

به نام خدا



تشخیص دهنده لبه

پردازش چندهسته‌ای - پروژه پایانی - دکتر هاجر فلاحتی

دانشکده مهندسی کامپیوتر

دانشگاه صنعتی شریف

تابستان ۱۴۰۱

نویسندگان:

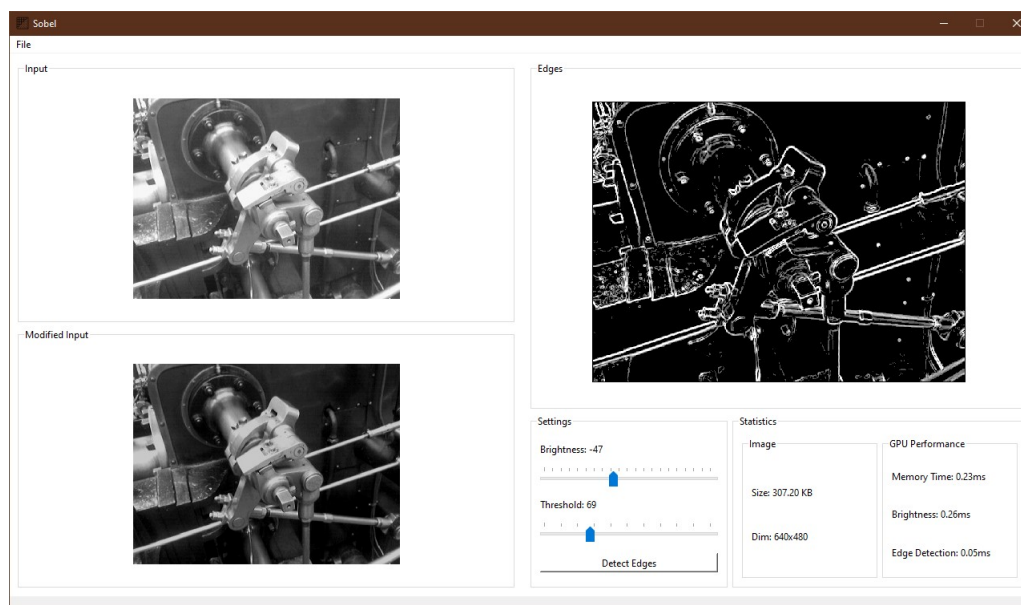
فاطمه خاشعی-۹۷۱۰۵۸۹۹

علی حاتمی تاجیک-۹۸۱۰۱۳۸۵



فهرست مطالب

۱	مقدمه	۲
۲	پیاده‌سازی	۲
۱.۲	پیاده‌سازی ساده	۲
۲.۲	پیاده‌سازی بهینه‌سازی شده	۳
۳.۲	افزودن Shared Memory	۳
۴.۲	تغییر روشنایی	۵
۵.۲	پیکربندی اجرای کرنل	۵
۳	بررسی و مقایسه عملکرد	۶
۱.۳	دلیل بیشتر بودن زمان هنگام استفاده از حافظه اشتراکی	۷
۲.۳	دلیل انتخاب سایز بلوک	۷
۴	رابط کاربری گرافیکی	۸
۵	جمع‌بندی	۸



شکل ۱: شمایل نهایی برنامه

۱ مقدمه

در این پروژه سعی شده است که یک تشخیص دهنده لبه^۱ با استفاده از Spbel Operator پیاده‌سازی شود که بر روی پردازنده‌های گرافیکی شرکت Nvidia (پلتفرم CUDA) اجرا شود. برای بهبود زمان اجرای برنامه از تکنیک‌هایی برای بالابردن بهره‌وری استفاده شده است و همچنین برای رابط کاربری بهتر از واسط کاربری گرافیکی استفاده شده است. شمایل نهایی در شکل ۱ آمده است.

۲ پیاده‌سازی

پیاده‌سازی انجام شده با استفاده از الگوریتمی که در [Wikipedia](https://en.wikipedia.org/wiki/Sobel_operator) آمده بود پیاده‌سازی ساده این الگوریتم را از کد متلبی که داده شده بود الهام گرفتیم و مطابق با آن الگوریتم را در C++ پیاده‌سازی کردیم.

