

Performance portability via Nim metaprogramming

Xiao-Yong Jin and James C. Osborn
Argonne Leadership Computing Facility, Argonne National Laboratory

A DOE ECP effort from Lattice QCD

August 23
DOE COE Performance Portability Meeting 2017
Denver, Colorado

“Complexity” seems to be a lot like “energy”: you can transfer it from the end user to one/some of the other players, but the total amount seems to remain pretty much constant for a given task.

—Ran, 5 Mar 2000

Outline

- Benchmark code in Nim
- Application domain: Lattice Gauge Theory
- Nim programming language
- onGPU
- Benchmark result
 - P100 + XEON
 - KNL

Benchmark

$$X_{ij}^l \leftarrow X_{ij}^l + Y_{ik}^l \times Z_{kj}^l$$

```
3: proc test(vecLen, memLen: static[int]; N: int) =
4:   var
5:     x = newColorMatrixArray(vecLen, memLen, N)
6:     y = newColorMatrixArray(vecLen, memLen, N)
7:     z = newColorMatrixArray(vecLen, memLen, N)

32:  threads:
33:     x := 0
34:     y := 1
35:     z := 2
36:     timeit "CPU": x += y * z
37:
38:     timeit "GPU5":
39:       onGpu(N, 32):
40:         x += y * z
41:     timeit "GPU6": onGpu(N, 64): x += y * z
42:     timeit "GPU7": onGpu(N, 128): x += y * z
43:
44:  threads: timeit "CPU": x += y * z
```

Define complex
3×3 matrix field

CPU threads

Set diagonal elements

Benchmark

Run statement block on GPU

Rerun with
different T/B

Back to CPU threads

CPU threads

- Takes a block of code statements
- Wraps in a function with lexically scoped thread local objects
- References to variables outside the code block are managed by Nim
- Runs the function under `omp parallel` directive
- A custom iterator over the array indices takes care of actual data parallel operations

`x += y * z`  `x[i] += (y * z)[i]`

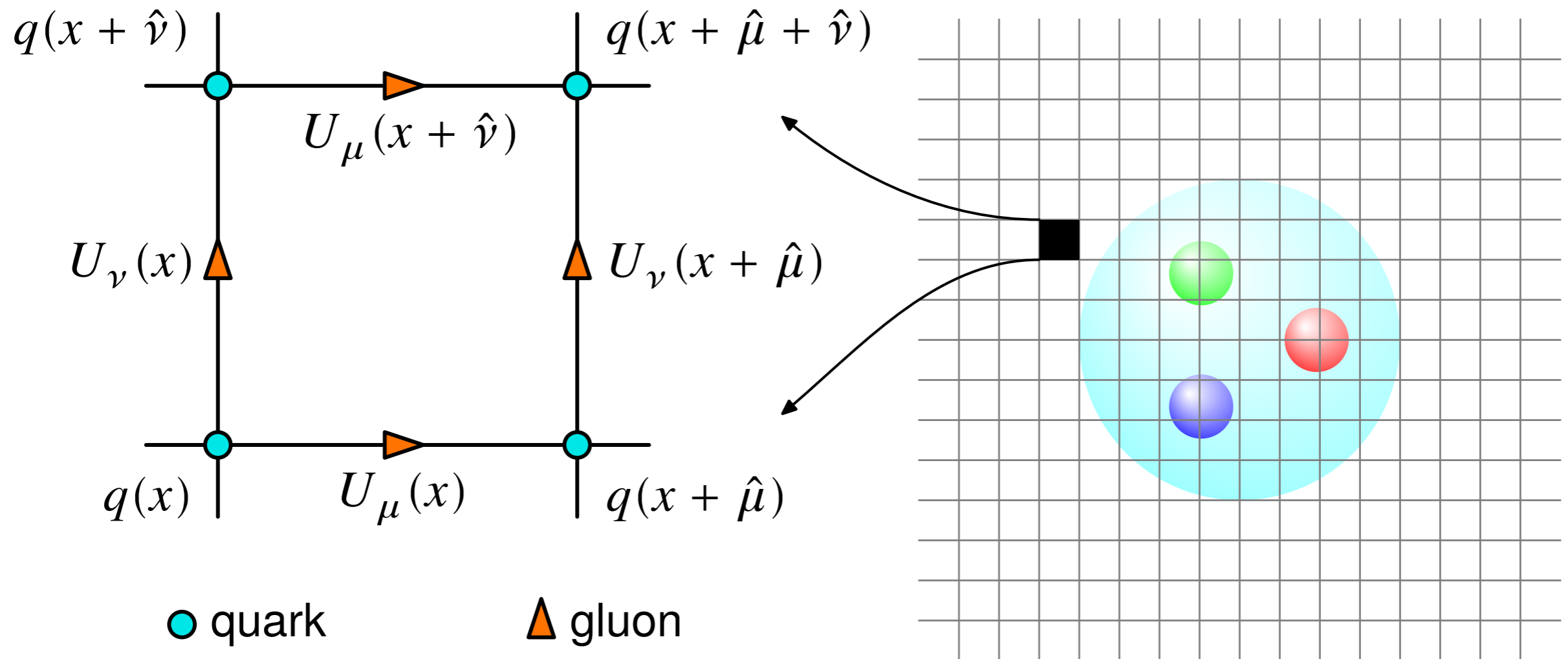
AST based overloading for data parallel ops

```
1: type ArrayIndex* = SomeInteger or ShortVectorIndex
2:
3: template indexArray*(x: ArrayObj, i: ArrayIndex): untyped =
4:   x.p[i]
5:
6: macro indexArray*(x: ArrayObj{call}, y: ArrayIndex): untyped =
7:   result = newCall(ident($x[0]))
8:   for i in 1..<x.len:
9:     let xi = x[i]
10:    result.add( quote do:
11:      indexArray(`xi`, `y`) )
12:
13: template []*(x: ArrayObj, i: ArrayIndex): untyped =
14:   indexArray(x, i)
```

