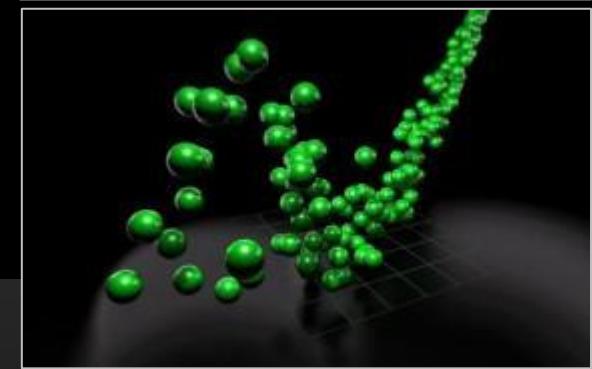
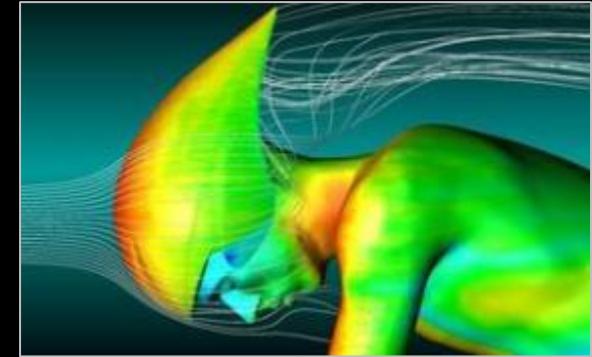


Gear Up for Speedup

Delivering 4x Speedups with ArrayFire

Julia Levites – NVIDIA

John Melonakos – AccelerEyes



Introduction

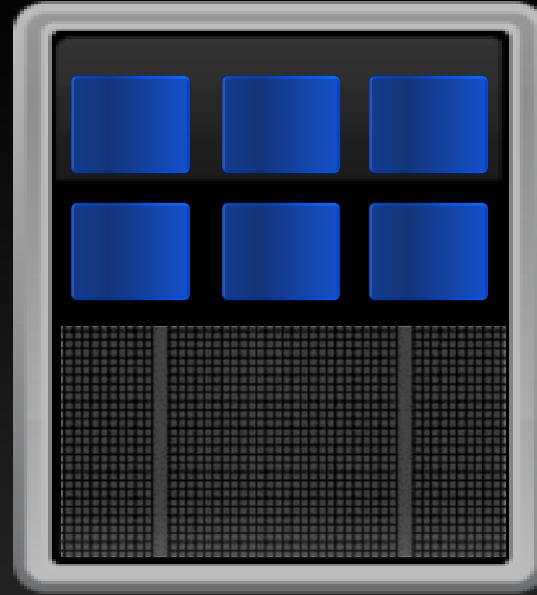
- What is GPU computing?
- GPU computing with ArrayFire

Gear Up for Speedup!

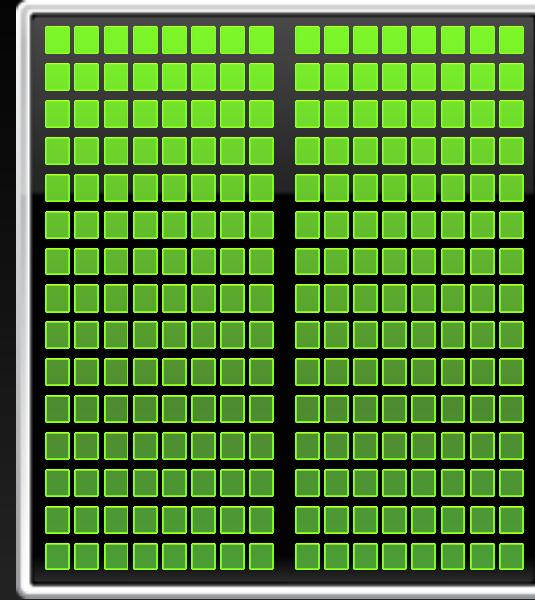
Experience 4X speed up with ArrayFire and GPUs



