنهادی UniTSA از دو ماژول تشکیل شده است:

الف) ماژول آموزش عامل عمومی

این بخش یک عامل RL عمومی را با استفاده از تقاطع‌های متفاوت آموزش می‌دهد. این ماژول از طراحی جدید عامل بهره می‌گیرد که قادر است تقاطع‌هایی با توپولوژی‌ها و طرح‌های سیگنالی متفاوت را با یک ساختار ثابت مدل کند. این امر با استفاده از ویژگی‌های حرکت‌ها و اعمال از نوع «فاز بعدی یا نه» و همچنین پنج روش نوآورانه افزایش حالت ترافیک محقق می‌شود. ماژول فاین‌تیون برای تقاطع‌های کلیدی در این مرحله، مدل حاصل از ماژول اول برای برخی تقاطع‌های خاص که اهمیت بیشتری دارند، به‌صورت تطبیقی اصلاح می‌شود. جزئیات هر مرحله در بخش‌های بعدی تشریح خواهد شد.

۴-۲ طراحی عامل و ماتریس تقاطع

وضعیت: تقاطع‌ها ممکن است تعداد خطوط متفاوتی داشته باشند که در صورت ثبت ویژگی‌ها در سطح خط، فضای حالت ابعاد متفاوتی خواهد داشت. همان‌طور که در بخش ۳ اشاره شد، مشاهده می‌شود که در یک چهارراه تنها هشت سیگنال حرکت معتبر وجود دارد، بدون توجه به تعداد خطوط.

با الهام از این مشاهده، ما پیشنهاد می‌کنیم که وضعیت تقاطع در زمان t با یک ماتریس تقاطع (به شکل زیر نمایش داده شود):

$$J_t = \begin{bmatrix} m_1^t \\ m_2^t \\ \vdots \\ m_8^t \end{bmatrix} \quad (1)$$

که در آن $[\cdot]^T$ ترانهاده را نشان می‌دهد و بردار m_i^t از طول ۸، اطلاعات استخراج شده از حرکت شماره i در زمان t را نشان می‌دهد.

اجزای بردار m_i^t :

برداری که هر حرکت شامل سه دسته اطلاعات است:

الف) ویژگی‌های ترافیکی

این ویژگی‌ها سطح تراکم حرکت را با استفاده از پارامترهای زیر توصیف می‌کنند:

جریان متوسط ترافیک $F_{i,t}$

بیشینه اشغال $O_{i,t}^{\max}$ (Occupancy)

میانگین اشغال $O_{i,t}^{\text{mean}}$

ب) ویژگی‌های حرکت

این بخش اطلاعات ذاتی حرکت را ارائه می‌کند:

جهت حرکت I_i^S مثلاً مستقیم بودن یا نبودن

تعداد خطوط اختصاص یافته به حرکت L_i

۳- ویژگی‌های سیگنال ترافیکی

این بخش شامل سه پارامتر باینری است:

- $I_{i,t}^{cg}$: آیا حرکت فعلاً چراغ سبز دارد یا نه
- $I_{i,t}^{ng}$: آیا حرکت در فاز بعد چراغ سبز خواهد داشت یا نه
- $I_{i,t}^{mg}$: آیا مدت زمان حداقل چراغ سبز فعلی به پایان رسیده یا نه

این ۸ ویژگی به‌سادگی قابل استخراج هستند، که مدل را برای استفاده عملی قابل اجرا می‌سازد.

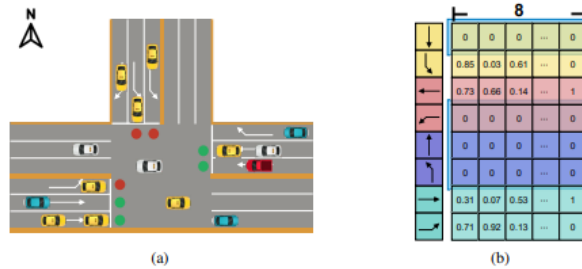
بردار m_i^t به‌صورت زیر تعریف می‌شود:

$$m_i^t = [F_i, t, O_i, tmax, O_i, tmean, I_{is}, L_i, I_i, tcg, I_i, tng, I_i, tmg]^T \quad (2)$$

۳- افزودن اطلاعات زمانی

مشاهده یک ماتریس I_t در یک زمان واحد برای درک کامل دینامیک ترافیک کافی نیست. برای رفع این مشکل، ما آخرین K مشاهده را وارد وضعیت می‌کنیم تا عامل بتواند الگوها و روندهای ترافیکی را بهتر تشخیص دهد. در نتیجه، وضعیت عامل در زمان t به‌صورت زیر تعریف می‌شود:

$$S_t \in \mathbb{R}^{K \times 8 \times 8} \quad (3)$$



شکل ۳) یک تقاطع سه راهی به همراه ماتریس متقاطع آن

در نهایت، شایان ذکر است که هنگامی که تعداد حرکت‌ها در یک تقاطع کمتر از ۸ باشد (برای مثال، یک تقاطع سه‌راهی)، از صفرپُرسی (Zero Padding) استفاده می‌شود. برای نمونه، شکل ۳ a یک تقاطع سه‌راهی رایج را نشان می‌دهد که تنها حرکت‌های شرق (E)، چرخش چپ شرق (EL)، غرب (W) و چرخش چپ جنوب (SL) فعال هستند. همان‌گونه که در شکل ۳ b دیده می‌شود، سطرهای مربوط به آن چهار حرکت بلااستفاده با صفر پر می‌شوند تا اندازه ماتریس، مشابه اندازه ماتریس در تقاطع چهارراهی باقی بماند.

عمل :

یک طراحی اقدام (Action) واقع‌بینانه و قابل اجرا باید ایمنی تمامی کاربران ترافیک را در نظر بگیرد. اگرچه طراحی «انتخاب فاز بعدی (choose next phase)» [۱۰ و ۱۲ و ۱۸] می‌تواند بازدهی تقاطع را به‌طور چشمگیری افزایش دهد، اما این روش توالی اصلی چراغ‌های راهنمایی را برهم می‌زند و در نتیجه، ایمنی رانندگان را تحت تأثیر قرار می‌دهد.

در مقابل، این پژوهش از طراحی اقدام «حفظ یا تغییر [۲۳ و ۲۵ و ۹] استفاده می‌کند. این طراحی به مفهوم چرخه‌ها وفادار می‌ماند و فازها را به‌صورت ترتیبی اجرا می‌کند (برای مثال: فاز ۱، فاز ۲، فاز ۳، فاز ۴، سپس دوباره فاز ۱، فاز ۲، و به همین ترتیب). عامل (agent) بر اساس وضعیت t تصمیم می‌گیرد که:

- فاز فعلی را حفظ کند، یا

- به فاز بعدی منتقل شود.

به دلیل وجود اطلاعات فاز فعلی I^{cg} و اطلاعات فاز بعدی I^{ng} در ماتریس تقاطع، در ادامه نشان داده می‌شود که این طراحی اقدام، مقیاس‌پذیری بسیار خوبی برای تقاطع‌هایی با طرح‌های سیگنال متفاوت ارائه می‌دهد.

پاداش:

منفی طول صف متوسط در هر حرکت q_i به‌عنوان پاداش انتخاب شده است. معیارهایی مانند زمان انتظار، زمان سفر و تأخیر استفاده نشده‌اند، زیرا به‌دست آوردن آن‌ها از تجهیزات واقعی تشخیص ترافیک عملی نیست. در نتیجه، تابع پاداش پیشنهادی به‌صورت زیر تعریف می‌شود:

$$r_t = -\left(\sum_{i=1}^8 q_i\right) - \frac{\mu}{\sigma + \varepsilon} \quad (4)$$

که در آن ε عددی کوچک است تا از تقسیم بر صفر جلوگیری شود. علاوه بر این، μ و σ به ترتیب میانگین و انحراف معیار اولین R مقدار پاداش هستند. این مقادیر به صورت زیر محاسبه می‌شوند:

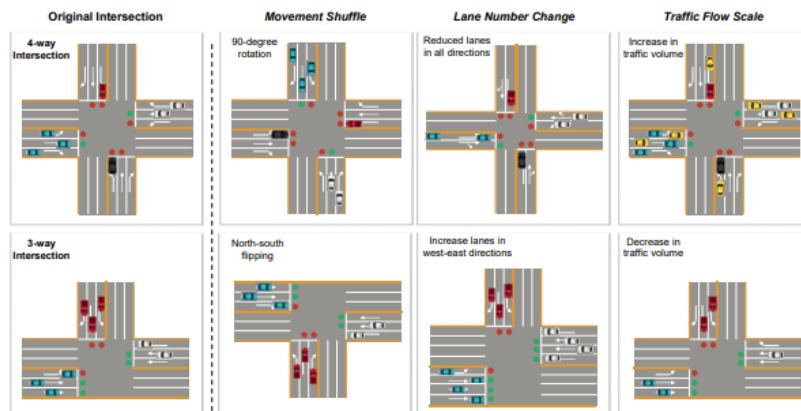
$$\mu = \frac{1}{R-1} \sum_{j=1}^{R-1} r_j \quad (5)$$

$$\sigma = \sqrt{\frac{1}{R-1} \sum_{j=1}^{R-1} (r_j - \mu)^2} \quad (6)$$

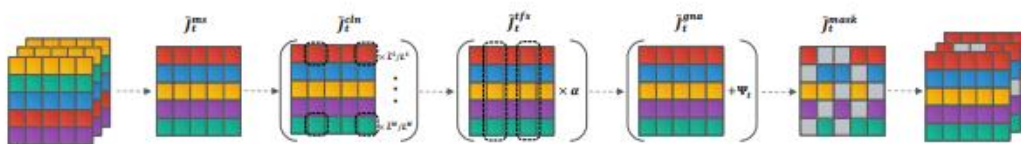
پاداش نرمال‌سازی می‌شود تا روند آموزش سریع‌تر و کارآمدتر گردد.

