**Final Project: GPU Accelerated Support Vector Machines via Quadratic Programming**
Armaan Kohli - ECE453 Advanced Computer Architecture
Spring 2020

*Remarks*

We implemented a fast support vector machine (SVM) classifier that leverages GPU architecture to achieve performance gains for large datasets. Specifically, our project formulates a SVM classifier as a quadratic program, which we can solve using the alternating direction method of multipliers (ADMM) with either LDL factorisation on the CPU or preconditioned conjugate gradient (PCG) on the GPU. Our project is an extension of the the `osqp` API [1], having added additional tools for SVM problem generation and linear system solving using PCG and CUDA to target an Nvidia Jetson Nano. We compare the results of our GPU implementation to the benchmarks provided by `osqp`.

*Support Vector Machines*

The support vector machine (SVM) is an example of a convex optimization problem w/ quadratic criteria. SVM is a technique for classifying data that may not be linearly separable. This is accomplished by performing a linear classification problem in a higher dimensional feature space where the data is linearly separable.

First, let's assume we have $N$ tuples, $(x_i, y_i)$, where $x_i \in R^m$ represent feature vectors, and $y_i$ represents the true class, $y_i \in \{-1, 1\}$. For instance, in the Fig 1, the blue dots correspond to $x_i$ from $y_i = -1$, and the red dots correspond to features from $y_i = 1$. Let's also assume we've already done some kind of projection into a space where the classes are linearly separable.

Now, we only have access to each value of $x_i$, the class they belong to, the corresponding $y_i$, is unknown. Let's define a hyperplane by:

$$\{x : f(x) = x^T \beta + 1 = 0\} \tag{1}$$

A classification rule induced by this hyperplane $f(x)$ is:

$$G(x) = sign[x^T \beta + 1] \tag{2}$$

Since we are in a space where the classes are linearly separable, we can find a function $f(x) = x^T \beta + 1$ with $y_i f(x_i) > 0 \ \forall \ i$. Hence, we are able to find a hyperplane that creates the largest margin between the two classes. This is what the SVM accomplishes.
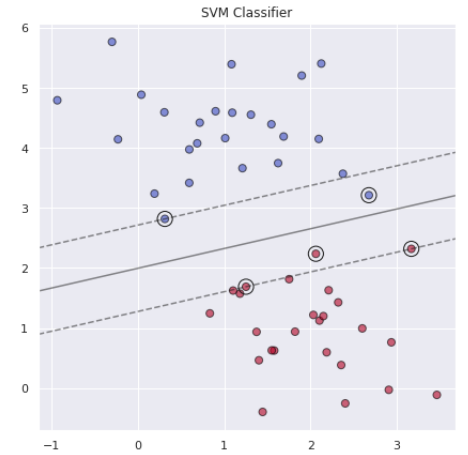


Figure 1: An illustrative example of an SVM classifier. Most of the construction of the SVM section is based on *The Elements of Statistical Learning* [2].

The circled points on the dotted lines are the so-called support vectors. The points between the dashed line and the decision boundary are within the margin. The SVM seeks to find the hyperplane that creates the largest margin, subject to the constraint to minimize the total distance of points on the wrong side of the margin.

We can formulate an SVM as a quadratic program using the method described in [2]. In this form, we can use an algorithm called alternating direction method of multipliers (ADMM).

## Implementation Details

The `osqp` paper presents an two new ways of solving ADMM, a direct method an indirect method. The algorithm in full-form is presented below:

---

**Algorithm 1:** ADMM algorithm as presented in [1]

---

*given: $x^0, z^0, y^0$ and parameters $\rho > 0$, $\sigma > 0$, $\alpha \in [0,2]$*

**while** *not terminated* **do**

$\quad (\tilde{x}^{k+1}, v^{k+1}) \leftarrow \begin{bmatrix} P + \sigma I & A^T \\ A & -\rho^{-1}I \end{bmatrix} \begin{bmatrix} \tilde{x}^{k+1} \\ v^{k+1} \end{bmatrix} = \begin{bmatrix} \sigma x^k - q \\ z^k - \rho^{-1}y^k \end{bmatrix}$

$\quad \tilde{z}^{k+1} \leftarrow z^k + \rho^{-1}(v^{k+1} - y^k)$

$\quad x^{k+1} \leftarrow \alpha \tilde{x}^{k+1} + (1 - \alpha)x^k$

$\quad z^{k+1} \leftarrow \prod \left( \alpha \tilde{z}^{k+1} + (1 - \alpha)z^k + \rho^{-1}y^k \right)$

$\quad y^{k+1} \leftarrow y^k + \rho(\alpha \tilde{z}^{k+1} + (1 - \alpha)z^k - z^{k+1})$

**end**

---

Using this algorithm, we can very quickly compute solutions to SVM problems. However, there is one big issue with the algorithm, namely the matrix inversion step, which must be done every iterations until convergence. This is what we'll exploit to get performance gains. The CPU-based algorithm that `osqp` presents uses LDL matrix factorization to solve this linear system, called the KKT matrix. LDL factorization can be costly as the size of the problem increases, so for large datasets, we can leverage parallelism, and implement an indirect solver for this KKT matrix using PCG. As opposed to a direct method such as KKT, PCG can be efficiently parallelized.