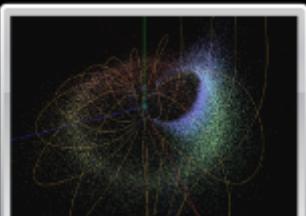
CPU



GPU



Add GPUs: Accelerate x86 Applications



100X

Astrophysics
RIKEN



130X

Quantum Chemistry
U of Illinois, Urbana



25X

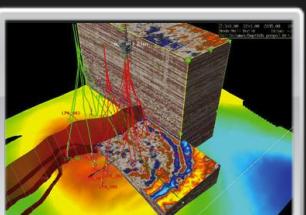
Fast Flow Simulation
FluiDyna



60X

Geospatial Imagery
PCI Geomatics

GPUs Accelerate HPC



12X

Kirchoff Time Migration
Acceleware



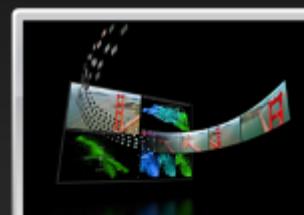
30X

Gene Sequencing
U of Maryland



149X

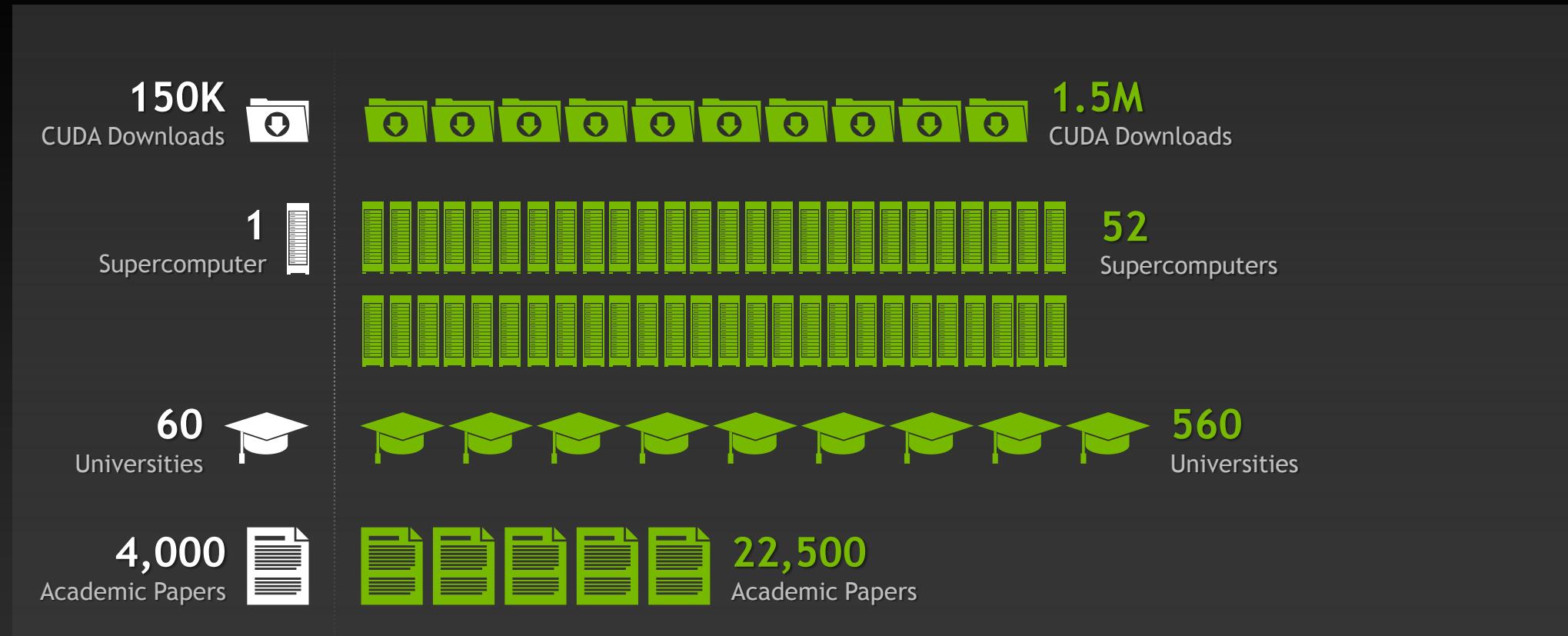
Financial Simulation
Oxford



18X

Video Transcoding
Elemental Tech

Explosive Growth of GPU Computing



2008

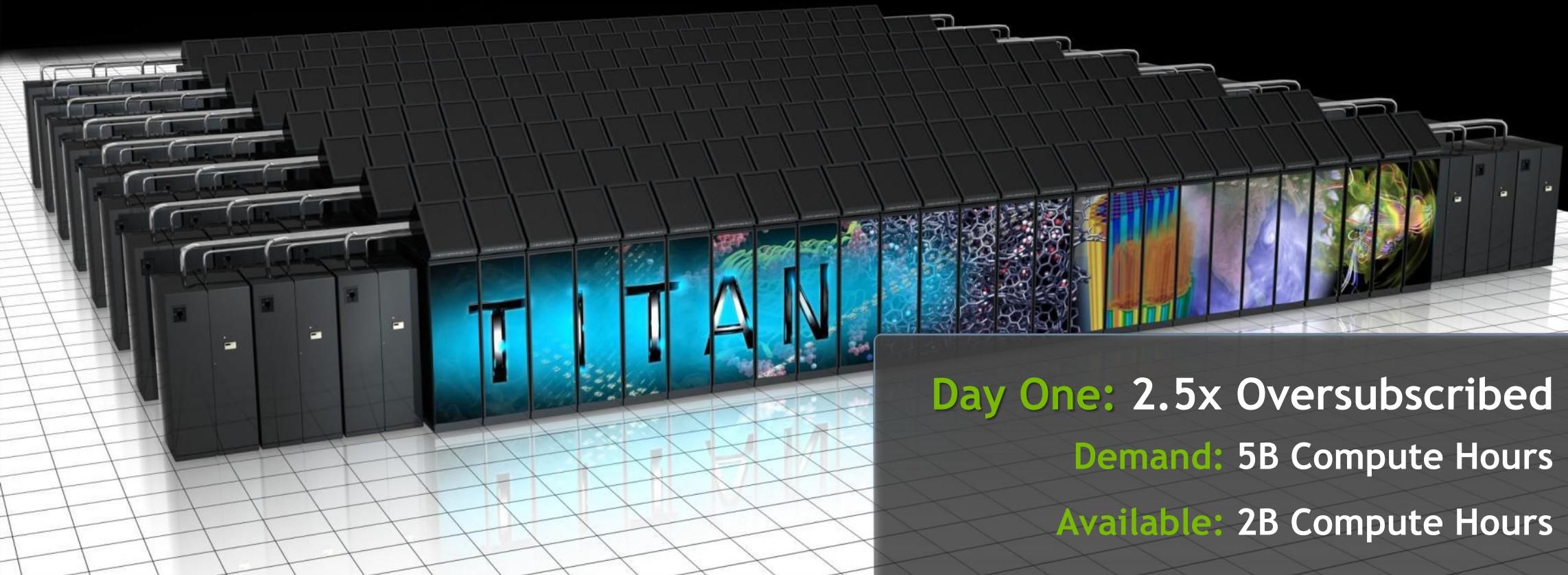
2012

Titan: World's #1 Open Science Supercomputer

18,688 Tesla Kepler GPUs

Fully GPU Accelerated Computing System

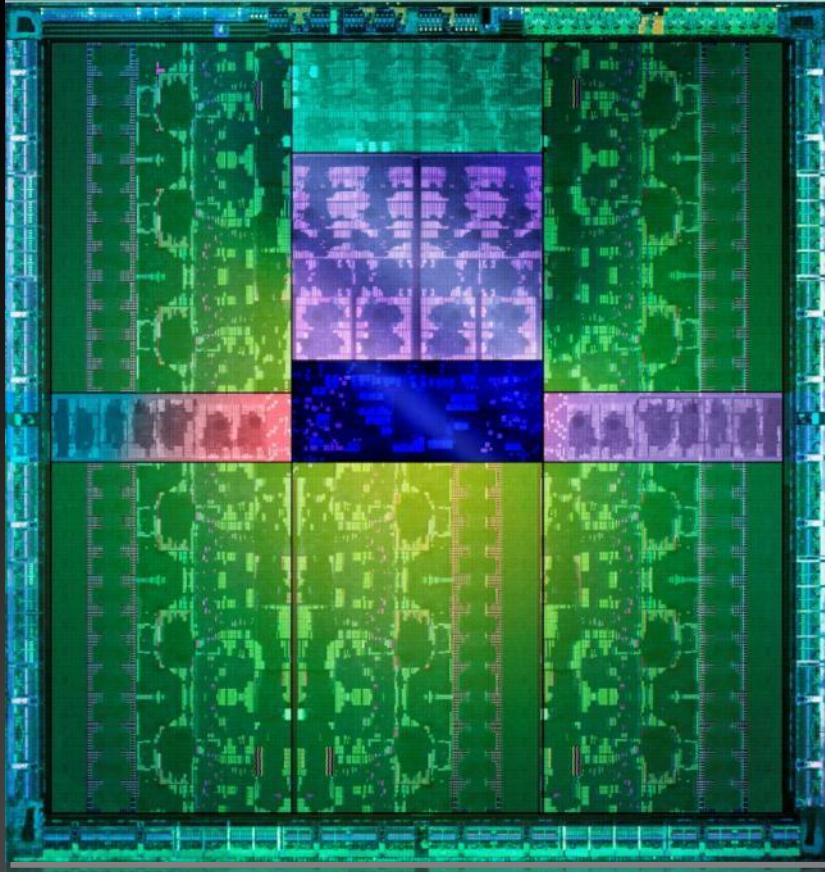
20+ Petaflops: 90% of Performance from GPUs



Day One: 2.5x Oversubscribed

Demand: 5B Compute Hours

Available: 2B Compute Hours



KEPLER

THE WORLD'S FASTEST, MOST
EFFICIENT HPC ACCELERATOR

- 3x Perf/Watt
- Max Utilization
- Easier to Program

Test Drive Kepler with ArrayFire

How can Gear Up help me?



Tonči Jukić

@toncijukic



 Follow

Today we have completed #ArrayFire implementation for our GPGPU farm and ported over 30% of code within only 3 days.
accelereyes.com/products/array...

 Reply

 Retweeted

 Favorite

{ 🔥 } ArrayFire

Fast, simple & comprehensive GPU library

- Hand-optimized GPU functions – Top speeds
