

# Code Specialisation of Auto Generated GPU Kernels

Troels Blum

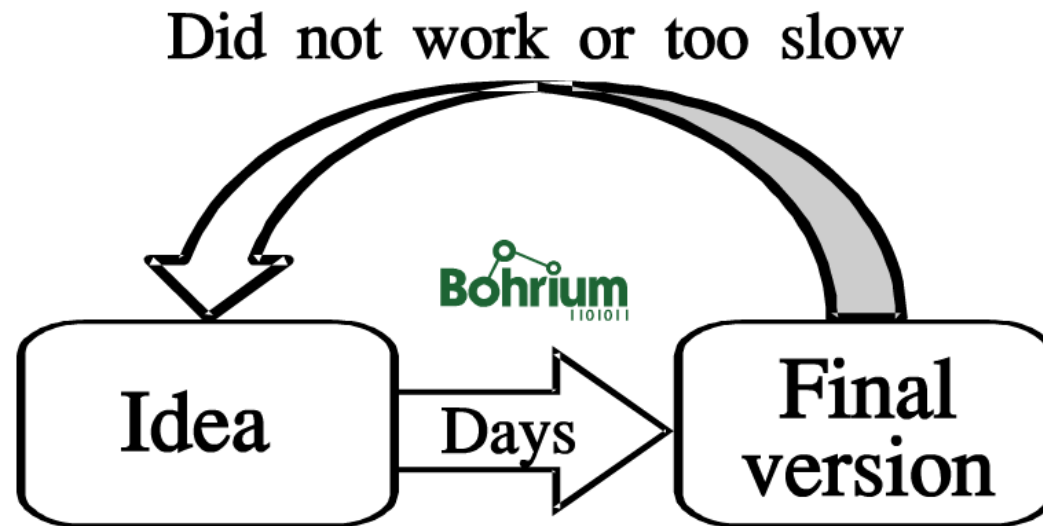
University of Copenhagen  
DENMARK



<http://bh107.org>

# Development Cycle

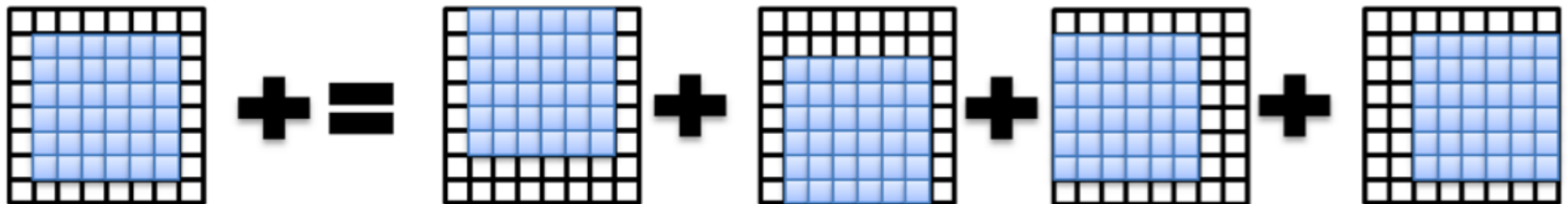
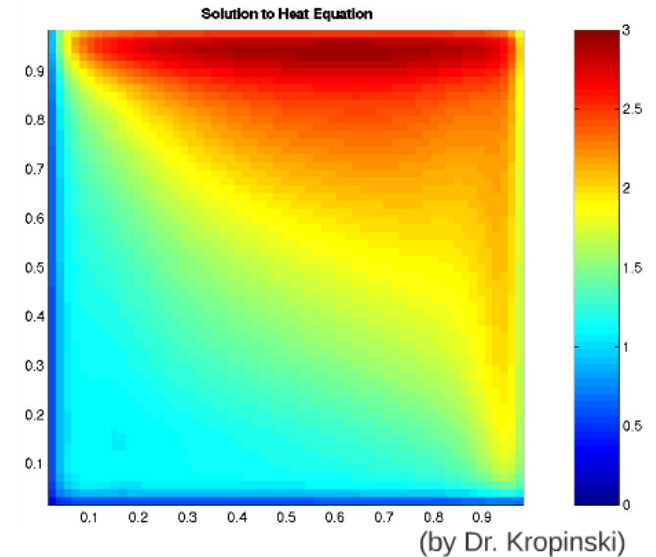
## High-Performance vs High-Productivity



