

بهبود کارایی تبدیل موجک گسسته دوبعدی با استفاده از تکنیک موازی سازی در سطح داده

عبدالصیر تیباش^۱، دانشجوی کارشناسی ارشد؛ اسدالله شاه بهرامی^۲، دانشیار

۱- گروه مهندسی کامپیوتر - دانشکده فنی - دانشگاه گیلان - رشت - ایران - BasirTibash@msc.guilan.ac.ir

۲- گروه مهندسی کامپیوتر - دانشکده فنی - دانشگاه گیلان - رشت - ایران - shahbahrami@guilan.ac.ir

چکیده: تبدیل موجک گسسته دوبعدی (2D-DWT) به صورت گسترده‌ای در کاربردهای مختلف پردازش داده‌های چندرسانه‌ای از جمله استانداردهای فشرده سازی تصاویر و ویدئو مورد استفاده قرار می‌گیرد. با این وجود، این تبدیل دارای پیچیدگی محاسباتی بالاتری نسبت به تبدیل‌های مرسوم مانند تبدیل گسسته کسینوسی و دیگر توابع موجود در استانداردهای فشرده سازی است و بیشترین درصد از زمان اجرا را به خود اختصاص می‌دهد. در این مقاله، برای بهبود کارایی 2D-DWT از مجموع دستورات فناوری‌های توسعه برداری پیشرفته AVX/AVX2 و جمع ضرب ترکیبی (FMA) که قابلیت پردازش ۲۵۶ بیت داده با استفاده از معماری یک دستورالعمل و چندین داده (SIMD) که توسط اکثر پردازشگرهای همه منظوره (GPP) پشتیبانی می‌گردد، پیشنهاد شده است. با استفاده از این فناوری‌ها قابلیت پردازش هشت داده ۳۲ بیتی برای اعداد اعشاری و شانزده داده ۱۶ بیتی برای اعداد صحیح شانزده بیتی در ثبات‌های SIMD یک GPP فراهم می‌گردد. علاوه بر این نداشت تبدیل‌های مختلف موجک به روش پردازش‌های سطر-ستونی که پردازش‌های سطر و ستونی را جداگانه انجام می‌دهد و مبتنی بر خط که هر دو، سطرها و ستون‌های تصویر را در یک حلقه پردازش می‌کند، استفاده شده است. نتایج پیاده سازی موازی سازی تبدیل‌های مختلف بر روی یک پلتفرم GPP نشان داد که کارایی، 2D-DWT به ازای اندازه تصاویر مختلف را می‌توان تا ۲۸/۸ برابر نسبت به پیاده سازی سریال بالا برد. همچنین نداشت مبتنی بر خط که باعث استفاده بهتر از ساختار سلسله مراتبی حافظه می‌گردد، کارایی را نسبت به نداشت سطر - ستونی بیشتر بهبود می‌دهد.

واژه‌های کلیدی: پردازشگرهای همه منظوره، پردازش موازی، تبدیل موجک گسسته دوبعدی، موازی سازی سطح داده، یک دستورالعمل چندین داده.

Performance Improvement of 2D Discrete Wavelet Transform using Data-Level Parallelism Technique

A. Tibash¹, MSc. student; A. Shahbahrami², Associate Professor

1- Department of Computer Engineering, Faculty of Engineering, University of Guilan, Rasht, Iran, Email: BasirTibash@msc.guilan.ac.ir

2- Department of Computer Engineering, Faculty of Engineering, University of Guilan, Rasht, Iran, Email: Shahbahrami@guilan.ac.ir

Abstract: The two-Dimensional Discrete Wavelet Transform (2D-DWT) is widely used in various applications for multimedia data processing, including image and video compression standards. However, this transform is computational intensive than conventional conversions, such as the discrete cosine transform. In this paper, in order to improve the performance of 2D-DWT, we use Single Instruction, Multiple Data (SIMD) set instructions including Advanced Vector Extensions (AVX), Fused Multiply-Add (FMA), and AVX2 supported by most General-Purpose Processors (GPP). These technologies capable to process 256-bit data located in SIMD registers. The AVX technology can process eight 32-bit floating point numbers, while AVX2 processes sixteen 16-bit fixed-point numbers. In other words, it is possible to exploit 8- and 16-way data-level parallelism. In addition, two different way of parallelism, Row Column Wavelet Transform (RCWT) which processes rows and columns separately and Line-Based Wavelet Transform (LBWT) that processes both rows and columns in a single loop are used. Experimental results of different wavelet transform with different image sizes on a GPP show that the speedups of up to 28.8x yield. Furthermore, LBWT approach improves performance more than RCWT. This is because it uses memory hierarchy structure more efficiently than RCWT approach.

Keywords: Data-Level Parallelism, Discrete Wavelet Transform, General-Purpose Processor, Parallel processing, Single Instruction, Multiple Data.

تاریخ ارسال مقاله: ۱۳۹۶/۱۰/۱۹

تاریخ اصلاح مقاله: ۱۳۹۷/۰۴/۰۲ و ۱۳۹۷/۰۶/۱۰

تاریخ پذیرش مقاله: ۱۳۹۷/۱۰/۰۳

نام نویسنده مسئول: اسدالله شاه بهرامی

نشانی نویسنده مسئول: گروه مهندسی کامپیوتر - دانشکده فنی - دانشگاه گیلان، رشت، ایران.

۱- مقدمه

و از مجموعه دستورالعمل‌های توسعه برداری پیشرفته (AVX)^۱، AVX2 و جمع ضرب ترکیبی (FMA)^۲ که امکان قابلیت پردازش بردارهای ۲۵۶ بیتی را فراهم می‌کنند بهره نبرده‌اند.

در این مقاله، جهت بهبود کارایی 2D-DWT از مفهوم موازی‌سازی سطح داده (DLP)^۳ با استفاده از فناوری‌های AVX/AVX2 و جمع ضرب ترکیبی (FMA) در SIMD توسعه‌یافته در GPPs پیشنهاد شده است. بدین منظور دو بردارسازی^۴ SIMD برای پنج تبدیل پرکاربرد هار (Haar) و Daub-4 با رویکرد پیچش^۵ و هار، CDF-5/3 و CDF-9/7 با رویکرد پلکانی^۶ بر اساس دو الگوریتم پیاده‌سازی RCWT و LBWT ارائه می‌شود. نتایج ارزیابی کارایی حاکی از آن است که بردارسازی از الگوریتم RCWT میزان کارایی تبدیل‌های مختلف را به ازای ابعاد مختلف تصاویر بین ۲/۷ الی ۲۳/۷ برابر بهبود می‌بخشد درحالی‌که میزان این افزایش کارایی برای بردارسازی SIMD از الگوریتم LBWT بیشتر و بین ۳/۳ الی ۲۸/۸ برابر است.

ادامه مقاله به صورت زیر ساختار بندی می‌شود. در بخش دوم ابتدا مروری کوتاه از مبانی نظری تبدیل موجک گسسته و SIMD ارائه شده و سپس پیشینه تحقیق مورد بررسی قرار گرفته است. در بخش سوم بردارسازی تبدیل موجک گسسته بر اساس دو روش پیاده‌سازی RCWT و LBWT ارائه شده است. بخش چهارم به ارزیابی نتایج آن‌ها بر روی یک پردازشگر با ریز معماری Haswell پرداخته است. در نهایت بخش آخر به نتیجه‌گیری اختصاص یافته است.

۲- مفاهیم اولیه و پیشینه تحقیق

۲-۱- تبدیل موجک گسسته

تبدیل موجک گسسته می‌تواند به عنوان یک روش مناسب جهت تجزیه یک سیگنال به اجزای فرکانسی بالا گذر و پایین گذر درک شود. این تبدیل ویژگی‌های بسیار مناسبی را ارائه می‌دهد که از جمله آن‌ها می‌توان به موارد زیر اشاره نمود [۸]:

- به اندازه کافی اطلاعات مورد نیاز برای تحلیل و بررسی موج اصلی ارائه می‌دهد.
- می‌توان توسط آن سیگنال اصلی را در فرکانس‌های مختلف و با وضوح‌های کاملاً متفاوت بررسی و تحلیل نمود.
- رویکردهای مختلفی برای اعمال DWT بر روی سیگنال از جمله رویکرد پیچش یا کانولوشن و پلکانی ارائه شده است. همچنین، به منظور پیاده‌سازی دوبعدی این تبدیل دو الگوریتم تبدیل موجک سطر-ستونی و تبدیل موجک مبتنی بر خط ارائه گردیده است. در این بخش به بیان این رویکردها و الگوریتم‌ها پرداخته شده است. در زیر بخش‌های ۲-۱-۱ و ۲-۱-۲ به ترتیب دو رویکرد مختلف پیچش و پلکانی برای پیاده‌سازی DWT تشریح شده‌اند. در زیر بخش ۲-۱-۳ الگوریتم‌های پیاده‌سازی DWT دوبعدی با استفاده DWT یک بعدی ارائه شده است.

تبدیل موجک گسسته (DWT)^۱ یک ابزار ریاضی است که امکان تجزیه و تحلیل زمان-فرکانسی سیگنال را فراهم می‌سازد. امروزه این تبدیل در حوزه پردازش تصویر، مخصوصاً در الگوریتم‌های فشرده‌سازی، استخراج ویژگی‌های یک تصویر، نشانه‌گذاری رقمی و تشخیص الگو به صورت گسترده‌ای مورد استفاده قرار می‌گیرد. به عنوان مثال، در استانداردهای فشرده‌سازی JPEG-2000 و MPEG-4 از DWT استفاده شده است. در نتیجه، مزیت تبدیل موجک نسبت به تبدیل‌های مرسوم مانند تبدیل فوریه در حال حاضر به خوبی شناخته شده است [۳-۱].

با این وجود DWT دوبعدی از نظر مدت زمان مورد نیاز برای اجرا، می‌تواند بسیار ناکارآمد باشد. به عنوان مثال، نشان داده شده که در استاندارد JPEG-2000 بیشترین درصد زمان اجرای پردازنده مربوط به تابع DWT دوبعدی است. همچنین، نشان داده شده که DWT به طور میانگین ۴۶ درصد از زمان کدگذاری فشرده‌سازی‌های بدون اتلاف و ۶۸ درصد از زمان کدگذاری فشرده‌سازی‌های با اتلاف را به خود اختصاص داده و پیچیدگی محاسباتی بالاتری نسبت به تبدیل کسینوسی گسسته (DCT)^۲ دوبعدی دارد [۴، ۵].