- C, C++, Fortran
- CUDA 5.0 and Kepler K20 support
- v2.0 (coming this month)

accelereyes.com/arrayfire

estimate π

estimate π

```
#include <stdio.h>
#include <arrayfire.h>
using namespace af;
int main() {
    // 20 million random samples
    int n = 20e6;
    array x = randu(n,1), y = randu(n,1);
    // how many fell inside unit circle?
    float pi = 4 * sum<float>(sqrt(x*x+y*y)<1) / n;
    printf("pi = %g\n", pi);
    return 0;
}
```

estimate π

```
#include <stdio.h>
#include <arrayfire.h>
using namespace af;

int main() {
    // 20 million random samples
    int n = 20e6;
    array x = randu(n,1), y = randu(n,1);
    // how many fell inside unit circle?
    float pi = 4 * sum<float>(sqrt(x*x+y*y)<1) / n;
    printf("pi = %g\n", pi);
    return 0;
}
```

estimate π

```
#include <stdio.h>
#include <arrayfire.h>
using namespace af;
int main() {
    // 20 million random samples
    int n = 20e6;
    array x = randu(n,1), y = randu(n,1);
    // how many fell inside unit circle?
    float pi = 4 * sum<float>(sqrt(x*x+y*y)<1) / n;
    printf("pi = %g\n", pi);
    return 0;
}
```

estimate π

```
#include <stdio.h>
#include <arrayfire.h>
using namespace af;
int main() {
    // 20 million random samples
    int n = 20e6;
    array x = randu(n,1), y = randu(n,1);
    // how many fell inside unit circle?
    float pi = 4 * sum<float>(sqrt(x*x+y*y)<1) / n;
    printf("pi = %g\n", pi);
    return 0;
}
```

