

محاسبه مقدار آسیب‌پذیری مدارهای مجتمع دیجیتال در برابر تروجان سخت‌افزاری با استفاده از شبکه‌های

عصبی پیچشی

هادی جهانی راد^{*۱}

گروه مهندسی برق، الکترونیک و مخابرات، دانشکده مهندسی، دانشگاه کردستان، سنندج، ایران

چکیده

با پیشرفت تراشه‌های مجتمع دیجیتال و پیاده‌سازی سیستم‌های پیچیده بر روی آن‌ها، مخاطراتی در رابطه با عملکرد آن‌ها ایجاد شده است. تروجان‌های سخت‌افزاری (HT) از مهمترین نوع مخاطرات هستند که سبب ایجاد خطا در عملکرد تراشه، افزایش توان مصرفی و نشت کردن اطلاعات ذخیره‌شده بر روی تراشه‌ها می‌شوند. در نتیجه، ارزیابی میزان آسیب‌پذیری تراشه‌ها در برابر انواع مختلف تروجان‌های سخت‌افزاری دارای اهمیت بسیار زیادی است. در این مقاله روشی دقیق در سطح چیدمان (layout)، بر مبنای استفاده از شبکه‌های عصبی پیچشی (CNN)، برای محاسبه میزان آسیب‌پذیری تراشه‌های دیجیتال در برابر تروجان‌های سخت‌افزاری ارائه شده است. عوامل اصلی مؤثر بر میزان خطرپذیری تراشه‌های دیجیتال شامل میزان فضاهای خالی در چیدمان، منابع مسیردهی استفاده نشده، فعالیت سیگنال‌های داخلی، و قابلیت آزمون-پذیری گیت‌های مدار می‌باشند. برای تولید مجموعه داده مناسب، چیدمان فیزیکی هر پیاده‌سازی از یک مدار دیجیتال با استخراج این عوامل، به یک تصویر دیجیتال تبدیل شده است. پس از تولید مجموعه داده مناسب که شامل ۱۰۰۰۰ تصویر است، فرآیند یادگیری شبکه عصبی پیچشی تکمیل می‌شود و شبکه آموزش یافته برای تعیین میزان آسیب‌پذیری در برابر تروجان‌های سخت‌افزاری مورد استفاده قرار می‌گیرد. با مشخص شدن میزان آسیب‌پذیری مدار پیاده‌سازی شده، طراح می‌تواند تغییرات لازم را برای مقاوم کردن تراشه در برابر تروجان‌های سخت‌افزاری اعمال نماید. نتایج شبیه‌سازی بر روی چیدمان مدارهای معیار (ISCAS 85, 89) نشان می‌دهد، میزان دقت رهیافت پیشنهادی ۹۲٪ می‌باشد. همچنین روش پیشنهادی، مشکل ناشی از عدم مدل‌سازی دقیق عوامل مؤثر در تعیین خطرپذیری در روش‌های پیشین را مرتفع کرده و دقت محاسبه آسیب‌پذیری را نسبت به بهترین روش موجود در مطالعات پیشین، ۱۷٪ افزایش می‌دهد.

مقاله پژوهشی

تاریخ دریافت:

۱۴۰۲/۱۲/۱۷

تاریخ پذیرش:

۱۴۰۳/۴/۱۶

کلیدواژه‌ها:

مدارهای مجتمع دیجیتال،
تروجان‌های سخت‌افزاری،
شبکه‌های عصبی پیچشی، کلاس
بندی، یادگیری ماشین

نویسنده مسئول:

h.jahanirad@uok.ac.ir

doi : 10.22034/ABMIR.2024.21351.1049

۱- مقدمه

می‌باشد. He و همکارانش در مرجع [۳] از تحلیل اثرات جانبی الکترومغناطیسی ناشی از تروجان‌های سخت‌افزاری برای آشکارسازی استفاده کرده‌اند.

در مرجع [۴] از اطلاعات فاز در اسکن حرارتی فروسرخ برای ایجاد الگوی حرارتی تراشه استفاده شده است. با بهره‌گیری از این روش، نقاط داغ با ابعاد کوچک (که می‌توانند ناشی از فعالیت تروجان‌ها باشند) قابل تشخیص هستند. Koushanfar و Mirhosseini در مرجع [۵] یک روش چند حالتی با ترکیب الگوهای حرارتی و توانی برای آشکارسازی توسعه داده‌اند. در مرجع [۶] از روش تحلیل الگوی دمایی (که با تبدیلات لازم از الگوی توانی استخراج می‌شوند) برای آشکارسازی استفاده شده است.

از پردازش تصویر برای آشکارسازی تروجان‌های سخت‌افزاری در پژوهش‌های مختلف استفاده شده است. تصویربرداری SEM از سطح تراشه و تحلیل آن برای آشکارسازی تروجان سخت‌افزاری براساس روش‌های پردازش تصویر، در مرجع [۷] ارائه شده است. Jahanirad و Rahimifar در مرجع [۸] از روش‌های پردازش تصویر برای تحلیل الگوی حرارتی مدار مجتمع در راستای آشکارسازی تروجان‌های سخت‌افزاری استفاده کرده‌اند.

با توجه به تنوع بسیار زیاد تروجان‌های سخت‌افزاری و عدم توانایی کامل روش‌های آشکارسازی برای تشخیص وجود آن‌ها در تراشه نهایی، ضروری است که طراحی مدار مجتمع به گونه‌ای انجام شود که امکان قراردادن تروجان‌های سخت‌افزاری در آن تراشه به حداقل ممکن برسد. برای رسیدن به این هدف، باید ارزیابی میزان آسیب‌پذیری چیدمان طراحی شده IC در مقابل تروجان سخت‌افزاری انجام گیرد. ویژگی‌های متعددی در این ارزیابی موثرند که مهمترین آن‌ها به شرح زیر است: معمولاً طراحان تروجان‌های سخت‌افزاری، هسته سخت‌افزاری خود را در فضاهای خالی از چیدمان تراشه، قرار می‌دهند. از طرف دیگر اتصالات لازم برای تریگر کردن و همچنین اعمال خروجی هسته‌های

جایگذاری تروجان‌های سخت‌افزاری^۱ در مدارهای مجتمع مدرن، برای ایجاد اختلال در عملکرد سیستم و یا مصرف توان بیشتر، به یکی از مخاطرات جدی تبدیل شده است. به‌طور کلی تروجان سخت‌افزاری، یک هسته یا ماژول سخت‌افزاری کوچک است که در مکان‌های خاصی از چیدمان IC^۲ قرار داده می‌شود. فعالسازی تروجان‌های سخت‌افزاری با یک تریگر که متصل به گره‌های با میزان فعالیت کم هستند، انجام می‌شود. آشکارسازی وجود تروجان‌های سخت‌افزاری در تراشه‌های مورد استفاده در سیستم‌های حساس (نظامی، مخابراتی و مانند آن) بسیار با اهمیت است. روش‌های آشکارسازی، به سه دسته اصلی تقسیم می‌شوند. دسته اول از روش‌های آشکارسازی، بر مبنای آزمون‌های منطقی عمل می‌کنند و برای آشکارسازی تروجان‌های همیشه-فعال، کارایی بیشتری دارند. دسته دوم از روش‌های آشکارسازی، به تحلیل اثرات جانبی (side channel) تروجان سخت‌افزاری، مانند بوجود آمدن نقاط داغ^۳ در اثر فعالیت تروجان، افزایش توان مصرفی تراشه و یا افزایش تأخیر مسیر بحرانی مدار، می‌پردازند. راهکارهایی که در این دسته قرار می‌گیرند، بر مبنای تحلیل الگوی حرارتی تراشه، الگوی توان مصرفی تراشه، و ابزارهای یادگیری ماشین استوارند. دسته سوم براساس مهندسی معکوس تراشه دیجیتال به آشکارسازی تروجان‌های سخت‌افزاری می‌پردازند. این دسته از روش‌ها، از جمله روش‌های مخرب بوده و تنها در موارد خاصی برای تولید تراشه‌های طلایی مورد استفاده قرار می‌گیرند.

تعدادی از پژوهش‌های مهم در زمینه آشکارسازی تروجان‌های سخت‌افزاری به شرح زیر است. در مرجع [۱] از روش مهندسی معکوس در سطح چیدمان و گیت، برای آشکارسازی استفاده شده است که بسیار زمان‌بر و پرهزینه است. روشی دیگر در مرجع [۲] بر مبنای قرار دادن حسگرهای حرارتی در بخش‌های مختلف IC ارائه شده است که در زمان اجرا (runtime) با استفاده از اطلاعات حسگرها وجود تروجان‌ها را آشکار می‌کند. این روش پرهزینه است و دقت آن به شدت وابسته به دقت بودن اطلاعات حسگرها

³ Hotspots

¹ Hardware Trojans

² Integrated Circuits

رئوس نوآوری‌های این پژوهش را می‌توان به صورت ذیل خلاصه کرد:

- تولید مجموعه داده مناسب و ارائه روش تولید آن برای مطالعات مربوط به بررسی خطرپذیری مدارهای دیجیتال.
 - ارائه راهکار مناسب برای استخراج ویژگی‌های اصلی و مؤثر در تعیین میزان خطرپذیری برای مدارهای دیجیتال.
 - توسعه یک چارچوب نرم‌افزاری برای تعیین کلاس خطرپذیری در برابر تروجان‌های سخت‌افزاری در مدارهای دیجیتال بر مبنای شبکه‌های عصبی پیچشی.
- در ادامه مقاله مطالب زیر ارائه می‌شوند. در بخش دوم مفاهیم اولیه برای ورود به بحث توضیح داده می‌شود. در بخش سوم، رهیافت پیشنهادی برای ارزیابی آسیب‌پذیری مدارهای مجتمع دیجیتال در سطح چیدمان با استفاده از شبکه‌های عصبی کانولوشن ارائه می‌شود. نتایج شبیه‌سازی و مقایسه با روش‌های [۹-۱۱] در بخش ۴ ارائه می‌شود و سپس در بخش نتیجه‌گیری جمع‌بندی مقاله انجام می‌گردد.

