# Parallel Computing with GPUs
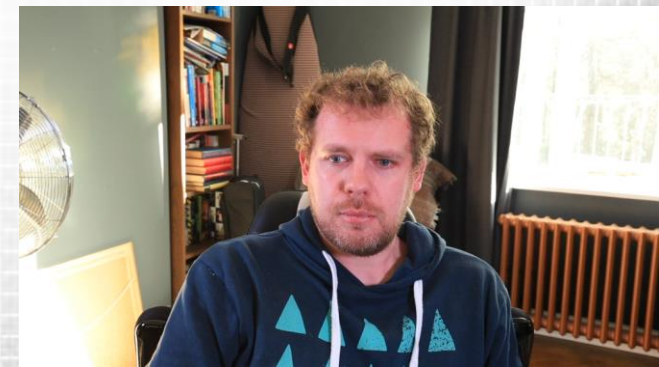
# Performance
# Part 2 – L1 Cache

Dr Paul Richmond

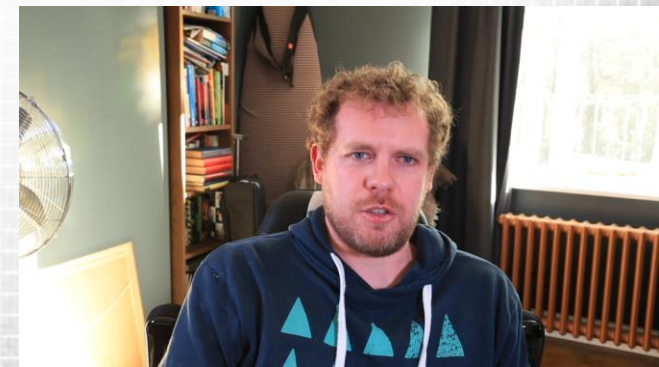http://paulrichmond.shef.ac.uk/teaching/COM4521/

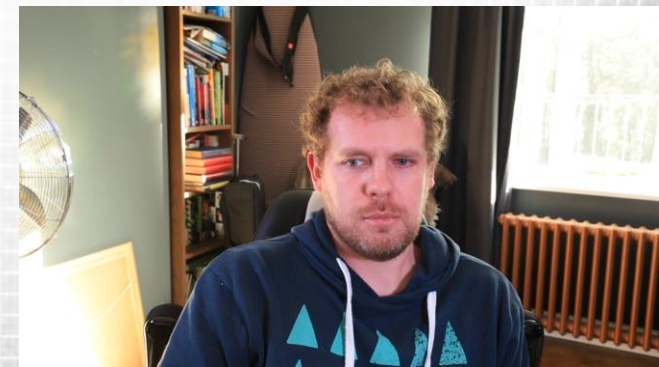# This Lecture (learning objectives)