Data types

f32

real single precision

c32

complex single precision

c64

complex double precision

array

container object

b8

boolean byte

u32

unsigned integer

s32

signed integer

```
array x = randu(n, f32);  
array y = randu(n, f64);  
array z = rand(n, u32);
```

f64

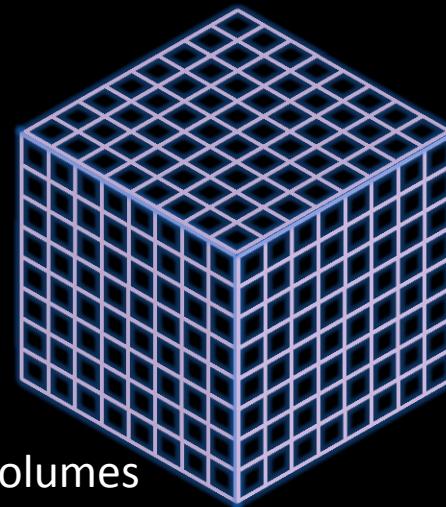
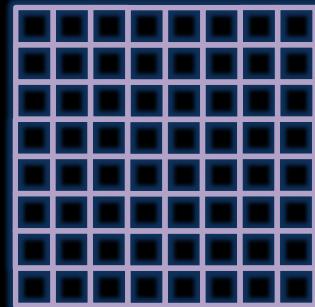
real double precision

Up to 4D Support



vectors

matrices



volumes

... 4D

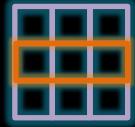
Subscripting

*ArrayFire Keywords: **end**, **span***

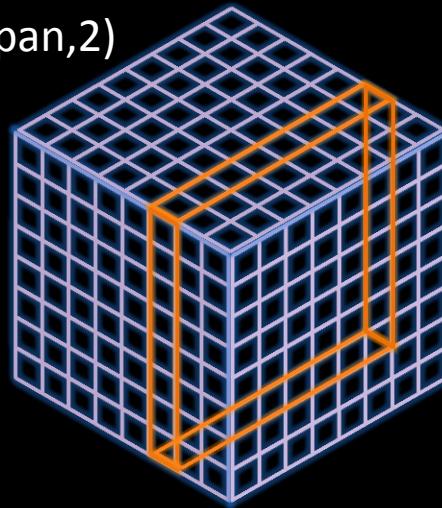
A(1,1)



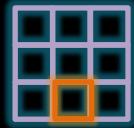
A(1,span)



A(span,span,2)

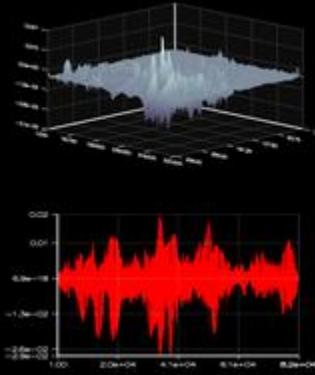
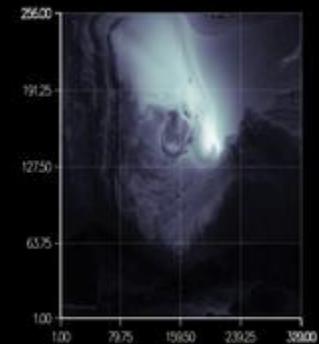
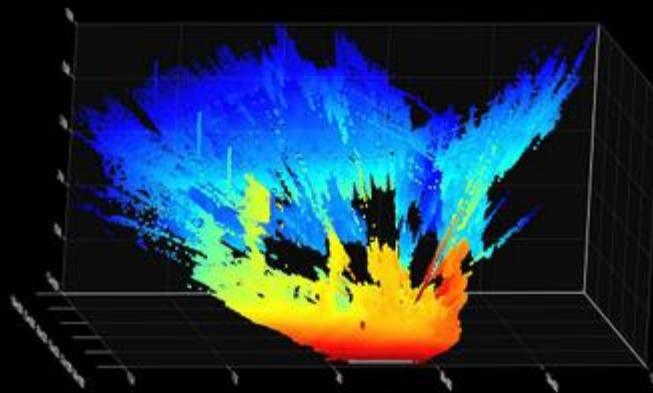


A(end,1)

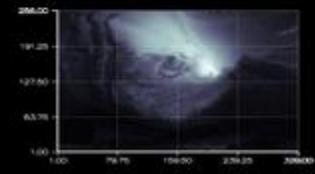
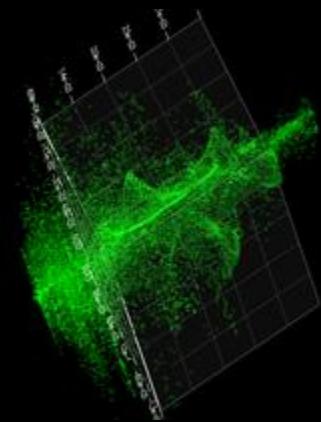
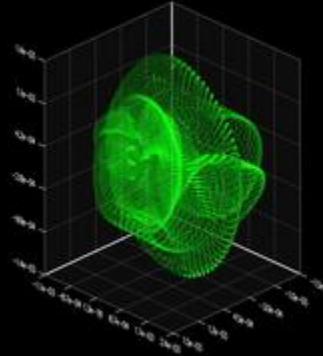
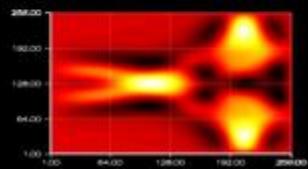
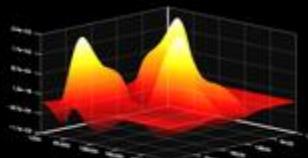


