

Easy, Effective, Efficient: GPU Programming in Python with PyOpenCL and PyCUDA

Andreas Klöckner

Courant Institute of Mathematical Sciences
New York University

PASI: The Challenge of Massive Parallelism
Lecture 4 · January 8, 2011

Outline

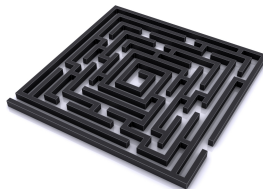
- 1 PyCUDA
- 2 Automatic GPU Programming
- 3 GPU-DG: Challenges and Solutions



Lab Solutions

Lab solutions:

- Lab 1 yesterday:
Sorry, posted wrong tarball
(I think)
- Will post lab solutions after second
lab today:
[http://tiker.net/tmp/
pasi-lab-solution.tar.gz](http://tiker.net/tmp/pasi-lab-solution.tar.gz)



Outline

- 1 PyCUDA
- 2 Automatic GPU Programming
- 3 GPU-DG: Challenges and Solutions



Whetting your appetite

```
1 import pycuda.driver as cuda
2 import pycuda.autoint, pycuda.compiler
3 import numpy
4
5 a = numpy.random.randn(4,4).astype(numpy.float32)
6 a_gpu = cuda.mem_alloc(a.nbytes)
7 cuda.memcpy_htod(a_gpu, a)
```

[This is `examples/demo.py` in the PyCUDA distribution.]

Whetting your appetite

```
1 mod = pycuda.compiler.SourceModule("""
2     __global__ void twice(float *a)
3     {
4         int idx = threadIdx.x + threadIdx.y*4;
5         a[idx] *= 2;
6     }
7     """)
8
9 func = mod.get_function("twice")
10 func(a_gpu, block=(4,4,1))
11
12 a_doubled = numpy.empty_like(a)
13 cuda.memcpy_dtoh(a_doubled, a_gpu)
14 print a_doubled
15 print a
```

Whetting your appetite

```

1 mod = pycuda.compiler.SourceModule("""
2     __global__ void twice(float *a)
3     {
4         int idx = threadIdx.x + threadIdx.y*4;
5         a[idx] *= 2;
6     }
7     """)
8
9 func = mod.get_function("twice")
10 func(a_gpu, block=(4,4,1))
11
12 a_doubled = numpy.empty_like(a)
13 cuda.memcpy_dtoh(a_doubled, a_gpu)
14 print a_doubled
15 print a

```

Compute kernel

Whetting your appetite, Part II

Did somebody say “Abstraction is good”?



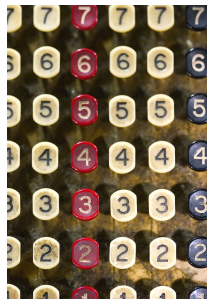
Whetting your appetite, Part II

```
1 import numpy
2 import pycuda.autoinit
3 import pycuda.gpuarray as gpuarray
4
5 a_gpu = gpuarray.to_gpu(
6     numpy.random.randn(4,4).astype(numpy.float32))
7 a_doubled = (2*a_gpu).get()
8 print a_doubled
9 print a_gpu
```

gpuarray: Simple Linear Algebra

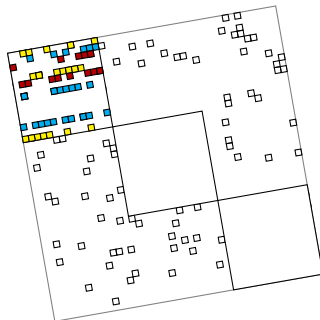
pycuda.gpuarray:

- Meant to look and feel just like numpy.
 - `gpuarray.to_gpu(numpy_array)`
 - `numpy_array = gpuarray.get()`
- `+`, `-`, `*`, `/`, `fill`, `sin`, `exp`, `rand`,
basic indexing, `norm`, inner product, ...
- Mixed types (`int32 + float32 = float64`)
- `print gpuarray` for debugging.
- Allows access to raw bits
 - Use as kernel arguments, textures, etc.