۱.۲ پیاده‌سازی ساده

ساده‌ترین پیاده‌سازی که به ذهن می‌رسد، پیاده‌سازی کانولوشن و انجام آن روی عکس و در نهایت گرفتن جذر جمع مجذور دو عدد و اعمال آستانه^۲ نهایی است. البته برای راحتی در پیاده‌سازی فیلترها به صورتی پیچیده شده اند که ضرب عضو به عضو آنها کفایت کند (چرخاندن فیلتر پیش از اجرای کرنل انجام گرفته است). در این پیاده‌سازی ساده هر ترد مسئول یک پیکسل از خروجی نهایی است و ۹ ضرب (سایز فیلتر) را انجام داده، نتایج را با یکدیگر جمع می‌کند و نتیجه کانولوشن را به ما می‌دهد. این ضرب‌ها با استفاده از دو حلقه حول مرکز فیلتر انجام می‌گیرد. کد کرنل پیاده‌سازی شده در ادامه آمده است.

```
1 __global__ void sobelCUDA(const uint8_t* image,
2   const int8_t* xKernel,
3   const int8_t* yKernel,
4   uint8_t* output,
5   int width,
6   int height,
7   int kernelDim,
```

^۱Edge Detector

^۲Threshold



```

8   int threshold)
9   {
10      int i = blockIdx.y * blockDim.y + threadIdx.y;
11      int j = blockIdx.x * blockDim.x + threadIdx.x;
12      int idx = i * width + j;
13
14      int center = (kernelDim - 1) / 2;
15      float S1 = 0, S2 = 0;
16      int jshift, ishift;
17      int out;
18
19      if (i >= center && j >= center &&
20          i < height - center && j < width - center)
21      {
22          for (int ii = 0; ii < kernelDim; ii++) {
23              for (int jj = 0; jj < kernelDim; jj++) {
24                  jshift = jj + j - center;
25                  ishift = ii + i - center;
26                  S1 += image[ishift * width + jshift] * xKernel[ii * kernelDim + jj];
27                  S2 += image[ishift * width + jshift] * yKernel[ii * kernelDim + jj];
28              }
29          }
30
31          out = sqrtf(S1 * S1 + S2 * S2);
32          out = out > 255 ? 255 : out;
33          output[idx] = out > threshold ? out : 0;
34      }
35  }

```

۲.۲ پیاده‌سازی بهینه‌سازی شده

مشکلی که در کد بالا وجود دارد این است که دسترسی به حافظه‌ای بالا وجود دارد و همینطور با علم به اینکه قرار است چند ضرب و جمع وجود داشته باشد اما با این حال از حلقه‌های تو در تو استفاده شده است که باعث می‌شود هم دسترسی به حافظه بیهوده داشته باشیم و هم دستورات مربوط به حلقه بار اضافی برای اجرا باشند. راه حل این است که ما فیلترها را از پیش داریم و میدانیم که از هر کدام سه ضرب و جمع بیهوده استفاده می‌شوند (چرا که عامل صفر را داریم ضرب می‌کنیم). برای همین از این روش استفاده می‌کنیم که فیلترها را به صورت کد-سخت شده^۳ درون کد کرنل می‌آوریم. با این کار هم دسترسی به حافظه که مربوط به فیلترها بوده است را کم کرده‌ایم، هم حلقه‌ها را از بین برده‌ایم، هم از صفر بودن عامل‌های فیلتر برای کم کردن تعداد دستورات و هم از معلوم بودن فیلتر برای کم کردن تعداد ضرب‌ها (تبدیل ضرب در ۱ و -۱ به جمع و تفریق) استفاده کرده‌ایم. همچنین از الگوی دسترسی منظمی نیز برای کم شدن miss در حافظه کش استفاده شده است. در کرنل بهینه‌سازی شده از قطعه کد زیر به جای حلقه‌های تو در تو استفاده شده است:

```

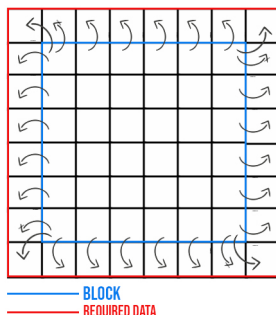
1  S1 = image[(i - 1) * width + (j + 1)] - image[(i - 1) * width + (j - 1)] +
2  (image[(i + 1) * width + (j + 1)] - image[(i + 1) * width + (j - 1)] +
3  2 * (image[i * width + (j + 1)] - image[i * width + (j - 1)]);
4
5  S2 = image[(i - 1) * width + (j - 1)] + (image[(i - 1) * width + (j + 1)] +
6  2 * (image[(i - 1) * width + j] - image[(i + 1) * width + j])
7  - image[(i + 1) * width + (j - 1)] - image[(i + 1) * width + (j + 1)]);

