



TinyML جعبه ابزاری برای به روزرسانی های امن بی سیم و ارزیابی عملکرد مدل های RIOT-ML:

ژائولان هوانگ^۱، کوئن زندبرگ^۱، کاسپار شلیزر^۱، امانوئل باجلی^{۱،۲}

دریافت: ۱۷ ژانویه ۲۰۲۴ / پذیرش: ۱۲ مه ۲۰۲۴ / انتشار آنلاین: ۲۲ مه ۲۰۲۴ © نویسنده(گان) ۲۰۲۴

چکیده

متخصصان حوزه TinyML تاکنون فاقد یک جعبه ابزار جامع و «شامل باتری» برای ساده سازی ادغام مداوم، استقرار مداوم و ارزیابی عملکرد اجرای مدل های مختلف یادگیری ماشین بر روی سخت افزارهای مختلف اینترنت اشیا کم مصرف بوده اند. مقاله ما با پرداختن به این شکاف، RIOT-ML را معرفی می کند، یک جعبه ابزار همه کاره که برای کمک به طراحان و محققان اینترنت اشیا در این وظایف ساخته شده است. برای این منظور، ما RIOT-ML را بر اساس ادغام مجموعه ای از قابلیت ها از یک سیستم عامل تعبیه شده کم مصرف، یک ترنسپایلر و کامپایلر مدل جهانی، یک جعبه ابزار برای اندازه گیری عملکرد TinyML و یک چارچوب به روزرسانی امن کم مصرف از طریق هوا طراحی کردیم - که همه آنها در یک بستر آزمایشی اینترنت اشیا با دسترسی آزاد در دسترس جامعه قابل استفاده هستند. پیاده سازی متن باز ما از RIOT-ML و آزمایش های اولیه ای که در مورد آنها گزارش می دهیم، کاربرد آن را در ارزیابی تجربی عملکرد مدل TinyML در میان ناوگان بردهای اینترنت اشیا کم مصرف تحت آزمایش در میدان، با طیف گسترده ای از معماری های میکروکنترلر ناهمگن و پیکربندی های اتصال شبکه ناوگان، نشان می دهد. وجود یک جعبه ابزار متن باز مانند RIOT-ML برای تسریع تحقیقات در زمینه ترکیب هوش مصنوعی و اینترنت اشیا و تقویت تحقق کامل پتانسیل محاسبات لبه ضروری است.

کلمات کلیدی: هوش مصنوعی، اینترنت اشیا، یادگیری ماشین، قدرت کم، میکروکنترلر، معیارها، به روزرسانی نرم افزار، MLOps

امقدمه

همچنان که هوش مصنوعی (AI) بیشتر و بیشتر در زندگی ما نفوذ می کند، مکانیسم هایی مانند شبکه های عصبی عمیق [۱] در مکان های بیشتر و بیشتری در سیستم های توزیع شده مختلف مورد استفاده قرار می گیرند (یا استقرار آنها برنامه ریزی شده است). به طور خاص، یک شبکه حسگر بی سیم (WSN) می تواند با استقرار هوش مصنوعی در گره های لبه، پوشش و اتصال خود را بهبود بخشد و مصرف انرژی و پهنای باند را کاهش دهد [۲].

خط لوله داده با هوش مصنوعی معمولاً نیاز به ایجاد و استفاده از ... دارد. مدل، یعنی یک ساختار لایه ای از الگوریتم های پیچیده (که به عنوان ... نیز شناخته می شود) پیرامورها که داده ها را تفسیر کرده و براساس آن داده ها تصمیم گیری می کنند. این مدل ابتدا باید آموزش داده شود (یادگیری فاز [۱])، قبل از اینکه بتواند به مرحله تولید برسد (مورد استفاده برای استنباط).

کارهای اخیر از یادگیری ماشین کوچک (TinyML) [۳، ۴] تلاش های جامعه برای بهینه سازی مدل ها برای تطبیق با بودجه های منابع کوچک تر (و در عین حال عملکرد کارآمد)

ب ژائولان هوانگ

zhaolan.huang@fu-berlin.de

^۱ دانشگاه آزاد برلین، برلین، آلمان اینریا ساکلا،

^۲ پاریس، فرانسه

میکروکنترلرهای کم مصرف در اینترنت اشیا (IoT). در نتیجه، هر دو امکان یادگیری و جایگذاری استنتاج برای در بر گرفتن ترمینال های فوق العاده کم مصرف گسترش می یابند.

با این حال، ابزارهای متن باز عمومی و راحت، فاقد طراحی هستند که بتوانند ترکیبی از هوش مصنوعی و اینترنت اشیا (AIoT) را به کار گیرند، طراحانی که ملزم به انجام موارد زیر باشند:

- عملکرد مدل های خود را هنگامی که در امتداد پیوستار ابر-لبه-ترمینال قرار می گیرند، ارزیابی کنند، به خصوص هنگامی که امکان قرارگیری آنها در دستگاه های مختلف مبتنی بر میکروکنترلر در نظر گرفته شود.
- مدل های خود را به دقت تنظیم کنند و تنگناهای عملکرد را در جزئیات لایه مدل، در میکروکنترلرهای مختلف شناسایی کنند.
- یک میکروکنترلر مناسب برای اجرای مدل خود، برای یک کار هدفمند که روی یک دستگاه کم مصرف که قرار است طراحی شود، اجرا می شود، انتخاب کنند.
- به روزرسانی مداوم و ایمن مدل های از پیش آماده شده روی ناوگانی از دستگاه های ناهمگن تحت آزمایش، از راه دور، از طریق شبکه
- نظارت بر عملکرد محاسباتی مدل های مستقر شده از راه دور، از طریق شبکه، مثلاً روی داده های جدید جمع آوری شده در محل

۲ کارهای مرتبط

کارهای اخیر بررسی کرده اند [۶] دامنه چارچوب ها، وظایف و معیارهای یادگیری ماشین، از جمله بررسی جامع پشته TinyML و خطلوله استقرار. برای تطبیق با بودجه های منابع ناچیز معمول میکروکنترلرها (کیلوبایت حافظه، مصرف برق بر حسب میلی وات، فرکانس واحد پردازش مرکزی (CPU) بر حسب مگاهرتز) و در عین حال حفظ عملکرد در سطح قابل قبول و حفظ قابلیت حمل به سخت افزارهای بسیار چندریختی در این دسته، باید با تعدادی از چالش ها روبرو شد.

پلتفرم های نرم افزاری اینترنت اشیا تعبیه شده سیستم های عامل متن باز اینترنت اشیا مختلفی برای ارائه انتزاع سخت افزار، اشتراک گذاری منابع اولیه و دسترسی راحت به لوازم جانبی (مانند حسگر/محرك، زیرسیستم شبکه) روی میکروکنترلرهای ناهمگن کم مصرف استفاده می شوند. کارهای قبلی مانند [۷] چنین سیستم های عاملی را بررسی می کند، که از جمله نمونه های برجسته آن می توان به RIOT [۸] اشاره کرد. و زفیر [۴] با این حال، تاکنون چنین پلتفرم های نرم افزاری هیچ پشتیبانی از چارچوب یادگیری ماشین ارائه نمی دهند - یا اگر ارائه می دهند، این پشتیبانی بسیار محدود است. علاوه بر این، پیشرفته ترین پشتیبانی ها تاکنون معمولاً سخت افزاری هستند - یا مختص فروشنده، مثلاً با کتابخانه هایی که توسط STM32CubeMx یا ARM CMSIS-NN ارائه می شوند.

بسترهای آزمایشی اینترنت اشیا کم مصرف بسترهای آزمایشی مختلف، دسترسی از راه دور به مجموعه ای از دستگاه های مبتنی بر میکروکنترلر قابل برنامه ریزی مجدد را ارائه می دهند. کارهای قبلی مانند [۱۰،۹] چنین بسترهای آزمایشی را بررسی کنید، که از جمله نمونه های برجسته آن می توان به آزمایشگاه اینترنت اشیا با دسترسی آزاد اشاره کرد [۱۱] که دسترسی از راه دور Bare-Metal (سریال از طریق پروتکل کنترل انتقال (TCP)) را به ناوگانی متشکل از صدها برد کم مصرف محبوب از انواع مختلف ارائه می دهد.

مجموعه های بنچمارک برای TinyML بنچمارک کردن یادگیری ماشینی روی سخت افزار کم مصرف، چالش های متعددی را به دنبال دارد [۱۲] کارهای قبلی مانند MLPerf Tiny [۱۳] یک مجموعه معیار استاندارد (مجموعه ای ثابت از وظایف نماینده یادگیری ماشین) برای ارزیابی عملکرد سخت افزار معین و یک پلتفرم آنلاین برای تولیدکنندگان جهت انتشار نتایج معیار مقایسه ای خود ارائه می دهد. در مقابل، RIOT-ML یک جعبه ابزار قدرتمندتر و قابل تنظیم تر برای انجام بررسی های امکان سنجی مدل های یادگیری ماشین تعریف شده توسط کاربر در دستگاه های کم مصرف، با درجه انعطاف پذیری و سفارشی سازی بیشتر ارائه می دهد.

بنچمارک های TinyML کارهای قبلی مانند [۱۴] بر مقایسه عملکرد چارچوب های مختلف یادگیری ماشین روی دو برد کم مصرف COTS (Arduino Nano BLE 33 و F401RE - STM32 NUCLEO) تمرکز دارد. به طور خاص،

بنابراین، این مقاله به معرفی موارد زیر می پردازد: *جعبه ابزار RIOT برای یادگیری ماشینی* (یک گردش کار را برای فشرده سازی، فلش کردن و ارزیابی خودکار مدل های دلخواه پیاده سازی می کند RIOT-ML. کم مصرف را ترکیب می کند IoT متن باز است که به طور تنگاتنگی یک کامپایلر مدل عمومی و یک سیستم عامل محبوب AIoT یک جعبه ابزار، RIOT-ML روی بردهای کم مصرف تجاری (COTS) دلخواه بر اساس معماری های مختلف میکروکنترلر (MCU) محبوب. با بهره گیری از یک پشته شبکه کم مصرف و پرکاربرد در ترکیب با مکانیسم های نرم افزاری امن اینترنت اشیا، RIOT-ML همچنین امکان کنترل، نظارت و به روزرسانی مدل ML را روی مجموعه هایی از چنین دستگاه هایی از راه دور، بر روی توپولوژی های شبکه ناهمگن فراهم می کند.

