

ابزار تولید و بهینه‌سازی شبکه‌های یادگیری عمیق برای پیاده‌سازی روی FPGA ها

نیما شیرین‌زاده^۱، مهدی امینیان^۲

۱- دانشجوی کارشناسی ارشد، گروه مهندسی کامپیوتر، دانشکده فنی و مهندسی، دانشگاه گیلان، رشت

nima.shirinzadeh@gmail.com

۲- استادیار، گروه مهندسی کامپیوتر، دانشکده فنی و مهندسی، دانشگاه گیلان، رشت

mahdi.aminian@guilan.ac.ir

خلاصه

یکی از انواع مهم شبکه‌های عصبی، شبکه‌های عصبی بازگشتی هستند. به طور کلی از این شبکه‌ها می‌توان برای پیدا کردن الگوها در داده‌هایی که توالی در آن‌ها وجود دارد، استفاده کرد. از جمله کاربردهای شبکه‌های عصبی بازگشتی می‌توان به تشخیص صوت، ترجمه‌گرها، تشخیص دست‌خط، تشخیص حالت‌های انسانی، کنترل ربات‌ها، تولید موسیقی و کارهای مختلف دیگر اشاره کرد. اما یکی از چالش‌های مهم برای این نوع الگوریتم‌ها پیاده‌سازی آن‌ها روی سخت‌افزار است. پیاده‌سازی این الگوریتم‌ها روی CPU و GPU و FPGA امکان پذیر است. پیاده‌سازی با استفاده از CPU و GPU به دلیل سهولت پیاده‌سازی و وجود راهکارهای نرم‌افزاری آسان، روش‌های پیاده‌سازی شناخته شده‌تر هستند. از طرفی، امروزه FPGAها به دلیل مصرف توان پایین و امکان طراحی بهینه و مناسب هر الگوریتم و امکان پیاده‌سازی موازی می‌تواند در بخش‌های مختلف جایگزین مناسبی برای CPU و GPU باشد. با توجه به ویژگی‌های ذکر شده FPGA می‌تواند ابزار مناسبی برای پیاده‌سازی شبکه‌های عصبی نیز باشد. با توجه به ماهیت پیچیده الگوریتم‌های شبکه‌های عصبی، پیاده‌سازی این نوع الگوریتم‌ها روی FPGA کاری زمان‌بر است و از طرفی برای پیاده‌سازی بهینه و مناسب الگوریتم‌های شبکه‌های عصبی به تجربه کافی در زمینه پیاده‌سازی سخت‌افزاری الگوریتم‌ها و همچنین تسلط کافی به ماژول‌های تشکیل دهنده یک شبکه‌ی عصبی نیاز است. اما برای رفع این مشکلات و امکان پیاده‌سازی سریع و بهینه شبکه‌های عصبی روی FPGA می‌توان با ترکیب راهکارهای سخت‌افزاری و نرم‌افزاری نظیر روش‌های بهینه‌سازی مختلف ماژول‌های تشکیل دهنده شبکه‌های عصبی این امکان را فراهم نمود تا شبکه‌ی عصبی پیاده شده در پلتفرم‌های نرم‌افزاری مانند Keras و Pytorch را به شبکه‌های سخت‌افزاری بهینه‌سازی شده و قابل پیاده‌سازی روی FPGA تبدیل کرد. یکی از خروجی‌های سخت‌افزاری ممکن، خروجی در زبان‌های سطح بالا مانند C یا C++ برای پیاده‌سازی از طریق HLS^۱ها و سنتز سطح بالا است. که به این ترتیب می‌توان فرایند زمان‌بر پیاده‌سازی و بهینه‌سازی سخت‌افزاری یک شبکه را در کوتاه‌ترین زمان ممکن طی نمود و همچنین این امکان فراهم می‌شود تا استفاده‌کنندگان حوزه‌های نرم‌افزاری هم بدون نیاز به دانش سخت‌افزاری بالا بتوانند به راحتی یک شبکه‌ی عصبی روی سخت‌افزار اجرا کنند.

کلمات کلیدی: شبکه‌های عصبی، سنتز سطح بالا، شتاب‌دهنده‌های سخت‌افزاری، بهینه‌سازی، LSTM^۱, Framework, HLS

۱. مقدمه

در گذر زمان و با پیشرفت تکنولوژی‌های سخت‌افزاری، در زمینه‌هایی مانند افزایش قدرت پردازشی و افزایش فضای ذخیره‌سازی، باعث شد تا الگوریتم‌های هوش مصنوعی^۲ مورد توجه قرار بگیرند و امکان استفاده‌ی عمومی آن‌ها فراهم شود، تا با استفاده از این الگوریتم‌ها امکان حل مسائلی که برای انسان غیرممکن بود فراهم شود. از جمله الگوریتم‌های مهم هوش مصنوعی و یادگیری ماشین، الگوریتم‌های یادگیری عمیق^۳ و شبکه‌های عصبی^۴ هستند که امروزه نقش مهمی را در بین الگوریتم‌های هوش مصنوعی بازی می‌کنند. یکی از نوآوری‌های مهم شبکه‌های عصبی، شبکه‌های عصبی بازگشتی^۵ هستند. این نوع شبکه‌ها به گونه‌ای طراحی شده‌اند که از فرایند یادگیری نورون‌های مغز انسان تقلید می‌کنند. به همین علت می‌توانند روی داده‌های توالی‌پذیر مانند صوت یا متن پردازش انجام دهند و مانند یک انسان به یادگیری و پیدا کردن الگوهای موجود در داده‌ها و پیش‌بینی بپردازند. از مهمترین کاربردهای شبکه‌های عصبی بازگشتی، پردازش زبان‌های طبیعی [1]، سیستم تشخیص گفتار [2]، سیستم تشخیص صوت [6] [3]، و سیستم تشخیص دست‌خط [4] است. امروزه از شبکه‌های عصبی بازگشتی ساده به دلیل حافظه بسیار کوتاه آن که باعث از بین رفتن گرادینان می‌شود استفاده نمی‌شود. [5] دو نوع مهم شبکه‌های عصبی که امروزه به صورت گسترده مورد استفاده قرار می‌گیرند شبکه‌های LSTM و GRU^۶ هستند. این شبکه‌های عصبی مشکل شبکه‌های عصبی ساده‌تر را ندارند و ارتباط بین داده‌های متوالی را حفظ می‌کنند.

برای پیاده‌سازی این شبکه‌های عصبی از GPU ها استفاده می‌شود. GPU ها با توجه به امکان اجرای موازی، بستر خوبی برای اجرای الگوریتم‌های شبکه‌های عصبی هستند. اما برای اجرای این الگوریتم‌ها می‌توان از FPGA ها نیز استفاده کرد. FPGA ها هم مانند GPU ها امکان اجرای موازی الگوریتم‌ها را دارند. و با توجه به اینکه FPGA ها مصرف توان کمتری دارند و این امکان را می‌دهند تا برای هر الگوریتم معماری مناسب و بهینه طراحی کنیم، آن‌ها را گزینه مناسب‌تری برای پیاده‌سازی شبکه‌های عصبی بازگشتی می‌کند. ولی به علت اینکه اجرای الگوریتم‌های پیچیده مانند الگوریتم‌های هوش مصنوعی و بهینه‌سازی آن‌ها کاری زمان‌بر و نیاز به تجربه کافی در زمینه پیاده‌سازی سخت‌افزاری الگوریتم‌ها دارد، عموماً برای پیاده‌سازی الگوریتم‌های شبکه‌های عصبی از بسترهای نرم‌افزاری موجود و GPU استفاده می‌شود.

