Testing Fine-Grained Parallelism for the ADMM on a Factor-Graph Nate Derbinsky

Wentworth Institute of Technology



Ning Hao Oracle



AmirReza Oghbaee Northeastern



Mohammad Rostami UPenn



José Bento Boston College



The Problem

- Large-scale optimization problems
 - Many hard/soft constraints
 - Many discrete/continuous variables
 - Often not smooth/convex
- Diverse applications: vision/imaging, robotics, machine learning, ...
- Thus, we need tools that balance
 - Problem independence
 - Ease of use
 - Efficiency/scalability



Our Framework: parADMM

 General optimization via the Alternating Direction Method of Multipliers (ADMM)

State-of-the-art performance on numerous tasks

- Automatically exploit CPU/GPU parallelism given user's serial code
 - Written in C; integrates CUDA, OpenMP
- Empirical GPU/CPU speedup results
 - Packing
 - Optimal control
 - Machine learning (SVM)



Optimization Problem



Testing Fine-Grained Parallelism for the ADMM on a Factor-Graph

V V

Example: Packing

Fit *n* circles of radius *r* in a square of side-length *s* without overlap (non-convex, NP-hard, ∞ solutions)







Packing Objective

 $\min_{x_1, x_2, x_3} : Box(x_1) + Box(x_2) + Box(x_3) +$

 $Collision(x_1, x_2) + Collision(x_2, x_3) + Collision(x_1, x_3)$





ADMM

Alternating Direction Method of Multipliers [Boyd et al. '11]

General

- Arbitrary objective functions, constraints, and variables
- Global minimum for *convex* problems (and demonstrably successful for non-convex as well)
- If converges, produces a *feasible* solution (all hard constraints met)

Interruptible

 Iterative algorithm; intermediate results can serve as heuristic start for complementary approaches

Scalable and Parallelizable

 Formulated as a *decomposition-coordination* problem; leads naturally to concurrency at multiple levels (e.g. MapReduce, multi-core, GPU)



PCO | IPDPS 2016 | Nate Derbinsky

Message-Passing ADMM

[Derbinsky et al. '13]



Testing Fine-Grained Parallelism for the ADMM on a Factor-Graph

May 23, 2016

Example parADMM Program

```
graph Cpu graph;
                          // bipartite graph in the CPU
graph Cpu Gpu graph;
                         // interface between CPU and GPU graph
                         // bipartite graph in the GPU
graph* Gpu graph;
cudaMalloc( (void **) &Gpu_graph , sizeof(graph) );
startG( &Cpu graph, number of dims per edge ); // initialize empty CPU graph
// add nodes to the bipartite graph
int index of variables 1[] = {1,2,3}; int number of variables 1 = 3;
int index of variables_2[] = {1,4,5}; int number_of_variables_2 = 3;
int index_of_variables_3[] = {2,5};
                                       int number of variables 3 = 2;
                                       int number_of_variables_4 = 1;
int index of variables 4[] = {5};
addNode(&Cpu graph, proximal operator 1, (void *) parameters 1,
                         number_of_variables_1, index_of_variables_1);
   size parameters 1,
addNode(&Cpu graph, proximal operator 2, (void *) parameters 2,
   size_parameters_2,
                         number_of_variables_2, index_of_variables_2);
addNode(&Cpu graph, proximal operator 3, (void *) parameters 3,
   size parameters 3,
                         number of variables 3, index of variables 3);
addNode(&Cpu graph, proximal operator 4, (void *) parameters 4,
   size parameters 4,
                         number of variables 4, index of variables 4);
// set the rhos and alpha values (all equal)
initialize_RHOS_APHAS(&Cpu_graph, rho, alpha);
// initialize the ADMM variables (at random between lower and upper bound)
initialize_X_N_Z_M_U_rand(&Cpu_graph, lower_bound, upper_bound, lower_bound,
   upper bound, lower bound, upper bound, lower bound, upper bound, lower bound,
   upper_bound);
// initialize the Cpu Gpu graph using the GPU graph
copyGraphFromCPUtoGPU(Cpu_graph, &Cpu_Gpu_graph, Gpu_graph);
// iterate the ADMM
for (int i = 0; i< numiterations; i++)</pre>
{
   updateXGPU<<< ... , ... >>>(GPU graph);
   updateMGPU<<< ... , ... >>>(GPU graph);
   updateZGPU<<< ... , ... >>>(GPU graph);
   updateUGPU<<< ... , ... >>>(GPU graph);
   updateNGPU<<< ... , ... >>>(GPU graph);
}
// copy the variables Z from the GPU graph to the CPU graph
cudaMemcpy(Cpu graph.Z , Cpu Gpu graph.Z , CPU graph.num dims *
   CPU graph.num vars * sizeof(double),cudaMemcpyDeviceToHost);
. . .
```