افزایش حالت ترافیکی:

در پژوهش‌های اخیر، تکنیک‌های افزایش داده اثربخشی خود را در بهبود توانایی تعمیم مدل‌های یادگیری تقویتی نشان داده‌اند [۳۸، ۳۹، ۴۰]. با آموزش عامل‌ها بر مجموعه‌ای متنوع‌تر از داده‌های افزایش‌یافته، مدل‌های یادگیری تقویتی قادر خواهند بود عملکرد بهتری در شرایط گوناگون و دیده‌نشده ارائه دهند.



شکل ۴) تصویرسازی اعمال سه روش افزایش وضعیت ترافیک که بر تقاطع‌های چهارراهی و سه‌راهی



شکل ۵) مراحل دقیق در بلوک افزایش وضعیت ترافیک

قابلیت آن‌ها در مواجهه با وظایف دیده‌نشده را افزایش می‌دهد. رایج‌ترین روش‌های افزایش داده که در ادبیات گزارش شده‌اند، شامل افزودن نویز گاوسی و ماسک‌گذاری هستند. در این پژوهش، سه روش نوآورانه دیگر برای افزایش وضعیت ترافیک، به‌طور خاص برای مسائل کنترل سیگنال ترافیکی مبتنی بر یادگیری تقویتی (RL-based TSC) می‌کنیم؛ یعنی جابجایی حرکت‌ها، تغییر تعداد و مقیاس‌بندی جریان ترافیک که در شکل ۴ نشان داده شده‌اند. شکل ۵ فرایند مرحله‌به‌مرحله اعمال شدن پنج روش افزایش داده را بر روی وضعیت ترافیکی K_t نشان می‌دهد که در نهایت به وضعیت افزوده‌شده نهایی \tilde{K}_t منتهی می‌شود. در طول آموزش، یک مینی‌بچ از داده‌ها به‌طور تصادفی از حافظهٔ بازپخش یا از مسیرهای افزوده‌شده اخیر نمونه‌برداری می‌شود. اگرچه افزایش داده به‌صورت تصادفی روی مینی‌بچ اعمال می‌شود، اما در میان فریم‌های پشت‌شده یکسان باقی می‌ماند. همچنین شایان ذکر است که این پنج روش افزایش وضعیت ترافیکی می‌توانند مستقیماً بر ماتریس تقاطع /اعمال شوند، که باعث می‌شود عامل بتواند سناریوها و ساختارهای مختلف تقاطع را بهتر یاد بگیرد و تطبیق یابد.

جابجایی حرکت‌ها

در این روش، سطرهای ماتریس تقاطع جابجا می‌شوند تا چرخش‌ها یا تقارن‌های مختلف یک تقاطع شبیه‌سازی شوند. به‌صورت شهودی، جابجایی حرکت‌ها را می‌توان نوعی «چرخش مؤثر» تقاطع در نظر گرفت (طبق شکل ۴). فرض اصلی پشت این روش این است که اقدامی که عامل اتخاذ می‌کند نباید پس از چرخش تقاطع تغییر کند.

به‌صورت ریاضی، این عملیات به صورت زیر مدل‌سازی می‌شود:

$$\tilde{J}_t^{ms} = P \cdot J_t, \quad (7)$$

که در آن، J_t ماتریس اصلی و \tilde{J}_t^{ms} ماتریس افزوده شده هستند. ماتریس P یک ماتریس جایگشتی است که سطرهای J_t را جابجا می کند.

با اعمال «جابجایی حرکتها» بر تمام ماتریسهای موجود در S_t ، حالت جدید به صورت زیر تعریف می شود:

$$\tilde{S}_t^{ms} = [\tilde{J}_{t-K+1}^{ms}, \tilde{J}_{t-K+2}^{ms}, \dots, \tilde{J}_t^{ms}]. \quad (8)$$

تغییر تعداد خطوط

برای قرار دادن عامل در معرض ترکیبهای متنوع تری از ساختارهای جاده، تعداد خطوط L_i در هر بردار حرکت \tilde{m}_i^{ms} به طور تصادفی تغییر داده می شود. این روش باعث می شود عامل در طول آموزش، با پیکربندیهای گوناگون خطوط آشنا شود و توانایی تطبیق با تقاطعهای مختلف را افزایش دهد. همچنین، ویژگیهای ترافیکی مانند جریان و اشغال، متناسب با ضریب برجای ماندن نسبتها مقیاس می شوند.

مثالی از این روش در شکل ۴ آمده است: در یک تقاطع ۴ راهی، تعداد خطوط هر جهت به ۲ کاهش یافته است؛ و در یک تقاطع ۳ راهی تعداد خطوط ورودی شرق و غرب به ۴ افزایش یافته است. در این فرایند، نسبت تعداد وسایل نقلیه در هر جهت حفظ می شود؛ بنابراین انتظار می رود تصمیم عامل پیش و پس از این تغییر یکسان بماند.

عملیات به صورت زیر تعریف می شود:

$$\tilde{m}_i^{cln} = f(\tilde{m}_i^{ms}, \tilde{L}_i), i = 1, 2, \dots, 8 \quad (9)$$

که در آن:

- L_i و \tilde{L}_i متغیرهای تصادفی یکنواخت هستند که تعداد خطوط اصلی و اصلاح شده را نشان می دهند،
- تابع f وظیفه تنظیم ویژگیهای ترافیکی را دارد.

تابع f به صورت زیر تعریف می شود:

$$f(\tilde{m}_i^{ms}, \tilde{L}_i) = \begin{cases} [\tilde{m}_i^{ms}]_k \cdot \frac{\tilde{L}_i}{L_i}, & k = 1, 2, 3, 5 \\ [\tilde{m}_i^{ms}]_k, & \text{otherwise} \end{cases} \quad (10)$$

- که در آن $[\cdot]_k$ ورودی k ام بردار را نشان می دهد.

- پس از اعمال این روش، حالت جدید برابر است با:

$$\tilde{S}_t^{cln} = [\tilde{j}_{t-K+1}^{cln}, \dots, \tilde{j}_t^{cln}] \quad (11)$$

که

$$\tilde{j}_t^{cln} = [\tilde{m}_1^{cln}, \tilde{m}_2^{cln}, \dots, \tilde{m}_8^{cln}]^T. \quad (12)$$

مقیاس بندی جریان ترافیک

- برای تغییر تمرکز عامل از مقادیر مطلق ترافیک به توزیع نسبی وسایل نقلیه، یک ضریب مقیاس بندی جریان معرفی می شود. در این روش، مقدار جریان و اشغال هر حرکت در ماتریس تقاطع در یک عدد تصادفی یکنواخت α ضرب می شود.

- عدد α برای تمام حرکت های یک حالت ثابت باقی می ماند تا نسبت های بین حرکت ها حفظ شوند.
- این روش به عامل کمک می کند تا اهمیت نسبی حرکت ها را یاد بگیرد و کمتر به مقدارهای مطلق وابسته باشد. لازم به ذکر است که این روش تعداد خطوط را تغییر نمی دهد.

- عملیات به صورت زیر تعریف می شود:

$$\tilde{m}_i^{tfs} = g(\tilde{m}_i^{cln}, \alpha) \quad (13)$$

که تابع g به صورت زیر تعریف می شود:

$$g(\tilde{m}_i^{cln}, \alpha) = \begin{cases} [\tilde{m}_i^{cln}]_k \cdot \alpha, & k = 1, 2, 3 \\ [\tilde{m}_i^{cln}]_k, & \text{otherwise} \end{cases} \quad (14)$$

حالت نهایی این مرحله \tilde{S}_t^{tfs} نامیده می شود.