در این مقاله، قصد داریم تا با ترکیب راهکارهای سخت‌افزاری و نرم‌افزاری پیاده‌سازی شبکه‌های عصبی این امکان را فراهم کنیم تا بتوان شبکه‌ی عصبی تولید شده در سطح نرم‌افزار را به شبکه عصبی بهینه‌سازی شده و قابل پیاده‌سازی روی سخت‌افزار تبدیل کنیم. بدین ترتیب علاوه بر تولید سریع شبکه، شبکه بهینه و مناسب برای پیاده‌سازی سخت‌افزاری

^۱ Long short-term Memory

^۲ Artificial Intelligence

^۳ Deep Learning

^۴ Neural Network

^۵ Recurrent Neural Network

^۶ Neuron

^۷ Gated Recurrent Units

تولید می‌گردد. همچنین از طریق این راهکار، این امکان فراهم می‌شود تا کاربران نرم‌افزاری هم بتوانند شبکه‌ی عصبی خود را با کمترین دانش سخت‌افزاری روی FPGA پیاده‌سازی کنند.

در ادامه این مقاله، بخش ۲ به بررسی کارهای پیشین در زمینه‌های مرتبط می‌پردازد. در بخش ۳، به مفاهیم اولیه‌ی شبکه‌های عصبی و سخت‌افزار و روش‌های به‌کار گرفته‌شده در پیاده‌سازی آن پرداخته می‌شود. در بخش ۴، چالش‌های پیش رو در این زمینه و ایده‌های حل مساله مطرح می‌شود. سپس در بخش ۵ روش پیشنهادی و چگونگی اجرای آن توضیح داده شده است؛ و در بخش ۶، ارزیابی و نتایج آن ارائه شده است. در نهایت، در بخش ۷، نتیجه‌گیری حاصل از اجرای روش پیشنهادی بیان می‌شود.

۲. کارهای پیشین

امروزه پیاده‌سازی GRU/LSTM روی FPGA یکی از موضوعات مناسب تحقیقاتی است. برای مثال در مقاله‌های [7][8] یک LSTM با استفاده از زبان C طراحی شده است که با استفاده از HLS^۹ به RTL^{۱۰} تبدیل شده است. در مقاله [9] و در مقاله‌ی [10]، LSTM در سطح RTL طراحی شده است و در پلتفرم Xilinx zynq پیاده‌سازی شده است. استفاده از این روش زمان بیشتری صرف نموده است ولی خروجی آن می‌تواند بهتر از خروجی HLS باشد. استفاده از HLS برای پیاده‌سازی شبکه‌های عصبی یک انتخاب مناسب است ولی استفاده از این ابزار به تنهایی نمی‌تواند فرایند پیاده‌سازی شبکه‌های عصبی روی FPGA را کامل کند. یک شبکه‌ی عصبی LSTM/GRU از المان‌های مختلفی مانند تابع‌های فعال‌ساز^{۱۱} و ضرب ماتریسی تشکیل شده است. برای پیاده‌سازی بهینه شبکه‌های عصبی روی سخت‌افزار باید این المان‌ها به گونه‌ای پیاده‌سازی شوند تا مناسب پیاده‌سازی روی سخت‌افزار باشد. در شبکه‌های عصبی تابع‌های فعال‌ساز Sigmoid، Tahn و Softmax تابع‌هایی هستند که بیشترین استفاده را در شبکه‌های عصبی بازگشتی دارند. اما این توابع به دلیل وجود محاسبات توانی^{۱۲} و تقسیم اعداد شناور^{۱۳} مقدار زیادی از منابع سخت‌افزاری را مصرف می‌کنند. برای پیاده‌سازی تابع‌های فعال‌ساز باید از پیاده‌سازی تقریبی این المان‌ها روی سخت‌افزار استفاده کرد تا با مصرف کمتر به خروجی مورد انتظار رسید. در مقاله [7] [11] از روش PLAN برای پیاده‌سازی Sigmoid استفاده شده است. در مقاله [12] برای پیاده‌سازی Tanh رابطه‌ای تقریبی ارائه شده است که برای محاسبه نتیجه از Sigmoid استفاده می‌کند. همچنین در مقاله [13] برای پیاده‌سازی یک لایه Softmax معماری جدیدی ارائه شده است که عملیات تقسیم اعداد شناور را از پیاده‌سازی حذف کرده است. در مقاله [14] برای بهینه‌سازی مصرف^{۱۳} BRAM از Tiling استفاده شده تا به این ترتیب این امکان فراهم شود تا شبکه‌ی بزرگ‌تری را بتوان روی FPGA پیاده‌سازی کرد.

در مقالات ذکر شده سعی شد با استفاده از روش‌های بهینه‌سازی مختلف، بخش‌های مختلف را به صورت کارآتر برای سخت‌افزار طراحی کنند. در روش پیشنهادی این مقاله با ترکیب روش‌های بهینه‌سازی موجود در قالب یک پلتفرم و

^۹ Register-transfer Level

^{۱۰} Activation Function

^{۱۱} Exponentiation

^{۱۲} Floating-point Divisions

^{۱۳} Block RAMs

علاوه بر آن، با بهره گیری از راهکارهای نرم افزاری سعی در تولید سریع و بهینه ی یک شبکه سخت افزاری از خروجی شبکه پیاده سازی شده از Keras با استفاده از Python و HLS را داریم.

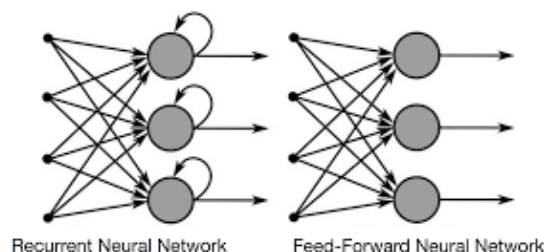
به طور خلاصه، این مقاله موارد زیر را ارائه می دهد:

- جامع بودن پلت فرم پیشنهادی نسبت به موارد موجود از نظر دسترسی به تابع های فعال ساز مختلف و نوع شبکه های قابل پیاده سازی متفاوت
- ایجاد شبکه ی عصبی از خروجی نرم افزاری و صرفه جویی در زمان
- بهینه سازی تابع های فعال ساز و مساحت و تاخیر آنها
- بهینه سازی مصرف BRAM و مساحت آنها

۳. مفاهیم اولیه

۳.۱. شبکه های عصبی بازگشتی

شبکه ی عصبی بازگشتی یک شبکه ی عصبی عمیق است. این نوع شبکه ها به گونه ای طراحی شده اند که می توانند با داده هایی که با هم ارتباط دارند کار کنند. بر خلاف شبکه های Feed-Forward [15] در این شبکه ها خروجی تولید شده در زمان قبلی را به عنوان ورودی در زمان حال استفاده می کنند. که بدین ترتیب می تواند توالی و ارتباط داده ها در یک مجموعه داده حفظ شود. در شکل ۱ ساختار شبکه ی RNN^5 و Feed-Forward نشان داده شده است.