معماری‌های مختلفی جهت بهبود کارایی 2D-DWT ارائه شده و بر روی سخت‌افزاری مختلفی از جمله پردازشگرهای همه‌منظوره (GPP)^۳ با بهره‌برداری از توسعه یک دستورالعمل و چندین داده (SIMD)^۴ پیاده‌سازی شده‌اند. بردارسازی‌های انجام شده از DWT دوبعدی بر روی GPPها با استفاده از مجموعه دستورالعمل‌های SIMD را می‌توان بر اساس معیارهای مختلفی از جمله الگوریتم‌های پیاده‌سازی دوبعدی، نوع پیاده‌سازی و مجموعه دستورالعمل‌های SIMD استفاده شده به دسته‌های مختلفی تقسیم کرد. از نظر الگوریتم پیاده‌سازی دوبعدی تبدیل موجک گسسته می‌توان آن‌ها را به دو دسته بردارسازی SIMD از الگوریتم تبدیل موجک سطر-ستونی (RCWT)^۵ و الگوریتم موجک مبتنی بر خط (LBWT)^۶ تقسیم کرد. در [۵، ۶] بردارسازی SIMD از تبدیل دابچیز-۴ (Daub-4)^۷ و کوهان-دابچیز-فوو-۵/۳ (CDF-5/3)^۸ با استفاده از این دو الگوریتم ارائه و نتایج حاصل از آن‌ها مقایسه شده‌اند. نتایج حاصل از پیاده‌سازی‌ها نشان می‌دهد اگرچه بردارسازی SIMD از الگوریتم RCWT آسان‌تر از الگوریتم LBWT است، اما بهبود کارایی بردارسازی SIMD با استفاده از الگوریتم LBWT بیشتر است. بر اساس نوع پیاده‌سازی نیز می‌توان آن‌ها را به دو دسته پیاده‌سازی صحیح و اعشاری تقسیم کرد. در [۷] نشان داده شده که می‌توان ضرایب پالایش اعشاری تبدیل‌های مختلف را به صورت صحیح درآورده و محاسبات را به صورت صحیح و با استفاده از مجموعه دستورالعمل‌های صحیح انجام داد. آن‌ها نشان داده‌اند که پیاده‌سازی صحیح به ازای خطایی ناچیزی می‌تواند کارایی را بهبود بخشد. با این وجود در اکثر پیاده‌سازی‌های انجام شده از مجموعه دستورالعمل‌های توسعه‌های چندرسانه‌ای (MMX)^۹ و جریان‌سازی توسعه SIMD (SSE)^{۱۰} به منظور بردارسازی DWT دوبعدی استفاده شده

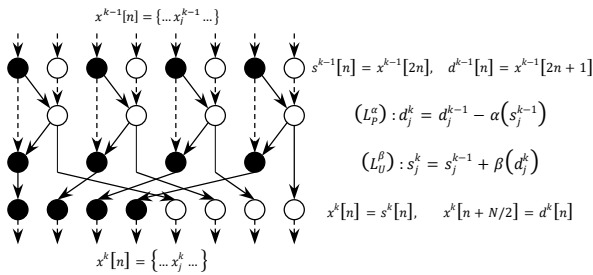
۱-۱-۲- رویکرد پیش

در روش تبدیل موجک از یک تابع پایه به نام موجک مادر و توابع دیگری که از بسط و گسترش موجک مادر تولید می‌شوند، استفاده می‌گردد [۹]. در رویکرد پیش ضرایب موجک رزولوشن‌های پایین از هر مرحله تبدیل موجک گسسته به صورت بازگشتی بر طبق معادلات (۱) و (۲) محاسبه می‌گردد:

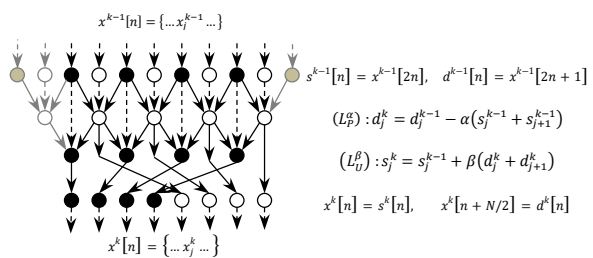
$$c_{j-1,k} = \sum c_{j,m} \cdot h[m - 2k] \quad (1)$$

$$d_{j-1,k} = \sum c_{j,m} \cdot g[m - 2k] \quad (2)$$

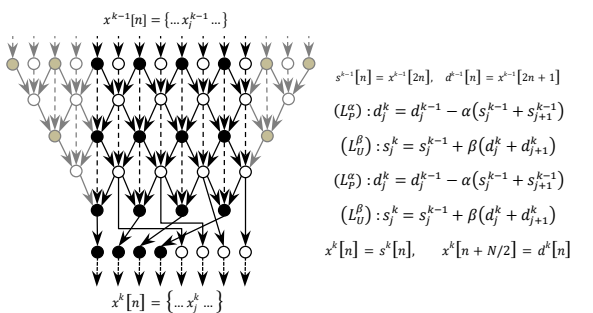
که در آن $c_{j,m}$ نشان‌دهنده زمین ضریب پایین‌گذر در سطح تبدیل m و $d_{j,m}$ نشان‌دهنده زمین ضریب بالاگذر در سطح تبدیل m است. شکل ۱ نمودار بلوکی برای DWT یک‌بعدی با استفاده از رویکرد پیش را برای سه سطح نشان می‌دهد. در این شکل $h[n]$ و $g[n]$ به ترتیب ضرایب پالایش بالا گذر و پایین گذر، $\downarrow 2$ نمونه گاهی با مقیاس ۲ را نشان می‌دهند [۱۰، ۱۱].



الف) تبدیل هار



ب) تبدیل CDF-5/3



ج) تبدیل CDF-9/7

شکل ۳: رویکرد پلکانی برای پیاده‌سازی سه تبدیل هار، CDF-5/3 و CDF-9/7 [۱۲].

ابتدا سیگنال ورودی به دو جریان از ضرایب تقسیم می‌شوند: آرایه‌های زوج و فرد، به‌طور معمول سری‌های تقریب و جزئیات نامیده شده و فرکانس‌های پایین و بالای طیف موجک را نشان می‌دهند. به عبارت دیگر، سیگنال ورودی چندگانه در داخل یک نمایش چندمرحله‌ای به صورت زیر تجزیه شده است:

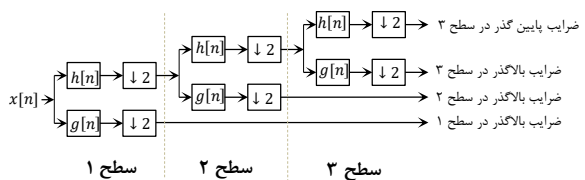
$$s^{k-1}[n] = x^{k-1}[2n], \quad d^{k-1}[n] = x^{k-1}[2n + 1] \quad (3)$$

سپس، قسمت جزئیات با استفاده از مقادیر فعلی هر دو جریان پیش‌بینی (معادله (۴)) و در مرحله بعد قسمت تقریب طبق معادله (۵) به صورت زیر به‌روزرسانی می‌گردد $s^k[n]$ و $d^k[n]$ به ترتیب ضرایب تقریب و جزئیات سطح k از سطح $k-1$ هستند):

$$d^k[n] = d^{k-1}[n] - L_p(s^{k-1}[n]) \quad (4)$$

$$s^k[n] = s^{k-1}[n] + L_U(d^k[n]) \quad (5)$$

به زبان ساده، L_U و L_p (با پیکان مشکی در شکل ۲ نشان داده شده‌اند) که عملگرهای پلکان پیش‌بینی و به‌روزرسانی در معادله (۴) و معادله (۵) هستند، می‌توانند به‌عنوان مجموع وزنی محلی ضرایب پلکان بیان



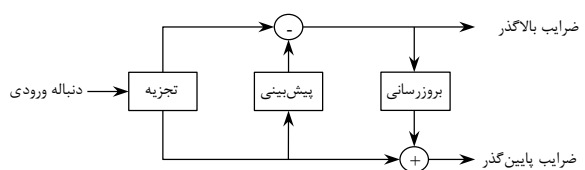
شکل ۱: نمودار بلوکی تحلیل DWT یک‌بعدی با استفاده از رویکرد پیش برای سه سطح تجزیه [۱۰].

پالایه‌ها و فیلترهای مختلفی برای پیاده‌سازی تبدیل موجک گسسته از جمله پالایه‌های دابجیز و CDF که دو تا از مهم‌ترین این پالایه‌ها هستند، ارائه گردیده‌اند. همچنین، در استانداردهای فشرده‌سازی JPEG-2000 و MPEG-4 از این پالایه‌ها استفاده شده است.

۲-۱-۲- رویکرد پلکانی

در این بخش به صورت خلاصه ایده اصلی رویکرد پلکانی ارائه گردیده و سه روش پیاده‌سازی افقی، عمودی و مورب آن بیان شده‌اند. رویکرد پلکانی شامل سه مرحله تجزیه، پیش‌بینی و به‌روزرسانی است که نمودار کلی آن در شکل ۲ نشان داده شده است [۶].

شکل ۳ توضیح مصوری از DWT یک‌بعدی با رویکرد پلکانی بر روی تبدیل‌های هار، CDF-5/3 و CDF-9/7 را نشان می‌دهد. همان‌طور که در این شکل مشاهده می‌گردد رویکرد پلکانی DWT شامل چندین مرحله پلکان (دو مرحله برای هار و CDF-5/3 و چهار مرحله برای CDF-9/7) و یک مرحله نرمال‌سازی اختیاری در پایان است [۱۲].



شکل ۲: رویکرد پلکانی برای پیاده‌سازی DWT دوبعدی [۶].

رویکرد پلکانی مبتنی بر بردار سازی افقی به چندین عمل خواندن و نوشتن نتایج میانی جهت محاسبه یک جفت ضرایب پلکانی d_n و s_n نیازمند است. با توجه به پیاده سازی GPP به ازای تصاویر بزرگ، این نتایج میانی چندین بار از حافظه نهان واحد پردازشگر مرکزی اخراج خواهند شد. در نتیجه، فقدان های حافظه نهان زیادی در چنین محاسباتی رخ خواهد داد [۱۳].

بردار سازی عمودی روشی دیگر جهت بردار سازی گراف جریان داده رویکرد پلکانی است که روش تک حلقه نیز نامیده شده است. گراف داده بردار سازی عمودی برای تبدیل CDF-9/7 با رویکرد پلکانی در شکل ۴ با ناحیه (۲) نشان داده شده است. در این روش، محاسبات تبدیل با رویکرد پلکانی به جای استفاده از حلقه های متعدد بر روی همه ضرایب در داخل یک حلقه انجام می شود؛ بنابراین، در هر تکرار حلقه یک جفت از ضرایب پلکان s_n و d_n محاسبه می گردند. محاسبات در درون ناحیه (۲) نشان داده شده در این شکل، به دلیل وجود وابستگی بین داده ها، نمی تواند به طوری موازی پیاده سازی شود. با این وجود، این روش به دلیل اینکه ضرایب فقط یک بار خوانده و نوشته می شوند بسیار سودمند است. چون این کار مانع از رخ دادن فقدان های حافظه نهان غیر ضروری می شود. در بردار سازی های DWT یک بعدی با استفاده SIMD، می توان چندین ناحیه از بردار سازی عمودی را به صورت موازی هم زمان پردازش کرد [۱۳].