## Preconditioned Conjugate Gradient

In cases where the number of datapoints or the dimensionality of the space is too large for LDL factorization to work effectively, we can instead use an indirect method for solving this KKT linear system. PCG was suggested as a solution in [3]. This is accomplished by re-writing the KKT constraints as in the following form:

$$(P + \sigma I + A^T R A)\tilde{x}^{k+1} = \sigma x^k - q + A^T(Rz^k - y^k) \tag{3}$$

Which can be solved using the PCG algorithm below.

## GPU Optimizations

In order to write an efficient GPU implementation of PCG, we used the following CUDA Toolkit libraries: `Thrust`, `cuBLAS` and `cuSPARSE`. `Thrust` provides a high-level interface for essential data parallel primitives, such as scan and sort operations. `cuBLAS` is a CUDA implementation of BLAS (Basic Linear Algebra Subprograms). We use only level-1 cuBLAS API functions that implement the inner product, axpyoperation, scalar-vector multiplication, and computation of norms. `cuSPARSE` is a CUDA library that contains a set of linear algebra subroutines for handling sparse matrices. Throughout the GPU codebase, we use a CSR format of sparse

**Algorithm 2:** PCG algorithm as presented in [3]

*initialise:* $r^0 = Kx^0 - b$, $y^0 = M^{-1}r^0$, $p^0 = -y^0$, $k = 0$

**while** $||r^k|| > \epsilon||b||$ **do**

$\qquad \alpha^k \leftarrow -\dfrac{(r^k)^T y^k}{(p^k)^T K p^k}$

$\qquad x^{k+1} \leftarrow x^k + \alpha^k p^k$

$\qquad r^{k+1} \leftarrow r^k + \alpha^k K p^k$

$\qquad y^{k+1} \leftarrow M^{-1} r^{k+1}$

$\qquad \beta^{k+1} \leftarrow -\dfrac{(r^{k+1})^T y^{k+1}}{(r^k)^T y^k}$

$\qquad p^{k+1} \leftarrow -y^{k+1} + \beta^{k+1} p^k$

$\qquad k \leftarrow k + 1$

**end**

matricies, since that is the format that `cuSPARSE` uses. The rest of the `osqp` api uses CSC representation, so additional code was needed to perform conversion.

For the kernel functions, we used 64 elements per thread and 1024 threads per block. We found that this gave the best performance empirically, through a more thorough analysis or parameter search would be required to test for optimal settings.

## Results & Discussion

In order to test the performance of our implementation, we generated random SVM problems of various dimensions and computed the run time. To create the CPU benchmarks, we use the CPU code provided by `osqp`. To test the GPU code, we used our own PCG implementation that we then interfaced with the `osqp` api. We generated the SVM problems using a simple `python` script and used the provided code generation tools provided by `osqp` to generate header files containing the data. Below in Fig. 2 is the performance of the SVM classifier according to the results from [3].
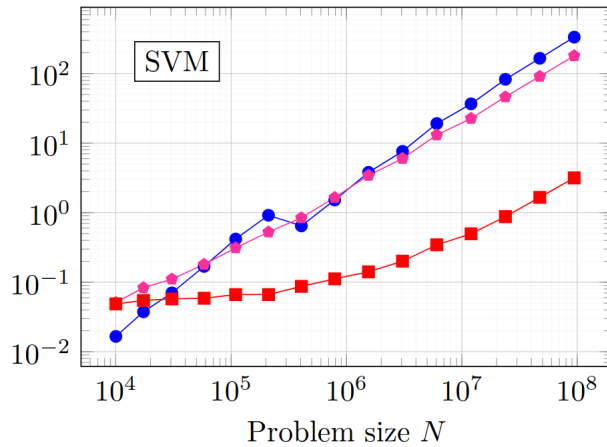


Figure 2: This a plot replicated from [3]. The red line on the right-hand graph shows the performance of the algorithm from [3] in terms of runtime measured in seconds versus problem size. The blue line is the CPU implementation from [1].