شکل ۱- ساختار شبکه ی عصبی RNN و Feed-Forward

رابطه های زیر نحوه ی عملکرد و محاسبات یک شبکه ی عصبی ساده را نشان می دهد. که در آن X بردار ورودی، H بردار خروجی در زمان قبلی، و Y بردار خروجی زمان فعلی است.

$$h_t = \sigma(w_{xh} + w_{hh} + b_h) \quad (1)$$

$$y_t = w_{hy} + b_y \quad (2)$$

w_{xh} ماتریس وزنی است که لایه ی خروجی را به لایه ی نهان مرتبط می کند، w_{hy} ماتریس وزنی است که لایه ی نهان را به لایه ی خروجی مرتبط می کند، و w_{hh} ماتریس وزن بازگشتی است که ورودی لایه ی نهان در زمان قبلی را به لایه ی

^{۱۴} Area

^{۱۵} Recurrent Neural Network

نهان در زمان فعلی مرتبط می‌کند. بردارهای b_h و b_y بردارهای بایاس هستند. همچنین، σ بیان‌گر یک تابع فعال‌ساز است که می‌تواند Sigmoid یا Tanh باشد.

۲.۳. شبکه‌های عصبی بازگشتی

در شبکه‌های عصبی ساده تنها این امکان وجود دارد که از ورودی زمان قبلی یادگیری انجام دهد و این امکان وجود ندارد تا ورودی زمان‌های دورتر نیز در تصمیم‌گیری در داده‌های مرتبط در نظر گرفته شود. همچنین در شبکه‌های عصبی ساده مشکل Vanishing Gradient نیز وجود دارد. این حالت ممکن است در هنگام آموزش یک شبکه با چند لایه رخ دهد و باعث شود که تغییر در وزن‌ها هنگام آموزش به سختی انجام شود. برای حل این مشکلات، معماری جدیدی برای شبکه‌های عصبی به نام LSTM [5] ارائه شد. در شکل ۲ معماری یک سلول LSTM نشان داده شده است. در زمان t خروجی‌های زیر از گیت‌های LSTM تولید می‌شود:

$$\dot{i}_t = \sigma(w_{ix}x_t + w_{iy}y_{t-1} + w_{ic}c_{t-1} + b_i) \quad (3)$$

$$f_t = \sigma(w_{fx}x_t + w_{fy}y_{t-1} + w_{fc}c_{t-1} + b_f) \quad (4)$$

$$g_t = \sigma(w_{gx}x_t + w_{gy}y_{t-1} + w_{gc}c_{t-1} + b_g) \quad (5)$$

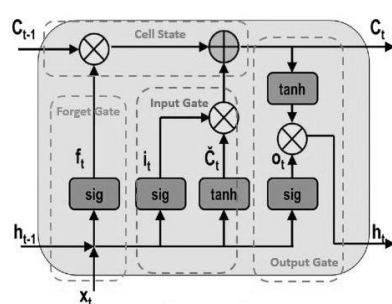
$$c_t = f_t \circ x_t + g_t \circ i_t \quad (6)$$

$$o_t = \sigma(w_{ox}x_t + w_{oy}y_{t-1} + w_{oc}c_{t-1} + b_o) \quad (7)$$

$$m_t = o_t \circ h(c_t) \quad (8)$$

$$y_t = w_{ym}m_t \quad (9)$$

در روابط (3) تا (9)، W و U ماتریس وزن را نشان می‌دهند و b نشان‌دهنده بردارهای بایاس هستند. پارامترهای I, f, O و c به ترتیب گیت‌های ورودی، گیت فراموشی، گیت خروجی و سلول فعالیت حافظه‌ی طولانی مدت را نشان می‌دهد. عملگر (نقطه تو خالی) ضرب عنصر به عنصر را نشان می‌دهد و عملگر $+$ جمع عنصر را نشان می‌دهد. علامت سیگما تابع Sigmoid را نشان می‌دهد و h تابع Tanh را نشان می‌دهد. هدف از تابع Sigmoid و Tanh نرمال‌سازی ورودی‌ها به منظور جلوگیری از تاثیر بیش از حد برخی ورودی‌ها نسبت به دیگر ورودی‌ها است. در میان تمام عملگرها ضرب ماتریس بردار بیشترین عملگر به کار رفته در حافظه و محاسبات است. ابعاد به کار رفته در x_t ، y_t و c_t یکسان هستند.



شکل ۲- معماری یک سلول LSTM

۳.۳. شبکه عصبی GRU

نوع دیگری از معماری شبکه‌های عصبی، GRU ها هستند. معماری سلول GRU مشابه معماری سلول LSTM است با این تفاوت که تعداد پارامترهای یک سلول GRU کمتر از یک سلول LSTM است. عملکرد GRU در برخی از کارها مانند مدل‌سازی موسیقی و پردازش زبان‌های طبیعی مانند عملکرد LSTM است. از این رو در این نوع کارها استفاده از GRU کمک به کاهش حجم محاسبات و فضای مورد استفاده می‌کند. [16] معماری یک سلول GRU در شکل ۳ نشان داده شده است. در زمان t خروجی‌های زیر از گیت‌های GRU تولید می‌شود:

$$z_t = \sigma(w_z x_t + u_z h_{t-1} + b_z) \quad (10)$$

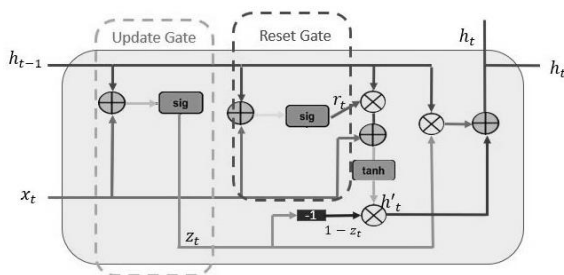
$$r_t = \sigma(w_r x_t + u_r h_{t-1} + b_r) \quad (11)$$

$$h'_t = \sigma(w_h x_t + u_h (r_t \circ h_{t-1}) + b_h) \quad (12)$$

$$h_t = \sigma((1 - z_t) \circ h_{t-1} + z_t \circ h'_t) \quad (13)$$

$$y_t = h_t \quad (14)$$

همانطور که مشخص است یک سلول GRU به LSTM شباهت دارد اما از تعداد گیت‌های کمتری تشکیل شده است. Γ و Z به ترتیب گیت به روزرسانی و گیت تنظیم مجدد را نشان می‌دهد.