در [۱۴]، بردار سازی مورب از رویکرد پلکانی ارائه شده است. گراف جریان داده بردار سازی مورب برای تبدیل CDF-9/7 با رویکرد پلکانی در شکل ۴ با ناحیه (۳) داده شده است. در این روش نیز مانند بردار سازی عمودی، محاسبات تبدیل پلکانی در داخل یک حلقه انجام شده و در هر تکرار حلقه یک جفت از ضرایب پلکانی s_n و d_n محاسبه گردند. با این تفاوت که در بردار سازی مورب برخلاف بردار سازی عمودی، وابستگی بین داده ها در ناحیه (۳) حذف شده است؛ بنابراین عملیات نشان داده شده در داخل این ناحیه می توانند به صورت موازی انجام شوند.

۲-۱-۳- الگوریتم های پیاده سازی DWT دوبعدی

تبدیل موجک گسسته می تواند به دوبعدی یا چند بعدی توسعه یابد. الگوریتم های مختلفی از جمله تبدیل موجک سطر-ستونی و تبدیل موجک مبتنی بر خط برای پیاده سازی DWT به صورت دوبعدی وجود دارد.

در پیاده سازی مبتنی بر سطر - ستونی (RCWT)، تبدیل موجک گسسته دوبعدی به دو DWT یک بعدی که پالایش افقی و پالایش عمودی نامیده شده اند، تقسیم شده است. پالایش افقی سطرهای تصویر ورودی را پردازش و ضرایب موجک را در ماتریس کمکی ذخیره می کند. سپس، پالایش عمودی ستون های ماتریس کمکی را، پردازش و نتایج را در ماتریس اصلی ذخیره می کند. به عبارت دیگر، در این الگوریتم همه سطرها قبل از شروع فرایند پالایش عمودی، پالایش می شوند. پیچیدگی محاسباتی پالایش افقی و عمودی یکسان و برابر $O(n^2)$ است. در این

شوند. رویکرد پلکانی دو موجک α و β استفاده می کنند در حالی که موجک CDF-9/7 از دو جفت وزن (α) و (β) و (γ) و (δ) که در جدول ۱ لیست شده اند استفاده می کنند. در نهایت، ضرایب نهایی در مرحله ادغام (معادله (۶)) که N نشان دهنده اندازه ابعاد تصویر است) به مکان های مناسب، با یا بدون نرمال سازی، فرستاده شده و یک سطح از تبدیل تکمیل می گردد. معکوس تبدیل را می توان با استفاده از اعمال مراحل بیان شده در جهت عکس و وزن علامت مخالف به دست آورد [۱۲].

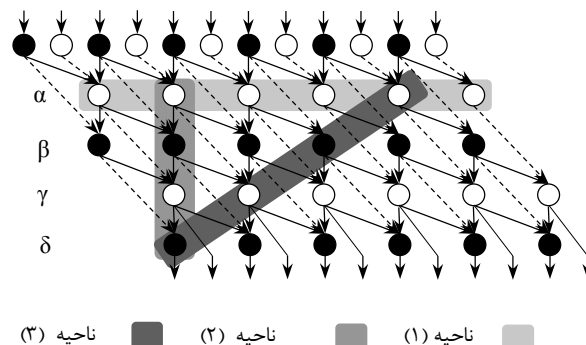
$$x^k[n] = s^k[n], \quad x^k[n + N/2] = d^k[n] \quad (6)$$

جدول ۱: ضرایب پلکان از خانواده های موجک α ، β ، γ و δ CDF-9/7

	CDF-9/7	CDF-5/3	هار	
α	۱/۵۸۶۱۳۴۳۴	۰/۵	۱/۰	
β	-۰/۰۵۲۹۸۰۱۲	۰/۲۵	۰/۵	
γ	-۰/۸۸۲۹۱۱۰۸	---	---	
δ	۰/۴۴۳۵۰۶۸۵	---	---	

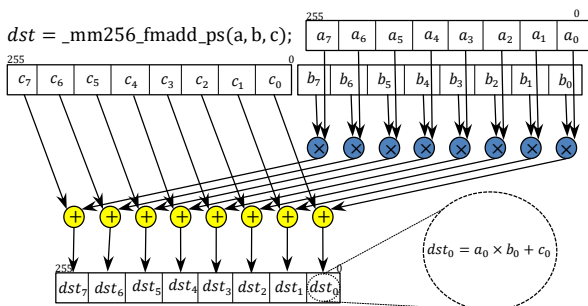
باید توجه داشت که تبدیل CDF-9/7 با تلاف است، بدین معنا که دنباله ورودی به تبدیل CDF-9/7 و دنباله حاصل از معکوس آن به دلیل ضرایب موجک غیر دقیق نشان داده شده در جدول ۱ دقیقاً یکسان نیستند. علاوه بر این، CDF-9/7 به دلیل داشتن چهار مرحله پلکانی دارای هزینه محاسباتی بیشتری در مقایسه با دو موجک α و β CDF-5/3 با دو مرحله پلکانی است. شکل ۳ همچنین نشان می دهد که L_p و L_H تبدیل موجک α و β می تواند به صورت محلی و بدون خواندن مقادیر همسایه در طول یک مرز بلاک اعمال شود چرا که هر دو مرحله پیش بینی و به روز رسانی فقط از مقادیر یک سمت استفاده می کنند؛ اما CDF-9/7 و CDF-5/3 نیاز به پیش بینی و به روز رسانی عنصر فعلی با استفاده از مقادیر همسایه هر دو طرف دارند [۱۲].

به منظور پیاده سازی رویکرد پلکانی سه روش مبتنی بر بردار سازی افقی، عمودی و مورب وجود دارد. در شکل ۴ این سه روش بردار سازی به ترتیب با ناحیه های (۱)، (۲) و (۳) نشان داده شده اند.



شکل ۴: رویکرد پلکانی با سه روش مبتنی بر بردار سازی افقی ناحیه ۱، عمودی ناحیه ۲ و مورب ناحیه ۳.

از بردارسازی خودکار می‌گردد. در [۱۸] دو روش به‌منظور کاهش فاصله^{۱۷} میان بردارسازی خودکار و بردارسازی صریح ارائه شده است. بااین‌وجود، بردارسازی صریح کارایی را به میزان بیشتری نسبت به بردارسازی خودکار بهبود می‌بخشد و برای همین برخی از الگوریتم‌ها به روش صریح بردار سازی می‌شوند. در [۱۹] از بردارسازی صریح SIMD به‌منظور کاهش پیچیدگی و افزایش کارایی در الگوریتم هش و در [۲۰] از مجموعه دستورالعمل‌های SSE4.2، AVX و AVX2 به‌منظور بردارسازی صریح الگوریتم‌های پردازش تصویر استفاده شده که در آن‌ها کارایی را به میزان چشمگیری بهبود بخشیده‌اند. در [۲۱] نیز معماری یک FMA سه سطحی طراحی شده که در عین حفظ کارایی میزان ناحیه انرژی مصرفی را به‌اندازه ۵۰ درصد کاهش داده است. به همین دلیل از بردارسازی صریح به‌عنوان روشی برای بهبود کارایی الگوریتم‌های مختلفی از جمله 2D-DWT استفاده شده که در بخش بعدی پیشینه تحقیق بردارسازی DWT با استفاده از دستورالعمل‌های SIMD مورد بررسی قرار گرفته است.



شکل ۵: ساختار دستورالعمل `_mm256_fmadd_ps`.

۲-۳- پیشینه تحقیق

در [۱۶، ۲۲] تبدیل CDF-9/7 با رویکرد پلکانی را با بهره‌برداری از مجموعه دستورالعمل‌های SIMD پیاده‌سازی کرده‌اند. آن‌ها از بردارسازی خودکار با استفاده از دستورالعمل‌های SSE بهره برده‌اند. در واقع در این کار، بر روی سلسله‌مراتب حافظه تمرکز شده و چندین تکنیک از جمله کاشی کاری جهت بهبود تمرکز موضوعی در زمان و مکان مراجعات در نظر گرفته شده‌اند. همچنین، به‌منظور پیاده‌سازی دوبعدی DWT از الگوریتم LBWT استفاده شده است.

در [۲۳] بر روی بردارسازی تبدیل CDF-9/7 با رویکرد پلکانی تمرکز شده است. در این کار، تبدیل CDF-9/7 با استفاده از روش بردارسازی عمودی و دستورالعمل‌های SSE بردارسازی شده است. در این پیاده‌سازی پالایش افقی و عمودی تبدیل دوبعدی در یک حلقه منفرد ترکیب شده‌اند. به همین دلیل، در این پیاده‌سازی بافری با ابعاد $16 \times M$ جهت ذخیره‌سازی نتایج میانی نیاز است. جایی که M برابر با تعداد ستون‌های تصویر ورودی است.

الگوریتم پالایش‌های افقی و عمودی به‌صورت جداگانه‌ای و به ترتیب به تصویر ورودی اعمال شده‌اند. هر یک از ماتریس‌های $N \times M$ نیازمند NMC بایت از حافظه هستند. جایی که c نشان‌دهنده تعداد بایت‌های موردنیاز برای بیان کردن یک ضریب موجک است [۱۵].

در روش تبدیل موجک مبتنی بر خط، پالایش عمودی بلافاصله پس‌ازاینکه تعداد سطرهای کافی که در رویکرد پیش‌توسط طول ضرایب پالایش مشخص می‌شود، به‌صورت افقی پالایش شدند شروع می‌گردد. الگوریتم LBWT از یک حلقه برای پردازش سطرها و ستون‌ها باهم استفاده می‌کند. این الگوریتم، DWT دوبعدی برای یک تصویر $N \times M$ را توسط مراحل زیر انجام می‌دهد. ابتدا پالایش L سطر از تصویر را پردازش و نتایج مقادیر بالا گذر و پایین گذر را به‌صورت بازه آمیخته در یک بافر $L \times M$ ذخیره می‌کند. سپس، ستون‌های این بافر کوچک توسط پالایش عمودی پردازش و نتایج در ماتریس نهایی ذخیره می‌گردد. در نهایت این مراحل جهت پردازش همه سطرها و ستون‌ها تکرار می‌گردد [۱۵].