مشارکت های مقاله ای کمک های ما به شرح زیر است:

- مایک جعبه ابزار جهانی برای ارزیابی داخلی و نظارت از راه دور مدل های RIOT-ML (TinyML) روی دستگاه های کم مصرف طراحی کردیم. RIOT-ML جعبه ابزار U-TOE را ادغام و گسترش می دهد [۵]، بررسی های امکان سنجی برای استفاده از مدل های دلخواه در طیف گسترده ای از پلتفرم های سخت افزاری اینترنت اشیا را ارائه می دهد. این به محققان و توسعه دهندگان اجازه می دهد تا تنگنای عملکرد یک مدل مشخص را در یک دستگاه هدف پیدا کنند. نتایج ارزیابی، امکان طراحی مشترک با سایر اجزا در سطح سیستم را فراهم می کند و به بهینه سازی مدل ها و پیکربندی های ML برای موارد استفاده خاص کمک می کند و امکان دسترسی به بهترین عملکرد ممکن در دستگاه های هدف را فراهم می کند.
- ما مکانیزم هایی را برای به روزرسانی های امن مدل از طریق هوا (OTA) و همچنین برای استقرار و مدیریت مداوم مدل های دلخواه بر روی دستگاه های با محدودیت منابع بر روی پیکربندی های شبکه دلخواه، که ممکن است شامل لینک های شبکه کم مصرف نیز باشد، طراحی، پیاده سازی و ارزیابی می کنیم.
- کدرا منتشر کردیم RIOT-ML تحت مجوز متن باز. این پیاده سازی امکان کامپایل، فلش کردن، ارزیابی و به روزرسانی های OTA امن شبکه های عصبی مختلف (مدل های محاسباتی مبتنی بر گراف) را از چارچوب های اصلی یادگیری ماشین بر روی بردهای کم مصرف مختلف مبتنی بر معماری های مجموعه دستورالعمل محبوب (RISC-V، ESP32، ARM Cortex-M)، فراهم می کند.
- ما معیارها و ارزیابی تجربی مقایسه ای را با استفاده از RIOT-ML ارائه می دهیم که هم در یک بستر آزمایشی اینترنت اشیا با دسترسی آزاد و هم در ایستگاه های کاری شخصی قابل تکرار است، که بینش هایی در مورد عملکرد استنتاج با مدل های مختلف روی سخت افزارهای کم مصرف مختلف ارائه می دهد و نشان می دهد که چگونه RIOT-ML می تواند توسط محققان و توسعه دهندگان تجربی TinyML برای تنظیم دقیق پیکربندی های اینترنت اشیا دوباره مورد استفاده قرار گیرد.

خروجی TensorFlow، PyTorch ...

مانند BBC: microbit، nrf52840dk، Arduino Zero، HiFive ...

دین TinyPART/RIOT-ML <https://github.com/TinyPART/RIOT-ML>

۴ ببینید <https://www.zephyrproject.org/>

معماری ها و پشته های شبکه. در نتیجه، TinyMLOps (و تا حد زیادی خود MLOps) الگویی است که هنوز در مراحل ابتدایی خود قرار دارد.

به روزرسانی های نرم افزاری برای اینترنت اشیا کم مصرف

به روزرسانی های نرم افزاری در زمینه اینترنت اشیا برای رفع آسیب پذیری های امنیتی، گسترش قابلیت ها و بهبود عملکرد دستگاه های کم مصرف در طول عمرشان بسیار مهم هستند. مطالعات تحقیقاتی مانند [۲۲-۲۴] بر مکانیسم های به روزرسانی امن، از جمله تکنیک های پوشش دهنده احراز هویت، بررسی های یکپارچگی و رمزگذاری، با هدف کاهش حملات ناشی از دشمنان بالقوه با سطوح مختلف قدرت محاسباتی، تا سطوح قدرت محاسبات کوانتومی، تمرکز داشت. جنبه های کارایی و قابلیت اطمینان، همراه با استراتژی هایی برای به حداقل رساندن بار شبکه و مدیریت توان، با هدف بهینه سازی فرآیند به روزرسانی OTA در محیط های اینترنت اشیا با منابع محدود شبکه و باتری، بررسی شده اند. [۲۵] تلاش های استانداردسازی در کارگروه مهندسی اینترنت (IETF)، مانند IPv6 روی شبکه های بی سیم شخصی کم مصرف (6LoWPAN) و پروتکل کاربردی محدود (CoAP) [۲۶، ۲۷]، فضای حافظه داخلی و هزینه های انتقال شبکه پروتکل اینترنت نسخه ۶ (IPv6) و تعامل شبیه به پروتکل انتقال ابرمتن (HTTP) را در شبکه کاهش داده اند، در حالی که استانداردهای جدیدتر مانند به روزرسانی های نرم افزاری برای اینترنت اشیا (SUIT) [۲۸] هدف بیشتر، مشخص کردن معماری های عمومی، مدل های داده و فراداده برای به روزرسانی های نرم افزار اینترنت اشیا با امنیت کم و مصرف انرژی کم از طریق پروتکل های انتقال مانند CoAP است.

موارد فوق در جدول خلاصه شده است.

۳ پیشینه عملکرد TinyML تحلیل

از یک سو، از آنجایی که فوری ترین محدودیت بودجه منابع در میکروکنترلرها مربوط به محدودیت های حافظه است، معمولاً در حد کیلوبایت، ارزیابی عملکرد TinyML معمولاً در درجه اول بر معیارهای اندازه گیری مصرف حافظه تمرکز دارد - در حالی که سرعت اجرا را نیز در نظر می گیرد - همانطور که در زیر توضیح داده شده است. از سوی دیگر، تجزیه و تحلیل عملکرد TinyML را می توان در سطوح مختلف جزئیات انجام داد: در سطح مدل جهانی یا در سطح اپراتور، برای جزئیات دقیق تر، همانطور که در ادامه توضیح داده شده است.

۳.۱ معیارهای عملکرد

معیارهای در نظر گرفته شده، بینش هایی در مورد امکان سنجی، کارایی و استفاده از منابع برای انتقال بار استنتاج مدل به دستگاه های کم مصرف ارائه می دهند. با تجزیه و تحلیل این معیارها، کاربران می توانند تصمیمات اولیه ای در مورد انتخاب مدل، تکنیک های بهینه سازی و پیکربندی های سخت افزاری بگیرند.

دو چارچوب TinyML، Tensorflow Lite برای میکروکنترلرها (TFLM) و X-CUBE-AI را در زمینه تشخیص حرکت و تشخیص کلمه بیدارباش، محک زدند. کارهای دیگری مانند [۱۵] مدل های TFLM را روی چندین برد مبتنی بر میکروکنترلر آزمایش کردند. در حالی که چنین مقالاتی مقایسه عملکرد چارچوب های خاص را روی بردهای خاص برای وظایف خاص ارائه می دهند، RIOT-ML انعطاف پذیری و عمومیت بیشتری ارائه می دهد و به توسعه دهندگان اجازه می دهد طیف وسیع تری از مدل های (مشخص شده توسط کاربر) را روی انواع بیشتری از دستگاه های کم مصرف ارزیابی کنند و به جزئیات اجرای مدل های ML بپردازند.

کامپایلر و ترنسپایلر مدل TinyML کامپایلرهای مانند ماشین مجازی تنسور (TVM) [۱۶] می توان از آن برای خودکارسازی ترجمه و کامپایل مدل های ارائه شده توسط چارچوب های اصلی یادگیری ماشین (Pytorch، TFML، و غیره) استفاده کرد تا روال های سطح پایین را در معرض نمایش قرار داده و آنها را برای اجرا بر روی ویژگی های خاص واحد پردازش CPU، واحد پردازش گرافیکی (GPU) و غیره) بهینه کند. اخیراً افزونه ای از TVM به نام uTVM معرفی شده است که اهداف سخت افزاری کوچکتری از جمله انواع MCU ها (واحدهای میکروکنترلر) را اضافه می کند.

پروفایلرهای مدل TinyML [۱۷] TVM-EXray توسعه دهندگان TinyML را قادر می سازد تا به جزئیات سطح لایه اجرای ML دسترسی پیدا کنند و مشکلات استقرار ابری تا لایه را تشخیص دهند. توسعه دهندگان می توانند خطوط لوله استقرار لایه را با قابلیت استفاده بالا، با استفاده از کمتر از ۱۵ خط کد برای بررسی کامل، تجزیه و تحلیل و اشکال زدایی کنند. با این حال، اتکا به Tensorflow Lite قابلیت تطبیق مدل ها از چارچوب های ML بیشتر و استقرار آنها در دستگاه های کم مصرف را محدود می کند. چارچوب های اصلی ML (Pytorch، TFML، و غیره) پروفایلر داخلی ارائه می دهند [۱۸] چنین ابزارهایی به توسعه دهندگان اجازه می دهند عملکرد مدل های خود را اندازه گیری کنند. می توان از آن ها برای جمع آوری معیارهایی مانند زمان استنتاج و میزان استفاده از حافظه استفاده کرد که سپس می توان آن ها را برای بهینه سازی عملکرد مدل تجزیه و تحلیل کرد. اگرچه می تواند جزئیات اجرا را در سطح لایه در اختیار ما قرار دهد، اما هنوز فاقد پشتیبانی از استقرار و ارزیابی روی دستگاه در دستگاه های مختلف اینترنت اشیا است، در حالی که RIOT-ML یک جعبه ابزار عمومی تر است که یک راه حل جامع در طیف وسیعی از دستگاه های کم مصرف ارائه می دهد.

عملیات یادگیری ماشین کوچک (TinyMLOps) مدل یادگیری ماشینی در دستگاه های با منابع محدود، گسترش می دهد. چندین چارچوب CI/CD را برای پشتیبانی از گردش های کاری یکپارچه سازی مداوم و استقرار مداوم TinyMLOps، MLOps است و هدف آن ارائه و نگهداری مدل های یادگیری ماشینی به صورت کارآمد و قابل اعتماد در این زمینه است. این امر مستلزم استقرار مداوم مدل از ابتدا تا انتها و نظارت بر عملکرد آن است (DevOps) عملیات توسعه) الگویی اقتباس شده از مهندسی نرم افزار MLOps [۱۹] -۱۲۱ اجزای MLOps را برای سرورهای معمولی یا خوشه های بزرگ ارائه می دهد، اما معمولاً دستگاه های TinyML را به دلیل چالش های منحصر به فردی که این دستگاه ها با آن مواجه هستند، به ویژه چندرخی شدید در سخت افزار، در بر نمی گیرد.

چارچوب	میکروکنترلر	نوع مدل	ارزیابی از راه دور/بالا بردن	گرانو	مدل برای سوار شدن به سول
MLPerf	بله	مشخص شده	خیر	مدل	خیر
ML-EXray	خیر	تی اف لایت	خیر	لایه	خیر
تی اف لایت	بله	تی اف لایت	خیر	لایه	خیر
پیتورچ	خیر	مشعل	خیر	لایه	خیر
بوتی وی ام	بله	یونیورسال	خیر	اپراتور	خیر
یو-تو	بله	یونیورسال	خیر	اپراتور	بله
شورش-ML	بله	یونیورسال	بله	اپراتور	بله

ای وی ال، ارزیابی، سول راه حل، بالا، به روزرسانی، گرانودانه بندی

عملکرد را به حداکثر رسانده و رد پای منابع را در دستگاه های کم مصرف به حداقل برسانید.

مصرف حافظه (RAM) این معیار، میزان فضای حافظه پویا (حافظه دسترسی تصادفی اولیه (RAM)) مصرف شده توسط مدل در طول استنتاج را اندازه گیری می کند. این معیار، میزان اشغال فضای حافظه توسط فعال سازی مدل را نشان می دهد و برای دستگاه های کم مصرف که منابع حافظه محدودی دارند، مهم است. استفاده کارآمد از حافظه، امکان استقرار مدل های بزرگ تر و پیچیده تر را در چنین دستگاه هایی فراهم می کند.

مصرف حافظه (فلش مموری) این معیار، میزان فضای ذخیره سازی مورد نیاز برای ذخیره دستورالعمل محاسباتی و پارامترهای مرتبط را، معمولاً بر حسب ناحیه حافظه فلش، تعیین می کند. این معیار، میزان فضای ذخیره سازی مدل را در دستگاه کم مصرف نشان می دهد. به حداقل رساندن مصرف فضای ذخیره سازی، امکان جای دادن چندین مدل در دستگاه یا هماهنگی با سایر برنامه های ضروری را فراهم می کند.

