# Open-source Tools For GPU Programming in Large Classrooms

**Abdul Dakkak, Carl Pearson, Cheng Li**
{dakkak,pearson,cli99}@illinois.edu

# WebGPU

Description   **Code**   Questions   Attempts   History

## Machine Problem Code (Past Deadline)

Compile & Run ▾

```
1   #include <wb.h>
2
3   #define wbCheck(stmt)                                                    \
4     do {                                                                   \
5       cudaError_t err = stmt;                                             \
6       if (err != cudaSuccess) {                                          \
7         wbLog(ERROR, "Failed to run stmt ", #stmt);                      \
8         wbLog(ERROR, "Got CUDA error ... ", cudaGetErrorString(err));    \
9         return -1;                                                        \
10      }                                                                   \
11    } while (0)
12
13  /// For simplicity, fix #bins=1024 so scan can use a single block and no padding
14  #define NUM_BINS 1024
15
16  /********************************************************************************
17   GPU main computation kernels
18   *******************************************************************************/
19
20  __global__ void gpu_normal_kernel(float *in_val, float *in_pos, float *out,
21                                    int grid_size, int num_in) {
22
23    //@@ INSERT CODE HERE
24
25    int outIdx = blockIdx.x * blockDim.x + threadIdx.x;
26
27    if (outIdx < grid_size) { // Boundary check
28
29      // Local accumulator
30      float out_reg = 0.0f;
31
32      // Loop over input elements and compute
33      for (int inIdx = 0; inIdx < num_in; ++inIdx) {
34        const float in_val_reg = in_val[inIdx];
35        const float dist = in_pos[inIdx] - (float)outIdx;
36        out_reg += (in_val_reg * in_val_reg) / (dist * dist);
37      }
38
39      // Commit final result
40      out[outIdx] += out_reg;
```

# Originally Designed for MOOC

➜ Around 100k students registered for Coursera's Heterogeneous Parallel Programming course
➜ Targeted weekly labs
➜ Labs auto-graded based with dataset

# Intro to CUDA

Around 200 students from UIUC
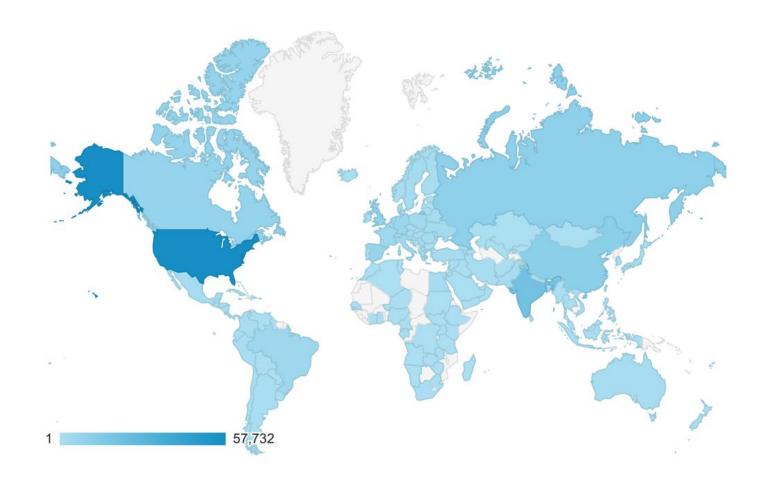
# Advanced CUDA

Around 100 students for UIUC and collaborating institutions

# Summer School

Around 100 students from all over the world

# Coursera HPP

Around 20,000 students worldwide

Students Per Offering

1                                    57,732

# Problem

# Restrictions with WebGPU

➔ Cannot modify programming environment
  ◆ Build scripts / libraries / dataset / …
  ◆ Cannot use profilers and debuggers
➔ User restricted within a sandboxed environment

# Intro and Advanced CUDA Project

➜ Develop a CUDA version of a CNN
➜ Given unoptimized sequential code
➜ Significant part of the total grade
➜ Around 4-6 weeks to complete
➜ Users should be "*root*"
➜ github.com/webgpu/ece408project
➜ github.com/webgpu/ece508-convlayer

# Pipeline

# Jupyter Notebook Interface to RAI

➔ Make it easy to develop interactive labs
➔ Built on top of Jupyter
➔ Implements a client/server that speaks the IPython protocol

