

libjacket

The World's Largest, Fastest GPU Library



The Matrix Companion to CUDA



CUDA

CUBLAS

CUFFT

Kernels

Easy GPU Acceleration in C++

```
#include <stdio.h>
#include <jacket.h>
using namespace jkt;
int main() {
    int n = 20e6; // 20 million random samples
    f32 x = f32::rand(n,1), y = f32::rand(n,1);
    // how many fell inside unit circle?
    float pi = 4.0 * sum_vector(sqrt(x*x + y*y) < 1) / n;
    printf("pi = %g\n", pi);
    return 0;
}
```

Easy GPU Acceleration in C++

```
#include <stdio.h>
#include <jacket.h>
using namespace jkt;

int main() {
    int n = 20e6; // 20 million random samples
    f32 x = f32::rand(n,1), y = f32::rand(n,1);
    // how many fell inside unit circle?
    float pi = 4.0 * sum_vector(sqrt(x*x + y*y) < 1) / n;
    printf("pi = %g\n", pi);
    return 0;
}
```

Easy GPU Acceleration in C++

```
#include <stdio.h>
#include <jacket.h>
using namespace jkt;
int main() {
    int n = 20e6; // 20 million random samples
    f32 x = f32::rand(n,1), y = f32::rand(n,1);
    // how many fell inside unit circle?
    float pi = 4.0 * sum_vector(sqrt(x*x + y*y) < 1) / n;
    printf("pi = %g\n", pi);
    return 0;
}
```

On GPU

Matrix Types

f32
real single precision

f64
real double precision

b8
boolean byte

c32
complex single precision

c64
complex double precision

Matrix Types: ND Support

f64

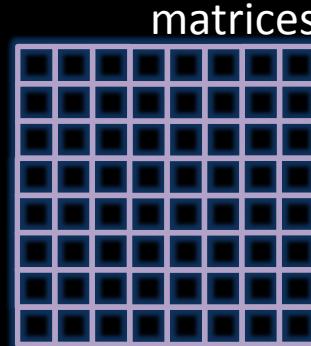
real double precision



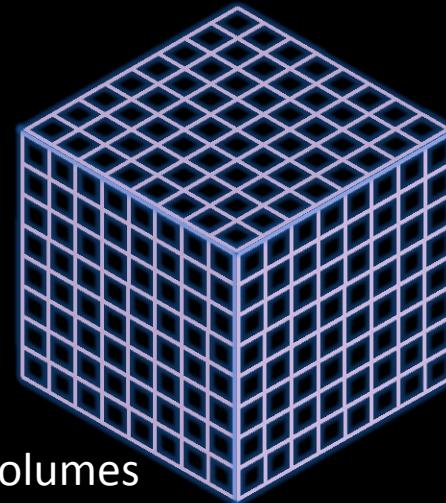
vectors

f32

real single precision



matrices



volumes

b8

boolean byte

... ND

c32

complex single precision

c64

complex double precision

Matrix Types: Easy Manipulation

f64

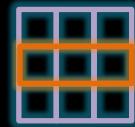
real double precision

Jacket Keywords: end, span

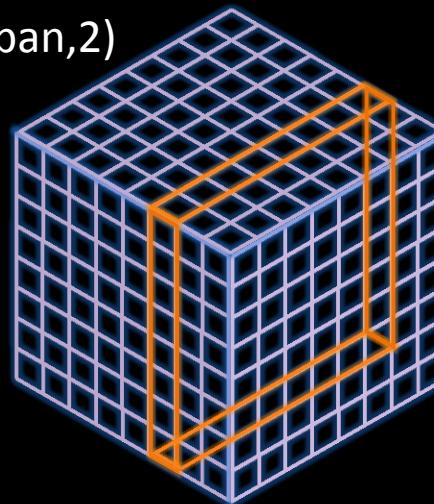
A(1,1)



A(1,span)



A(span,span,2)



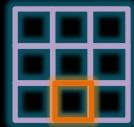
b8

boolean byte

f32

real single precision

A(end,1)



A(end,span)



c32

complex single precision

c64

complex double precision

Easy GPU Acceleration in C++

```
#include <stdio.h>
#include <jacket.h>
using namespace jkt;
int main() {
    int n = 20e6; // 20 million random samples
    f32 x = f32::rand(n,1), y = f32::rand(n,1);
    // how many fell inside unit circle?
    float pi = 4.0 * sum_vector(sqrt(x*x + y*y) < 1) / n;
    printf("pi = %g\n", pi);
    return 0;
}
```

Result on CPU

Python, C, and Fortran

```
# python  
import jacket as jkt  
b = jkt.rand(4,4)
```

```
! Fortran  
type(f32)  
dst = rand(4,4)
```

```
// C  
float *d_x = NULL;  
cudaMalloc(&d_x, 64);  
jkt_rand_S(d_x, 16);
```



gfor: what is it?

