

RPC 2020



Virtual Research Presentation Conference

GPU Computing for Fast, Massively Parallel Simulations of Manufacturing Processes

Principal Investigator: Richard Otis (357)
Co-Is: Hamsa Shwetha Venkataram (1740)
Program: Spontaneous Concept



RPC-127



Jet Propulsion Laboratory
California Institute of Technology

Tutorial Introduction

Abstract

JPL is increasingly relying on manufacturing simulation software, including for additive manufacturing process design. As our manufacturing simulation capability scales up to whole parts, the demand for parallel processing will increase faster than what our present algorithms can support. The objective of this work was to explore moving certain types of engineering computations to the GPU, for the purpose of performing billions of computations in parallel with the potential upside being multiple orders-of-magnitude improvement in calculation times.

Empirical evidence was found during this study that the orders-of-magnitude (**100x+**) speedups enabled by GPUs, reported in the literature, are possible to achieve on problems relevant to JPL and NASA. The speedups enabled by GPU-aware simulation codes may constitute a game-changing technology opportunity for certain classes of physical simulation problems.

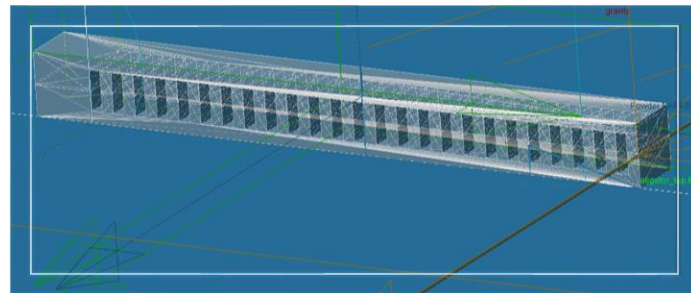
Problem Description

JPL is experiencing computational “growing pains” relating to our increasing infusion of advanced manufacturing technology, particularly 3D printing. Internal customers are demanding results from more advanced manufacturing processes with less iteration time.

SOA process modeling codes do not provide the just-in-time accuracy needed to deliver quick-turn solutions, resulting in time-consuming experimental trials.

GPU-based algorithms have been reported as an approach for achieving orders-of-magnitude speedups in performance for certain classes of problems. Within the last 12 months the software tools for programming GPUs have matured significantly in terms of accessibility to non-experts.

If the reported speedups can be replicated in modeling codes used by NASA and JPL, this would mean running high-fidelity process models, which currently take *weeks* or *months*, in a matter of hours.



Dr. Cornelia Altenbuchner/JPL/CIF

```
# Each thread computes one element in the result matrix.
# The dot product is chunked into dot products of TPB-long vectors.
tmp = 0.
for i in range(bpg):
    # Preload data into shared memory
    sA[tx, ty] = A[x, ty + i * TPB]
    sB[tx, ty] = B[tx + i * TPB, y]

    # Wait until all threads finish preloading
    cuda.syncthreads()

    # Computes partial product on the shared memory
    for j in range(TPB):
        tmp += sA[tx, j] * sB[j, ty]

    # Wait until all threads finish computing
    cuda.syncthreads()
```

Methodology

Several GPU-enabled array computing packages were considered during this study, including:

- Tensorflow: <https://tensorflow.org/>
- JAX: <https://jax.readthedocs.io/en/latest/>
- Numba: <http://numba.pydata.org/>
- CuPy: <https://cupy.dev/>
- Scikit-cuda: <https://scikit-cuda.readthedocs.io/en/latest/>
- CUDA.jl: <https://juliagpu.org/>