افزودن نویز گاوسی:

برای افزودن تصادفی بودن به داده‌های آموزشی، نویز گاوسی مستقیماً به ماتریس تقاطع افزوده می‌شود. این نویز می‌تواند روی تمام اجزای ماتریس از جمله ویژگی‌های ترافیکی، ویژگی‌های حرکت و ویژگی‌های سیگنال تأثیر بگذارد. این روش باعث می‌شود عامل در برابر شرایط ناواضح و نامطمئن مقاوم‌تر شود.

عملیات به صورت زیر است:

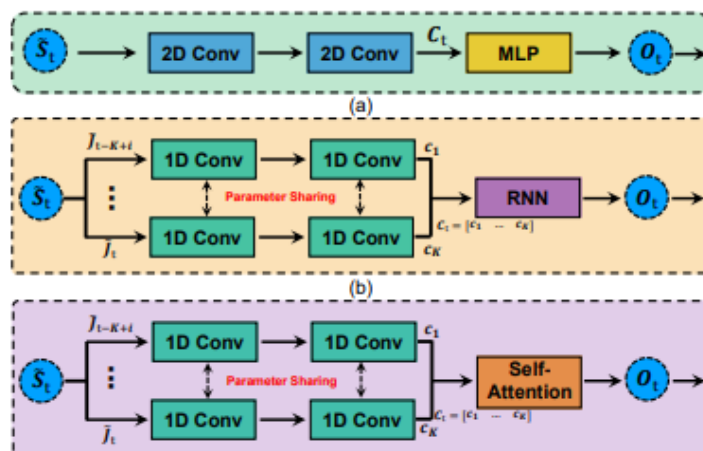
$$\tilde{J}_t^{gna} = \tilde{J}_t^{tfs} + \Psi_t, \quad (15)$$

که در آن $\Psi_t \sim N(0, I)$ نویز گاوسی است.

پس از این مرحله، حالت جدید \tilde{S}_t^{gna} نامیده می‌شود.

نقاب گذاری:

برای تشویق عامل به یادگیری تغییرات جریان ترافیک، به طور تصادفی برخی از مقدارهای موجود در ماتریس تقاطع در یک لحظه زمانی خاص صفر می‌شوند. ماسک گذاری عامل را مجبور می‌کند که برای استنتاج دینامیک‌های ترافیک، به اطلاعات قبل و بعد از بازه ماسک شده تکیه کند. این روش توانایی عامل را در درک نوسانات ترافیکی افزایش می‌دهد.



شکل ۶) سه نوع بلوک استخراج ویژگی تقاطع (الف) ساختار مبتنی بر CNN؛ (ب) ساختار مبتنی بر RNN؛ (ج) ساختار مبتنی بر ترانسفورماتور

۴-۴ استخراج ویژگی‌های تقاطع

در این بخش سه ساختار شبکه عصبی برای استخراج اطلاعات تقاطع از حالات ترافیکی افزوده شده به کار گرفته می‌شوند، که عبارت‌اند از: شبکه کانولوشنی چندبعدی (Convolutional Neural Network — CNN)، شبکه بازگشتی (Recurrent Neural Network — RNN) و ترنسفورمرها (Transformers).

ساختار مبتنی بر CNN

همان‌طور که در شکل ۶ (a) نشان داده شده است، از یک ساختار مبتنی بر CNN مجهز به دو لایه کانولوشن دوبعدی استفاده می‌شود تا اطلاعات سری‌زمانی موجود درون یک گره جاده استخراج شود. نمایش پنهان حاصل را C_t می‌نامیم که به صورت زیر تعریف می‌شود:

$$C_t = \text{ReLU} \left[(W_2^{2d} W_1^{2d} \tilde{S}_t) \right], \quad (16)$$

که در آن W_1^{2d} و W_2^{2d} پارامترهای قابل آموزش دو بلوک کانولوشن دوبعدی هستند. به طور مشخص، بلوک اول W_1^{2d} اطلاعات مربوط به حرکت‌های ترافیکی را استخراج می‌کند و سپس بلوک دوم W_2^{2d} بر اساس آن، اطلاعات خاص تقاطع را می‌گیرد. همچنین $\text{ReLU}(\cdot)$ نشان‌دهنده تابع فعال‌سازی ReLU است.

سپس C_t از طریق لایه (های) چندلایه‌ای پرسپترون (MLP) عبور داده شده و بردار ویژگی O_t تولید می‌شود:

$$O_t = W_u C_t + b_u, \quad (17)$$

که W_u و b_u پارامترهای قابل آموزش لایه MLP هستند.

ساختار مبتنی بر RNN

بر خلاف روش مبتنی بر CNN، ساختار مبتنی بر RNN که در شکل ۶ (b) نشان داده شده از یک لایه کانولوشن ۱D با به اشتراک گذاری پارامترها برای استخراج اطلاعات از هر ماتریس تقاطع استفاده می‌کند. لایه کانولوشن یکبعدی بر هر ماتریس تقاطع \tilde{A}_{t-K+i} برای $i = 1, 2, \dots, K$ اعمال می‌شود که این ماتریس همان حالت

افزوده شده در یک گام زمانی خاص در پنجره تاریخی است. خروجی های حاصل که با c_1, c_2, \dots, c_K نشان داده می شوند، اطلاعات تقاطع را در K زمان مختلف ضبط می کنند و به صورت زیر تعریف می گردند:

$$c_i = \text{ReLU} \left(W_2^{1d} W_1^{1d} \tilde{J}_{t-K+i} \right), i = 1, 2, \dots, K, \quad (18)$$

که در آن W_1^{1d} و W_2^{1d} پارامترهای قابل آموزش دو لایه کانولوشن یک بعدی هستند. این خروجی ها سپس به ماژول RNN خورنده می شوند:

$$h_i = \tanh (W_x c_i + h_{i-1} W_h + b_h), i = 1, 2, \dots, K, \quad (19)$$

که W_x, W_h, b_h وزن ها و بایاس لایه مخفی در ساختار RNN هستند. حالت مخفی نهایی h_K که از خروجی آخرین گام RNN به دست می آید برای محاسبه ویژگی کل تقاطع در بازه زمانی مورد نظر به کار می رود:

$$O_t = W_v h_K + b_v, \quad (20)$$

که W_v و b_v وزن ها و بایاس لایه خروجی در ساختار RNN می باشند.

ساختار مبتنی بر ترنسفورمر (Transformer)

شکل ۶ (c) روش مبتنی بر ترنسفورمر را نمایش می دهد. در اینجا، ابتدا از یک شبکه CNN با به اشتراک گذاری وزن ها برای استخراج ویژگی از ماتریس تقاطع در هر گام زمانی استفاده می شود (مشابه معادله (18)). اما به جای بلوک RNN، از یک رمزگذار ترنسفورمر برای آشکارسازی وابستگی های زمانی بین توالی ویژگی ها استفاده می شود. توالی ویژگی ها را $\tilde{C}_t = [c_1, c_2, \dots, c_K]$ تعریف می کنیم. برای وارد کردن اطلاعات زمان بندی، یک تعبیه قابل آموزش با نماد c_{class} به \tilde{C}_t اضافه می کنیم به طوری که:

$$C_t = \langle c_{\text{class}}, \tilde{C}_t \rangle = [c_{\text{class}}, c_1, c_2, \dots, c_K]. \quad (21)$$

خروجی رمزگذار ترنسفورمر بر پایه این توالی ورودی اصلاح شده، به عنوان نمایش حالت ترافیکی O_t به کار می رود. در بلوک رمزگذار ترنسفورمر، مکانیزم Self-attention به صورت زیر محاسبه می شود:

ابتدا پروجکشن های Query، Key و Value را محاسبه می کنیم:

$$Q_C = C_t W_Q, K_C = C_t W_K, V_C = C_t W_V, \quad (22)$$

و سپس خروجی attention را به صورت:

$$Z_t = \varphi \left(\frac{Q_C K_C^T}{\sqrt{C_d}} \right) V_C, \quad (23)$$

که در آن W_Q, W_K, W_V پارامترهای پروجکشن مربوطه هستند، C_d بعدی است که برای نرمال سازی ضرب در دنباله attention استفاده می شود و $\varphi(\cdot)$ عملگر SoftMax است. خروجی نهایی O_t بردار نخست Z_t معمولاً مرتبط با توکن c_{class} (می باشد).

۵-۴ آموزش RL و فاین تیون

برای آموزش مدل UniTSA از الگوریتم Proximal Policy Optimization (PPO) استفاده می شود. همان طور که در شکل ۲ نشان داده شده، عامل در تعامل با سناریوهای ترافیکی مختلف از تقاطع های با ساختارهای گوناگون، مسیرهایی متشکل از مشاهدات، اقدامات و پاداش ها را جمع آوری می کند. این مسیرها (trajectories) مبنای محاسبه تابع های هزینه سیاست (policy loss) و ارزش (value loss) هستند که برای به روز رسانی وزن های شبکه های Actor و Critic به کار می روند.