A(end,span)





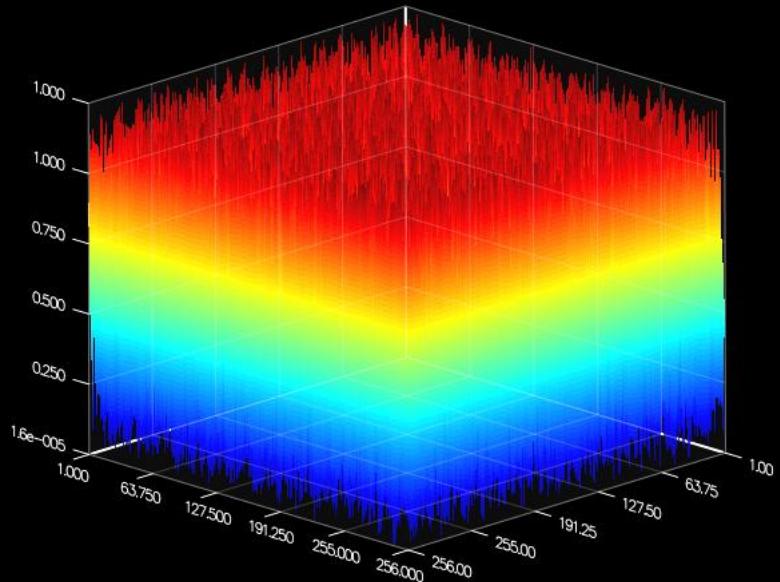
viz



Compute + Render

- asynchronous
- non-blocking
- throttled at 35 Hz

```
#include <arrayfire.h>
using namespace af;
int main() {
    // Infinite number of random 3d surfaces
    const int n = 256;
    while (1) {
        array x = randu(n,n);
        // 3d surface plot
        surface(x);
    }
    return 0;
}
```



GFOR -- data parallel for-loops

Example: Matrix Multiply

Serial matrix multiplications (3 kernel launches)

```
for (i = 0; i < 3; i++)
    C(span,span,i) = A(span,span,i) * B;
```

iteration i = 1

$$C(:,:,1) = A(:,:,1) * B$$

Example: Matrix Multiply

Serial matrix multiplications (3 kernel launches)

```
for (i = 0; i < 3; i++)
    C(span,span,i) = A(span,span,i) * B;
```

iteration i = 1

$$C(:,:,1) = A(:,:,1) * B$$

iteration i = 2

$$C(:,:,2) = A(:,:,2) * B$$

Example: Matrix Multiply

Serial matrix multiplications (3 kernel launches)

```
for (i = 0; i < 3; i++)
    C(span,span,i) = A(span,span,i) * B;
```

iteration i = 1

$$C(:,:,1) = A(:,:,1) * B$$

iteration i = 2

$$C(:,:,2) = A(:,:,2) * B$$

iteration i = 3

$$C(:,:,3) = A(:,:,3) * B$$

GFOR: data parallel for-loop

- Good for lots of small problems

Serial matrix multiplications (3 kernel launches)

```
for (i = 0; i < 3; i++)
    C(span,span,i) = A(span,span,i) * B;
```

Parallel matrix multiplications (1 kernel launch)

```
gfor (array i, 3)
    C(span,span,i) = A(span,span,i) * B;
```

Example: Matrix Multiply

Parallel matrix multiplications (1 kernel launch)

```
gfor (array i, 3)
    C(span, span, i) = A(span, span, i) * B;
```

simultaneous iterations i = 1:3

$$\begin{matrix} \text{C}(:,:,1) \\ \text{A}(:,:,1) \\ \text{B} \end{matrix} = \begin{matrix} \text{C}(:,:,1) \\ \text{A}(:,:,1) \\ \text{B} \end{matrix} * \begin{matrix} \text{C}(:,:,1) \\ \text{A}(:,:,1) \\ \text{B} \end{matrix}$$

$$\begin{matrix} \text{C}(:,:,2) \\ \text{A}(:,:,2) \\ \text{B} \end{matrix} = \begin{matrix} \text{C}(:,:,2) \\ \text{A}(:,:,2) \\ \text{B} \end{matrix} * \begin{matrix} \text{C}(:,:,2) \\ \text{A}(:,:,2) \\ \text{B} \end{matrix}$$

$$\begin{matrix} \text{C}(:,:,3) \\ \text{A}(:,:,3) \\ \text{B} \end{matrix} = \begin{matrix} \text{C}(:,:,3) \\ \text{A}(:,:,3) \\ \text{B} \end{matrix} * \begin{matrix} \text{C}(:,:,3) \\ \text{A}(:,:,3) \\ \text{B} \end{matrix}$$

Example: Matrix Multiply

Parallel matrix multiplications (1 kernel launch)

```
gfor (array i, 3)
    C(span, span, i) = A(span, span, i) * B;
```

Think of GFOR as compiling 1 stacked kernel with all iterations.

simultaneous iterations i = 1:3

$$\begin{matrix} \text{C}(,,1:3) \\ \text{A}(,,1:3) \\ \text{B} \end{matrix} = \begin{matrix} \text{C}(,,1:3) \\ \text{A}(,,1:3) \\ \text{B} \end{matrix} * \begin{matrix} \text{C}(,,1:3) \\ \text{A}(,,1:3) \\ \text{B} \end{matrix} *$$