تأخیر محاسباتی این معیار، میزان مصرف زمان انجام استنتاج برای هر نمونه ورودی، چه در سطح مدل و چه در سطح عملگرهای منفرد درون مدل را اندازه گیری می کند. این معیار، سرعت استنتاج مدل را روی دستگاه کم مصرف نشان می دهد و نقش حیاتی در برنامه های بلند رنگ یا حساس به تأخیر ایفا می کند. فرکانس کلاک هسته، استراتژی های حافظه پنهان و تأخیر ارتباط بین حافظه و هسته در حال کار، تأثیر زیادی بر این شاخص دارند.

قیمت تراشه این معیار، هزینه سیستم-آن-چیپ (SoC) مورد استفاده در دستگاه کم مصرف را در نظر می گیرد. قیمت SoC بر مقرون به صرفه بودن کلی و امکان پذیری استقرار مدل در سیستم های توزیع شده در مقیاس بزرگ تأثیر می گذارد. SoC های کم هزینه ترمی توانند استقرار را در دسترس تر و مقرون به صرفه تر کنند.

۳.۲ جزئیات اندازه گیری

در مورد تحلیل عملکرد یادگیری ماشین در حوزه های دیگر، عملکرد TinyML را می توان در سطوح مختلف جزئیات اندازه گیری کرد:

ارزیابی هر مدل در این سطح کلی، عملکرد مدل به عنوان یک کل، یعنی میزان منابع مصرفی ناشی از اجرای مدل شامل تمام لایه ها و عملگرهای آن، اندازه گیری می شود. به عنوان مثال، این امر امکان ارزیابی مصرف منابع برای استنتاج با کد آماده تولید، در یک پیکربندی سخت افزار صنعتی خاص را فراهم می کند.

ارزیابی به ازای هر اپراتور در این سطح، عملکرد یک یا چند اپراتور (یعنی یک یا چند جزء از مدل) به طور جداگانه اندازه گیری می شود. این اندازه گیری به ازای هر اپراتور می تواند به شناسایی اپراتورهای خاصی که در عملکرد یا ناکارآمدی نقش دارند، در بهینه سازی کارایی مدل و شناسایی تنگناهای احتمالی کمک کند.

۴ پیشینه به روزرسانی مدل TinyML

یک مدل مستقر و در حال اجرا بر روی یک دستگاه می تواند در سطوح مختلف به روزرسانی شود. بسته به این سطح، انعطاف پذیری و هزینه های ترافیک شبکه می تواند متفاوت باشد.

به روزرسانی های میان افزار هرچند ممکن است تعجب آور باشد، اما رویکرد غالب در حوزه به روزرسانی های نرم افزار اینترنت اشیا کم مصرف در دستگاه های مبتنی بر میکروکنترلر، انجام به روزرسانی میان افزار است، یعنی انتقال، تأیید و (دوباره) نصب کل نرم افزار در حال اجرا روی دستگاه (به جز یک بوت لودر حداقلی، در برخی موارد). بنابراین، یک رویکرد ساده برای به روزرسانی یک مدل از راه دور از طریق شبکه، انتشار و استقرار به روزرسانی میان افزار حاوی مدل جدید است. این سطح مستلزم انتقال بیشترین حجم داده به دستگاه است که می تواند بار شبکه و نیازهای حافظه بافر را افزایش دهد.

به روزرسانی های کامل مدل یک رویکرد اصلاح شده تر برای به روزرسانی نرم افزار از راه دور از طریق شبکه، انتشار و استقرار به روزرسانی های جزئی نرم افزار به جای به روزرسانی های میان افزار است. با این رویکرد، کاربر فقط کد دودویی مدل، یعنی عملگرهای آن و پارامترهای مرتبط (وزن ها، بایاس و غیره) را به روزرسانی می کند. انتقال شبکه کاهش می یابد، در حالی که انعطاف پذیری بالا باقی می ماند، زیرا کاربر می تواند عملگرها را با جریان های اجرایی بهینه شده به روزرسانی کند، یا حتی یک مدل را با معماری کاملاً متفاوت با عملکرد بهتر در دقت یا استفاده از منابع جایگزین کند.

مکانیسم هایی که ما طراحی می کنیم بنیادی هستند. آن ها باید امکان گسترش به سمت کاهش حملات پیچیده تر را فراهم کنند.

۶ چارچوب RIOT-ML

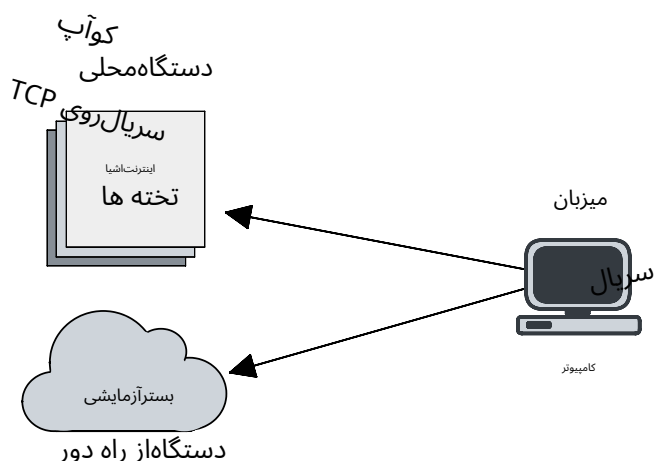
برای ارزیابی مدل های دلخواه تولید شده توسط چارچوب های یادگیری ماشین رایج روی بردهای کم مصرف مختلف استفاده می کنند U-TOE را برای انجام کامپایل مدل، بهینه سازی و فلش کردن، به روزرسانی امن مدل و مدیریت مدل ادغام می کند و از RIOT و uTVM، RIOT-ML

همانطور که در شکل نشان داده شده است. کاربران می توانند از یک رایانه شخصی میزبان (PC) (ترکیبی از لینوکس و زنجیره ابزار RIOT-ML) برای استقرار، ارزیابی و به روزرسانی مدل های خود در دستگاه های اینترنت اشیا محلی (مثلاً متصل از طریق گذرگاه سریال جهانی (USB) به رایانه شخصی خود) یا روی دستگاه های اینترنت اشیا از راه دور از طریق یک بستر آزمایشی که دسترسی فیزیکی به انواع مختلف MCU ها را از طریق شبکه ارائه می دهد، استفاده کنند.

۶.۱ طراحی معماری

همانطور که در شکل نشان داده شده است. ۲. این جعبه ابزار از اجزای کلیدی زیر تشکیل شده است:

- **کامپایلر مدل.** را به عنوان ورودی دریافت می کند، سپس کارایی مدل ها را برای میکروکنترلرها افزایش می دهد و امکان اجرای آنها را روی دستگاه های کم مصرف فراهم می کند. (Tensor-Flow، PyTorch مانند) کارآمد استفاده می کند. این کامپایلر، خروجی چارچوب های معمول یادگیری ماشین C برای تبدیل مدل های دلخواه شبکه عصبی به کد uTVM از کامپایلر RIOT-ML
- **محیط سیستم عامل و پشتیبانی از سخت افزار.** یک سیستم عامل همه منظوره برای دستگاه های اینترنت اشیا کم مصرف است که برای ارائه یک محیط زمان اجرای سبک برای اجرای مدل و ارزیابی روی میکروکنترلرها انتخاب شده است. این پایه، قابلیت توسعه و پشتیبانی از طیف گسترده را برای بردهای کم مصرف ناهمگن فراهم می کند RIOT



شکل ۱. راه اندازی سخت افزار RIOT-ML. کاربران می توانند بردهای محلی را از طریق سریال به کامپیوتر میزبان متصل کنند یا از سرویس برد از راه دور در بستر آزمایشی اینترنت اشیا استفاده کنند.

به روزرسانی های جزئی مدل یک رویکرد دقیق تر، به روزرسانی تنها زیرمجموعه ای از پارامترهای مشخص شده از یک مدل از پیش مستقر شده است. به عنوان مثال، کاربر می تواند یک لایه واحد (یا حتی یک پارامتر واحد) را برای به روزرسانی مشخص کند که می تواند ترافیک شبکه و مصرف برق مرتبط را به حداقل برساند. چنین سطح دقیقی ممکن است برای گره های اینترنت اشیا کم مصرف با همان مدل که پارامترها را مبادله می کنند (یادگیری فدرال/توزیع شده) یا پارامترهای جدیدتر را از یک سرور مرکزی دریافت می کنند، مناسب باشد. با این حال، توجه داشته باشید که انعطاف پذیری کاهش می یابد، زیرا چنین به روزرسانی نمی تواند معماری یک مدل را به طور قابل توجهی تغییر دهد.

۵ مدل تهدید و امنیت TinyML به روزرسانی های مدل

کاربران به مدل های یادگیری ماشینی برای ارائه پیش بینی های دقیق و بی طرفانه اعتماد دارند. برای حفظ قابلیت اعتماد مدل، حداقل بررسی های یکپارچگی و احراز هویت در طول فرآیند به روزرسانی مدل برای اطمینان از قابلیت اطمینان سیستم های TinyML و کاهش حملاتی مانند مسمومیت مدل، که در آن مهاجمان سعی در دستکاری داده های آموزشی یا تزریق بیت های مخرب به پارامترهای مدل دارند، ضروری است. در این مقاله، ما به بررسی دفاع های پیچیده در برابر مسمومیت مدل نمی پردازیم. در عوض، پنج بردار حمله اساسی را در طول فرآیند به روزرسانی مدل بررسی می کنیم:

به روزرسانی مدل دستکاری شده: یک مهاجم مخزن مدل را در اختیار دارد و سعی می کند مدل های معیوب را آپلود کند، که توسط دستگاه های اینترنت اشیا برای به روزرسانی مدل دریافت می شوند.

به روزرسانی پارامتر دستکاری شده: این مورد مشابه مورد فوق است، اما در جزئیات تغییر پارامتر (مخرب).

به روزرسانی مدل غیرمجاز: یک مهاجم غیرمجاز تلاش می کند تا از دستگاه اینترنت اشیا درخواست کند تا مدل های اصلاح شده را دریافت و مستقر کند.

به روزرسانی غیرمجاز پارامتر: این همان مورد فوق است، اما در جزئیات تغییرات پارامتر (غیرمجاز).

نقض محرمانگی: به روزرسانی های مدل ها یا پارامترها باید رمزگذاری شوند تا فقط نگهدارنده مجاز و دستگاه به روزرسانی های مدل رمزگشایی شده دسترسی داشته باشند.

توجه داشته باشید که این مدل تهدید محدودیت هایی دارد: این مدل مواردی را که ریشه اعتماد، یعنی نگهدارنده مجاز مدل، خود به خطر افتاده باشد، پوشش نمی دهد. به عنوان مثال، خود نگهدارنده مجاز می تواند سرکش شود، یا به گونه ای دیگر فریب بخورد تا به روزرسانی های مدل را با بدافزار ترکیب کند، یا پارامترهای مدل را به طور مخرب تغییر دهد. با این وجود، هر مدل تهدیدی باید بردارهای فوق را پوشش دهد، و از این نظر، امنیت

یک برنامه ریزی مناسب^۵ این کتابخانه تولید شده توسط uTVM سپس به طور مشترک با RIOT، یک سرور RPC و یک کارگر اندازه گیری در یک میان افزار اجرایی کامپایل می شود که سپس به طور خودکار روی دستگاه (از طریق USB یا از راه دور از طریق شبکه) فلش می شود.