- Data-Parallel for loop, e.g.

Serial matrix multiplications (3 kernel launches)

```
for i = 1:3  
    C(:, :, i) = A(:, :, i) * B;
```

Parallel matrix multiplications (**1** kernel launch)

```
gfor i = 1:3  
    C(:, :, i) = A(:, :, i) * B;
```

example: matrix multiply

- Data-Parallel for loop, e.g.

Serial matrix multiplications (3 kernel launches)

```
for i = 1:3  
    C(:, :, i) = A(:, :, i) * B;
```

iteration i = 1

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

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

example: matrix multiply

- Data-Parallel for loop, e.g.

Serial matrix multiplications (3 kernel launches)

```
for i = 1:3
    C(:, :, i) = A(:, :, i) * B;
```

iteration i = 1

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

iteration i = 2

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

example: matrix multiply

- Data-Parallel for loop, e.g.

Serial matrix multiplications (3 kernel launches)

```
for i = 1:3
    C(:, :, i) = A(:, :, i) * B;
```

iteration i = 1

$$\begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} * \begin{bmatrix} & \\ & \\ & \end{bmatrix}$$

$C(:, :, i)$ $A(:, :, i)$ B

iteration i = 2

$$\begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} * \begin{bmatrix} & \\ & \\ & \end{bmatrix}$$

$C(:, :, i)$ $A(:, :, i)$ B

iteration i = 3

$$\begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} * \begin{bmatrix} & \\ & \\ & \end{bmatrix}$$

$C(:, :, i)$ $A(:, :, i)$ B

example: matrix multiply

Parallel matrix multiplications (**1** kernel launch)

```
gfor i = 1:3  
    C(:,:,i) = A(:,:,i) * B;
```

simultaneous iterations i = 1:3

$$\begin{array}{ccc} \begin{matrix} \text{C(:,:,1)} \\ \text{A(:,:,1)} \end{matrix} & = & \begin{matrix} \text{B} \end{matrix} \\ \begin{matrix} \text{C(:,:,2)} \\ \text{A(:,:,2)} \end{matrix} & = & \begin{matrix} \text{B} \end{matrix} \\ \begin{matrix} \text{C(:,:,3)} \\ \text{A(:,:,3)} \end{matrix} & = & \begin{matrix} \text{B} \end{matrix} \end{array}$$

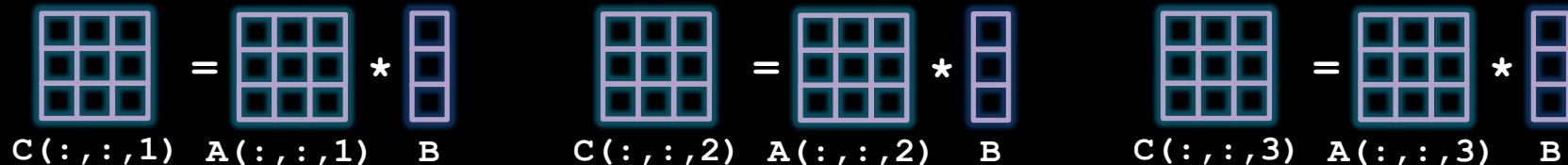
example: matrix multiply

Parallel matrix multiplications (1 kernel launch)

```
gfor i = 1:3  
    C(:, :, i) = A(:, :, i) * B;
```

Think “**gfor i=1:3**” as “extend 1:3 in the next dimension”

```
simultaneous iterations i = 1:3
```

$$\begin{array}{ccc} \text{C}(:,:,1) & = & \text{A}(:,:,1) * \text{B} \\ \text{C}(:,:,2) & = & \text{A}(:,:,2) * \text{B} \\ \text{C}(:,:,3) & = & \text{A}(:,:,3) * \text{B} \end{array}$$


example: matrix multiply

Parallel matrix multiplications (1 kernel launch)

```
gfor i = 1:3  
    C(:, :, i) = A(:, :, i) * B;
```