❑The L1 Cache

    ❑Describe how to enable global L1 caching

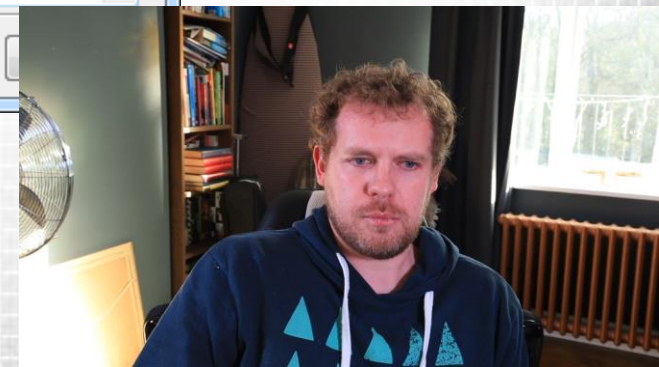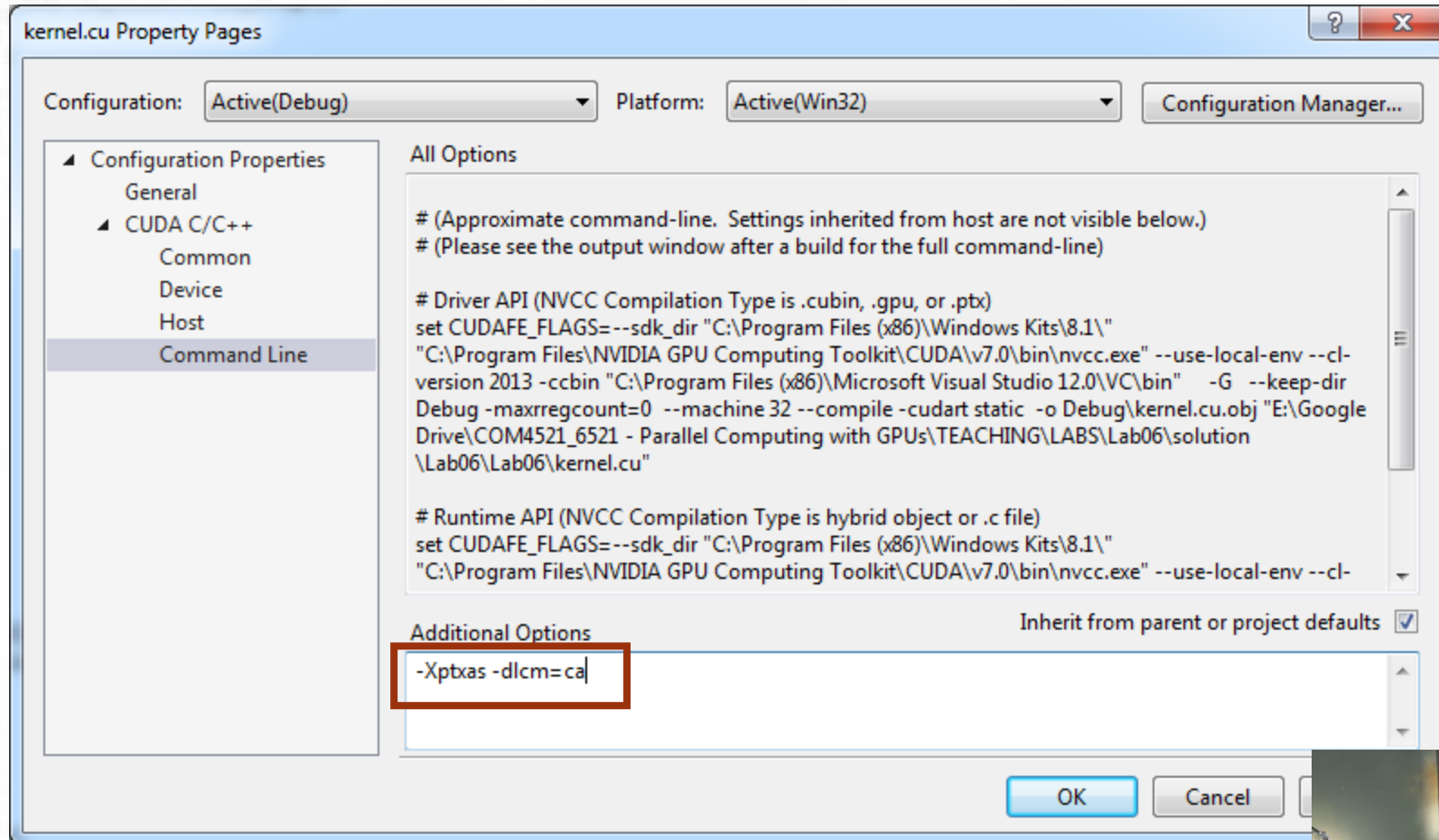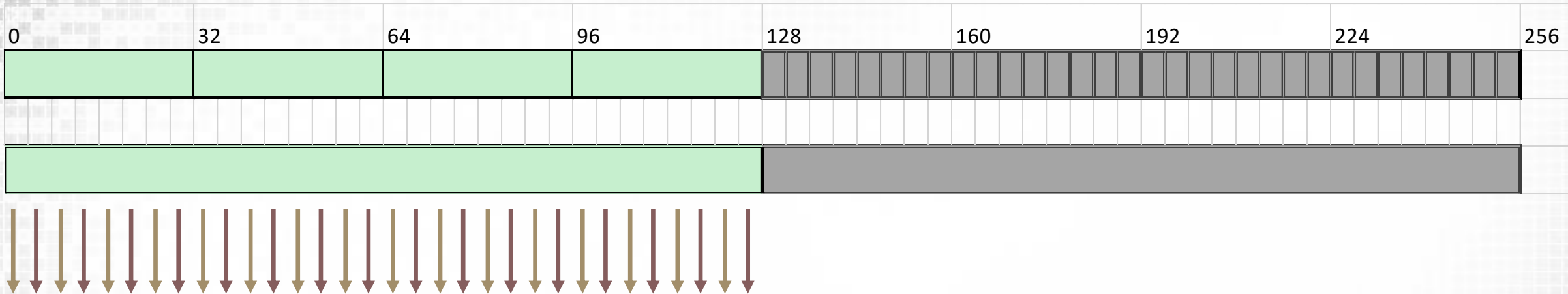    ❑Analyse the performance implications of example access patterns