۲-۲- یک دستورالعمل و چندین داده

در چند سال اخیر، پردازشگرهای همه‌منظوره به‌منظور افزایش کارایی برنامه‌های کاربردی چندرسانه‌ای، معماری مجموعه دستورالعمل‌های خودشان را با توسعه SIMD از جمله مجموعه دستورالعمل‌های توسعه‌های چندرسانه‌ای، جریان‌سازی توسعه SIMD، توسعه برداری پیشرفته، AVX2 و جمع ضرب ترکیبی گسترش داده‌اند. هر یک از آن‌ها مجموعه‌ای از دستورالعمل‌ها را جهت انجام عملیات SIMD بر روی عناصر داده صحیح و یا اعشاری بسته‌بندی شده فراهم می‌کنند.

توسعه MMX قابلیت پردازش بردارهای ۶۴ بیتی در دامنه صحیح و توسعه SSE قابلیت پردازش بردارهای ۶۴ و ۱۲۸ بیتی در دامنه صحیح و اعشاری را فراهم می‌کنند. درحالی‌که توسعه AVX، AVX2 و FMA قابلیت پردازش بردارهای ۲۵۶ بیتی را فراهم و ویژگی‌های پیشرفته‌تری نسبت به نسل‌های پیشین از توسعه SIMD فراهم می‌کنند.

توسعه FMA، توسعه AVX را بخصوص در محاسبات عددی اعشاری پیشرفت می‌دهد. به‌عنوان مثال در شکل ۵ ساختار دستورالعمل `_mm256_fmadd_ps` که یک دستورالعمل FMA است نشان داده شده است. این دستورالعمل عناصر ۳۲ بیتی اعشاری بسته‌بندی شده در a و b را ضرب، سپس نتایج میانی را با عناصر بسته‌بندی شده در c جمع و نتایج نهایی را در dst ذخیره می‌کند. در واقع امکان انجام شانزده عمل ضرب و جمع توسط یک دستورالعمل فراهم می‌سازد. مجموعه دستورالعمل‌های AVX2 انجام محاسبات در دامنه صحیح و اعشاری با استفاده از ثبات‌های برداری ۲۵۶ بیتی را فراهم می‌کنند [۱۶].

در مقالات زیادی از بردارسازی به‌منظور بهبود کارایی در الگوریتم‌های مختلف استفاده شده است. در [۱۷] بررسی و تحلیل‌ها نشان داده که استراتژی‌های برنامه‌نویسی، الگوریتم‌ها و ساختار داده‌ها مورد استفاده در اکثر برنامه‌های کامپیوتری مانع از بهره‌برداری کامپایلرها

الگوریتم پیاده‌سازی دوبعدی (RCWT یا LBWT)، تبدیل‌های موجک استفاده‌شده، نوع پیاده‌سازی (صحیح یا اعشاری)، مجموعه دستورالعمل‌های استفاده‌شده، سایر خصوصیات این پیاده‌سازی‌ها را نشان می‌دهد. همان‌طور که در این جدول مشاهده می‌شود در اکثر کارها از مجموعه دستورالعمل‌های MMX و SSE به‌منظور پیاده‌سازی 2D-DWT استفاده‌شده است.

۳- روش پیشنهادی

در این بخش بردارسازی SIMD از پنج موجک پرکاربرد از جمله تبدیل‌های، هار و Daub-4 با رویکرد پیچش و تبدیل‌های هار، CDF-5/3 و CDF-9/7 با رویکرد پلکانی با استفاده از فناوری‌های AVX/AVX2 و جمع ضرب ترکیبی (FMA) پیشنهاد شده است. این بردارسازی‌ها بر اساس نوع الگوریتم پیاده‌سازی تبدیل موجک گسسته دوبعدی به دودسته بردارسازی الگوریتم RCWT و بردارسازی الگوریتم LBWT تقسیم شده است. این دو بردارسازی در ادامه این بخش به‌صورت جداگانه‌ای بحث شده‌اند.

۳-۱- بردارسازی الگوریتم تبدیل موجک سطر-ستونی

در این بخش پیاده‌سازی‌های AVX/AVX2 از الگوریتم RCWT برای دو تبدیل هار و Daub-4 با رویکرد پیچش و سه تبدیل هار، CDF-5/3 و CDF-9/7 با رویکرد پلکانی ارائه‌شده و نتایج کارایی آن‌ها بحث شده است. جهت ساده‌سازی بردارسازی SIMD الگوریتم RCWT تا انتهای بخش ۳-۱، تنها بردارسازی SIMD نامیده شده است. این بخش به‌صورت زیر سازمان‌دهی شده است.

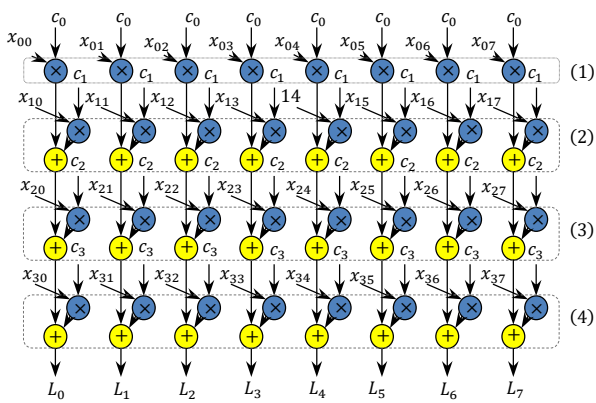
در کارهای زیادی دو تبدیل Daub-4 اعشاری با رویکرد پیچش و CDF-5/3 صحیح با رویکرد پلکانی را به ترتیب با استفاده از دستورالعمل‌های SSE و MMX پیاده‌سازی کرده‌اند. در [۲۴، ۲۵] بردارسازی SIMD از DWT دوبعدی را با استفاده از الگوریتم تبدیل موجک سطر-ستونی ارائه کرده‌اند. در این پیاده‌سازی‌ها نشان داده‌شده که بردارسازی SIMD از پالایش عمودی به‌مراتب سراسر از پالایش افقی است. چراکه، در پالایش عمودی عناصری که باید پردازش شوند به‌صورت متوالی در حافظه ذخیره‌شده‌اند، درحالی‌که بردارسازی SIMD از پالایش افقی نیازمند بازآرایی^{۱۸} عناصر ورودی است. همچنین در [۲۶، ۵] بردارسازی SIMD از این دو تبدیل با استفاده از الگوریتم تبدیل موجک مبتنی بر خط ارائه و نتایج حاصل از آن با بردارسازی با استفاده از الگوریتم RCWT مقایسه شده است. نتایج حاصل از پیاده‌سازی‌ها نشان می‌دهد اگرچه بردارسازی SIMD از الگوریتم RCWT آسان‌تر از الگوریتم LBWT است، اما بهبود کارایی بردارسازی SIMD از تبدیل موجک گسسته دوبعدی با استفاده از الگوریتم LBWT بیشتر است.

در [۲۷] تبدیل CDF-9/7 با رویکرد پلکانی را با بهره‌برداری از خصوصیات مجموعه دستورالعمل‌های SIMD پیاده‌سازی کرده‌اند. این پیاده‌سازی‌ها بر مبنای بردارسازی افقی، عمودی و مورب ارائه‌شده‌اند. نتایج ارزیابی کارایی نشان می‌دهد که بردارسازی SIMD از تبدیل CDF-9/7 با رویکرد پلکانی بر مبنای بردارسازی عمودی به‌مراتب کارا تر از بردارسازی آن بر مبنای بردارسازی افقی و مورب است.

جدول ۲ یک دید کلی از پیاده‌سازی‌های مختلف DWT دوبعدی بر روی پردازشگرهای همه‌منظوره باقابلیت بهره‌برداری از توسعه SIMD بر اساس معیارهای مختلف از جمله رویکرد پیاده‌سازی (پیچش یا پلکانی)،

جدول ۲: یک دید کلی از پیاده‌سازی‌های مختلف DWT دوبعدی بر روی پردازشگرهای همه‌منظوره.

منبع	رویکرد	الگوریتم پیاده‌سازی دوبعدی	تبدیل	نوع پیاده‌سازی	مجموعه دستورالعمل	سایر خصوصیات
[۲۸]	پلکانی	LBWT	CDF-9/7	اعشاری	SSE	برای ذخیره نتایج میانی نیازمند یک بافر $M \times 16$ است. استفاده از تکنیک خط لوله.
[۲۶]	پلکانی	LBWT	CDF-9/7	اعشاری	SSE	برای ذخیره نتایج میانی نیازمند دو بافر است. استفاده از تکنیک بلاک بندی.
[۲۳، ۱۶]	پیچش و پلکانی	RCWT	Daub-4	اعشاری	SSE	استفاده از تکنیک موازی‌سازی سطح داده. بازآرایی عناصر ورودی برای پالایش افقی.
			CDF-5/3	صحیح	MMX	
[۲۴، ۵]	پیچش و پلکانی	LBWT	Daub-4	اعشاری	SSE	برای ذخیره نتایج میانی نیازمند یک بافر $M \times L$ است. استفاده از تکنیک موازی‌سازی سطح داده.
			CDF-5/3	صحیح	MMX	
[۲۹]	پلکانی	LBWT	CDF-9/7	اعشاری	AVX	برای ذخیره نتایج میانی نیازمند یک بافر $M \times 4$ است. استفاده از تکنیک بلاک بندی.



شکل ۶: گراف جریان داده پالایش عمودی برای تبدیل Daub-4

شکل ۷ نحوه محاسبه هشت ضریب موجک خروجی برای پالایش عمودی با استفاده از دستورالعمل‌های AVX/FMA را نشان می‌دهد. همان‌طور که در شکل دیده می‌شود برای عملیات ضرب و جمع‌های متوالی (کادرهای خط‌چین (۲)، (۳) و (۴)) از دستورالعمل `_mm256_fmadd_ps` که یک دستورالعمل FMA است، استفاده گردیده است. در جدول ۳ کارایی آن را در مقایسه با کارایی دو دستورالعمل `_mm256_mul_ps` و `_mm256_add_ps` نشان داده شده است.

بردارسازی پالایش افقی به‌مراتب سخت‌تر از بردارسازی پالایش عمودی است. گراف جریان داده از پالایش افقی در شکل ۸ نشان داده شده است. همان‌طور که این شکل نشان می‌دهد چهار نمونه ورودی مختلف با چهار نمونه مختلف از ضرایب ضرب شده و نتایج میانی در یک عمولند میانه انباشته شده است.