Think “**gfor i=1:3**” as “extend 1:3 in the next dimension”

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 i = 1:3  
    C(:, :, i) = A(:, :, i) * B;
```

Think “**gfor i=1:3**” as “extend 1:3 in the next dimension”

simultaneous iterations i = 1:3

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

example: printing to screen

- gfor runs its body only once on the CPU

Loop body executed three times

```
for i = 1:3  
    disp(i);
```

Loop body executed one time

```
gfor i = 1:3  
    disp(i);
```

example: printing to screen

- gfor runs its body only once on the CPU

Loop body executed three times

```
for i = 1:3  
    disp(i);
```

→ >> go;
1
2
3

Loop body executed one time

```
gfor i = 1:3  
    disp(i);
```

→ >> go;
1

example: summing over columns

- Think of gfor as “syntactic sugar” to write vectorized code in an iterative style.

Three passes to sum all columns of B

```
for i = 1:3  
    A(i) = sum(B(:,i));
```

Both equivalent to “**sum(B)**”,
but latter is faster (more
explicitly written)



One pass to sum all columns of B

```
gfor i = 1:3  
    A(i) = sum(B(:,i));
```

Easy Multi GPU Scaling

```
f32 *y = new f32[n];
for (int i = 0; i < n; ++i) {
    device(i);                      // change GPUs
    f32 x = f32::rand(5,5);         // add work to GPU's queue
    y[i] = fft(x);                  // more work in queue
}