# Command line Interface

```
1   rai:
2     version: 0.2 # this is required
3     #   image: gcc:6.3.0
4     image: ppc64le/gcc
5   resources:
6     cpu:
7       architecture: ppc64le
8     network: false
9     # gpu:
10    #   count: 1
11  commands:
12    build:
13      - echo "Building project"
14      - gcc /src/main.c
15      - ./a.out
16
```

Submission Spec

```
* Checking your athentication credentials.
* Preparing your project directory for upload.
* Uploading your project directory. This may take a few minutes.
 358 B / 358 B [████████████████████████████]          100.00% 5.23 KiB/s 0s
* Folder uploaded. Server is now processing your submission.
* Your job request has been posted to the queue.
* Server has accepted your job submission and started to configure the container.
* Downloading your code.
* Using ppc64le/gcc as container image.
* Starting container.
* Running echo "Building project"
Building project
* Running gcc /src/main.c
* Running ./a.out
Hello Universe!!
* * The build folder has been uploaded to http://s3.amazonaws.com/files.rai-project.com/userdata/buil
d-377d8ae0-64da-441c-80fb-bff5e717e13f.tar.tar.gz. The data will be present for only a short duration
 of time.
* Server has ended your request.
```

```
1
2   #include "stdio.h"
3
4   int main() {
5       printf("Hello Universe!!\n");
6       return 0;
7   }
```

User Program

https://asciinema.org/a/6k5e96itnqu6ekbji60c3kgy4
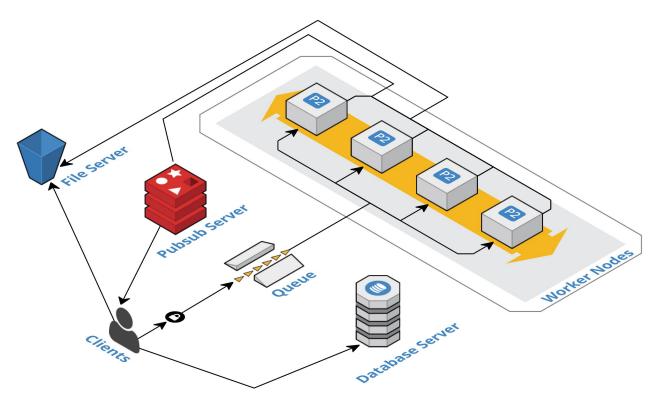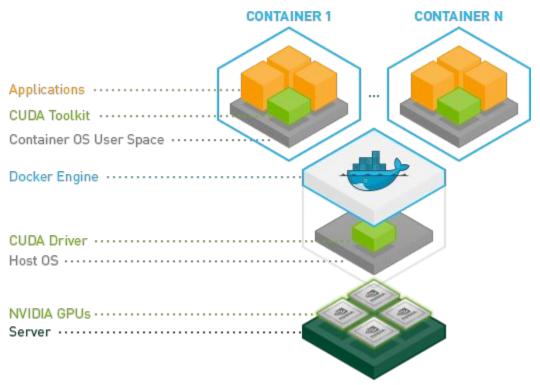
# Demo

# Architecture

# Current Deployment Setup

# Docker Layer



Wrote our own docker volume plugin

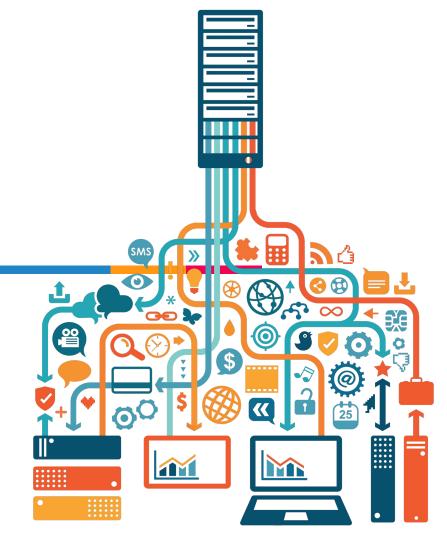# Not Just Project Submission

▷ A set of reusable components serving as a runtime

▷ Submission specific code is contained and small (<2KLoc)
  ○ Client logic is ~400 lines of code
  ○ Server logic is ~800 lines of code

| Service | Available Backends |
| --- | --- |
| Authentication | Secret, **Auth0** |
| Queue | NSQ, **SQS**, Redis, Kafka, NATS |
| Database | RethinkDB, MongoDB, MySQL, Postgres, SQLite, … |
| Registry | Etcd, Consul, BoltDB, Zookeeper |
| Config | **Yaml**, Toml, JSON, Environment |
| PubSub | EC, **Redis**, GCP, NATS, SNS |
| Tracing | XRay, Zipkin, StackDriver |
| Logger | **StackDriver**, **JournalD**, Syslog, Kinesis |
| Store | **S3**, Minio |
| Container | **Docker** |
| Serializer | BSON, **JSON** |