۶.۲ به روزرسانی و مدیریت مدل از طریق هوا

مرباط ها و ماژول هایی را هم روی دستگاه های اینترنت اشیا و هم روی میزبان (مدیر ناوگان) ارائه می دهیم تا به روزرسانی و مدیریت OTA امن را امکان پذیر کنیم.

معماری از یک سو، اجزای نرم افزاری که روی دستگاه های اینترنت اشیا نصب شده اند و می توانند به دستورات مدیریتی پاسخ دهند و بر اساس اعلان های به روزرسانی عمل کنند، بر روی پشته شبکه کم مصرف RIOT IPv6 و پیاده سازی SUIT پیاده سازی شدند. رابط های به روزرسانی و مدیریت به عنوان ماژول های جداگانه RIOT پیاده سازی شدند. از سوی دیگر، میزبان، به عنوان مدیر ناوگان، مسئول ساخت و هماهنگی به روزرسانی ها، ارسال دستورات و اعلان ها و نگهداری مخزن به روزرسانی (سرور) با رجیستری (CoAP) است که بارهای به روزرسانی را ذخیره می کند.

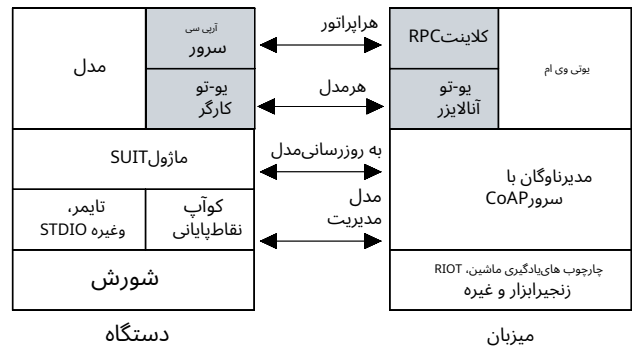
مدیریت مدل برای کنترل از راه دور، نگهداری و نظارت بر عملکرد و کارایی مدل ML به طور کارآمد و مداوم، RIOT-ML چندین نقطه پایانی CoAP را در دستگاه های مدیریت شده IoT در معرض نمایش قرار می دهد، همانطور که در جدول نشان داده شده است.^۲ این نقاط پایانی می توانند توسط تعمیرکار مجاز از طریق ماژول مدیریت ناوگان قابل دسترسی باشند.

روش به روزرسانی مدل این روش شامل دو مرحله است. در مرحله اول، که در شکل نشان داده شده است.^۴ یک نگهدارنده مجاز که با جفت کلید عمومی/خصوصی شناسایی می شود (پای، سی) فایل باینری به روزرسانی و فراداده های مرتبط که طبق SUIT (مانیفست SUIT) ایمن شده اند را تولید و ارسال می کند. در مرحله دوم، که در شکل نشان داده شده است.^۵، مسئول نگهداری (با استفاده از ماژول مدیریت ناوگان) به دستگاه های اینترنت اشیا مدیریت شده اطلاع می دهد که یک به روزرسانی جدید باید دریافت، تأیید و نصب شود، که سپس این کار انجام می شود.

توجه داشته باشید که طبق مشخصات SUIT، می توان از طرح های مختلفی برای امضای دیجیتال، هشینگ و رمزگذاری استفاده کرد (مقایسه عملکرد در کارهای قبلی مانند [موجود است ۲۴] در زیر، فرض می کنیم از الگوریتم هش امن ۲۵۶ بیتی (SHA256) استفاده می شود [۲۹] برای هش کردن، و الگوریتم Ed25519 [۳۰] برای امضای دیجیتال.

در ادامه، مراحل به روزرسانی با جزئیات بیشتر توضیح داده شده است:

هیک برنامه، بهینه سازی سطح پایین را برای اجرای حلقه مشخص می کند، که باعث افزایش میزان دسترسی به حافظه پنهان و حافظه می شود. برنامه بهینه با مشخصات مدل و دستگاه تعیین شده و توسط الگوریتم های جستجوی اکتشافی بر اساس اندازه گیری های روی دستگاه شناسایی می شود [۱۶].

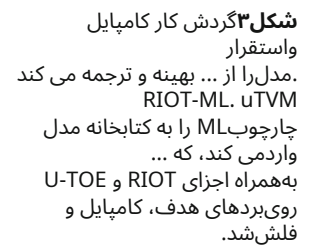


شکل ۲ معماری نرم افزار و اجزای چارچوب RIOT-ML. به رنگ خاکستری: اجزای جعبه ابزار U-TOE

- **ماژول ارزیابی مدل** این کامپوننت، ابزار U-TOE را ادغام می کند^۵ برای انجام ارزیابی مدل در محل. بخش^۹ معماری و رویه های اندازه گیری آن را با جزئیات شرح می دهد.
- **ماژول به روزرسانی مدل امن** این ماژول پیاده سازی درخواست نظرات (RFC) [9019] SUIT را ادغام می کند.^{۲۲} برای ارائه به روزرسانی های نرم افزاری امن و کم مصرف در دستگاه های اینترنت اشیا. این ابزار از SUIT manifests که حاوی اطلاعات نسخه مدل، هش های payloadها و امضای دیجیتال از سوی نگهدارندگان است، برای اطمینان از یکپارچگی و اصالت در طول فرآیند به روزرسانی استفاده می کند. ما آن را در RIOT-ML تطبیق داده و ادغام کرده ایم تا از به روزرسانی مدل OTA امن، همانطور که در بخش توضیح داده شده است، پشتیبانی کند.^{۶.۲}
- **نقاط پایانی مدیریت شبکه** را در دستگاه های اینترنت اشیا که این چارچوب را اجرا می کنند، در معرض نمایش قرار می دهد و به عنوان رابط هایی برای مدیریت مدل عمل می کند. کاربران می توانند ابر داده ها و وضعیت مدل (نام، نسخه، نتایج ارزیابی و غیره) را جستجو کنند، رفتارهای مدل را کنترل کنند یا به روزرسانی مدل را فعال کنند، همانطور که در بخش توضیح داده شده است CoAP نقاط پایانی RIOT-ML.^{۶.۲}
- **ماژول مدیریت ناوگان** مدیریت شده را در مورد در دسترس بودن به روزرسانی های مدل مطلع می سازد. همچنین یک رابط کاربر پسند ارائه می دهد IoT ارسال می کند و دستگاه های CoAP رجیستری را تولید و امضا می کند، آنها را به یک مخزن SUIT مدل های یادگیری ماشین روی دستگاه خودکار می کند. این مدل، بارهای به روزرسانی را می سازد، مانیفست های CI/CD مدیریت ناوگان را روی میزبان برای به روزرسانی های از راه دور مدل ها پیاده سازی می کند و یک خط لوله ساخت را برای پشتیبانی از RIOT-ML

علاوه بر این، RIOT-ML فراهم می کند که تکنیک برای یک بستر آزمایشی اینترنت اشیا مبتنی بر ابر که تعامل یکپارچه با بردهای از راه دور را با استفاده از سریال روی TCP امکان پذیر می کند.

گردش کار ما در شکل نشان می دهیم.^۳ نمای سطح بالا از گردش کار معمول با RIOT-ML. در مرحله اولیه، RIOT-ML ابتدا مشخصات دستگاه هدف را جمع آوری می کند تا گزینه های کامپایل برای uTVM و RIOT را تعیین کند. سپس، uTVM یک "کتابخانه" مدل غیر بهینه تولید می کند. سپس برخی از استراتژی های بهینه سازی استاتیک در این مرحله، به ویژه بر اساس نوع دستگاه هدف، اعمال می شوند تا تعیین کنند



۷. دستگاه، فایل باینری به روزرسانی (بار داده) تعیین شده در مانیفست SUIT را دریافت کرده و بررسی یکپارچگی را انجام می دهد.

نقاط پایانی CoAP	اظهارات
(برای مانیتور)	
مدل/وضعیت	دریافت وضعیت کاری مدل،
مدل/نام	دریافت نام مدل
مدل/پارامترها/اطلاعات	دریافت اطلاعات پارامترها
مدل/eval_result	دریافت آخرین نتایج ارزیابی (تأخیر محاسباتی، از دست دادن داده و غیره)
(برای کنترل)	
مدل/توقف	مدل را متوقف کنید
مدل/اجرا	مدل را اجرا کنید
مدل/اجرا-ارزیابی	مدل را ارزیابی کنید
مدل/پارامترها/به روزرسانی	به روزرسانی جزئی را فعال کنید
(برای کت و شلوار)	
/active/in/slot/suit	تعداد اسلات های فعال یا غیرفعال فریمور را دریافت کنید
کت و شلوار/ماشه	به روزرسانی میان افزار را فعال کنید
کت و شلوار/نسخه	دریافت نسخه فعلی سیستم عامل

شکل ۴: تولید و ارسال مانیفست و پیلودهای SUIT برای به روزرسانی مدل امن با RIOT-ML. جفت کلید (پی، سر) مسئول نگهداری مجاز را احراز هویت می‌کند. سرور CoAP هیچ اطلاعاتی از کلید مخفی ندارد. سر و فقط به عنوان مخزن آثار باستانی عمل می‌کند.

تحلیلگر روی میزبان اجرا می شود، معیارهای آپلود شده از دستگاه را آماری می کند و یک رابط کاربری قابل خواندن توسط انسان برای کاربران فراهم می کند.

گردش کار ارزیابی با U-TOE: پس از استقرار برنامه روی دستگاه هدف، یک کانال دو طرفه بین میزبان و دستگاه برقرار می شود، همانطور که در شکل ۱ نشان داده شده است. ۲. عامل اندازه گیری شروع به جمع آوری معیارهای عملکرد در سطح مشخص شده توسط کاربر می کند و داده های معیارها را به تحلیلگر آپلود می کند. در نهایت، کاربران می توانند آمار کلی معیارهای مدل را به دست آورند یا گلوگاه عملکرد را با جزئیات اجرای هر اپراتور پیدا کنند. تمام داده های خام معیارها در یک فایل گزارش برای تجزیه و تحلیل بیشتر و سفارشی سازی شده توسط کاربر ذخیره می شوند.

۷.۲ روش اندازه گیری

مادو رویه اندازه گیری برای پشتیبانی از ارزیابی در جزئیات مختلف طراحی کردیم. این رویه ها در چندین مؤلفه از جعبه ابزار اجرایی شوند و بیشتر حجم کار عمدتاً روی برد هدف قرار دارند. مراحل زیر روال اندازه گیری را پس از کامپایل یک برنامه اجرایی شرح می دهند. مراحل مشخص شده با **عدد پررنگ** روی تخته هدف اجرا می شوند.