# Array Programming

## 2D 5-point stencil

```
def heat_eq(grid, epsilon=0.005):  
    delta = epsilon + 1  
    while delta > epsilon:  
        work = (grid[1:-1,1:-1] + grid[0:-2,1:-1] +  
               grid[1:-1,2:] + grid[1:-1,0:-2] +  
               grid[2:,1:-1]) * 0.2  
        delta = numpy.sum(numpy.absolute(work - grid[1:-1,1:-1]))  
        grid[1:-1,1:-1] = work  
    return grid
```

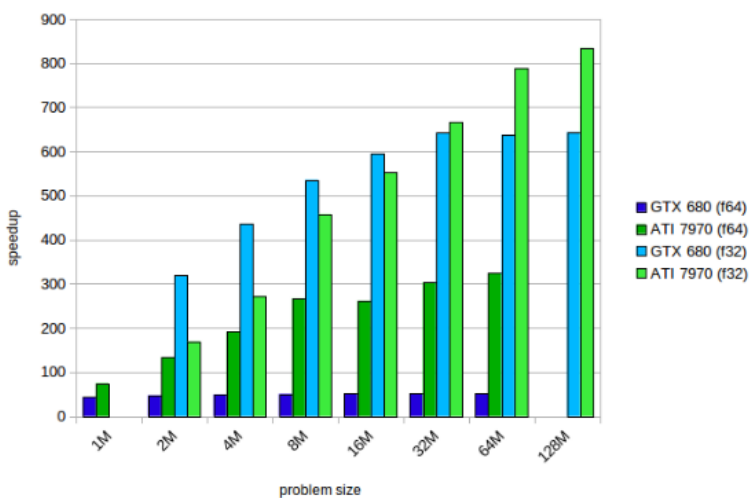


$A[1:-1,1:-1] += A[:-2, 1:-1] + A[2:, 1:-1] + A[1:-1,0:-2] + A[2:,1:-1]$

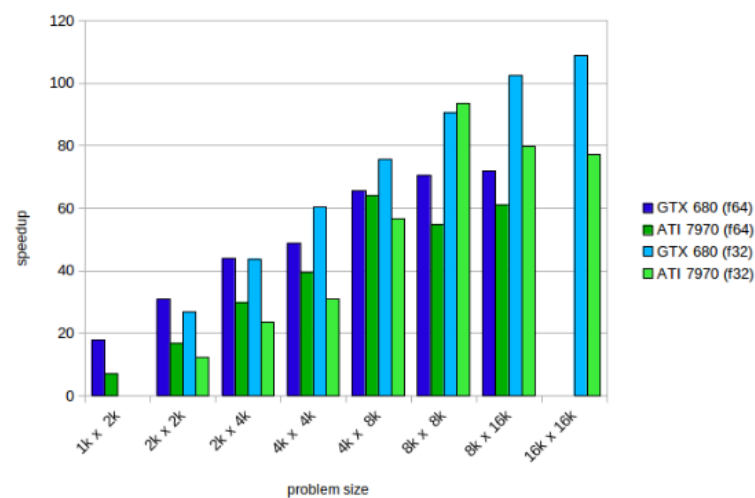
# Benchmarks

Relative speedup compared to native NumPy

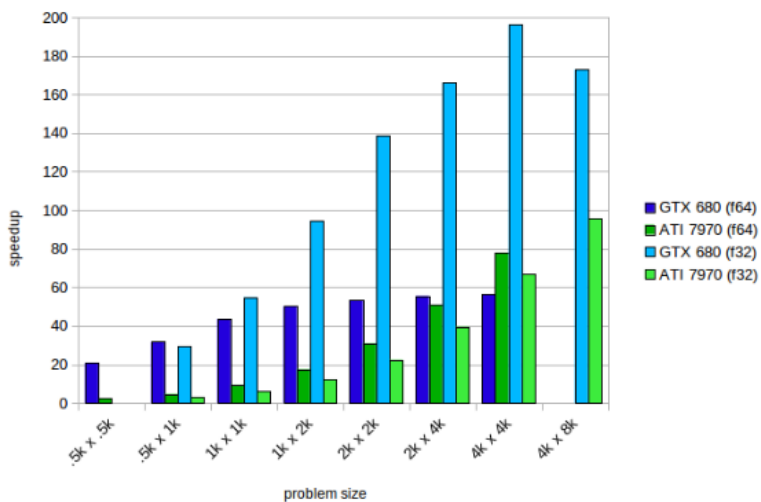
## Black-Scholes



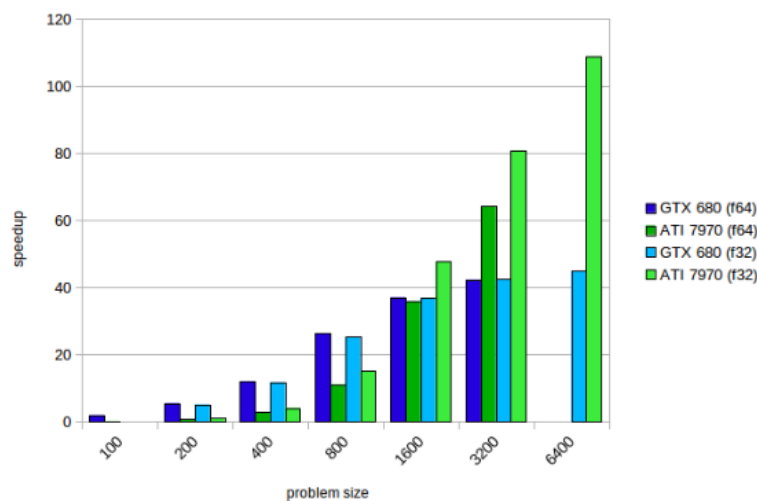
## Successive Over-relaxation (SOR)



## Shallow Water



## N-body



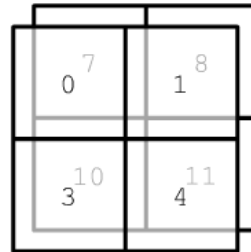
# Bohrium Byte-code

Data structure

type	float64		
ndim	3		
start	0		
shape	2	2	2
stride	7	3	1
data	*		



Data layout



Seen 3d-array

```

ADD      t1, v1, v2
ADD      t2, t1, v3
DIV      t3, t2, 3.0
SUB      t4, t3, v1
ABS      t5, t4
ADD_REDUCE t6, t5
    
```

# Generated GPU kernels

```
__kernel heat_eq (...)  
{  
    const size_t gidx = get_global_id(0);  
    if (gidx >= ds0)  
        return;  
    const size_t gidy = get_global_id(1);  
    if (gidy >= ds1)  
        return;  
    double v1 = a1[gidy*v1s2 + gidx*v1s1 + v1s0];  
    double v2 = a1[gidy*v2s2 + gidx*v2s1 + v2s0];  
    double v4 = a1[gidy*v4s2 + gidx*v4s1 + v4s0];  