تابع هزینه سیاست (Policy Loss)

تابع هزینه سیاست اختلاف بین سیاست فعلی و سیاست به روز رسانی شده بر پایه ۰ مسیرهای جمع آوری شده را اندازه گیری می کند. این تابع عامل را تشویق می کند تا احتمال اقداماتی را که منجر به پاداش های بالاتر می شوند افزایش دهد و احتمال اقدامات منفی را کاهش دهد. به صورت ریاضی تابع هزینه ۰ سیاست چنین فرموله می شود:

$$L_{pf}(\theta) = \mathbb{E}_t[\min(r_t A_t, \text{clip}(r_t, 1 - \epsilon, 1 + \epsilon) A_t)], \quad (24)$$

که \mathbb{E}_t نشان دهنده اپراتور امید ریاضی تخمینی است و r_t نسبت میان سیاست جدید $\pi_\theta(a_t | s_t)$ و سیاست قدیمی $\pi_{\theta_{old}}(a_t | s_t)$ است:

$$r_t = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)}. \quad (25)$$

علاوه بر این، $A_t = r_t + \gamma V(s_{t+1}) - V(s_t)$ تابع مزیت است و استفاده از تابع $\text{clip}(\cdot)$ موجب به‌روزرسانی‌های پایدارتر در سیاست می‌شود.

تابع هزینه ارزش (Value Loss)

تابع هزینه ارزش اختلاف بین مقدار تخمینی تابع ارزش و پاداش‌های واقعی به‌دست‌آمده را اندازه می‌گیرد تا تابع ارزش بهتر بتواند مجموع پاداش‌های آتی مورد انتظار را تقریب بزند. این تابع به‌صورت زیر تعریف می‌شود:

$$L_{vf}(\theta) = \mathbb{E}_t[(V_{\theta}(s_t) - \hat{R}_t)^2], \quad (26)$$

که در آن $V_{\theta}(s_t)$ مقدار تخمینی تابع ارزش تحت سیاست θ است و $\gamma^k r_{t+k}$ و $\hat{R}_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$ نمایانگر «پاداش تا سرانجام (reward-to-go)» می‌باشد.

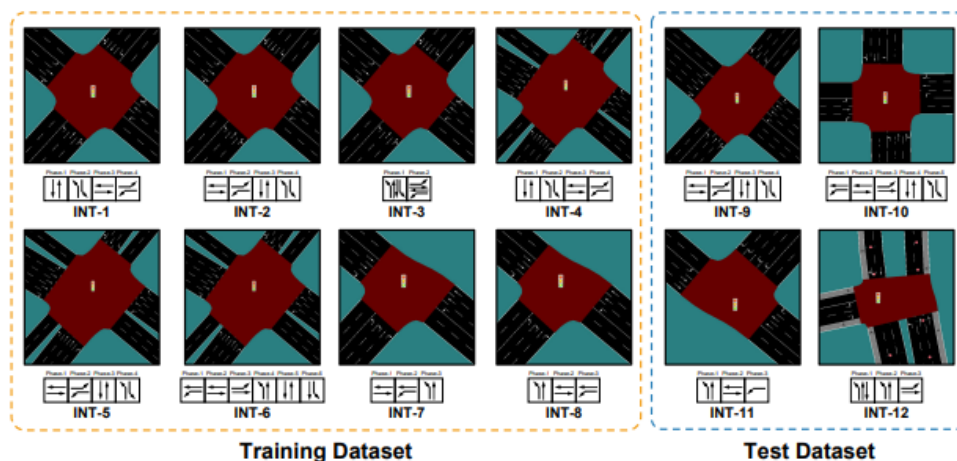
تابع هدف نهایی

پس از محاسبه تابع هزینه سیاست و تابع هزینه ارزش، تابع هدف نهایی به‌صورت زیر تعریف می‌شود:

$$L(\theta) = -L_{pf}(\theta) + \lambda L_{vf}(\theta), \quad (27)$$

که در آن λ ضریب وزن‌دهی به خسارت ارزش است.

این پژوهش یک مدل جهانی مؤثر برای کنترل سیگنال‌های ترافیکی پیشنهاد می‌کند که با کمینه‌سازی تابع هدف فوق از طریق PPO بهینه می‌گردد. علاوه بر این، از آن‌جا که برخی تقاطع‌ها در عمل از اهمیت بیشتری برخوردارند، برای فاین‌تیون عملکرد روی این تقاطع‌های حیاتی از روش LoRA استفاده می‌شود. به‌طور مشخص، ماژول‌های LoRA به وزن‌های لایه‌های چگال (dense) در شبکه Actor و Critic افزوده می‌شوند (همان‌طور که در شکل ۲ نشان داده شده است). در فرایند فاین‌تیون، وزن‌های پیش‌آموزش‌دیده اصلی ثابت نگه داشته می‌شوند و تنها ماژول‌های LoRA به‌روزرسانی می‌گردند. این طراحی اجازه می‌دهد مدل بدون افزایش قابل توجه در تعداد پارامترها به ویژگی‌های ویژه هر تقاطع سازگار شود و در عین حال کارایی آموزش حفظ شود.



شکل ۷) تمام توپولوژی‌های تقاطع به همراه فازهای موجودشان که برای آموزش و آزمایش استفاده شده‌اند

۵) آزمایش‌ها

۵-۱ تنظیمات آزمایش

برای ارزیابی مدل پیشنهادی، آزمایش‌های گسترده‌ای با استفاده از بسته نرم‌افزاری SUMO انجام شد. SUMO یک شبیه‌ساز ترافیک میکروسکوپی متن‌باز است که برای مدیریت شبکه‌های بزرگ طراحی شده است. این نرم‌افزار رابط Traffic Control Interface (TraCI) را ارائه می‌دهد که امکان کنترل چراغ‌های ترافیکی و دریافت اطلاعات شرایط ترافیکی تقاطع‌ها را فراهم می‌کند. جریان ترافیک و اشغال (occupancy) هر حرکت با تحلیل موقعیت‌ها و مسیرهای خودروها محاسبه می‌شود. برای شبیه‌سازی شرایط واقعی، تنها خودروهایی که در فاصله ۱۵۰ متر از تقاطع قرار دارند، در نظر گرفته شده‌اند. همچنین پس از چراغ سبز، چراغ زرد به مدت ۳ ثانیه نمایش داده شد تا ایمنی راننده حفظ شود. زمان انتظار هر خودرو به عنوان معیار عملکرد جهت ارزیابی اثربخشی روش‌های مختلف استفاده شده است؛ مقدار کمتر این معیار نشان‌دهنده عبور سریع‌تر خودروها از تقاطع است. برای آموزش مدل، از پیاده‌سازی الگوریتم Proximal Policy Optimization (PPO) ارائه شده توسط کتابخانه Stable Baselines3 استفاده شد. برای افزایش سرعت آموزش، از ۳۰ پردازش موازی بهره گرفته شد و تعداد کل گام‌های محیط آموزش برابر با ۱۰ میلیون تعیین گردید. نمایش حالت شامل ۸ نمونه قبلی ماتریس تقاطع بود و فاصله زمانی بین دو اقدام متوالی ۵ ثانیه در نظر گرفته شد.

جدول ۱) تنظیمات ابر پارامترها

Hyper-parameter	Value
Learning rate	0.0001
Trajectory memory size	3000
Clipping range ϵ	0.2
Discount factor γ	0.99
Value function coefficient λ	0.9
Scale hyperparameter α	1
Rank of LoRA module	8

ابر پارامترها نیز مطابق جدول ۱ تنظیم شدند. شبکه‌های Actor و Critic به صورت دو لایه کاملاً متصل طراحی شدند. ابعاد ورودی $\{۳۲, ۶۴\}$ و ابعاد خروجی $\{۲, ۳۲\}$ برای Actor و $\{۱, ۳۲\}$ برای Critic تنظیم شد.