$Vec_A0 = \begin{bmatrix} x_{07} & x_{06} & x_{05} & x_{04} & x_{03} & x_{02} & x_{01} & x_{00} \end{bmatrix}$	$Vec_C0 = \begin{bmatrix} c_0 & c_0 & c_0 & c_0 & c_0 & c_0 & c_0 & c_0 \end{bmatrix}$
$Vec_A1 = \begin{bmatrix} x_{17} & x_{16} & x_{15} & x_{14} & x_{13} & x_{12} & x_{11} & x_{10} \end{bmatrix}$	$Vec_C1 = \begin{bmatrix} c_1 & c_1 & c_1 & c_1 & c_1 & c_1 & c_1 & c_1 \end{bmatrix}$
$Vec_A2 = \begin{bmatrix} x_{27} & x_{26} & x_{25} & x_{24} & x_{23} & x_{22} & x_{21} & x_{20} \end{bmatrix}$	$Vec_C2 = \begin{bmatrix} c_2 & c_2 & c_2 & c_2 & c_2 & c_2 & c_2 & c_2 \end{bmatrix}$
$Vec_A3 = \begin{bmatrix} x_{37} & x_{36} & x_{35} & x_{34} & x_{33} & x_{32} & x_{31} & x_{30} \end{bmatrix}$	$Vec_C3 = \begin{bmatrix} c_3 & c_3 & c_3 & c_3 & c_3 & c_3 & c_3 & c_3 \end{bmatrix}$

- (1) `Vec_L = _mm256_mul_ps(Vec_A0, Vec_C0);`
- (2) `Vec_L = _mm256_fmadd_ps(Vec_A1, Vec_C1, Vec_L);`
- (3) `Vec_L = _mm256_fmadd_ps(Vec_A2, Vec_C2, Vec_L);`
- (4) `Vec_L = _mm256_fmadd_ps(Vec_A3, Vec_C3, Vec_L);`

$Vec_L = \begin{bmatrix} L_7 & L_6 & L_5 & L_4 & L_3 & L_2 & L_1 & L_0 \end{bmatrix}$
--

شکل ۷: محاسبه هشت ضریب موجک خروجی برای پالایش عمودی با استفاده از دستورالعمل‌های AVX/FMA

جدول ۳: کارایی دستورالعمل `_mm256_fmadd_ps` در مقایسه با کارایی دو دستورالعمل `_mm256_mul_ps` و `_mm256_add_ps`

توان عملیاتی	تأخیر	دستورالعمل
۱	۳	<code>_mm256_add_ps</code>
۰/۵	۵	<code>_mm256-mul_ps</code>
۰/۵	۵	<code>_mm256_fmadd_ps</code>

پیاده‌سازی‌های SIMD از دو تبدیل هار و Daub-4 با رویکرد پیش‌در زیر بخش ۳-۱-۱ ارائه شده است. سپس، در زیر بخش ۳-۱-۲ پیاده‌سازی‌های SIMD از تبدیل‌های هار، CDF-5/3 و CDF-9/7 با رویکرد پلکانی مورد بحث قرار گرفته است.

۳-۱-۱-۳ پیاده‌سازی SIMD رویکرد پیش

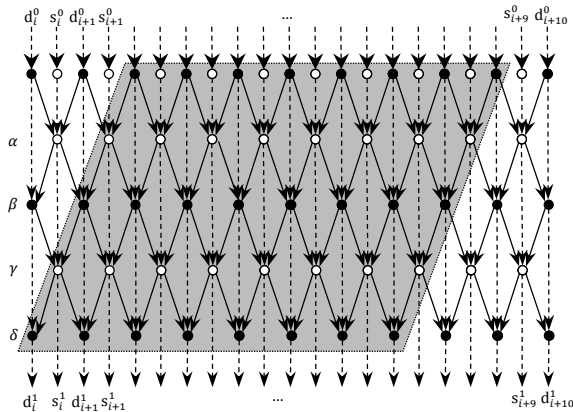
هر دو تبدیل هار و Daub-4 با رویکرد پیش برای محاسبه ضرایب بالا گذر و پایین گذر از ضرایب پالایش اعشاری ۳۲ بیتی استفاده می‌کنند؛ در نتیجه، پیاده‌سازی SIMD این دو تبدیل بسیار مشابه‌اند. به همین دلیل، در ادامه این بخش پیاده‌سازی‌های SIMD این دو تبدیل باهم مورد بحث قرار می‌گیرند.

با توجه به اعشاری بودن ضرایب پالایش این دو تبدیل بسیار مشخص است که بردارسازی آن‌ها با استفاده از دستورالعمل‌های AVX و استخراج هشت مسیر موازی‌سازی $8 \times 32 = 256$ انجام می‌شود. همچنین، همان‌طور که قبلاً اشاره شد بردارسازی SIMD پالایش عمودی نسبت به پالایش افقی سراسر تر است. چون در پالایش عمودی عناصری که می‌توانند هم‌زمان پردازش شوند به‌صورت متوالی در حافظه ذخیره شده‌اند. به‌عنوان مثال، تبدیل Daub-4 را در نظر بگیرید که $x_{i,j}$ نمونه‌های ورودی، c_0, c_1, c_2, c_3 ضرایب پایین گذر و $L_{i,j}$ مقادیر پایین گذر باشند. آنگاه پالایش عمودی از مقادیر پایین گذر را می‌توان با استفاده از معادله (۷) به دست آورد.

$$\begin{aligned}
 & (L_{i,j} L_{i,j+1} L_{i,j+2} L_{i,j+3} L_{i,j+4} L_{i,j+5} L_{i,j+6} L_{i,j+7}) = \\
 & ((c_0 c_0 c_0 c_0 c_0 c_0 c_0 c_0) \times \\
 & (x_{2i,j} x_{2i,j+1} x_{2i,j+2} x_{2i,j+3} x_{2i,j+4} x_{2i,j+5} x_{2i,j+6} x_{2i,j+7})) + \\
 & ((c_1 c_1 c_1 c_1 c_1 c_1 c_1 c_1) \times \\
 & (x_{2i+1,j} x_{2i+1,j+1} x_{2i+1,j+2} x_{2i+1,j+3} x_{2i+1,j+4} x_{2i+1,j+5} x_{2i+1,j+6} x_{2i+1,j+7})) + \\
 & ((c_2 c_2 c_2 c_2 c_2 c_2 c_2 c_2) \times \\
 & (x_{2i+2,j} x_{2i+2,j+1} x_{2i+2,j+2} x_{2i+2,j+3} x_{2i+2,j+4} x_{2i+2,j+5} x_{2i+2,j+6} x_{2i+2,j+7})) + \\
 & ((c_3 c_3 c_3 c_3 c_3 c_3 c_3 c_3) \\
 & \times (x_{2i+3,j} x_{2i+3,j+1} x_{2i+3,j+2} x_{2i+3,j+3} x_{2i+3,j+4} x_{2i+3,j+5} x_{2i+3,j+6} x_{2i+3,j+7})) \quad (7)
 \end{aligned}$$

در این معادله، عملگر \times نشان‌دهنده ضرب چند عنصری برداری است. چون قابلیت موازی‌سازی هشت داده با استفاده از یک دستور SIMD وجود دارد، هر ضریب به‌صورت یک عدد ۳۲ بیتی در یک ثبات SIMD ۲۵۶ بیتی، هشت بار تکرار شده است. معادله مشابهی نیز می‌توان برای مقادیر بالا گذر رسم کرد. شکل ۶ گراف جریان داده را برای پالایش عمودی نشان می‌دهد. همان‌طور که در شکل نشان داده شده، هشت نمونه ورودی از هر سطر به‌صورت هم‌زمان در یک ضریب پالایش (هر ضریب پالایش باید در هر هشت زیر باند یک ثبات رسانه تکرار شود) ضرب شده‌اند. بعد از هشت عملیات ضرب در چهار سطر متوالی با ضرایب مختلف، نتایج هر ستون با یکدیگر جمع شده‌اند.

پلکان بعدی (به ترتیب مرحله بروز رسانی یک، مرحله پیش‌بینی دوم و مرحله بروز رسانی دوم) نیاز است را نمی‌توان در یک ثابت نگه داشت، دسترسی به آن‌ها نیازمند بارگذاری مجدد مقادیر خروجی مرحله پلکان قبلی است. بردارسازی پالایش عمودی تبدیل CDF-9/7 نسب به پالایش افقی ساده‌تر است. چراکه، به دلیل قرار گرفتن دنباله‌های زوج و فرد در سطرها کاملاً جداگانه، پیاده‌سازی آن نیازمند بازآرایی دنباله‌های زوج و فرد ورودی نیست. همانند پیاده‌سازی تبدیل‌های با رویکرد پیچش، بردارسازی پالایش عمودی آسان‌تر بردارسازی پالایش افقی است. در این مورد، دنباله‌های زوج و فرد نیاز نیست که از هم جدا شوند چون آن‌ها در سطرها مختلف قرار دارند؛ اما همانند پالایش افقی به دلیل انجام شدن عملیات جمع در مرحله بروز رسانی، دو کپی از دنباله فرد که یکی از d_{i+1}^0 و دیگری از d_i^0 شروع می‌شود نیاز است. در شکل ۱۰ پیاده‌سازی پالایش عمودی تبدیل CDF-5/3 با رویکرد پلکانی با استفاده از دستورالعمل‌های AVX2 برای شانزده ضریب پایین‌گذر و شانزده ضریب بالا‌گذر نشان داده شده است.



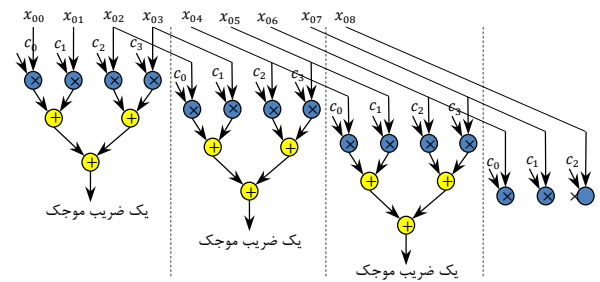
شکل ۹: گراف جریان داده از تبدیل CDF-9/7.

```
Vec_A1 = _mm256_loadu_si256 ((__m256i *)&temp[i][j]);
Vec_A2 = _mm256_loadu_si256 ((__m256i *)&temp[i + 1][j]);
Vec_A3 = _mm256_loadu_si256 ((__m256i *)&temp[i + 2][j]);

//-----Predict 1-----
Vec_A3 = _mm256_add_epi16(Vec_A1, Vec_A3);
Vec_A3 = _mm256_srai_epi16(Vec_A3, 1);
Vec_A3 = _mm256_sub_epi16(Vec_A2, Vec_A3);
int k = i / 2 + N/2;
_mm256_store_si256((__m256i *)&image[k][j], Vec_A3);

//-----Update 1-----
Vec_A2 = _mm256_loadu_si256 ((__m256i *)&image[k - 1][j]);
Vec_A2 = _mm256_add_epi16(Vec_A3, Vec_A2);
Vec_A2 = _mm256_srai_epi16(Vec_A2, 2);
Vec_A2 = _mm256_add_epi16(Vec_A2, Vec_A1);
_mm256_store_si256((__m256i *)&image[i / 2][j], Vec_A2);
```

شکل ۱۰: پیاده‌سازی پالایش عمودی تبدیل CDF-5/3 با رویکرد پلکانی با استفاده از دستورالعمل‌های AVX2 برای شانزده ضریب پایین‌گذر و شانزده ضریب بالا‌گذر.