شکل ۳- معماری یک سلول GRU

۴.۳. شتاب دهنده‌های سخت افزاری

یکی از بسترهای مناسب برای پیاده‌سازی شبکه‌های عصبی GPU است. GPU به دلیل داشتن تعداد هسته‌های زیاد می‌تواند قابلیت اجرای موازی داشته باشد به همین منظور یک گزینه‌ی خوب برای پیاده‌سازی شبکه‌های عصبی است. اما نوع دیگری از سخت‌افزار که مناسب پیاده‌سازی شبکه‌های عصبی است FPGAها هستند. این نوع سخت‌افزار هم به دلیل امکان اجرای موازی گزینه‌ی مناسبی برای پیاده‌سازی شبکه‌های عصبی است.

اما FPGA تفاوت‌هایی با GPU دارد که آن را به گزینه‌ی مناسب‌تری برای پیاده‌سازی شبکه‌های عصبی تبدیل می‌کند. FPGAها با فرکانس پایین‌تری کار می‌کنند، از این رو مصرف توان کمتر را با عملکرد یکسان بخاطر قابلیت

موازی سازی بالاتر دارند. همچنین به دلیل وجود بلوک های داخلی حافظه (BRAM) می توان به بیشترین حد موازی سازی در FPGA ها رسید. نکته ی مهم دیگر انعطاف پذیر بودن FPGA ها است که باعث می شود بتوان برای هر الگوریتمی، مداری بهینه طراحی کرد. به بیان دیگر امکان تغییر سریع معماری، یا حتی تغییر بخشی از آن در حین کار باقی بخش ها، با توجه به کاربرد وجود دارد تا الگوریتم ها به کارایی بالاتری برسند. [17] اما یکی از معایب پیاده سازی یک شبکه ی عصبی روی FPGA زمان بر بودن پیاده سازی الگوریتم است. همچنین پیاده سازی بهینه الگوریتم نیاز به تجربه کافی در پیاده سازی سخت افزاری الگوریتم ها مورد نیاز دارد و به دلیل اینکه این کار به صورت دستی انجام می شود باعث افزایش زمان پیاده سازی و تست یک شبکه روی FPGA می شود.

۳.۵. روش های پیاده سازی سخت افزاری

یکی از راه های پیاده سازی الگوریتم ها روی سخت افزار استفاده از HDL ها^{۱۶} (زبان های توصیف سخت افزار) است. اما برای پیاده سازی با این روش بسیاری از جزئیات پیاده سازی نیز باید توسط طراح پیاده سازی شود که این زمان به زمان پیاده سازی الگوریتم اصلی و بهینه سازی اضافه می شود. روش دیگر پیاده سازی استفاده از HLS ها (زبان های سنتز سطح بالای مدارهای دیجیتال) است که از زبان های سطح بالا مانند C یا C++ برای پیاده سازی استفاده می کنند. مزیت پیاده سازی با HLS کاهش زمان پیاده سازی و تست است. با استفاده از HLS حدود ۲ برابر سریع تر می توان الگوریتم ها را پیاده سازی کرد. اما مشکل استفاده از HLS این است که کارایی آن نسبت به کارایی خروجی HDL کمتر است، اما با این حال می تواند تا حد قابل قبولی به کارایی روش دستی HDL در زمان بسیار کوتاه تر نزدیک شود. بنابراین HLS در بسیاری از موارد می تواند گزینه مناسبی برای پیاده سازی بسیاری از الگوریتم ها روی FPGA باشد. در این مقاله نیز برای پیاده سازی شبکه های عصبی از HLS استفاده شده است. تا علاوه بر اینکه پیاده سازی و تست الگوریتم سریع تر انجام شود این امکان برای کاربران نرم افزاری نیز فراهم شود تا شبکه ی خود را روی سخت افزار پیاده سازی و تست کنند. [18]

۴. طرح مسئله

GPU و FPGA به دلیل دارا بودن قابلیت اجرای موازی از جمله بسترهای مناسب برای اجرای الگوریتم های شبکه های عصبی هستند. اما FPGA با توجه به انعطاف پذیری بالا و مصرف توان کمتر گزینه بهتری نسبت به GPU هاست. اما برای پیاده سازی شبکه ی عصبی روی GPU ها رابط های نرم افزاری متنوعی وجود دارد که این امکان را به کاربر می دهد تا در سریع ترین زمان و به آسانی بتواند یک شبکه را پیاده سازی کند و همچنین از قابلیت اجرای موازی آن استفاده کند. به همین دلایل GPU ها گزینه ی پر استفاده تری هستند. با وجود اینکه FPGA ها یکی از مناسب ترین ابزارها برای پیاده سازی شبکه ی عصبی هستند، اما پیاده سازی شبکه های عصبی روی FPGA ها با چالش هایی همراه است و به راحتی پیاده سازی روی GPU یا CPU ها نیست. یکی از

^{۱۶} Hardware Description Language

^{۱۷} Performance

مشکلات پیاده سازی شبکه های عصبی روی FPGAها نبود پلتفرم هایی مانند Keras یا Pytorch است تا بتوان سریعاً یک شبکه برای سخت افزار تولید کرد.

برای این منظور باید جزییات شبکه های عصبی را برای رسیدن به یک پلتفرم پیاده سازی سخت افزاری کاملاً بررسی نمود. یک شبکه ی عصبی از ماژول ها و بخش های مختلفی تشکیل شده است، و به دلیل این که ماژول های شبکه های عصبی منابع زیادی برای انجام محاسبات خود مصرف می کنند، امکان پیاده سازی یک شبکه ی عصبی بزرگ روی FPGA بصورت همزمان کار غیرممکنی است. بنابراین تنها تسلط روی جزییات شبکه برای پیاده سازی کافی نیست و ماژول ها باید بهینه سازی نیز شوند. برای بهینه کردن هر ماژول باید از روش های متفاوت و مختلفی با توجه به نوع ماژول استفاده کرد. بنابراین، یکی دیگر از مشکلات، علاوه بر زمان مورد نیاز برای پیاده سازی، پیچیدگی و زمان مورد نیاز بالا برای بهینه سازی یک شبکه ی عصبی برای اجرا روی FPGA است.

برای بهینه سازی شبکه های عصبی دو چالش مهم وجود دارد. یکی از این چالش ها، وجود عملیات هایی مانند توان رسانی و تقسیم اعداد شناور در توابع فعال سازی پر استفاده است که منابع محاسباتی زیاد از جمله^{۱۸} DSP،^{۱۹} LUT و^{۲۰} FFها را مصرف می کند. چالش مهم دیگر، مصرف بالای حافظه BRAM در شبکه های عصبی است. یک شبکه ی عصبی از پارامترهای فراوانی در قالب ماتریس تشکیل شده است که ذخیره سازی آن ها با استفاده از بلوک های BRAM انجام می شود. مقدار داده های ورودی یک شبکه ی عصبی بسیار بیشتر از تعداد بلوک های BRAM موجود روی FPGA است، بنابراین پیاده سازی اکثر شبکه ها روی FPGA امکان پذیر نیست.