جدول ۲) تمام تنظیمات تقاطع‌ها

Intersection ID	Training Dataset								Test Dataset			
	INT-1	INT-2	INT-3	INT-4	INT-5	INT-6	INT-7	INT-8	INT-9	INT-10	INT-11	INT-12
roads	4	4	4	4	4	4	3	3	4	4	3	3
lanes per road	(3,3,3,3)	(3,3,3,3)	(3,3,3,3)	(3,4,4,5)	(3,4,4,5)	(3,4,4,5)	(3,3,3)	(3,3,3)	(3,4,3,4)	(3,3,3,3)	(4,3,3)	(2,3,2)
phases	4	4	2	4	4	6	3	3	4	5	3	3

۲-۵ مجموعه داده‌ها

این مطالعه تقاطع‌هایی با ساختارهای متنوع در نظر گرفته است تا یک مدل جهانی واحد برای پیش‌بینی اقدامات در تمامی تقاطع‌ها آموزش داده شود. به‌طور مشخص، ۱۲ تقاطع با تعداد فازها، خطوط عبور و جاده‌های ورودی متفاوت (۳-راه یا ۴-راه) ساخته شده و برای آزمایش استفاده شد. از این ۱۲ تقاطع، هشت تقاطع برای آموزش و چهار تقاطع برای آزمون کنار گذاشته شدند. توپولوژی و فازهای ۱۲ تقاطع در شکل ۷ نشان داده شده و جدول ۲ خلاصهٔ پیکربندی‌های آن‌ها را ارائه می‌دهد. برای مثال، INT-4 شامل چهار جاده دوطرفه است با سه خط در جهت شمال-جنوب، پنج خط در جهت غرب-شرق و چهار خط در هر یک از دو جهت دیگر INT-4. شامل چهار فاز است که هر فاز دو سیگنال حرکت مختلف را ترکیب می‌کند. بنابراین، پیکربندی INT-4 در جدول ۲ شامل چهار جاده با تعداد خطوط مشخص (۳, ۴, ۴, ۵) به ترتیب جهت عقربه‌های ساعت و چهار فاز می‌باشد.

در مجموعه داده آموزش، سه توپولوژی متفاوت وجود دارد:

- INT-1، INT-2، و INT-3 تقاطع ۴-راه معمولی با هر جاده شامل سه خط.
 - INT-4، INT-5، و INT-6 تقاطع ۴-راه بزرگ با بیش از چهار خط در هر جاده.
 - INT-7، INT-8، و INT-9 تقاطع ۳-راه.
- برای تقاطع‌های با توپولوژی مشابه، تقاطع‌های جدید با تغییر ترتیب یا تعداد فازها ساخته شد. به‌طور مثال، INT-1 و INT-2 دارای پیکربندی یکسان هستند اما ترتیب فازها متفاوت است. همچنین INT-1 و INT-3 در تعداد فازها متفاوت هستند.
- برای ارزیابی عملکرد مدل در تقاطع‌های دیده‌نشده، چهار سناریوی آزمون شامل INT-9، INT-10، INT-11، و INT-12 طراحی شد.
- INT-9 و INT-10: تقاطع‌های ۴-راه با پیکربندی خطوط و فازهای متفاوت نسبت به مجموعه آموزش.
 - INT-11: تغییر خطوط و فازها بر اساس INT-8.
 - INT-12: شبیه‌سازی ترافیک واقعی شهر Ingolstadt، آلمان [43].
- علاوه بر تنوع ساختار، برای هر تقاطع ۱۰۰ مسیر منحصر به فرد تولید شد که سه چهارم آن‌ها برای آموزش و یک چهارم برای ارزیابی استفاده شد. هر مسیر دارای مدت ۳۰,۰۰۰ ثانیه (~ ۸ ساعت) بود.

5.3 روش‌های مقایسه‌ای (Compared Methods)

برای ارزیابی عملکرد UniTSA، مدل نهایی با روش‌های کلاسیک و پیشرفته RL برای کنترل ترافیک مقایسه شد:

- **Fix Time [2]** کنترل زمان ثابت با چرخه و فازهای از پیش تعیین‌شده. دو نسخه FixTime-30 و FixTime-40 هر فاز ۳۰ و ۴۰ ثانیه
- **Webster [3]** تعیین طول چرخه و تقسیم فاز بر اساس حجم ترافیک در یک بازه زمانی. می‌تواند زمان سفر را کاهش یا ظرفیت تقاطع را افزایش دهد. برای عدالت، در این مطالعه بر اساس ترافیک لحظه‌ای اعمال شد.

- **[4] SOTL** کنترل سیگنال خودسازمانده که طول سیگنالها را بر اساس آستانه تعداد خودروهای منتظر تنظیم می کند.
- **[10] MPLight** استفاده از ساختار FRAP و مکانیسم پاداش مبتنی بر فشار، همچنین MLP با اشتراک پارامتر برای افزایش سازگاری با توپولوژی های مختلف.
- **[18] AttendLight** مدل توجه برای آموزش مدل جهانی که قادر به مدیریت تقاطع های با ساختار و جریان متفاوت است. شامل دو مدل توجه: اول برای تعداد جاده ها و خطوط و دوم برای تصمیم گیری بین تقاطع ها با تعداد فاز متفاوت.

روش های متغیر: UniTSA

- **UniTSA (Single)** آموزش در یک محیط، استفاده از ساختار RNN برای استخراج اطلاعات حالت ترافیک.
- **UniTSA (Multi)** آموزش همزمان در چند محیط با ساختارهای مختلف شبکه CNN، RNN، (Transformer).
- **UniTSA (Multi+TSA)** ترکیب پنج روش تقویت حالت ترافیک (Traffic State Augmentation) با UniTSA (Multi) برای افزایش تنوع داده ها و بهبود عملکرد.

جدول ۳) نتایج کمی (میانگین زمان انتظار به ازای هر وسیله نقلیه) تقاطع های آموزشی برای مدل های جهانی. مقدار کمتر نشان دهنده عملکرد بهتر است و کمترین مقادیر با حروف پررنگ مشخص شده اند.

	INT-1	INT-2	INT-3	INT-4	INT-5	INT-6	INT-7	INT-8
Fix-30 [2]	39.458	38.862	8.323	35.123	35.031	51.923	18.920	18.595
Fix-40 [2]	50.696	52.108	10.667	44.888	45.540	61.743	23.411	24.759
Webster [3]	26.466	26.889	5.751	25.413	24.541	40.128	12.755	13.574
SOTL [4]	16.048	16.561	2.764	28.169	27.902	27.404	8.639	7.925
MPLight [10]	19.111	14.659	4.469	16.067	19.925	19.115	6.654	7.523
AttendLight [18]	16.483	13.893	3.860	16.903	18.915	20.795	6.532	8.104
UniTSA (Multi+CNN)	13.776	13.580	3.265	14.790	15.437	18.751	6.592	6.894
UniTSA (Multi+RNN)	13.692	13.437	3.495	15.198	14.896	18.612	6.393	6.625
UniTSA (Multi+Trans)	14.976	20.459	2.811	16.489	16.481	23.793	7.552	7.175
UniTSA (Multi+CNN+TSA)	14.071	14.150	3.036	15.990	16.472	24.383	6.329	6.328
UniTSA (Multi+RNN+TSA)	13.450	13.578	3.007	14.470	14.462	18.689	6.242	6.164
UniTSA (Multi+Trans+TSA)	13.335	13.314	3.311	14.648	14.566	18.579	6.643	6.571

جدول ۴) نتایج کمی تقاطع‌های آزمون برای مدل‌های جهانی

	INT-9	INT-10	INT-11	INT-12
Fix-30	40.341	38.360	17.136	17.440
Fix-40	60.043	50.580	23.504	20.570
Webster	26.191	27.766	11.981	13.280
SOTL	23.031	28.066	7.452	8.070
MPLight	23.674	21.475	8.447	15.047
AttendLight	18.080	18.501	7.393	12.982
UniTSA (Multi+CNN)	17.877	18.244	7.063	12.820
UniTSA (Multi+RNN)	16.771	15.450	6.807	10.140
UniTSA (Multi+Trans)	17.640	16.269	6.598	9.630
UniTSA (Multi+CNN+TSA)	13.075	14.245	6.794	7.550
UniTSA (Multi+RNN+TSA)	12.677	14.208	5.914	6.730
UniTSA (Multi+Trans+TSA)	13.054	16.186	5.983	6.190

۴-۵ نتایج تقاطع‌های آموزشی

عملکرد UniTSA با روش‌های Fix Time، Webster، SOTL، MPLight و AttendLight بر روی تقاطع‌های آموزش مقایسه شد. نتایج متوسط زمان انتظار هر خودرو در جدول ۳ ارائه شده است.

- روش‌های مبتنی بر RL در اکثر تقاطع‌ها عملکرد بهتری نسبت به روش‌های کلاسیک نشان دادند.
- SOTL در INT-2 کمترین زمان انتظار را داشت اما نیاز به آستانه‌های دستی داشت که تعمیم‌پذیری آن را محدود می‌کند.
- بین روش‌های RL جهانی، UniTSA نسبت به MPLight و AttendLight بهبود قابل توجهی داشت؛ به‌طور میانگین، بهبود ۱۵٪ و ۱۲٪ نسبت به MPLight و AttendLight در هشت تقاطع مشاهده شد.

توضیح عملکرد:

- UniTSA علاوه بر اشتراک پارامتر، با جایگزینی بلوک RNN یا Transformer به جای MLP، اطلاعات زمانی وضعیت ترافیک را استخراج می‌کند.
- پنج روش تقویت حالت ترافیک به عامل اجازه می‌دهد در طول آموزش با تنوع بیشتری از وضعیت‌های تقاطع مواجه شود.
- بررسی ساختار شبکه‌ها نشان می‌دهد که RNN نسبت به CNN عملکرد بهتری دارد.