```

۳.۲ افزودن Shared Memory

در این قسمت برای کم کردن دسترسی به حافظه از Shared Memory استفاده می‌کنیم که پیش از عملیات اصلی آن قسمتی از عکس که در اجرای کد اصلی استفاده می‌شود را در حافظه اشتراکی بریزیم. در قسمت ۳ به این مسئله می‌پردازیم که این کار نه تنها باعث افزایش سرعت نمی‌شود بلکه از سرعت می‌کاهد. در هر بلوک ۱۰۲۴ ریسه در چیدمان ۳۲ در ۳۲ کار می‌کنند و هر کدام مسئول یک پیکسل از خروجی هستند که هر کدام به پیکسل‌های اطراف خودشان نیاز دارند که با این حساب در هر بلوک از داده‌های یک قسمت ۳۴ در ۳۴ از عکس اصلی استفاده می‌شود.

³Hard-Coded



شکل ۲: خواندن داده‌های لبه

برای این کار ابتدا هر ترد پیکسل منطبق بر پیکسل خروجی را از ورودی به حافظه اشتراکی می‌آورد. سپس ریسه‌هایی که در لبه‌های قسمت ۳۲ در ۳۲ عکس اصلی قرار دارند پیکسل‌های مجاور خود که در خارج از محدوده ۳۲ در ۳۲ قرار دارند می‌آورند. شکل ۲ نشان می‌دهد که مسئول خواندن هر پیکسل که در حافظه اشتراکی قرار می‌گیرد با کدام ترد متناظر است (تنها برای پیکسل‌های کناری در شکل نشان داده شده‌اند). حافظه اشتراکی به وسیله قطعه کد زیر که پیش از کد اصلی قرار می‌گیرد پر می‌شود:

```
1  __shared__ uint8_t sdata[34 * 34];
2  if (i < height && j < width) {
3      // Main data
4      sdata[(y + 1) * 34 + x + 1] = image[idx];
5
6      // Fix Boudnries
7      if (y == 0 && blockIdx.y != 0) {
8          sdata[x + 1] = image[(i - 1) * width + j];
9          if (x == 0 && blockIdx.x != 0) {
10             sdata[0] = image[(i - 1) * width + (j - 1)];
11         }
12     }
13     if (x == 0 && blockIdx.x != 0) {
14         sdata[y34 + 34] = image[(i) * width + j - 1];
15         if (y == 31 && i != height - 1) {
16             sdata[33 * 34] = image[(i + 1) * width + j - 1];
17         }
18     }
19     if (y == 31 && i != height - 1) {
20         sdata[33 * 34 + x + 1] = image[(i + 1) * width + j];
21         if (x == 31 && j != width - 1) {
22             sdata[33 * 34 + 33] = image[(i + 1) * width + j + 1];
23         }
24     }
25     if (x == 31 && j != width - 1) {
26         sdata[y34 + 34 + 33] = image[i * width + j + 1];
27         if (y == 0 && blockIdx.y != 0) {
28             sdata[33] = image[(i - 1) * width + j + 1];
29         }
30     }
31 }
32
33 // waits untill shared data is completed
34 __syncthreads();
```

حال قطعه کد اصلی به جای اینکه از حافظه داده خود را بردارد از حافظه اشتراکی این کار را انجام می‌دهد.



۴.۲ تغییر روشنایی

برای تغییر روشنایی از یک کرنل ساده استفاده می‌شود که کد آن در ادامه آمده است. پس از افزودن مقدار مورد نظر به هر پیکسل، اگر از باند مربوط به یک پیکسل گذشته بود آنرا به محدوده برمی‌گردانیم.

```
1 __global__ void changeBrightnessCUDA(uint8_t* input, const int width,  
2     const int height, const int brightness)  
3 {  
4     int val;  
5  
6     int i = blockIdx.y * blockDim.y + threadIdx.y;  
7     int j = blockIdx.x * blockDim.x + threadIdx.x;  
8     int idx = i * width + j;  
9  
10    if (i < height && j < width) {  
11        val = input[idx] + brightness;  
12        // Truncate the result (0..255)  
13        if (val > 255) {  
14            val = 255;  
15        }  
16        else if (val < 0) {  
17            val = 0;  
18        }  
19        input[idx] = val;  
20    }  
21 }
```