همچنین، برای پیاده سازی یک شبکه روی FPGA نیاز به یک رابط نرم افزاری جامع وجود دارد تا این امکان را فراهم کند که بتوان در سریع ترین زمان به یک خروجی قابل پیاده سازی رسید که عملکرد بهینه ای داشته باشد. در نتیجه باید راهکاری ارائه کرد که پیاده سازی شبکه ی عصبی روی FPGA را آسان تر کند، تا استفاده از آن برای کسانی که تجربه کافی در زمینه پیاده سازی سخت افزاری الگوریتم ها را ندارند نیز امکان پذیر باشد. این مقاله قصد دارد با ارائه یک واسط سخت افزاری نرم افزاری، و ایجاد امکانات پیاده سازی سخت افزاری ماژول های مهم شبکه های عصبی، قابلیت تطبیق و تبدیل شبکه های مختلف را برای پیاده سازی روی FPGA ایجاد کند. این رابط، شبکه های توصیف شده در پلتفرم های نرم افزاری مانند KERAS را به کدهای سخت افزاری تبدیل می کند، علاوه بر بهینه سازی کد، امکان زمان بندی اجرای لایه های مختلف روی بستر سخت افزاری را فراهم می کند.

۵. روش پیشنهادی

در روش پیشنهادی هدف فراهم آوردن امکان پیاده سازی بهینه و سریع شبکه های عصبی شامل LSTM، GRU و DENSE با استفاده از شبکه ی آموزش دیده شده در Keras است. برای این منظور باید بخش های مهم این شبکه ها به صورتی پیاده سازی شوند که مناسب پیاده سازی روی سخت افزار باشند. در شبکه های عصبی توابع فعال ساز و ضرب های ماتریسی بخش هایی هستند که بیشترین منابع را مصرف می کنند. بنابراین در ابتدا راهکارهایی ارائه می دهیم تا مصرف منابع برای توابع فعال ساز مهم بهینه شود. سپس مصرف منابع فرآیند کلی یک شبکه ی عصبی و ضرب ماتریسی را بهینه می کنیم و در آخر راهکاری ارائه می دهیم تا شبکه ی تولید شده در ابزار نرم افزاری Kears را با استفاده از واسط

سخت افزاری نرم افزاری و ماژول های بهینه سازی شده ی شبکه های عصبی برای سخت افزار، به شبکه ی بهینه ی قابل پیاده سازی روی سخت افزار تبدیل کنیم.

۱.۵. پیاده سازی تابع های فعال ساز

یکی از ماژول های مهم شبکه های عصبی تابع های فعال ساز غیر خطی هستند. در اکثر این نوع ماژول ها عملیات توان رسانی و تقسیم اعداد شناور وجود دارد. پیاده سازی این دسته از عملیات روی سخت افزار منابع زیادی مصرف می کند و همچنین باعث افزایش تاخیر می شود. بنابراین، برای پیاده سازی این ماژول ها باید از روش های تقریبی استفاده کرد. در پیاده سازی با روش های تقریبی علاوه بر بهبود تاخیر، به اندازه ی قابل ملاحظه ای در منابع مصرفی صرفه جویی می شود. برای این پلتفرم توابع Sigmoid [11] و Tanh [12] در قالب یک ماژول تقریبی پیاده سازی شده است. همچنین Softmax به صورت یک لایه Softmax تقریبی ارائه شده است. از Relu به دلیل اینکه ماهیتا مصرف منابع کمی دارد بدون تغییر در پیاده سازی استفاده شده است.

۱.۱.۵. پیاده سازی Sigmoid

یکی از انواع پر کاربرد توابع فعال ساز، تابع فعال ساز Sigmoid است. رابطه (15)، تابع Sigmoid را نشان می دهد که شامل توان رسانی و تقسیم اعداد شناور است. برای پیاده سازی این ماژول از روش PLAN [11] استفاده شده است. در این روش عملیات های پرهزینه مانند توان رسانی و تقسیم اعداد اعشاری حذف شده است و این تابع با استفاده از μx ها و Shift Registerها پیاده سازی شده است.

در شکل ۴ پیاده سازی Sigmoid با روش PLAN نشان داده شده است. در این روش در بخش مثبت ۷ نقطه ی کلی روی نمودار انتخاب می شود. و بین هر دو نقطه یک معادله خط تشکیل می شود و در نهایت ۴ پاره خط بر اساس نقاط روی نمودار اصلی بدست می آیند که به صورت تقریبی می تواند به ازای هر ورودی، خروجی نزدیک به نمودار اصلی را تولید کند. نقاط به گونه ای انتخاب شده اند که معادلات خط بدست آمده تنها با استفاده از یک عمل شیفت و یک عمل جمع می تواند خروجی تقریبی تابع Sigmoid را تولید کند. تفاوت مصرف منابع و تاخیر با استفاده از روش PLAN نسبت به پیاده سازی اصلی در جدول ۲ نشان داده شده است.

همانطور که در شکل ۵ مشخص شده نمودار Sigmoid یک نمودار متقارن است و خط تقارن در نقطه ی $x=0$ قرار دارد. بنابراین کافی است محاسبات فقط برای بخش مثبت انجام شود. برای بدست آوردن مقدار خروجی برای ورودی های منفی کافی است از خروجی محاسبه شده با استفاده از $|x|$ ، یک واحد کم کنیم.

$$\text{sigmoid}(x) = \frac{1}{1+e^{-x}} \quad (15)$$

```
const Output sigmoid(Input input)
{
    Input y = 1;
    Input twoscomp = -input;
    bool signbit = input[0];
    if (signbit){
        input = twoscomp;
    }
    if((Input) 2.375 <= input && input < (Input) 5){
        y = (input >> 5) + Input(0.84375);
    }else if((Input) 1 <= input && input < (Input) 2.375){
        y = (input >> 3) + Input(0.625);
    }else if( (Input) 0 <= input && input < (Input) 1){
        y = (input >> 2) + Input(0.5);
    }

    if(signbit) {
        y = Input(1) - y;
    }

    return Output(y);
}
```

شکل ۴- پیاده سازی sigmoid با استفاده از روش PLAN در HLS

۵.۱.۲. پیاده سازی Tanh

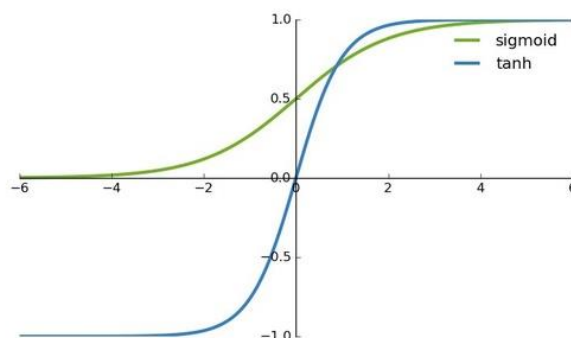
یکی دیگر از توابع فعال ساز مهم در شبکه های عصبی بازگشتی، تابع فعال ساز Tanh است. همانطور که در رابطه (16) مشخص است این تابع از ۴ عملیات توان رسانی اعداد شناور تشکیل شده است. بنابراین پیاده سازی این تابع نسبت به تابع Sigmoid در رابطه (15) از تاخیر بیشتری برخوردار است و مصرف منابع به مراتب بالاتری در حدود ۳ برابر بیشتر از Sigmoid را داراست.

