# Global memory and Coalescing
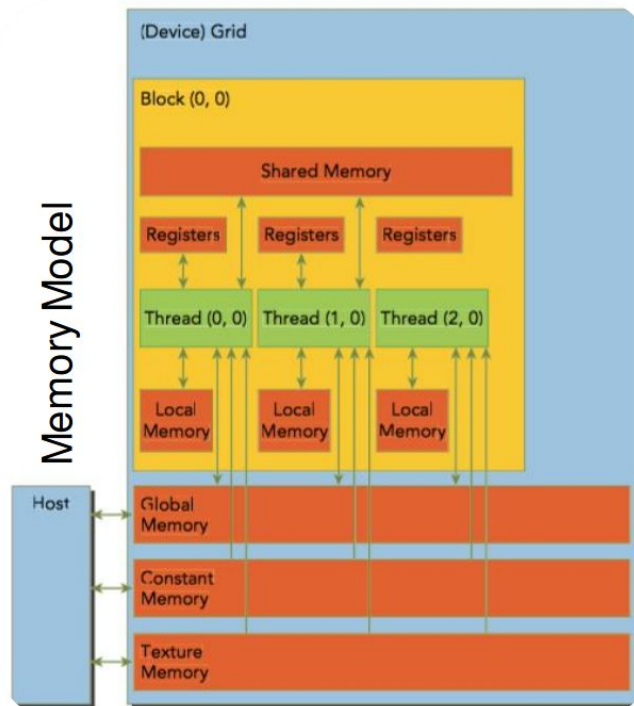
A click through example
By Nicolas Bohm Agostini

Memory Model
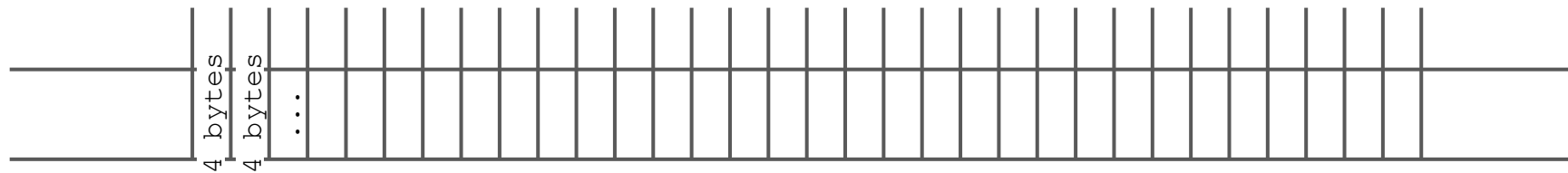
(Device) Grid

Block (0, 0)

Shared Memory

Registers | Registers | Registers

Thread (0, 0) | Thread (1, 0) | Thread (2, 0)

Local Memory | Local Memory | Local Memory

Host

Global Memory

Constant Memory

Texture Memory

Memory Model

(Device) Grid

Block (0, 0)

Shared Memory

Registers    Registers    Registers

Thread (0, 0)    Thread (1, 0)    Thread (2, 0)

Local Memory    Local Memory    Local Memory
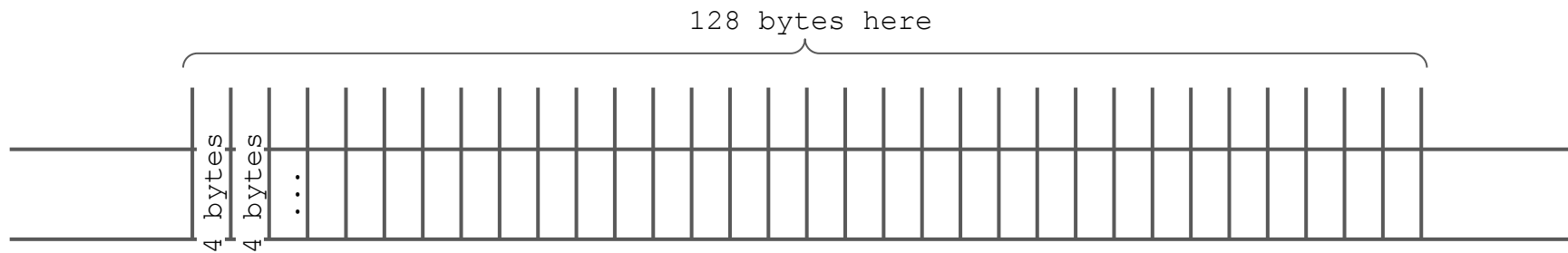
Host