    double v6 = a1[gidy*v6s2 + gidx*v6s1 + v6s0];  
    double v8 = a1[gidy*v8s2 + gidx*v8s1 + v8s0];  
    double v0;  
    v0 = v1 + v2;  
    double v3;  
    v3 = v0 + v4;  
    double v5;  
    v5 = v3 + v6;  
    double v7;  
    v7 = v5 + v8;  
    double v9;  
    v9 = v7 * s0;  
    double v10;  
    v10 = v9 - v1;  
    double v11;  
    v11 = fabs(v10);  
    a5[gidy*v9s2 + gidx*v9s1 + v9s0] = v9;  
    a7[gidy*v11s2 + gidx*v11s1 + v11s0] = v11;  
}
```

```
__kernel heat_eq(...)  
{  
    const size_t gidx = get_global_id(0);  
    if (gidx >= 1998)  
        return;  
    const size_t gidy = get_global_id(1);  
    if (gidy >= 1998)  
        return;  
    double v1 = a1[gidy*2000 + gidx*1 + 2001];  
    double v2 = a1[gidy*2000 + gidx*1 + 1];  
    double v4 = a1[gidy*2000 + gidx*1 + 2002];  
    double v6 = a1[gidy*2000 + gidx*1 + 2000];  
    double v8 = a1[gidy*2000 + gidx*1 + 0];  
    double v0;  
    v0 = v1 + v2;  
    double v3;  
    v3 = v0 + v4;  
    double v5;  
    v5 = v3 + v6;  
    double v7;  
    v7 = v5 + v8;  
    double v9;  
    v9 = v7 * s0;  
    double v10;  
    v10 = v9 - v1;  
    double v11;  
    v11 = fabs(v10);  
    a5[gidy*1998 + gidx*1 + 0] = v9;  
    a7[gidy*1998 + gidx*1 + 0] = v11;  
}
```

# Assumptions

- Specialization will enhance kernel run-time
- Compilation of kernels is time consuming

# Benchmark applications

- Stencil applications
  - 1D, 3 points
  - 2D, 9 points
  - 3D, 27 points
  - 4D, 81 points
- Gaussian elimination
- LU decomposition