- در غیاب تقویت حالت ترافیک، UniTSA (Multi+CNN) و UniTSA (Multi+RNN) عملکرد بهتری نسبت به UniTSA (Multi+Trans) داشتند، زیرا ترنسفورمر به داده‌های بیشتری نیاز دارد.
- با اعمال تقویت حالت ترافیک، مدل می‌تواند با تقاطع‌های متنوع‌تری تعامل داشته باشد و UniTSA (Multi+Trans+TSA) در بسیاری از سناریوها از UniTSA (Multi+RNN) و UniTSA (Multi+RNN+TSA) عملکرد بهتری نشان داد.

فرمول LoRA در UniTSA

اگر ماتریس وزن پیش‌آموزش‌دیده $W \in \mathbb{R}^{n \times m}$ با ماژول $\Delta W = W_A W_B^T$ LoRA داشته باشیم، با $W_B \in \mathbb{R}^{m \times d}$ ، $W_A \in \mathbb{R}^{n \times d}$ و $d \ll n$ خروجی لایه به شکل زیر است:

$$z = Wx + \Delta Wx = Wx + \alpha r W_A W_B^T x, \quad (28)$$

که W_B و W_A به ترتیب با ماتریس صفر و ماتریس گوسی با میانگین صفر مقداردهی اولیه شده‌اند، α ضریب مقیاس و r رتبه ماژول LoRA است. با ترکیب آموزش RL با PPO و فاین تیون LoRA، مدل جهانی پیشنهادی قادر است چالش‌های تقاطع‌های مهم با ساختارهای متنوع را به‌طور مؤثر مدیریت کند.

۵-۵ نتایج تقاطع‌های آزمون

در مرحله بعد، ویژگی کلیدی UniTSA مورد ارزیابی قرار گرفت تا بررسی شود آیا می‌تواند برای تقاطع‌های دیده‌نشده (unseen intersections) نیز مورد استفاده قرار گیرد. چهار تقاطع شامل INT-9، INT-10، INT-11 و INT-12 به‌طور ویژه برای اهداف آزمایشی انتخاب شدند. همان‌طور که در بخش ۵.۲ بیان شد، این تقاطع‌ها از نظر تعداد خطوط و تعداد فازها با تقاطع‌های مجموعه آموزش متفاوت هستند.

در زمینه کاهش زمان سفر متوسط خودروها در عبور از تقاطع‌ها، مدل UniTSA (Multi+RNN+TSA) نشان داد که زمان سفر متوسط خودروها به ترتیب تقریباً ۳۲.۹٪ و ۴۱.۳٪ نسبت به MPLight و AttendLight کاهش یافته است. این نتایج نشان‌دهنده اثربخشی روش پیشنهادی در بهینه‌سازی کنترل چراغ‌های ترافیکی (TSC) و بهبود کارایی جریان ترافیک می‌باشد.

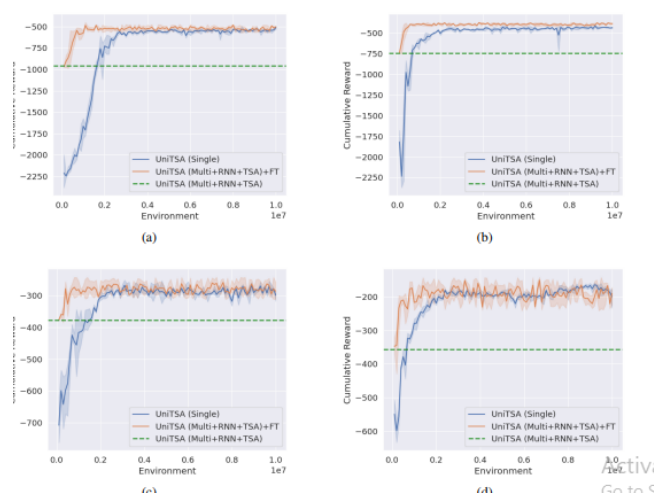
قابل توجه است که استفاده از تکنیک‌های تقویت حالت ترافیک موجب بهبود عملکرد بین واریانت‌های مختلف UniTSA می‌شود. به‌طور مثال:

- UniTSA (Multi+RNN+TSA)
- UniTSA (Multi+RNN+TSA)
- UniTSA (Multi+Trans+TSA)

به ترتیب بهبودهای ۲۳.۴٪، ۱۹.۸٪ و ۱۷.۹٪ را در مقایسه با UniTSA (Multi+RNN) ، UniTSA (Multi+RNN) و UniTSA (Multi+RNN) نشان دادند. این بهبود به دلیل افزوده شدن تنوع بیشتر در سناریوهای تقاطع در داده‌های آموزشی است که از طریق تکنیک‌های تقویت حالت ترافیک، مانند روش “تغییر تعداد خطوط ایجاد شده و امکان تولید ترکیب‌های متنوعی از پیکربندی خطوط را فراهم می‌کند.

جدول ۴ عملکرد روش‌های مختلف، شامل روش‌های پایه و واریانت‌های مختلف مدل UniTSA را جمع‌بندی می‌کند. مشابه نتایج به‌دست آمده در مجموعه آموزش، الگوریتم‌های مبتنی بر RL عملکرد قابل توجهی نسبت به روش‌های سنتی کنترل ترافیک نشان دادند.

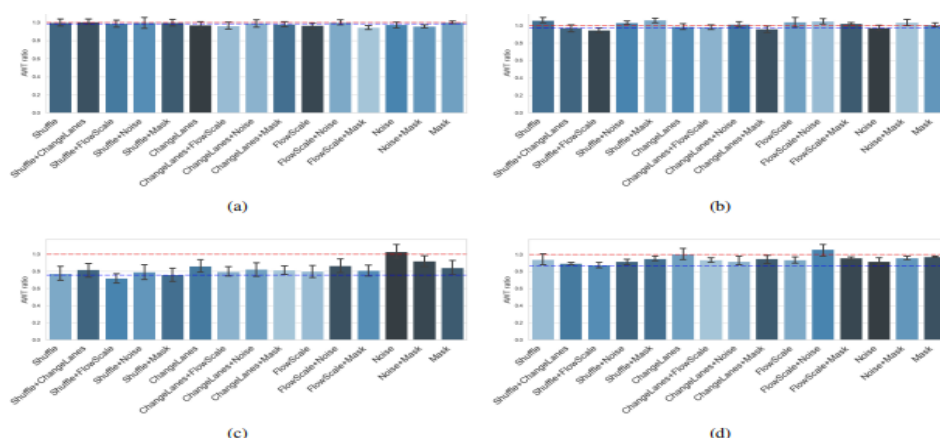
در بین روش‌های جهانی مبتنی بر RL، مدل‌های UniTSA در تمامی تقاطع‌های آزمون عملکرد برتری داشتند. به ویژه، UniTSA (Multi+RNN+TSA) و UniTSA (Multi+Trans+TSA) در میان تمامی روش‌ها بهترین عملکرد را ارائه کردند.



شکل ۸) مراحل محیطی روش‌های مختلف در چهار تقاطع آزمون. (الف) INT-9. (ب) INT-10. (ج) INT-11. (د) INT-12.

جدول ۵) نتایج کمی تنظیم دقیق در تقاطع‌های آزمایشی

	INT-9	INT-10	INT-11	INT-12
UniTSA (Multi+RNN)	16.771	15.450	6.807	10.140
UniTSA (Multi+RNN+TSA)	12.677	14.208	5.914	6.730
<i>1M Environment Steps</i>				
UniTSA (Single)	23.683	14.863	8.417	5.024
UniTSA (Multi+RNN+TSA) + FT	10.995	12.553	4.475	3.534
<i>10M Environment Steps</i>				
UniTSA (Single)	10.353	12.254	4.359	3.089
UniTSA (Multi+RNN+TSA) + FT	10.292	11.081	4.153	3.110



شکل ۹) تحلیل مقایسه‌ای روش‌های تقویت وضعیت ترافیک در تقاطع‌های آموزشی و آزمایشی انتخاب‌شده (الف) INT-1. (ب) INT-7. (ج) INT-9. (د) INT-11

۵-۶ نتایج فاین تیونینگ در تقاطع‌های آزمون

در سناریوهای عملی، برخی تقاطع‌ها به دلیل اهمیت‌شان نیاز به توجه ویژه دارند. برای رسیدگی به این موضوع، ما با مدل جهانی آموزش‌دیده توسط UniTSA (Multi+RNN) شروع می‌کنیم. ساختار RNN برای UniTSA انتخاب شد زیرا عملکرد بهتری در اغلب تقاطع‌ها نسبت به ساختارهای CNN و Transformer دارد. در مقایسه با UniTSA (Single)، که روی یک سناریوی واحد آموزش داده شده است، مدل حاصل می‌تواند به سرعت به عملکرد مدل تک‌محیطی برسد یا حتی آن را پس از تنها چند مرحله آموزش پشت سر بگذارد.