They found that using a GPU implementation for SVM solving didn't give massive improvements. However the CPU implementation scales much faster than the GPU implementation as the problem size increases. Furthermore, latency between the host and target devices meant that for smaller problems, the CPU algorithm was superior to the GPU algorithm. But as the problem size increased, the overhead was overcome and the GPU achieved better performance.

For our implementation, we were unable to generate problems greater than $10^4$. However, we did see similar performance improvements, our results are pictured in Fig. 3. As was the case for the official implementation, we found that for small problem size, the overhead made the CPU more efficient than the GPU. But, as the problem size increased, the GPU gave superior performance. And, as the problem size increased, the GPU runtime increased at a much slower rate than the CPU runtime.

Looking at absolute performance however, and not trends, tells us that our GPU implementation is actually superior, since it surpasses CPU performance for smaller problems faster than the official version. GPU performance beats CPU performance around for $N = 3 \times 10^4$ in Fig. 2, whereas this cross over happens much earlier for our implementation, at around $N = 5 \times 10^3$. This could be due to the choice of GPU specific parameters, such as the number of threads per block or elements per thread, or even the physical GPU itself.
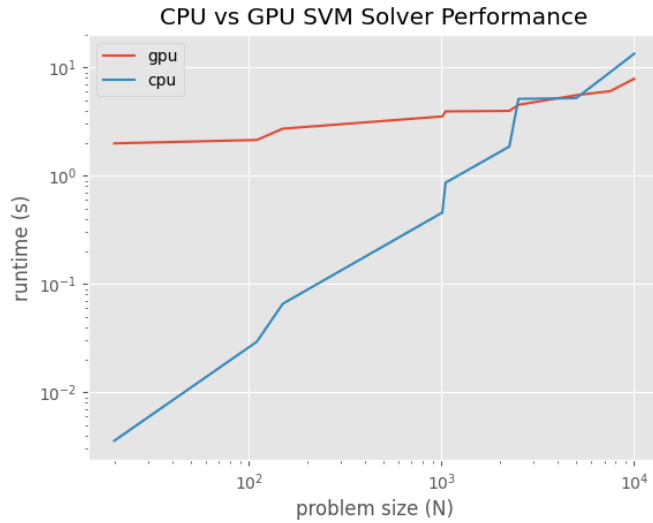


Figure 3: This a plot of our SVM solver results. Unlike the results presented by [3], we were not able to generate problems larger than $10^4$. However, the trends in the graphs are the same. For larger problem sizes, GPU outperforms CPU, and GPU performance scales better for larger problems. These metrics were compiled by generating random SVM problems of different sizes, then averaging over several problems. All experiments were run on a Jetson Nano.

*Conclusion*

In conclusion, we were able to develop our own version of a GPU implementation of ADMM for SVM classification using the `osqp` api. We also successfully verified the results from [3] a regime of relatively scale problems. Furthermore, our GPU implementation achieves higher performance than reported in the paper in the regime of small scale problems.

*References*

[1] B. Stellato, G. Banjac, P. Goulart, A. Bemporad, and S. Boyd, "OSQP: An operator splitting solver for quadratic programs," *Mathematical Programming Computation*, 2020. [Online]. Available: https://doi.org/10.1007/s12532-020-00179-2

[2] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, ser. Springer Series in Statistics. New York, NY, USA: Springer New York Inc., 2001.

[3] M. Schubiger, G. Banjac, and J. Lygeros, "GPU acceleration of ADMM for large-scale quadratic programming," *arXiv:1912.04263*, 2019.

## Appendix A: Code

The code below is `generate_problem.py`. It is used to generate SVM classification problems

```python
"""
generate random svm problems
"""

import numpy as np
import scipy as sp

from scipy import sparse
import utils.codegen_utils as cu

sp.random.seed(1234)
n = 50
m = 100
N = int(m / 2)
gamma = 1.0
b = np.hstack([np.ones(N), -np.ones(N)])
A_upp = sparse.random(N, n, density=0.5)
A_low = sparse.random(N, n, density=0.5)
Ad = sparse.vstack(
    [
        A_upp / np.sqrt(n) + (A_upp != 0.0).astype(float) / n,
        A_low / np.sqrt(n) - (A_low != 0.0).astype(float) / n,
    ],
    format="csc",
)

# osqp data
Im = sparse.eye(m)
P = sparse.block_diag([sparse.eye(n), sparse.csc_matrix((m, m))], format="csc")
q = np.hstack([np.zeros(n), gamma * np.ones(m)])
A = sparse.vstack(
    [
        sparse.hstack([sparse.diags(b).dot(Ad), -Im]),
        sparse.hstack([sparse.csc_matrix((m, n)), Im]),
    ],
    format="csc",
)
l = np.hstack([-np.inf * np.ones(m), np.zeros(m)])
u = np.hstack([-np.ones(m), np.inf * np.ones(m)])

cu.generate_problem_data(P, q, A, l, u, "svm")
```