۲- محاسبه ویژگی‌های چیدمان تراشه

همان‌گونه که در بخش پیش توضیح داده شد، گام اول در ارزیابی آسیب‌پذیری مدار مجتمع در برابر تروجان‌های سخت‌افزاری، محاسبه ویژگی‌های مؤثر در قراردادان هسته‌های تروجان سخت‌افزاری در چیدمان تراشه است. این ویژگی‌ها شامل فضاهای خالی چیدمان، تراکم منابع مسیره‌دهی، کنترل‌پذیری و فعالیت سیگنال مربوط به گره‌های مدار و مسیرهای غیربحرانی (از لحاظ تأخیر) می‌باشد. در زیربخش‌های آینده، نحوه محاسبه هرکدام از این ویژگی‌ها توضیح داده می‌شود.

۲-۱ فضاهای خالی در چیدمان تراشه

برای تولید چیدمان مدار مجتمع، ابتدا کد HDL^۱ مدار در سطح گیت سنتز می‌شود و پس از آن نرم‌افزار Cadence Design System، با استفاده از گیت‌های معیار (کتابخانه‌ای) مدار سنتز شده را به چیدمان بهینه تبدیل می‌کند. به صورت کلی، چیدمان نهایی مدارهای مجتمع به یک شبکه منظم $M \times M$ از واحدهای مربع شکل

تروجان‌های سخت‌افزاری به بخش‌های دیگر مدار مجتمع، از طریق منابع استفاده نشده در شبکه مسیره‌دهی تراشه انجام می‌گیرد. تریگر کردن هسته‌های سخت‌افزاری تروجان‌ها از طریق اتصال به گره‌هایی از مدار که فعالیت سیگنالی پایین دارند، انجام می‌شود. این کار به طراحان تروجان‌های سخت‌افزاری کمک می‌کند که از روش‌های آشکارسازی بر مبنای آزمون‌های منطقی، مصون بمانند. همچنین برای عبور از آزمون‌های بر مبنای تأخیر مسیر، تروجان‌ها باید به‌گونه‌ای در تراشه قرار گیرند که تأخیر مسیر بحرانی را افزایش ندهند. بنابراین وجود مسیرهای غیر بحرانی زیاد در ساختار مدار مجتمع می‌تواند آن را در برابر تروجان‌های سخت‌افزاری آسیب‌پذیرتر کند.

با توجه به وجود عوامل متعدد در ارزیابی آسیب‌پذیری مدارهای مجتمع دیجیتال در برابر تروجان‌های سخت‌افزاری، توسعه دادن یک مدل فراگیر و کارا که بتواند انواع مختلف حالت‌های ممکن را پوشش دهد، بسیار پیچیده است. در روش‌های پیشین [۹-۱۱] ابتدا مقادیر فضای خالی، تراکم منابع مسیره‌دهی، فعالیت سیگنال در گره‌های داخلی مدار محاسبه می‌شود. سپس با استفاده از روابط ساده ریاضی (ضرب کردن مقادیر بدست آمده در مرحله اول) میزان آسیب‌پذیری بخش‌های مختلف مدار محاسبه می‌شود. مهم‌ترین اشکال این روش‌ها آن است که ترکیب معیارهای مختلف، به صورت ناقص انجام شده است و همان‌طور که در بخش‌های بعدی نشان داده می‌شود این امر به نتایج نادرست می‌انجامد.

در این مقاله روشی کارا برای ارزیابی آسیب‌پذیری مدارهای مجتمع دیجیتال در برابر تروجان‌های سخت‌افزاری ارائه می‌شود که در آن ابتدا سطح چیدمان مدار به صورت یک شبکه توری (با سوراخ‌های مربعی شکل) بخش‌بندی می‌شود. سپس برای هرکدام از بخش‌ها معیارهای اصلی محاسبه می‌شود و ترکیب آن‌ها منجر به تولید یک تصویر دیجیتال برای اعمال به شبکه عصبی پیچشی می‌گردد. استفاده از ظرفیت یادگیری بالای شبکه‌های عصبی کانولوشن منجر به یادگیری روابط و قضایای پیچیده بین عوامل مؤثر در میزان خطرپذیری تراشه‌های دیجیتال خواهد شد و در نتیجه نقص‌های مدلسازی در روش‌های پیشین به صورت مناسب برطرف می‌گردد.

^۱ Hardware Description Language

۲-۲ تراکم منابع مسیره‌دهی در چیدمان تراشه

همان‌گونه که پیشتر اشاره شد، هسته‌های سخت‌افزاری مربوط به تروجان‌های سخت‌افزاری برای فعال شدن و تأثیر مخرب بر عملکرد مدار نیازمند استفاده از منابع مسیره‌دهی هستند. اطلاعات مربوط به منابع مسیره‌دهی موجود در چیدمان تراشه، از فایل‌های خروجی نرم‌افزار Layout Editor مانند Design Exchange Format و Library Exchange Format استخراج می‌شوند. در نهایت تراکم هر کدام از واحدهای مربع شکل در چیدمان نهایی تراشه، با استفاده از رابطه ۱ محاسبه می‌شود. در این رابطه $U_r(m,n)$ و $A_r(m,n)$ به ترتیب مقدار منابع مسیره‌دهی استفاده‌شده و تعداد منابع مسیره‌دهی موجود در واحد (m,n) هستند. هرچه تعداد سیم‌های بخش مسیره‌دهی در فضای بالایی واحد مربع شکل بیشتر باشد، امکان پیدا کردن منابع استفاده نشده برای تریگر کردن و اتصال خروجی‌ها بیشتر می‌شود.

$$Cong(m,n) = \frac{U_r(m,n)}{A_r(m,n)} \quad (1)$$

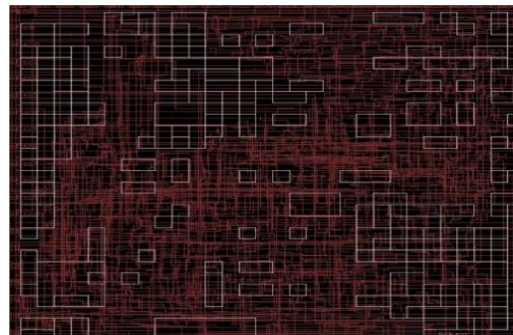
۲-۳ فعالیت سیگنال در گره‌های مدار

تریگر کردن هسته‌های تروجان سخت‌افزاری با استفاده از گره‌هایی از مدار انجام می‌شود که فعالیت پایین‌تری دارند. منظور از فعالیت سیگنال در یک گره، تعداد متوسط 0 و 1 شدن‌های سیگنال در یک بازه زمانی مشخص است. در تحلیل مدارهای مجتمع دیجیتال معمولاً برای بیان فعالیت سیگنال از مفهوم چگالی انتقال سیگنال استفاده می‌شود که طبق رابطه ۲ تعریف می‌شود. در این رابطه $n_x(T)$ تعداد انتقال از 0 به 1 و از 1 به 0 در گره x از مدار در مدت زمان T می‌باشد.

$$D(x) = \lim_{T \rightarrow \infty} \frac{n_x(T)}{T} \quad (2)$$

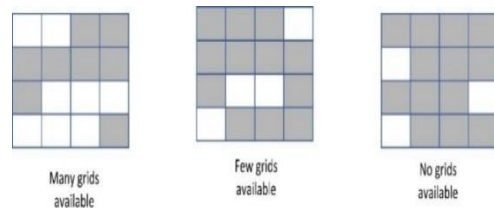
در روش پیشنهادی این مقاله، محاسبه چگالی انتقال سیگنال مطابق روند نشان داده در شکل (۳) انجام می‌شود. در مرحله اول، کد Verilog مدار موردنظر وارد سنتزگر ODIN II می‌شود که

تقسیم‌بندی می‌شود. به‌عنوان مثال چیدمان بهینه مدار معیار b15 یک سطح مربعی $300\mu m \times 300\mu m$ است. بنابراین هرکدام از واحدهای شبکه با انتخاب $M=10$ ، یک سطح مربعی $30 \times 30 \mu m$ می‌باشد. با فرض اینکه حداقل مساحت خالی که می‌تواند بعنوان فضای خالی در نظر گرفته شود، برابر مساحت چیدمان یک گیت وارونگر (مثلاً INVX0 در کتابخانه SAED-EDK90nm) باشد، می‌توان توزیع فضاهای خالی را در چیدمان نهایی استخراج کرد. در شکل ۱، این توزیع در مورد مدار معیار EthernetMAC10GE نشان داده شده است. معمولاً هسته‌های سخت‌افزاری تروجان‌های سخت‌افزاری مساحتی بیشتر از مساحت یک واحد فضای خالی (Unit-Space) نیاز دارد. بعنوان مثال فرض کنید یک هسته سخت‌افزاری تروجان، به مساحتی ۵ برابر Unit-Space نیاز دارد. در شکل ۲، سه واحد مربع شکل از چیدمان تراشه نشان داده شده است که در آن‌ها فضاهای خالی به صورت‌های مختلف توزیع شده‌اند. با توجه به تعداد فضاهای خالی مورد نیاز برای پیاده‌سازی تروجان‌های سخت‌افزاری، تنها واحد مربع شکل در شکل سمت چپ برای قرار دادن تروجان‌های سخت‌افزاری موردنظر قابل استفاده می‌باشد.



شکل (۱): فضاهای خالی در چیدمان مدار

EthernetMAC10GE



شکل (۲): میزان سختی وارد کردن تروجان به مدار در شرایط مختلف

۲-۵ کنترل‌پذیری گره‌های مدار

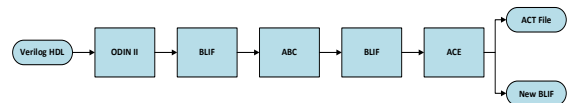
تریگر کردن تروجان و اتصال آن به سایر بخش‌های مدار، باید در گره‌هایی صورت گیرد که کنترل‌پذیری و مشاهده‌پذیری کمی داشته باشند. این امر سبب می‌شود که روش‌های آشکارسازی برمبنای آزمون، با احتمال کمتری امکان تشخیص هسته‌های تروجان را داشته باشند. برنامه SCOAP^۳ که به صورت گسترده در آزمون مدارهای دیجیتال به کار برده می‌شود، برای محاسبه میزان کنترل‌پذیری گره‌های داخلی مدار، مورداستفاده قرار می‌گیرد [۱۲].