# L1 Cache

❑Can utilise L1 via unified L1/Texture/Read-only memory

❑Can also be enabled globally

  ❑What affect does this have on performance of memory movement?

   ❑Can be good in certain circumstances

    ❑Coalesced access with adjacent warps reading same data

   ❑Can also be bad

    ❑Un-coalesced access performance is worse

    ❑Increases over-fetch

  ❑Does my card support global L1 Caching?

   ❑Check globalL1CacheSupported and localL1CacheSupported CUDA device properties

    ❑Maxwell 5.2 reports globalL1CacheSupported false when in fact true!

  ❑Enabling L1 caching of global loads

   ❑Pass the `-Xptxas -dlcm=ca` flag to nvcc at compile time

    ❑-dlcm=cg can be used to disable L1 on devices which use it by default
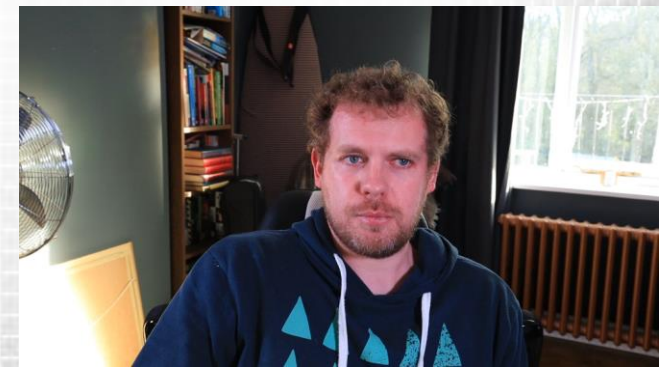
# Enabling L1 Cache in Visual Studio

# L1 Coalesced Memory Access
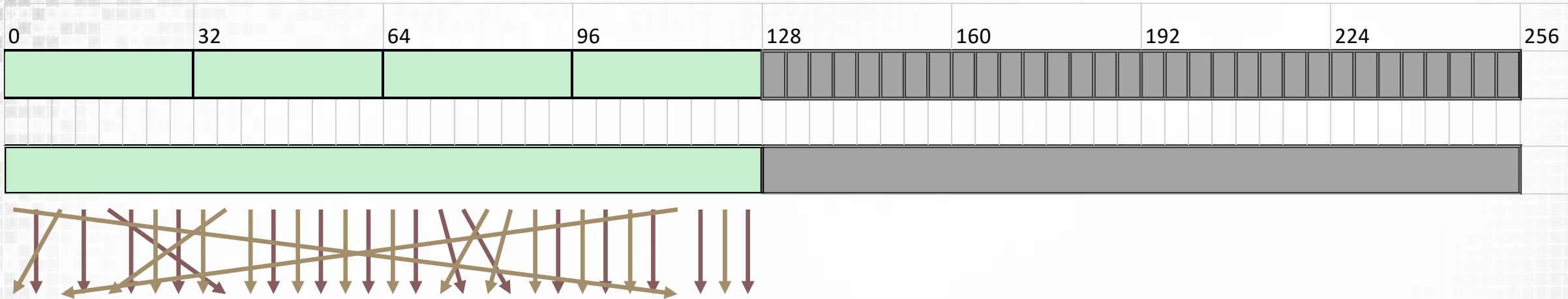
| 0 | 32 | 64 | 96 | 128 | 160 | 192 | 224 | 256 |

```
__global__ void copy(float *odata, float* idata)
    int xid = blockIdx.x * blockDim.x + threadIdx.x;
    odata[xid] = idata[xid];
}
```