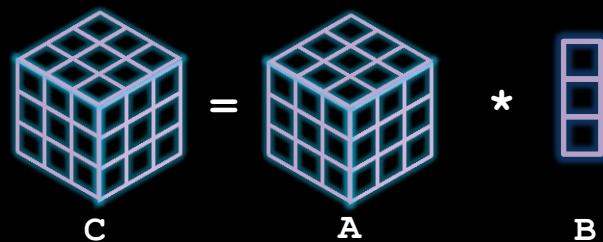
Example: Matrix Multiply

Parallel matrix multiplications (1 kernel launch)

```
gfor (array i, 3)
    C(span, span, i) = A(span, span, i) * B;
```

Think of GFOR as compiling 1 stacked kernel with all iterations.

simultaneous iterations i = 1:3

$$\text{C} = \text{A} * \text{B}$$


Easy Multi GPU Scaling

```
array *y = new array[n];
for (int i = 0; i < n; ++i) {
    deviceset(i);                      // change GPUs
    array x = randu(5,5);              // add work to GPU's queue
    y[i] = fft(x);                    // more work in queue
}

// all GPUs are now computing simultaneously
```

Hundreds of Functions...

reductions

- sum, min, max, count, prod
- vectors, columns, rows, etc

dense linear algebra

- LU, QR, Cholesky, SVD, Eigenvalues, Inversion, Solvers, Determinant, Matrix Power

convolutions

- 2D, 3D, ND

FFTs

- 2D, 3D, ND

image processing

- filter, rotate, erode, dilate, morph, resize, rgb2gray, histograms

interpolate & scale

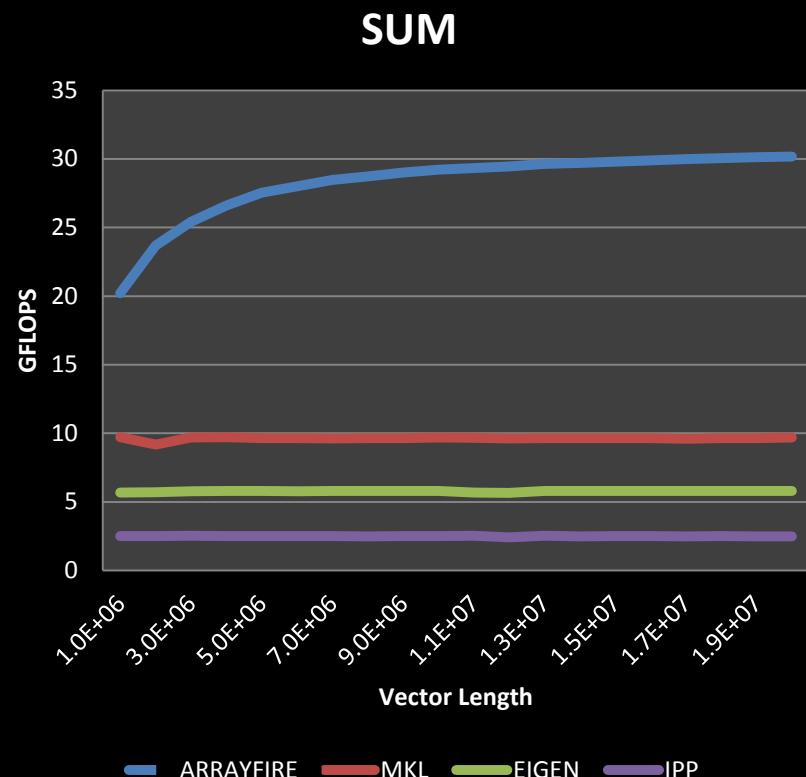
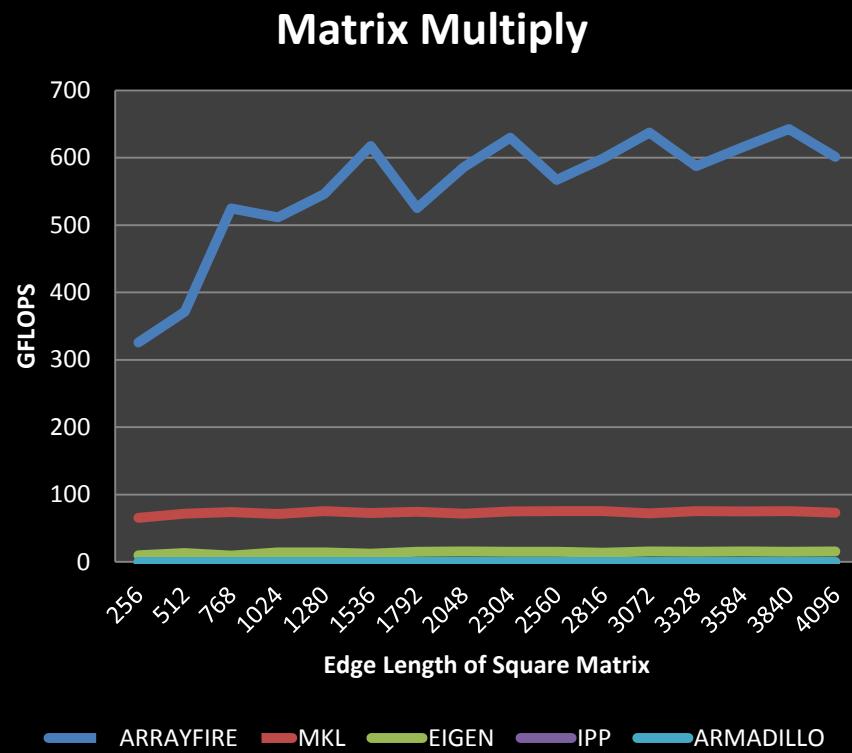
- vectors, matrices
- rescaling