در این روش، مدار سطح‌بندی می‌شود و با شروع از ورودی‌های اصلی مدار محاسبه کنترل‌پذیری گیت‌های منطقی با توجه به مقدار کنترل‌پذیری ورودی‌های آن گیت انجام می‌شود. منظور از کنترل‌پذیری 0 (CC0) و کنترل‌پذیری 1 (CC1) آن است که ایجاد منطق 0 و 1 در خروجی گیت با استفاده از تنظیم منطق ورودی‌های گیت، به چه اندازه مشکل است. در جدول ۱ نحوه محاسبه CC0 و CC1 برای تعدادی از گیت‌های پایه‌ای ارائه شده است. بعد از اعمال SCOAP بر روی مدار دیجیتال و تعیین مقدار کنترل‌پذیری گره‌های مختلف مدار، می‌توان برای واحدهای مربع شکل چیدمان تراشه معیار کنترل‌پذیری را با میانگین گرفتن از مقادیر کنترل‌پذیری گره‌های موجود در آن واحد محاسبه کرد.

جدول (۱): نحوه محاسبه کنترل‌پذیری چهار گیت منطقی

	CC0	CC1
AND	$\min(\text{CC0}(x_1), \text{CC0}(x_2))+1$	$[\text{CC1}(x_1)+\text{CC1}(x_2)]+1$
OR	$[\text{CC0}(x_1)+\text{CC0}(x_2)]+1$	$\min(\text{CC1}(x_1), \text{CC1}(x_2))+1$
XOR	$\min([\text{CC0}(x_1)+\text{CC0}(x_2), (\text{CC1}(x_1)+\text{CC1}(x_2))]+1$	$\min([\text{CC0}(x_1)+\text{CC1}(x_2), (\text{CC1}(x_1)+\text{CC0}(x_2))]+1$
BUFFER	$\text{CC0}(x)+1$	$\text{CC1}(x)+1$

خروجی آن یک فایل در قالب BLIF^۱ می‌باشد. در مرحله دوم، فایل تولیدشده، توسط نرم‌افزار ABC^۲ مورد پردازش قرار می‌گیرد که هدف از آن انجام بهینه‌سازی و نگاشت مدار سنتز شده بر روی تراشه با استفاده از تکنولوژی ساخت موردنظر می‌باشد. خروجی مرحله دوم که در قالب فایل BLIF می‌باشد با استفاده از نرم‌افزار ACE 2 جهت استخراج مقدار چگالی انتقال (فعالیت) گره‌های داخلی مدار، مورد ارزیابی قرار می‌گیرد. نتایج این مرحله در فایل ACT ذخیره می‌گردد.



شکل (۳): مراحل به دست آوردن فایل مربوط به فعالیت سویچ‌ها

۲-۴ محاسبه تأخیر مسیرها

برای محاسبه تأخیر مسیرهای موجود از ورودی‌ها به خروجی‌های مدار، ابتدا گراف جهت‌دار مدار تولیدشده و با استفاده از ویژگی‌های الکترونیکی گیت‌ها، به هرکدام از یال‌های گراف مقدار عددی تأخیر انتساب داده می‌شود. تأخیر هر مسیر از ورودی به خروجی، با جمع کردن تأخیر یال‌های متعلق به آن مسیر (d_e) محاسبه می‌شود (رابطه ۳). در رابطه (۳)، P نمایانگر مسیر موردنظر و d_e بیانگر تأخیر مربوط به جزء e از مسیر P می‌باشد.

$$\text{delay}(P) = \sum_{e \in P} d_e \quad (3)$$

وقتی تروجان سخت‌افزاری به یک مسیر از مدار اضافه می‌شود، مقدار تأخیر ورودی به خروجی هسته تروجان سخت‌افزاری به تأخیر مسیر موردنظر اضافه می‌شود. در صورتی که این افزایش تأخیر، مسیر غیر بحرانی را به مسیر بحرانی مدار تبدیل کند، آنگاه وجود تروجان بوسیله روش‌های آزمون برمبنای تأخیر آشکار می‌شود. لذا طراح تروجان‌های سخت‌افزاری سعی در شناسایی مسیرهای غیر بحرانی مدار دارد. محاسبه معیار تأخیر مسیر در هرکدام از واحدهای مربع شکل، با شمارش تعداد مسیرهای غیر بحرانی در آن واحد انجام می‌شود.

⁴ Combinational Controlability-0

⁵ Combinational Controlability-1

¹ Berkeley Logic Interchange Format

² A System for Sequential Synthesis and Verification

³ Sandia Controllability Analysis Program

قرار نگیرد ($UR(r)=0$) و یا گره‌ای از مدار که داخل واحد است دارای فعالیت سیگنالی بیش از آستانه موردنظر باشد ($N_{LP}(r)=0$)، مقدار آسیب‌پذیری این نوع واحدها 0 تلقی می‌شود. این نتیجه به صورت آشکار درست نیست: طراح تروجان سخت‌افزاری می‌تواند هسته تروجان‌های سخت‌افزاری را در فضای خالی موجود در واحد موردنظر قرار دهد و برای تریگر کردن آن از گره‌های مدار در واحدهای کناری که دارای فعالیت سیگنال پایین‌تری هستند، استفاده کند. بنابراین وجود فضای خالی بزرگ در یک واحد می‌تواند معیار خطرپذیری بالایی باشد و صرفاً بدلیل عدم وجود تریگر مناسب نمی‌توان از خطر وجود چنین واحدهایی صرفنظر کرد.

خطای دیگری که در روش بر مبنای معادله اتفاق می‌افتد، مربوط به واحدهای پرتراکم در چیدمان تراشه می‌باشد. این واحدها که معمولاً در مرکز چیدمان قرار می‌گیرند، دارای منابع مسیره‌ی زیادی هستند و همچنین به دلیل قرار گرفتن تعداد بیشتری گره مدار در آنها، احتمال وجود گره‌هایی با فعالیت سیگنال پایین (جهت تریگر کردن تروجان‌های سخت‌افزاری) زیاد است. در این صورت، حتی اگر فضای کوچکی از سطح آنها خالی باشد ($WS(r)$ کوچک)، آسیب‌پذیری نهایی این واحدها به دلیل بسیار بزرگ بودن سایر معیارها تعداد منابع مسیره‌ی بلااستفاده ($UR(r)$) و ($N_{LP}(r)$)، مقدار قابل ملاحظه‌ای است. این درحالی است که بدلیل عدم وجود فضای خالی کافی در چنین واحدهایی، اصولاً امکان قراردادن هسته‌های تروجان‌های سخت‌افزاری در آنها وجود ندارد و در نتیجه این واحدها آسیب‌پذیر نیستند.

در رهیافت ارائه شده در [۹] برای محاسبه میزان خطرپذیری تراشه‌های دیجیتال در برابر تروجان‌های سخت‌افزاری، نقشه خطرپذیری تروجان^۱ (TVM) برای سطح مشبک چیدمان فیزیکی^۲ استخراج شده است. TVM در واقع یک تصویر خاکستری با تعداد پیکسل‌های برابر با تعداد مربع‌های واحد موجود در چیدمان فیزیکی مشبک است. میزان خاکستری بودن هر پیکسل نمادی از

۳- روش پیشنهادی برای محاسبه میزان خطرپذیری تراشه‌های دیجیتال در برابر تروجان‌های سخت‌افزاری

با توجه به وجود عوامل متعدد، تعیین میزان آسیب‌پذیری یک مدار مجتمع در برابر تروجان‌های سخت‌افزاری امری بسیار پیچیده است. شاخص‌ترین روش‌های موجود برای تعیین میزان خطرپذیری، در مراجع [۹-۱۱] ارائه شده‌اند. در مرجع [۱۰]، ابتدا طبق رابطه ۴، مقدار خام خطرپذیری برای هر واحد مربعی شکل از سطح تراشه ($V(r)$) با ضرب کردن مقدار نرمالیزه شده فضای خالی ($WS(r)$) و تعداد منابع مسیره‌ی بلااستفاده ($UR(r)$)، محاسبه می‌شود. سپس آسیب‌پذیری واحد مربع شکل r براساس اینکه تروجان سخت‌افزاری نسبت به تأخیر مقاوم باشد ($V_{Td}(r)$) یا نسبت به توان دینامیک مقاوم باشد ($V_{LP}(r)$) و یا هر دو عامل قرار گرفته ($V_{Td\&LP}(r)$) با استفاده از روابط ۵-۷ محاسبه می‌شود.

$$V(r) = WS(r) \times UR(r) \quad (4)$$

$$V_{Td}(r) = V(r) \times N_{NC}(r) \quad (5)$$

$$V_{LP}(r) = V(r) \times N_{LP}(r) \quad (6)$$

$$V_{Td\&LP}(r) = V(r) \times N_{NC\&LP}(r) \quad (7)$$

در روابط (۵) تا (۷)، $N_{NC}(r)$ تعداد مسیره‌های غیربحرانی از لحاظ تأخیر در واحد مربعی r $N_{LP}(r)$ تعداد گره‌هایی است که مقدار نرخ سوئیچینگ آنها از یک مقدار آستانه کمتر است و $N_{NC\&LP}(r)$ تعداد گره‌هایی هستند که مقدار فعالیت سوئیچینگ آنها کمتر از مقدار آستانه می‌باشد و در عین حال بر روی مسیره‌های غیربحرانی قرار گرفته‌اند.

این روش دارای اشکالاتی اساسی در تعیین درست میزان آسیب‌پذیری مدار مجتمع می‌باشد. یکی از مهمترین اشکالات در این روش مربوط به واحدهایی است که یا کاملاً خالی هستند و یا سطح بسیار کمی از آنها اشغال شده است. این واحدها معیار فضای خالی را بخوبی ارضا می‌کنند و مکان مناسبی برای قراردادن هسته‌های سخت‌افزاری تروجان‌های سخت‌افزاری هستند. اما در روش [۱۰]، در صورتی که منبع مسیره‌ی بلااستفاده‌ای در آن واحد

² Layout

¹ Trojan Vulnerability Map

طراح تروجان به این گره‌ها را ایجاد می‌کنند، به شدت کاهش می‌یابند. معیار سوم به فاصله بین گره‌های دارای فعالیت سیگنال پایین با فضاهای خالی موجود در سطح چیدمان فیزیکی مربوط می‌شود. بدیهی است که هرچه این فاصله کمتر باشد، امکان ساختن هسته تروجان و اتصال آن به تریگر آسان‌تر خواهد بود. از جمله اشکالات این رهیافت عبارت است از عدم مدلسازی دقیق فضاهای خالی، مدلسازی ناقص تریگرها و عدم شناسایی دقیق مسیرهای ممکن برای اتصال تریگر به هسته تروجان.