## Kernel types

Dynamic

Fixed

Selected



```

def heat_eq(grid, epsilon=0.005):
    delta = epsilon + 1
    while delta > epsilon:
        work = (grid[1:-1,1:-1] + grid[0:-2,1:-1] +
                grid[1:-1,2:] + grid[1:-1,0:-2] +
                grid[2:,1:-1]) * 0.2
        delta = numpy.sum(numpy.absolute(work - grid[1:-1,1:-1]))
        grid[1:-1,1:-1] = work
    return grid

```

1D - 4D  
stencil

```

def gauss(a):
    for c in xrange(1,a.shape[0]):
        a[c:,c-1:] = a[c:,c-1:] -
                    (a[c:,c-1]/a[c-1,c-1:c])[:,None] *
                    a[c-1,c-1:]
    a /= numpy.diagonal(a)[:,None]
    return a

```

Gauss  
elimination

```

def lu(a):
    u = a.copy()
    l = numpy.zeros like(a)
    numpy.diagonal(l)[:] = 1.0
    for c in xrange(1,u.shape[0]):
        l[c:,c-1] = u[c:,c-1] / u[c-1,c-1:c]
        u[c:,c-1:] = u[c:,c-1:] -
                    l[c:,c-1][:,None] * u[c-1,c-1:]
    return (l,u)

```

1D - 4D  
stencil

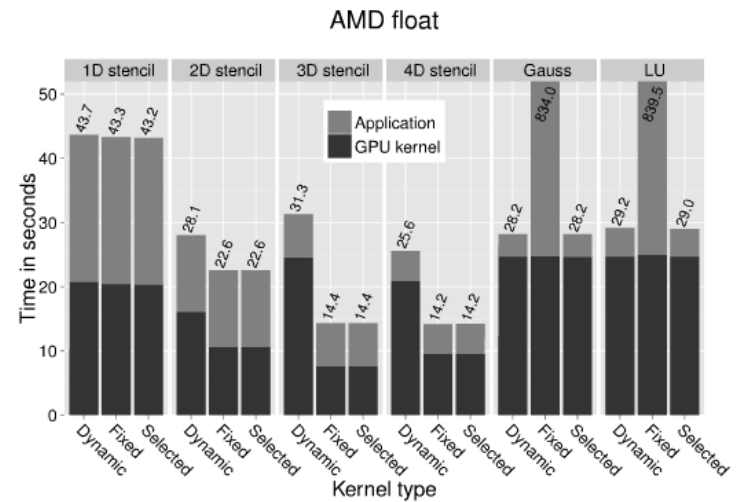
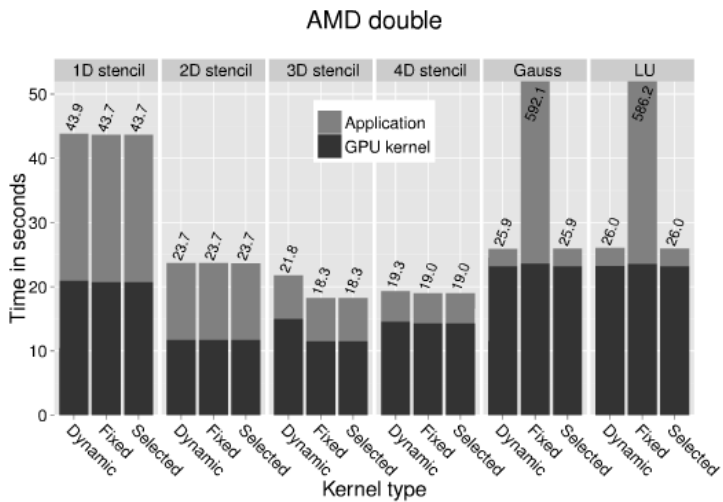
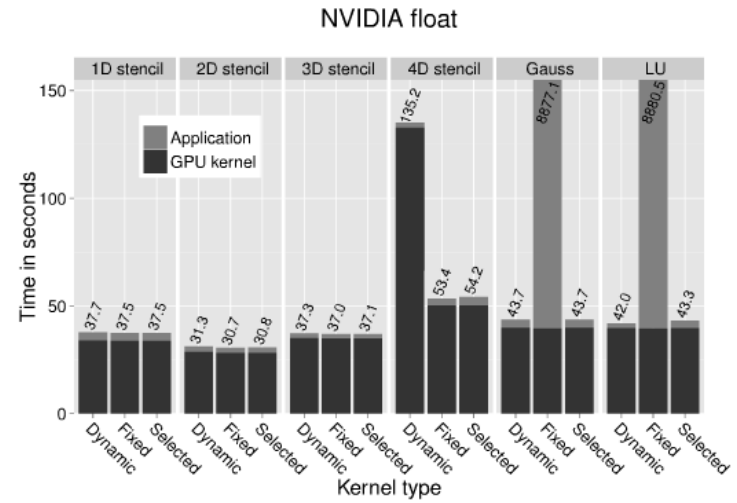
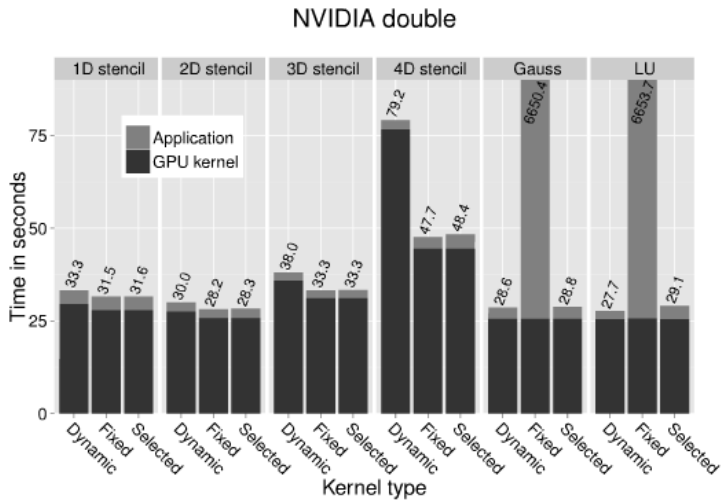
# GPU Specs