This is my Global memory

Constant Memory

Texture Memory

This is my Global memory

This is my Global memory

4 bytes
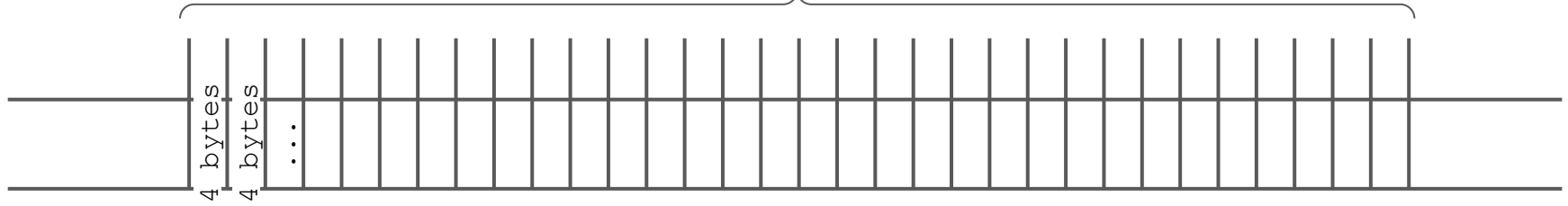
4 bytes

...

These are groups of
4 bytes
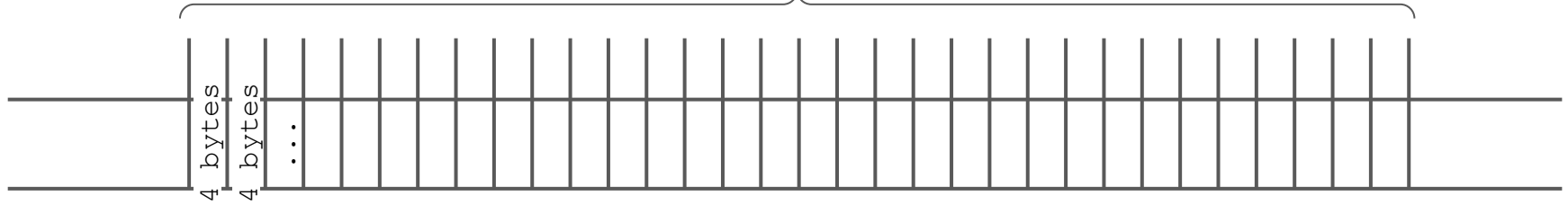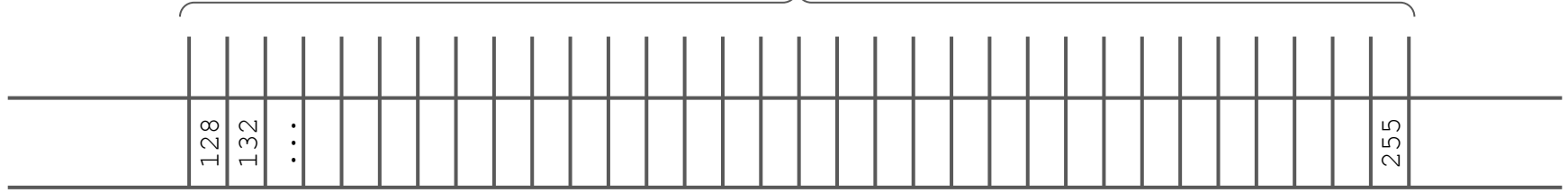
128 bytes here

4 bytes 4 bytes ...

These are groups of
4 bytes

128 bytes here. 32 Groups of 4 bytes

4 bytes 4 bytes ...