۳-۱ تولید مجموعه داده مناسب از مدارهای دیجیتال

برای تولید مجموعه داده موردنظر این پژوهش، از مجموعه مدارهای معیار $ISCAS\ 85^2$ و $ISCAS\ 89$ استفاده شده است که شامل ۲۵ مدار دیجیتال با اندازه و ساختارهای مختلف می‌باشند. هر کدام از مدارهای موجود در این مجموعه از تعداد متنوعی گیت منطقی (AND, OR, DFF, ...) تشکیل شده است که در قالب یک گراف چند لایه به هم متصل شده‌اند. همچنین ارتباط اجزاء داخلی مدار از طریق ورودی‌های اصلی (PI) و خروجی‌های اصلی (PO) انجام می‌پذیرد. در هنگام پیاده‌سازی مدارهای دیجیتال بر روی سطح تراشه‌ها، ابتدا الگوی جانشانی مناسبی برای قراردادن گیت‌های منطقی و PI ها و PO ها تولید می‌شود. سپس در مرحله مسیره‌دهی ارتباط بین اجزاء داخلی با استفاده از منابع مسیره‌دهی موجود در تراشه (سیم‌ها، کنتاکت‌ها و Via ها و سوئیچ‌ها) تحقق پیدا می‌کند.

با توجه به تعداد محدود مدارهای دیجیتال موجود در مجموعه‌های $ISCAS\ 85$ و $ISCAS\ 89$ ، برای هر کدام از مدارها، چهارصد پیاده‌سازی مختلف تولید شده است، که در نهایت منجر به تولید ۱۰۰۰۰ پیاده‌سازی با ویژگی‌های متفاوت شده است.

در فرآیند ایجاد مجموعه‌داده موردنیاز این پژوهش، برای تولید پیاده‌سازی‌های مختلف از یک مدار معیار، مشخصه‌های مختلفی به صورت اتفاقی انتخاب شده‌اند. ابتدا، مساحت تراشه دیجیتال برای هر پیاده‌سازی از مقدار کمینه مورد نیاز تا ده برابر این مقدار کمینه، تغییر می‌کند. این رویکرد منجر به تولید فضاهای خالی با اندازه‌ها و الگوهای متفاوت در پیاده‌سازی‌های یک مدار می‌گردد. سپس،

میزان خطرپذیری آن است. مقدار سطح خاکستری هر پیکسل بر اساس معیار $TISF^1$ و مطابق رابطه (۸) محاسبه می‌شود.

در این رابطه BWS نمایانگر نسبت فضای خالی در واحد مورد نظر، TWS نسبت فضای خالی کل چیدمان فیزیکی، BFF و TFF به ترتیب بیانگر تعداد DFF در واحد مورد نظر و در کل چیدمان فیزیکی می‌باشند. WSD بیانگر توزیع فضاهای خالی است که با استفاده از رابطه (۹) محاسبه می‌شود. در این رابطه (x_i, y_i) و (x_m, y_m) به ترتیب بیانگر تعداد واحدهای فضای خالی موجود در واحد مربعی مورد بررسی، مختصات مرکز واحد مربعی و مختصات مرکز واحد فضای خالی می‌باشند.

$$TISF = \frac{BWS}{WSD} \times \frac{BFF}{TFF} \quad (8)$$

$$WSD = \sum_{i=1}^n \frac{1}{n} [(x_i - x_m)^2 + (y_i - y_m)^2] \quad (9)$$

اشکالات مربوط به مدلسازی فضاهای خالی روش ارائه شده در مرجع [۱۰]، در این رهیافت نیز وجود دارد. همچنین تمرکز این روش بر روی معیارهای فضای خالی و تعداد فلیپ‌فلاپ‌های داده (DFF)، منجر به عدم بررسی سایر معیارهای مهم از جمله منابع مسیره‌دهی، میزان فعالیت سیگنال و غیره شده است و بنابراین این رهیافت دقت متوسطی در تعیین میزان خطرپذیری چیدمان فیزیکی دارد.

در مرجع [۱۱] رهیافت دیگری برای تعیین میزان خطرپذیری در سطح چیدمان فیزیکی ارائه شده است که از سه معیار مختلف استفاده می‌کند. معیار اول بر مبنای استخراج هیستوگرام فضاهای خالی چیدمان فیزیکی می‌باشد که در آن ارتباط بین واحدهای فضای خالی با استفاده از روش ۴ همسایگی شناسایی می‌شود. معیار دوم مربوط به توانایی بلوکه کردن دسترسی به گره‌های با مقدار فعالیت پایین است. همان‌طور که پیشتر اشاره شد، این گره‌ها می‌توانند به‌عنوان تریگر برای هسته تروجان مورد استفاده قرارگیرند. نویسندگان معتقدند که با افزایش تراکم مسیرها در اطراف این گره‌ها، منابع مسیره‌دهی بلااستفاده که امکان دسترسی

ترکیب شده و با یک عدد ۸ بیتی به مؤلفه سبز هر پیکسل از تصویر رنگی نگاشته می‌شود. برای این کار، ابتدا مقدار CO ، $CC1$ ، $CC0$ و فعالیت سیگنال (act) هر کدام از گره‌های مدار محاسبه و نرمالیزه می‌شوند. برای نرمالیزه کردن هر مشخصه مقدار آن در واحد مربع موردنظر را بر ماکزیمم مقدار آن مشخصه در کل چیدمان فیزیکی تقسیم می‌شود تا یک عدد در بازه $[0,1]$ حاصل شود. این مقادیر نرمالیزه شده با استفاده از رابطه (۱۰) با هم ترکیب می‌شوند و حاصل عملیات با استفاده از تبدیل اعداد ممیز ثابت^۱ به یک عدد ۸ بیتی تبدیل می‌گردد. این مقدار نهایی به عنوان مؤلفه سبز رنگ پیکسل معادل از تصویر RGB ذخیره می‌شود.

$$RR_{norm} = \sqrt{(CC0 + CC1 + CO)^2 + act^2} \quad (10)$$

در رابطه (۱۰)، CO ، $CC1$ ، $CC0$ و act به ترتیب بیانگر میزان کنترل‌پذیری ترکیبی ۰، کنترل‌پذیری ترکیبی ۱، مشاهده‌پذیری ترکیبی و مقدار فعالیت سوئیچینگ گره‌های مدار می‌باشند. با توجه به تعریف مقدار فضای خالی در هر مربع واحد که برابر نسبت مساحت خالی به کل مساحت مربع واحد می‌باشد، عدد مربوطه همواره در بازه ۰ تا ۱ بوده و بنابراین نیاز به نرمالیزه کردن، ندارد. نسبت به‌دست‌آمده با استفاده از تبدیل عدد اعشار به عدد ممیز ثابت، به یک عدد ۸ بیتی تبدیل می‌شود و به مؤلفه قرمز رنگ در پیکسل معادل نسبت داده می‌شود. به عنوان مثال اگر $3/4$ مساحت یک مربع از سطح تراشه فضای خالی و مابقی با چیدمان فیزیکی گیت‌های منطقی پر شده باشد، مقدار عدد معادل ممیز ثابت برابر ۱۹۲ می‌باشد که به صورت عدد باینری 11000000 ذخیره خواهد شد.

مؤلفه آبی پیکسل‌ها به مقدار ویژگی تراکم منابع مسیردهی و مسیرهای غیربحرانی در هر مربع واحد، مربوط می‌شود. برای محاسبه عدد نرمالیزه معادل، ابتدا تعداد سیم‌هایی که برای مسیردهی سیگنال‌های مدار (از مسیرهای غیربحرانی) از مربع موردنظر عبور می‌کند، محاسبه می‌شود. سپس این عدد بر تعداد کل سیم‌هایی که قابلیت عبور از سطح مربع واحد موردنظر را دارند، تقسیم می‌گردد. عدد نرمالیزه حاصل با استفاده از روش تبدیل اعداد ممیز ثابت به

برای هر پیاده‌سازی، جانشانی منحصر به فردی در نظر گرفته می‌شود. به عبارت دیگر در ۴۰۰ پیاده‌سازی مختلف یک مدار معیار، ۴۰۰ جانشانی متفاوت تولید می‌شود. به دلیل متفاوت بودن مقدار فعالیت سیگنال، کنترل‌پذیری و مشاهده‌پذیری گیت‌های مختلف، این رویکرد موجب تولید ویژگی‌های مختلف برای یک موقعیت هندسی خاص از سطح تراشه خواهد شد. در نهایت، در هر پیاده‌سازی از یک مدار معیار، انتخاب نوع خاصی از طراحی برای پیاده‌سازی گیت‌های منطقی از میان طرح‌های مختلف (طراحی مکملی، نسبی، ترانزیستور گذر و دروازه انتقالی) صورت می‌گیرد. با توجه به متفاوت بودن اندازه و چیدمان فیزیکی گیت‌های منطقی در این روش‌های طراحی، توزیع و الگوی فضاهای خالی و میزان تراکم در فرایند مسیردهی در مکان‌های مختلف در سطح تراشه، از تنوع زیادی برخوردار خواهد شد. مرحله دوم از فرایند تولید مجموعه داده، شامل تبدیل هر کدام از پیاده‌سازی‌های تولیدشده برای مدارهای معیار، به یک تصویر رنگی (RGB) جهت استفاده در مراحل یادگیری و آزمون شبکه‌های عصبی عمیق می‌باشد. این تصاویر دارای ابعاد 128×128 پیکسل هستند. مرحله اول در این تبدیل عبارت است از مشبک کردن سطح تراشه به واحدهای مربع شکل به‌گونه‌ای که هر کدام از این واحدها به یک پیکسل از تصویر نگاشته می‌شود. در تصویر رنگی، هر کدام از پیکسل‌ها با سه مؤلفه قرمز، سبز و آبی مشخص می‌شود و با یک عدد N بیتی نمایش داده می‌شود. به عنوان مثال اگر $N=8$ باشد، آنگاه هر کدام از مؤلفه‌ها با یک عدد صحیح بین ۰ تا ۲۵۵ نمایش داده می‌شوند.