Array arithmetic (addition, multiplication) and linear algebra capabilities such as matrix inversion were also investigated for each of the packages.

```

JAX

In [8]: npX = np.load('input_X.npy')
npGM = np.load('input_GM.npy')
npchemical_potentials = np.array([-1000, -1000, -1000, -1000, -1000])

In [12]: def driving_force_jax(inp_x, inp_gm, inp_mu):
return jnp.dot(inp_x, inp_mu).block_until_ready() - inp_gm

## warm cache
driving_force_jax(npX, npGM, npchemical_potentials)

Out[12]: DeviceArray([[[[ 44967.168 ,  34403.746 ,  35703.016 , ...,
-3631.798 , -6328.6255, -24201.709 ]]]], dtype=float32)

In [14]: npX = device_put(npX)
npGM = device_put(npGM)

In [2]: import tensorflow as tf
import numpy as np
print("Num GPUs Available: ", len(tf.config.experimental.list_physical_devices('GPU')))

Out[14]: Num GPUs Available: 1

In [3]: with tf.device('/CPU:0'):
X = tf.constant(np.load('input_X.npy'), dtype=tf.float32)
GM = tf.constant(np.load('input_GM.npy'), dtype=tf.float32)
chemical_potentials = tf.constant([-1000, -1000, -1000, -1000, -1000], dtype=tf.float32)

@tf.function(experimental_compile=False)
def driving_force_cpu(inp_x, inp_gm, inp_mu):
return tf.tensordot(inp_x, inp_mu, axes=(4,0)) - inp_gm

In [4]: # Warm cache
driving_force_cpu(X, GM, chemical_potentials)

Out[4]: <tf.Tensor: shape=(1, 1, 1, 26461866), dtype=float32, numpy=
array([[[[ 44967.168 ,  34403.746 ,  35703.016 , ..., -3631.798 ,
-6328.6255, -24201.709 ]]]], dtype=float32)>

In [5]: %time driving_force_cpu(X, GM, chemical_potentials)
CPU times: user 110 ms, sys: 7.02 ms, total: 117 ms
Wall time: 113 ms

Out[5]: <tf.Tensor: shape=(1, 1, 1, 26461866), dtype=float32, numpy=
array([[[[ 44967.168 ,  34403.746 ,  35703.016 , ..., -3631.798 ,
-6328.6255, -24201.709 ]]]], dtype=float32)>

In [6]: with tf.device('/CPU:0'):
X = tf.constant(np.load('input_X.npy'), dtype=tf.float32)
GM = tf.constant(np.load('input_GM.npy'), dtype=tf.float32)
chemical_potentials = tf.constant([-1000, -1000, -1000, -1000, -1000], dtype=tf.float32)

@tf.function(experimental_compile=True)
def driving_force_cpu_xla(inp_x, inp_gm, inp_mu):
return tf.tensordot(inp_x, inp_mu, axes=(4,0)) - inp_gm

In [7]: # Warm cache
driving_force_cpu_xla(X, GM, chemical_potentials)

Out[7]: <tf.Tensor: shape=(1, 1, 1, 26461866), dtype=float32, numpy=
array([[[[ 44967.168 ,  34403.746 ,  35703.016 , ..., -3631.798 ,
-6328.6255, -24201.709 ]]]], dtype=float32)>

In [8]: %time driving_force_cpu_xla(X, GM, chemical_potentials)
  
```

Results

Goals: Demonstrate 100x+ speedup from GPU on a JPL-relevant simulation problem

Accomplishments

- Demonstrated **3-4x** speedups from GPU at JPL
- Observed **500x** speedup due to GPU by a NASA-funded University team for a relevant simulation

Next Steps: Pursue the solution of a more focused problem using the Numba toolkit to achieve the observed speedups in a “real world” scenario.

The Julia programming language is also an interesting option for new simulation software packages.

```
59 function hyperplane(compositions::MatrixType,  
60     energies::FloatVector, target_composition::FloatVector,  
61     chemical_potentials::FloatVector, total_moles::FloatType,  
62     fixed_chempot_indices::IndexVector, fixed_comp_indices::IndexVector,  
63     result_fractions::FloatVector, result_simplex::IndexVector)  
64     num_components = size(compositions)[2]  
65     num_points = length(energies)  
66     num_fixed_chempots = size(fixed_chempot_indices)[1]  
67     simplex_size = num_components - num_fixed_chempots  
68     # composition index of -1 indicates total number of moles, i.e., N=1 condition  
69     included_composition_indices = fixed_comp_indices  
70     best_guess_simplex = sort(setdiff!(collect(1:num_components), fixed_chempot_indices))  
71     free_chempot_indices = best_guess_simplex[:]  
72     candidate_simplex = best_guess_simplex[:]  
73     trial_simplices = Array{IntType, 2}(undef, simplex_size, simplex_size)  
74     fractions = Array{FloatType, 2}(undef, simplex_size, simplex_size)  
75     driving_forces = Array{FloatType, 1}(undef, num_points)  
76     for i in 1:simplex_size  
77         trial_simplices[i, :] = best_guess_simplex  
78     end  
79     trial_matrix = Array{FloatType, 3}(undef, simplex_size, simplex_size, simplex_size)  
80     candidate_tieline = Array{FloatType, 2}(undef, simplex_size, simplex_size)  
81     candidate_energies = Array{FloatType, 1}(undef, simplex_size)  
82     candidate_potentials = Array{FloatType, 1}(undef, simplex_size)  
83     smallest_fractions = Array{FloatType, 1}(undef, simplex_size)  
84     saved_trial = 0  
85  
86     max_iterations = 1000  
87     for _ in 1:max_iterations  
88         for trial_idx in 1:simplex_size  
89             for comp_idx in 1:simplex_size  
90                 ici = included_composition_indices[comp_idx]  
91                 for simplex_idx in 1:simplex_size  
92                     if ici > 0  
93                         trial_matrix[comp_idx, simplex_idx, trial_idx] = compositions[trial_simplices[trial_idx, simplex_idx], ici]  
94                     else  
95                         # ici = -1, refers to N=1 condition  
96                         trial_matrix[comp_idx, simplex_idx, trial_idx] = 1 # 1 mole-formula per formula unit of a phase  
97                     end # if  
98                 end # for  
99             end # for  
100         end # for
```

Publications and References

<https://rapids.ai/>

<https://cloud.google.com/tpu/docs/tpus>

<https://hackage.haskell.org/package/accelerate>

<http://numba.pydata.org/>

<https://scikit-cuda.readthedocs.io/en/latest/>

<https://github.com/AdditiveModeling/pyphasefield/blob/a089672a9659b36a888f2be22c913d2235bd3a30/pyphasefield/pyphasefield/Engines/NCGPU.py>