ارزیابی هر مدل این حالت بر عملکرد مدل در یک محیط تولید واقعی تمرکز دارد. روال اندازه گیری مربوطه به شرح زیر است:

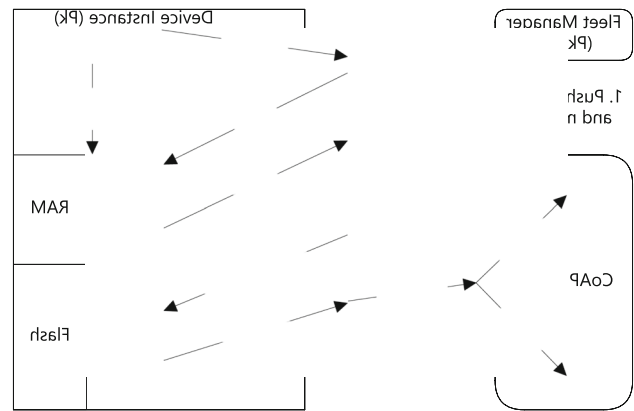
۱. محاسبه ی میزان مصرف حافظه و فضای ذخیره سازی مدل بر اساس فایل (Executable and Linkable Format) ELF. ما تخصیص حافظه ی پویا را غیرفعال می کنیم تا تحلیل استاتیک ردپای حافظه را فعال کنیم.
۲. برنامه اجرایی را روی برد اینترنت اشیا محلی یا از راه دور مستقر کنید.
۳. استنتاج مدل را بر اساس تعداد آزمایش های خاص کاربر با ورودی تصادفی روی توزیع یکنواخت تکرار کنید.

۴. تأخیر محاسباتی هر آزمایش را ثبت کنید.

۵. سوابق را برای تجزیه و تحلیل بیشتر و بایگانی به دستگاه میزبان آپلود کنید.

در پایان ارزیابی، نتایج به همراه آمار (مثلاً فاصله اطمینان ۹۵٪، میانه، حداکثر و حداقل) روی دستگاه میزبان ارائه می شوند، از جمله تأخیر محاسباتی و مصرف حافظه و فضای ذخیره سازی.

ارزیابی به ازای هر اپراتور در مقابل ارزیابی مدل برتر، این حالت بر کارایی و میزان مصرف منابع هر اپراتور تمرکز دارد و امکان کشف گلوگاه های عملکرد درون مدل ها را فراهم می کند. به لطف ارزیاب زمان درون مکانیسم uTVM RPC، می توانیم تأخیر محاسباتی را در سطح اپراتور اندازه گیری کنیم. انتزاع بالای تایمر و سریال در RIOT به ما امکان می دهد پیاده سازی ... را یکپارچه کنیم.



شکل ۵: اطلاع رسانی، دریافت، تأیید و نصب به روزرسانی های مدل با ML-RIOT. به صورت ایمن از طریق شبکه

محرمانگی به روزرسانی مدل مشخص می کند (به [مراجعه کنید] payload SUIT همچنین مکانیسم های رمزگذاری (اختیاری) را برای فایل های باینری IETF [۳]). برای این کار، از یک کلید رمزگذاری محتوای متقارن استفاده می شود، که این کلید یا از قبل به اشتراک گذاشته شده است یا به صورت آنی ایجاد می شود، به عنوان مثال از طریق یک تبادل دیفی-هلمن ایستا و زودگذر (EDHOC) [۳۲] در صورت استفاده، این مکانیسم محرمانگی را برای به روزرسانی های مدل ها فراهم می کند.

۹ جعبه ابزار U-TOE

رابرایی انجام کامپایل و ارزیابی مدل، روی یک دستگاه مبتنی بر میکروکنترلر، ادغام می کند uTVM و RIOT ابزاری است که ما طراحی و پیاده سازی کرده ایم و U-TOE

۷.۱ طراحی معماری

این جعبه ابزار از اجزای کلیدی زیر تشکیل شده است:

- **ماژول RPC** برای ارزیابی مصرف منابع در سطح اپراتور، U-TOE از مکانیسم فراخوانی رویه از راه دور (RPC) در uTVM استفاده می کند. مکانیسم RPC امکان آپلود و راه اندازی توابع روی برد های IoT را از طریق سریال فراهم می کند. این امر برای آزمایش و پروفایل سازی از راه دور مفید است و U-TOE را قادر می سازد تا اپراتورهای مدل را برای اندازه گیری تأخیر محاسباتی و میزان استفاده از حافظه، در خود جای دهد. این سیستم از یک کلاینت روی میزبان و یک سرور روی دستگاه هدف تشکیل شده است که دستورات و دستورالعمل های اجرایی را از میزبان دریافت می کند.

- **ماژول ارزیابی** شامل دو واحد است: اندازه گیری کننده و تحلیل گر. همانطور که در شکل نشان داده شده است، کارگر اندازه گیری برای دستیابی به معیارهای عملکرد در سطح مدل یا اپراتور در MCU مستقر شده است. این کارگر علاوه بر انجام اندازه گیری ردپای منابع، مسئول تصادفی سازی ورودی مدل و گزارش داده های معیارها به دستگاه میزبان است.

اندازه گیری زمان و ارتباط RPC روی بردهای اینترنت اشیا دلخواه. در اینجا روال اندازه گیری مربوطه آمده است:

۱. با استفاده از API داخلی TVM، میزان حافظه اشغال شده توسط دستگاه را در سطح اپراتور تجزیه و تحلیل کنید.
۲. برنامه اجرایی را روی برد اینترنت اشیا محلی یا از راه دور مستقر کنید.
۳. سرور RPC را روی برد اینترنت اشیا راه اندازی کنید.
۴. کلاینت RPC را برای ارزیابی و ثبت عملکرد اجرای هر اپراتور راه اندازی کنید.

لازمه ذکر است که ساختار عملگر ساخته شده توسط uTVM معمولاً با نسخه دست ساز در چارچوب ML مغایرت دارد. دلیل این امر آن است که uTVM به عنوان کامپایلر مدل، بهینه سازی مدل (یعنی ادغام عملگرها) را اعمال می کند و جزئیات اجرا (یعنی حساب کوانتیزاسیون) را در طول تبدیل و کامپایل وارد می کند، که به طور بالقوه چندین عملگر را در یک عملگر ادغام می کند یا عملگرهای اضافی را وارد می کند. با این وجود، ما عملگرهای uTVM را با پارامترهای مدل (وزن ها، بایاس ها و غیره) حاشیه نویسی می کنیم، به طوری که کاربران می توانند یک عملگر خاص را به لایه مربوطه مرتبط کنند.

۸. آزمایش با RIOT-ML

در ادامه قابلیت های RIOT-ML را نشان می دهیم. دو دسته از آزمایش های انجام شده را گزارش می دهیم. در این بخش، ابتدا تنظیمات آزمایشی برای هر دسته از آزمایش ها را شرح می دهیم، سپس بخش ۹ نتایج جمع آوری شده را ارائه و تحلیل خواهیم کرد.

۸.۱ ارزیابی عملکرد مدل

ما آزمایش هایی را برای اعتبارسنجی عملکرد و سازگاری RIOT-ML در سمت مدل (پشتیبانی جهانی از ساختار مدل و چارچوب های ML) و در سمت دستگاه (پشتیبانی طیف گسترده برای دستگاه های IoT) انجام دادیم. از این رو، ما به دو جهت متعامد پرداختیم: برای پشتیبانی از دستگاه، یک LeNet-5 کوانتیزه شده را روی بردهای مختلف IoT ارزیابی کردیم؛ برای سازگاری مدل، چندین مدل را روی یک برد اکتشافی محلی STM32F746G ارزیابی کردیم.

انتخاب مدل ما مدل های از پیش آموزش دیده و کوانتیزه شده را از مخازن متن باز انتخاب کردیم^۶ که وظایف معمول TinyML (طبقه بندی تصویر، تشخیص کلمات کلیدی، ردیابی بصری کلمات، حذف نویز و تشخیص ناهنجاری) را هدف قرار می دهند. وزن ها و فعال سازی های مدل به صورت عدد صحیح ۸ بیتی کوانتیزه شده اند، با این حال ورودی ها و خروجی ها در IEEE 754 با دقت ممیزشاور باقی می مانند [۳۳] قالب.

بهینه سازی مدل ما فقط از بهینه سازی مبتنی بر قانون داخلی در uTVM استفاده کردیم. بنابراین، تمام استراتژی های بهینه سازی اکتشافی مانند زمان بندی مدل غیرفعال شدند.

پیکربندی MCU ما حافظه پنهان داده ها و دستورالعمل ها را غیرفعال کردیم تا اثر «دیوار حافظه» را در مدل ML مشاهده کنیم. فرکانس کلک هسته توسط کد مقداردهی اولیه CPU در RIOT از پیش تنظیم شده بود و نتایج آزمایش در بخش ارائه شده است. ۹.

استقرار ترکیبی این آزمایش ها هم بر روی بردهای اینترنت اشیا محلی و هم از راه دور ارائه شده توسط FIT IoT-LAB انجام شد. لازم به ذکر است که برای هر ارزیابی، تعداد آزمایش ها را از قبل روی ده تنظیم کرده ایم تا خطای تصادفی را در نظر بگیریم.

۸.۲ به روزرسانی های OTA مدل امن

ما آزمایش هایی را برای به روزرسانی مدل هایی که روی برد کیت توسعه nRF52840dk در بستر آزمایشی FIT IoT-LAB اجرا می شوند، انجام دادیم و میزان مصرف منابع اجزای مهم و سربار انتقال شبکه را اندازه گیری کردیم. اندازه گیری های ما از مدل 5-LeNet به عنوان یک مورد استفاده استفاده کردند. ما جزئیات مختلفی را برای به روزرسانی در نظر می گیریم: یا به روزرسانی میان افزاری در عوض، یک به روزرسانی جزئی که فقط مربوط به وزن های لایه نهایی (به عنوان طبقه بندی کننده) مدل 5-LeNet کوانتیزه شده است.

راه اندازی شبکه برد nRF52840 از طریق یک لینک دسترسی رادیویی کم مصرف IEEE 802.15.4 (با استفاده از پشته شبکه عمومی (GNRC) در RIOT) به شبکه متصل می شود که ترافیک 6LoWPAN و IPv6 می تواند روی آن جریان یابد. از طریق یک روتر مرزی IPv6 میانی، پیام های CoAP می توانند به/از دستگاه منتقل شوند و به مدیر ناوگان در رایانه شخصی کاربر متصل به اینترنت میزبانی می شود. براساس استفاده از این رویکرد (6LoWPAN و IPv6)، ما اطمینان حاصل می کنیم که این رویکرد می تواند به صورت سرتاسری از طریق اینترنت و بر روی مجموعه ای دلخواه از لینک ها اجرا شود که می تواند نه تنها شامل لینک های معمولی اترنت و وای فای، بلکه برخی از لینک های رادیویی کم مصرف (مانند بلوتوث کم مصرف (BLE)، برد بلند (LoRa)، IEEE 802.15.4 نیز باشد.

پیکربندی رمزنگاری ما قبل از راه اندازی ماژول های به روزرسانی مدل، کلیدهای عمومی و مخفی را با استفاده از الگوریتم Ed25519 تولید کردیم. امضا و رمزگذاری شیء COSE (CBOR [۳۴] طبق دستورالعمل SUIT برای امضای مانیفست و بارهای داده CoAP استفاده شد. کتابخانه های رمزنگاری که ما در عمل استفاده کردیم عبارتند از vlibcose و c25519، که سیستم های با حافظه کم را هدف قرار می دهند و به ترتیب COSE و Ed25519 را پیاده سازی می کنند.