شکل ۸: گراف جریان داده از پالایش افقی برای تبدیل Daub-4.

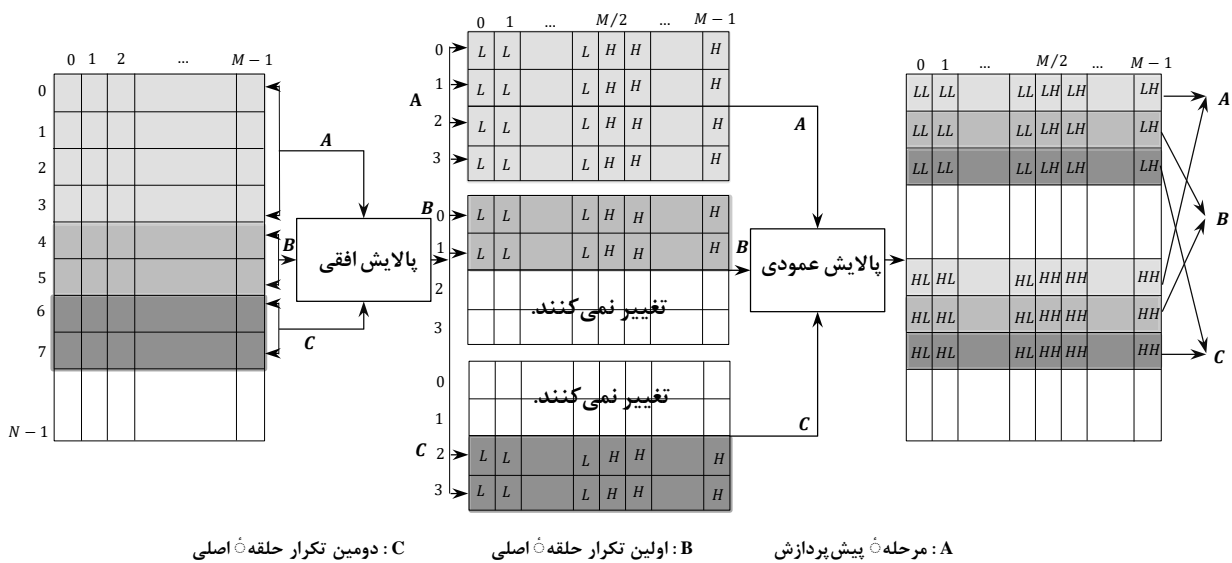
پیاده‌سازی پالایش افقی به روش پیاده‌سازی پالایش عمودی، نیازمند بازآرایی عناصر ورودی است. برای بازآرایی عناصر ورودی از دستورالعمل `_mm256_setr_ps` (تنظیم عناصر ۳۲ بیتی اعشاری بسته‌بندی شده در `dst` با عناصر عرضه‌شده به ترتیب عکس) استفاده شده است.

۳-۲- بردارسازی SIMD رویکرد پلکانی

به دلیل ساختار رویکرد پلکانی تبدیل‌های هار، CDF-5/3 و CDF-9/7 این رویکرد به روش کاملاً متفاوتی نسبت به تبدیل‌های هار و Daub-4 با رویکرد پیچش بردارسازی شده‌اند. همچنین، پیاده‌سازی‌های SIMD از تبدیل‌های هار و CDF-5/3 با رویکرد پلکانی بسیار متفاوت با پیاده‌سازی SIMD از تبدیل CDF-9/7 با رویکرد پلکانی است. اولاً، تبدیل‌های هار و CDF-5/3 با رویکرد پلکانی از محاسبات صحیح استفاده می‌کنند و به همین دلیل پیاده‌سازی‌های SIMD آن‌ها دستورالعمل‌های AVX2 را بکار می‌برد. درحالی‌که پیاده‌سازی SIMD تبدیل CDF-9/7 با رویکرد پلکانی به دلیل اعشاری بودن محاسبات از دستورالعمل‌های AVX استفاده می‌کند. دوماً، با توجه به ضرایب پالایش‌های هار و CDF-5/3 نشان داده‌شده در جدول ۱ می‌توان این ضرایب پالایش را پیاده‌سازی SIMD تبدیل‌های هار، CDF-5/3 و CDF-9/7 با رویکرد پلکانی از بردارسازی عمودی بهره‌گرفته شده است.

در شکل ۹ گراف جریان داده پالایش افقی تبدیل CDF-9/7 نشان داده شده است. همان‌طور که در شکل دیده می‌شود این تبدیل دارای چهار مرحله پلکانی است. به راحتی می‌توان محاسبات مراحل پلکانی این تبدیل را به صورت متوالی با استفاده از دستورالعمل‌های AVX نگاهت کرد و در هر تکرار حلقه هشت ضریب بالا‌گذر و هشت ضریب پایین‌گذر (ناحیه خاکستری‌رنگ) را محاسبه کرد؛ اما بردارسازی پالایش افقی تبدیل CDF-9/7 همانند تبدیل‌های هار و CDF-5/3 نیازمند بازآرایی دنباله‌های زوج و فرد است.

این بردارسازی با استفاده از دستورالعمل `_mm256_setr_ps` انجام شده است. همچنین، چون مقادیر خروجی مرحله پلکان قبلی (به ترتیب مقادیر بالا‌گذر مرحله پیش‌بینی یک، مقادیر پایین‌گذر مرحله به‌روزرسانی یک و مقادیر بالا‌گذر مرحله پیش‌بینی دوم) که برای مرحله



A: مرحله پیش پردازش B: اولین تکرار حلقه اصلی C: دومین تکرار حلقه اصلی

شکل ۱۱: مرحله پیش پردازش و دو تکرار حلقه اصلی پیاده‌سازی SIMD الگوریتم LBWT برای تبدیل Daub-4 با رویکرد پیش با استفاده از موازی سازی هشت-مسیره.

سطرهای ۰ و ۱ بافر ذخیره می‌شوند. چراکه این دو سطر در مرحله قبلی به صورت عمودی پالایش شده و دیگر مورد استفاده قرار نمی‌گیرند. سپس، پالایش عمودی به بافر اعمال می‌گردد. در دومین تکرار حلقه اصلی، مراحل موجود در اولین تکرار حلقه با این تفاوت که نتایج حاصل از پالایش افقی در سطرهای ۲ و ۳ از بافر ذخیره می‌گردد، تکرار می‌شوند. در نهایت این دو مرحله به صورت متوالی بروی کل تصویر ورودی اعمال می‌گردند.

جدول ۴: خصوصیات بافرهای استفاده شده برای بردار سازی الگوریتم LBWT پنج تبدیل مدنظر.

رویکرد	تبدیل	تعداد بافر	اندازه بافر
پیچش	هار	۱	$2 \times M$
	Daub-4	۱	$4 \times M$
پلکانی	هار	۱	$2 \times M$
	CDF-5/3	۱	$3 \times M$
	CDF-9/7	۲	$5 \times M$

جدول ۴ خصوصیات بافرهای استفاده شده برای بردار سازی الگوریتم LBWT دو تبدیل هار و Daub-4 با رویکرد پیچش و سه تبدیل هار، CDF-9/7 و CDF-5/3 با رویکرد پلکانی را نشان می‌دهد. به این نکته باید توجه کرد که به دلیل وجود چهار مرحله پلکانی در تبدیل CDF-9/7 این تبدیل نیازمند یک بافر اضافی برای ذخیره سازی نتایج حاصل از دو مرحله اول پلکانی است. چون این نتایج در تکرار بعدی حلقه مورد استفاده قرار می‌گیرند. همچنین، تبدیل‌های Daub-4 و هار با رویکرد پیچش و CDF-9/7 با رویکرد پلکانی به صورت اعشاری و با استفاده از دستور العمل‌های AVX/FMA پیاده‌سازی شده‌اند. در حالی که

۲-۳- بردار سازی الگوریتم تبدیل موجک مبتنی بر خط

همان‌طور که در بخش ۲-۱-۳ بیان شد، الگوریتم تبدیل موجک مبتنی بر خط برای پردازش سطرها و ستون‌ها تنها از یک حلقه استفاده می‌کند. در این بخش ابتدا بردار سازی الگوریتم LBWT برای دو تبدیل هار و Daub-4 با رویکرد پیچش و سه تبدیل هار، CDF-9/7 و CDF-5/3 با رویکرد پلکانی در زیر بخش ۳-۲-۱ ارائه شده است. سپس، نتایج ارزیابی کارایی برای بردار سازی الگوریتم LBWT در زیر بخش ۳-۲-۲ ارائه گردیده است.

۳-۲-۱- بردار سازی الگوریتم LBWT برای رویکرد پیچش و پلکانی

بردار سازی الگوریتم LBWT برای دو تبدیل هار و Daub-4 با رویکرد پیچش و سه تبدیل هار، CDF-9/7 و CDF-5/3 با رویکرد پلکانی بسیار مشابه است. به منظور بردار سازی الگوریتم LBWT بافر صف حلقوی به اندازه $L \times M$ استفاده می‌شود، جایی که M برابر با عرض تصویر ورودی و L برای تبدیل‌های مختلف متفاوت است. به عنوان مثال، اندازه L برای Daub-4 با رویکرد پیچش برابر طول ضرایب یعنی چهار است.

در شکل ۱۱ مرحله پیش پردازش و دو تکرار حلقه اصلی پیاده‌سازی SIMD الگوریتم LBWT برای تبدیل Daub-4 با رویکرد پیچش با استفاده از موازی سازی هشت-مسیره نشان داده شده است. در مرحله پیش پردازش ابتدا چهار سطر از تصویر ورودی به صورت افقی پردازش و نتایج در چهار سطر از بافر ذخیره می‌شود. سپس، بافر پر شده به صورت عمودی پالایش شده و ضرایب موجک محاسبه شده در یک ماتریس $N \times M$ اضافی ذخیره می‌گردد. در اولین تکرار حلقه اصلی، ابتدا دو سطر بعدی از تصویر ورودی به صورت افقی پردازش و نتایج دست‌آمده در

داده و دستورالعمل را کاهش می‌دهد [۳۰-۳۲]. برای ارزیابی کارایی پیاده‌سازی‌ها شش تصویر محک‌زن^{۱۹} استاندارد در پردازش تصویر که در شکل ۱۲ نشان داده شده‌اند، با ابعاد متفاوت از 128×128 تا 4096×4096 مورد استفاده قرار گرفته‌اند و متوسط کارایی به ازای هر اندازه تصویر در شکل‌های ۱۳ تا ۱۵ نشان داده شده است. تمامی پیاده‌سازی‌ها، کدهای سریال و موازی توسط یک برنامه‌نویس و بر روی یک سیستم مورد آزمون و ارزیابی قرار گرفته است.