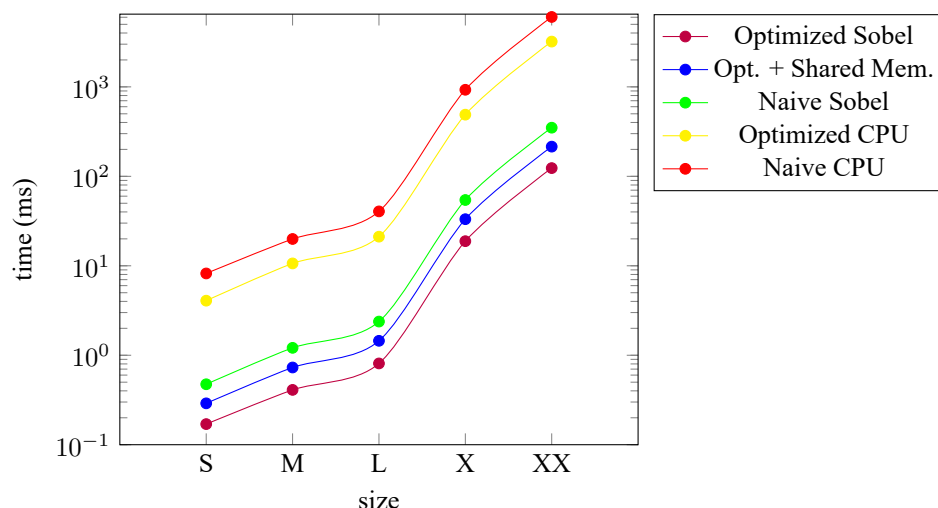
۵.۲ پیکربندی اجرای کرنل

اجرای کرنل‌ها روند کلی زیر را طی می‌کند:

۱. حافظه عکس اصلی تخصیص داده می‌شود.
۲. عکس در حافظه کپی می‌شود.
۳. کرنل تغییر روشنایی اجرا می‌شود.
۴. صبر می‌شود تا تغییرات صورت بگیرد.
۵. نتیجه به صورت آسنکرون روی استریمی جدا به هاست منتقل می‌شود.
۶. حین جایجایی عکسی که روشنایی آن تغییر کرده است کرنل تشخیص لبه اجرا می‌شود.
۷. نتیجه نهایی به هاست منتقل می‌شود.

در ادامه کد اجرایی آمده است. البته برای خوانایی گزارش به صورت خلاصه آمده است.

```
1 dim3 block(BLOCK_DIM, BLOCK_DIM);  
2 dim3 grid(width/BLOCK_DIM + (width%BLOCK_DIM!=0),  
3     height/BLOCK_DIM + (width % BLOCK_DIM != 0));  
4  
5 cudaMalloc((void**)&dev_input, imageSize);  
6 cudaMalloc((void**)&dev_edge, imageSize);  
7 cudaMemcpy(dev_input, input, imageSize, cudaMemcpyHostToDevice);  
8  
9 changeBrightnessCUDA <<<grid, block, 0, main>>> (dev_input, width, height, brightness);  
10 cudaGetLastError();  
11 cudaDeviceSynchronize();  
12  
13 cudaMemcpyAsync(bright, dev_input, imageSize, cudaMemcpyDeviceToHost, mem);
```



شکل ۳: نتایج بنچ‌مارک انجام شده روی متدهای مختلف

```

14
15  sobelOptimizedShCUDA <<<grid, block, 0, main>>> (dev_input, dev_edge, width, height,
16      threshold);
17  cudaEventRecord(stop);
18  cudaGetLastError();
19  cudaDeviceSynchronize();

```

در اینجا با فرستادن غیرهمگام تصویری که روشنایی آن تغییر کرده است، از اورلپ صورت گرفته را در زمان نهایی ذخیره می‌کنیم و **Throughput** حافظه را افزایش می‌دهیم. دلیل نحوه انتخاب این سایز برای بلوک در بخش ۳ توضیح داده شده است.