شکل ۸ تغییرات مجموع پاداش‌ها را طی مراحل آموزش برای مدل‌های مختلف در تقاطع‌های آزمون نشان می‌دهد. خط سبز نقطه‌چین نتایج اعمال مستقیم مدل جهانی بر تقاطع‌های جدید بدون هیچ فاین تیونینگ را نشان می‌دهد. قابل توجه است که مدل حتی بدون هرگونه فاین تیونینگ یا یادگیری انتقالی نیز نتایج امیدوارکننده‌ای ارائه می‌دهد.

خط آبی نشان‌دهنده مدل آموزش‌دیده از صفر و خط نارنجی نشان‌دهنده مدل فاین تیون شده بر اساس مدل جهانی است. مشاهده می‌شود که مدل تک‌محیطی تقریباً در ۳ میلیون مرحله آموزش همگرا می‌شود، در حالی که مدل فاین تیون شده با تنها حدود ۱ میلیون مرحله آموزش به عملکرد مشابه دست می‌یابد و در نتیجه حدود ۶۶٪ کاهش در زمان محاسبات با حفظ عملکرد مشابه به دست می‌آید.

جدول ۵ تحلیل دقیقی از عملکرد مدل پس از فاین تیونینگ ارائه می‌دهد. در ۱ میلیون مرحله آموزش، مدل فاین تیون شده بهبود عملکرد متوسط ۳۶٪ نسبت به UniTSA (Single) در چهار تقاطع آزمون نشان داد. حتی پس از ۱۰ میلیون مرحله آموزش، مدل‌های فاین تیون شده همچنان ۳٪ بهتر از مدل‌های آموزش‌دیده از صفر عمل کردند. این ویژگی برای کاربردهای بلادرنگ بسیار جذاب است، زیرا در شبکه‌های جاده‌ای با بیش از ۱۰۰۰ تقاطع می‌توان تعداد تعاملات با محیط را به طور قابل توجهی کاهش داد و در عین حال عملکرد مشابهی را حفظ کرد، که به طور چشمگیری کارایی آموزش را افزایش می‌دهد.

۷-۵ تحلیل مقایسه‌ای روش‌های تقویت حالت ترافیک

در این بخش، اثربخشی روش‌های تقویت حالت ترافیک در بهبود عملکرد مدل UniTSA (Multi+RNN) تحلیل شده است. آزمایش‌ها با استفاده از ترکیب‌های مختلف دو تکنیک تقویت حالت ترافیک روی تقاطع‌های آموزش و آزمون انجام شد. برای ارزیابی تأثیر این روش‌ها، نسبت زمان انتظار متوسط (AWT ratio) بین مدل‌های با و بدون تقویت حالت ترافیک محاسبه شد. مقدار کمتر از ۱ نشان‌دهنده بهبود عملکرد ناشی از روش‌های تقویت است.

شکل ۹ نتایج تحلیل مقایسه‌ای را برای تقاطع‌های INT-1، INT-7، INT-9 و INT-11، که نماینده ساختارهای معمول تقاطع در سناریوهای واقعی هستند، نشان می‌دهد. هر میله در نمودار مربوط به نسبت متوسط AWT برای ترکیب خاصی از روش‌های تقویت است و نوارهای خطا، فاصله اطمینان ۹۵٪ را نشان می‌دهند. خط آبی نقطه‌چین نسبت AWT مدل UniTSA (Multi+RNN+TSA) را که از تمامی روش‌های تقویت حالت ترافیک استفاده می‌کند، نشان می‌دهد.

شکل ۹ (a) و ۹ (b) نتایج به دست آمده از تقاطع‌های آموزش INT-1 و INT-7 را نشان می‌دهد. بررسی این نتایج نشان می‌دهد که بهبود متوسط عملکرد از طریق تقویت حالت ترافیک حدود ۲٪ بوده است. همچنین، استفاده از روش‌های Noise و Mask در INT-7 باعث کاهش عملکرد شد، که می‌تواند ناشی از این باشد که مدل قبلاً الگوها و ویژگی‌های پایه‌ای تقاطع‌های آموزش را تا حد زیادی یاد گرفته است و افزودن تغییرات اضافی از طریق تقویت حالت ترافیک ممکن است مزایای قابل توجهی نداشته باشد.

با این حال، روش‌های تقویت حالت ترافیک هنگام مواجهه با ساختارهای تقاطع دیده‌نشده عملکرد قابل توجهی ارائه دادند. شکل ۹ (c) و ۹ (d) نسبت‌های AWT برای تقاطع‌های آزمون INT-9 و INT-11 را نشان می‌دهد. نتایج به وضوح نشان می‌دهد که بیشتر روش‌های تقویت حالت ترافیک عملکرد سیاست پایه را در تقاطع‌های آزمون بهبود دادند. نمونه‌های متنوع آموزشی ایجاد شده از طریق تقویت حالت ترافیک باعث بهبود عملکرد مدل شدند. از میان تکنیک‌ها، Movement Shuffle و Traffic Flow Scale به‌طور ویژه‌ای مؤثر بودند، زیرا به مدل اجازه می‌دهند تا از طیف گسترده‌تری از سناریوها یاد بگیرد و خود را تطبیق دهد و در نتیجه عملکرد در مجموعه آزمون افزایش یابد.

۶ نتیجه‌گیری

در این مقاله، یک چارچوب کنترل چراغ ترافیک مبتنی بر RL به صورت جهانی (UniTSA) برای ساختارهای مختلف تقاطع در محیط‌های V2X پیشنهاد شد. به طور مشخص، UniTSA امکان آموزش یک عامل RL جهانی را با استفاده از ماتریس تقاطع برای توصیف وضعیت تقاطع فراهم می‌کند. برای مقابله با تقاطع‌های دیده‌نشده، روش‌های جدید تقویت حالت ترافیک پیشنهاد شد تا داده‌های جمع‌آوری شده توسط عامل غنی شود و عملکرد و قابلیت تعمیم مدل در تقاطع‌های ناشناخته بهبود یابد. نتیجه این است که UniTSA نیاز به سفارشی‌سازی و توسعه مجدد برای هر تقاطع منفرد را از بین می‌برد و در عین حال پیاده‌سازی ساده، کارآمد و متن‌باز ارائه می‌دهد، که آن را به چارچوبی ارزشمند برای تحقیقات آینده در حوزه کنترل چراغ ترافیک مبتنی بر RL با داده کارآمدی و قابلیت تعمیم بالا در محیط‌های V2X تبدیل می‌کند. نتایج تجربی گسترده نشان داد که UniTSA کوتاه‌ترین زمان انتظار متوسط را در میان انواع مختلف تقاطع‌ها به دست می‌آورد، عملکرد بهتری نسبت به روش‌های موجود ارائه می‌دهد و حتی از مدل‌های آموزش دیده از صفر با فاین تیونینگ پیشی می‌گیرد.

قدردانی (Acknowledgments)

این تحقیق توسط برنامه ملی کلیدی تحقیقات و توسعه چین تحت گرنت شماره 2020YFB1807700 و برنامه Shanghai Pujiang تحت گرنت شماره 21PJD092 حمایت شده است.

References

- [1] Fehda Malik, Hasan Ali Khattak, and Munam Ali Shah. Evaluation of the impact of traffic congestion based on sumo. In 2019 25th International Conference on Automation and Computing (ICAC), pages 1–5. IEEE, 2019.
- [2] Alan J Miller. Settings for fixed-cycle traffic signals. *Journal of the Operational Research Society*, 14(4):373–386, 1963.
- [3] Thomas Urbanik, Alison Tanaka, Bailey Lozner, Eric Lindstrom, Kevin Lee, Shaun Quayle, Scott Beaird, Shing Tsoi, Paul Ryus, Doug Gettman, et al. Signal timing manual, volume 1. Transportation Research Board Washington, DC, 2015.
- [4] Carlos Gershenson. Self-organizing traffic lights. arXiv preprint nlin/0411066, 2004.
- [5] Ishu Tomar, S Indu, and Neeta Pandey. Traffic signal control methods: Current status, challenges, and emerging trends. *Proceedings of Data Analytics and Management: ICDAM 2021, Volume 1*, pages 151–163, 2022.
- [6] Wang Tong, Azhar Hussain, Wang Xi Bo, and Sabita Maharjan. Artificial intelligence for vehicle-to-everything: A survey. *IEEE Access*, 7:10823–10843, 2019.
- [7] Tamás Wágner, Tamás Ormándi, Tamás Tettamanti, and István Varga. Spat/map v2x communication between traffic light and vehicles and a realization with digital twin. *Computers and Electrical Engineering*, 106:108560, 2023.
- [8] Yit Kwong Chin, Lai Kuan Lee, Nurmin Bolong, Soo Siang Yang, and Kenneth Tze Kin Teo. Exploring q-learning

optimization in traffic signal timing plan management. In 2011 third international conference on computational

intelligence, communication systems and networks, pages 269–274. IEEE, 2011.