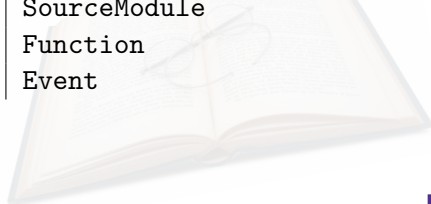
Sparse Matrix-Vector on the GPU

- New feature in 0.94:
Sparse matrix-vector multiplication
- Uses “packeted format”
by Garland and Bell (also
includes parts of their code)
- Integrates with `scipy.sparse`.
- Conjugate-gradients solver
included
 - Deferred convergence
checking



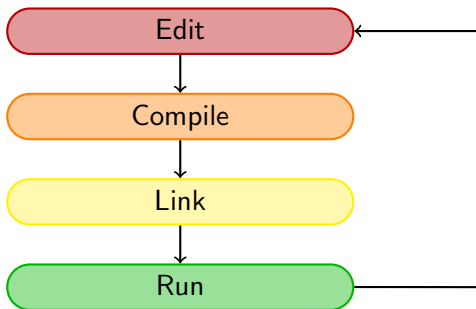
PyOpenCL \leftrightarrow PyCUDA: A (rough) dictionary

PyOpenCL	PyCUDA
Context	Context
CommandQueue	Stream
Buffer	mem_alloc / DeviceAllocation
Program	SourceModule
Kernel	Function
Event (eg. enqueue_marker)	Event



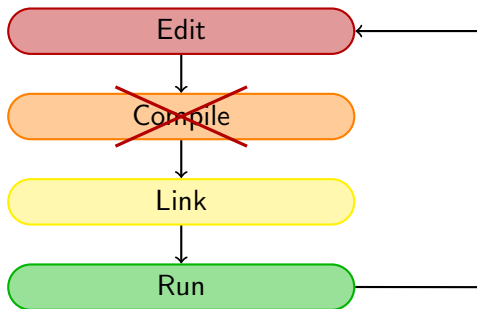
Scripting: Interpreted, not Compiled

Program creation workflow:



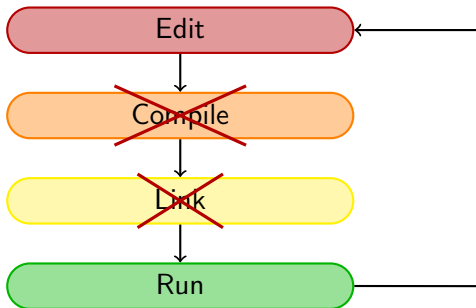
Scripting: Interpreted, not Compiled

Program creation workflow:

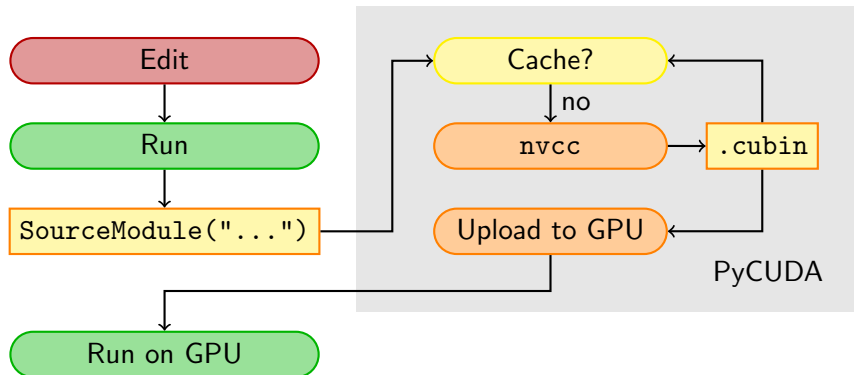


Scripting: Interpreted, not Compiled

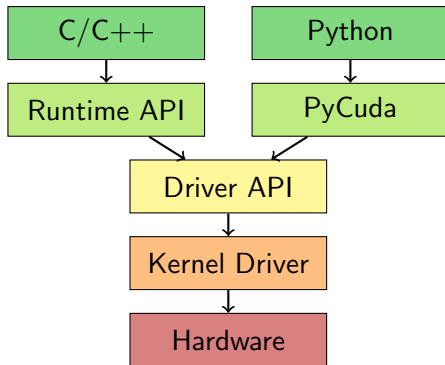
Program creation workflow:



PyCUDA: Workflow



PyCUDA in the CUDA ecosystem



CUDA has two Programming Interfaces:

- “Runtime” high-level (separate install)
- “Driver” low-level (`libcuda.so`, comes with GPU driver)

PyCUDA: Vital Information

- <http://mathematician.de/software/pycuda>
- Complete documentation
- X Consortium License
(no warranty, free for all use)
- Convenient abstractions
Array, Fast Vector Math, Reductions
- Requires: numpy, Python 2.4+
(Win/OS X/Linux)



Outline

- 1 PyCUDA
- 2 Automatic GPU Programming
- 3 GPU-DG: Challenges and Solutions



Automating GPU Programming

GPU programming can be time-consuming, unintuitive and error-prone.

- Obvious idea: Let the computer do it.
- One way: Smart compilers



Automating GPU Programming

GPU programming can be time-consuming, unintuitive and error-prone.

- Obvious idea: Let the computer do it.
- One way: Smart compilers
 - GPU programming requires complex tradeoffs
 - Tradeoffs require heuristics
 - Heuristics are fragile



Automating GPU Programming

GPU programming can be time-consuming, unintuitive and error-prone.

- Obvious idea: Let the computer do it.
- One way: Smart compilers
 - GPU programming requires complex tradeoffs
 - Tradeoffs require heuristics
 - Heuristics are fragile
- Another way: Dumb enumeration
 - Enumerate loop slicings
 - Enumerate prefetch options
 - Choose by running resulting code on actual hardware



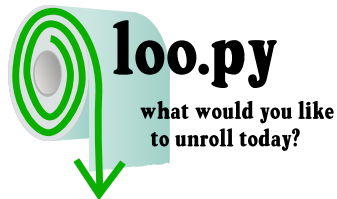
Loo.py Example

Empirical GPU loop optimization:

```
a, b, c, i, j, k = [var(s) for s in "abcijk"]
n = 500
k = make_loop_kernel([
    LoopDimension("i", n),
    LoopDimension("j", n),
    LoopDimension("k", n),
], [
    (c[i+n*j], a[i+n*k]*b[k+n*j])
])

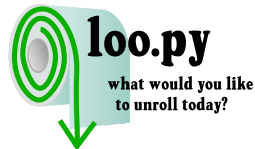
gen_kwargs = {
    "min_threads": 128,
    "min_blocks": 32,
}
```

→ Ideal case: Finds 160 GF/s kernel
without human intervention.



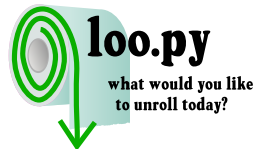
Loo.py Status

- Limited scope:
 - Require input/output separation
 - Kernels must be expressible using “loopy” model
(i.e. indices decompose into “output” and “reduction”)
 - Enough for DG, LA, FD, ...



Loo.py Status

- Limited scope:
 - Require input/output separation
 - Kernels must be expressible using “loopy” model
(i.e. indices decompose into “output” and “reduction”)
 - Enough for DG, LA, FD, ...
- Kernel compilation limits trial rate
- Non-Goal: Peak performance
- Good results currently for dense linear algebra and (some) DG subkernels



Outline

1 PyCUDA

2 Automatic GPU Programming

3 GPU-DG: Challenges and Solutions

- Introduction
- Challenges
- Benefits of Metaprogramming
- GPU-DG: Performance and Generality
- Viscous Shock Capture



Outline

1 PyCUDA

2 Automatic GPU Programming

3 GPU-DG: Challenges and Solutions

- Introduction
- Challenges
- Benefits of Metaprogramming
- GPU-DG: Performance and Generality
- Viscous Shock Capture



Discontinuous Galerkin Method

Let $\Omega := \bigcup_i D_k \subset \mathbb{R}^d$.



Discontinuous Galerkin Method

Let $\Omega := \bigcup_i D_k \subset \mathbb{R}^d$.



Goal

Solve a *conservation law* on Ω :

$$u_t + \nabla \cdot F(u) = 0$$

Discontinuous Galerkin Method

Let $\Omega := \bigcup_i D_k \subset \mathbb{R}^d$.