Vendor:	AMD	NVIDIA
Model:	HD 7970	GTX 680
Driver version:	1214.3 (VM)	331.38
#Cores:	2048	1536
Clock:	1000 MHz	1006 MHz
Memory:	3GB GDDR5	2GB DDR5
Memory bandwidth:	288 GB/s	192 GB/s
Peak performance:	4096 GFLOPS	3090 GFLOPS

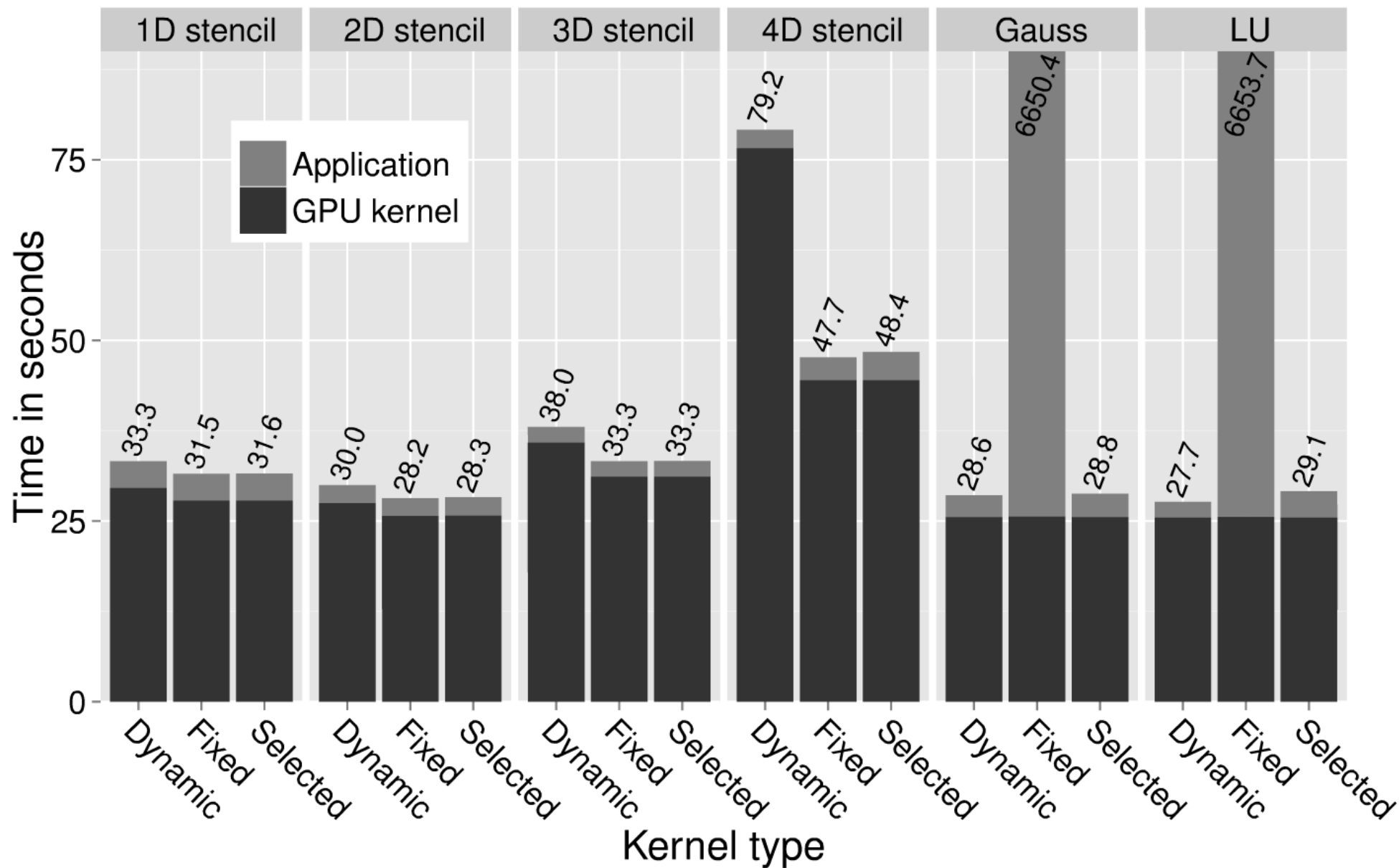
# Compile times

	NVIDIA	AMD	ratio
1D stencil	1445ms	117ms	<b>12.3</b>
2D stencil	1460ms	125ms	<b>11.7</b>
3D stencil	1503ms	151ms	<b>9.9</b>
4D stencil	3993ms	482ms	<b>8.3</b>
Gauss	2199ms	186ms	<b>11.8</b>
LU	2929ms	241ms	<b>12.2</b>
<i>Average</i>			<b><i>10.4</i></b>