These are groups of
4 bytes

128 bytes here. 32 Groups of 4 bytes
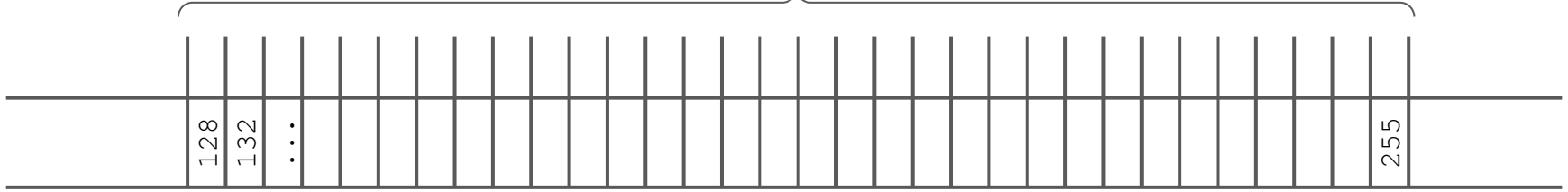
4 bytes | 4 bytes | . . .

These are groups of
4 bytes

128 bytes here. 32 Groups of 4 bytes

128 132 ... 255

This chunk of data got allocated starting on byte 128 of my
Global Memory

128 bytes here. 32 Groups of 4 bytes

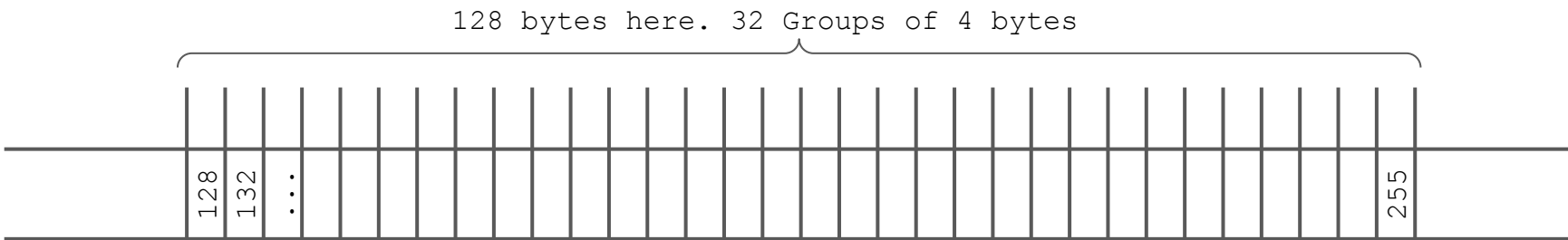| 128 | 132 | : : | | | | | | | | | | | | | | | | | | | | | | | | | | | | | 255 |

```
__device__ void add1(int* a, int vector_size) {
    id = blockIdx.x*blockDim + threadIdx.x;

    if (id<vector_size)
        a[id]=a[id]+1;
}

int main() {
    // do setup

    add1<<<4,256>>>(dev_a,vector_size);

    // do cleanup
}
```
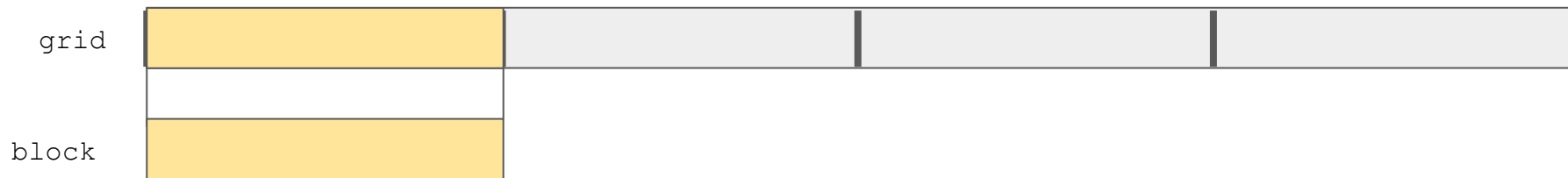
128 bytes here. 32 Groups of 4 bytes

| 128 | 132 | . . | | | | | | | | | | | | | | | | | | | | | | | | | | | | 255 |

```
int main() {
        // do setup

        add1<<<4,256>>>(dev_a,vector_size);

        // do cleanup
}
```

