

به نام خدا



دانشگاه تهران

پردیس دانشکده‌های فنی

دانشکده برق و کامپیوتر



گزارش کار و مستندات پروژه کارشناسی

پیاده‌سازی سخت‌افزاری الگوریتم شبکه عصبی برای تشخیص تصویر
مجموعه داده MNIST بر روی FPGA

استاد داور

دکتر بیژن علیزاده

استاد راهنما

دکتر زین‌العابدین نوابی

علی شایان پور

۸۱۰۱۹۸۵۳۲

بهار ۱۴۰۲

۱- چکیده.....	4
۲- مقدمه.....	5
۳- شرح پروژه.....	6
۳.۱- شبکه عصبی MLP برای تشخیص اعداد دست‌نوشته.....	6
۳.۱.۱- مقدمه‌ای بر شبکه MLP.....	6
۳.۱.۲- پیاده‌سازی.....	7
۳.۲- شبکه عصبی CNN برای تشخیص اعداد دست‌نوشته.....	8
۳.۲.۱- مقدمه‌ای بر شبکه عصبی CNN.....	8
۳.۲.۲- پیاده‌سازی.....	9
۳.۳- استخراج ضرایب و پیاده‌سازی شبکه CNN بدون استفاده از کتابخانه.....	10
۳.۳.۱- لایه Convolution.....	10
۳.۳.۲- لایه Max Pooling.....	11
۳.۳.۳- لایه Dense یا Fully Connected.....	11
۳.۴- اضافه کردن ماژول دوربین.....	11
۳.۵- پردازش به صورت Real Time.....	12
۳.۶- اجرا کردن شبکه CNN بر روی بخش PS برد PYNQ.....	12
۳.۶.۱- برد PYNQ-Z2.....	12
۳.۶.۲- پیاده‌سازی.....	13
۳.۷- ساخت سیستم‌های ساده و ارتباط PS و PL برد PYNQ.....	13
۳.۷.۱- ساختار PS و PL برد PYNQ.....	14
۳.۷.۲- پیاده‌سازی.....	15
۳.۸- پروتکل AXI-4 Stream.....	17
۳.۹- طراحی ماژول General Purpose AXI-4 Stream Interface.....	19
۳.۱۰- ماژول Dense و ترکیب آن با GP AXIS IF.....	19
۴- نتایج.....	21
۴.۱- آموزش شبکه MLP برای تشخیص اعداد دست‌نوشته MNIST.....	21
۴.۲- درستی‌سنجی مدل train شده MLP بر روی داده‌های واقعی دست‌نوشته.....	22
۴.۲.۱- اعداد دست‌نوشته با رنگ سیاه.....	22
۴.۲.۲- اعداد دست‌نوشته با رنگ آبی، ضخامت نازک.....	23
۴.۲.۳- اعداد دست‌نوشته با رنگ آبی، ضخامت کلفت.....	23

23.....	۴.۲.۴- نتیجه‌گیری
24.....	۴.۳- آموزش شبکه CNN برای تشخیص اعداد دست‌نوشته MNIST
24.....	۴.۴- درستی‌سنجی مدل train شده CNN بر روی داده‌های واقعی دست‌نوشته
24.....	۴.۴.۱- اعداد دست‌نوشته با رنگ سیاه
25.....	۴.۴.۲- اعداد دست‌نوشته با رنگ آبی، ضخامت نازک
25.....	۴.۴.۳- اعداد دست‌نوشته با رنگ آبی، ضخامت کلفت
25.....	۴.۴.۴- نتیجه‌گیری
26.....	۴.۵- اضافه کردن ماژول دوربین
27.....	۴.۶- شبیه‌سازی ماژول GP AXIS Interface
27.....	۴.۷- شبیه‌سازی ماژول Dense یا Fully Connected
28.....	۴.۸- شبیه‌سازی ماژول ReLU
28.....	۴.۹- زمان‌بندی‌های نرم‌افزاری بر روی پردازنده PC
29.....	۴.۱۰- زمان‌بندی‌های نرم‌افزاری بر روی قسمت PS برد PYNQ
30.....	۴.۱۱- زمان‌بندی‌های سخت‌افزاری
31.....	۴.۱۲- نتیجه زمان‌بندی‌ها و مقایسه قسمت PS و PL برد PYNQ
32.....	۴.۱۲- نتیجه زمان‌بندی‌ها و مقایسه قسمت PL برد PYNQ با پردازنده PC
33.....	۴.۱۳- زمان‌بندی اجرا Real Time بر روی پردازنده PC
35.....	۵- جمع بندی
36.....	۶- مراجع

پیشرفت روزافزون و فراگیر هوش مصنوعی بر کسی پوشیده نیست. از طرفی پیاده‌سازی سخت‌افزاری این الگوریتم‌ها گام مهمی در تحقق استفاده‌های هر چه بیشتر از این الگوریتم‌هاست. هدف از این پروژه پیاده‌سازی الگوریتم شبکه عصبی برای کاربرد پردازش تصویر به منظور تشخیص تصویر است. پروسه پردازش و اجرای الگوریتم‌های هوش مصنوعی با توجه به پیچیدگی‌ها و لایه‌های زیاد، باعث افت سرعت و زمان محاسبه طولانی می‌شود و هدف از پروژه دریافت و پردازش این اطلاعات با سرعت بسیار بالا می‌باشد.

امروزه با توسعه فراگیر و سریع هوش مصنوعی و پیچیدگی مدل‌ها، کتابخانه‌ها و روش‌های زیادی برای ماژولار کردن و ساده‌سازی عملیات‌های هوش مصنوعی به وجود آمده است. این ساده‌سازی‌ها در زبان‌های سطح بالا مانند پایتون در نهایت منجر به از دست دادن سرعت اجرا در محاسبات می‌شود. هر چند کتابخانه‌ها با انجام موازی‌سازی‌هایی در سطح سیستم، باعث تسریع این فرآیند شده‌اند اما این سرعت باز هم برای مدل‌های بسیار پیچیده و مخصوصاً برای پردازش‌های Real Time بسیار پایین است. راهکارهایی که برای این موضوع وجود دارد استفاده حداکثری از پردازش‌های موازی و GPU است اما استفاده از این سخت‌افزارها بسیار هزینه‌بر است. در این پروژه به پیاده‌سازی سخت‌افزاری شبکه عصبی MLP و CNN پرداخته‌ایم. زمانبندی لایه‌های مختلف محاسبه شده و با یک‌دیگر مقایسه شده است. در این پروژه از پروتکل‌های عمومی و استاندارد استفاده شده است تا پیاده‌سازی آن جامع و ساده باشد. همچنین به دلیل flexibility زیادی که سخت‌افزار دارد می‌توان با اضافه کردن یا کم کردن واحدهای محاسباتی، به سرعت بیشتری (به بهای هزینه بیشتر) رسید. تمامی پیاده‌سازی‌ها در این پروژه با استفاده از حداقل سخت‌افزار بوده که باعث هزینه کمتر و مصرف توان کمتر می‌شود و با وجود کمترین واحدهای محاسباتی همچنان سرعت به طرز چشمگیری بیشتر از نرم‌افزار می‌باشد.

در این بخش به شرح فعالیت‌ها و روند کار می‌پردازیم. همچنین روش‌های انجام آن به صورت کامل شرح داده می‌شود. بخش‌ها به ترتیب روند اجرایی می‌باشد و نتیجه‌های این بخش‌ها در بخش ۴ (نتایج) موجود است.

۳.۱- شبکه عصبی MLP برای تشخیص اعداد دست‌نوشته

۳.۱.۱- مقدمه‌ای بر شبکه MLP

شبکه‌های عصبی چند لایه (Multilayer Perceptron یا MLP) یکی از انواع رایج شبکه‌های عصبی است و برای بسیاری از وظایف مورد استفاده قرار می‌گیرد. این شبکه‌ها برای مجموعه وسیعی از مسائل قابل استفاده است و دارای ویژگی‌های مناسبی هستند که آن‌ها را در برخی از موارد مورد ترجیح قرار می‌دهد. در زیر به برخی از دلایل استفاده از MLP در شبکه‌های عصبی اشاره خواهیم کرد:

۱- **قدرت تقریبی:** MLP قدرت تقریبی بالایی دارد، به این معنی که می‌تواند توابع پیچیده را با دقت بالا تقریب بزند. با استفاده از لایه‌های مخفی (hidden layers) و با تعداد مناسب نورون‌ها، MLP می‌تواند تقریبی دقیق از توابع غیرخطی را توسط مجموعه‌ای از توابع خطی ارائه دهد.

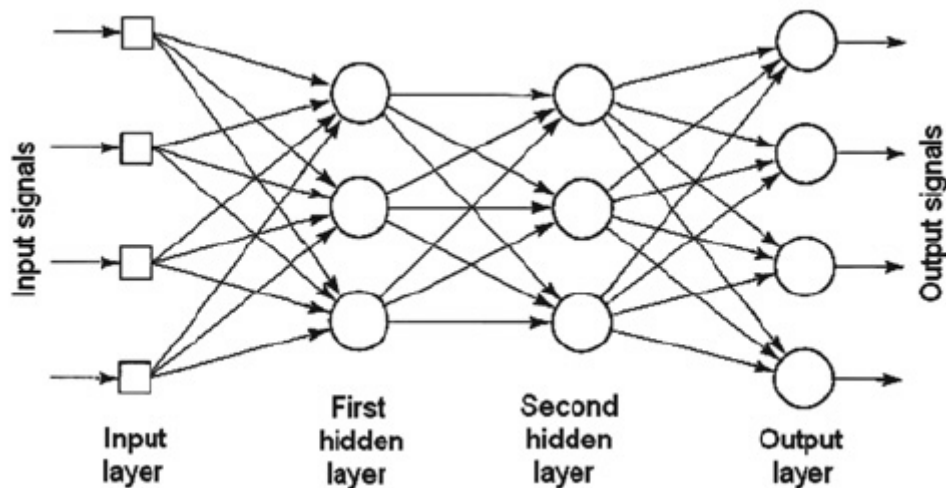
۲- **آموزش قابل تعمیم:** MLP معمولاً با استفاده از الگوریتم Backpropagation و روش‌های بهینه‌سازی، قابلیت آموزش را دارند. با این روش، شبکه می‌تواند از رویکرد نمونه‌های آموزشی تعمیم بیاموزد و قادر به پیش‌بینی برای نمونه‌های جدید باشد.

۳- **قابلیت استفاده در وظایف تشخیص الگو:** MLP به خوبی برای وظایف تشخیص الگو مناسب است. با ترکیب چندین لایه عصبی در MLP، شبکه قادر است ویژگی‌های پیچیده‌تر را از داده‌ها استخراج کند و الگوهای پنهان را شناسایی کند.

۴- **پیچیدگی قابل تنظیم:** تعداد لایه‌ها و تعداد نورون‌ها در هر لایه قابل تنظیم است، به این معنی که با تغییر تعداد لایه‌ها و نورون‌ها می‌توان انعطاف پذیری بیشتری در شبکه‌های عصبی داشت و آن‌ها را با دقت و توانایی مناسب برای وظایف مورد نیاز پیکربندی کرد.

۵- **قابلیت استفاده در وظایف غیرخطی:** MLP قادر است توابع غیرخطی را تقریب بزند. این امکان به MLP می‌دهد که در وظایفی که روابط پیچیده‌تری وجود دارد، مفید باشد.