- When an ArrayObj is indexed, if the object is a function call, the indexing goes inside the call

$x[i] += (y * z)[i] \rightarrow x[i] += y[i] * z[i]$

Lattice gauge theory



- Large 4D (5D) grid of small vectors/matrices with homogeneous stencil operations — large sparse linear algebra

Nim

- Modern (since 2008) language
- “Efficient Expressive Elegant”
- Statically typed systems language (full access to low-level objects & code) with type inference
- Generates C or C++ code & compile with any compiler
- Integrated build system (no Makefile necessary): copy main program, modify, compile
- <https://nim-lang.org>

Nim—both low-level and high-level

- Low-level efficiency
 - Can manually manage memory instead of GC
 - Cross module inlining and constant unfolding
 - Whole program dead code elimination
- High-level wrappers & libraries
 - gmp, bignum, nimblas, linalg(LAPACK), ...
 - bindings to GTK2, the Windows API, the POSIX API, OpenGL, SDL, Cairo, Python, Lua, TCL, X11, libzip, PCRE, libcurl, MySQL, SQLite, ...
 - `exportC` to create static/dynamic libraries
- NimScript: shell-like scripting
 - Used in compiler for compile-time evaluation
 - Available to plug in to application and can interface with rest of application

Nim—metaprogramming

- **Templates:** in-line code substitutions, also allows overloading, completely hygienic (if desired)
- **Generics:** applies to types, procedures, templates, and macros also allows type-classes, concepts
- **Macros:** similar to Lisp: syntax tree of arguments passed to macro at compile time (type checked or untyped)
- **AST based overloading:** allows specialization based on the AST of the arguments

New framework: QEX (Quantum EXpressions)

- Data parallel library for tensor objects on a lattice including shifts, reductions
- Mostly in Nim, with USQCD SciDAC C libraries
- High level interface in development
- Available on <https://github.com/jcosborn/qex>
- Performance portability study: cudanim
 - Supports arrays on both CPU and GPU
 - Checkout <https://github.com/jcosborn/cudanim>