❑ All addresses fall in one 128B cache line
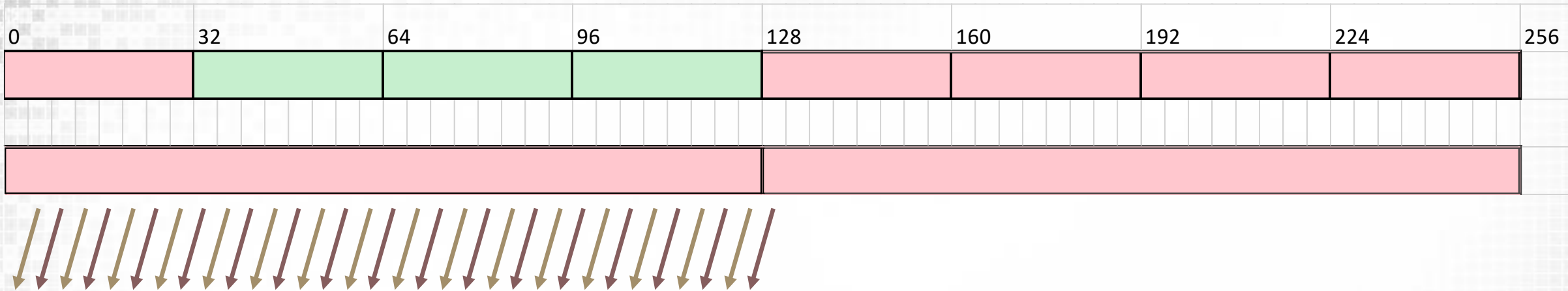   ❑ Single transaction
❑ 100% bus utilisation

# L1 Permuted Memory Access

| 0 | 32 | 64 | 96 | 128 | 160 | 192 | 224 | 256 |
|---|---|---|---|---|---|---|---|---|

❑Any thread within the warp can permute access

❑Same as L2

# L1 Offset Memory Access

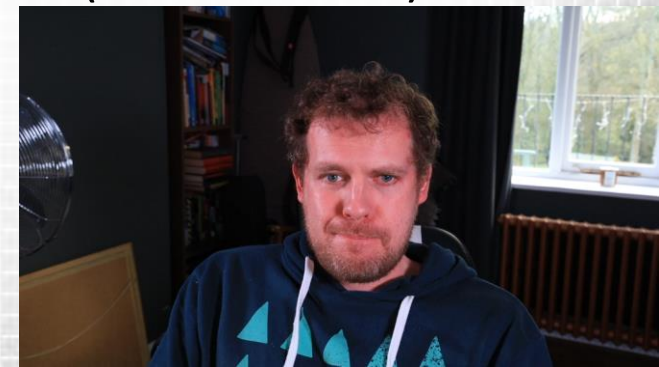| 0 | 32 | 64 | 96 | 128 | 160 | 192 | 224 | 256 |
|---|----|----|----|-----|-----|-----|-----|-----|

```
__global__ void copy(float *odata, float* idata)
    int xid = blockIdx.x * blockDim.x + threadIdx.x + OFFSET;
    odata[xid] = idata[xid];
}
```

❑ If memory accesses are offset then parts of the cache line will be unused (shown in red) e.g.

   ❑ 2 transactions of 256B of which 128B is required: 50% utilisation

❑ For strided and random access performance is much worse with L1

# Summary

❑The L1 Cache

    ❑Describe how to enable global L1 caching

    ❑Analyse the performance implications of example access patterns

❑Next Lecture: Occupancy