IMPACT

# Usage / Pedigree from Last Semester

➜ Around 170 students had to use the system for submission
➜ Students were using Linux, OSX, Windows, and WLS
➜ Students uploaded and generated around 100GB of data



**Used 25 Workers**

# Currently

➜ Running on the 2 IBM Minsky machines
➜ Used by around 100 people in the 508 class (UIUC and Minnesota)
   ◆ For the last lab
   ◆ For open-ended projects
➜ Students developed their own containers solving anything from Matrix factorization (for recommender systems) to Molecular simulations

# CarML

# CarML - Deploy ML Artifacts w/RAI

➜ Make it easy to deploy ML artifacts
➜ Makes it possible for people to test tools / ML models without investing time in installing software dependencies and getting HW resources

# Resources

# GPU TEACHING KIT FOR ACCELERATED COMPUTING
## Breaking the Barriers to GPU Education in Academia
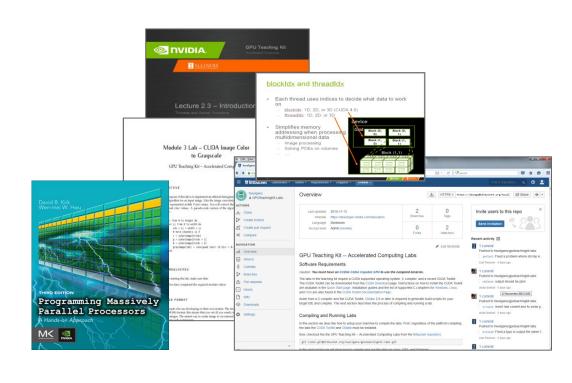
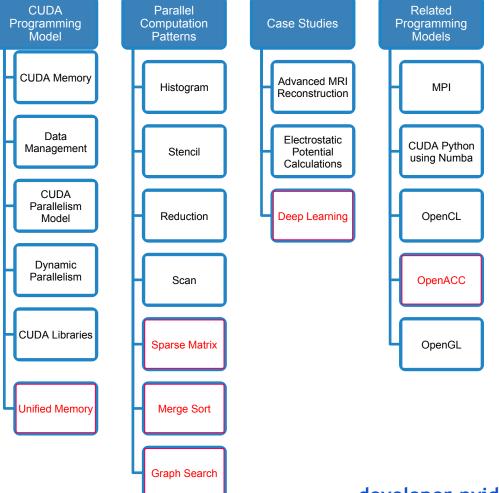Co-developed by UIUC and NVIDIA for educators

Comprehensive teaching materials

3rd Ed. PMPP E-book by Hwu/Kirk
Lecture slides and notes
Lecture videos
Hands-on labs/solutions
Larger coding projects/solutions
Quiz/exam questions/solution

GPU compute resources

NVIDIA online free Qwiklab credits
AWS credits

**developer.nvidia.com/teaching-kits**

| CUDA Programming Model | Parallel Computation Patterns | Case Studies | Related Programming Models |
|---|---|---|---|
| CUDA Memory | Histogram | Advanced MRI Reconstruction | MPI |
| Data Management | Stencil | Electrostatic Potential Calculations | CUDA Python using Numba |
| CUDA Parallelism Model | Reduction | Deep Learning | OpenCL |
| Dynamic Parallelism | Scan | | OpenACC |
| CUDA Libraries | Sparse Matrix | | OpenGL |
| Unified Memory | Merge Sort | | |
| | Graph Search | | |

**developer.nvidia.com/teaching-kits**

# Questions, Criticisms, and Concerns?

# Thank you

**Abdul Dakkak, Carl Pearson, Cheng Li**

{dakkak,pearson,cli99}@illinois.edu