Goal

Solve a *conservation law* on Ω :

$$u_t + \nabla \cdot F(u) = 0$$

Example

Maxwell's Equations: EM field: $E(x, t)$, $H(x, t)$ on Ω governed by

$$\partial_t E - \frac{1}{\varepsilon} \nabla \times H = -\frac{j}{\varepsilon},$$

$$\nabla \cdot E = \frac{\rho}{\varepsilon},$$

$$\partial_t H + \frac{1}{\mu} \nabla \times E = 0,$$

$$\nabla \cdot H = 0.$$

Discontinuous Galerkin Method

Multiply by test function, integrate by parts:

$$\begin{aligned} 0 &= \int_{D_k} u_t \varphi + [\nabla \cdot F(u)] \varphi \, dx \\ &= \int_{D_k} u_t \varphi - F(u) \cdot \nabla \varphi \, dx + \int_{\partial D_k} (\hat{n} \cdot F)^* \varphi \, dS_x, \end{aligned}$$

Substitute in basis functions, introduce elementwise stiffness, mass, and surface mass matrices S , M , M_A :

$$\partial_t u^k = - \sum_{\nu} D^{\partial_{\nu}, k} [F(u^k)] + L^k [\hat{n} \cdot F - (\hat{n} \cdot F)^*]|_{A \subset \partial D_k}.$$

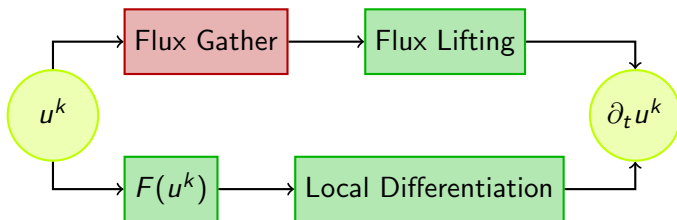
For straight-sided simplicial elements:

Reduce $D^{\partial_{\nu}}$ and L to reference matrices.



Decomposition of a DG operator into Subtasks

DG's execution decomposes into two (mostly) separate branches:



Green: Element-local parts of the DG operator.

DG on GPUs: Possible Advantages



DG on GPUs: Why?

- GPUs have deep Memory Hierarchy
 - The majority of DG is local.
- Compute Bandwidth \gg Memory Bandwidth
 - DG is arithmetically intense.
- GPUs favor dense data.
 - Local parts of the DG operator are dense.

DG on the GPU: What are we trying to achieve?

Objectives:

- *Main: Speed*
Reduce need for compute-bound clusters
- *Secondary: Generality*
Be applicable to many problems
- *Tertiary: Ease-of-Use*
Hide complexity of GPU hardware

Setting (for now):

- Specialize to straight-sided simplices
- Optimize for (but don't specialize to) tetrahedra (ie. 3D)
- Optimize for “medium” order (3...5)