sorting

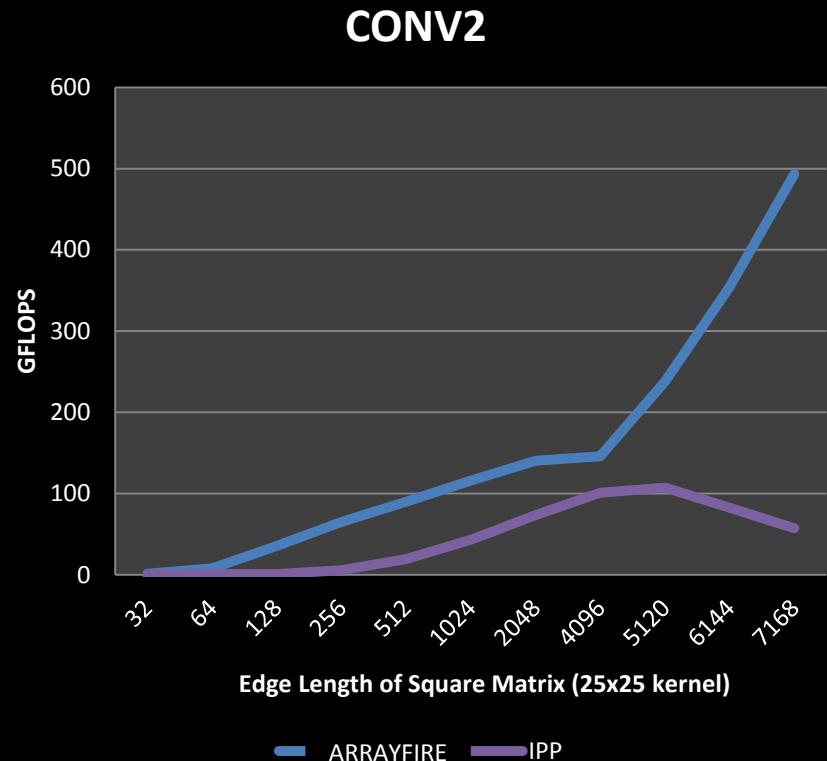
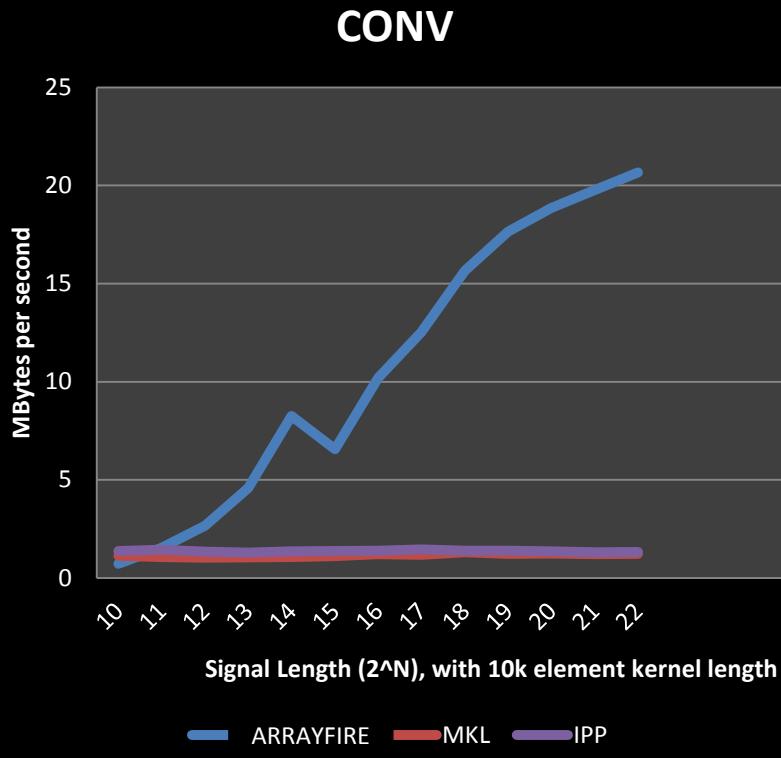
- along any dimension
- sort detection

and many more...