^۷ ببینید <https://github.com/bergzand/libcose>

^۸ ببینید <https://www.dlbeer.co.nz/oss/c25519.html>

^۶ دیدن <https://github.com/ARM-software/ML-zoo> و <https://mlcommons.org>

جدول ۳: ارزیابی مدل LeNet-5 روی بردهای مختلف اینترنت اشیا

مکس	حداقل	میان	تأخیر حاسباتی (میلی ثانیه) 95%-CI	فشار محوری سازه (کیلوپایت)	حافظه (کیلوپایت)	هسته	برد/MCU
۲۶۲.۲۱۶	۲۶۱.۳۵۰	۲۶۲.۱۸۷	[۲۶۲.۱۳۹, ۲۶۱.۸۲۹]	۶۴.۳۴۰	۱۱.۲۸۸	در ۳۳ مگاهرتز M0+	rpi-pico/RP2040
۱۷۶.۹۷۵	۱۷۶.۹۲۴	۱۷۶.۹۵۸	[۱۷۶.۹۶۵, ۱۷۶.۹۳۶]	۶۵.۱۶۸	۱۱.۲۰۸	در ۴۸ مگاهرتز M0+	ATSAMD21G18/آردوینو-نیزو
۱۸۲.۰۹۸	۱۸۲.۰۵۱	۱۸۲.۰۶۸	[۱۸۲.۰۸۲, ۱۸۲.۰۶۱]	۶۴.۹۴۰	۱۱.۲۹۲	در ۴۸ مگاهرتز M0+	samr30-xpro/ATSAMR30G18A
۷۰.۱۵۱	۷۰.۰۹۱	۷۰.۱۱۷	[۷۰.۱۳۰, ۷۰.۱۰۸]	۶۵.۱۷۲	۲۸.۷۰۴	در ۱۲۵ مگاهرتز M0+	b-1072z-lrwan1/STM32L072CZ
۲۰۰.۴۰۴	۲۰۰.۳۲۳	۲۰۰.۳۶۷	[۲۰۰.۳۸۴, ۲۰۰.۳۳۷]	۶۶.۰۸۰	۱۱.۱۰۰	در ۳۳ مگاهرتز M3	openmote-b/CC2538SF53
۹۷.۷۶۴	۹۷.۷۳۳	۹۷.۷۵۱	[۹۷.۷۵۷, ۹۷.۷۴۰]	۶۲.۲۶۰	۱۱.۲۹۶	در ۷۲ مگاهرتز M3	M3/STM32F103RE آزمايشگاه اینترنت اشيا
۹۸.۶۷۹	۹۸.۶۳۷	۹۸.۶۶۱	[۹۸.۶۶۸, ۹۸.۶۴۹]	۶۳.۱۸۰	۱۱.۲۸۸	در ۴۸ مگاهرتز M4	nrf52840dk/nRF52840
۶۶.۱۶۳	۶۶.۰۸۷	۶۶.۰۸۸	[۶۶.۱۱۲, ۶۶.۰۷۸]	۶۱.۳۳۲	۱۱.۲۴۸	در ۶۴ مگاهرتز M4	nucleo-wl55jc/STM32WL55JC
۵۲.۹۰۲	۵۲.۹۰۰	۵۲.۹۰۱	[۵۲.۹۰۱, ۵۲.۹۰۰]	۶۱.۶۰۴	۱۱.۲۸۸	در ۸۰ مگاهرتز M4	b-1475e-10t01a/STM32L475VG
۳۹.۶۰۴	۳۹.۵۹۹	۳۹.۶۰۱	[۳۹.۶۰۲, ۳۹.۶۰۰]	۶۴.۷۱۲	۱۱.۰۷۶	مگاهرتز 216 @ M7	stm32f746g-disco/STM32F746NG
۸۵.۵۸۴	۸۵.۵۷۶	۸۵.۵۸۲	[۸۵.۵۸۳, ۸۵.۵۸۰]	۱۵۷.۷۱۹	۱۱۵.۹۵۸	پایه کانس ۸۰ مگاهرتز ESP32	hifive1b/SiFive FE310-G002
۵۴.۹۶۱	۵۴.۹۳۸	۵۴.۹۵۳	[۵۴.۹۵۷, ۵۴.۹۴۷]	۲۲۴.۲۷۴	۲۵۸.۸۷۴	در ۸۰ مگاهرتز V-RISC	speedlongan-nano/GD32VF103GBT6
۳۷.۷۹۱	۳۷.۷۷۹	۳۷.۷۸۹	[۳۷.۷۸۹, ۳۷.۷۸۳]	۱۰۶.۴۲۲	۱۰۳.۱۰۸	در ۳۴ مگاهرتز V-RISC	esp32-c3-devkitc-esp32-c3FN4
۱۵۴.۹۲۸	۱۵۴.۷۱۷	۱۵۴.۷۴۷	[۱۵۴.۱۶۶, ۱۵۴.۶۳۱]	۶۶.۴۹۲	۶۰.۸۸۴	در ۱۰ مگاهرتز V-RISC	esp32-wroom-32/ESP32-D0WDQ6

جدول ۱۴: ارزیابی مدل های کوانتیزه شده مختلف روی برد stm32f746-disco

مدل (#پارامترها)	طبقه	حافظه (کیلوبایت)	تعداد پیکسل های ورودی (کیلوبایت)	تأخیر محاسباتی (میلی ثانیه) 95%-CI	میان	میان	مکس
MobileNetV1-0.25x (۵۰ کیلوبایت)	کلمات پیدایش بصری	۱۸۵.۳۵۲	۴۹۱.۶۶۸	[۱۴۳۵.۹۳۸, ۱۴۳۵.۹۳۷]	۱۴۳۵.۹۳۸	۱۴۳۵.۹۳۸	۱۴۳۵.۹۳۹
رمگذاری خودکار عمیق (۲۶۴ کیلوبایت)	تشخیص نااهنجاری	۶.۵۳۲	۳۹۲.۶۹۶	[۳۵.۶۳۸, ۳۵.۶۳۷]	۳۵.۶۳۸	۳۵.۶۳۸	۳۵.۶۳۹
RNNoise (۸۷ هزار)	سرکوب نویز	۴.۶۸۸	۱۱۹.۶۵۲	[۱۲.۱۵۷, ۱۲.۱۵۱]	۱۲.۱۵۴	۱۲.۱۴۸	۱۲.۱۶۰
لی نت-۵ (۴ هزار)	طبقه بندی تصویر	۱۲.۰۶۸	۶۵.۸۵۱	[۳۹.۶۰۳, ۳۹.۵۹۹]	۳۹.۶۰۱	۳۹.۵۹۸	۳۹.۶۰۵
کوچک DS-CNN (۲۲ هزار)	تشخیص کلمات کلیدی	۶۸.۹۹۲	۷۱.۷۹۶	[۴۶۱.۳۹۶, ۴۶۱.۳۹۵]	۴۶۱.۳۹۶	۴۶۱.۳۹۶	۴۶۱.۳۹۷

تعداد پارامترهای مدل در کنار نام مدل ها ارائه شده است

این مدل ها از مخزن بنچمارک های MTPerf Tiny در ... سرچشمه می گیرند. <https://github.com/mlcommons/tiny>

این مدل ها از ARM Model Zoo سرچشمه می گیرند. <https://github.com/ARM-software/ML-zoo> همه مدل ها

در INT8 توسط TFLite از قبل آموزش داده شده و کوانتیزه شدند. به جز LeNet-5 که توسط Pytorch انجام شد.

تحلیل اندازه گیری های RIOT-ML

در این بخش، نتایج آزمایش هایی را که در بخش [مطالب] توضیح داده ایم، ارائه می دهیم. ۸.

۹.۱ ارزیابی عملکرد مدل

ارزیابی هر مدل میز ۳ مصرف منابع مدل LeNet-5 را روی بردهای مختلف اینترنت اشیا، که با ارزیابی هر مدل تولید شده اند، نشان می دهد. MCUها بر اساس خانواده گروه بندی شده و به ترتیب صعودی فرکانس کلک در هر گروه مرتب شده اند. MCUهای سری ARM Cortex-M هیچ تفاوت قابل توجهی در استفاده از حافظه و فضای ذخیره سازی نشان ندادند و با افزایش فرکانس هسته، تأخیر محاسباتی کاهش یافت. یک پرت برد rpi-pico مبتنی بر RP2040 است، اما 16 کیلوبایت رم اضافی در واقع برای اشکال زدارزرزرو شده است. با پشتیبانی کامل از پردازش سیگنال دیجیتال (DSP) و مجموعه دستورالعمل های Thumb-2، میکروکنترلرهای Cortex-M3 و M4- با فرکانس کلک هسته یکسان، عملکرد بهتری نسبت به Cortex-M0+ دارند. یکی دیگر پرت روی SiFive RISC-V MCU کشف شد. با بالاترین فرکانس کلک هسته، این MCU از نظر تأخیر محاسباتی و استفاده از حافظه، نامطلوب ترین رتبه را کسب کرد. این MCU از یک فلش خارجی رابط جانبی سریال NOR (SPI) برای ذخیره سازی داده ها و برنامه ها استفاده می کند که باعث افت شدید عملکرد در هنگام غیرفعال کردن حافظه پنهان می شود.

میز ۴ نتایج مدل های مختلف ML را روی وظایف نماینده TinyML روی بردهای IoT منفرد ارائه می دهد و پشتیبانی جهانی از چارچوب های مختلف ML و ساختارهای مدل را به نمایش می گذارد. به جز LeNet-5 که روی یک دستگاه میزبان محلی با Pytorch آموزش دیده است، بقیه از باغ وحش های مدل متن باز آمده اند. ستون های حافظه و ذخیره سازی به میزان مصرف منابع آنها اشاره دارند.

همانطور که انتظار می رود، مشاهده می کنیم که چگونه مصرف فضای ذخیره سازی متناسب با کاهش تعداد پارامترهای مدل کاهش می یابد. یک مورد پرت آشکار این است که کوچک $DS-CNN$ که فضای ذخیره سازی بیشتری نسبت به ... مصرف می کند. بی نت-۵ اگرچه تقریباً دو برابر پارامترهای مدل دارد. با این حال، تجزیه و تحلیل بیشتر نشان داد که کوچک $DS-CNN$ تقریباً سه برابر فاکتورهای مقیاس بندی بیشتری (در نقطه شناور) در مقایسه با بی نت-۵ این عوامل مقیاس بندی به عنوان پارامترهای مدل محسوب نمی شوند، اما حجم زیادی از فضای ذخیره سازی را اشغال می کنند. همانطور که انتظار می رفت، پیچیده ترین مدل، *موبایل نت وی*، بیشترین حافظه (RAM) و منابع محاسباتی را مصرف کرد.

جدول ۵ خروجی ارزیابی هر اپراتور از مدل سینوس TFlite روی برد stm32f746-disco

اپراتورها	زمان (میکروها)	زمان (%)	پارامز	حافظه	ذخیره سازی
add_nn_relu	۸.۸۵۶	۱۵.۲۲%	ص ۰، ص ۱	۰.۱۲۸	۰.۱۲۸
add_nn_relu_1	۴۶.۶۸۲	۸۰.۲۳%	ص ۲، ص ۳	۰.۱۲۸	۱.۰۸۸
اضافه کردن	۲.۶۴۶	۴.۵۴%	صفحه ۴، صفحه ۵	۰.۰۶۸	۰.۰۶۸

پیشوند تولید شده خودکار `uTVM default_fused_nn_dense` نام عملگر برای شفاف سازی ارائه نشده است. میزان مصرف حافظه و فضای ذخیره سازی بر حسب کیلوبایت ارائه شده است.