grid

block

128 bytes here. 32 Groups of 4 bytes

128  132  ⋅⋅⋅  255

```
int main() {
    // do setup

    add1<<<4,256>>>(dev_a,vector_size);

    // do cleanup
}
```
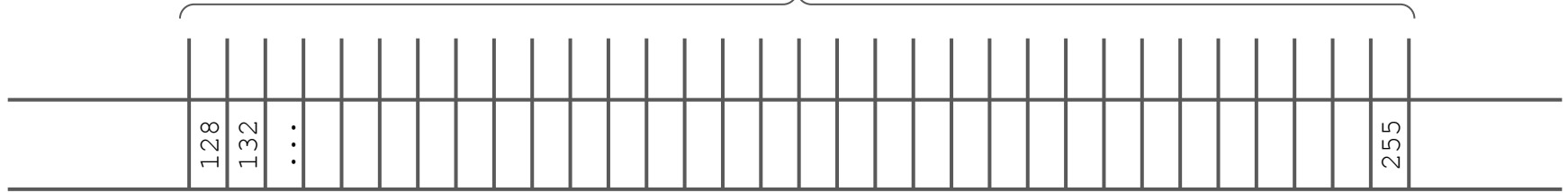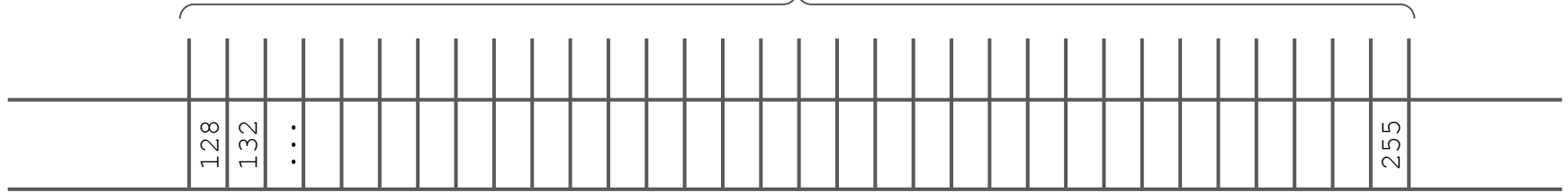
…                                    block                                    ...

128 bytes here. 32 Groups of 4 bytes

128 132 .. 255
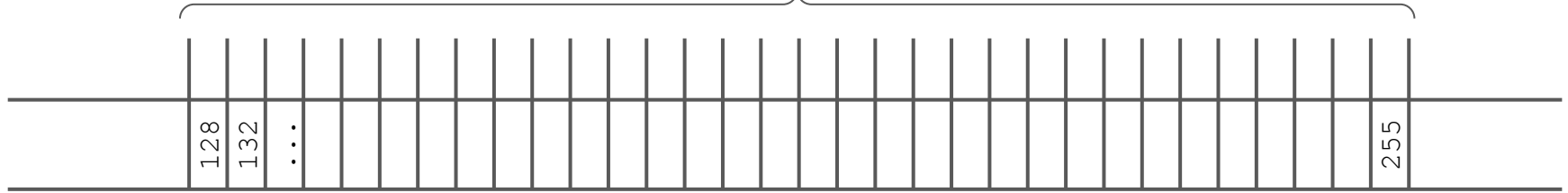
```
int main() {
    // do setup

    add1<<<4,256>>>(dev_a,vector_size);

    // do cleanup
}
```

… block with 256/32 warps ...

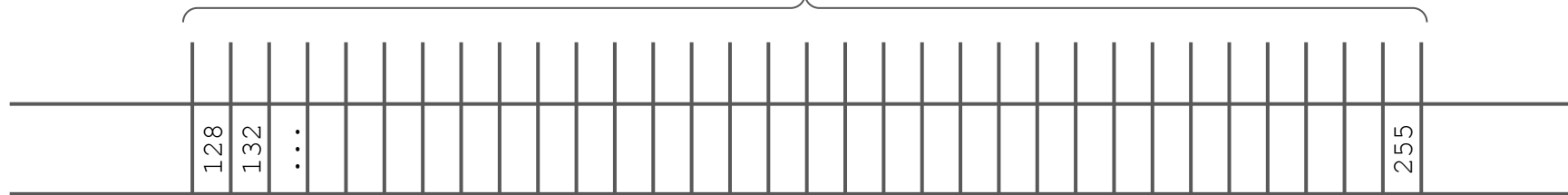128 bytes here. 32 Groups of 4 bytes

128 132 ∷ 255

```
int main() {
    // do setup

    add1<<<4,256>>>(dev_a,vector_size);

    // do cleanup
}
```

… block with 8 warps ...