همان‌گونه که پیشتر ذکر شد، پنج ویژگی اصلی در تعیین خطرپذیری مدارهای دیجیتال در برابر تروجان‌های سخت‌افزاری موثرند (فضاهای خالی، تراکم منابع مسیردهی و مسیرهای غیربحرانی، فعالیت سیگنالی گره‌های مدار، میزان مشاهده‌پذیری و مقدار کنترل‌پذیری گره‌های مدار). این ویژگی‌ها باید به صورت مناسبی به سه مؤلفه قرمز، سبز و آبی در تصاویر رنگی نگاشته شوند. باتوجه به ارتباط تنگاتنگ سه ویژگی فعالیت سیگنالی، مشاهده‌پذیری و کنترل‌پذیری گره‌های مدار باهم، این سه ویژگی

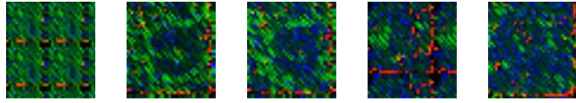
³ Switching Activity

¹ Fixed point

² Combinational Observability

پیاده‌سازی شده در مجموعه داده با یک برچسب که تعیین‌کننده کلاس خطرپذیری آن است، مشخص می‌شود.

در شکل ۴ چند نمونه از تصاویر مربوط به کلاس‌های مختلف از مجموعه داده تولیدشده نشان داده شده است.



شکل (۴): چند تصویر از مجموعه داده تولیدشده

۲-۳ شبکه‌های عصبی پیچشی

شبکه‌های عصبی کانولوشن، جزء پرکاربردترین راهکارهای یادگیری عمیق در زمینه‌های مختلف می‌باشند [۱۳]. از لحاظ معماری، این شبکه‌ها یک نوع شبکه عصبی Feed-forward می‌باشند که داده‌ها را معمولاً به صورت ماتریس‌های دوبعدی دریافت می‌کنند. بخش‌های اصلی شبکه عصبی کانولوشن شامل لایه ورودی، لایه پیچشی، لایه ادغام^۱، لایه سطح‌کننده، لایه کاملاً متصل و لایه خروجی می‌باشد. یک نمونه از شبکه‌های عصبی کانولوشن در شکل ۵ نشان داده شده است.

وظیفه اصلی لایه پیچشی استخراج ویژگی‌های مختلف از ورودی شبکه است. بعضی از لایه‌های پیچشی (لایه‌های نزدیک به لایه ورودی) وظیفه استخراج ویژگی‌های سطح پایین مانند لبه‌ها، خط‌ها و گوشه‌ها را برعهده دارند. لایه‌های بعدی برای استخراج ویژگی‌های پیچیده‌تر، با ترکیب ویژگی‌های سطح پایین استخراج‌شده در لایه‌های ابتدایی، به‌کار می‌روند. با وجود توانایی استخراج ویژگی‌ها توسط لایه‌های پیچشی، برای انجام تصمیم‌گیری صحیح در مورد کلاس تصویر ورودی توسط لایه کاملاً متصل، لازم است که تعداد ویژگی‌های استخراج‌شده به صورت هوشمندانه‌ای کاهش یابند. این امر با استفاده از لایه ادغام محقق می‌شود. لایه ادغام با نمونه‌برداری از ویژگی‌های استخراج‌شده به وسیله لایه پیچشی، به کاهش تعداد پارامترهای مؤثر در شبکه کمک می‌کند. پرکاربردترین نوع نمونه‌برداری در لایه ادغام استفاده از عملیات ماکزیمم (max pooling) و میانگین (average pooling) می‌باشد. در روش ادغام ماکزیمم، مقدار بیشینه از پیکسل‌هایی که توسط پنجره ادغام پوشش داده می‌شوند انتخاب می‌گردد و در روش میانگین مقدار متوسط پیکسل‌ها به عنوان خروجی در نظر گرفته می‌شود. در هر دو روش یک پیکسل به نمایندگی از کل پیکسل‌های پنجره به خروجی لایه ادغام می‌رود که این امر منجر به کاهش تعداد ویژگی‌ها برای لایه بعدی می‌گردد.

یک عدد باینری ۸ بیتی تبدیل می‌شود و به عنوان مولفه آبی پیکسل موردنظر در تصویر RGB ذخیره می‌گردد. چنین محاسبه‌ای از آن جهت منطقی می‌باشد که فرد طراح تروجان سخت‌افزاری می‌تواند از ظرفیت استفاده نشده برای سیم‌کشی در سطح مربع واحد برای ایجاد مسیرهای موردنیاز جهت اتصال تریگر به هسته تروجان سخت‌افزاری، استفاده کند.

بستر نرم‌افزاری توسعه یافته در این پژوهش قادر است که مدارهای دیجیتال را از نظر میزان خطرپذیری به ده دسته با خطرپذیری‌های مختلف، کلاس‌بندی کند. چون آموزش شبکه عصبی عمیق نیازمند داده‌های برچسب دار است، بنابراین در مجموعه داده تولیدشده باید هر کدام از مدارهای پیاده‌سازی شده با برچسب کلاس‌های ده گانه برچسب‌گذاری شوند. برای تعیین این برچسب‌ها، انواع مختلفی تروجان سخت‌افزاری تولیدشده و امکان قراردادن آن‌ها در مدار پیاده‌سازی شده، بررسی شده است.

جهت محاسبه میزان سختی قراردادن هر تروجان در مدار پیاده‌سازی شده، مراحل ذیل طی می‌شود. مرحله اول شامل یافتن فضای خالی مناسب جهت قراردادن هسته تروجان موردنظر است. در صورتی که چنین فضایی وجود نداشته باشد، فرآیند قراردادن تروجان شکست می‌خورد و میزان خطرپذیری مدار دیجیتال کاهش می‌یابد. در صورت یافتن جای مناسب، تعداد امکانات لازم برای تریگر کردن آن تروجان و وجود منابع مسیره‌دهی لازم برای اتصال آن تریگر به هسته تروجان بررسی می‌شود. برای هر تریگر، امتیازی بر اساس میزان آشکارسازی و سختی استفاده از آن تریگر محاسبه می‌شود و در نهایت برای تمام تروجان‌های قابل پیاده‌سازی، امتیازات جمع شده و به عنوان امتیاز کلی آن مدار ثبت می‌گردد.

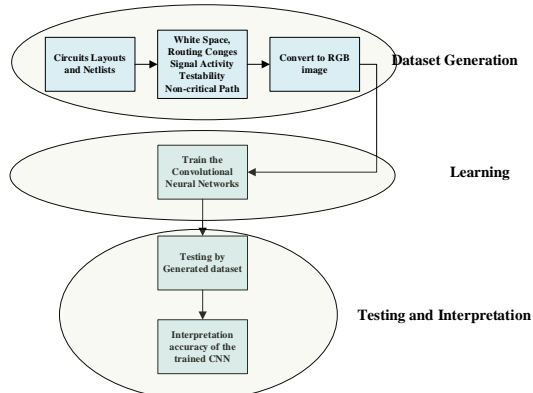
پس از محاسبه امتیاز تمام ۱۰۰۰۰ مدار پیاده‌سازی شده، تصاویر براساس امتیازها به صورت نزولی مرتب‌سازی می‌شوند. بالاترین امتیاز مربوط به مدارهای با خطرپذیری بسیار زیاد و پایین‌ترین امتیاز مربوط به مدارها با خطرپذیری کم است. این بازه به ۱۰ کلاس مختلف طبقه‌بندی می‌شوند که در آن شماره کلاس بیشتر، نمایانگر خطرپذیری بیشتر است. در نهایت هرکدام از مدارهای

¹ Pooling

۳-۳ روند پیشنهادی کلاس‌بندی مدارهای دیجیتال با

استفاده از شبکه عصبی کانولوشن

در شکل (۶) روند پیشنهادی این مقاله نشان داده شده است. در مرحله اول، داده‌های مورد نیاز (مجموعه داده) برای اجرای عملیات یادگیری شبکه عصبی پیچشی و ارزیابی صحت نتایج تولید می‌شود. توضیحات مربوط به این مرحله در بخش (۳-۱) ذکر گردیده است.



شکل (۶): جزئیات گام‌های مختلف در روند پیشنهادی

در مرحله دوم، معماری شبکه عصبی پیچشی مورد نظر انتخاب می‌شود و براساس مجموعه داده تولیدشده آموزش می‌بیند. در مرحله آموزش ۷۰٪ داده‌های کل به صورت تصادفی انتخاب می‌شوند که در میان این داده‌ها نیز، ۲۰٪ برای اعتبارسنجی کنارگذاشته می‌شوند.

در مرحله سوم، با استفاده از ۳۰٪ داده انتخاب شده به عنوان داده آزمون، میزان دقت شبکه آموزش یافته با استفاده از معیارهای مختلف سنجیده و تحلیل می‌شود.

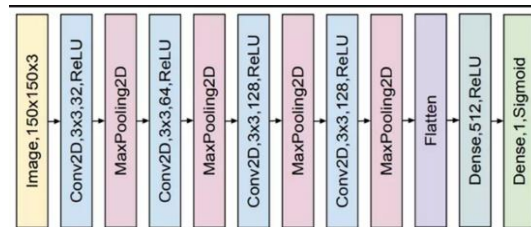
مساله بسیار مهم در فرآیند سه مرحله‌ای پیشنهادی، نحوه تعیین مقدار بهینه ابرمشخصه‌های^۳ می‌باشد. یکی از راهکارهای پرکاربرد در این زمینه، استفاده از روش جستجوی حریصانه است. در این روش ابتدا یک لیست از ابرمشخصه‌های مورد نیاز شبکه (نرخ یادگیری، اندازه دسته، تعداد epoch و ...) تهیه می‌شود و گستره مقدار هرکدام از آن‌ها، به ۲۵ مقدار مختلف تقسیم‌بندی می‌شود. در ابتدا مقادیر اولیه هرکدام از ابرمشخصه‌ها تعیین می‌شوند، سپس مرحله آموزش شبکه با استفاده از این ابرمشخصه‌ها تکمیل می‌شود و عملکرد شبکه ارزیابی می‌شود. در مرحله بعد، مقدار جدیدی از ابرمشخصه‌ها انتخاب می‌شود و شبکه بازآموزی می‌شود و اگر شبکه با ابرمشخصه‌های جدید عملکرد بهتری نسبت به بهترین حالت قبلی آن‌ها داشته باشد، ابرمشخصه‌های جدید به عنوان حالت بهینه انتخاب

از آنجاییکه لایه تماماً متصل، یک بردار از داده‌ها را به صورت ورودی می‌پذیرد، لذا خروجی آخرین لایه پیچشی که به صورت ماتریس‌های دوبعدی می‌باشد، به وسیله لایه مسطح کننده به یک بردار تبدیل می‌شود. لایه تماماً متصل از شبکه عصبی پیچشی، ساختاری شبیه شبکه عصبی^۱ MLP دارد که در آن خروجی نرون‌های موجود در هر لایه به ورودی تمام نرون‌های لایه بعد متصل می‌شوند. وظیفه اصلی این بخش از شبکه عصبی پیچشی، تشخیص کلاس بردار خروجی لایه مسطح کننده می‌باشد. لازم به ذکر است که لایه آخر در ساختار بخش تماماً متصل، از تابع softmax برای تشخیص کلاس داده ورودی مورد نظر استفاده می‌کند (رابطه ۱۱).

$$\text{softmax}(x_i) = \frac{e^{-\beta x_i}}{\sum_{j=1}^K e^{-\beta x_j}} \quad (11)$$

در رابطه (۱۱)، فرض بر این است که بردار ورودی تابع softmax شامل K عنصر می‌باشد (vector=(x1, x2,..., xK) و مشخصه (با توجه به شرایط مساله تعیین می‌شود).