۳ بررسی و مقایسه عملکرد

در ابتدا به مقایسه عملکرد کرنل‌های متفاوت می‌پردازیم. نتیجه **Benchmark** اجرا شده روی کدهای مختلف در نمودار شکل ۳ آمده است. (دقت شود که اندازه محور زمان لگاریتمی است). در این شکل به وضوح زمان بهتر اجرا بر روی GPU آمده است. برای این تست از سخت‌افزار زیر استفاده شده است:

CPU 11th Gen Intel®Core™i7-1165G7 @ 2.80GHz

GPU NVIDIA GeForce MX450

سایز تصاویر به شرح زیر است:

S 729 x 480

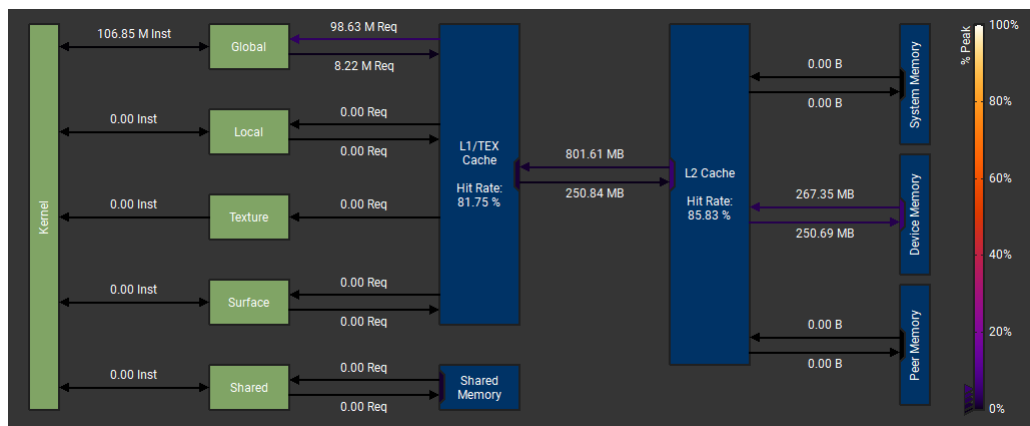
M 1166 x 768

L 1639 x 1080

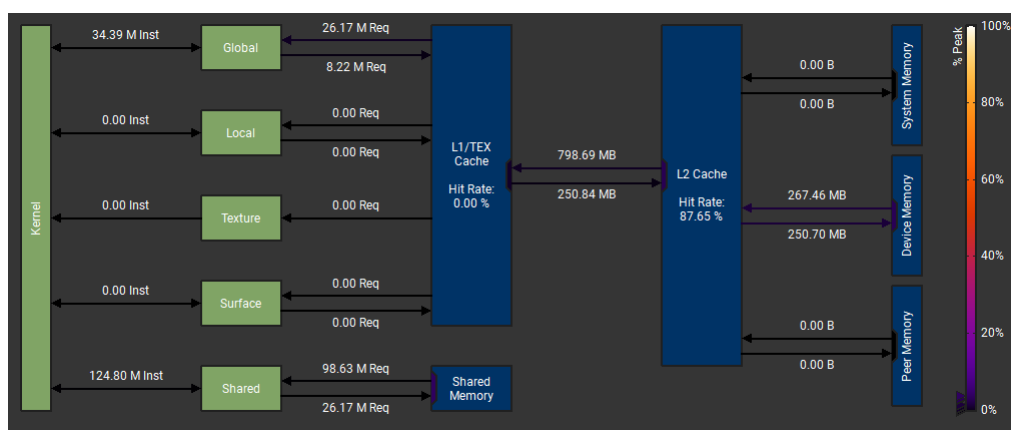
X 7881 x 5192

XX 20000 x 13176

حال به برخی توضیحات می‌پردازیم.



