

اطلاعیه دفاع

| | | | |
|--|--|--|--|
| نام دانشجو: علیرضا تقوی زاده | | نام استاد راهنما: دکتر دارا رحمتی | |
| مقطع: کارشناسی ارشد | | رشته: مهندسی کامپیوتر | |
| نوع دفاع: | | تاریخ: ۱۴۰۳/۰۷/۳۰ | |
| <ul style="list-style-type: none"> • دفاع پروپوزال <input type="checkbox"/> • دفاع پایان نامه <input checked="" type="checkbox"/> • دفاع رساله دکترا <input type="checkbox"/> | | ساعت: ۱۷:۰۰-۱۸:۰۰ | |
| | | مکان: کلاس ۱۱۷ | |
| عنوان: پیاده‌سازی مدارهای فعال ساز شبکه‌های عصبی با استفاده از محاسبات تصادفی | | | |
| داوران خارجی: دکتر شاهین حسابی | | داوران داخلی: دکتر مهدیانی و دکتر عطارزاده | |
| <p>چکیده: شبکه‌های انتقالی انقلابی در یادگیری عمیق را بنا نهادند که به برتری در پردازش زبان طبیعی، بینایی کامپیوتر و سایر حوزه‌ها منتج شد. با این حال، عملکرد آن‌ها به بهای نیازهای محاسباتی قابل توجه، به ویژه در سازوکار توجه است. پیچیدگی محاسباتی بالا و نیازهای حافظه مانع کارایی آن‌ها در سخت‌افزارهای همه منظوره میشود که منجر به چالش‌هایی در برنامه‌های کاربردی بلادرنگ و دارای محدودیت منابع می‌شود. این پایان‌نامه به بررسی توانایی شتابدهی سخت‌افزاری برای افزایش عملکرد مدل‌های انتقالی با حفظ دقت می‌پردازد. با توجه به وابستگی حجم محاسبه و حافظه در سازوکار توجه در این شبکه به ورودی از درجه دو است و همچنین افزونگی بسیار بالای تعداد وزن‌های این شبکه که شامل تعداد قابل توجهی وزن نزدیک به صفر و کم ارزش است، نیاز است که با رویکردهای متفاوت در ابتدا به کاهش حجم محاسبات بپردازیم که در این پایان‌نامه از دو رویکرد کوانتیزاسیون که برای کاهش عرض بیت محاسبات و شتابدهی کاملاً براساس عدد صحیح و همچنین هرس تدریجی کلمات برای کاهش حجم محاسبات اضافی بهره بردیم. علاوه بر این به دلیل وجود تابع سافت‌مکس در محاسبات سازوکار توجه که باعث تحمیل افزایش حجم محاسبات میگردد از یک تقریب بر اساس معادله درجه دوم استفاده کردیم که این محاسبات نیز کاملاً براساس عدد صحیح بنا نهاده شده‌است. برای افزایش کارایی سخت‌افزاری شتابدهنده پیشنهادی رویکردی پیاده‌سازی به گونه‌ای است که در آن سخت‌افزار سازوکار توجه به شکل یک خط لوله طراحی و پیاده‌سازی شده است که باعث افزایش توان عملیاتی گردد. در نهایت با توجه به مقایسه انجام شده با کارهای گذشته در ارزیابی‌های انرژی مصرفی، مساحت، مساحت مؤثر و حجم حافظه اشغالی بر روی تراشه به ترتیب به کاهش ۲.۳۵، ۳.۷۷، ۳.۲۵ و ۴ برابری دست یافتیم.</p> | | | |