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (16)$$

همانطور که در شکل ۵ نشان داده شده است نمودار دو تابع Sigmoid و Tanh مشابه یک دیگر است. [12] بنابراین با توجه به قوانین انتقال نمودارها می توان نقطه بدست آمده روی نمودار Sigmoid را به نقطه ای روی نمودار Tanh انتقال داد. با توجه به رابطه (17) می توان با استفاده از تابع Sigmoid پیاده سازی شده با روش PLAN و عملیات Shift، تابع فعال ساز Tanh را به صورت تقریبی پیاده سازی کرد.

در رابطه (17) زمانی که ورودی x دو برابر می شود نمودار Sigmoid جمع تر شده و در واقع شیب نمودار بیشتر شده و به نمودار Tanh نزدیک تر می شود. برای انتقال نمودار در جهت محور y به سمت پایین یک واحد از تابع Sigmoid کم شده است و برای این که بازه تابع در جهت محور y دو برابر شود نیاز است تا خروجی تابع Sigmoid دو برابر شود. به این ترتیب خروجی تقریبی تابع Tanh با استفاده اعمال شیفت و تفریق روی تابع Sigmoid قابل محاسبه است. همچنین در جدول ۳ تفاوت در مقدار منابع مصرفی و تاخیر تابع تقریبی پیاده سازی شده نسبت به پیاده سازی اولیه نشان داده شده است.

$$\tanh(x) = 2\sigma(2x) - 1 \quad (17)$$



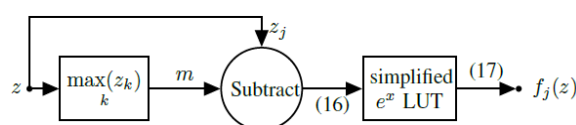
شکل ۵- مقایسه نمودار Sigmoid و Tanh

۳.۱.۵. پیاده سازی Softmax

برای این ابزار لایه Softmax در قالب یک لایه تقریبی [13] پیاده‌سازی شده است که برای عمل طبقه‌بندی^{۲۲} در انتهای شبکه مورد استفاده قرار می‌گیرد.

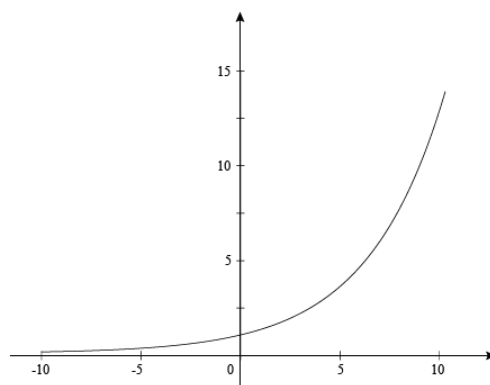
$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum e^{x_j}} \quad (18)$$

با در نظر گرفتن رابطه (18) برای بهبود منابع مصرفی این تابع از معماری ارائه شده در شکل ۶ استفاده شده است در این معماری عملیات تقسیم شناور حذف شده و تنها به یک عملیات توان رسانی اعداد شناور نیاز است. برای محاسبه Softmax در این روش ابتدا بیشترین مقدار در بین اعداد پیدا می‌شود. سپس، بعد از تفریق کردن بیشترین مقدار ورودی از اعداد ورودی، این اعداد در عملیات توان رسانی عدد e شرکت داده می‌شوند. باید در نظر داشت که معماری ارائه شده در صورتی که در لایه‌های میانی استفاده شود احتمال اینکه خطای شبکه را بالا ببرد وجود دارد. اما در لایه آخر، با توجه به اینکه مقدار خطا ناچیز است این معماری می‌تواند یکی از مناسب‌ترین پیاده‌سازی‌ها برای تابع Softmax باشد.



شکل ۶- معماری پیشنهادی یک لایه softmax [13]

وظیفه اصلی لایه Softmax دادن امتیاز بیشتر به ورودی بزرگتر است و معمولاً در لایه آخر برای تعیین نتیجه نهایی استفاده می‌شود. در معماری ارائه شده ابتدا در بین ورودی‌هایی که به لایه Softmax وارد می‌شوند بزرگترین ورودی انتخاب می‌شود. سپس ورودی بزرگ از تمام ورودی‌ها کم می‌شود. اعداد به دست آمده بعد از عمل تفریق در بازه بین صفر تا منفی بی‌نهایت قرار می‌گیرند. با اعمال یک تابع \exp روی این ورودی‌ها می‌توان خروجی نهایی را به دست آورد. با توجه به شکل ۷ ورودی‌ها بعد از عمل تفریق در سمت چپ محور y قرار می‌گیرند. بنابراین خروجی‌ها در بازه ۰ تا ۱ قرار دارند و بزرگترین مقدار ممکن عدد ۱ به ازای ورودی صفر است. در این معماری این حالت زمانی اتفاق می‌افتد که بزرگترین ورودی از خودش کم می‌شود. به این ترتیب به ورودی‌ها از بزرگ به کوچک امتیاز اختصاص می‌یابد و ورودی بزرگتر تعیین می‌شود.



شکل ۷- نمودار یک تابع \exp

۵.۲. روش Tiling

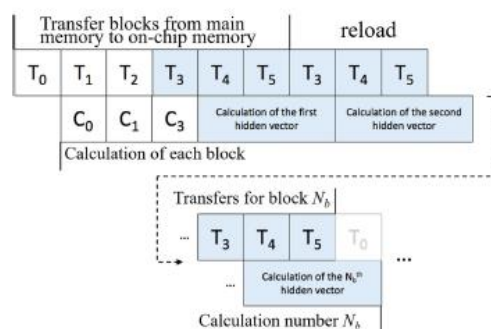
یکی از منابع مهم برای پیاده‌سازی شبکه‌های عصبی روی سخت‌افزار بلوک‌های BRAM هستند. از این منابع نسبت به LUTها، FFها تعداد کمتری وجود دارد و همین‌طور فضای ذخیره‌سازی این حافظه نیز نسبت به حافظه‌های خارجی کمتر است. همان‌طور که در بخش ۳ دیدیم یک لایه LSTM یا GRU از تعداد زیاد ماتریس وزن، بایاس و ورودی تشکیل شده است. بنابراین اگر تعداد نوروں‌های یک لایه یا تعداد ویژگی‌های ماتریس ورودی زیاد شود اندازه ماتریس وزن و بایاس نیز تغییر می‌کند. که در این حالت به دلیل تعداد کم و فضای ذخیره‌سازی محدود BRAMها امکان پیاده‌سازی شبکه‌های بزرگ روی FPGA وجود ندارد. بنابراین یکی از چالش‌های مهم، بهینه‌سازی مصرف BRAM در پیاده‌سازی سخت‌افزاری شبکه است.