جدول ۶: پارامترهای معماری سکوی پیاده‌سازی.

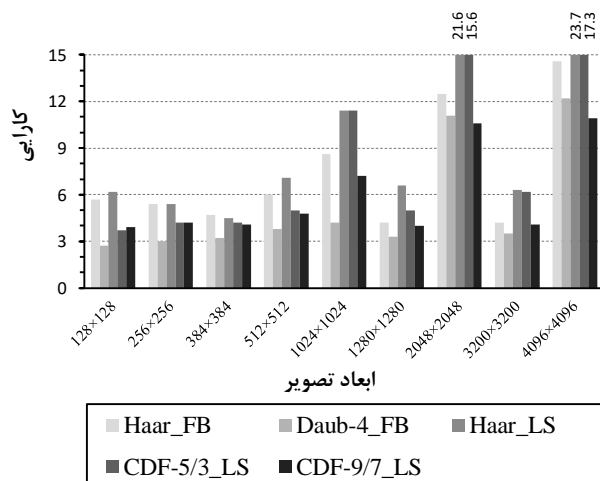
واحد پردازشگر	Intel® Core™ i7-4500U CPU @1.80 GHz
حافظه نهان L1	128 KB
حافظه نهان L2	512 KB
حافظه نهان L3	4 MB
حافظه اصلی	6 GB



شکل ۱۲: تصاویر محک‌زن استاندارد استفاده شده برای ارزیابی کارایی

۴-۳- ارزیابی کارایی بردارسازی الگویتیم RCWT

در این بخش پیاده‌سازی‌های سریال تبدیل‌ها به‌عنوان پیاده‌سازی پایه در نظر گرفته شده و نتایج پیاده‌سازی‌های SIMD موازی شده برای تمامی تبدیل‌ها با آن‌ها مقایسه شده‌اند. در شکل ۱۳ کارایی پیاده‌سازی‌های SIMD برای تبدیل‌های هار و Daub-4 با رویکرد پیچش و تبدیل‌های هار، CDF-5/3 و CDF-9/7 با رویکرد پلکانی نسبت به پیاده‌سازی سریال الگویتیم‌های متناظر نشان داده شده است.



بردارسازی تبدیل‌های هار و CDF-5/3 با رویکرد پلکانی به‌صورت صحیح و با استفاده از دستورالعمل‌های AVX2 انجام شده است.

۴-۴- ارزیابی کارایی

۴-۱- مجموعه داده‌ها

به‌منظور ارزیابی کارایی بردارسازی SIMD بر اساس روش پیشنهادی برای پنج موجک پرکاربرد از جمله تبدیل‌های هار و Daub-4 با رویکرد پیچش که در ادامه به ترتیب با Haar_FB و Daub-4_FB نشان داده می‌شوند و تبدیل‌های هار، CDF-5/3 و CDF-9/7 با رویکرد پلکانی نیز به ترتیب با Haar_LS، CDF-5/3_LS و CDF-9/7_LS نشان می‌شوند ارائه شده است. جدول ۵ یک دید کلی از خصوصیات پیاده‌سازی انجام شده از این تبدیل‌ها را بر اساس معیارهای مختلف از جمله رویکرد پیاده‌سازی (پیچش یا پلکانی)، نوع پیاده‌سازی (صحیح یا اعشاری)، مجموعه دستورالعمل‌های استفاده شده و سطح موازی‌سازی به‌دست‌آمده را نشان می‌دهد. در ستون آخر این جدول، ماکزیمم تعداد داده‌هایی که به‌طور موازی در هر یک از تبدیل‌ها می‌توانند پردازش شوند را نشان داده شده است.

جدول ۵: یک دید کلی از خصوصیات پیاده‌سازی پنج تبدیل موجک.

تبدیل	رویکرد	نوع پیاده‌سازی	مجموعه دستورالعمل‌های SIMD	سطح موازی‌سازی قابل استخراج
هار	پیچش	اعشاری	AVX	۸
Daub-4	پیچش	اعشاری	AVX/FMA	۸
هار	پلکانی	صحیح	AVX2	۱۶
CDF-5/3	پلکانی	صحیح	AVX2	۱۶
CDF-9/7	پلکانی	اعشاری	AVX/FMA	۸

۴-۲- محیط ارزیابی کارایی و تصاویر محک‌زن

همه پیاده‌سازی‌ها بر روی پردازشگر Intel® Core™ i7-4500U با ریز معماری Haswell اجرا شده و پارامترهای معماری این سیستم در جدول ۶ به‌صورت خلاصه نشان داده شده است. همه نسخه‌ها با استفاده از gcc با بهینه‌سازی سطح O2- کامپایل و بر روی یک سیستم با بار کاری سبک اجرا شده‌اند. کارایی با استفاده از نرخ تعداد چرخه اجرا محاسبه شده است. تعداد کلی چرخه‌ها با استفاده از شمارنده چرخه به‌دست‌آمده است. شمارنده چرخه یک ابزار بسیار دقیق برای اندازه‌گیری زمان مصرف شده بین دو نقطه مختلف در اجرای برنامه است. برای حذف اثرات تعویض متن و فقدان‌های حافظه نهان اجباری، از رویکرد اندازه‌گیری K- بهترین و گرم کردن حافظه نهان استفاده شده است. این بدان معناست که تابع K بار اجرا شده و سریع‌ترین زمان در نظر گرفته شده است. اجرای تابع حداقل یک‌بار قبل از شروع اندازه‌گیری تأثیر فقدان‌های حافظه نهان

همچنین شکل ۱۵ کارایی موازی سازی LBWT تبدیلها را نسبت به پیاده سازی سریال نشان می دهد. با توجه به شکل، بردار سازی LBWT میزان کارایی را برای تبدیل های مختلف با ابعاد تصویر مختلف بین ۳/۳ الی ۲۸/۸ برابر بهبود می بخشد. در دو موجک Haar (Haar_LS) و CDF-5/3 این دو تبدیل از مجموعه دستورالعمل های AVX2، اعداد صحیح با قابلیت موازی سازی شانزده داده در یک دستور SIMD استفاده شده است. در حالی که برای سه تبدیل دیگر از مجموعه دستورالعمل های AVX، اعداد اعشاری با قابلیت موازی سازی هشت داده در یک دستورالعمل استفاده شده است. همچنین میزان بهبود کارایی برنامه نویسی SIMD برای تصاویری که دارای ابعاد توانی از دو، 2^n هستند بدلیل فقدان های حافظه نهان کمتر نسبت به سایر ابعاد دیگر بیشتر است. با توجه به نتایج ارائه شده در شکل های ۱۳ و ۱۵ می توان نتیجه گرفت که فناوری های جدید SIMD، کارایی تبدیل های مختلف 2D-DWT را به طور قابل ملاحظه ای افزایش می دهد و هزینه این افزایش فقط روش برنامه نویسی و نگاشت از برنامه سریال به برنامه موازی است. همچنین با توجه به شکل ۱۴ می توان نتیجه گرفت که کارایی الگوریتم برداری سازی LBWT نسبت به RCWT بیشتر است.

۵- نتیجه گیری

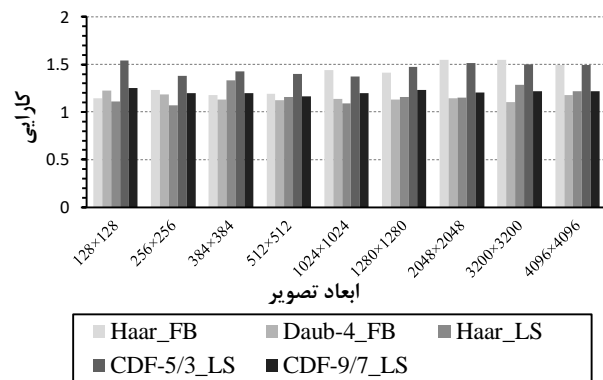
تبدیل موجک گسسته دوبعدی (2D-DWT) یک ابزار قوی در پردازش سیگنال های چندرسانه ای است. ولی این تبدیل از نظر محاسباتی وقت گیر و حجم محاسبات آن زیاد است. به منظور بهبود کارایی این تبدیل از موازی سازی سطح داده با استفاده از مجموع دستورات توسعه یافته، یک دستورالعمل و چند داده (SIMD) با استفاده از فناوری های AVX/AVX2 و FMA که تقریباً در تمامی پردازنده های همه منظوره جدید موجود هستند و قابلیت پردازش ثابت های ۲۵۶ بیتی را دارند پیشنهاد شد. به منظور موازی سازی از دو روش بردار سازی پردازش سطری-ستونی (RCWT) که سطرها و ستون های تصویر جداگانه پردازش می شوند و مبتنی بر خط (LBWT) که هر دو سطرها و ستون ها در یک حلقه پردازش می شوند استفاده گردید. پیچیدگی پیاده سازی برداری الگوریتم LBWT به دلیل انجام پالایش افقی و پالایش عمودی در یک حلقه، نسبت به بردار سازی الگوریتم RCWT بیشتر است؛ اما کارایی موازی سازی آن نسبت به RCWT بیشتر است. به طوری که حداکثر کارایی قابل استخراج با استفاده از RCWT به ازای تبدیل های مختلف موجک و رویکردهای پیچش و پلکانی و اندازه تصاویر متفاوت تا ۴۰۹۶، ۲۳/۷ در حالی که با استفاده از LBWT ۲۸/۸ برابر است. کلیه پیاده سازی ها به روش صریح بوده که توسط یک برنامه نویس صورت گرفته است. برای کارهای آتی پیشنهاد می گردد که عملیات بردار سازی با استفاده از بهینه سازهای کامپایلرها به طور خودکار، ضمنی صورت گیرد و میزان کارایی آن را، با بردار سازی صریح مورد مقایسه و ارزیابی قرارداد و محدودیت های موجود در کامپایلرها به دلیل عدم موقعیت در برخی موارد برای بردار سازی شناسایی گردد.

شکل ۱۳: کارایی پیاده سازی های SIMD تبدیل های مختلف نسبت به پیاده سازی سریال.