Outline

1 PyCUDA

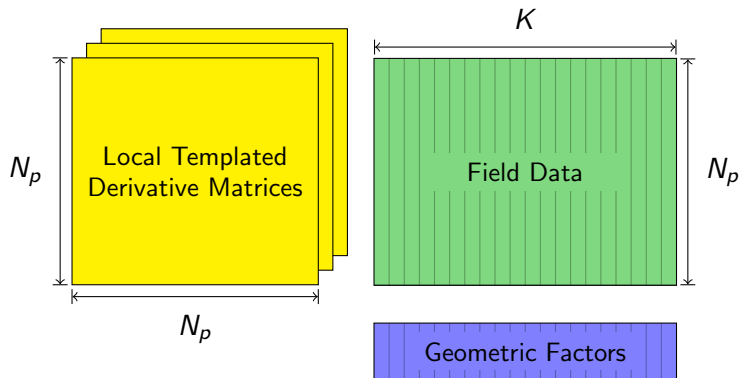
2 Automatic GPU Programming

3 GPU-DG: Challenges and Solutions

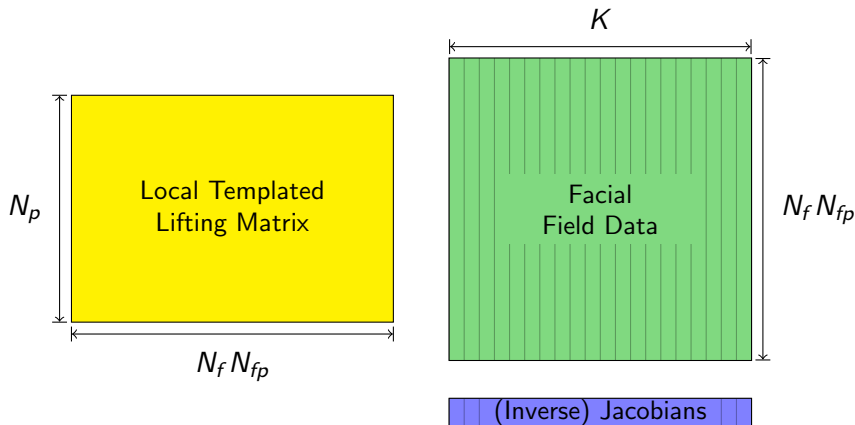
- Introduction
- **Challenges**
- Benefits of Metaprogramming
- GPU-DG: Performance and Generality
- Viscous Shock Capture



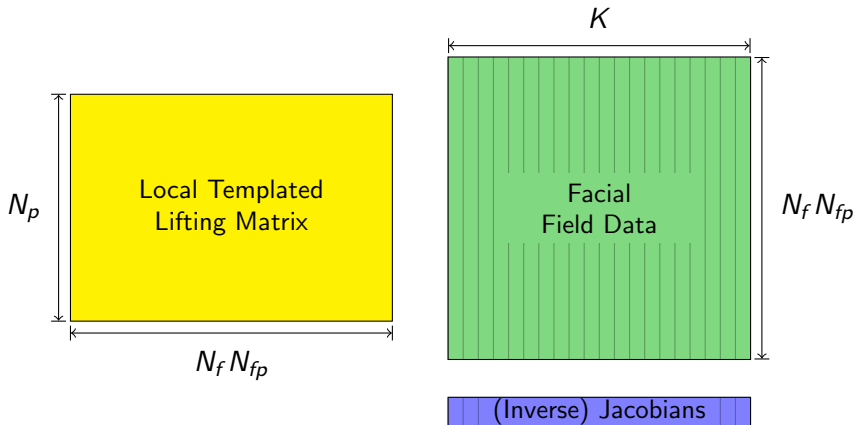
Element-Local Operations: Differentiation