آموزش شبکه‌های پیچشی که شامل تنظیم پارامترهای مختلف شبکه (وزن‌ها، مقادیر بایاس و ...) است به صورت آموزش با نظارت است. در این نوع یادگیری از یک مجموعه داده که در آن بروی هرکدام از داده‌ها برچسب مناسب قرارداده شده است، استفاده می‌شود. در ابتدای مرحله یادگیری شبکه، مقادیر پارامترها به صورت تصادفی انتخاب می‌شوند. سپس داده‌های موجود در مجموعه داده به صورت تصادفی به شبکه اعمال می‌شوند. برای هرکدام از داده‌ها، مقدار خروجی شبکه استخراج می‌شود و سپس با توجه به مقدار خطای تولیدشده (بر اساس یک تابع هزینه)، در مرحله پس انتشار^۲، مقدار پارامترهای قابل تنظیم شبکه اصلاح می‌گردند. پس از تکرار چندین باره الگوریتم یادگیری، شبکه عصبی پیچشی در بهترین وضعیت خود قرار می‌گیرد. در نهایت با استفاده از داده‌های آزمون مقدار دقت شبکه عصبی سنجیده می‌شود.



شکل (۵): نمونه‌ای از معماری شبکه عصبی کانولوشن

³ Hyper-parameters

¹ Multilayer Perceptron

² Back propagation

FN^۷: مواردی از مجموعه داده که برای آن‌ها، شبکه عصبی به اشتباه تشخیص می‌دهد تصویر مورد نظر متعلق به کلاس تحت بررسی نیست.

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN}) \quad (۱۲)$$

$$\text{TPR} = \text{TP} / (\text{TP} + \text{FN}) \quad (۱۳)$$

$$\text{TNR} = \text{TN} / (\text{TN} + \text{FP}) \quad (۱۴)$$

با توجه به تعاریف بالا و روابط (۱۲-۱۴) سه معیار زیر برای محاسبه کارایی شبکه‌های عصبی آموزش یافته مطابق روند پیشنهادی، مورد استفاده قرار می‌گیرند.

۱- صحت: طبق رابطه (۱۲)، نسبت تصاویری از مجموعه داده که به صورت صحیح توسط شبکه عصبی عمیق کلاس بندی می‌شود به کل تصاویر.

۲- TPR^۸: طبق رابطه (۱۳)، نسبت تصاویری که در یک کلاس قرار دارند و به درستی کلاس بندی شده‌اند به کل تصاویر موجود در آن کلاس.

۳- TNR^۹: طبق رابطه (۱۴)، نسبت تعداد تصاویری که شبکه عصبی عمیق به درستی تشخیص داده است که متعلق به کلاس تحت بررسی نیست به کل تصاویری که در آن کلاس قرار نمی‌گیرند.

در جدول (۲)، ساختار شبکه عصبی پیچشی و مقادیر ابرمشنه‌ها که پس از اعمال الگوریتم جستجوی حرصانه به دست آمده است، گزارش شده است.

به دلیل اینکه روش‌های پیشین بر مبنای یادگیری ماشین، تنها برای تشخیص وجود تروجان‌های سخت‌افزاری بعد از تولید تراشه توسعه پیدا کرده‌اند [۱۴-۱۷] و به ارزیابی میزان خطرپذیری تراشه‌ها در برابر تروجان‌های سخت‌افزاری در مرحله طراحی نپرداخته‌اند، استخراج نتایج برای طبقه بندی کننده‌های^{۱۱} معمولی بر روی مجموعه داده این پژوهش، به صورت مستقل انجام شده است. در جدول (۳) مقایسه میزان دقت به دست آمده توسط شبکه عصبی پیچشی با سایر روش‌های طبقه بندی کننده خطی مشهور گزارش شده است. لازم به ذکر است که نتایج مربوط به طبقه بندی کننده‌های معمولی (طبقه بندی کننده‌های مرکب^{۱۲}، ماشین

می‌شوند و در غیر این صورت مقدار بهینه قبلی حفظ می‌شود. هرچند این روش نسبت به دیگر الگوریتم‌ها (جستجوی اتفاقی و بهینه سازی Bayesian) بیشتر طول می‌کشد، لیکن ابرمشنه‌های انتخاب شده عملکرد بسیار خوبی دارند.

ذکر این نکته ضروری به نظر می‌رسد که در مدلسازی با استفاده از شبکه‌های عصبی پیچشی، باید مساله انطباق بیش از حد^۱ به طریق مناسبی حل شود. برای اینکار روش‌های مختلفی همچون استفاده از روش‌های انتظام^۲، توقف زودهنگام^۳، حذف کردن^۴ و غیره ارائه شده است. L2-Regulization و ۲۰٪ حذف کردن به همراه توقف زودهنگام راهکارهای جلوگیری از انطباق بیش از حد مورد استفاده در مرحله آموزش در شبکه عصبی پیچشی پیشنهادی این پژوهش می‌باشند.

۴- شبیه سازی‌ها و نتایج

در این بخش به بررسی دقت روش پیشنهادی در کلاس بندی خطرپذیری تراشه‌های دیجیتال در برابر تروجان‌های سخت‌افزاری پرداخته می‌شود. شبیه سازی‌ها در محیط MATLAB و Python (کتابخانه Tensorflow مربوط به یادگیری عمیق) بر روی سیستم کامپیوتری با پردازنده نسل ۶ از Core-i7 با حافظه RAM ۸ گیگابایت و فرکانس ۲/۵ گیگاهرتز و کارت گرافیکی GeForce 940MX انجام شده‌اند. در مرحله یادگیری شبکه‌های عصبی عمیق، از بهینه ساز GCD، نرخ یادگیری ۰/۰۰۲، تعداد epoch برابر ۱۰۰ و تابع هزینه categorical cross-entropy استفاده شده است. همان گونه که در بخش پیش ذکر شد، مقادیر بهینه ابرمشنه‌ها با استفاده از روش جستجوی حرصانه بر روی شبکه کامل هایپر پارامترها استخراج شده‌اند. این پژوهش جهت اندازه گیری میزان صحت روند پیشنهادی از مفاهیم ذیل بهره می‌برد:

TP^۵: مواردی از مجموعه داده که برای آن‌ها، کلاس پیش بینی شده توسط شبکه عصبی عمیق با کلاس واقعی تصویر یکسان است.

TN^۶: مواردی از مجموعه داده که برای آن‌ها، شبکه عصبی عمیق به درستی تشخیص می‌دهد تصویر مورد نظر متعلق به کلاس تحت بررسی نیست.

FP^۷: مواردی از مجموعه داده که برای آن‌ها، شبکه عصبی به اشتباه تشخیص می‌دهد تصویر مورد نظر متعلق به کلاس تحت بررسی است.

⁷ False Positive

⁸ False Negative

⁹ True Positive Rate

¹⁰ True Negative Rate

¹¹ Classifiers

¹² Ensemble Classifier

¹ Overfitting

² Regularization

³ Early Stopping

⁴ Dropout

⁵ True Positive

⁶ True Negative

جدول (۳) میزان دقت طبقه‌بندی در روش‌های مختلف

Method	Accuracy%
Ensemble	75.7
Tree	67.5
SVM	72.4
KNN	68.2
[9]	62.4
[10]	57.8
[11]	54.8
Proposed CNN	92.6

مقایسه معیار دقت نتایج شبکه عصبی کانولوشن و بهترین الگوریتم‌های کلاس‌بندی خطی مورداستفاده در شبیه‌سازی‌ها (طبقه‌بندی کننده مرکب) بر مبنای ماتریس Confusion در شکل (۷) نشان داده شده است. برای استخراج این ماتریس‌ها از یک مجموعه ۳۰۰۰ تایی از تصاویر مربوط به کلاس‌های ده‌گانه موجود در مجموعه داده که به صورت اتفاقی انتخاب شده‌اند، استفاده شده است. در ماتریس Confusion، هر سطر بیانگر کلاس تشخیص داده شده توسط طبقه‌بندی کننده و هر ستون بیانگر کلاسی است که داده مورد آزمایش واقعاً به آن تعلق دارد. برای یک طبقه‌بندی کننده ایده‌آل مقدار خانه‌هایی که روی قطر اصلی قرار دارند برابر ۱ است. مقادیر درصد‌های داخل هر خانه به جز سطر آخر و همچنین ستون سمت راست، از تقسیم عدد مربوط به آن خانه بر مجموع تمام اعداد داخل جدول (کل نمونه‌های استفاده شده) به دست می‌آید. اعداد درصد سبزرنگ مربوط به خانه‌های ستون سمت راست از تقسیم عدد موجود در مربع سبزرنگ سطر مربوطه بر تمام اعداد موجود در آن سطر به دست می‌آید و عدد درصد قرمز رنگ از تفریق درصد سبزرنگ از عدد ۱۰۰ به دست می‌آید. رویکرد مشابهی برای محاسبه درصد‌های سبز و قرمز رنگ در خانه‌های سطر آخر اتخاذ شده است. بنابراین، اعداد سبزرنگ ستون سمت راست بیانگر آن هستند که طبقه‌بندی کننده در مورد داده‌هایی که تشخیص داده شده است کلاس آن‌ها در ردیف خاصی قرار دارد، چه میزان صحیح عمل کرده است و همچنین درصد سبزرنگ هر ستون در سطر آخر بیان می‌کند که طبقه‌بندی کننده در مورد تشخیص کلاس داده‌های هر کدام از کلاس‌های ده‌گانه به درستی عمل کرده است. به طور مثال، در مورد داده‌هایی که متعلق به کلاس شماره ۴ هستند، شبکه عصبی پیشنهادی دارای دقت ۸۹/۱٪ و طبقه‌بندی مرکب دارای دقت ۵۷/۷٪ می‌باشد. از مقایسه جداول Confusion این نتیجه حاصل می‌شود که شبکه عصبی کانولوشن پیشنهادی، نسبت به دیگر الگوریتم‌ها از دقت بالاتری برای

بردار پشتیبان^۱، درخت^۲، k-همسایه نزدیک^۳ با آموزش آن‌ها در Toolbox مربوطه از نرم‌افزار MATLAB صورت گرفته است. برای آموزش این طبقه‌بندی کننده‌ها از مجموعه داده تولید شده در این پژوهش استفاده شده است و مقدار ابرمشخصه‌ها با استفاده از روش جستجوی حریصانه موجود در Toolbox استخراج شده است.