در شکل ۳.۱.۱.۱ یک نمونه شبکه MLP با دو لایه میانی که هر کدام ۳ hidden unit دارند را مشاهده می‌کنیم.



شکل ۳.۱.۱.۱- نمونه یک شبکه MLP با دو لایه میانی

۳.۱.۲- پیاده‌سازی

ابتدا برای مدل و train کردن به سراغ ساده‌ترین معماری برای تشخیص تصویر می‌رویم. این معماری شامل دو لایه میانی می‌باشد. برای انتخاب اینکه از چند hidden unit استفاده کنیم، باید بر اساس دقتی که به دست می‌آید عمل کنیم. برای همین کار از لایه‌های کم به لایه‌های زیاد با استفاده از زبان پایتون، کدی می‌نویسیم که از تعداد لایه ۱ شروع کرده و یک مدل train کند و دقت را حساب کند و در یک فایل csv یادداشت کند؛ پس از این کار آن تعداد لایه‌ها را یکی زیاد کرده و این کار را از اول تکرار کند. این کار تا زمانی که دقت بیشتر از ۹۸ درصد بشود ادامه پیدا می‌کند.

ورودی شبکه یک عکس ۲۸ در ۲۸ بوده و خروجی ۱۰ تا عدد که احتمال هر کدام از ۰ تا ۹ بودن را مشخص می‌کند. در نهایت بر اساس نتایج قسمت ۴.۱ انتخاب می‌شود که از ۱۰۰ تا hidden unit استفاده شود تا هم دقت خوبی داشته باشیم و هم Resource بیش از حد استفاده نکنیم.

پس از انتخاب تعداد لایه، مدل train شده شبکه را در یک فایل ذخیره کرده تا بعدتر از آن استفاده کنیم. سپس به بررسی و تست این شبکه روی داده‌های واقعی و چیزی به جز مجموعه داده MNIST می‌پردازیم. این مدل دقت خوبی در تشخیص اعداد نداشته و مجبور به تغییر مدل خود می‌شویم.

۳.۲- شبکه عصبی CNN برای تشخیص اعداد دست‌نوشته

۳.۲.۱- مقدمه‌ای بر شبکه عصبی CNN

شبکه‌های عصبی CNN یا Convolutional Neural Networks برای پردازش داده‌های دارای ساختار شبکه‌ای، مانند تصاویر و سیگنال‌های صوتی، استفاده می‌شوند. این شبکه‌ها به خاطر ویژگی‌های خاص خود، در بسیاری از وظایف بینایی ماشین و پردازش تصویر عملکرد بسیار خوبی دارند. در زیر به برخی از مزایا، معایب و موارد استفاده از شبکه‌های عصبی CNN خواهیم پرداخت:

مزایا:

- استخراج ویژگی‌های سلسله مراتبی: شبکه‌های عصبی CNN، با استفاده از لایه‌های Convolution و لایه‌های Pooling، توانایی استخراج ویژگی‌های سلسله مراتبی از داده‌ها را دارند. این به معنی آن است که شبکه‌ها می‌توانند ویژگی‌های ساده‌تر مانند خطوط و لبه‌ها را در لایه‌های اولیه شناسایی کرده و به ویژگی‌های پیچیده‌تر مانند الگوها و اشیاء در لایه‌های بالاتر پردازند.
- اشتراک پارامترها: یکی از ویژگی‌های مهم CNN، اشتراک پارامترها است. این به معنی آن است که وزن‌ها و پارامترهای استفاده شده در یک لایه Convolution برای تمام نقاط داده مشترک هستند. این ویژگی باعث می‌شود که تعداد قابل تنظیم پارامترها در CNN به صورت قابل تحملی کاهش یابد و در نتیجه مدل قابلیت یادگیری و تعمیم بیشتری داشته باشد.
- کاهش تعداد پارامترها: استفاده از لایه‌های Pooling در CNN به کاهش تعداد پارامترها و حجم داده ورودی کمک می‌کند. این موضوع به افزایش سرعت آموزش و پیش‌بینی مدل کمک می‌کند.

معایب:

- نیاز به داده آموزش بزرگ: استفاده از CNN معمولاً نیاز به مجموعه داده‌های آموزش بزرگی دارد. برای آموزش شبکه‌های عصبی CNN، ممکن است نیاز به مجموعه داده‌های بزرگی باشد که ممکن است در برخی برنامه‌ها محدودیت‌هایی ایجاد کند.
- پیچیدگی محاسباتی: شبکه‌های عصبی CNN، معمولاً پیچیدگی محاسباتی بالایی دارند، به ویژه در موارد استفاده از شبکه‌های عمیق. آموزش و استفاده از این شبکه‌ها ممکن است زمان و منابع محاسباتی زیادی را مصرف کند.

موارد استفاده:

- تشخیص الگو: به خاطر قابلیت استخراج ویژگی‌های تصویری، در وظایف تشخیص الگو مانند تشخیص چهره‌ها، تشخیص اشیاء و تشخیص افراد استفاده می‌شوند.
- دسته‌بندی تصاویر: برای دسته‌بندی تصاویر در مواردی مانند تشخیص شیء در تصاویر، تشخیص بیماری‌ها از تصاویر پزشکی و تشخیص اشیاء در تصاویر مورد استفاده قرار می‌گیرند.
- ترجمه ماشینی: در وظایف ترجمه ماشینی بر روی متن‌ها و جملات مورد استفاده قرار می‌گیرند.

- تشخیص سیگنال‌های صوتی: برای تشخیص الگوها و ویژگی‌های صوتی در مواردی مانند تشخیص سیگنال‌های صوتی و تشخیص سیگنال‌های زبانی استفاده می‌شوند.

۳.۲.۲- پیاده‌سازی

با توجه به ویژگی‌های مذکور در بخش قبلی بهتر است که برای تشخیص تصویر از مدل CNN استفاده کنیم. برای این کار یک شبکه عصبی با لایه‌ها و مشخصات زیر train می‌کنیم:

۱- لایه Convolution 2D + تابع فعال‌سازی ReLU (ورودی ۲۸ در ۲۸؛ کرنل ۳ در ۳؛ ۳۲ کانال)

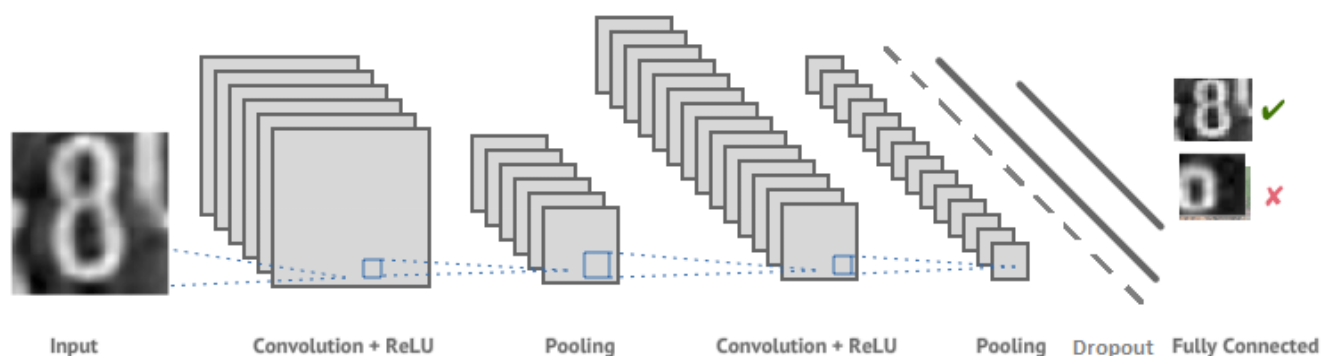
۲- لایه Max Pooling 2D (ورودی ۲۶ در ۲۶؛ ۳۲ کانال)

۳- لایه Convolution 2D + تابع فعال‌سازی ReLU (ورودی ۱۳ در ۱۳؛ کرنل ۳ در ۳؛ ۶۴ کانال)

۴- لایه Max Pooling 2D (ورودی ۱۱ در ۱۱؛ ۶۴ کانال)

۵- لایه Dense یا Fully Connected + تابع فعال‌سازی Softmax (ورودی ۱۶۰۰؛ خروجی ۱۰)

شکل لایه‌ها در شکل ۳.۲.۲.۱ قابل مشاهده است.



شکل ۳.۲.۲.۱- شبکه عصبی CNN طراحی شده برای تشخیص اعداد دست‌نوشته از روی تصویر

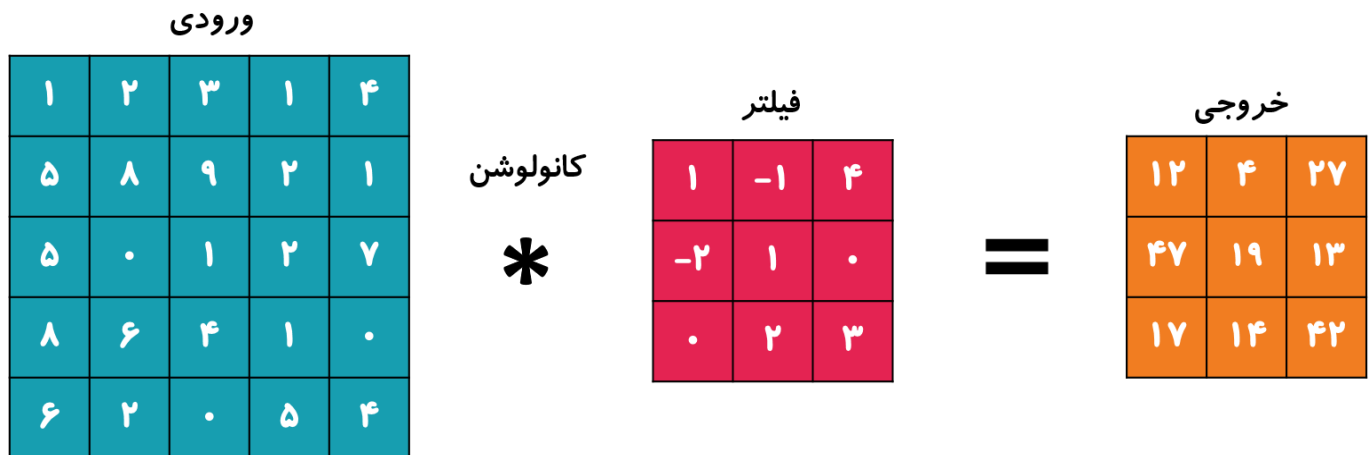
پس از پیاده‌سازی این مدل و train کردن این شبکه عصبی و ذخیره این مدل به عنوان فایل برای استفاده‌های بعدی، باید تست‌های خارج از مجموعه داده MNIST را (مانند بخش قبل) انجام دهیم؛ این بار (برعکس بخش قبل) دقت بسیار خوبی گرفته و می‌توانیم ادامه راه را با همین مدل ادامه دهیم (نتایج در بخش ۴.۴).

۳.۳- استخراج ضرایب و پیاده‌سازی شبکه CNN بدون استفاده از کتابخانه

با توجه به اینکه در پیاده‌سازی سخت‌افزاری به کتابخانه‌های هوش مصنوعی دسترسی نداشته و تنها می‌توانیم از ضرب‌کننده، جمع‌کننده و بقیه عملیات‌های پایه استفاده کنیم، نیاز است تا تمامی لایه‌ها بدون استفاده از کتابخانه و تنها با استفاده از حلقه، ضرب و جمع، انجام شود؛ به همین منظور در این بخش در ابتدا مقدار ضرایب از فایل مدل train شده که قبلاً ذخیره کرده بودیم، استخراج می‌کنیم و از آن به منظور انجام عملیات محاسباتی استفاده می‌کنیم.