onGpu

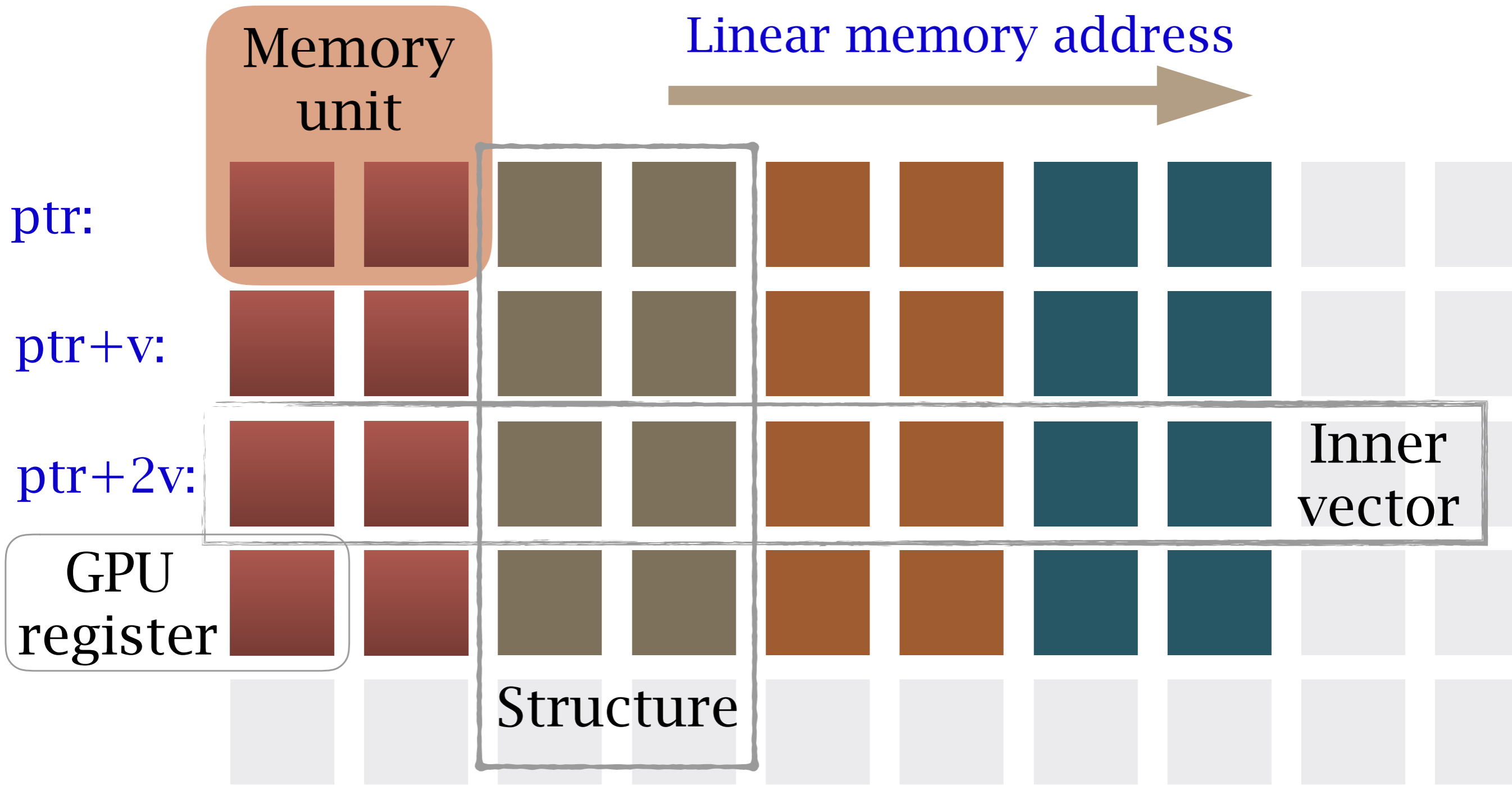
```
1: template onGpu*(nn, tpb: untyped, body: untyped): untyped =
2:   block: Copy memory to GPU
3:     var v = packVars(body, getGpuPtr) if necessary
4:     type ByCopy { .bycopy. } [T] = object collect pointers in v
5:     d: T
6:     proc kern(xx: ByCopy[type(v)]) { .cudaGlobal. } =
7:       template deref(k: int): untyped = xx.d[k]
8:       substVars(body, deref) Generate kernel function
9:     let ni = nn.int32
10:    let threadsPerBlock = tpb.int32
11:    let blocksPerGrid = (ni+threadsPerBlock-1) div threadsPerBlock
12:    cudaLaunch(kern, blocksPerGrid, threadsPerBlock, v)
13:    discard cudaDeviceSynchronize()
14: template onGpu*(nn: untyped, body: untyped): untyped =
15:   onGpu(nn, 64, body)
16: template onGpu*(body: untyped): untyped =
17:   onGpu(512*64, 64, body)
```