همان‌گونه که ملاحظه می‌شود، میزان دقت شبکه پیشنهادی (بالای ۹۲٪) نسبت به سایر روش‌ها بسیار بهتر است. یکی از دلایل این عملکرد بهتر را می‌توان در توانایی فوق‌العاده لایه‌های پیچشی شبکه عصبی پیشنهادی ما برای استخراج ویژگی‌های اساسی مؤثر در میزان خطرپذیری تراشه‌های دیجیتال دانست. این لایه‌ها قادرند تعداد زیادی ویژگی مختلف را در تصویر ورودی پردازش کنند و در نهایت آن‌ها را به گونه‌ای به لایه‌های تماماً متصل ارائه دهند که بیشترین تفاوت را در کلاس‌های مختلف داشته باشند. از سوی دیگر، الگوریتم فراگیر و پیچیده شبکه‌های عصبی پیچشی امکان تحلیل ویژگی‌های اصلی را در تصاویر مجموعه داده فراهم می‌کند، که این امکان با استفاده از الگوریتم‌های یادگیری ساده‌تر سایر طبقه‌بندی کننده‌ها وجود ندارد. همان‌گونه که پیشتر ذکر شد، ویژگی‌های تعیین کننده میزان خطرپذیری تراشه‌های دیجیتال، دارای رفتار غیرخطی در نواحی مختلف سطح تراشه می‌باشند و در نتیجه ترکیب اثرات این ویژگی‌ها و اتخاذ تصمیم در مورد میزان خطرپذیری کل تراشه فرآیندی بسیار پیچیده است که نیازمند استفاده از مدلی غیرخطی با توانایی یادگیری بالا دارد.

جدول (۲) مقادیر بهینه ابرمشخصه‌ها

مقدار بهینه	ابرمشخصه
5	تعداد لایه‌های پیچشی
5	تعداد لایه‌های ادغام
3	تعداد لایه‌های تماماً متصل
0.9	مومتوم
0.1	نرخ یادگیری اولیه
0.09	ضریب کاهش نرخ یادگیری
0.95	نرخ تضعیف میانگین متحرک
200	تعداد دورها در هر تکرار
0.0005	تضعیف وزن‌ها
128	اندازه دسته

³ K-nearest neighbors (KNN)

¹ Support-Vector Machine (SVM)

² Tree

Confusion Matrix

Output Class	0	1	2	3	4	5	6	7	8	9	Accuracy	Loss
0	1145 14.9%	25 0.3%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	3 0.0%	0 0.0%	5 0.1%	97.1% 2.9%	
1	52 0.7%	824 10.7%	31 0.4%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	90.6% 9.4%	
2	5 0.1%	52 0.7%	847 11.0%	40 0.5%	3 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	89.3% 10.7%	
3	1 0.0%	4 0.1%	44 0.6%	682 8.9%	30 0.4%	3 0.0%	1 0.0%	0 0.0%	0 0.0%	2 0.0%	88.9% 11.1%	
4	0 0.0%	0 0.0%	2 0.0%	28 0.4%	476 6.2%	16 0.2%	2 0.0%	0 0.0%	1 0.0%	2 0.0%	90.3% 9.7%	
5	2 0.0%	0 0.0%	0 0.0%	3 0.0%	21 0.3%	358 4.7%	16 0.2%	1 0.0%	3 0.0%	1 0.0%	88.2% 11.8%	
6	1 0.0%	0 0.0%	0 0.0%	2 0.0%	2 0.0%	30 0.4%	289 3.8%	16 0.2%	1 0.0%	9 0.1%	82.6% 17.4%	
7	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	2 0.0%	10 0.1%	216 2.8%	6 0.1%	11 0.1%	87.8% 12.2%	
8	3 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	6 0.1%	171 2.2%	41 0.5%	77.0% 23.0%	
9	15 0.2%	0 0.0%	0 0.0%	0 0.0%	2 0.0%	0 0.0%	1 0.0%	3 0.0%	9 0.1%	2106 27.4%	98.8% 1.4%	
	93.5% 6.5%	91.0% 9.0%	91.5% 8.5%	90.1% 9.9%	89.1% 10.9%	87.3% 12.7%	90.3% 9.7%	88.2% 11.8%	89.5% 10.5%	96.6% 3.4%	92.5% 7.5%	

الف

Confusion Matrix

Output Class	0	1	2	3	4	5	6	7	8	9	Accuracy	Loss
0	1129 14.7%	53 0.7%	4 0.1%	3 0.0%	1 0.0%	1 0.0%	0 0.0%	3 0.0%	0 0.0%	9 0.1%	93.8% 6.2%	
1	60 0.8%	728 9.5%	104 1.4%	10 0.1%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	80.7% 19.3%	
2	7 0.1%	114 1.5%	661 8.6%	135 1.8%	16 0.2%	1 0.0%	1 0.0%	0 0.0%	0 0.0%	1 0.0%	70.6% 29.4%	
3	1 0.0%	9 0.1%	143 1.9%	482 6.3%	106 1.4%	13 0.2%	2 0.0%	0 0.0%	0 0.0%	1 0.0%	63.7% 36.3%	
4	3 0.0%	1 0.0%	10 0.1%	107 1.4%	308 4.0%	78 1.0%	15 0.2%	3 0.0%	1 0.0%	2 0.0%	58.3% 41.7%	
5	2 0.0%	0 0.0%	4 0.0%	15 0.2%	78 1.0%	230 3.0%	68 0.9%	12 0.2%	3 0.0%	4 0.0%	55.3% 44.7%	
6	2 0.0%	0 0.0%	0 0.0%	2 0.0%	17 0.2%	69 0.9%	168 2.2%	63 0.8%	14 0.2%	12 0.2%	48.4% 51.6%	
7	0 0.0%	0 0.0%	0 0.0%	1 0.0%	3 0.0%	11 0.1%	44 0.6%	103 1.3%	43 0.6%	17 0.2%	46.4% 53.6%	
8	3 0.0%	0 0.0%	0 0.0%	0 0.0%	2 0.0%	2 0.0%	14 0.2%	44 0.6%	81 1.1%	69 0.9%	37.7% 62.3%	
9	17 0.2%	0 0.0%	0 0.0%	0 0.0%	2 0.0%	3 0.0%	5 0.1%	8 0.1%	49 0.6%	2064 26.8%	95.3% 4.7%	
	92.2% 7.8%	80.4% 19.6%	71.4% 28.6%	63.7% 36.3%	57.7% 42.3%	56.1% 43.9%	52.5% 47.5%	42.0% 58.0%	42.4% 57.6%	94.7% 5.3%	77.4% 22.6%	

ب

شکل (۷): الف) ماتریس confusion شبکه عصبی پیچشی (ب)

ماتریس confusion طبقه‌بندی کننده مرکب

مقایسه نتایج مربوط به سه رهیافت مراجع [۹]، [۱۰] و [۱۱] نشان می‌دهد که عملکرد مرجع [۹] از دو روش دیگر بهتر است، که دلیل اصلی آن ارائه مدلی نسبتاً دقیق‌تر از فضاهای خالی چیدمان فیزیکی در این روش می‌باشد. لیکن عدم مدلسازی دقیق مسیره‌دهی تریگر به هسته تروجان و عدم شناخت دقیق فرآیند تریگر، دقت آن را تا ۶۲٪ کاهش داده است. از طرفی دیگر، در نظر گرفتن معیارهای آزمون‌پذیری (کنترل‌پذیری و مشاهده‌پذیری) در رهیافت مرجع [۱۰] سبب شده است که این روش دقت بالاتری نسبت به رهیافت مرجع [۱۱] داشته باشد. اشکالات دیگر

تشخیص کلاس مربوط به هر مدار و همچنین دقت بالاتری در تشخیص عدم تعلق یک مدار به یک کلاس خاص برخوردار است.

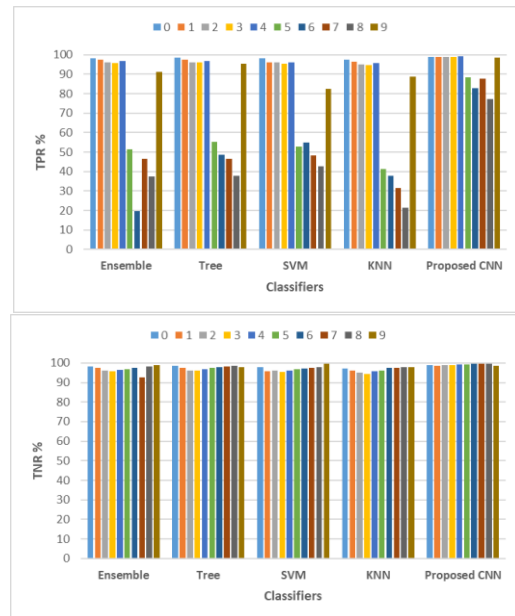
در شبیه‌سازی دیگری که نتایج آن در شکل (۸) نشان داده شده است، مقادیر TPR و TNR مربوط به شبکه عصبی کانولوشن پیشنهادی و دیگر الگوریتم‌های طبقه‌بندی کننده، گزارش شده است. همان‌گونه که پیشتر ذکر شد، TPR بیانگر میزان توانایی یک طبقه‌بندی کننده در تشخیص صحیح داده‌هایی است که به یک کلاس خاص تعلق دارند. این معیار در شبکه‌عصبی پیچشی پیشنهادی در همه کلاس‌ها نسبت به سایر طبقه‌بندی کننده‌ها برتری محسوس دارد و این برتری مخصوصاً در کلاس‌های بالاتر (که میزان خطرپذیری بیشتری نسبت به تروجان‌های سخت‌افزاری دارند) به‌خوبی قابل مشاهده است. به‌عبارت دیگر، در سایر طبقه‌بندی کننده‌ها به‌احتمال بسیار بیشتر ممکن است یک مدار دیجیتال که میزان خطرپذیری بالایی دارد به‌عنوان مداری تشخیص داده شود که دارای خطرپذیری کم و یا متوسط است و این امر منجر به خطرپذیری بیشتر در قراردادن تروجان‌های سخت‌افزاری در این تراشه‌ها می‌شود.