با این حال، همانطور که در جدول نشان داده شده است ۴ سربار اجرا، مصرف حافظه و تأخیر محاسباتی را نمی توان صرفاً از روی مقدار پارامترهای مدل به طور کلی به طور قابل اعتمادی پیش بینی کرد. ساختار مدل و الگوی محاسباتی مرتبط با آن، و همچنین اثرات فشرده سازی، نیز به شدت بر سربار اجرا تأثیر می گذارند (از این رو، RIOT-ML به عنوان یک ابزار آزمایشی برای محک زنی مفید است!).

ارزیابی به ازای هر اپراتور ما در اینجا از یک مدل کوچک با تنها سه لایه از TFlite به عنوان نمونه استفاده کردیم تا از پیچیدگی غیر ضروری در نمایش جلوگیری کنیم و نتایج خروجی در جدول ارائه شده است. ۵ گلوگاه های محاسباتی در عملگر قرار دارند `add_nn_relu_1` و با بالاترین میزان مصرف حافظه و فضای ذخیره سازی نیز همراه است. می توانیم لایه های متناظر مدل اصلی را با نشانه هایی از پارامترهای مرتبط، که وزن ها، بایاس یا سایر پارامترهای قابل آموزش مدل هستند، ردیابی کنیم و اعمال استراتژی های بهینه سازی را روی لایه های هدفمند امکان پذیر سازیم.

۹.۲ به روزرسانی های OTA مدل امن

در جدول نشان می دهیم ۶ تفکیک مصرف منابع. از آنجایی که به روزرسانی های کامل و جزئی از یک پشته نرم افزاری مشترک استفاده می کنند، مصرف آنها یکسان باقی می ماند. مشاهده می کنیم که به دلیل تعداد زیاد پارامترهای مدل ثابت (وزن ها و بایاس)، مدل ML مؤلفه ای است که بیشترین فضای ذخیره سازی (حافظه فلش) را استفاده می کند. در مورد استفاده از حافظه (RAM)، پشته شبکه و ماژول SUIT به دلیل نیاز به بافر زیاد برای مدیریت بسته های شبکه و بارهای SUIT، از حافظه نسبتاً بیشتری استفاده می کنند. همچنین مشاهده می کنیم که پیاده سازی های رمزنگاری به کار رفته توسط SUIT (یعنی libcose و c25519) حداقل مصرف حافظه و فضای ذخیره سازی را دارند.

۹.۳ ارزیابی اولیه هزینه های انتقال شبکه

ما علاوه بر این در جدول اندازه گیری کردیم ۶ اندازه مانیفست SUIT و بارهای باینری به روزرسانی که باید برای انجام فرآیند به روزرسانی مدل از طریق شبکه منتقل شوند. اینها

نهاد	حافظه	درصد	ذخیره سازی (انتقال)	درصد
میکروکنترلر	۲۵۶,۰۰۰	-	۱۰۲۴,۰۰۰	-
میان افزار	۴۳,۸۶۴	۱۰۰,۰۰	۱۲۶,۳۷۴	۱۰۰,۰۰
مدل یادگیری ماشینی	۱۲,۴۲۰	۲۸,۳۱	۴۹,۸۵۸	۳۹,۴۵
پشته شبکه	۱۳,۷۱۴	۳۱,۲۶	۳۵,۱۲۱	۲۷,۷۹
کت و شلوار	۱۳,۹۳۳	۳۱,۷۶	۱۵,۳۵۱	۱۲,۱۴
کریپتو	۱,۳۴۰	۳,۰۵	۵,۸۲۶	۴,۶۱
مانیفست SUIT	-	-	۰,۴۷۱	۰,۳۷
لایه نهایی	-	-	۰,۸۴۰	۰,۶۷

در این اندازه گیری ها، ما رمزگذاری SUIT را لحاظ نکردیم، فقط تأیید امضا و بررسی یکپارچگی را در نظر گرفتیم. مصرف حافظه و فضای ذخیره سازی بر حسب کیلوبایت ارائه شده است.

از یک طرف، محققان و متخصصانی که سخت افزار اینترنت اشیا دارند که توسط سیستم عامل متن باز RIOT پشتیبانی می شود (در حال حاضر بیش از ۲۵۰ نوع برد، با استفاده از بیش از ۶۰ نوع CPU) می توانند از RIOT-ML به صورت آماده و مستقیماً روی بردهای خود استفاده کنند. از سوی دیگر، همراه با استفاده از بستر آزمایشی رایگان و آزاد IoT-Lab، حتی محققان و متخصصانی که چنین سخت افزاری را در محل خود ندارند، می توانند با استفاده از RIOT-ML، کمپین های ارزیابی تجربی در مقیاس بزرگ انجام دهند.

دیدگاه ها همچنان که پشتیبانی از بردها و پردازنده های RIOT به مرور زمان گسترش و بهبود می یابد، و همچنان که uTVM نیز به طور موازی پشتیبانی از سایر معماری ها را گسترش می دهد (هر دو جامعه متن باز بسیار فعال هستند)، RIOT-ML می تواند در مدت زمان کوتاهی پشتیبانی خود را برای موارد استفاده جدید گسترش دهد و به طور خودکار پشتیبانی uTVM را برای بردهای جدید و پشتیبانی RIOT را برای مدل های جدید اضافه کند. به این ترتیب، RIOT-ML ممکن است به صورت ارگانیک رشد کند و به یک پیوند مفید بین این دو جامعه تبدیل شود. علاوه بر این، اگرچه کار روی RIOT-ML در این مقاله فقط بر استنتاج روی میکروکنترلرهای تک هسته ای متمرکز بوده است، اما پتانسیل بالایی برای گسترش جعبه ابزار ارائه شده توسط RIOT-ML برای پشتیبانی از سناریوهای یادگیری روی دستگاه و بهینه سازی بهره برداری از میکروکنترلرهای چند هسته ای وجود دارد.

انتیجه گیری

این مقاله RIOT-ML را معرفی می کند، یک جعبه ابزار جدید که ما برای ساده سازی تلاش های متخصصان AIoT طراحی کرده ایم. RIOT-ML ادغام مداوم مدل TinyML، استقرار مداوم ایمن و ارزیابی عملکرد را از راه دور بر روی شبکه، روی مجموعه ای از دستگاه های IoT ناهمگن تحت آزمایش، تسهیل می کند. به طور دقیق تر، با RIOT-ML، کاربران می توانند خروجی مدل Zoo را توسط چارچوب های مختلف یادگیری ماشین سنتی (مانند TensorFlow، PyTorch) دریافت کرده و خودکارسازی کنند.

۹۰ دیدن <https://github.com/RIOT-OS/RIOT/tree/master/boards>

۱۰ ببینید <https://www.iot-lab.info/>

جدول ۶: مصرف فضای ذخیره سازی و هزینه های انتقال به روزرسانی های مدل روی برد nrf52840dk با مدل 5- LeNet مستقر شده

اندازه گیری های ما می توانند به سنجش بار وارده بر شبکه کمک کنند. در مورد بارهای داده، اندازه ی میان افزار سریال شده و وزن لایه ی نهایی به ترتیب ۱۳۳.۴۴ کیلوبایت و ۰.۸۴ کیلوبایت است. در مقام مقایسه، اندازه ی مانیفست SUIT برابر با ۰.۴۷۱ کیلوبایت است که نسبت تقریبی ۲۶۰:۲۰۱ را نشان می دهد. این به ویژه بدان معناست که سربار بار شبکه ناشی از استفاده از مکانیسم های امنیتی SUIT در RIOT-ML (برای یکپارچگی، احراز هویت و مجوز) هنگام استفاده از رویکرد به روزرسانی میان افزار برای به روزرسانی مدل ها، بسیار کمتر از ۱٪ باقی می ماند.

بحث در مورد جزئیات به روزرسانی مدل ساده ترین رویکرد برای به روزرسانی مدل از طریق هوا در دستگاه های اینترنت اشیا، به روزرسانی های میان افزار است، زیرا سیستم تعبیه شده از پیچیدگی های ناشی از بارگذاری پویا روی سخت افزار ناهمگن رهایی می یابد. با این حال، این راحتی با هزینه بالایی از نظر هزینه های انتقال شبکه همراه است. اندازه گیری ها در جدول ۶ اشاره می شود که به روزرسانی فقط مدل (۴۹ کیلوبایت) به جای کل میان افزار (۱۲۶ کیلوبایت) بیش از ۶۰٪ در هزینه های انتقال شبکه صرفه جویی می کند. فراتر از این، اگر در برخی موارد به روزرسانی های مدل بتوانند از یادگیری انتقالی عبور کنند، یعنی به روزرسانی فقط بخش هایی از مدل، هزینه های انتقال شبکه نسبتاً ناچیز می شود. به عنوان مثال، در اندازه گیری های ما، ما فقط آخرین لایه را به روزرسانی کردیم که بیش از ۹۹٪ از هزینه های انتقال شبکه را حذف کرد. با این حال، با این رویکرد، بهبود دقت استنتاج با به روزرسانی ممکن است محدودتر باشد. بنابراین، کاربران موظفند هنگام طراحی طرح به روزرسانی مدل خود، بین سربارهای شبکه، پیچیدگی فنی و دقت مدل تعادل برقرار کنند.

قابل تکرار و سفارشی RIOT-ML ۱۰ آزمایش ها

ماکد منبع کامل ابزار RIOT-ML را در Github به آدرس زیر منتشر کردیم. <https://github.com/TinyPART/RIOT-ML> مجوز متن باز LGPL نسخه ۳، برای جزئیات بیشتر در مورد نحوه شروع کار عملی با RIOT-ML، خواننده به راهنمای جامع ارجاع داده می شود. [Readme.md](#) در مخزن.