[9] Guanjie Zheng, Yuanhao Xiong, Xinshi Zang, Jie Feng, Hua Wei, Huichu Zhang, Yong Li, Kai Xu, and Zhenhui

Li. Learning phase competition for traffic signal control. In Proceedings of the 28th ACM international conference

on information and knowledge management, pages 1963–1972, 2019.

[10] Chacha Chen, Hua Wei, Nan Xu, Guanjie Zheng, Ming Yang, Yuanhao Xiong, Kai Xu, and Zhenhui Li. Toward a

thousand lights: Decentralized deep reinforcement learning for large-scale traffic signal control. In Proceedings of

the AAAI Conference on Artificial Intelligence, volume 34, pages 3414–3421, 2020.

[11] Xinshi Zang, Huaxiu Yao, Guanjie Zheng, Nan Xu, Kai Xu, and Zhenhui Li. Metalight: Value-based metareinforcement learning for traffic signal control. In Proceedings of the AAAI Conference on Artificial Intelligence,

volume 34, pages 1153–1160, 2020.

[12] Enming Liang, Zicheng Su, Chilin Fang, and Renxin Zhong. Oam: An option-action reinforcement learning framework for universal multi-intersection control. In Proceedings of the AAAI Conference on Artificial Intelligence,

volume 36, pages 4550–4558, 2022.

[13] Azzedine Boukerche, Dunhao Zhong, and Peng Sun. A novel reinforcement learning-based cooperative traffic

signal system through max-pressure control. IEEE Transactions on Vehicular Technology, 71(2):1187–1198,

2022.

[14] Liang Zhang, Qiang Wu, Jun Shen, Linyuan Lü, Bo Du, and Jianqing Wu. Expression might be enough:

representing pressure and demand for reinforcement learning based traffic signal control. In International

Conference on Machine Learning, pages 26645–26654. PMLR, 2022.

[15] Yuanhao Xiong, Guanjie Zheng, Kai Xu, and Zhenhui Li. Learning traffic signal control from demonstrations.

In Proceedings of the 28th ACM International Conference on Information and Knowledge Management, pages 2289–2292, 2019.

[16] Stefano Giovanni Rizzo, Giovanna Vantini, and Sanjay Chawla. Time critic policy gradient methods for traffic

signal control in complex and congested scenarios. In Proceedings of the 25th ACM SIGKDD International

Conference on Knowledge Discovery and Data Mining, KDD ’19, page 1654–1664, New York, NY, USA, 2019.

Association for Computing Machinery.

[17] Tianshu Chu, Jie Wang, Lara Codecà, and Zhaojian Li. Multi-agent deep reinforcement learning for large-scale

traffic signal control. IEEE Transactions on Intelligent Transportation Systems, 21(3):1086–1095, 2019.

[18] Afshin Oroojlooy, Mohammadreza Nazari, Davood Hajinezhad, and Jorge Silva. Attendlight: Universal attentionbased reinforcement learning model for traffic signal control. Advances in Neural Information Processing Systems, 33:4079–4090, 2020.

[19] Zian Ma, Chengcheng Xu, Yuheng Kan, Maonan Wang, and Wei Wu. Adaptive coordinated traffic control for

arterial intersections based on reinforcement learning. In 2021 IEEE International Intelligent Transportation

Systems Conference (ITSC), pages 2562–2567. IEEE, 2021.

[20] Seyed Sajad Mousavi, Michael Schukat, and Enda Howley. Traffic light control using deep policy-gradient and

value-function-based reinforcement learning. IET Intelligent Transport Systems, 11(7):417–423, 2017.

[21] Mohammad Aslani, Mohammad Saadi Mesgari, Stefan Seipel, and Marco Wiering. Developing adaptive traffic

signal control by actor–critic and direct exploration methods. In Proceedings of the Institution of Civil EngineersTransport, volume 172, pages 289–298. Thomas Telford Ltd, 2019.

[22] Haoran Su, Yaofeng D Zhong, Joseph YJ Chow, Biswadip Dey, and Li Jin. Emvlight: A multi-agent reinforcement learning framework for an emergency vehicle decentralized routing and traffic signal control system.

Transportation Research Part C: Emerging Technologies, 146:103955, 2023.

[23] Elise Van der Pol and Frans A Oliehoek. Coordinated deep reinforcement learners for traffic light control.

Proceedings of learning, inference and control of multi-agent systems (at NIPS 2016), 8:21–38, 2016.

[24] Patrick Mannion, Jim Duggan, and Enda Howley. An experimental review of reinforcement learning algorithms

for adaptive traffic signal control. Autonomic road transport support systems, pages 47–66, 2016.

[25] Hua Wei, Guanjie Zheng, Huaxiu Yao, and Zhenhui Li. Intellilight: A reinforcement learning approach for

intelligent traffic light control. In Proceedings of the 24th ACM SIGKDD International Conference on Knowledge

Discovery & Data Mining, pages 2496–2505, 2018.

[26] Lun-Hui Xu, Xin-Hai Xia, and Qiang Luo. The study of reinforcement learning for traffic self-adaptive control

under multiagent markov game environment. Mathematical Problems in Engineering, 2013, 2013.

17

[27] Mohammad Aslani, Mohammad Saadi Mesgari, and Marco Wiering. Adaptive traffic signal control with actorcritic methods in a real-world traffic network with different traffic disruption events. Transportation Research Part

C: Emerging Technologies, 85:732–752, 2017.

- [28] Mohammad Aslani, Stefan Seipel, Mohammad Saadi Mesgari, and Marco Wiering. Traffic signal optimization through discrete and continuous reinforcement learning with robustness analysis in downtown tehran. *Advanced Engineering Informatics*, 38:639–655, 2018.
- [29] Halit Bugra Tulay and Can Emre Koksak. Road state inference via channel state information. *IEEE Transactions on Vehicular Technology*, pages 1–14, 2023.
- [30] Mohamed MG Farag, Hesham A Rakha, Emadeldin A Mazied, and Jayanthi Rao. Integration large-scale modeling framework of direct cellular vehicle-to-all (c-v2x) applications. *Sensors*, 21(6):2127, 2021.
- [31] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [32] Pablo Alvarez Lopez, Michael Behrisch, Laura Bieker-Walz, Jakob Erdmann, Yun-Pang Flötteröd, Robert Hilbrich, Leonhard Lücken, Johannes Rummel, Peter Wagner, and Evamarie Wießner. Microscopic traffic simulation using sumo. In *2018 21st international conference on intelligent transportation systems (ITSC)*, pages 2575–2582. IEEE, 2018.
- [33] Peter Koonce and Lee Rodegerdts. Traffic signal timing manual. Technical report, United States. Federal Highway Administration, 2008.
- [34] Arthur G Sims and Kenneth W Dobinson. The sydney coordinated adaptive traffic (scat) system philosophy and benefits. *IEEE Transactions on vehicular technology*, 29(2):130–137, 1980.
- [35] Pravin Varaiya. The max-pressure controller for arbitrary networks of signalized intersections. *Advances in*

dynamic network modeling in complex transportation systems, pages 27–66, 2013.

[36] Hua Wei, Guanjie Zheng, Vikash Gayah, and Zhenhui Li. A survey on traffic signal control methods. arXiv

preprint arXiv:1904.08117, 2019.

[37] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural networks: A

survey. IEEE transactions on pattern analysis and machine intelligence, 44(9):5149–5169, 2021.

[38] Misha Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. Reinforcement

learning with augmented data. Advances in neural information processing systems, 33:19884–19895, 2020.

[39] Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. arXiv preprint arXiv:2004.13649, 2020.

[40] Nicklas Hansen and Xiaolong Wang. Generalization in reinforcement learning by soft data augmentation. In 2021

IEEE International Conference on Robotics and Automation (ICRA), pages 13611–13617. IEEE, 2021.

[41] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization

algorithms. arXiv preprint arXiv:1707.06347, 2017.

[42] Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stablebaselines3: Reliable reinforcement learning implementations. The Journal of Machine Learning Research,

22(1):12348–12355, 2021.

[43] James Ault and Guni Sharon. Reinforcement learning benchmarks for traffic signal control. In Thirty-fifth

Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1), 2021.

[44] Seung-Bae Cools, Carlos Gershenson, and Bart D’Hooghe. Self-organizing traffic lights: A realistic simulation.

Advances in applied self-organizing systems, pages 45–55, 2013.

[45] Hua Wei, Chacha Chen, Guanjie Zheng, Kan Wu, Vikash Gayah, Kai Xu, and Zhenhui Li. Presslight: Learning

max pressure control to coordinate traffic signals in arterial network. In Proceedings of the 25th ACM SIGKDD

International Conference on Knowledge Discovery & Data Mining, pages 1290–1298, 2019