The code below is `cuda_pcg.cu`. It contains all of the main CUDA routines for the PGC algorithm.

```
/* cuda pcg algorithm */

#include "cuda_pcg.h"
#include "csr_type.h"
#include "cuda_handler.h"
#include "cuda_malloc.h"
#include "cuda_lin_alg.h"
#include "cuda_wrapper.h"
#include "helper_cuda.h"

#ifdef __cplusplus
extern "C" {
extern CUDA_Handle_t *CUDA_handle;
}
#endif

__global__ void scalar_division_kernel(c_float *res, const c_float *num, const c_float *den)
{
    *res = (*num) / (*den);
}

/* computes:  d_y = (P + sigma*I + A'*R*A) * d_x */
static void mat_vec_prod(cudapcg_solver *s, c_float *d_y, const c_float  *d_x, c_int device)
{

    c_float *sigma;
    c_float H_0 = 0.0;
    c_float H_1  = 1.0;
    c_int n = s->n;
    c_int m = s->m;
    csr *P  = s->P;
    csr *A  = s->A;
    csr *At = s->At;

    if (device)
    {
        sigma = s->d_sigma;
    }
    else
    {
        sigma = s->h_sigma;
    }

    /* d_y = d_x */
    checkCudaErrors(cudaMemcpy(d_y, d_x, n * sizeof(c_float), cudaMemcpyDeviceToDevice));

    /* d_y *= sigma */
    checkCudaErrors(cublasTscal(CUDA_handle->cublasHandle, n, sigma, d_y, 1));

    /* d_y += P * d_x */
```

```
51      checkCudaErrors(cusparseCsrmv(CUDA_handle->cusparseHandle, P->alg, P->m, P->n, P->nnz, &H_1, P->
         MatDescription, P->val, P->row_ptr, P->col_ind, d_x, &H_1, d_y, P->buffer));
52
53      if (m == 0) return;
54
55      if (!s->d_rho_vec)
56      {
57          /* d_z = rho * A * d_x */
58          checkCudaErrors(cusparseCsrmv(CUDA_handle->cusparseHandle, A->alg, A->m, A->n, A->nnz, s->h_rho,
             A->MatDescription, A->val, A->row_ptr, A->col_ind, d_x, &H_0, s->d_z, A->buffer));
59      }
60      else
61      {
62          /* d_z = A * d_x */
63          checkCudaErrors(cusparseCsrmv(CUDA_handle->cusparseHandle, A->alg, A->m, A->n, A->nnz, &H_1, A->
             MatDescription, A->val, A->row_ptr, A->col_ind, d_x, &H_0, s->d_z, A->buffer));
64
65          /* d_z = diag(d_rho_vec) * dz */
66          cuda_vec_ew_prod(s->d_z, s->d_z, s->d_rho_vec, m);
67      }
68
69      /* d_y += A' * d_z */
70      checkCudaErrors(cusparseCsrmv(CUDA_handle->cusparseHandle, At->alg, At->m, At->n, At->nnz, &H_1, At
         ->MatDescription, At->val, At->row_ptr, At->col_ind, s->d_z, &H_1, d_y, A->buffer));
71  }
72
73  /* pcg algorithm */
74  c_int cuda_pcg(cudapcg_solver *s, c_float eps, c_int max_niter)
75  {
76
77      c_float *ptr_tmp;
78      c_int niter = 0;
79      c_int n = s->n;
80      c_float H_m_1 = -1.0;
81
82      /* set up problem */
83
84      if (!s->warm_start)
85      {
86          /* d_x = 0 */
87          checkCudaErrors(cudaMemset(s->d_x, 0, n * sizeof(c_float)));
88      }
89
90      /* d_p = 0 */
91      checkCudaErrors(cudaMemset(s->d_p, 0, n * sizeof(c_float)));
92
93      /* d_r = K * d_x */
94      mat_vec_prod(s, s->d_r, s->d_x, 0);
95
96      /* d_r -= d_rhs */
97      checkCudaErrors(cublasTaxpy(CUDA_handle->cublasHandle, n, &H_m_1, s->d_rhs, 1, s->d_r, 1));
```

```
98
99      /* h_r_norm = |d_r| */
100     s->vector_norm(s->d_r, n, s->h_r_norm);
101
102     /* need to change CUBLAS mode */
103     cublasSetPointerMode(CUDA_handle->cublasHandle, CUBLAS_POINTER_MODE_DEVICE);
104
105     if (s->precondition)
106     {
107         /* d_y = M \ d_r */
108         cuda_vec_ew_prod(s->d_y, s->d_diag_precond_inv, s->d_r, n);
109     }
110
111     /* d_p = -d_y */
112     checkCudaErrors(cublasTaxpy(CUDA_handle->cublasHandle, n, s->D_MINUS_ONE, s->d_y, 1, s->d_p, 1));
113
114     /* rTy = d_r' * d_y */
115     checkCudaErrors(cublasTdot(CUDA_handle->cublasHandle, n, s->d_y, 1, s->d_r, 1, s->rTy));
116
117     /* synchronize for timing */
118     cudaDeviceSynchronize();
119
120     /* Run the PCG algorithm */
121     while ( *(s->h_r_norm) > eps && niter < max_niter )
122     {
123
124         /* d_Kp = K * d_p */
125         mat_vec_prod(s, s->d_Kp, s->d_p, 1);
126
127         /* pKp = d_p' * d_Kp */
128         checkCudaErrors(cublasTdot(CUDA_handle->cublasHandle, n, s->d_p, 1, s->d_Kp, 1, s->pKp));
129
130         /* alpha = rTy / pKp */