2. شارما اچ، حق ای، بلاجرگ اف (2021) یادگیری ماشین در شبکه های حسگر بی سیم برای شهرهای هوشمند: یک بررسی. الکترونیک 10(9): 1012.
3. سانچز-ایورا آر، اسکارمتا ای اف (2020) اشیاء هوشمند مقرون به صرفه با قابلیت TinyML: چالش ها و فرصت ها. مجله مدارهای 18-4(3): 20 IEEE Syst
4. ری پی پی (2022) مروری بر TinyML: جدیدترین ها و چشم اندازها. مجله علوم کامپیوتر دانشگاه ملک سعود 34(4): 1595-1623
5. هوانگ زد، زندهرگ کی، شلازیر کی، باچلی ای (2023) U-TOE: جعبه ابزار ارزیابی جهانی TinyML برای اینترنت اشیا کم مصرف. در سال ۲۰۲۳، دوازدهمین کنفرانس بین المللی ifip/ieee در مورد ارزیابی و مدل سازی عملکرد در شبکه های سیمی و بی سیم (pewwn)، IEEE، صفحات ۱ تا ۶
6. ساها اس اس، ساندها اس اس، سریواستوا ام (2022) یادگیری ماشین برای سخت افزار کلاس میکروکنترلر - بررسی. IEEE Sensors 22(22): 21362-21390
7. هام او، باچلی ای، پیترسن اچ، تیسفتس ان (2015) سیستم های عامل برای دستگاه های ارزان قیمت در اینترنت اشیا: یک بررسی. IEEE Internet Things J 3(5): 720
8. باچلی ای، گوندوغان سی، هام او، کیتزمن پی، لندرز ام اس، پترسن اچ، والیش ام (2018) RIOT: یک سیستم عامل متن باز برای دستگاه های تعبیه شده سطح پایین در اینترنت اشیا. IEEE Internet Things J 5(6): 4428-4440
9. لیما ال ای، کیمورا بی ال، راست وی (2019) محیط های آزمایشی برای اینترنت اشیا: یک بررسی. IEEE Sens J 19(9): 3203-3211
10. Gluhak A, Krco S, Nati M, Pfisterer D, Mitton N, Razafindralambo T (2011) نظرسنجی در مورد امکانات برای تحقیقات آزمایشی اینترنت اشیا (2011 IEEE Commun Mag 49): 58-67
11. Harter G, Mitton N, Noel T, Vandaele J (2015) مقیاس بزرگ: دز دومین انجمن جهانی IoT یک بستر آزمایشی باز: LAB-FIT IoT و همکاران (2015) Adjih C, Baccelli E, Fleury E
12. بنبری سی آر، ردی وی جی، لام ام، فو دلیو، فاضل ای، هولمن جی، لوخمو توف ای و همکاران (۲۰۲۰) ارزیابی سیستم های TinyML: چالش ها و مسیرها. پیش از چاپ در arXiv:2003.04821
13. Banbury C, Reddi VJ, Torelli P, Holleman J, Jeffries N, Kiraly S, پائو دی و همکاران (2021) معیار کوچک MLPerf. پیش چاپ در arXiv:2106.07597
14. عثمان آ، عبید یو، جما ال، پروتو ام، برنولی دی (2022) بنچمارک پلتفرم های TinyML: در کاربردهای الکترونیک در صنعت، محیط زیست و جامعه: Applepies 2021، Springer: صفحات 139-148
15. Sudharsan B, Salerno S, Nguyen DD, Yahya M, Yada V, علی ام آی (۲۰۲۱) بنچمارک TinyML: اجرای شبکه های عصبی کاملاً متصل روی میکروکنترلرهای معمولی. در: هفتمین مجمع جهانی IEEE در مورد اینترنت اشیا (wf-iot) در سال ۲۰۲۱، IEEE، صفحات ۸۸۳-۸۸۴
16. چن تی، مورو تی، جیانگ زد، ژنگ ال، یان ای، شن اچ، چز ال و همکاران (۲۰۱۸) TVM: یک کامپایلر بهینه سازی سرتاسری خودکار برای یادگیری عمیق، ۵۷۸-۵۹۴
17. کیو اچ، واولیدو آی، لی جی، پرگامنت ای، واردن پی، چینچالی اس، کاتی اس (2022) ML-EXray: قابلیت مشاهده در استقرار ml در لبه. Proc Mach Learn Syst 4: 337
18. یوسف زاده اصل-میانداوب ای، روبروک تی، توزون پی (2023) پروفایل بندی و نظارت بر وظایف آموزشی یادگیری عمیق. در: مجموعه مقالات سومین کارگاه آموزشی یادگیری ماشین و سیستم ها، صفحات 18-25. <https://doi.org/10.1145/3578356.3592589>
19. Kreuzberger D, Kühl N, Hirschl S (2023) عملیات یادگیری ماشین (MLops): مرور کلی، تعریف و معماری. IEEE Access 11: 31866-31879. <https://doi.org/10.1109/ACCESS.2023.326213819>
20. لی ام تی، آریل جی (2023) TinyMLops برای MCUهای فوق العاده کم مصرف پلادرنگه برای طبقه بندی رویدادهای مبتنی بر فریم اعمال می شوند. در: مجموعه مقالات سومین کارگاه آموزشی یادگیری ماشین و سیستم ها، نیویورک، نیویورک، ایالات متحده آمریکا: انجمن ماشین آلات محاسباتی، صفحات 148-153

تطبیق پذیری، استقرار و ارزیابی عملکرد آنها هنگام اجرا بر روی طیف گسترده ای از بردهای اینترنت اشیا و کیت های توسعه مبتنی بر انواع مختلف میکروکنترلرهای کم مصرف (RISC-V، ARM Cortex-M، ESP32)، فعال سازی و تسهیل چنین ماتریس های آزمایشی بزرگی در واقع برای پیشرفت حوزه AIoT بسیار مهم است. برای این منظور، ما همچنین یک پیاده سازی متن باز بسیار قابل استفاده مجدد، مستندسازی شده و قابل تنظیم از RIOT-ML منتشر کردیم که از جوامع متن باز پویا مرتبط با RIOT و uTVM بهره می برد. در آخر، ما با ارائه نتایج ارزیابی تجربی اولیه در یک بستر آزمایشی با دسترسی آزاد، کاربرد عملی RIOT-ML را به نمایش می گذاریم.

تقدیرنامه هانویسندگان ما بلند از سد ریک آچی، نجیب آچیر و فلیکس بیسمن برای بحث ها و پیشنهادهای مفیدشان تشکر کنند.

سهم نویسندگان نرم افزار را توسعه داده و آزمایش کردند. همه نویسندگان مقاله را بررسی کردند KS و KZ، متن اصلی مقاله را مفهوم سازی و نوشتند: EB و ZH

بودجه بودجه دسترسی آزاد توسط Projekt DEAL فراهم و سازماندهی شده است. تحقیقاتی که منجر به این نتایج شده است، تا حدودی از برنامه امنیت سایبری آلمانی /فرانسوی MESRI-BMBF تحت توافقنامه اعطای کمک هزینه شماره CYAL-0005 و ANR-20-16KIS1395K بودجه دریافت کرده است.

در دسترس بودن داده هادر طول مطالعه فعلی هیچ مجموعه داده ای تولید یا تجزیه و تحلیل نشد.

اعلامیه ها

تضاد منافع نویسندگان هیچ گونه تضاد منافی را اعلام نمی کنند.

سلب مسئولیت این مقاله صرفاً منعکس کننده دیدگاه های نویسندگان است. MESRI و BMBF مسئولیتی در قبال هرگونه استفاده ای که ممکن است از اطلاعات موجود در آن صورت گیرد، ندارند.

دسترسی آزاد این مقاله تحت مجوز بین المللی Attribution 4.0 Creative Commons منتشر شده است که استفاده، اشتراک گذاری، اقتباس، توزیع و تکثیر در هر رسانه یا قالبی را مجاز می داند، مادامی که به نویسنده (گان) اصلی و منبع، اعتبار کافی داده شود، پیوندی به مجوز Creative Commons ارائه شود و در صورت ایجاد تغییرات، مشخص شود که آیا تغییراتی ایجاد شده است یا خیر. تصاویر یا سایر مطالب شخص ثالث در این مقاله در مجوز Creative Commons مقاله گنجانده شده اند، مگر اینکه در خط اعتباری مطلب، خلاف آن ذکر شده باشد. اگر مطلبی در مجوز Creative Commons مقاله گنجانده نشده باشد و استفاده مورد نظر شما طبق مقررات قانونی مجاز نباشد یا از میزان مجاز تجاوز کند، باید مستقیماً از دارنده حق چاپ اجازه بگیرید. برای مشاهده نسخه ای از این مجوز، به <http://creativecommons.org/licenses/by/4.0/>

منابع

1. سزه وی، چن وای اچ، یانگ تی جی، امر جی اس (2017) پردازش کارآمد شبکه های عصبی عمیق: یک آموزش و بررسی. Proc IEEE 12(2295-2329): 105

28. موران بی، شوفنیگ آچ، براون دی، مریاک ام (2021) معماری به روزرسانی میان افزار برای اینترنت اشیا. RFC 9019. <https://doi.org/10.17487/RFC9019>
۲۹. هانسن تی سوم، دی ای ای (۲۰۰۶) الگوریتم های هش امن ایالات متحده (SHA و HMACSHA). RFC ۴۶۳۴. <https://doi.org/10.17487/RFC4634>
- EdDSA (RFC8032). (۲۰۱۷) I. <https://doi.org/10.17487/RFC8032>. Josefsson S, Liusvaara
۳۱. تسوفنیگ آچ، هاسلی آر، موران بی، براون دی، تاکایاما کی (۲۰۲۳) بارهای رمزگذاری شده در مانیفست های SUIT (شماره پیش نویس اینترنت، 18-draft-ietf-suit-firmwareencryption). کار در حال انجام. گروه ویژه مهندسی اینترنت. برگرفته از <https://datatracker.ietf.org/doc/draft-ietf-suit-firmware-encryption/18/>
۳۲. سلندر جی، ماتسون جی پی، پالومبینی اف (۲۰۲۴) دیفی-هلمن زودگذر بر فراز (RFC ۹۵۲۸). <https://doi.org/10.17487/RFC9528>. COSE (EDHOC)
33. کاهان دلیو (1996) استاندارد IEEE 754 برای حساب ممیز شناور دودویی. یادداشت های سخنرانی وضعیت 11: 1776-94720 (IEEE 754)
34. Schaad J (2017) (CBOR رمزگذاری شی). <https://doi.org/10.17487/RFC8152>. COSE (RFC8152)
21. دیاز-دآرکایا جی، توره-باستیدا ای آی، زاراتی جی، میون آر، آلمیدا ای (2023) مطالعه مشترک چالش ها، فرصت ها و نقشه راه MLOps و AIops: یک بررسی سیستماتیک. ACM Comput Surv 56(4):1-30. <https://doi.org/10.1145/3625289>
22. زندبرگ ک، شلایزر ک، آکوستا ف، شوفنیگ ه، باجلی ای (2019) به روزرسانی های امن میان افزار برای دستگاه های اینترنت اشیا محدود با استفاده از استانداردهای باز: بررسی واقعیت. IEEE Access 7:71907-71920
23. منتوا ان اس، تارویری پی، ابو-محمود ای ام، آدیگان ام او (2019) به روزرسانی های امن میان افزار در اینترنت اشیا: یک بررسی. در کنفرانس بین المللی فناوری اطلاعات و مهندسی چندرشته ای 2019 (imitec). IEEE, صفحات 7-1
24. بانگاس جی، زندبرگ کی، باجلی ای، هرمان ای، اسمیت بی (2022) امنیت به روزرسانی نرم افزار مقاوم در برابر کوانتوم در دستگاه های تعبیه شده شبکه ای کم مصرف. در کنفرانس بین المللی رمزنگاری کاربردی و امنیت شبکه، اسپرینگر، صفحات 872-891
- 41-35(2) IEEE Commun Mag 58(2):35-41. به روزرسانی های نرم افزار هوایی در اینترنت اشیا: مروری بر اصول کلیدی (2020) I, De Poorter E
25. Bauwens J, Ruckebusch P, Giannoulis S, Moerman
26. بورمن سی، کاستلانی ای پی، شلیبی زد (2012) CoAP: یک پروتکل کاربردی برای میلیاردها گره کوچک اینترنتی. Comput 16(2):62-67. IEEE Internet
- برای نظرات IETF درخواست CoAP(پروتکل برنامه محدود
27. ShelbyZ, HartkeK, BormannC (2014) RFC7252:

یادداشت ناشر: اسپرینگر نیچر در مورد ادعاهای مربوط به صلاحیت قضایی در نقشه های منتشر شده و وابستگی های سازمانی بی طرف باقی می ماند.