speed

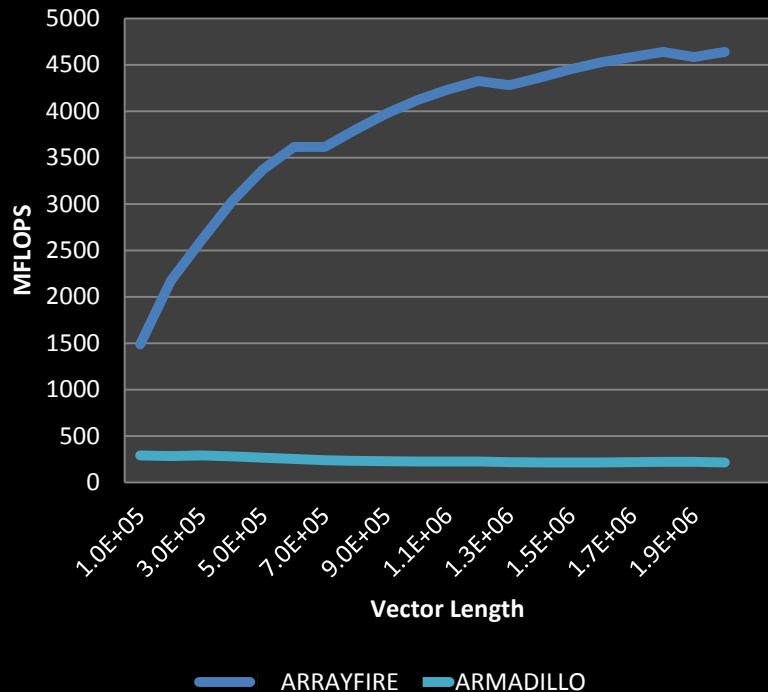


speed

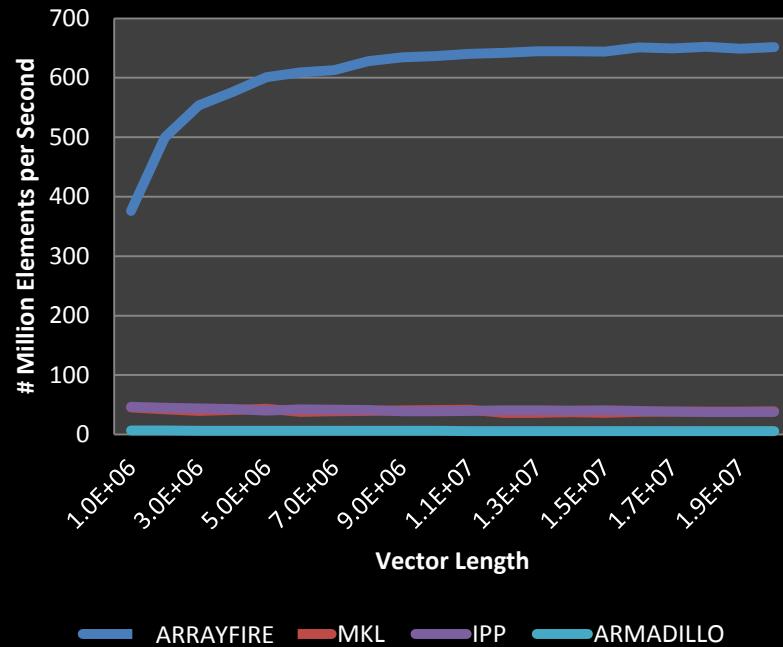


speed

Cumulative Sum



SORT



— ARRAYFIRE — MKL — IPP — ARMADILLO

Easy To Maintain

- Write your code once
 - Each new release improves the speed of your code
 - Tesla, Fermi, Kepler, ..
- Behind the scenes...we update



Licensing

- Free
 - Full performance, single GPU, internet connection to license server
- Pro
 - Multi-GPU, full linear algebra, offline



Get Help



Custom development

- Port CPU code
- Extend existing code to use the GPU
- Numerical, image, signal, algebra, ..

Training courses

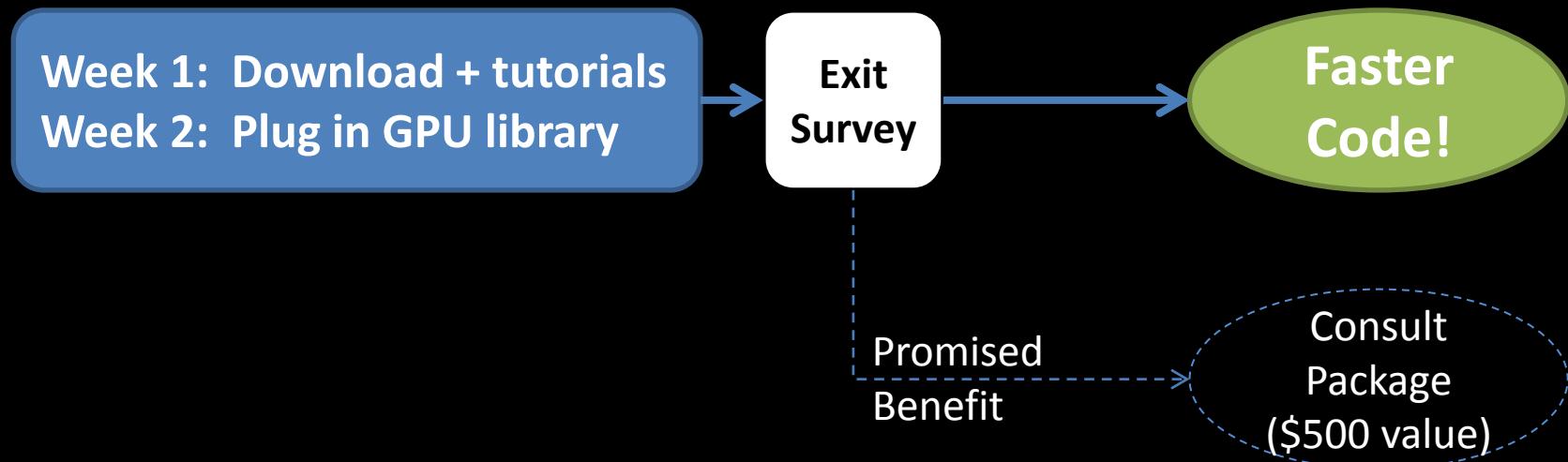
- 1-4 days
- CUDA, OpenCL
- Hands on work with your code

Gear Up for Speedup

- NVIDIA and AccelerEyes guarantee 4X code improvements in just 2 weeks
- Use ArrayFire GPU library for CUDA!



How it works...



What others are saying...



Dominic

“At Catsaras Investments Ltd, ArrayFire gave us 50X speedups on our financial models and is really easy-to-use. With ArrayFire, we avoid the time-sink of writing and optimizing GPU kernels by hand.”

Gear Up for Speedup!