131         scalar_division_kernel<<<1,1>>>(s->alpha, s->rTy, s->pKp);
132
133         /* d_x += alpha * d_p */
134         checkCudaErrors(cublasTaxpy(CUDA_handle->cublasHandle, n, s->alpha, s->d_p, 1, s->d_x, 1));
135
136         /* d_r += alpha * d_Kp */
137         checkCudaErrors(cublasTaxpy(CUDA_handle->cublasHandle, n, s->alpha, s->d_Kp, 1, s->d_r, 1));
138
139         if (s->precondition)
140         {
141             /* d_y = M \ d_r */
142             cuda_vec_ew_prod(s->d_y, s->d_diag_precond_inv, s->d_r, n);
143         }
144
145         /* Swap pointers to rTy and rTy_prev */
146         ptr_tmp = s->rTy_prev;
147         s->rTy_prev = s->rTy;
148         s->rTy = ptr_tmp;
```

```
149
150          /* rTy = d_r' * d_y */
151          checkCudaErrors(cublasTdot(CUDA_handle->cublasHandle, n, s->d_y, 1, s->d_r, 1, s->rTy));
152
153          /* Update residual norm */
154          s->vector_norm(s->d_r, n, s->d_r_norm);
155          checkCudaErrors(cudaMemcpyAsync(s->h_r_norm, s->d_r_norm, sizeof(c_float),
       cudaMemcpyDeviceToHost));
156
157          /* beta = rTy / rTy_prev */
158          scalar_division_kernel<<<1,1>>>(s->beta, s->rTy, s->rTy_prev);
159
160          /* d_p *= beta */
161          checkCudaErrors(cublasTscal(CUDA_handle->cublasHandle, n, s->beta, s->d_p, 1));
162
163          /* d_p -= d_y */
164          checkCudaErrors(cublasTaxpy(CUDA_handle->cublasHandle, n, s->D_MINUS_ONE, s->d_y, 1, s->d_p, 1))
       ;
165
166          cudaDeviceSynchronize();
167          niter++;
168
169      } /* End of the PCG algorithm */
170
171      /* change CUBLAS pointer mode back */
172      cublasSetPointerMode(CUDA_handle->cublasHandle, CUBLAS_POINTER_MODE_HOST);
173
174      return niter;
175  }
176
177  /* update preconditioning  */
178  void cuda_pcg_update_precond(cudapcg_solver *s, c_int P_updated, c_int A_updated, c_int R_updated)
179  {
180
181      void    *buffer;
182      c_float *mem_tmp;
183      c_int    n  = s->n;
184      csr      *At = s->At;
185
186      size_t buff_size = n * (sizeof(c_float) + sizeof(c_int));
187
188      if (!P_updated && !A_updated && !R_updated) return;
189
190      if (P_updated)
191      {
192          /* Update d_P_diag_val */
193          checkCudaErrors(cusparseTgthr(CUDA_handle->cusparseHandle, n, s->P->val, s->d_P_diag_val, s->
       d_P_diag_ind, CUSPARSE_INDEX_BASE_ZERO));
194      }
195
196      if (A_updated || R_updated)
```

```
197    {
198        /* Allocate memory */
199        cuda_malloc((void **) &mem_tmp, At->nnz * sizeof(c_float));
200        cuda_malloc((void **) &buffer, buff_size);
201
202        /* Update d_AtRA_diag_val */
203        if (!s->d_rho_vec)
204        {    /* R = rho*I  -->  A'*R*A = rho * A'*A */
205
206            if (A_updated)
207            {
208                /* Update d_AtA_diag_val */
209                cuda_vec_ew_prod(mem_tmp, At->val, At->val, At->nnz);
210                cuda_vec_segmented_sum(mem_tmp, At->row_ind, s->d_AtA_diag_val, buffer, n, At->nnz);
211            }
212
213            /* d_AtRA_diag_val = rho * d_AtA_diag_val */
214            cuda_vec_add_scaled(s->d_AtRA_diag_val, s->d_AtA_diag_val, NULL, *s->h_rho, 0.0, n);
215        }
216        else
217        {    /* R = diag(d_rho_vec)  -->  A'*R*A = A' * diag(d_rho_vec) * A */
218            cuda_mat_rmult_diag_new(At, mem_tmp, s->d_rho_vec);   /* mem_tmp = A' * R */
219            cuda_vec_ew_prod(mem_tmp, mem_tmp, At->val, At->nnz);      /* mem_tmp = mem_tmp * A */
220            cuda_vec_segmented_sum(mem_tmp, At->row_ind, s->d_AtRA_diag_val, buffer, n, At->nnz);
221        }
222
223        cuda_free((void **) &mem_tmp);
224        cuda_free((void **) &buffer);
225    }
226
227    /* d_diag_precond = sigma */
228    cuda_vec_set_sc(s->d_diag_precond, *s->h_sigma, n);
229
230    /* d_diag_precond += d_P_diag_val + d_AtRA_diag_val */
231    cuda_vec_add_scaled3(s->d_diag_precond, s->d_diag_precond, s->d_P_diag_val, s->d_AtRA_diag_val, 1.0,
        1.0, 1.0, n);
232
233    /* d_diag_precond_inv = 1 / d_diag_precond */
234    cuda_vec_reciprocal(s->d_diag_precond_inv, s->d_diag_precond, n);
235 }
```