۳.۳.۱- لایه Convolution

این لایه از یک کرنل (فیلتر) n در n و چند کانال تشکیل شده است. فرآیند آن به این صورت است که در کرنل روی تصویر قرار گرفته و تمام نقاط نظیر به هم در ورودی ضرب شده و در نهایت همه با هم جمع می‌شوند و خروجی را تشکیل می‌دهند (شکل ۳.۳.۱.۱).

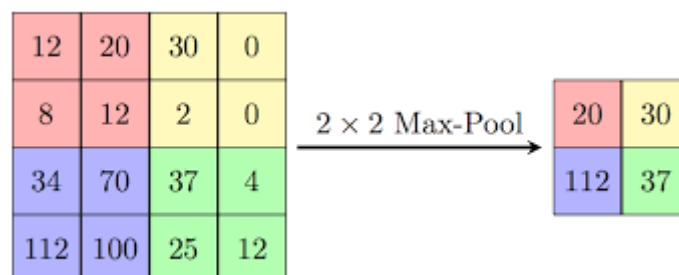


شکل ۳.۳.۱.۱

شکل ۳.۳.۱.۱ یکی از فیلترها بود و در لایه اصلی k تا فیلتر داریم که هر کدام را یک کانال می‌گوییم. برای مثال در لایه اول ۳۲ تا کانال داریم که یعنی در مجموع ۳۲ فیلتر داشته که با اعمال روی ورودی ۳۲ تصویر خروجی خواهیم داشت.

۳.۳.۲- لایه Max Pooling

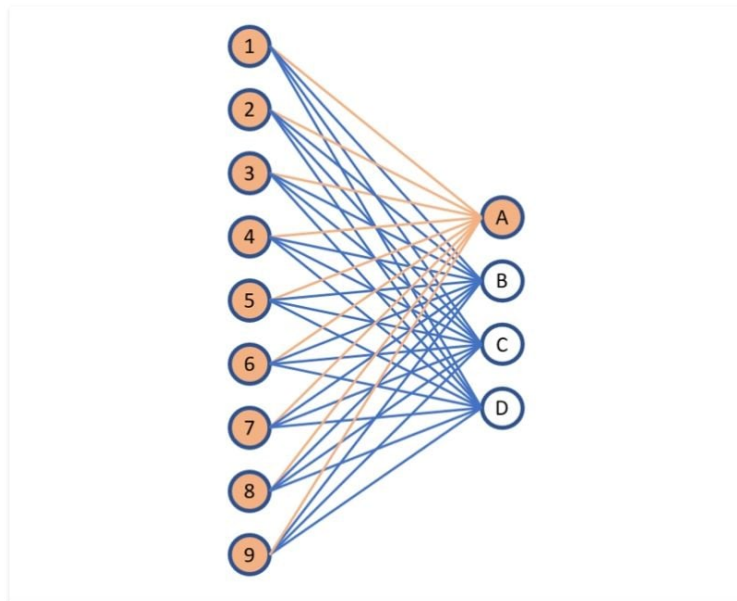
در این لایه یک پنجره ۲ در ۲ وجود دارد که با پیمایش تصاویر ورودی در آن پنجره ۲ در ۲، بیشترین مقدار را انتخاب می‌کند و بعد از آن طوری حرکت می‌کند که همپوشانی با ورودی‌های قبلی نداشته باشد.



شکل ۳.۳.۲.۱- نمونه‌ای از اعمال Max Pooling بر روی ورودی

۳.۳.۳- لایه Dense یا Fully Connected

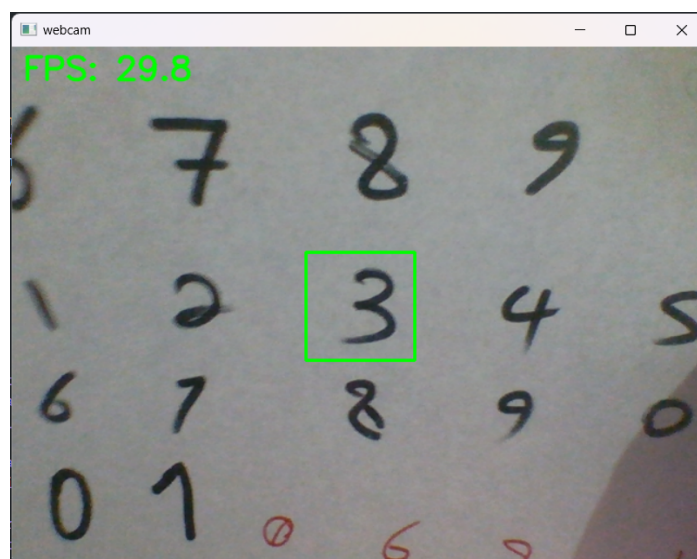
در این لایه، یک تعداد ورودی و یک تعداد خروجی با یک سایز مشخص داریم؛ تمامی ورودی‌ها و خروجی‌ها با یک یال به هم متصل شده‌اند و هر کدام از این یال‌ها دارای وزن می‌باشند. این وزن در ورودی ضرب شده و با ضرب ورودی بقیه یال‌ها در ورودی‌های دیگر (به ازای تمام ورودی‌ها) جمع می‌شود.



شکل ۳.۳.۳- یک نمونه شبکه Dense

۳.۴- اضافه کردن ماژول دوربین

در اینجا برای اینکه ورودی به صورت ثابت نباشد بتوانیم داده‌ها را از بیرون و توسط ماژول دوربین جمع‌آوری کنیم، از ماژول دوربین استفاده می‌کنیم؛ به این صورت که با استفاده از کتابخانه OpenCV یک پنجره باز کرده و با دسترسی به ماژول دوربین تصاویر آن را به صورت frame دریافت کرده و آن را روی ویندوز باز شده نمایش می‌دهیم. سپس یک کادر به اندازه کافی بزرگ با دور سبز رنگ ایجاد می‌کنیم تا یک مرزی برای عکسی که ذخیره می‌کنیم قرار دهیم (به دلیل بزرگ بودن بیش از حد تصویر هر frame)؛ با این کار تنها کافیسیت که عدد خود را درون و در مرکز این مربع قرار دهیم و سپس با زدن دکمه enter در زمان مناسب تصویر پردازش شده و خروجی داده شود. در شکل ۳.۴.۱ می‌توان نمونه‌ای از فرآیند گفته شده را مشاهده کرد.



شکل ۳.۴.۱- نمونه‌ای از اتصال ماژول دوربین به مانیتور لپ‌تاپ

۳.۵- پردازش به صورت Real Time

برای اجرای بهتر و داشتن دید بهتر به نحوه انجام پردازش، به جای زدن دکمه enter به این صورت عمل می‌کنیم که به برنامه در قدم اول یک آپشن Real Time و یک آپشن Manual اضافه می‌کنیم. پس از این کار به این صورت عمل می‌کنیم که هر زمان که برنامه توانست شروع به کار کند و تصویر را ذخیره و پیش‌بینی‌های لازم را روی آن انجام دهد؛ هر زمان که خروجی حاضر شد آن را چاپ کرده و تصویر بعدی را ذخیره کرده و مجدد کارهای بالا را تکرار کند.

این پیش‌بینی به دو صورت (۱) با استفاده از کتابخانه و (۲) بدون استفاده از کتابخانه قابل انجام است. در حالتی که از کتابخانه استفاده می‌کنیم، FPS بهتری می‌گیریم و این به دلیل موازی‌سازی‌هایی است که در کتابخانه Keras انجام می‌پذیرد. اما پیاده‌سازی بدون استفاده از کتابخانه به دلیل موازی نبود کارها، FPS کمتری می‌گیریم.

۳.۶- اجرا کردن شبکه CNN بر روی بخش PS برد PYNQ

۳.۶.۱- برد PYNQ-Z2

برد PYNQ یک برد توسعه مبتنی بر فریمورک Xilinx Zynq است که برای توسعه و برنامه‌ریزی سیستم‌های الکترونیکی و رایانه‌ای استفاده می‌شود. این برد با همکاری بین شرکت Xilinx و دانشگاه کالیفرنیا در سانتاکروز توسعه داده شده است.

برد PYNQ دارای سخت‌افزاری قدرتمند است که از پردازنده ARM Cortex-A9 و FPGA بر پایه زینک Xilinx تشکیل شده است. این ترکیب اجازه می‌دهد تا برنامه‌هایی با قابلیت‌های پیشرفته و پردازش سریع را اجرا کرده و سخت‌افزارهای قابل برنامه‌ریزی را پیاده‌سازی کند.

این برد با استفاده از محیط توسعه PYNQ، که بر پایه فریمورک Jupyter Notebook است، به برنامه‌نویسان امکان می‌دهد تا به سادگی کدهای Python برای کنترل سخت‌افزار و اجرای الگوریتم‌های پردازشی مستقل از سیستم عامل بنویسند. این محیط توسعه به کاربران اجازه می‌دهد تا از ویژگی‌های PYNQ بهره‌برداری کنند، از جمله پیکربندی سخت‌افزار FPGA، ارتباط با واسطه‌های دیگر مانند GPIO و I2C و اتصال به شبکه‌های بی‌سیم مانند Wi-Fi و Bluetooth.

با استفاده از برد PYNQ، می‌توان به سادگی پروژه‌های الکترونیکی پیچیده را پیاده‌سازی کرد و از قدرت پردازشی و قابلیت‌های FPGA بهره‌برداری کرد. این برد مناسب برای کاربران حرفه‌ای و آموزشی است و می‌تواند در زمینه‌های مختلفی مانند سیستم‌های مبتنی بر هوش مصنوعی، اینترنت اشیا، پردازش تصویر و صوت، رباتیک و بسیاری از برنامه‌های الکترونیکی دیگر مورد استفاده قرار بگیرد.

۳.۶.۲- پیاده‌سازی

اتصال به پایتون این نرم‌افزار به این طریق است که کابل LAN را به برد وصل می‌کنیم. همچنین یک کابل LAN دیگر را که به همان شبکه وصل است متصل می‌کنیم. با این کار برد PYNQ در شبکه یک سرور به صورت لوکال می‌سازد تا وسیله‌های دیگر از جمله PC بتوانند به طریق آن و با استفاده از مرورگر به آن متصل شده و وارد محیط Jupyter Notebook شوند.

برای پیدا کردن IP مورد نظر جهت اتصال کافی است که با اتصال PC به پورت Micro USB برد PYNQ و اجرا نرم‌افزاری که از پورت COM پشتیبانی کند (مانند Tera Term) به آن متصل شویم. برای دریافت و ارسال درست اطلاعات باید baud rate را در تنظیمات برابر ۱۱۵۲۰۰ قرار دهیم. سپس با دستور `ifconfig` اولین IP را کپی کرده و در مرورگر paste می‌کنیم و پس از آن به محیط Jupyter دسترسی داریم.