Element-Local Operations: Lifting

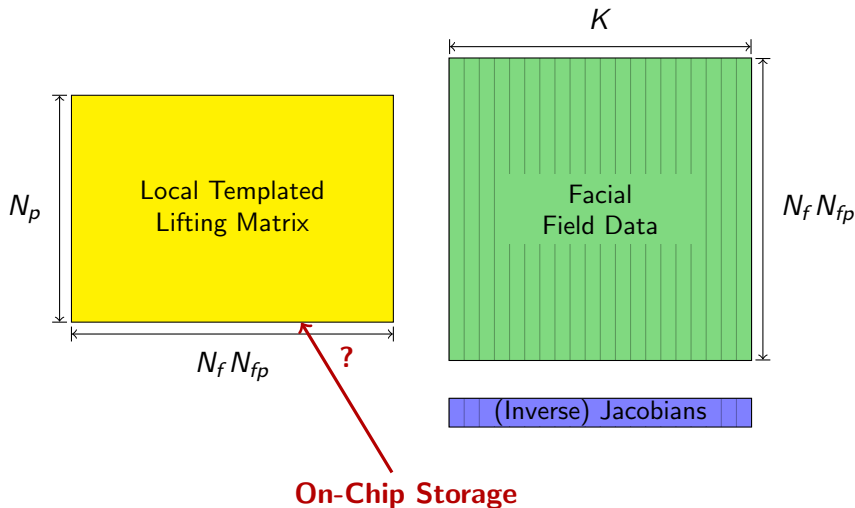


Element-Local Operations: Lifting

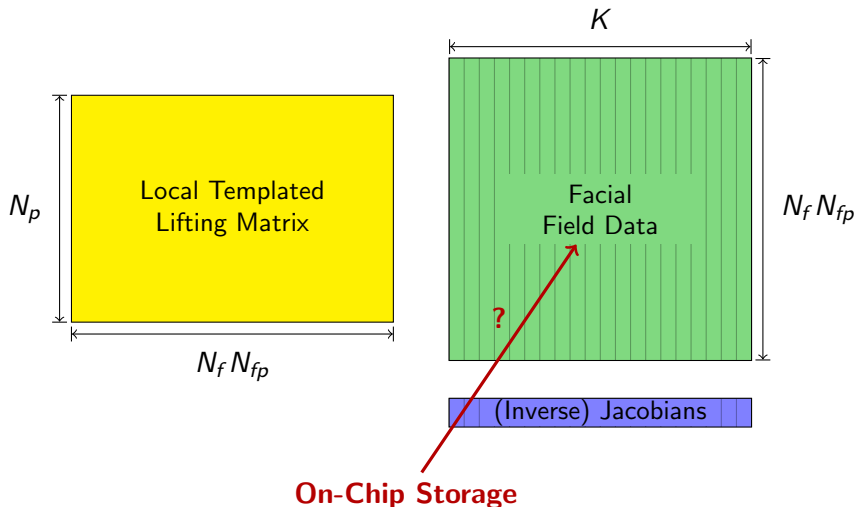


On-Chip Storage

Element-Local Operations: Lifting



Element-Local Operations: Lifting



Best use for on-chip memory?

Basic Problem

On-chip storage is scarce. . .

...and will be for the foreseeable future.

Possible uses:

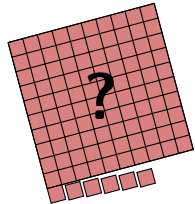
- Matrix/Matrices
- Part of a matrix
- Field Data
- Both

How to decide? Does it matter?



Work Partition for Element-Local Operators

Natural Work Decomposition:
One Element per Block

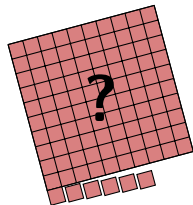


Work Partition for Element-Local Operators

Natural Work Decomposition:

One Element per Block


- ➕ Straightforward to implement
- ➕ No granularity penalty
- ➖ Cannot fill wide SIMD: unused compute power for small to medium elements
- ➖ Data alignment: Padding wastes memory
- ➖ Cannot amortize cost of preparation steps (e.g. fetching)



Loop Slicing for element-local parts of GPU DG

Per Block: K_L element-local mat.mult. + matrix load




Question: How should one assign work to threads? 

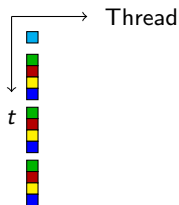
Loop Slicing for element-local parts of GPU DG

Per Block: K_L element-local mat.mult. + matrix load



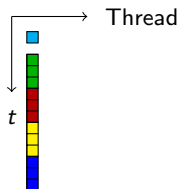
Question: How should one assign work to threads? 

w_s : in sequence



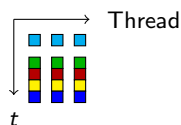
(amortize preparation)

w_i : "inline-parallel"



(exploit register space)

w_p : in parallel



Best Work Partition?

Basic Problem

Additional tier in parallelism offers additional choices. . .
...but very little in the way of guidance.

Possible work partitions:

- One or multiple elements per block?
- One or multiple DOFs per thread?
 - In parallel?
 - In sequence?
 - In-line?

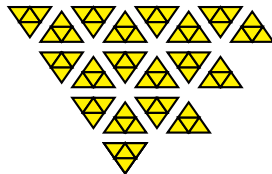


How to decide? Does it matter?

Work Partition for Surface Flux Evaluation

Granularity Tradeoff:

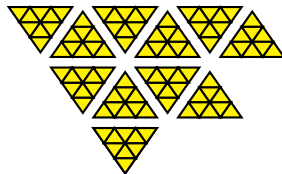
- Large Blocks:
 - ➕ More Data Reuse
 - ➖ Less Parallelism
 - ➖ Less Latency Hiding
- Block Size limited by two factors:
 - Output buffer size
 - Face metadata size
- Optimal Block Size:
not obvious ?