128 bytes here. 32 Groups of 4 bytes

128 132 :: 255

```
int main() {
    // do setup

    add1<<<4,256>>>(dev_a,vector_size);

    // do cleanup
}
```

1 WARP: 32 threads

… block with 8 warps ...

128 bytes here. 32 Groups of 4 bytes

128 132 .. 255

```
int main() {
    // do setup

    add1<<<4,256>>>(dev_a,vector_size);

    // do cleanup
}
```
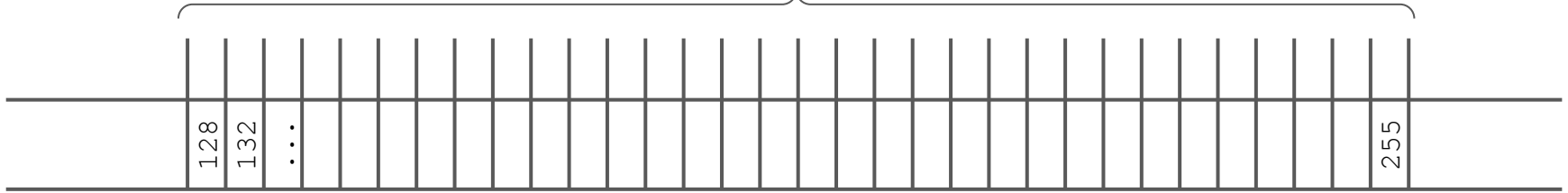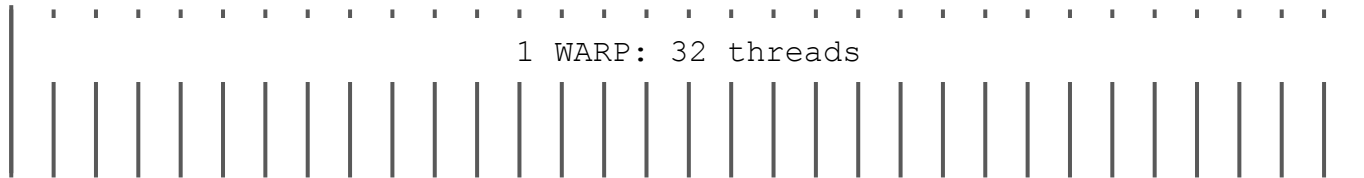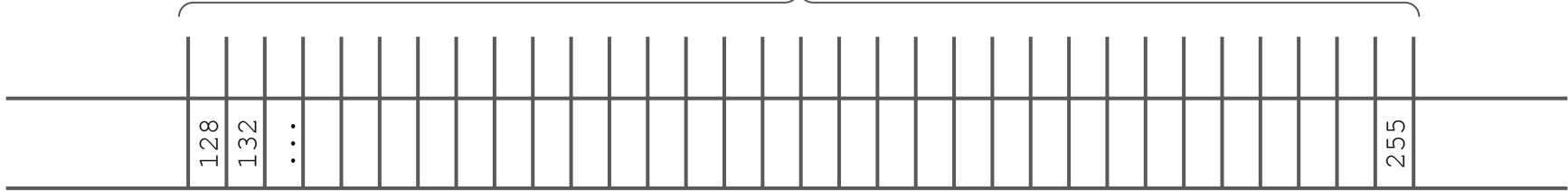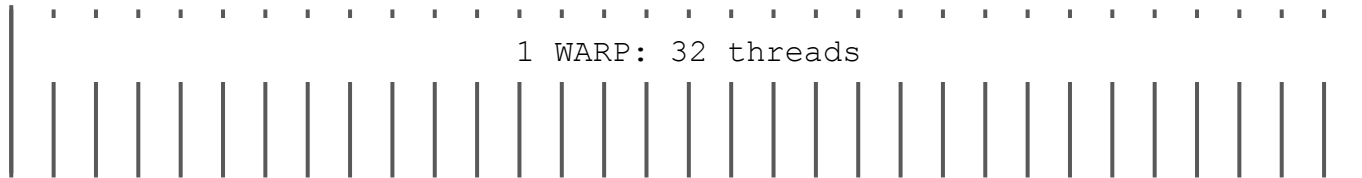
1 WARP: 32 threads