۳.۷- ساخت سیستم‌های ساده و ارتباط PS و PL برد PYNQ

در ابتدا برای اطمینان از کارکرد درست ماژول‌ها و اتصال درست PS و PL یک طراحی ساده انجام می‌دهیم تا قسمت PS و PL به هم متصل شده و قسمت PS یک سری داده را به DMA فرستاده و DMA آن را به یک FIFO بریزد و FIFO آن را به DMA برگرداند تا در نهایت DMA دوباره آن را در یک قسمت دیگر در حافظه نوشته و PS آن را بخواند. در صورت دیزاین درست، PS باید همان داده‌هایی که فرستاده است را بخواند. قبل از این که این ماژول را پیاده‌سازی کنیم، لازم است در ابتدا با ساختار PL و PS و نحوه ارتباط این دو قسمت آشنا شویم.

۳.۷.۱- ساختار PL و PS برد PYNQ

در برد PYNQ، ساختار PL (برنامه‌پذیر قابل برنامه‌ریزی) و PS (سیستم پردازنده) بر پایه فریمورک Xilinx Zynq قرار دارند.

PS یا همان Processing System یک بخش سخت‌افزاری است که شامل پردازنده ARM Cortex-A9 است. این پردازنده دارای سیستم‌عامل Linux است و می‌تواند برنامه‌ها را اجرا کند. PS وظیفه کنترل کلی سیستم را بر عهده دارد و مسئولیت‌هایی مانند مدیریت حافظه، کنترل پورت‌های ورودی/خروجی، و مدیریت منابع سیستمی را بر عهده دارد.

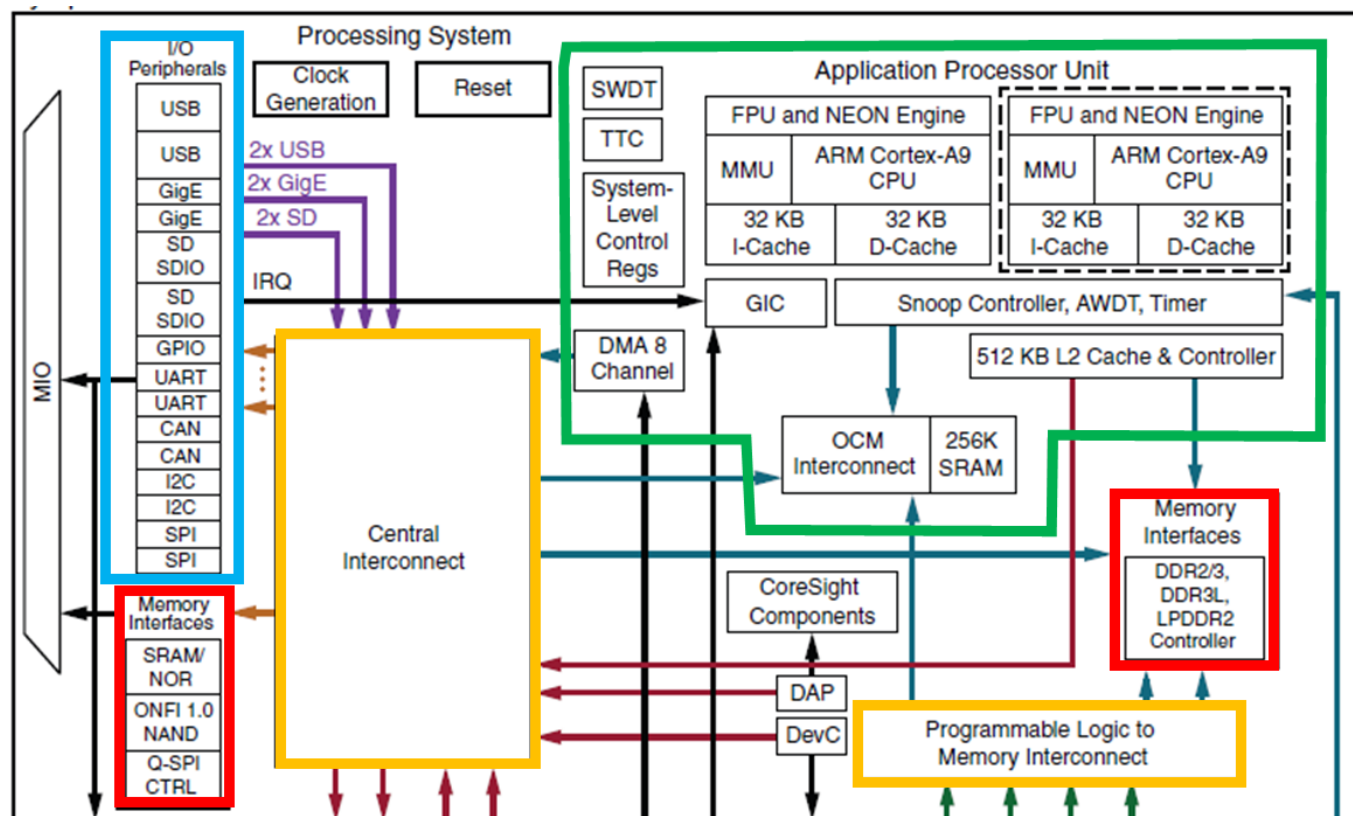
PL یا همان Programmable Logic یک بخش قابل برنامه‌ریزی است که شامل FPGA است. FPGA قابلیت برنامه‌ریزی و تنظیم مجدد مدارهای منطقی را دارد. با استفاده از زبان‌های سخت‌افزاری مانند VHDL یا Verilog، می‌توان سخت‌افزارهای خاصی را در FPGA پیاده‌سازی کرد. این امکان را به ما می‌دهد تا عملکرد سخت‌افزار را به طور دقیق و سفارشی تنظیم کرده و برنامه‌های خاصی را اجرا کنیم.

برای ارتباط بین PL و PS در برد PYNQ، وجود رابط‌های ارتباطی متنوعی از قبیل AXI یا همان Advanced eXtensible Interface و رابط‌های DMA یا همان Direct Memory Access امکان‌پذیر است. این رابط‌ها به PS اجازه می‌دهند با PL ارتباط برقرار کند و داده‌ها را به صورت مستقیم، منتقل کند. با استفاده از رابط‌های

ارتباطی این دو بخش، می‌توان داده‌ها و سیگنال‌ها را بین PS و PL انتقال داد و در عملکرد سیستم تغییری ایجاد کرد.

در مجموع، ارتباط بین PS و PL در برد PYNQ از طریق رابط‌های ارتباطی قابل برنامه‌ریزی برقرار می‌شود، که به PS امکان کنترل و ارسال داده به PL را می‌دهد و از طرف دیگر، با برنامه‌ریزی FPGA در PL، می‌توان سخت‌افزارهای خاصی را پیاده‌سازی کرده و با PS ارتباط برقرار کرد.

در شکل ۳.۷.۱.۱ می‌توان موارد گفته شده را مشاهده کرد.

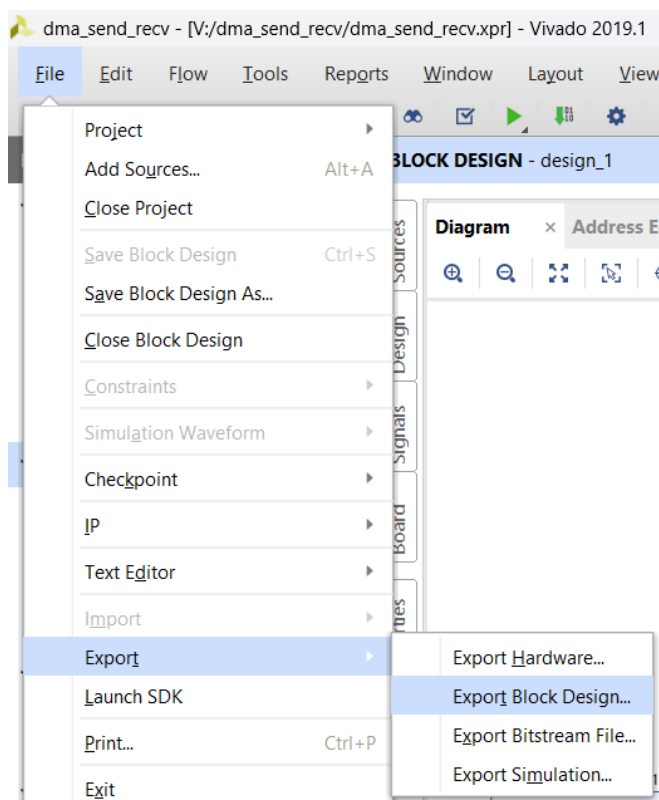


شکل ۳.۷.۱.۱- ساختار درونی برد PYNQ (معماری کلی ZYNQ)

۳.۷.۲- پیاده‌سازی

ابتدا باید دیزاین سخت‌افزاری خود را درون نرم‌افزار Vivado انجام دهیم. برای این کار ابتدا یک Block Design جدید درست می‌کنیم؛ سپس به آن ZYNQ Processing System را اضافه کرده و پورت‌های High Performance را فعال می‌کنیم (HP0 و HP2 هر کدام در حالت ۶۴ بیت). سپس DMA را افزوده و کانفیگ آن را مطابق شکل ۳.۷.۲.۱ انجام می‌دهیم.

پس از انجام این کار یک wrapper برای دیزاین خود درست کرده و در نهایت Generate Bitstream را می‌زنیم. این عملیات ممکن است زمان‌گیر باشد. پس از اتمام انجام کار و تولید شدن Bitstream، از گزینه‌های مشخص شده در شکل ۳.۷.۲.۳، گزینه‌های Export Block Design و Export Bitstream File را انتخاب می‌کنیم.



شکل ۳.۷.۲.۳- نحوه خروجی گرفتن دیزاین سخت‌افزار

حال فایل‌های خروجی را در Jupyter Notebook آپلود کرده و با اجرا کد پایتون مربوطه، دیزاین خود را تست می‌کنیم.

۳.۸- پروتکل AXI-4 Stream

پروتکل AXI-4 Stream یک پروتکل ارتباطی است که برای انتقال داده‌های پیوسته بین ماژول‌ها در سیستم‌های مبتنی بر FPGA استفاده می‌شود. این پروتکل از سری پروتکل‌های AXI شرکت Xilinx است و برای انتقال داده‌هایی با حجم بالا و نرخ انتقال سریع طراحی شده است. معماری AXI-4 Stream به صورت یکپارچه و ساده‌تری نسبت به دیگر نسخه‌های پروتکل AXI عمل می‌کند و تمرکز اصلی آن بیشتر بر روی انتقال داده‌ها و کمتر برای کنترل و کانفیگ کردن دارد. ویژگی‌های کلیدی AXI-4 Stream عبارتند از:

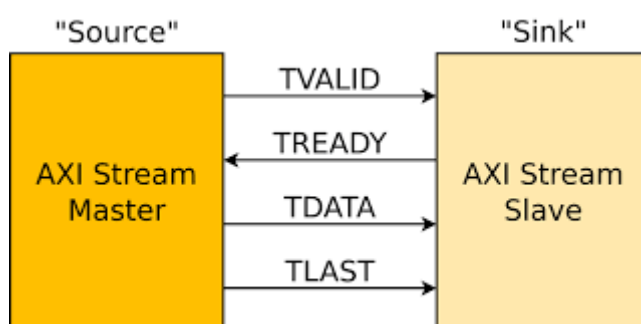
۱- **ارتباط پیوسته:** AXI-4 Stream برای انتقال داده‌ها از روی یک سیگنال پیوسته استفاده می‌کند، به این معنی که داده‌ها به صورت پیوسته و بدون وقفه انتقال می‌یابند.