The code below is `cuda_pcg_interface.c`. It wraps `cuda_pcg.cu` functions into a `lin_sys` solver module specified by the `osqp` api so that we can actually solve a linear system via PCG.

```c
/* interface for pcg_solver with main osqp api */

#include "cuda_pcg_interface.h"
#include "cuda_pcg.h"

#include "cuda_lin_alg.h"
#include "cuda_malloc.h"

#include "glob_opts.h"

static c_float compute_tolerance(cudapcg_solver *s, c_int admm_iter)
{
    c_float eps;
    c_float rhs_norm;

    /* compute the norm of RHS of the linear system */
    s->vector_norm(s->d_rhs, s->n, &rhs_norm);

    if (s->polish) return c_max(rhs_norm * CUDA_PCG_POLISH_ACCURACY, CUDA_PCG_EPS_MIN);

    switch (s->eps_strategy)
    {

    /* SCS strategy */
    case SCS_STRATEGY:
        eps = rhs_norm * s->start_tol  / pow((admm_iter + 1), s->dec_rate);
        eps = c_max(eps, CUDA_PCG_EPS_MIN);
        break;

    /* residual strategy */
    case RESIDUAL_STRATEGY:
        if (admm_iter == 1)
        {
            /* In case rhs = 0.0 we don't want to set eps_prev to 0.0 */
            if (rhs_norm < CUDA_PCG_EPS_MIN)
                s->eps_prev = 1.0;
            else
                s->eps_prev = rhs_norm * s->reduction_factor;
            /* Return early since scaled_pri_res and scaled_dua_res are meaningless before the first
    ADMM iteration */
            return s->eps_prev;
        }

        if (s->zero_pcg_iters >= s->reduction_threshold) {
            s->reduction_factor /= 2;
            s->zero_pcg_iters = 0;
        }