Example Proximal Operator

```
__device__ double d_max_alpha_radius(double n, double rho, double alpha)
 6
     {
             const int newton_iterations = 10;
 8
 9
             double r = d_{abs}d(n / 2.0);
10
11
             for ( int i=0; i<newton_iterations; i++ )</pre>
12
             {
13
                     const double num = -pow( r, alpha-1 ) + ( rho / alpha )*( r - n );
                     const double denom = -( alpha - 1 ) * pow( r, alpha-2 ) + ( rho / alpha );
14
15
                     r -= num / denom;
16
17
                     if (r < 0)
18
19
                     {
20
                              return 0.0;
                     }
21
22
             }
23
24
             return r;
25
     }
26
27
      _device__ void d_maximize_sum_of_squared_disk_radii_alpha(double *x, double *n, double *rhos, int numvars, int numdims, void *params)
     {
28
             double* paramsCasted = (double*) params;
29
30
             for ( int i=0; i<numvars; i++ )</pre>
31
             {
32
33
                     double rho_i = rhos[i];
34
                     double rad_i = d_max_alpha_radius(n[ i*numdims ], rho_i, paramsCasted[0]);
35
                     for ( int j=0; j<numdims; j++ )</pre>
36
                     {
37
                              x[ i*numdims + j ] = rad_i;
38
                     }
39
40
             }
41 }
```



Example Packing Solution



Fine-Grained Hypothesis

- When applying ADMM, problem decomposition is user-defined
 - Typically involves relatively few & large [with parallelism within factors], which may be useful in MapReduce-like context
- Our hypothesis: efficient and automatic parallelization via many & small factors (1 thread per factor)
 - Bonus: the user needs only implement factor logic using <u>serial</u> code, typically quite simple



Empirical Evaluation

- AMD Opteron Abu Dhabi 6300 @ 2.8GHz
 - Up to 32 cores, 128GB memory
 - pragma openmp parallel for
- NVIDIA Tesla K40
- Ubuntu
- Numerical results in <u>three</u> domains
 - Focus: speedup with problem size
 - No change to factor code
 - Minimal change to graph-creation code



Packing Circles in a Triangle



Testing Fine-Grained Parallelism for the ADMM on a Factor-Graph

*** ***

Phase Breakdown



Testing Fine-Grained Parallelism for the ADMM on a Factor-Graph

V V

Model Predictive Control (MPC)



Testing Fine-Grained Parallelism for the ADMM on a Factor-Graph

VVV

Binary Classification via SVM in \mathbb{R}^2



Testing Fine-Grained Parallelism for the ADMM on a Factor-Graph

V V V

Conclusions

- We achieve problem-independent scaling with user-supplied **serial** code
 - Multi-core CPU (5-6x with 32-cores)
 - GPU (10-18x; comparable with other GPUaccelerated libraries)
- Future work
 - Improved work-scheduling/topology
 - Multiple GPU/computer, asynchronous ADMM

VVV

Additional Details in the Paper

- Classic ADMM formulation
- parADMM internals

- Includes CUDA, OpenMP parameterization

- Survey of related tools/frameworks
- Empirical evaluation
 - Problem formulation
 - Analysis of results



The End :) Questions?



Ning Hao Oracle



AmirReza Oghbaee Northeastern



Mohammad Rostami UPenn



José Bento Boston College

https://github.com/parADMM/engine



Testing Fine-Grained Parallelism for the ADMM on a Factor-Graph

May 23, 2016