۲- **نبود handshake:** در AXI-4 Stream، نیاز به مکانیزم handshake برای تایید دریافت یا ارسال داده وجود ندارد. بنابراین، انتقال داده‌ها مستقل از سرعت یا نرخ انتقال برقرار می‌شود.

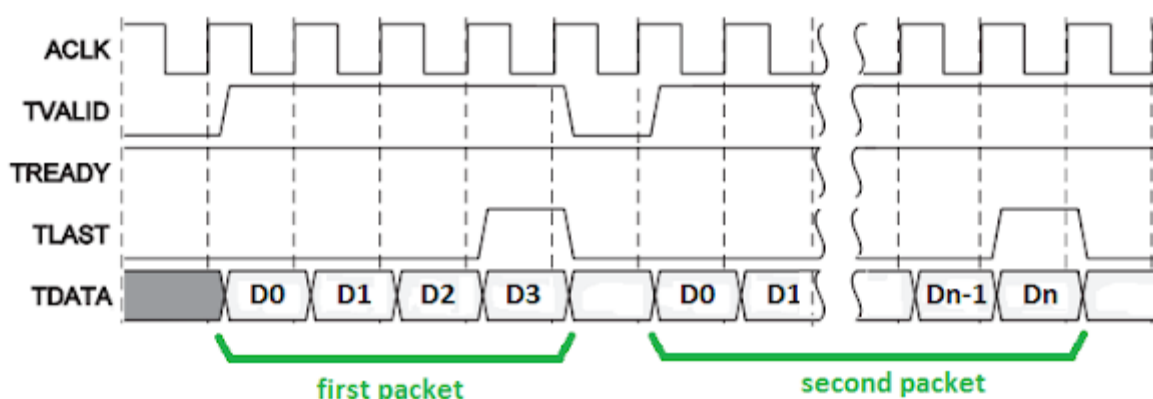
۳- **نبود پشتیبانی از چند کانال:** AXI-4 Stream تنها یک کانال انتقال داده را پشتیبانی می‌کند و از امکانات چندکاناله پیشرفته AXI-4 و AXI-4 Lite صرف نظر می‌کند.

۴- **ساختار ساده:** AXI-4 Stream با ساختاری ساده‌تر نسبت به AXI-4 عمل می‌کند، که باعث کاهش پیچیدگی و مصرف منابع سخت‌افزاری می‌شود.

AXI-4 Stream مناسب برای انتقال داده‌های پیوسته و پرسرعت در سیستم‌هایی است که نیاز به پردازش و انتقال داده‌های بزرگ و پیوسته دارند، مانند فرآیندهای پردازش تصویر، پردازش سیگنال‌های دیجیتال، پردازش سیگنال‌های صوتی و غیره.



شکل ۳.۸.۱- نحوه اتصال دو ماژول Master و Slave



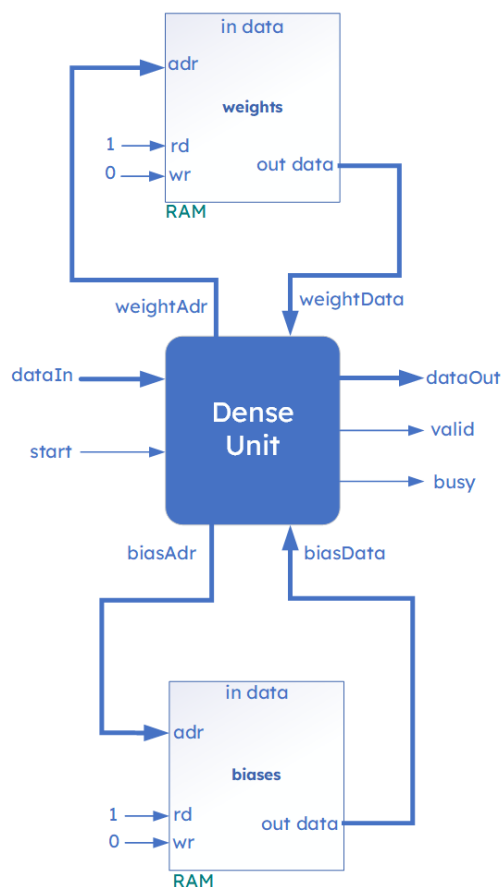
شکل ۳.۸.۲- سیگنال‌ها و شکل موج پروتکل AXIS

۳.۹- طراحی ماژول General Purpose AXI-4 Stream Interface

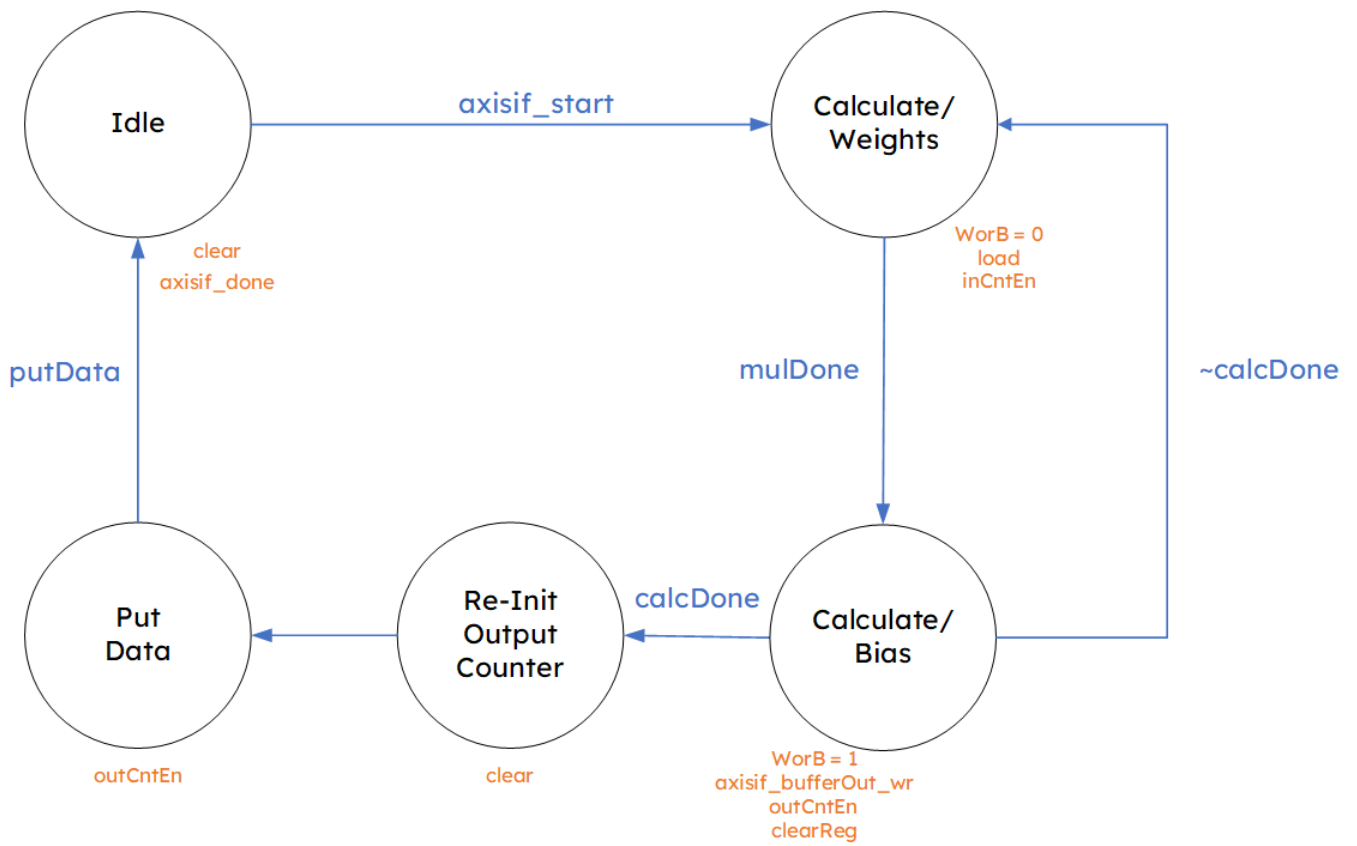
به منظور ارتباط ساده‌تر با ماژول‌های دیگر به گونه‌ای که این ماژول عمومیت ساده‌ای داشته و وفق دادن ماژول‌ها با آن ساده باشد و اینکه ما را از پیچیدگی‌های پروتکل AXI-4 Stream دور نگه دارد، این ماژول طراحی و معماری شد. این interface ابتدا حالت Slave دارد؛ به گونه‌ای که تمامی اطلاعات ورودی را درون یک buffer نوشته و در اختیار ماژول متصل قرار می‌دهد و سپس آن ماژول را با یک سیگنال axisif_start فعال می‌کند. این ماژول تمام داده‌ها را در داخل buffer خروجی gp axis interface نوشته و با یک سیگنال axisif_done آن را مطلع می‌کند. پس از آن این interface وارد حالت Master شده و داده‌ها را خروجی می‌دهد.

۳.۱۰- ماژول Dense و ترکیب آن با GP AXIS IF

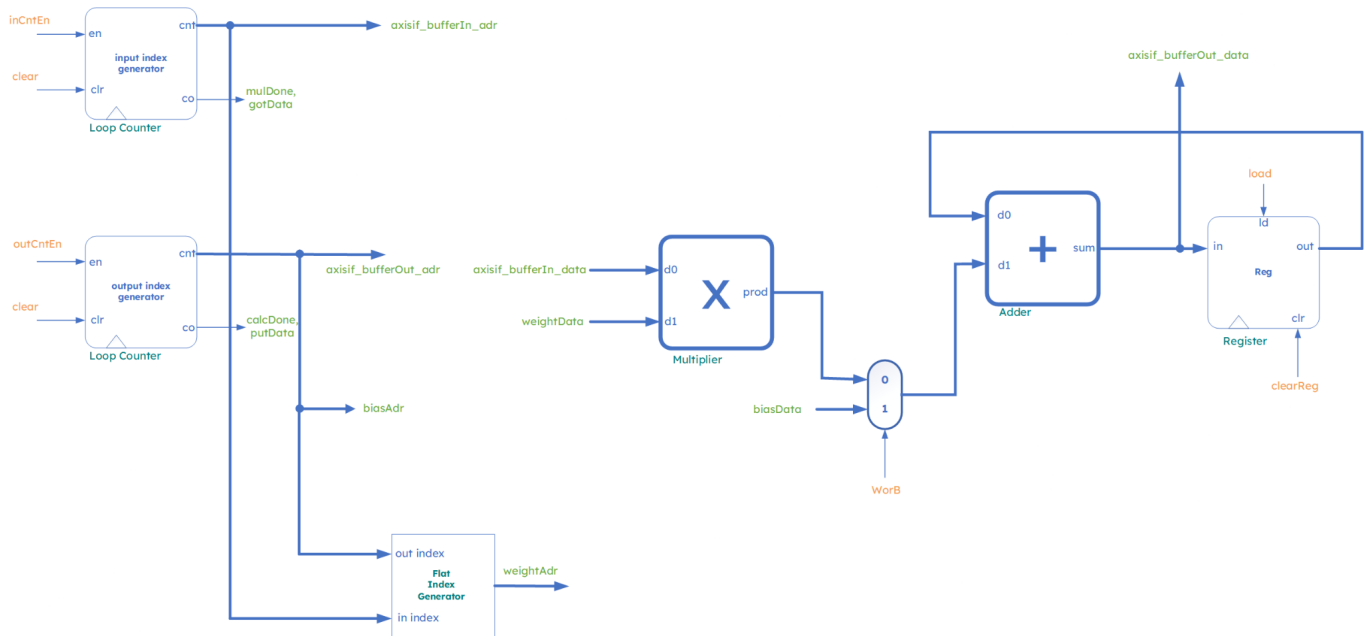
این ماژول تنها یک ضرب‌کننده و جمع‌کننده برای ضرب کردن وزن‌ها در ورودی و جمع زدن آن‌ها با هم و همینطور جمع زدن با bias استفاده می‌کند.