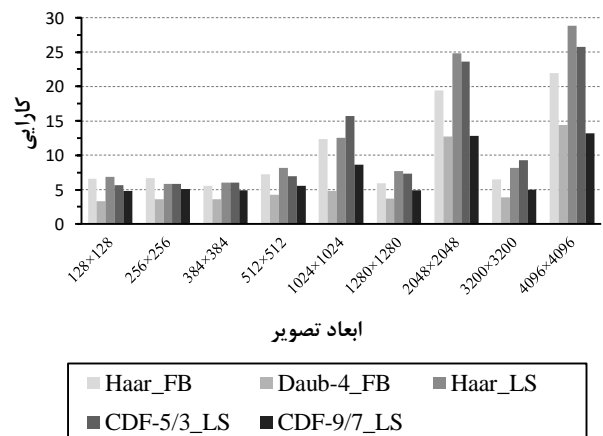
همان طور که در شکل ۱۳ مشاهده می شود بالاترین میزان افزایش کارایی مربوط به تبدیل های Haar_LS (23.7x) و CDF-5/3_LS (17.3x) با رویکرد پلکانی است. چراکه در این تبدیل ها به دلیل صحیح بودن ضرایب از دستورالعمل های AVX2 برای بردار سازی استفاده شده است. به عبارت دیگر تعداد داده های که به طور موازی پردازش می شوند همان طوری که در جدول ۵ نشان داده شده است ۱۶ است. برای سایر تبدیل های دیگر دارای میزان افزایش کارایی بین 2.6x تا 14.6x برای ابعاد مختلف هستند. همچنین با افزایش ابعاد تصاویر میزان افزایش کارایی برای تمامی تبدیل های افزایش یافته است.

۴-۴- ارزیابی کارایی بردار سازی الگوریتم LBWT

در شکل ۱۴ کارایی بردار سازی الگوریتم LBWT برای تبدیل های هار و Daub-4 با رویکرد پیچش و تبدیل های هار، CDF-5/3 و CDF-9/7 نسبت به بردار سازی الگوریتم RCWT نشان داده شده است. همان طور که در شکل دیده می شود برای تمامی تبدیل ها بردار سازی الگوریتم LBWT کارایی را نسبت به بردار سازی الگوریتم RCWT بهبود بخشیده است. چراکه در تمام نمودارها کارایی بیشتر از یک است.



شکل ۱۴: کارایی بردار سازی الگوریتم LBWT برای تبدیل های مورد نظر نسبت به بردار سازی الگوریتم RCWT.



شکل ۱۵: کارایی بردار سازی الگوریتم LBWT برای تبدیل های مدنظر نسبت به پیاده سازی سریال.

مراجع

- [16] D. Chaver, C. Tenllado, L. Piñuel, M. Prieto, and F. Tirado, "2-D Wavelet Transform Enhancement on General-Purpose Microprocessors: Memory Hierarchy and SIMD Parallelism Exploitation," In Proc. Int. Conf. on the High Performance Computing, vol. 2552, pp. 9-21, 2002.
- [17] F. Yazdanpanah, "An Approach for Analyzing Auto-vectorization Potential of Emerging Workloads", Microprocessors and Microsystems, vol. 49, pp. 139-149, 2017.
- [18] H. Saito, S. Preis, N. Panchenko, and X. Tian, *Reducing the Functionality Gap between Auto-vectorization and Explicit Vectorization*, Int. workshop on OpenMP, pp. 173-186, 2016.
- [19] T. Behrens, V. Rosenfeld, J. Traub, S. Breb, and V. Markl, "Efficient SIMD Vectorization for Hashing in OpenCL", Int. Conf. On Extending Database Technology, vol. 3, no. 4, 2018.
- [20] O. Reiche, C. Kobylko, F. Hanniq, J. Teich, "Auto-vectorization for Image Processing DSLs", Proc. in Int. Conf. on Languages, Compilers, and Tools for Embedded Systems, vol. 52, no. 5, pp. 21-30, 2017.
- [21] V. Arunachalam, A.N. J. Raj, N. Hampannavar, C. B. Bidul, "Efficient Dual-precision Floating-point Fused-multiply-add Architecture," Microprocessors and Microsystems, vol. 57, no. 1, pp. 23-31, 2018
- [22] D. G. Chaver, C. Tenllado, L. Pinuel, M. Prieto, and F. Tirado, *Vectorization of the 2D Wavelet Lifting Transform using SIMD Extensions*, IEEE Workshop on Parallel and Distributed Image Processing, Video Processing, and Multimedia, pp. 8-12, 2003.
- [23] R. Kutil, *A Single-Loop Approach to SIMD Parallelization of 2D Wavelet Lifting*, In IEEE Proc. of the 14th Euromicro Int. Conf. on Parallel, Distributed and Network-Based Processing, pp. 413-420, 2006.
- [24] A. Shahbahrami, B. Juurlink, and S. Vassiliadis, "Implementing the 2-D Wavelet Transform on SIMD-Enhanced General-Purpose Processors," IEEE Tran. on Multimedia, vol. 10, no. 1, pp. 43-51, 2008.
- [25] A. Shahbahrami, B. Juurlink, and S. Vassiliadis, *Performance Comparison of SIMD Implementations of the Discrete Wavelet Transform*, In 16th IEEE Int. Conf. on Application-Specific Systems, Architecture Processors, pp. 393-398, 2005.
- [26] A. Shahbahrami, and B. Juurlink, *A Comparison of Two SIMD Implementations of the 2D Discrete Wavelet Transform*, In Proc. 18th Annual Workshop on Circuits, Systems and Signal Processing, pp. 169-177, 2007.
- [27] D. Barina, and P. Zemcik, *Vectorization and Parallelization of 2-D Wavelet Lifting*, Journal of Real-Time Image Processing, pp. 1-13, 2015.
- [28] H. B. Prajapati, and S.K. Vij, *Analytical Study of Parallel and Distributed Image Processing*, In IEEE Int. Conf. on Image Information Processing, pp. 1-6, 2011.
- [29] P.P. Dang, and P.M. Chau, *Integer Fast Wavelet Transform and its VLSI Implementation for Low Power Applications*, In IEEE Workshop on Signal Processing Systems, pp. 93-98, 2002.
- [30] Intel Corporation, Intel® 64 and IA-32 Architectures. Software Developer's Manual, System Programming Guide, vol. 3A, Part 1, 2010.
- [31] D.B. Stewart, *Measuring Execution Time and Real-Time Performance*, In Embedded Systems Conf., vol. 141, 2001.
- [32] Intel Corporation, Intel® Architecture Code Analyzer, User Guide, Document Number: 321356-001US, 2009.
- [1] M. Rabbani, and R. Joshi, "An Overview of the JPEG2000 Still Image Compression Standard," Signal Processing: Image Communication, vol. 17, no. 1, pp. 3-48, 2002.
- [2] S. Battista, F. Casalino, and C. Lande, "MPEG-4: A Multimedia Standard for the Third Millenium, part 1," IEEE Multimedia, vol. 6, no. 4, pp.74-83, 1999.
- [3] S. Battista, F. Casalino, and C. Lande, "MPEG-4: A Multimedia Standard for the Third Millenium, part 2," IEEE Multimedia, vol. 6, no. 4, pp.76-84, 2000.
- [4] M.D. Adams, and R.K. Ward, "JasPer: A Portable Flexible Open-Source Software Tool Kit for Image Coding Processing," In Proc. IEEE, Int. Conf. on Acoustics, Speech, and Signal Processing, vol. 5, pp. 241-244, 2004.
- [5] A. Shahbahrami, "Improving the Performance of 2D Discrete Wavelet Transform using Data-Level Parallelism", In IEEE Int. Conf. on High Performance Computing and Simulation, pp. 362-368, 2011.
- [6] A. Shahbahrami, "Algorithms and Architectures for 2D Discrete Wavelet Transform," The Journal of Supercomputing, vol. 62, no. 2, pp. 1045-1064, May 2012.
- [7] عبدالصیر تیباش و اسدالله شاه بهرامی، «مقایسه کارایی پیاده‌سازی صحیح و اعشاری فیلترهای مختلف تبدیل موجک گسسته دوبعدی»، هشتمین کنفرانس فناوری اطلاعات و دانش، دانشگاه بوعلی سینا، ۱۷ و ۱۸ شهریور ۱۳۹۵.
- [8] فرشته صادقی، ابوالفضل جلیلود و سیدهادی حسینی، «ارائه یک روش ترکیبی مبتنی بر تبدیل موجک گسسته برای پیشبینی بارالکتريکی با استفاده از یک مدل دوبعدی»، مجله مهندسی برق دانشگاه تبریز، جلد ۴۵، شماره ۳، صفحه ۶۸-۷۸، پاییز ۹۴.
- [9] حسین شایقی و علی قاسمی، «پیشبینی قیمت روزانه برق با شبکه عصبی بهبودیافته مبتنی بر تبدیل موجک و روش آشوبناک جستجوی گرانشی»، مجله مهندسی برق دانشگاه تبریز، جلد ۴۵، شماره ۴، صفحه ۱۰۳-۱۱۳، پاییز ۹۴.
- [10] Z. N. Li, M.S. Drew, and J. Liu, *Fundamentals of Multimedia*, NJ: Prentice Hall, 2014.
- [11] H. Liao, M.K. Mandal, and B.F. Cockburn, "Efficient Architectures for 1-D and 2-D Lifting-Based Wavelet Transforms," IEEE Trans. on Signal Processing, vol. 52, no. 5, pp. 1315-1326, 2004.
- [12] T. M. Quan, and W.K. Jeong, "A Fast Discrete Wavelet Transform using Hybrid Parallelism on GPUs," IEEE Trans. on Parallel and Distributed Systems, vol. 27, no. 11, pp. 3088-3100, 2016.
- [13] D. Barina, and P. Zemcik, "Minimum Memory Vectorisation of Wavelet Lifting," In Int. Conf. on Advanced Concepts for Intelligent Vision Systems, Springer, pp. 91-101, 2013.
- [14] D. Barina, and P. Zemcik, "Diagonal Vectorisation of 2-D Wavelet Lifting," In IEEE Int. Conf. on Image Processing, pp. 2978-2982, 2014.
- [15] A. Shahbahrami, and B.H. Juurlink, *SIMD Architectural Enhancements to improve the Performance of the 2D Discrete Wavelet Transform*, 12th Euromicro Conference on Digital System Design, Architectures, Methods and Tools, pp. 497-504, 2009.

زیر نویس ها

⁹ MultiMedia eXtensions

¹⁰ Streaming SIMD Extensions

¹¹ Advanced Vector Extensions

¹² Fused Multiply-Add

¹³ Data Level Parallelism

¹⁴ Vectorization

¹⁵ Convolution

¹⁶ Lifting scheme

¹⁷ Gap

¹⁸ Rearrangement

¹⁹ Benchmark

¹ Discrete Wavelet Transform

² Discrete Cosine Transform

³ General-Purpose Processor

⁴ Single Instruction, Multiple Data

⁵ Row-Column Wavelet Transform

⁶ Line-Based Wavelet Transform

⁷ Daubechies-4

⁸ Cohan-Daubechies-Feauveau-5/3