128 bytes here. 32 Groups of 4 bytes

128  132  : :  255

```
__device__ void add1(int* a, int vector_size) {
    id = blockIdx.x*blockDim + threadIdx.x;

    if (id<vector_size)
        a[id]=a[id]+1;
}
```

Each thread reads one element of array **a**

1 WARP: 32 threads

128 bytes here. 32 Groups of 4 bytes

128 132 ... 255

Each thread reads one element
of array **a**

1 WARP: 32 threads

128 bytes here. 32 Groups of 4 bytes

128 132 .. 255

If this happens, all these **32 accesses**
will be serviced by **1 transaction**

1 WARP: 32 threads

128 bytes here. 32 Groups of 4 bytes

128 132 ·· 255

If this happens, all these **32 accesses**
will be serviced by **1 transaction**

As opposed to 32 transactions

1 WARP: 32 threads

128 bytes here. 32 Groups of 4 bytes

128 132 .. 255

If this happens, all these **32 accesses**
will be serviced by **1 transaction**

This is called **Memory Coalescing**!

1 WARP: 32 threads

Another kernel… that has some non-contiguous access based on tid

128 132 ⋮ ⋮ 255

This is also OK on "new" architectures kepler+

1 WARP: 32 threads

Yet Another kernel… that has some non-contiguous access based on tid

128 132 ... 255 256

This is **NOT** ok

**2 transactions** will be necessary

1 WARP: 32 threads

```
Typedef struct myCoordinates {
     int x;
     int y;
     int z;
} Coordinates;


int main () {
     // do stuff

     coordinates * a = malloc(sizeof(Coordinates)*vector_size);

     // do more stuff
}
```
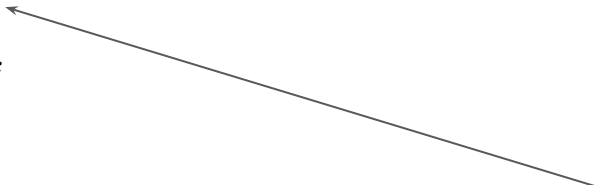
**In your lab today!**

```
Typedef struct myCoordinates {
    int x;
    int y;
    int z;
} Coordinates;


int main () {
    // do stuff

    coordinates * a = malloc(sizeof(Coordinates)*vector_size);

    // do more stuff
}
```

Size of 3 integers:
12 bytes

**In your lab today!**

```
Typedef struct myCoordinates {
    int x;
    int y;
    int z;
} Coordinates;


int main () {
    // do stuff

    coordinates * a = malloc(sizeof(Coordinates)*vector_size);

    // do more stuff
}
```
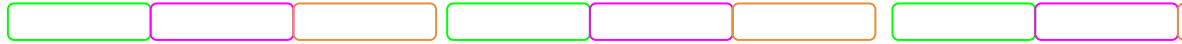
```
Typedef struct myCoordinates {
    int x;
    int y;
    int z;
} Coordinates;


int main () {
    // do stuff

    coordinates * a = malloc(sizeof(Coordinates)*vector_size);

    // do more stuff
}
```

In memory 1 Coordinates object would look like this

**In your lab today!**

```
Typedef struct myCoordinates {
    int x;
    int y;
    int z;
} Coordinates;


int main () {
    // do stuff

    coordinates * a = malloc(sizeof(Coordinates)*vector_size);

    // do more stuff
}
```

In memory Several Coordinates object would look like this

**In your lab today!**

```
Typedef struct myCoordinates {
    int x;
    int y;
    int z;
} Coordinates;


int main () {
    // do stuff

    coordinates * a = malloc(sizeof(Coordinates)*vector_size);

    // do more stuff
}
```
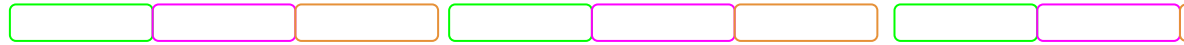
So, accessing only **x (the greens)** would be a strided access

**In your lab today!**

```
Typedef struct myCoordinates {
    int x;
    int y;
    int z;
} Coordinates;


int main () {
    // do stuff

    coordinates * a = malloc(sizeof(Coordinates)*vector_size);

    // do more stuff
}
```

What we want is:

Good luck! ¯\_(ツ)_/¯