شکل ۳.۱۰.۱- ساختار کلی ماژول Dense



شکل ۳.۱۰.۲- ساختار Controller ماژول Dense



شکل ۳.۱۰.۳- ساختار Datapath ماژول Dense

در این بخش به نتایج گرفته شده از پروژه و مقایسه و تحلیل آن می‌پردازیم.

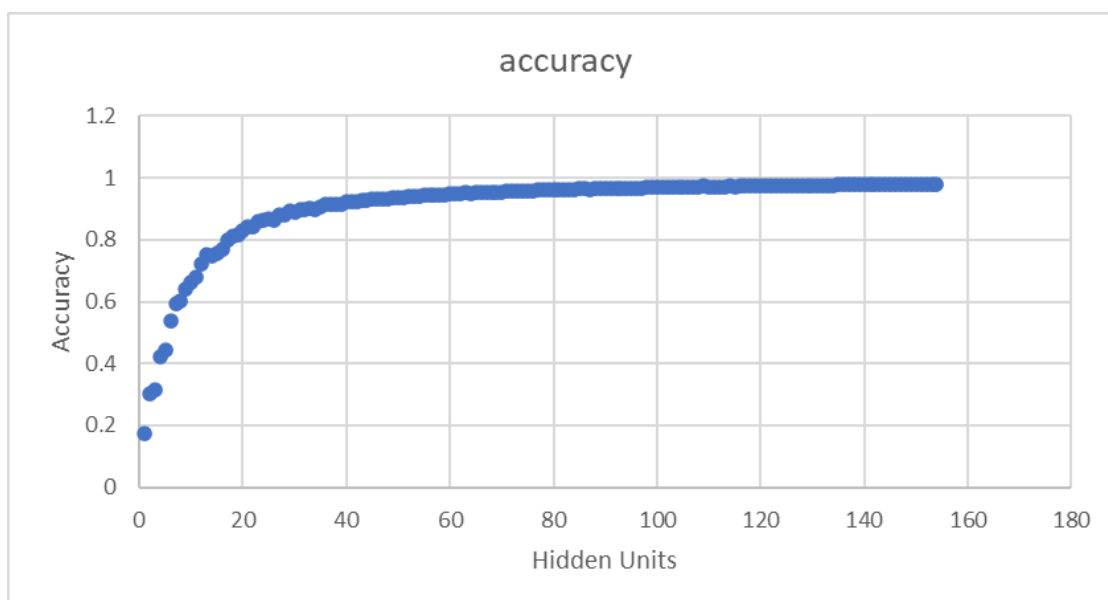
۴.۱- آموزش شبکه MLP برای تشخیص اعداد دست‌نوشته MNIST

در شکل ۴.۱.۱، می‌توانیم معماری و نحوه چینش لایه‌ها را مشاهده کنیم.

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 100)	78500
activation (Activation)	(None, 100)	0
dropout (Dropout)	(None, 100)	0
dense_1 (Dense)	(None, 100)	10100
activation_1 (Activation)	(None, 100)	0
dropout_1 (Dropout)	(None, 100)	0
dense_2 (Dense)	(None, 10)	1010
activation_2 (Activation)	(None, 10)	0
Total params: 89,610		
Trainable params: 89,610		
Non-trainable params: 0		

شکل ۴.۱.۱

در این شبکه دقت بر اساس hidden units بررسی شد و نتایج train آن در نمودار ۴.۱.۱ مشخص شده است:



نمودار ۴.۱.۱- دقت اعداد دست‌نوشته MNIST بر اساس تعداد hidden unit های شبکه MLP

در این نمودار مشاهده می‌شود که با افزایش دقت از یک تعداد hidden unit به بعد، دیگر دقت افزایش پیدا نمی‌کند پس برای استفاده کمتر از resource ها، عدد ۱۰۰ را انتخاب کردیم. این شبکه برای این تعداد hidden unit دارای دقت 0.9736 است.

۴.۲- درستی‌سنجی مدل train شده MLP بر روی داده‌های واقعی دست‌نوشته

۴.۲.۱- اعداد دست‌نوشته با رنگ سیاه

عدد	0	1	2	3	4	5	6	7	8	9
حدس	2	1	3	3	6	5	6	2	3	3
دقت	0.991	0.787	0.999	1.000	1.000	0.524	0.927	0.997	0.509	0.873
نمونه	0	1	2	3	4	5	6	7	8	9

جدول ۴.۲.۱.۱- نمونه تصاویر دست‌نوشته و دقت پیش‌بینی مدل

با توجه به جدول ۴.۲.۱.۱ مشاهده می‌کنیم که اعداد زیادی به اشتباه تشخیص داده می‌شوند؛ این در صورتی است که دقت train و test در مجموعه داده MNIST برابر با حدود ۹۷ درصد است و این موضوع احتمال overfit شدن را به ما می‌دهد.

در بخش‌های بعدی برای اعدادی با شکل، رنگ و ضخامت مختلف نیز این تست را اجرا می‌کنیم.

۴.۲.۲- اعداد دست‌نوشته با رنگ آبی، ضخامت نازک

عدد	0	1	2	3	4	5	6	7	8	9
حدس	5	5	5	5	5	5	5	5	5	5
دقت	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
نمونه	0	1	2	3	4	5	6	7	8	9

جدول ۴.۲.۲.۱- نمونه تصاویر دست‌نوشته و دقت پیش‌بینی مدل

مشاهده می‌شود که در تمام ستون‌ها، جواب ۵ و دقت ۱ بوده است.

۴.۲.۳- اعداد دست‌نوشته با رنگ آبی، ضخامت کلفت

عدد	0	1	2	3	4	5	6	7	8	9
حدس	5	5	3	5	5	5	5	3	3	5
دقت	1.000	1.000	0.995	1.000	1.000	1.000	1.000	1.000	0.934	1.000
نمونه	0	1	2	3	4	5	6	7	8	9

جدول ۴.۲.۳.۱- نمونه تصاویر دست‌نوشته و دقت پیش‌بینی مدل

همچنان ایراد بخش ۴.۲.۲ در این بخش نیز وجود دارد.

۴.۲.۴- نتیجه‌گیری

با توجه به نتایج این قسمت، معماری MLP برای توصیف کردن و پردازش این تصاویر مناسب نیست. علاوه بر آن زیاد کردن تعداد لایه‌های میانی نیز کار را بهتر نکرد.

۴.۳- آموزش شبکه CNN برای تشخیص اعداد دست‌نوشته MNIST

در شکل ۴.۳.۱ نحوه چینش لایه‌های CNN قابل مشاهده است.

```
x_train shape: (60000, 28, 28, 1)
60000 train samples
10000 test samples
Model: "sequential"
-----
)

conv2d_1 (Conv2D)          (None, 11, 11, 64)      18496
max_pooling2d_1 (MaxPooling (None, 5, 5, 64)      0
2D)

flatten (Flatten)          (None, 1600)            0
dropout (Dropout)          (None, 1600)            0
dense (Dense)               (None, 10)              16010

=====
Total params: 34,826
Trainable params: 34,826
Non-trainable params: 0
```

شکل ۴.۳.۱

در نهایت پس از train کردن به دقت 0.9908 رسیدیم.

۴.۴- درستی‌سنجی مدل train شده CNN بر روی داده‌های واقعی دست‌نوشته

۴.۴.۱- اعداد دست‌نوشته با رنگ سیاه

عدد	0	1	2	3	4	5	6	7	8	9
حدس	0	1	2	3	4	5	5	7	0	8
دقت	1.000	1.000	1.000	1.000	0.699	1.000	0.498	0.707	0.677	0.963
نمونه	0	1	2	3	4	5	6	7	8	9

جدول ۴.۴.۱.۱- نمونه تصاویر دست‌نوشته و دقت پیش‌بینی مدل

۴.۴.۲- اعداد دست‌نوشته با رنگ آبی، ضخامت نازک

عدد	0	1	2	3	4	5	6	7	8	9
حدس	2	2	2	2	2	2	2	2	2	2
دقت	0.390	0.392	0.411	0.509	0.456	0.358	0.404	0.322	0.399	0.448
نمونه	0	1	2	3	4	5	6	7	8	9

جدول ۴.۴.۲.۱- نمونه تصاویر دست‌نوشته و دقت پیش‌بینی مدل

۴.۴.۳- اعداد دست‌نوشته با رنگ آبی، ضخامت کلفت

عدد	0	1	2	3	4	5	6	7	8	9
حدس	0	1	2	3	4	5	5	2	0	9
دقت	0.999	0.963	0.999	1.000	0.852	0.991	0.508	0.685	0.707	0.822
نمونه	0	1	2	3	4	5	6	7	8	9

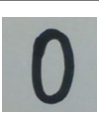

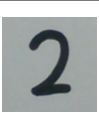
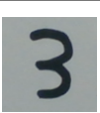
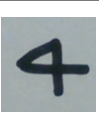
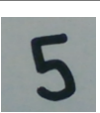
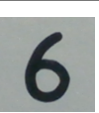
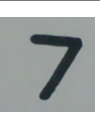

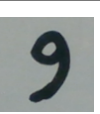
جدول ۴.۴.۳.۱- نمونه تصاویر دست‌نوشته و دقت پیش‌بینی مدل

۴.۴.۴- نتیجه گیری

در این بخش مشاهده کردیم که نتایج با معماری CNN به مراتب بهتر از MLP بود؛ در بعضی از تست‌کیس‌های CNN نیز تعداد خیلی کمی خطا وجود داشت که البته درصد دقت غالب آن‌ها کمتر از بقیه اعداد بود. در بخش ۴.۳.۲ و ۴.۴.۲ که اعداد دست‌نوشته آبی نازک بود، مشاهده کردیم که اعداد به خوبی پیش‌بینی نشدند و این به دلیل واضح نبودن اعداد بود. از طرفی دقت پیش‌بینی در معماری CNN برای این اعداد آبی نازک به مراتب پایین‌تر بود در حالی که این اعداد برای معماری MLP بسیاری بالا بود. به بیان دیگر در معماری MLP با درصد دقت بالایی، عدد را اشتباه حدس می‌زد اما در معماری CNN با درصد پایینی اشتباه تخمین زده می‌شد که این موضوع خود دلیلی برای استفاده از معماری CNN بود.

۴.۵- اضافه کردن ماژول دوربین

در جدول ۴.۵.۱ می‌توان اعداد داخل کادر را به عنوان عکس ثبت شده توسط ماژول دوربین مشاهده کرد. در این اعداد یک سری نکات وجود دارد که به آن‌ها خواهیم پرداخت.