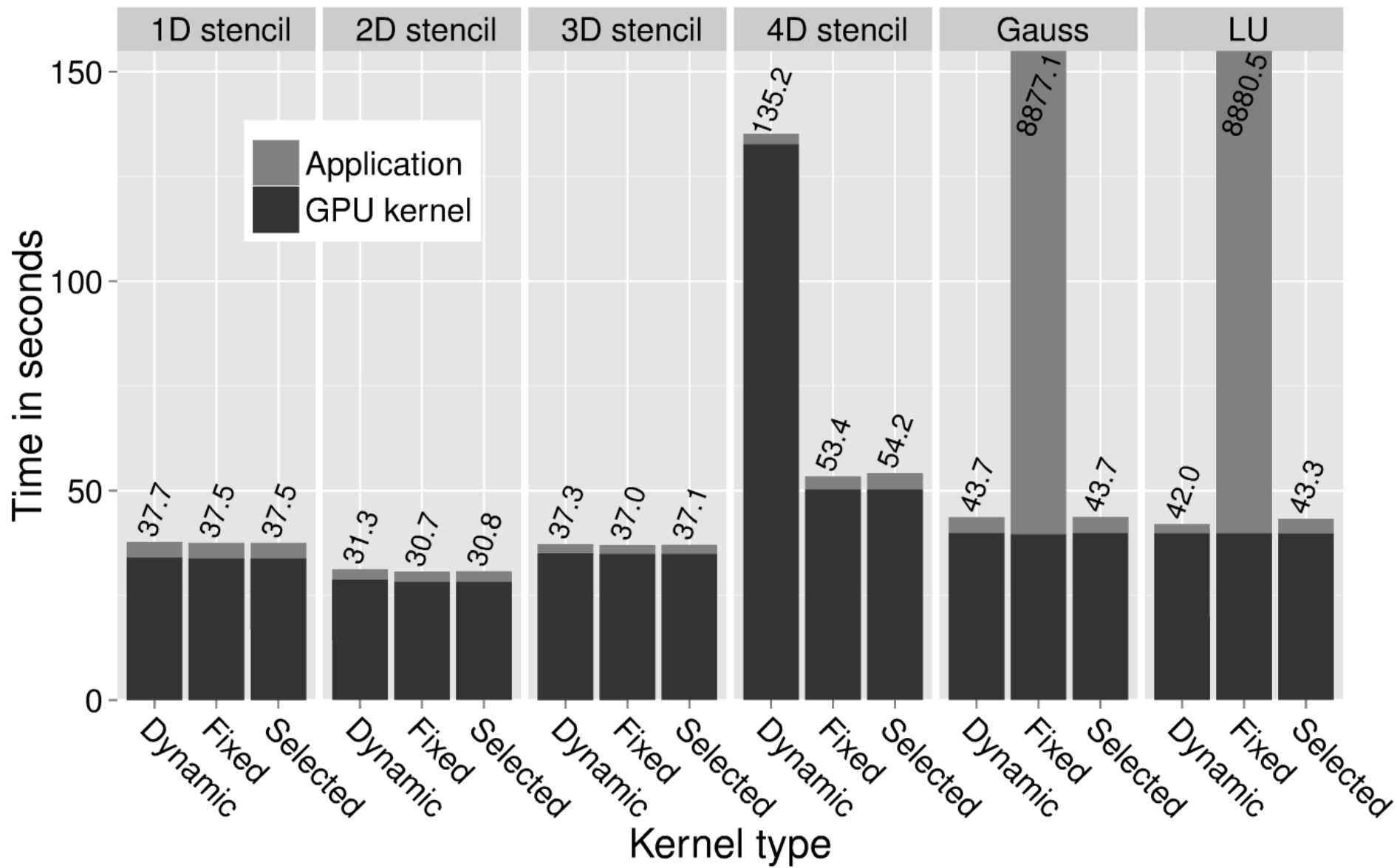
# Benchmarks



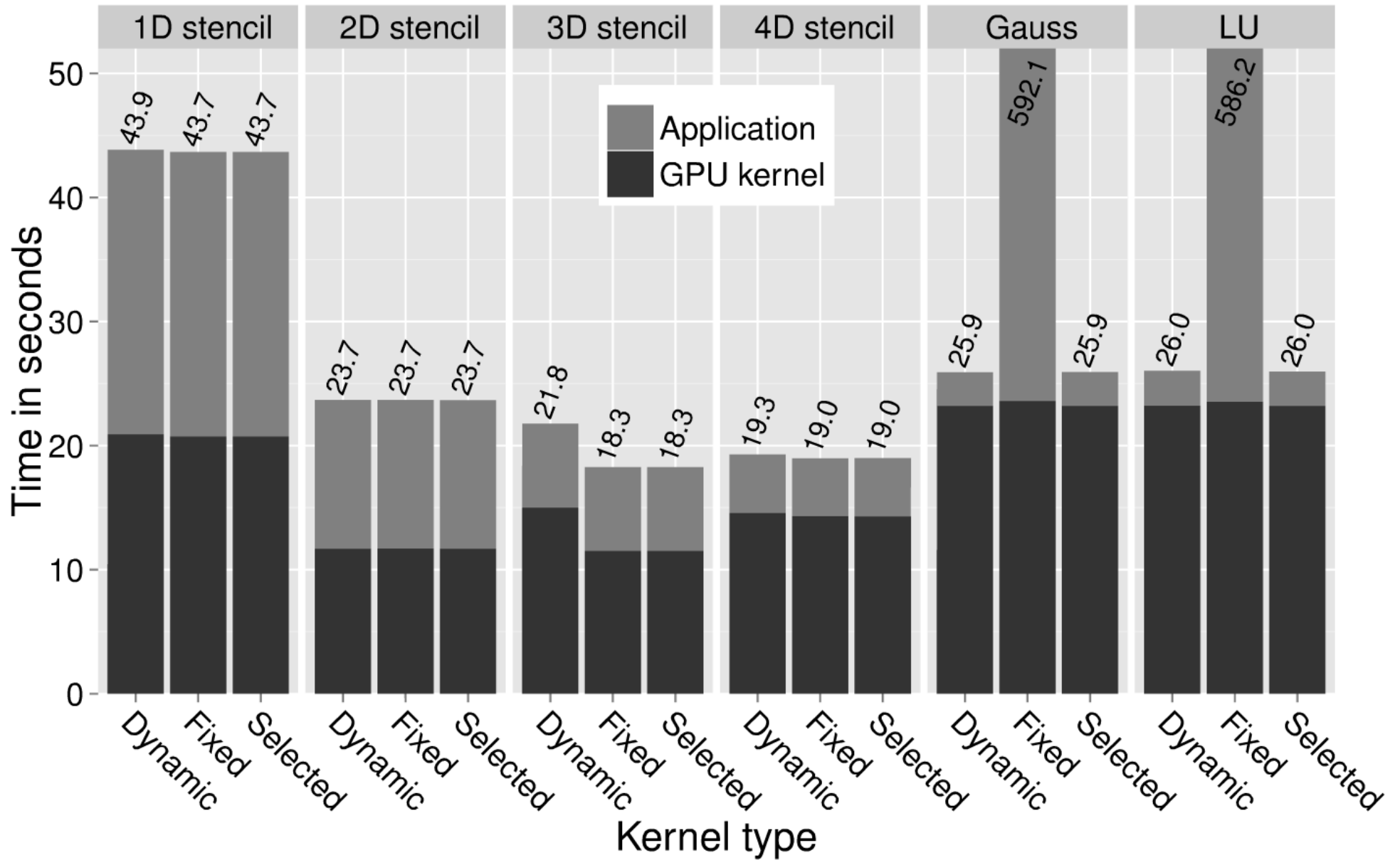