یکی از روش‌های موجود و پرکاربرد برای بهبود مصرف BRAM، استفاده از روش Tiling [14] در کنار روش Streaming ورودی‌ها از حافظه خارجی و Pipelining عملیات خواندن از حافظه خارجی و عملیات محاسباتی است. با ترکیب این روش‌ها می‌توان وزن‌ها و بایاس‌ها را روی حافظه خارجی DDR ذخیره‌سازی کرد و برای انتقال سریع آن‌ها به مدار از AXIS^{۲۳} استفاده کرد. اما به دلیل اینکه خواندن از حافظه خارجی به تاخیر عملیات محاسباتی اضافه می‌کند، همان‌طور که در شکل ۸ نشان داده شده است، بخشی که شامل خواندن ورودی از حافظه خارجی و سپس اجرای عملیات

^{۲۳}AXI4-Stream interface

اصلی است به صورت Pipeline پیاده‌سازی می‌شود. در این روش، در حین اجرای عملیات اصلی با ورودی‌های خوانده شده در t ، ورودی‌ها برای زمان $t+1$ نیز خوانده می‌شود تا تاخیر خواندن و زمان اجرای عملیات همپوشانی پیدا کرده و از ایجاد تاخیر بخاطر خواندن از حافظه خارجی جلوگیری شود.

با استفاده از این روش تنها بخشی از ورودی‌ها که برای محاسبات است خوانده می‌شود. این ورودی از یک حافظه خارجی انتقال می‌یابد. همچنین این امکان وجود دارد تا تعداد ورودی‌های لازم برای هر خواندن نیز تعیین شود تا به این ترتیب علاوه به صرفه جویی در مصرف BRAM، تاخیر مدار نیز مطابق نیاز تنظیم شود.



شکل ۸- نحوه عملکرد روش Tiling

۵.۳. واسط سخت‌افزاری نرم‌افزاری

در بخش‌های قبل توابع فعال‌ساز پرکاربرد شبکه‌های عصبی با استفاده از روش‌های تقریبی به گونه‌ای طراحی شد که با مصرف کمترین منابع، بهبود در تاخیر و دارا بودن دقت کافی، مناسب برای پیاده‌سازی روی سخت‌افزار و قابل استفاده برای شبکه‌های عصبی باشد. همچنین با استفاده از روش Streaming، Tiling و Pipelining فرایند کلی محاسبات در شبکه‌های عصبی به گونه‌ای طراحی شد که باعث کاهش مصرف BRAM در پیاده‌سازی سخت‌افزاری شبکه‌های عصبی گردد. بهینه‌سازی توابع فعال‌ساز و بهینه‌سازی مصرف BRAM این امکان را فراهم کرد تا شبکه‌های بزرگ نیز با تاخیر مناسب و مساحت کمتر روی FPGA قابل پیاده‌سازی باشند.

با وجود ماژول‌های قابل پیاده‌سازی شبکه‌های عصبی روی FPGA می‌توان با استفاده از یک واسط سخت‌افزاری و نرم‌افزاری، شبکه‌های تولید شده در پلتفرم‌های نرم‌افزاری Keras یا Pytorch را به شبکه‌های قابل پیاده‌سازی برای FPGA تبدیل کرد. واسط نرم‌افزاری ارائه شده با استفاده از زبان Python پیاده‌سازی شده است و این امکان را دارد تا لایه‌های LSTM، GRU و DENSE را از شبکه‌ی نرم‌افزاری آموزش داده شده تولید کند، که برای این کار از ماژول‌های بهینه شده مانند Sigmoid و باقی‌ماژول‌ها و روش‌های توضیح داده شده استفاده می‌کند. خروجی تولید شده کد C++ هر لایه برای HLS است که این امکان را می‌دهد تا هر لایه به صورت IP^{24} جداگانه با استفاده از vivado hls پیاده‌سازی شود و ورودی و خروجی هر لایه قابل استفاده و ذخیره سازی باشد. با استفاده از فایل‌های C++ تولید شده

در vivado hls می توان IP های مورد نیاز هر لایه را بدست آورد، سپس با استفاده از ابزار vivado ارتباط بین IP ها و پردازنده ایجاد می شود.

در ابتدا این واسط نرم افزاری فایل خروجی پلتفرم نرم افزاری را بررسی کرده، و به ازای هر لایه، تعداد ورودی، تعداد نوروها و نوع تابع های فعال ساز را تعیین می کند. سپس با استفاده از اطلاعات بدست آمده اقدام به تولید کد HLS هر لایه با توجه تعداد ورودی و نوروها و نوع تابع فعال ساز می کند.

این واسط به ازای هر IP، ورودی های وزن و بایاسی که باید بر اساس روش Tiling به سیستم وارد شود را نیز تولید می کند که می توان از این فایل ها برای تست شبکه و در زمان پیاده سازی به عنوان ورودی استفاده کرد. بدین ترتیب در سریع ترین زمان ممکن شبکه قابل پیاده سازی روی FPGA تولید شده که قابلیت تست و تغییرات سریع را نیز به کاربر می دهد.

نرم افزار فعلی تولید شبکه های عصبی شامل لایه های شبکه های LSTM، GRU و DENSE است. در تولید IP، هر لایه نرم افزار از توابع فعال ساز مورد نیاز استفاده می کند و بر اساس اندازه های هر لایه، اندازه های TILING و ضرب ماتریسی شبکه را تنظیم می کند. همچنین این امکان وجود دارد که ۴ نوع لایه DENSE با استفاده از ۴ تابع فعال ساز ارائه شده تولید کند که بنابراین در واسط سخت افزاری نرم افزاری فعلی امکان تولید ۶ نوع لایه مختلف بهینه شده وجود دارد.

۶. ارزیابی روش پیشنهادی

برای بررسی نتایج این مقاله از برد pynq-z1 به عنوان برد مبنا استفاده شده و برای ارائه نتایج از آمار و ارزیابی های خروجی نرم افزار HLS استفاده شده است. در نتایج بدست آمده مشخص شد که بدون روش های بهینه سازی اعمال شده، فقط پیاده سازی لایه با تعداد ۲۰۰ ورودی و ۵۱۲ نورو روی برد pynq-z1 امکان پذیر است و با این اندازه شبکه، حدود ۲۶۴ عدد از ۲۸۰ عدد BRAM موجود مصرف می شود. اما استفاده از روش Tiling در فرایند کلی محاسباتی شبکه ی عصبی این امکان را فراهم کرد تا بتوان لایه ای با ۳۰۰ برابر ورودی و ۱۲۰ برابر نورو بیشتر نیز قابل پیاده سازی باشد. همچنین پیاده سازی تابع های فعال ساز با روش های تقریبی کمک کرد تا در استفاده از DSP، FF و LUT به مقدار قابل توجهی صرفه جویی شود. در جدول ۱ منابع موجود روی برد pynq-z1 نشان داده شده است.

جدول ۱- منابع موجود روی برد pynq-z1

BRAM_18k	DSP48E	FF	LUT
280	220	106400	53200

همانطور که در جدول ۲ و جدول ۳ مشخص شده است استفاده از روش های تقریبی برای پیاده سازی تابع های فعال ساز Sigmoid و Tanh باعث بهبودی ۹۰ درصدی در مصرف منابع این دو تابع شده است. در جدول ۴ مقایسه منابع مصرفی لایه Softmax نشان داده شده است که با توجه به آمار حدود ۱۶ درصد در مصرف منابع بهبود داده شده است.

جدول ۲- مقایسه منابع مصرفی تابع فعال ساز Sigmoid با روش تقریبی پیشنهادی (Proposed Approx.) و بدون روش تقریبی (no Approx.)