شکل ۴: آنالیز مموری حالت بدون حافظه اشتراکی



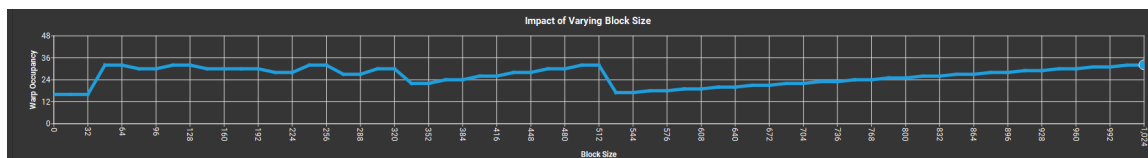
شکل ۵: آنالیز مموری حالت حافظه اشتراکی

۱.۳ دلیل بیشتر بودن زمان هنگام استفاده از حافظه اشتراکی

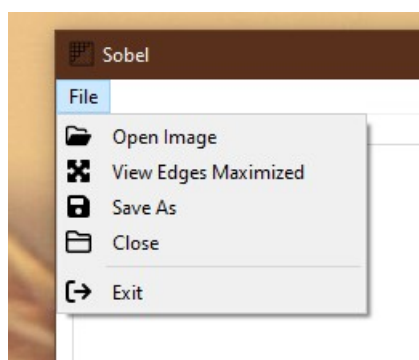
پس از اینکه متوجه این مورد شدیم که زمان اجرای هنگامی که از حافظه اشتراکی استفاده می‌شود بر خلاف انتظارات بیشتر از حالتی است که تماماً از حافظه اصلی می‌خوانیم با استفاده از نرم‌افزار Nsight Compute به پروفایلینگ تست انجام شده پرداختیم. پس از بررسی‌های انجام شده به این نتیجه رسیدیم که دلیل بیشتر شدن زمان این موضوع است که در کد *divergant*هایی ایجاد شده است که کارایی را پایین آورده است. همچنین در حالتی که از حافظه اصلی استفاده می‌شود Hit Rate بالایی در کش مرحله اول داریم که باعث می‌شود بدون حافظه اشتراکی دسترسی به حافظه همانند حالت اشتراکی باشد. در شکل ۵ و ۴ آنالیزهای مربوط به حافظه آمده است. همچنین آنالیز کلی صورت گرفته در محل rep.llvm.org/resources/profile.ncu موجود است.

۲.۳ دلیل انتخاب ساینز بلوک

با توجه به آنالیزی که روی کد انجام گرفته است، اندازه ۱۰۲۴ برای بلوک مناسب‌ترین سایز است و چون سر و کارمان با عکسهاست، اندازه ۱۰۲۴ در ۱ می‌تواند پرتی زیادی در ابعاد ایجاد کند برای همین از بلوک‌های مربعی ۳۲ در ۳۲ استفاده کرده‌ایم که کمترین بلوک ممکن در مساحت استفاده شود. در شکل ۶ آنالیز تعداد بلوک آمده است.



شکل ۶: آنالیز تعداد بلوک حافظه



شکل ۷: منو فایل

۴ رابط کاربری گرافیکی

همانطور که از شکل ۱ پیداست، رابط گرافیکی شامل پنج بخشی کلی است:

- نمایش ورودی
- نمایش ورودی تغییر یافته (روشنایی)
- نمایش خروجی
- انجام تنظیمات قبل از اجرای الگوریتم
- آمار مربوط به اجرا

تنظیمات قبل از اجرا شامل تنظیمات روشنایی و آستانه می‌شود. در منوی فایل همچنین باز کردن، ذخیره کردن، نمایش بزرگتر تصویر خروجی (زمانی که تصاویر بزرگ باشند، خروجی که نمایش داده می‌شود از دقت بالایی برخوردار نیست) و بستن برنامه است (شکل ۷). هنگام بستن و یا تعویض ورودی اگر تغییرات ذخیره نشده‌ای وجود داشته باشد از کاربر تاییدی برای ذخیره کردن یا نکردن گرفته می‌شود. برای پیاده‌سازی این رابط کاربری از فریم‌ورک Qt استفاده شده است.

۵ جمع‌بندی

در این پروژه یک تشخیص‌دهنده لبه برای CUDA پیاده‌سازی شده است که دارای ویژگی‌های زیر است:

- استفاده از مکانیزم‌هایی برای بهینه‌سازی الگوریتم:

- حذف دسترسی حافظه برای فیلتر
- حذف کردن حلقه ^۴

⁴Loop Unrolling



- در نظر نگرفتن عامل‌های صفر
- استفاده از الگوی دسترسی به حافظه منظم
- راهکاری برای استفاده از حافظه اشتراکی با تداخل بانک صفر (به قصد افزایش سرعت)
- استفاده از اندازه بلوک بهینه
- استفاده از کپی غیرهمگام برای همپوشانی کار مموری و محاسبات
- پشتیبانی از اندازه‌های بزرگ (8k)
- پیاده‌سازی آستانه
- پیاده‌سازی GUI