عدد	0	1	2	3	4	5	6	7	8	9
حدس	0	1	2	3	4	5	6	7	8	9
دقت	0.936	0.848	0.967	0.967	0.874	0.976	0.515	0.674	0.800	0.922
نمونه										

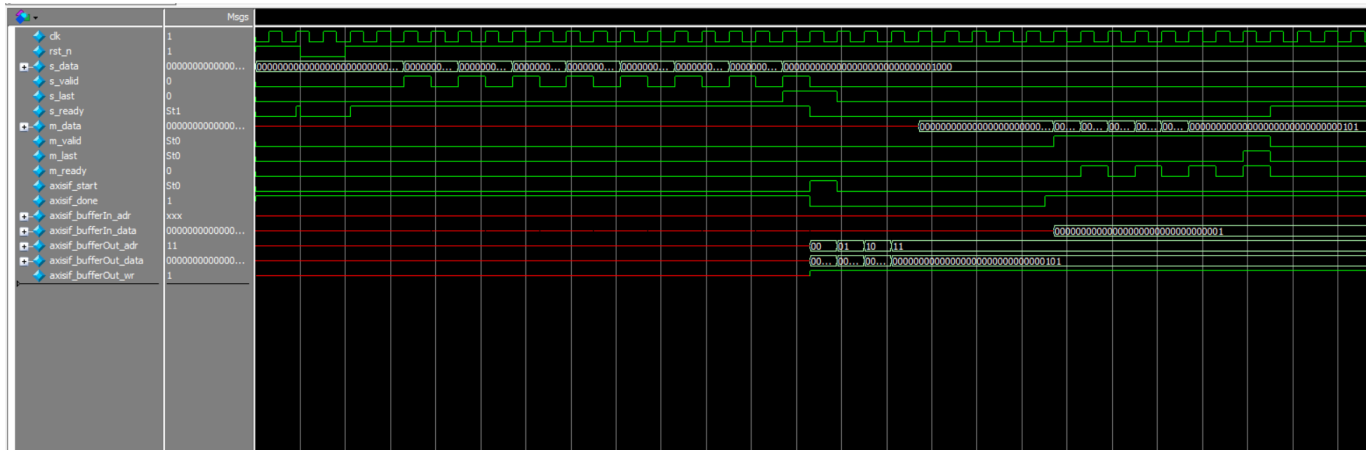
جدول ۴.۵.۱- عدد دست‌نوشته ثبت شده توسط ماژول دوربین و پیش‌بینی آن‌ها

تصاویر موجود در جدول ۴.۵.۱ از لحاظ ابعاد مانند داده‌های train مجموعه‌داده MNIST که ۲۸ در ۲۸ پیکسل است، نیست؛ پس برای این موضوع باید عملیات resize بر روی تصاویر ایجاد شود. نکته دیگری که وجود دارد این است که background تصاویر MNIST باید سیاه بوده و عدد باید با رنگ سفید باشد؛ به همین دلیل باید عملیات invert کردن روی تصاویر صورت بگیرد. از آنجایی که background این تصاویر کاملاً سفید نیست و هنگام invert کردن نیز کاملاً سفید نخواهد بود ممکن است روی دقت تاثیر بگذارد. از طرفی این اعداد باید کاملاً در مرکز عکس باشند تا شبیه مجموعه‌داده MNIST باشد.

با رعایت این نکات می‌توان به حداکثر دقت در عکس رسید.

۴.۶- شبیه‌سازی ماژول GP AXIS Interface

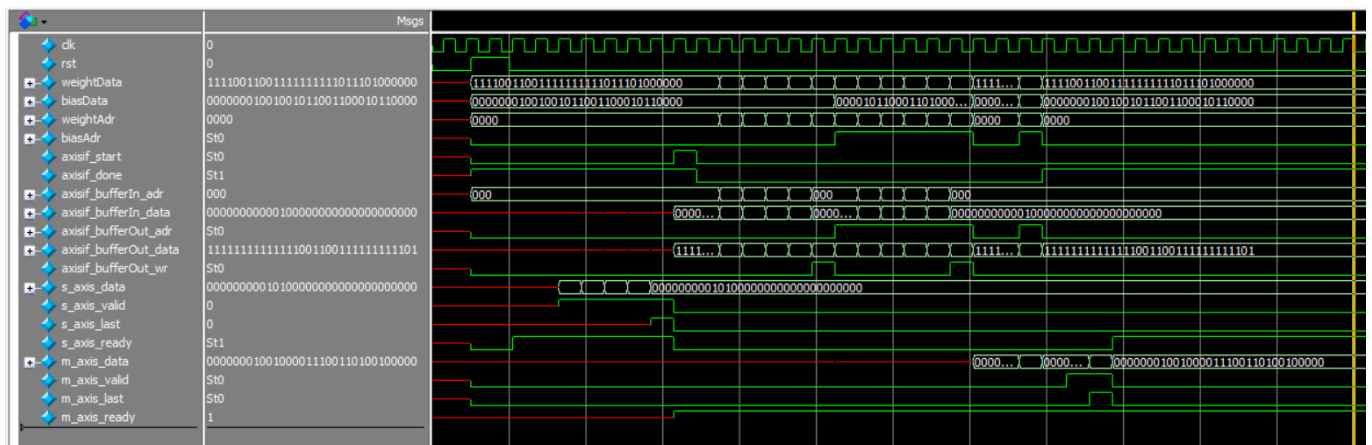
این ماژول به عنوان تست ۸ داده را با برقراری پروتکل AXI-4 Stream گرفته و ۴ داده با همان پروتکل خروجی می‌دهد؛ علاوه بر آن در میان فازهای slave و master، چند سیکل را به محاسبه اختصاص می‌دهد تا رفتار ماژولی که از آن استفاده می‌کند را شبیه‌سازی کند.



شکل ۴.۶.۱- شبیه‌سازی و شکل‌موج ماژول GP AXIS Interface

۴.۷- شبیه‌سازی ماژول Dense یا Fully Connected

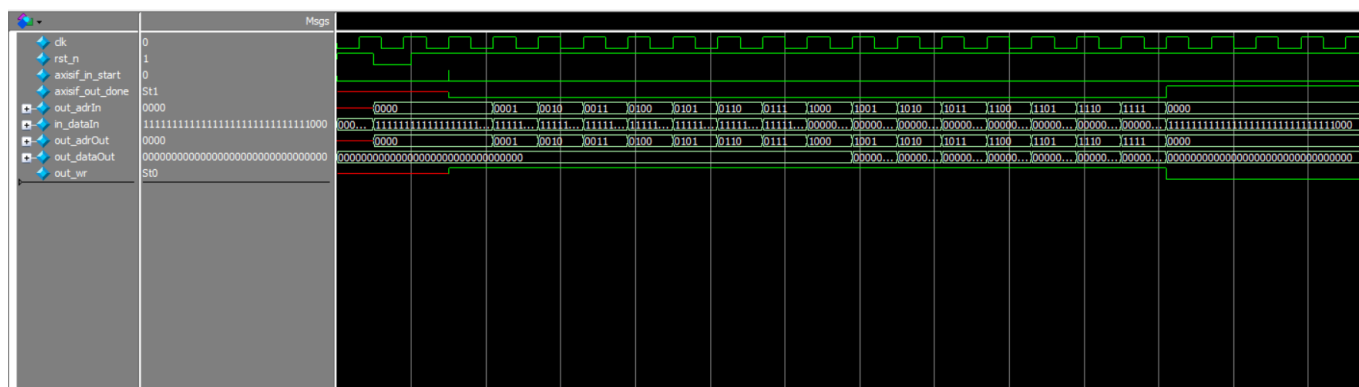
با توجه به تعداد زیاد ورودی برای ماژول Dense -تعداد ۱۶۰۰ تا داده ورودی- زمان simulation نرم‌افزار Modelsim بسیار طولانی شد و باید از یک LUT فیک استفاده کنیم و تعداد ورودی را برابر ۵ و خروجی را برابر ۲ در نظر می‌گیریم. با استفاده از این موضوع ماژول Dense را درستی سنجی می‌کنیم و در شکل ۴.۷.۱ مشاهده می‌شود که ابتدا ۵ داده ورودی از پورت s_axis_data وارد شده و پس از انجام محاسبات لازم آن را در بافر خروجی می‌نویسد و پس از اتمام تمامی محاسبات بر روی پورت m_axis_data برای خروجی ارسال می‌شود.



شکل ۴.۷.۱- شبیه‌سازی و شکل‌موج ماژول Dense

۴.۸- شبیه‌سازی ماژول ReLU

این ماژول به عنوان استفاده عمومی به عنوان activation function طراحی شده است و برای ارتباط با stage های بعدی، پروتکل AXI-4 Stream را پشتیبانی می‌کند. در شکل ۴.۸.۱ نحوه کار این ماژول مشاهده می‌شود.



شکل ۴.۸.۱- شبیه‌سازی و شکل‌موج ماژول ReLU

۴.۹- زمان‌بندی‌های نرم‌افزاری بر روی پردازنده PC

Layer	Name	Time (ms)
layer#0	Convolution 2D	396.3614
layer#0 activation	ReLU	7.0078
layer#1	MaxPooling 2D	6.0443
layer#2	Convolution 2D	1520.4479
layer#2 activation	ReLU	2.9997
layer#3	MaxPooling 2D	2.9995
layer#4	Flatten	0
layer#6	Dense	5.0367
layer#6 activation	Softmax	0

جدول ۴.۹.۱- زمان‌بندی اجرای لایه‌های مختلف به صورت نرم‌افزاری روی PC

۴.۱۰- زمان‌بندی‌های نرم‌افزاری بر روی قسمت PS برد PYNQ

در جدول ۴.۱۰.۱ تمامی زمان‌بندی‌ها قابل مشاهده است. لازم به ذکر است که تمامی زمان‌های نوشته شده به ثانیه است.

Layer	Name	Time (s)
layer#0	Convolution 2D	14.1355
layer#0 activation	ReLU	0.1697
layer#1	MaxPooling 2D	0.1646
layer#2	Convolution 2D	36.4254
layer#2 activation	ReLU	0.0842
layer#3	MaxPooling 2D	0.0595
layer#4	Flatten	0.0011
layer#6	Dense	0.0888
layer#6 activation	Softmax	0.0001

جدول ۴.۱۰.۱- زمان‌بندی اجرای لایه‌های مختلف به صورت نرم‌افزاری

۴.۱۱- زمان‌بندی‌های سخت‌افزاری

زمان انتقال داده‌ها از Memory واسطه DMA به یک ماژول AXIS برابر با 0.001171 ثانیه معادل 1.171ms می‌باشد و زمان دریافت داده‌ها نیز برابر با 0.001438 ثانیه معادل 1.438ms می‌باشد. زمان اول برای خواندن از حافظه و زمان دوم برای نوشتن در حافظه است که به همین دلیل بیشتر طول می‌کشد. لازم به ذکر است که فرکانس کاری مدار 100MHz می‌باشد.

Layer	Name	Dimension	Time (ms)
layer#0	Convolution	in: 784, out: 194688	1.95472
layer#0 activation	ReLU	in: 21632, out: 21632	0.43264
layer#1	MaxPooling	in: 21632, out: 21632	0.43264
layer#2	Convolution	in: 5408, out: 2230272	22.3568
layer#2 activation	ReLU	in: 7744, out: 7744	0.15488
layer#3	MaxPooling	in: 7744, out: 6400	0.14144
layer#4	Flatten	in: 1600, out: 1600	-
layer#6	Dense	in: 1600, out: 10	0.176
layer#6 activation	Softmax	in: 10: out: 10	TODO

جدول ۴.۱۱.۱- زمان‌بندی اجرای لایه‌های مختلف به صورت نرم‌افزاری

۴.۱۲- نتیجه زمان‌بندی‌ها و مقایسه قسمت PS و PL برد PYNQ

در قسمت ۴.۹ مشاهده شد که اجرا بر روی پردازنده PC خیلی سریع‌تر از پردازنده PYNQ بود. از طرفی اجرا کردن به صورت سخت‌افزاری روی PL بسیار سریع‌تر از PS بود و در جدول ۴.۱۲.۱ این زمان‌ها با هم مقایسه شده و speedup آن محاسبه شده است.

