## CUDA

| | |
|---|---|
| Kernel execution | `kernel Nam e<< <nu m_b locks, num_th reads, shared, stream>> (args);` |
| Kernel definition | `__global__ kernel Nam e(a rgs){…}` |
| Create memory on device | `cudaMa llo c(m emPtr, sizeIn Bytes);` |
| Free device memory | `cudaFr ee( mem Ptr);` |
| Create **pinned** memory on host | `cudaMa llo cHo st( memPtr, size_i n_b ytes);` |
| Free pinnen memory | `cudaFr eeH ost (me mPtr);` |
| Create stream | `cudaSt re a m Cr eat e( s t ream) ;` |
| Destroy stream | `cudaSt rea mDe str oy( str eam);` |
| Wait for Device | `cudaDe vic eSy nch ron ize();` |
| Wait for Stream | `cudaSt rea mSy nch ron ize (st ream);` |
| Copy data | `cudaMe mcp y(dst, src, numEle ments, cudaMe mcp yDe fautl);` |
| Async mem copy | `cudaMe mcp yAs ync (dst, src, numEle ments, cudaMe mcp yDe fautl, stream);` |
| Static device variable | `__device__ type name = value;` |
| Static shared mem | `__shared__ type name[s ize];` |
| Dynamic shared mem | `extern __shared__ type* name;` |
| Constant device mem | `__cons tant__ type name;` |

## CUDA (cont)

| | |
|---|---|
| Copy to mem constant | `cudaMe mcp yTo Sym bol (dst, src, sizeIn Byt` |
| Documentation | https://docs.nvidia.com/cuda/cuda-runtime-api/index.html |
| Device-only function | `__device__ auto functi onN ame (args) -> ret` |
| Host-only function | `__host__ auto functi onN ame (args) -> retur` |
| Host and device function | `__device__ __host__ auto functi onN ame (arg e{…}` |

## Thrust

| | |
|---|---|
| Universal Vector | `thrust ::u niv ers al_ vec tor <ty pe> nam e(s ize);` |

By **mavieth** (mavieth)
cheatography.com/mavieth/

Not published yet.
Last updated 27th February, 2026.
Page 1 of 1.