The generated kern

```
1: proc kern(xx: ByCopy[type(v)])
2:   { .codegenDecl: "__global__ $# $$$#" .} =
   # Some definitions omitted
13: inlineProcs: ← Inlines all
14:   template deref(k: int): untyped = procedure calls
15:     xx.d[k]
16:   substVars((x += y * z), deref)
```

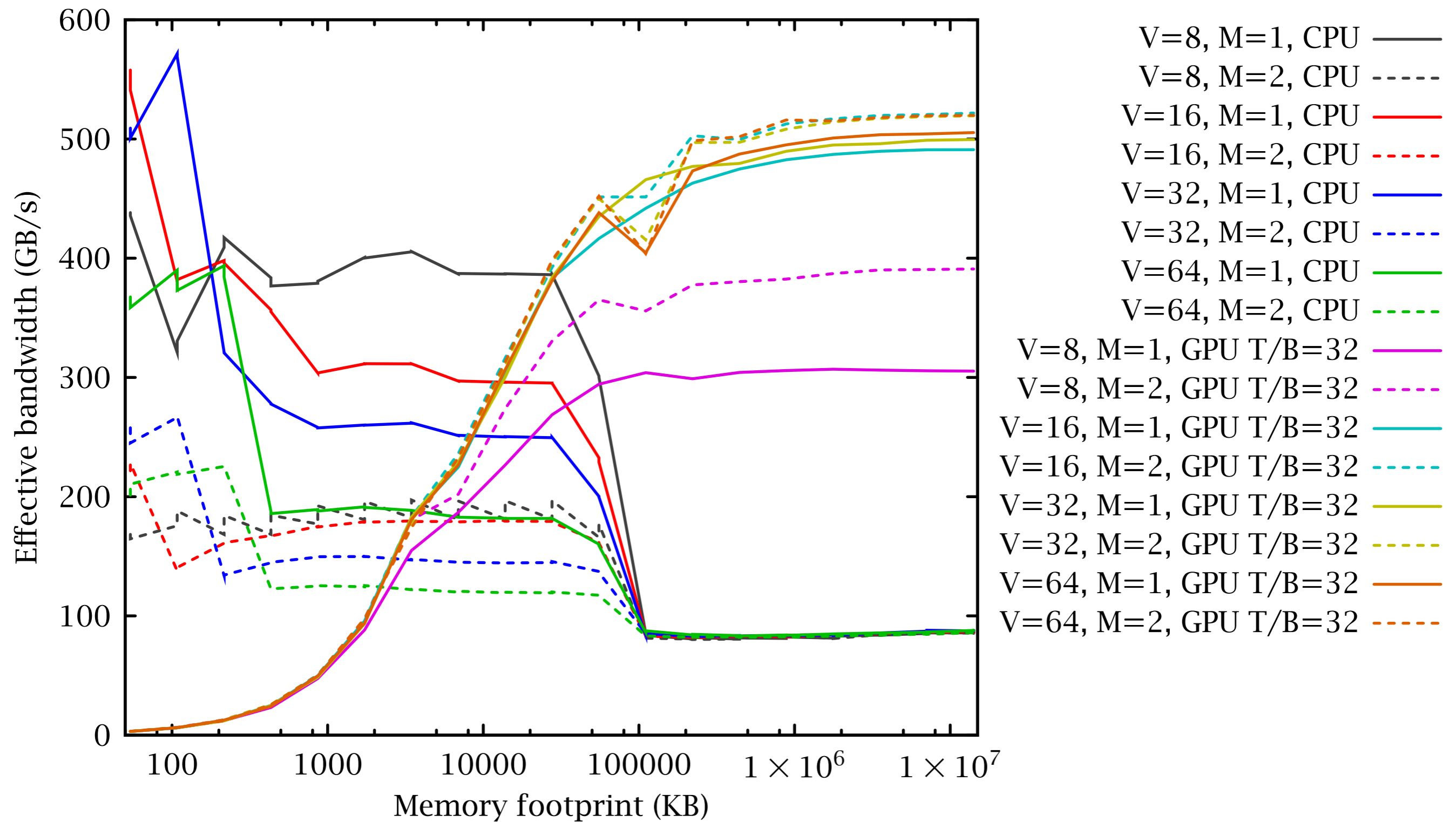
- xx is an object holding pointers to GPU memory
- substVars with the help of deref transforms