```

```
48          eps = s->reduction_factor * sqrt((*s->scaled_pri_res) * (*s->scaled_dua_res));
49          eps = c_max(c_min(eps, s->eps_prev), CUDA_PCG_EPS_MIN);
50          s->eps_prev = eps;
51          break;
52      }
53      return eps;
54  }
55
56  /* d_rhs = d_b1 + A' * rho * d_b2 */
57  static void compute_rhs(cudapcg_solver *s, c_float *d_b)
58  {
59
60      c_int n = s->n;
61      c_int m = s->m;
62
63      /* d_rhs = d_b1 */
64      cuda_vec_copy_d2d(s->d_rhs, d_b, n);
65
66      if (m == 0)
67          return;
68
69      /* d_z = d_b2 */
70      cuda_vec_copy_d2d(s->d_z, d_b + n, m);
71
72      if (!s->d_rho_vec)
73      {
74          /* d_z *= rho */
75          cuda_vec_mult_sc(s->d_z, *s->h_rho, m);
76      }
77      else
78      {
79          /* d_z = diag(d_rho_vec) * d_z */
80          cuda_vec_ew_prod(s->d_z, s->d_z, s->d_rho_vec, m);
81      }
82
83      /* d_rhs += A' * d_z */
84      cuda_mat_Axpy(s->At, s->d_z, s->d_rhs, 1.0, 1.0);
85  }
86
87  /* api fcns as perscribed by the osqp documentation */
88  c_int init_linsys_solver_cudapcg(cudapcg_solver **sp,
89                                   const OSQPMatrix *P,
90                                   const OSQPMatrix *A,
91                                   const OSQPVectorf *rho_vec,
92                                   OSQPSettings *settings,
93                                   c_float *scaled_pri_res,
94                                   c_float *scaled_dua_res,
95                                   c_int polish) {
96
97      c_int n, m;
98      c_float H_m_1 = -1.0;
```

```
99
100     /* allocate linsys solver structure */
101     cudapcg_solver *s = c_calloc(1, sizeof(cudapcg_solver));
102     *sp = s;
103
104     /* assign type and the number of threads */
105     s->type = CUDA_PCG_SOLVER;
106     s->nthreads = 1; // dummy, this changes for target device
107
108     /* dimensions */
109     n = OSQPMatrix_get_n(P);
110     m = OSQPMatrix_get_m(A);
111     s->n = n;
112     s->m = m;
113
114     /* pcg states */
115     s->polish = polish;
116     s->zero_pcg_iters = 0;
117
118     /* default norm and tolerance strategy */
119     s->eps_strategy = RESIDUAL_STRATEGY;
120     s->norm = CUDA_PCG_NORM;
121     s->precondition = CUDA_PCG_PRECONDITION;
122     s->warm_start_pcg = CUDA_PCG_WARM_START;
123     s->max_iter = (polish) ? CUDA_PCG_POLISH_MAX_ITER : CUDA_PCG_MAX_ITER;
124
125     /* tolerance strategy parameters */
126     s->start_tol = CUDA_PCG_START_TOL;
127     s->dec_rate = CUDA_PCG_DECAY_RATE;
128     s->reduction_threshold = CUDA_PCG_REDUCTION_THRESHOLD;
129     s->reduction_factor = CUDA_PCG_REDUCTION_FACTOR;
130     s->scaled_pri_res = scaled_pri_res;
131     s->scaled_dua_res = scaled_dua_res;
132
133     /* set pointers settings and data */
134     s->A = A->S;
135     s->At = A->At;
136     s->P = P->S;
137     s->d_P_diag_ind = P->d_P_diag_ind;
138     if (rho_vec)
139         s->d_rho_vec = rho_vec->d_val;
140     if (!polish)
141     {
142         s->h_sigma = &settings->sigma;
143         s->h_rho = &settings->rho;
144     }
145     else
146     {
147         s->h_sigma = &settings->delta;
148         s->h_rho = (c_float*) c_malloc(sizeof(c_float));
149         *s->h_rho = 1. / settings->delta;
```

```
150      }
151
152      /* allocate pcg iterates */
153      cuda_calloc((void **) &s->d_x, n * sizeof(c_float));
154      cuda_malloc((void **) &s->d_p, n * sizeof(c_float));
155      cuda_malloc((void **) &s->d_Kp, n * sizeof(c_float));
156      cuda_malloc((void **) &s->d_y, n * sizeof(c_float));
157      cuda_malloc((void **) &s->d_r, n * sizeof(c_float));
158      cuda_malloc((void **) &s->d_rhs, n * sizeof(c_float));
159
160      if (m != 0) cuda_malloc((void **) &s->d_z, m * sizeof(c_float));
161
162      /* allocate scalar in host memory that is page-locked and accessible to target */
163      cuda_malloc_host((void **) &s->h_r_norm, sizeof(c_float));
164
165      /* allocate target-side scalar values. This way scalars are packed in target memory */
166      cuda_malloc((void **) &s->d_r_norm, 8 * sizeof(c_float));
167      s->rTy = s->d_r_norm + 1;
168      s->rTy_prev = s->d_r_norm + 2;