# NVIDIA double



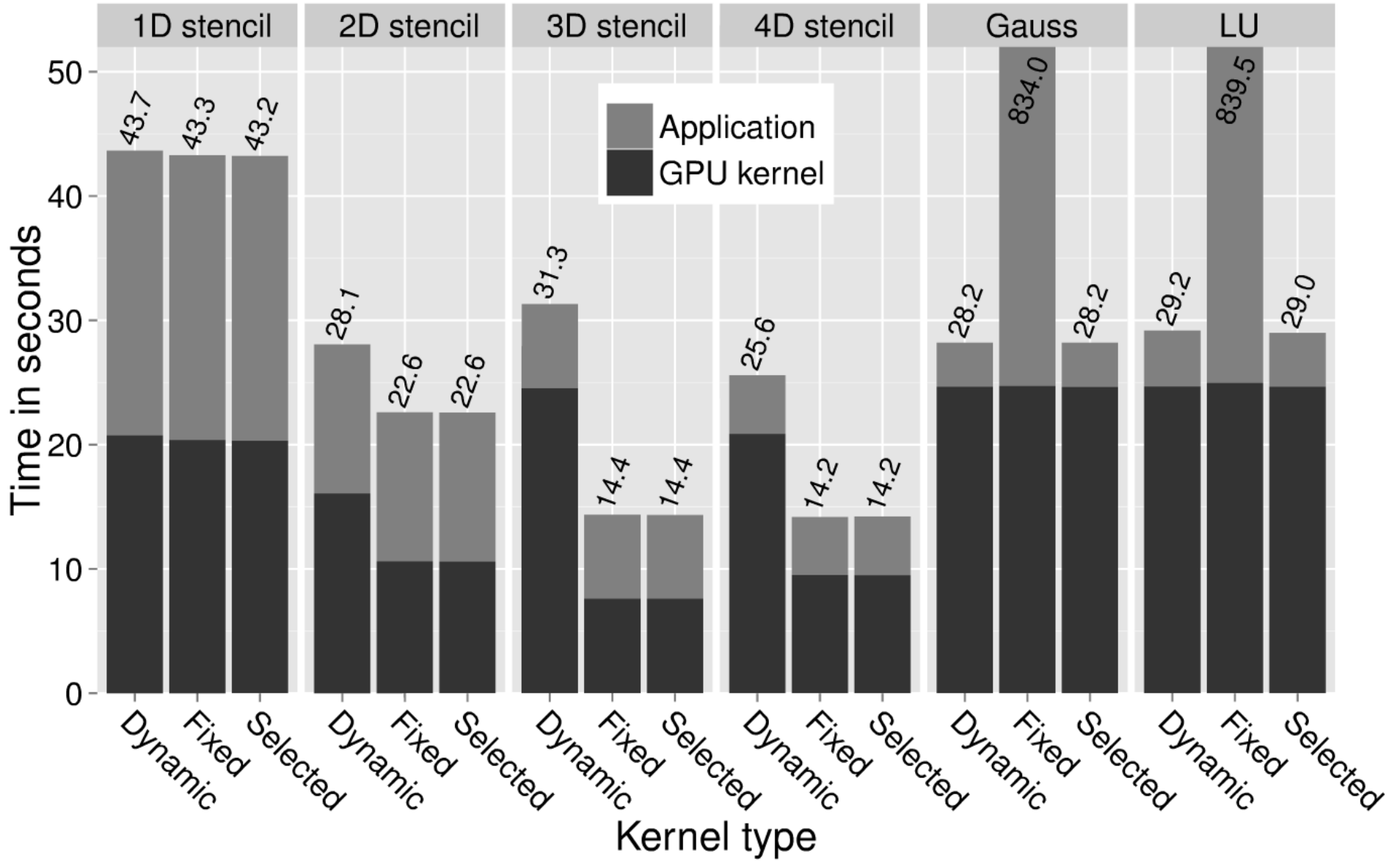
# NVIDIA float



# AMD double



# AMD float





# Thank You



<http://bh107.org>