// all GPUs are now computing simultaneously, until done
```

Hundreds of Functions...

reductions

- sum, min, max, any, all, nnz, prod
- vectors, columns, rows, etc

dense linear algebra

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

convolutions

- 2D, 3D, ND

FFTs

- 2D, 3D, ND

image processing

- filter, rotate, erode, dilate, bwmorph, resize, rgb2gray
- hist, histeq

interp and rescale

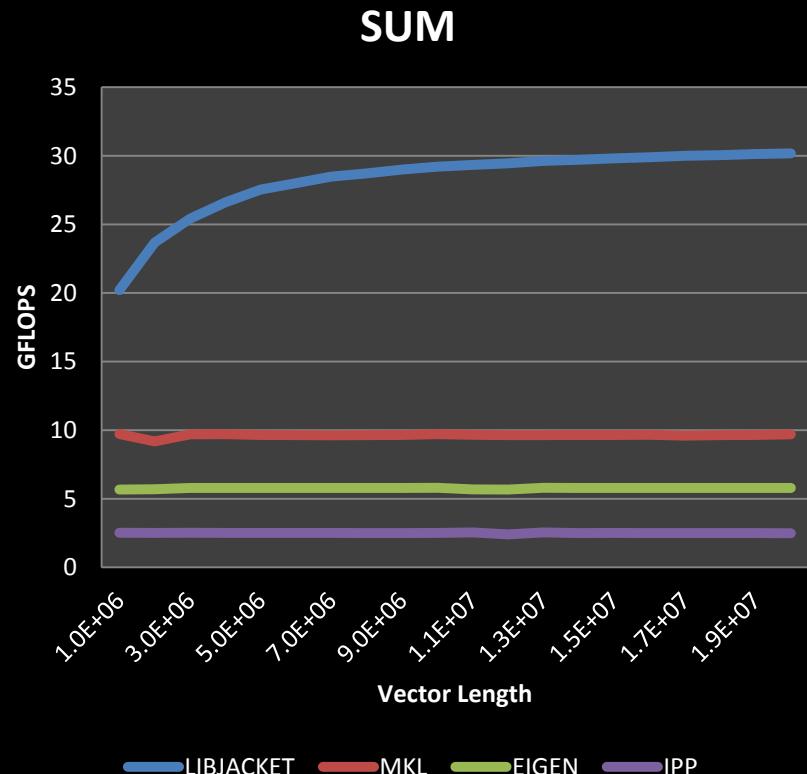
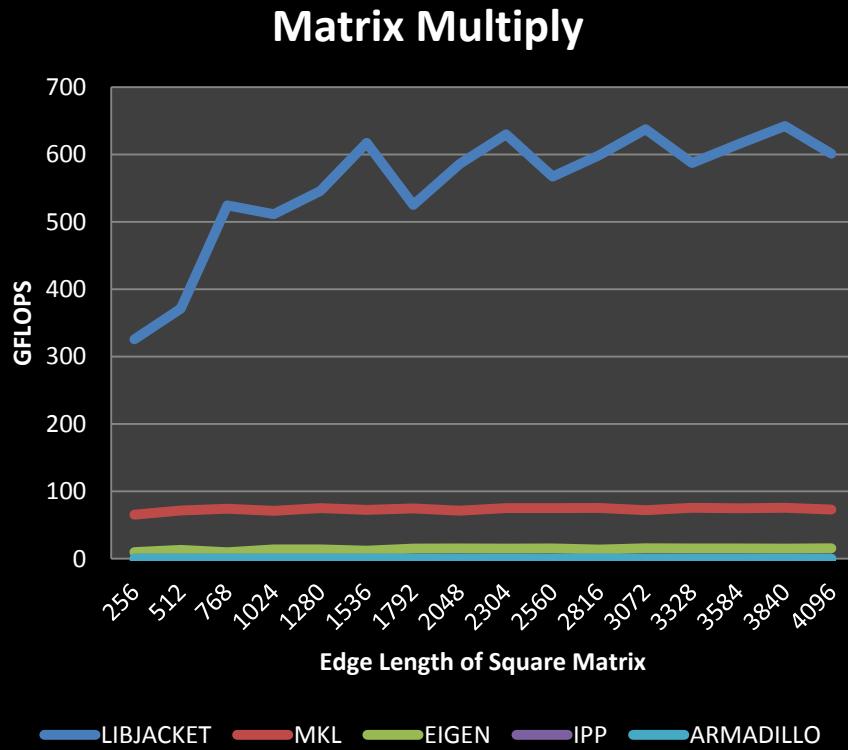
- vectors, matrices
- rescaling

sorting

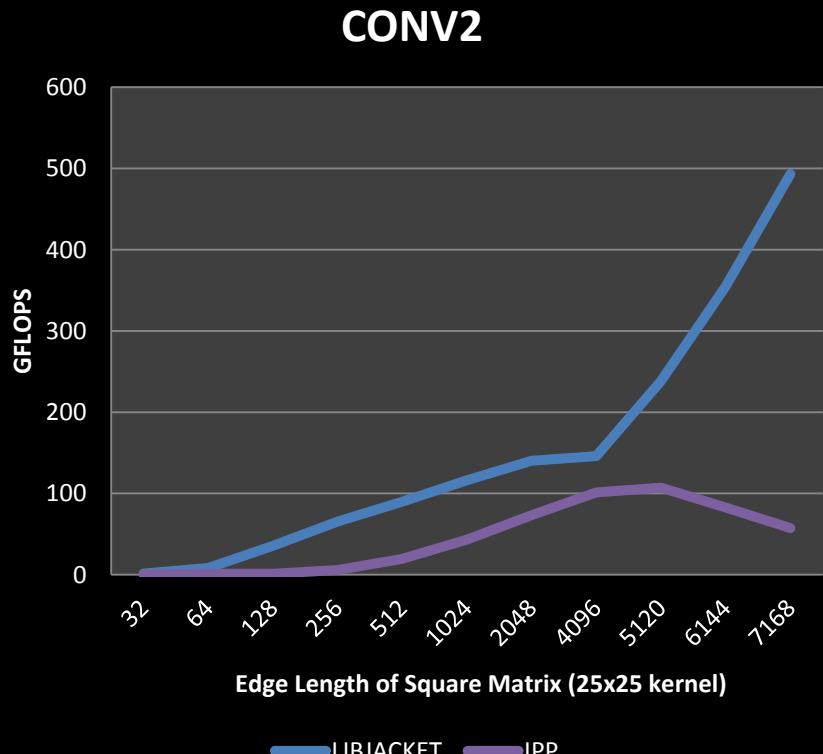
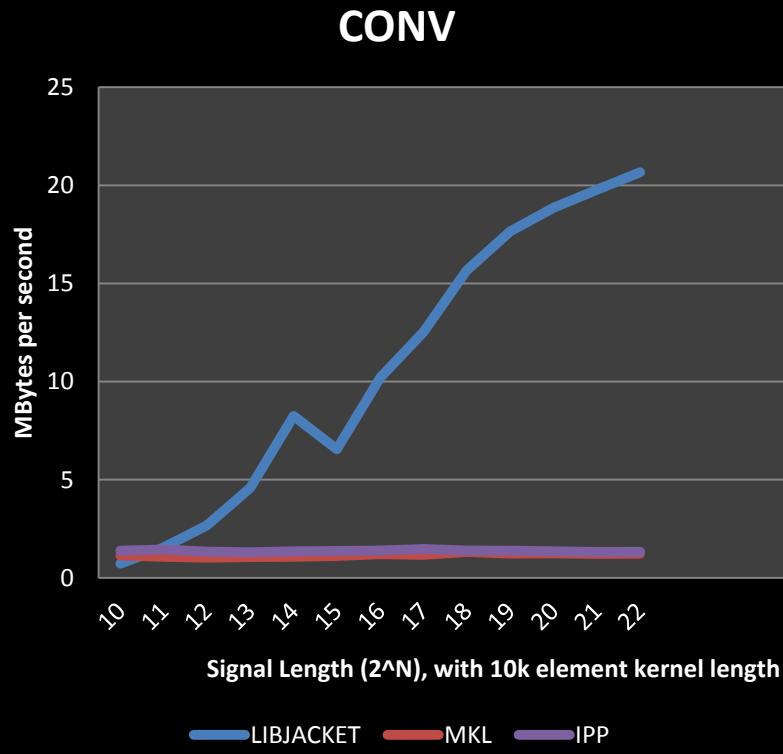
- along any dimension
- sort detection

and many more...

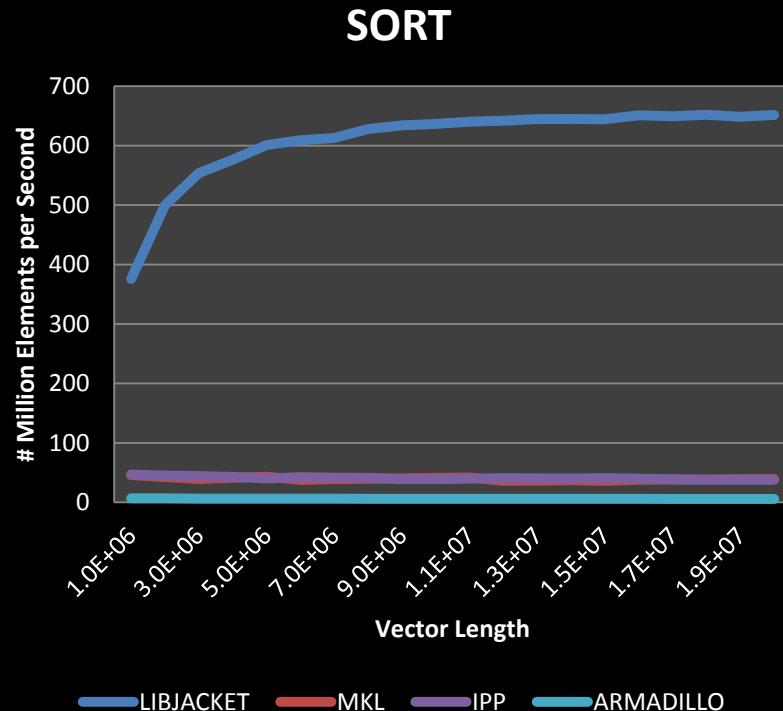
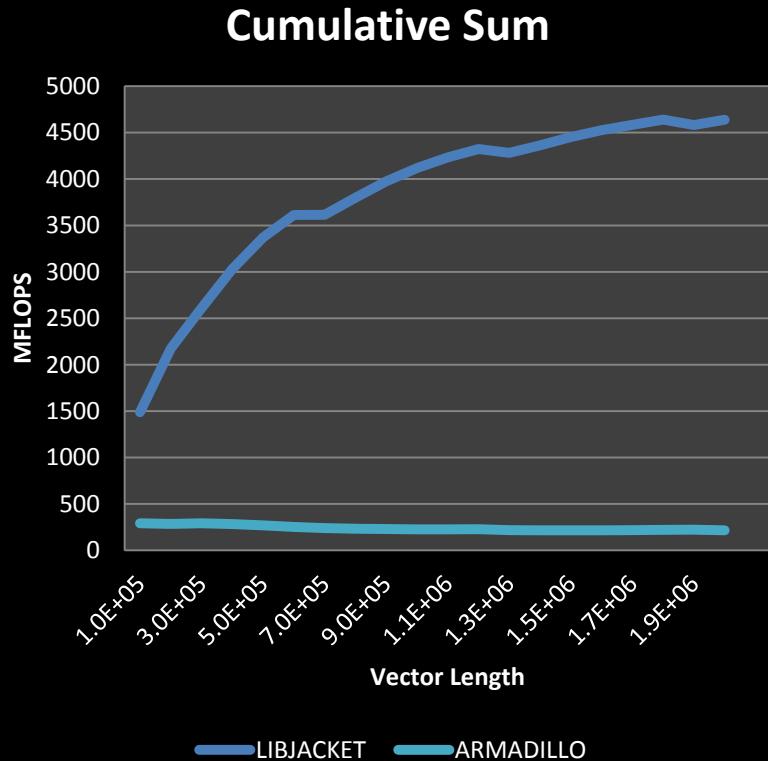
...Plenty of Speed



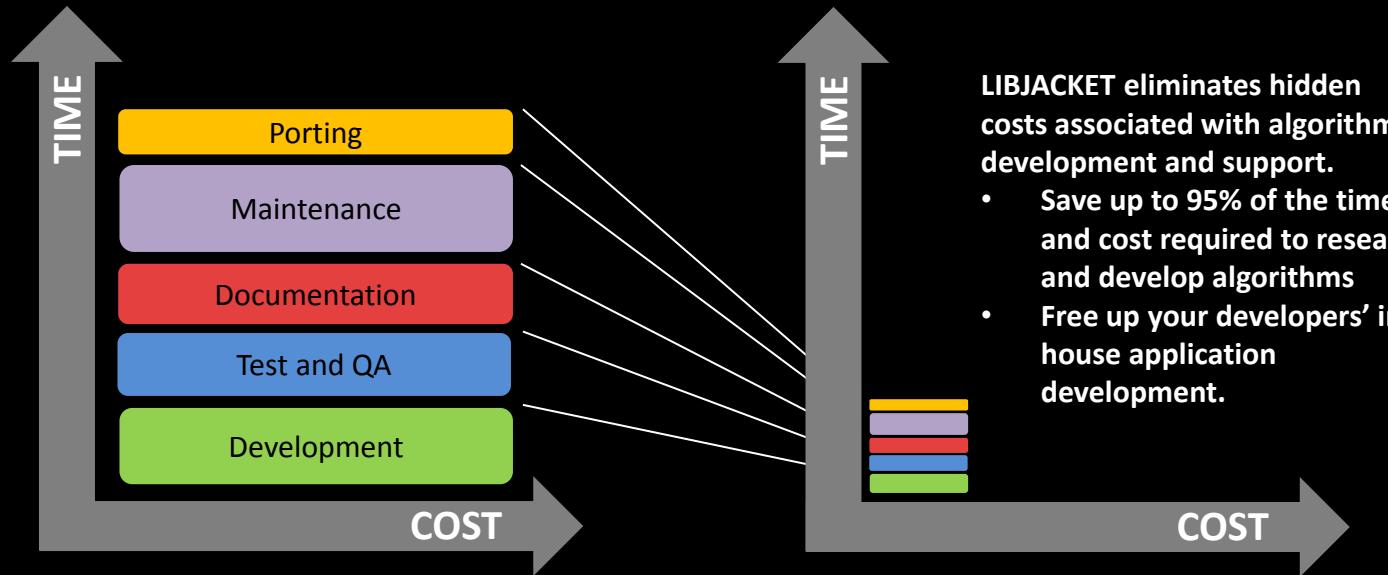
...Plenty of Speed



...Plenty of Speed



Build vs Buy



Easy To Maintain

- Write your code once and let LibJacket carry you through the coming hardware evolution.
 - Each new LibJacket release improves the speed of your code, without any code modification.
 - Each new LibJacket release leverages latest GPU hardware (e.g. Fermi), without any code modification.



Get libjacket

The World's Largest, Fastest GPU Library



www.accelereyes.com