169      s->alpha = s->d_r_norm + 3;
170      s->beta = s->d_r_norm + 4;
171      s->pKp = s->d_r_norm + 5;
172      s->D_MINUS_ONE = s->d_r_norm + 6;
173      s->d_sigma = s->d_r_norm + 7;
174      cuda_vec_copy_h2d(s->D_MINUS_ONE, &H_m_1, 1);
175      cuda_vec_copy_h2d(s->d_sigma, s->h_sigma, 1);
176
177      /* allocate memory for PCG preconditioning */
178      if (s->precondition)
179      {
180          cuda_malloc((void **) &s->d_P_diag_val, n * sizeof(c_float));
181          cuda_malloc((void **) &s->d_AtRA_diag_val, n * sizeof(c_float));
182          cuda_malloc((void **) &s->d_diag_precond, n * sizeof(c_float));
183          cuda_malloc((void **) &s->d_diag_precond_inv, n * sizeof(c_float));
184          if (!s->d_rho_vec) cuda_malloc((void **) &s->d_AtA_diag_val, n * sizeof(c_float));
185      }
186
187      /* Set the vector norm */
188      switch (s->norm) {
189      case 0:
190          s->vector_norm = &cuda_vec_norm_inf;
191          break;
192
193      case 2:
194          s->vector_norm = &cuda_vec_norm_2;
195          break;
196      }
197
198      s->solve = &solve_linsys_cudapcg;
199      s->warm_start = &warm_start_linsys_solver_cudapcg;
200      s->free = &free_linsys_solver_cudapcg;
```

```
201        s->update_matrices = &update_linsys_solver_matrices_cudapcg;
202        s->update_rho_vec  = &update_linsys_solver_rho_vec_cudapcg;
203
204        /* init PCG preconditioner */
205        if (s->precondition) cuda_pcg_update_precond(s, 1, 1, 1);
206
207        return 0;
208    }
209
210    /* main driver function for pcg */
211    c_int solve_linsys_cudapcg(cudapcg_solver *s, OSQPVectorf *b,c_int admm_iter)
212    {
213
214        c_int pcg_iter;
215        c_float eps;
216
217        /* compute the RHS of the reduced KKT system */
218        compute_rhs(s, b->d_val);
219
220        /* compute the required solution precision */
221        eps = compute_tolerance(s, admm_iter);
222
223        /* solve the linear system with PCG */
224        pcg_iter = cuda_pcg(s, eps, s->max_iter);
225
226        /* copy the first part of the solution to b->d_val */
227        cuda_vec_copy_d2d(b->d_val, s->d_x, s->n);
228
229        /* solution polishing */
230        if (!s->polish)
231        {
232            /* Compute d_z = A * d_x */
233            if (s->m) cuda_mat_Axpy(s->A, s->d_x, b->d_val + s->n, 1.0, 0.0);
234        }
235        else
236        {
237            /* Compute yred = (A * d_x - b) / delta */
238            cuda_mat_Axpy(s->A, s->d_x, b->d_val + s->n, 1.0, -1.0);
239            cuda_vec_mult_sc(b->d_val + s->n, *s->h_rho, s->m);
240        }
241
242        if (pcg_iter == 0)
243            s->zero_pcg_iters++;
244
245        return 0;
246    }
247
248    void warm_start_linsys_solver_cudapcg(cudapcg_solver *s, const OSQPVectorf *x)
249    {
250        cuda_vec_copy_d2d(s->d_x, x->d_val, x->length);
251    }
```

```
void free_linsys_solver_cudapcg(cudapcg_solver *s)
{

    if (s)
    {
        if (s->polish)
            c_free(s->h_rho);

        /* PCG iterates */
        cuda_free((void **) &s->d_x);
        cuda_free((void **) &s->d_p);
        cuda_free((void **) &s->d_Kp);
        cuda_free((void **) &s->d_y);
        cuda_free((void **) &s->d_r);
        cuda_free((void **) &s->d_rhs);
        cuda_free((void **) &s->d_z);

        /* free host memory */
        cuda_free_host((void **) &s->h_r_norm);

        /* target-side scalar values */
        cuda_free((void **) &s->d_r_norm);

        /* pcg preconditioner */
        cuda_free((void **) &s->d_P_diag_val);
        cuda_free((void **) &s->d_AtA_diag_val);
        cuda_free((void **) &s->d_AtRA_diag_val);
        cuda_free((void **) &s->d_diag_precond);
        cuda_free((void **) &s->d_diag_precond_inv);

        c_free(s);
    }
}

c_int update_linsys_solver_matrices_cudapcg(cudapcg_solver *s, const OSQPMatrix *P, const OSQPMatrix *A)
{

    if (s->precondition) cuda_pcg_update_precond(s, 1, 1, 0);
    return 0;
}

c_int update_linsys_solver_rho_vec_cudapcg(cudapcg_solver *s, const OSQPVectorf *rho_vec, c_float rho_sc
    )
{

    if (s->precondition) cuda_pcg_update_precond(s, 0, 0, 1);
    return 0;
}
```