در نهایت، برای مقایسه با کارهای پیشین نتایج صحت شبکه عصبی پیچشی پیشنهادی و الگوریتم‌های کلاس‌بندی با سه روش مراجع [۹]، [۱۰] و [۱۱] در جدول (۳) مقایسه شده‌اند. روش‌های مذکور بر اساس یادگیری ماشین نیستند و نتایج را برای مجموعه مدارهای دیگر ارائه کرده‌اند. لذا جهت مقایسه صحیح با روش‌های مبتنی بر یادگیری ماشین، نتایج آن‌ها (با کدنویسی روش‌های مذکور) برای مجموعه داده این پژوهش، بازتولید شده‌است. همان‌گونه که در بخش (۳) ذکر گردید، این روش‌ها هرکدام تعدادی از مشخصه‌های اصلی دخیل در میزان خطرپذیری را در نظر نگرفته‌اند و یا مدل‌سازی آن‌ها از این مشخصه‌ها جامع و کامل نمی‌باشد. از نتایج جدول (۳) نیز می‌توان دریافت که اختلاف زیادی در میزان دقت این روش‌ها با روش‌های مبتنی بر الگوریتم‌های هوشمند وجود دارد. دلیل اصلی این امر را می‌توان در توانایی فوق‌العاده شبکه‌های عصبی عمیق پیچشی در استخراج ویژگی‌های مؤثر در تعیین میزان خطرپذیری تراشه‌های دیجیتال دانست. این ویژگی‌ها دارای روابط پیچیده با یکدیگر هستند که امکان مدلسازی جامع و دقیق آن‌ها هم با استفاده از روش‌های تحلیلی (closed form) و هم با استفاده از طبقه‌بندی کننده‌های خطی معمولی غیر ممکن است.

References

- [1] Bao, C., Forte, D., & Srivastava, A. (2015). On reverse engineering-based hardware Trojan detection. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 35(1), 49-57.
- [2] Bhunia, S., Abramovici, M., Agrawal, D., Bradley, P., Hsiao, M.S., Plusquellic, J. and Tehranipoor, M., 2013. Protection against hardware trojan attacks: Towards a comprehensive solution. *IEEE Design & Test*, 30(3), pp.6-17.
- [3] He, J., Guo, X., Ma, H., Liu, Y., Zhao, Y. and Jin, Y., 2020, July. Runtime trust evaluation and hardware trojan detection using on-chip em sensors. In *2020 57th ACM/IEEE Design Automation Conference (DAC)* (pp. 1-6). IEEE.
- [4] Cozzi, M., Galliere, J.M. and Maurine, P., 2018, August. Exploiting Phase Information in Thermal Scans for Stealthy Trojan Detection. In *2018 21st Euromicro Conference on Digital System Design (DSD)* (pp. 573-576). IEEE.
- [5] Koushanfar, F. and Mirhoseini, A., 2010. A unified framework for multimodal submodular integrated circuits trojan detection. *IEEE Transactions on Information Forensics and Security*, 6(1), pp.162-174.
- [6] Chen, K., Arias, O., Guo, X., Deng, Q., & Jin, Y. (2022). IP-Tag: Tag-Based Runtime 3PIP Hardware Trojan Detection in SoC Platforms. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 42(1), 68–81
- [7] Q. Shi, N. Vashistha, H. Lu, H. Shen, B. Tehranipoor, D. L. Woodard, and N. Asadizanjani, "Golden gates: A new hybrid approach for rapid hardware trojan detection using testing and imaging," in *2019 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*. IEEE, 2019, pp. 61–71.
- [8] Rahimifar, M. M., & Jahanirad, H. (2020). Employing Image Processing Techniques for Hardware Trojans Detection. In *2020 10th International Conference on Computer and Knowledge Engineering (ICCKE)* (pp. 187–192). IEEE.
- [9] Trippel, T., Shin, K. G., Bush, K. B., & Hicks, M. (2020). ICAS: An extensible framework for estimating the susceptibility of ic layouts to additive trojans. In *2020 IEEE Symposium on Security and Privacy (SP)* (pp. 1742–1759). IEEE.

روش مرجع [۱۱] عدم مدل‌سازی مناسب سیگنال‌های تریگر و همچنین غیرقابل اعمال آن برای مدارهای ترکیبی می‌باشد.



شکل (۸): مقایسه TPR و TNR برای طبقه‌بندی‌کننده‌های مختلف

۵- نتیجه‌گیری

در این مقاله، به مدل‌سازی میزان خطرپذیری تراشه‌های دیجیتال در برابر حملات امنیتی ناشی از تروجان‌های سخت‌افزاری پرداخته شده است. در ابتدا مجموعه داده مناسبی از پیاده‌سازی‌های متنوع تراشه‌های دیجیتال ارائه شده و بر اساس روش قرارداد تروجان‌های سخت‌افزاری میزان خطرپذیری آن‌ها کلاس‌بندی شده‌اند. سپس مجموعه داده تولیدشده جهت آموزش شبکه عصبی کانولوشن مورد استفاده قرار گرفته و میزان توانایی شبکه پیشنهادی در استخراج ویژگی‌های مختلف، منجر به کلاس‌بندی با دقت ۹۲٪ شده است. مقایسه نتایج دقت کلاس‌بندی با سایر طبقه‌بندی‌کننده‌ها و همچنین روش‌های پیشین میزان بهبود ۱۷٪ را نشان می‌دهد. از مدل بدست آمده می‌توان در مراحل اولیه طراحی تراشه‌های دیجیتال برای استخراج میزان خطرپذیری مدارها و امکان به کارگیری رهیافت‌هایی برای ارتقاء امنیت تراشه‌ها بهره برد. با توجه به دقت بالای شبکه عصبی کانولوشن پیشنهادی در عین کم بودن سربار محاسباتی و حافظه، این امکان وجود دارد که از رهیافت پیشنهادی برای ارزیابی میزان خطرپذیری تراشه‌های دیجیتال در مرحله طراحی، در نرم‌افزارهای مربوط به طراحی با کمک کامپیوتر استفاده شود.



- on Computer-Aided Design of Integrated Circuits and Systems 37, no. 7 (2017): 1370-1383.
- [15] Kok, Chee Hoo, Chia Yee Ooi, Michiko Inoue, Mehrdad Moghbel, Sreedharan Baskara Dass, Hau Sim Choo, Nordinah Ismail, and Fawnizu Azmadi Hussin. "Net classification based on testability and netlist structural features for hardware trojan detection." In 2019 IEEE 28th Asian Test Symposium (ATS), pp. 105-1055. IEEE, 2019.
- [16] Hasegawa, Kento, Masao Yanagisawa, and Nozomu Togawa. "Trojan-feature extraction at gate-level netlists and its application to hardware-Trojan detection using random forest classifier." In 2017 IEEE International Symposium on Circuits and Systems (ISCAS), pp. 1-4. IEEE, 2017.
- [17] [17] Hasegawa, Kento, Youhua Shi, and Nozomu Togawa. "Hardware trojan detection utilizing machine learning approaches." In 2018 17th IEEE International Conference On Trust, Security And Privacy In Computing And Communications/12th IEEE International Conference On Big Data Science And Engineering (TrustCom/BigDataSE), pp. 1891-1896. IEEE, 2018.
- [10] Salmani, H., & Tehranipoor, M. M. (2016). Vulnerability analysis of a circuit layout to hardware Trojan insertion. *IEEE Transactions on Information Forensics and Security*, 11 (6), 1214–1225.
- [11] Bakhshizadeh, M., & Jahanian, A. (2014). Trojan Vulnerability Map: An efficient metric for modeling and improving the security level of hardware. *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, 97(11), 2218–2226.
- [12] Zamiri Azar, Kimia, Hadi Mardani Kamali, Farimah Farahmandi, and Mark Tehranipoor. "Basics of VLSI Testing and Debug." In *Understanding Logic Locking*, pp. 25-46. Cham: Springer International Publishing, 2023.
- [13] Cong, Shuang, and Yang Zhou. "A review of convolutional neural network architectures and their optimizations." *Artificial Intelligence Review* 56, no. 3 (2023): 1905-1969.
- [14] Chen, Xiaoming, Qiaoyi Liu, Song Yao, Jia Wang, Qiang Xu, Yu Wang, Yongpan Liu, and Huazhong Yang. "Hardware trojan detection in third-party digital intellectual property cores by multilevel feature analysis." *IEEE Transactions*

A convolutional Neural Networks-based Approach for Vulnerability Classification of Integrated Circuits against Hardware Trojans

Hadi Jahanirad*

Department of Electronics and Communication Engineering, University of Kurdistan, Sanandaj, Iran

Article Information

Original Research Paper

Received:

2024 March 7

Accepted:

2024 August 27

Keywords:

Digital integrated circuits,
Hardware Trojans,
Convolutional neural networks,
Classification, Machine
learning.

Corresponding Author*:

h.jahanirad@uok.ac.ir

Abstract

The vulnerability of digital integrated circuits against the Hardware Trojans (HT) has increased in recent decades due to the implementation of more complex systems on them. HTs could become a source of errors or apply to steal important information embedded in the implemented circuits. So, analyzing the vulnerability of digital integrated circuits in the early stages of production is of great merit. In this paper, a novel vulnerability classification method is introduced based on the deep convolutional neural networks (CNN) wherein five major effective features of vulnerability assessment are utilized (white space distribution, unutilized routing resources, signal activity of circuit nodes, delay of the circuit paths and, controllability of circuit nodes). In the proposed framework, first of all, a dataset containing 10000 images is generated using various digital circuit implementations. Then, a deep CNN is trained using the generated dataset meanwhile the most appropriate CNN's hyperparameters are achieved using a greedy optimization method. The simulation results reveal 92% accuracy of vulnerability classification which shows a 17% improvement in comparison with the best linear classifier and analytical methods.

 : 10.22034/ABMIR.2024.21351.1049

E-ISSN: [2821-2037](https://doi.org/10.22034/ABMIR.2024.21351.1049) /© 2023. Published by Yazd University This is an open access article under the CC BY 4.0 License (<https://creativecommons.org/licenses/by/4.0/>).