Sigmoid	BRAM_18k	DSP48E	FF	LUT
No Approx.	0	7	907	2519
Proposed Aprox.	0	0	17	229
Improvement	0	100%	98%	99%

جدول ۳- مقایسه منابع مصرفی تابع فعال ساز Tanh با روش تقریبی پیشنهادی (Proposed Approx.) و بدون روش تقریبی (No Approx.)

Tanh	BRAM_18k	DSP48E	FF	LUT
No approx.	0	14	1742	4719
Proposed Aprox.	0	0	25	398
Improvement	0	100%	99%	99%

جدول ۴- مقایسه منابع مصرفی یک لایه softmax با روش تقریبی پیشنهادی (Proposed Approx.) و بدون روش تقریبی (No Approx.)

Softmax	BRAM_18k	DSP48E	FF	LUT
No Approx.	0	7	1384	3964
Proposed Aprox.	0	6	1058	3564
Improvement	0	15%	25%	10%

قابل ذکر است که پیاده‌سازی شبکه با HLS نسبت به HDL دو برابر سریع‌تر انجام می‌شود، و زمان تست و پیاده‌سازی نیز به مقدار چشم‌گیری پایین است. پیاده‌سازی شبکه با استفاده از این واسط سخت‌افزاری نرم‌افزاری زمان مورد نیاز برای پیاده‌سازی دستی را حذف کرده و در سریع‌ترین زمان ممکن یک شبکه تولید می‌کند. علاوه بر این، این امکان را فراهم می‌کند که شبکه‌ی بدست آمده کارایی بالاتری داشته باشد و همچنین قابلیت تغییر شبکه، توسط طراح نیز همچنان وجود داشته باشد. علاوه بر طراحان سخت‌افزاری، امکان استفاده از این پلتفرم برای طراحان نرم‌افزاری هم وجود دارد تا بدون نیاز به دانش سخت‌افزاری بالا شبکه‌ی خود را روی سخت‌افزار پیاده‌سازی کنند.

۷. نتیجه‌گیری

در این مقاله، واسطی سخت‌افزاری نرم‌افزاری ارائه شد که می‌تواند با بهره‌گیری از روش‌های تقریبی و روش‌های بهینه‌سازی حافظه on-chip، شبکه‌های بهینه و قابل پیاده‌سازی روی FPGA را از شبکه‌های تولید شده در پلتفرم‌های نرم‌افزاری تولید کند. نتایج بررسی‌های واسط پیشنهادی مشخص کرد که بیشتر شبکه‌های تولید شده در پلتفرم‌های نرم‌افزاری با روش‌های بهینه‌سازی به کار گرفته شده در این واسط قابل پیاده‌سازی روی FPGA هستند. همچنین با توجه به اینکه تولید کد سخت‌افزاری در یک زبان سطح بالا و توسط ماشین انجام می‌شود، علاوه بر بالا بردن سرعت پیاده‌

سازی، امکان تغییرات بیشتر برای استفاده کننده را نیز در حین کار فراهم می کند و همچنین طراحان و کاربران نرم افزاری هم می توانند بدون داشتن دانش بالا در زمینه پیاده سازی سخت افزاری از این ابزار استفاده کنند.

۸. مراجع

1. Chowdhury, Gobinda G. "Natural language processing." Annual review of information science and technology 37, no. 1 (2003): 51-89.
2. Chiu, Chung-Cheng, Tara N. Sainath, Yonghui Wu, Rohit Prabhavalkar, Patrick Nguyen, Zhifeng Chen, Anjuli Kannan et al. "State-of-the-art speech recognition with sequence-to-sequence models." In 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 4774-4778. IEEE, 2018.
3. Aleksic, Petar, and Pedro J. Moreno Mengibar. "Voice recognition system." U.S. Patent 10,049,666, issued August 14, 2018.
4. Wigington, Curtis, Chris Tensmeyer, Brian Davis, William Barrett, Brian Price, and Scott Cohen. "Start, follow, read: End-to-end full-page handwriting recognition." In Proceedings of the European Conference on Computer Vision (ECCV), pp. 367-383. 2018.
5. Hochreiter, Sepp, and Jürgen Schmidhuber. "Long short-term memory." Neural computation 9, no. 8 (1997): 1735-1780.
6. Hasim Sak, Andrew Senior, and Françoise Beaufays. Long short term memory recurrent neural network architectures for large scale acoustic modeling. In Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), 2014.
7. Guan, Y., Yuan, Z., Sun, G., Cong, J.: FPGA-based accelerator for long short-term memory recurrent neural networks. In: 2017 22nd Asia and South Pacific Design Automation Conference (ASP-DAC), pp. 629-634 (2017).
8. Y. Guan, H. Liang, N. Xu, W. Wang, S. Shi, X. Chen, G. Sun, W. Zhang, and J. Cong, "FP-DNN: An automated framework for mapping deep neural networks onto FPGAs with RTL-HLS hybrid templates," in Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM), 2017, pp. 152-159.
9. Andre Xian Ming Chang, Berin Martini, and Eugenio Culurciello. Recurrent neural networks hardware implementation on fpga. arXiv preprint arXiv:1511.05552, 2015.
10. J. C. Ferreira, J. Fonseca, "An fpga implementation of a long short-term memory neural network" in ReConFigurable Computing and FPGAs (ReConFig) 2016 International Conference on, IEEE, pp. 1-8, 2016.
11. Amin, H., et al: "Piecewise linear approximation applied to nonlinear function of a neural network", *IEE Proc. Circuits, Devices Sys.*, 144, (6), pp. 313-317J, 1997.
12. Pedro Manuel Afonso Costa " Customized Hardware for Long-Short Term Memory Networks in Embedded Systems", IEEE, 2020.
13. I. Kouretas, V. Paliouras "Simplified Hardware Implementation of the Softmax Activation Function", 8th International Conference on Modern Circuits and Systems Technologies, 2019.
14. Zhiqiang Que, Yongxin Zhu, Hongxiang Fan, Hongxiang Fan, Jiuxi Meng, Xinyu Niu & Wayne Luk, "Mapping Large LSTMs to FPGAs with Weight Reuse", Journal of Signal Processing Systems, 2020.



مهندسی برق ایران



20th
IRANIAN STUDENT CONFERENCE
ON ELECTRICAL ENGINEERING



15. D. Svozil, V. Kvasnicka, J. Pospichal "Introduction to multi-layer feedforward neural networks", Chemometrics Intell Lab Systems, 39 (1997), pp. 43-62.

16. Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, Yoshua Bengio, "Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling", Presented in NIPS 2014 Deep Learning and Representation Learning Workshop, 2014.

17. Asano, Shuichi, Tsutomu Maruyama, and Yoshiki Yamaguchi. "Performance comparison of FPGA, GPU and CPU in image processing." In 2009 international conference on field programmable logic and applications, pp. 126-131. IEEE, 2009.

18. Tétrault Marc-André, "Two FPGA Case Studies Comparing High Level Synthesis and Manual HDL for HEP applications", IEEE NPSS Real Time Conference Record, 2018.