Layer	Name	PS Time (ms)	PL Time (ms)	Speed up
layer#0	Convolution 2D	14135.5	1.95472	723147%
layer#0 activation	ReLU	169.7	0.43264	39224%
layer#1	MaxPooling 2D	164.6	0.43264	38045%
layer#2	Convolution 2D	36425.4	22.3568	162927%
layer#2 activation	ReLU	84.2	0.15488	54364%
layer#3	MaxPooling 2D	59.5	0.14144	42067%
layer#4	Flatten	1.1	-	-
layer#6	Dense	88.8	0.176	50454%
layer#6 activation	Softmax	0.1	0.005	2000%

جدول ۴.۱۲.۱- مقایسه زمان‌بندی‌های PS و PL

مشاهده می‌شود که بیشترین speedup ها برای لایه‌های Convolution بوده است. در واقع لایه ۰ حدود ۷۲۰۰ برابر سریع‌تر و لایه ۲ حدود ۱۵۰۰ برابر سریع‌تر شده است. در مجموع کل زمان برای قسمت PS برد PYNQ برابر با 51.129s بوده و برای قسمت PL برابر با 25.65ms بوده است.

۴.۱۲- نتیجه زمان‌بندی‌ها و مقایسه قسمت PL برد PYNQ با پردازنده PC

با توجه به پایین بودن سرعت PS برد PYNQ نسبت به پردازنده PC، مقایسه بین پردازنده PC یک لپ‌تاپ با قسمت سخت‌افزاری قابل توجه می‌باشد.

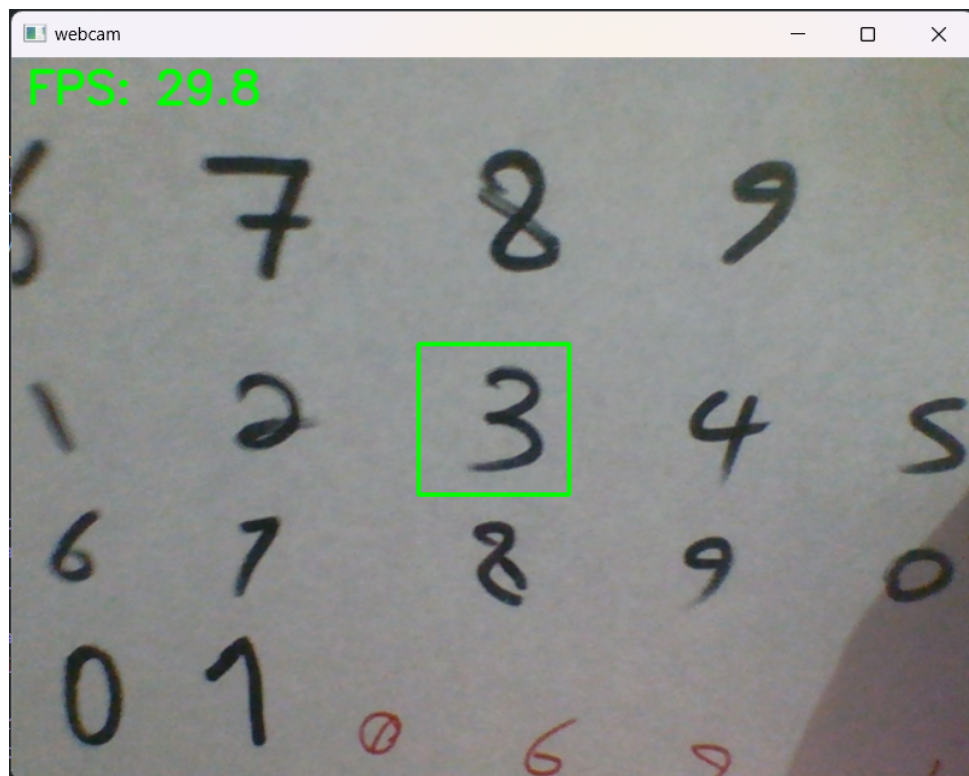
Layer	Name	PC Time (ms)	PL Time (ms)	Speed up
layer#0	Convolution 2D	396.3614	1.95472	20277%
layer#0 activation	ReLU	7.0078	0.43264	1620%
layer#1	MaxPooling 2D	6.0443	0.43264	1397%
layer#2	Convolution 2D	1520.4479	22.3568	6800%
layer#2 activation	ReLU	2.9997	0.15488	1937%
layer#3	MaxPooling 2D	2.9995	0.14144	2121%
layer#4	Flatten	0	-	-
layer#6	Dense	5.0367	0.176	2862%
layer#6 activation	Softmax	0	0.005	-

جدول ۴.۱۳.۱- مقایسه زمان‌بندی‌های PL و پردازنده PC

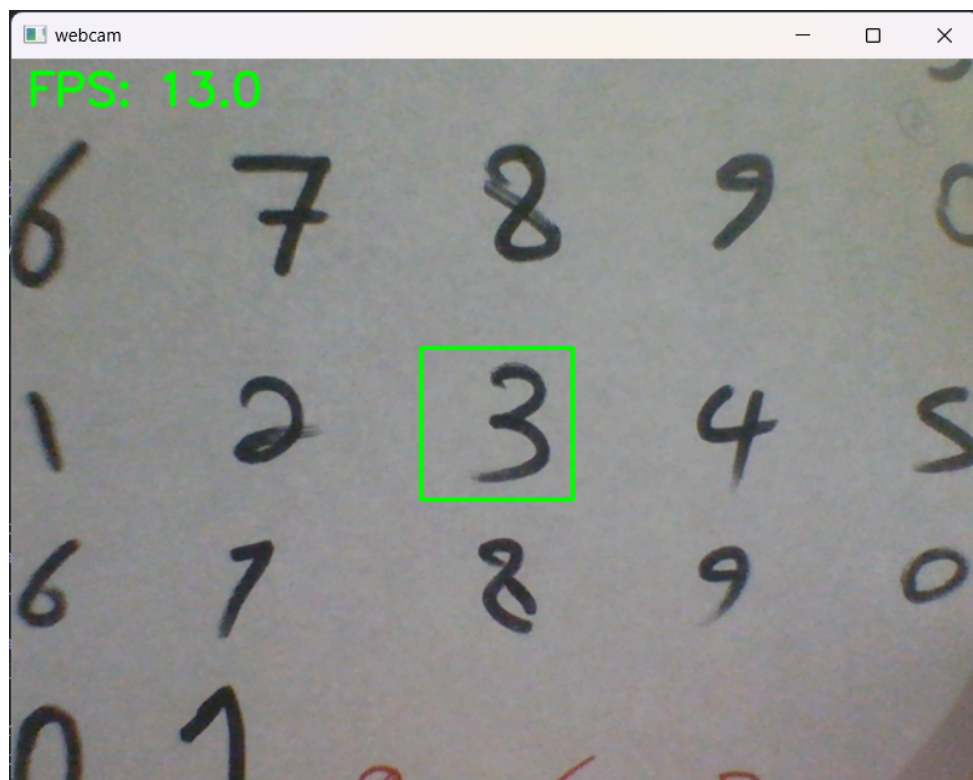
در مجموع مشاهده می‌شود که زمان کلی اجرای PC برابر با 1940.8973ms می‌باشد. بیشترین speed up برای لایه‌های Convolution بوده است که در لایه ۰ حدود ۲۰۰ برابر و در لایه ۲ حدود ۷۰ برابر سریعتر شده است.

قابل ذکر است که این تست‌ها صرفاً برای یک نمونه عکس می‌باشد و در صورتی که به صورت realtime کار کنیم، نتایج باز هم به نفع سخت‌افزار باز خواهد گشت زیرا در آنجا قابلیت pipeline هم وجود داشته و می‌توان به صورت همزمان، هنگامی که لایه بعدی در حال کار کردن است لایه قبلی نیز شروع به کار کند اما در نرم‌افزار باید تمام لایه‌ها اجرا شده و سپس محاسبات بر روی فریم بعدی صورت بگیرد.

۴.۱۳- زمان بندی اجرا Real Time بر روی پردازنده PC



شکل ۴.۱۳.۱- پردازش Real Time بر روی تصاویر با استفاده از کتابخانه keras و موازی سازی



شکل ۴.۱۳.۲- پردازش Real Time بر روی تصاویر با استفاده از حلقه های for بدون استفاده از کتابخانه و موازی سازی

همانطور که در شکل‌های ۴.۱۳.۱ و ۴.۱۳.۲ مشاهده می‌شود، میزان FPS با استفاده از کتابخانه keras و موازی‌سازی‌هایی که در این کتابخانه در نظر گرفته شده است از ۳۰ فریم بر ثانیه به ۱۳ فریم بر ثانیه هنگامی که برای محاسبات لایه‌ها از هیچ کتابخانه و موازی‌سازی‌ای استفاده نمی‌کنیم و صرفاً با استفاده از چندین حلقه for تودرتو و ضرایب این محاسبات را انجام می‌دهیم، کاهش می‌یابد. یعنی موازی‌سازی در پردازنده PC بین ۲ تا ۳ برابر می‌تواند کمک کند.

در این پروژه دیدیم که پیاده‌سازی سخت‌افزاری به جای نرم‌افزاری یک شبکه عصبی CNN چقدر می‌تواند مفید باشد و همچنین یک Design Space Exploration درباره زمانبندی‌های مختلف برای لایه‌های مختلف داشتیم. حال با استفاده از این داده‌ها می‌توان یافت که چه لایه‌هایی Bottleneck طراحی ما بوده و بهتر است آن را به صورت سخت‌افزاری پیاده کرد و چه لایه‌ای را می‌توان به صورت نرم‌افزاری پیاده‌سازی کرد تا از پیچیدگی زیاد هم کاست تا در نهایت به یک دیزاین efficient هم از لحاظ زمانی، هم از لحاظ Resource استفاده شده و هم از لحاظ پیچیدگی، رسید. به طور خاص در این پروژه با این معماری CNN دیدیم که یکی از Bottleneck های اصلی سیستم، هر دو لایه Convolution بود و زمان بسیار زیادی صرف محاسبات می‌کرد. در آینده می‌توان علاوه بر قسمت Evaluation، قسمت Train کردن را نیز به صورت سخت‌افزاری پیاده‌سازی کرد؛ زیرا یک بخش عظیمی از زمان در train کردن سپری می‌شود و می‌توان از این موضوع به خوبی استفاده کرد و زمان train کردن را فوق‌العاده کاهش داد.

- [1] [PYNQ](#)
- [2] [Zynq](#)
- [3] [XUP PYNQ-Z2](#)
- [4] [PYNQ Introduction – Python productivity for Zynq \(Pynq\) v1.0](#)
- [5] [Verilog AXI stream components for FPGA implementation](#)
- [6] [AXI4-Lite Interface Wrapper for Custom RTL in Vivado 2021.2 - Hackster.io](#)
- [7] [AXI4-Stream Interfaces](#)