$x += y * z \rightarrow \text{deref}(0) += \text{deref}(1) * \text{deref}(2)$
 $\rightarrow \text{xx.d}[0] += \text{xx.d}[1] * \text{xx.d}[2]$



Coalesced in-memory data layout

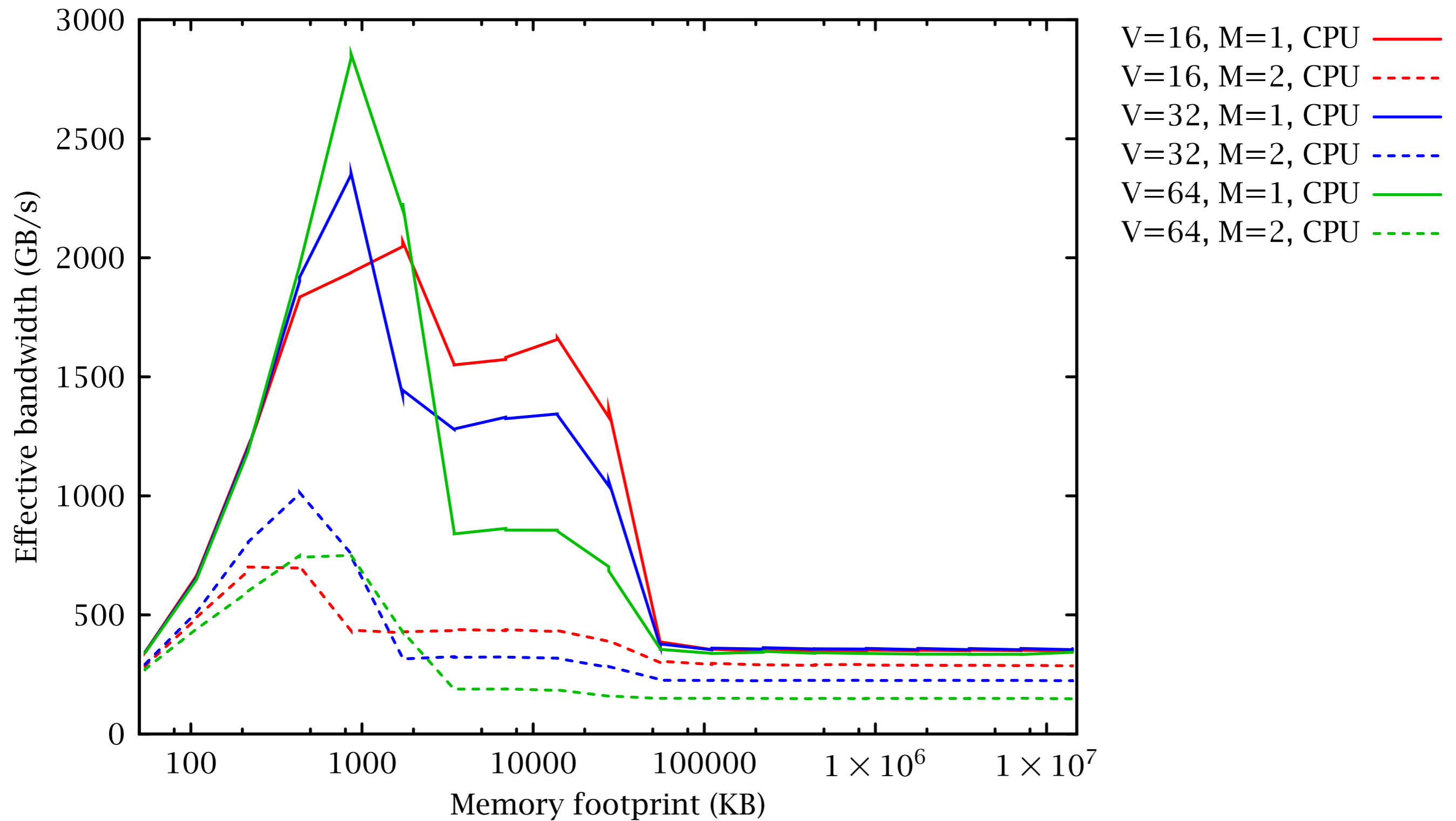
AoSoAoS



Tesla P100 +
2x Xeon E5-2687WV2

$$X_{ij}^l \leftarrow X_{ij}^l + Y_{ik}^l \times Z_{kj}^l$$

CUDA 8.0.61 + GCC 4.8.5



KNL (Xeon Phi 7210)
Flat MCDRAM

$$X_{ij}^l \leftarrow X_{ij}^l + Y_{ik}^l \times Z_{kj}^l$$

GCC 7.1.0

Summary & Outlook

- Nim metaprogramming helps hiding architecture differences under a unified data parallel API
- Toy benchmark saturates GPU bandwidth
- Considering API to use both CPU & GPU in a heterogeneous setting
- Many possibilities of using AST
 - Apply AST based optimizations (inlining, loop unrolling, temporary variable elimination), across multiple statements
 - Specialization at different levels
 - Craft application specific AST transformations
 - Help general purpose compiler with application specifics

Backup slides

GPU offloading (1 of 4)

```
1: template cudaDefs(body: untyped): untyped {.dirty.} =
2:   var gridDim{.global,importC,noDecl.}: CudaDim3
3:   var blockIdx{.global,importC,noDecl.}: CudaDim3
4:   var blockDim{.global,importC,noDecl.}: CudaDim3
5:   var threadIdx{.global,importC,noDecl.}: CudaDim3
6:   template getGridDim: untyped {.used.} = gridDim
7:   template getBlockIdx: untyped {.used.} = blockIdx
8:   template getBlockDim: untyped {.used.} = blockDim
9:   template getThreadIdx: untyped {.used.} = threadIdx
10:  template getThreadNum: untyped {.used.} = blockDim.x * blockIdx.x + threadIdx.x
11:  template getNumThreads: untyped {.used.} = gridDim.x * blockDim.x
12:  bind inlineProcs
13:  inlineProcs:
14:    body
```

- Overloaded template definitions
- Inline all Nim procedures in `body` to prepare for generating kernel function

GPU offloading (2 of 4)

```
34: macro cuda*(s,p: untyped): auto =
35:   let ss = s.strVal
36:   p.expectKind nnkProcDef
37:   result = p
38:   result.addPragma parseExpr("{.codegenDecl:\\"&ss&" $# $$$\".}") [0]
39:   result.body = getAst(cudaDefs(result.body))
40:   var s1 = newStmtList()
41:   s1.add( quote do:
42:     {.push checks: off.}
43:     {.push stacktrace: off.} )
44:   s1.add result
45:   result = s1
46: template cudaGlobal*(p: untyped): auto = cuda("__global__",p)
```

- Convert a procedure definition, p, with the overloaded templates in `cudaDefs`
- Turn it in to a `__global__` kernel

GPU offloading (3 of 4)

```
16: template cudaLaunch*(p: proc; blocksPerGrid, threadsPerBlock: SomeInteger;
17:                        arg: varargs[pointer, dataAddr]) =
18:   var pp: proc = p
19:   var gridDim, blockDim: CudaDim3
20:   gridDim.x = blocksPerGrid
21:   gridDim.y = 1
22:   gridDim.z = 1
23:   blockDim.x = threadsPerBlock
24:   blockDim.y = 1
25:   blockDim.z = 1
26:   var args: array[arg.len, pointer]
27:   for i in 0..
```

- Calls the cuda function, `cudaLaunchKernel`, with a passed in procedure, `p`

GPU offloading (4 of 4)

```
48: template onGpu*(nn,tpb: untyped, body: untyped): untyped =
49:   block:
50:     var v = packVars(body, getGpuPtr)
51:     type ByCopy {bycopy.} [T] = object
52:     d: T
53:     proc kern(xx: ByCopy[type(v)]) {.cudaGlobal.} =
54:       template deref(k: int): untyped = xx.d[k]
55:       substVars(body, deref)
56:       let ni = nn.int32
57:       let threadsPerBlock = tpb.int32
58:       let blocksPerGrid = (ni+threadsPerBlock-1) div threadsPerBlock
59:       cudaLaunch(kern, blocksPerGrid, threadsPerBlock, v)
60:       discard cudaDeviceSynchronize()
61: template onGpu*(nn: untyped, body: untyped): untyped = onGpu(nn, 64, body)
62: template onGpu*(body: untyped): untyped = onGpu(512*64, 64, body)
```

- Take a body of code chunk and put it in a kernel definition, kern
- kern calls cudaGlobal to setup other definitions, and takes care of syncing CPU memory to GPU