Work Partition for Surface Flux Evaluation

Granularity Tradeoff:

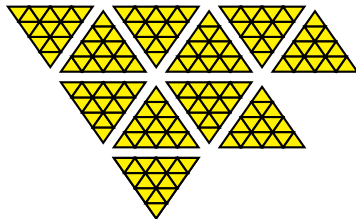
- Large Blocks:
 - ➕ More Data Reuse
 - ➖ Less Parallelism
 - ➖ Less Latency Hiding
- Block Size limited by two factors:
 - Output buffer size
 - Face metadata size
- Optimal Block Size:
not obvious ?



Work Partition for Surface Flux Evaluation

Granularity Tradeoff:

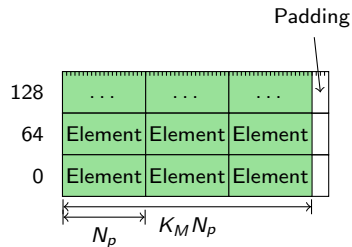
- Large Blocks:
 - ➕ More Data Reuse
 - ➖ Less Parallelism
 - ➖ Less Latency Hiding
- Block Size limited by two factors:
 - Output buffer size
 - Face metadata size
- Optimal Block Size:
not obvious ?



More than one Granularity

Different block sizes introduced so far:

- Differentiation
- Lifting
- Surface Fluxes



More than one Granularity

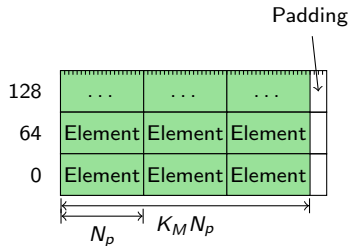
Different block sizes introduced so far:

- Differentiation
- Lifting
- Surface Fluxes

Idea

Introduce another, smaller block size to satisfy SIMD width and alignment constraints. (“*Microblock*”)

- And demand other block sizes be a multiple of this new size



More than one Granularity

Different block sizes introduced so far:

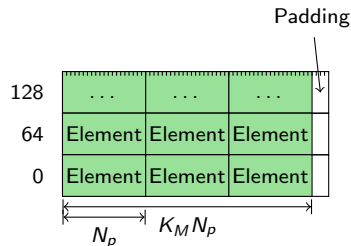
- Differentiation
- Lifting
- Surface Fluxes

Idea

Introduce another, smaller block size to satisfy SIMD width and alignment constraints. (“*Microblock*”)

- And demand other block sizes be a multiple of this new size

How big? Not obvious. ?



DG on GPUs: Implementation Choices

- Many difficult questions
- Insufficient heuristics
- Answers are hardware-specific and have no lasting value



DG on GPUs: Implementation Choices

- Many difficult questions
- Insufficient heuristics
- Answers are hardware-specific and have no lasting value



Proposed Solution: Tune automatically for hardware at computation time, cache tuning results.

- Decrease reliance on knowledge of hardware internals
- Shift emphasis from tuning *results* to tuning *ideas*

Outline

1 PyCUDA

2 Automatic GPU Programming

3 GPU-DG: Challenges and Solutions

- Introduction
- Challenges
- **Benefits of Metaprogramming**
- GPU-DG: Performance and Generality
- Viscous Shock Capture



Metaprogramming for GPU-DG

- Specialize code for user-given problem:
 - Flux Terms



Metaprogramming for GPU-DG

- Specialize code for user-given problem:
 - Flux Terms
- Automated Tuning:
 - Memory layout
 - Loop slicing
 - Gather granularity



Metaprogramming for GPU-DG

- Specialize code for user-given problem:
 - Flux Terms
- Automated Tuning:
 - Memory layout
 - Loop slicing
 - Gather granularity
- Constants instead of variables:
 - Dimensionality
 - Polynomial degree
 - Element properties
 - Matrix sizes



Metaprogramming for GPU-DG

- Specialize code for user-given problem:
 - Flux Terms
- Automated Tuning:
 - Memory layout
 - Loop slicing
 - Gather granularity
- Constants instead of variables:
 - Dimensionality
 - Polynomial degree
 - Element properties
 - Matrix sizes
- Loop Unrolling

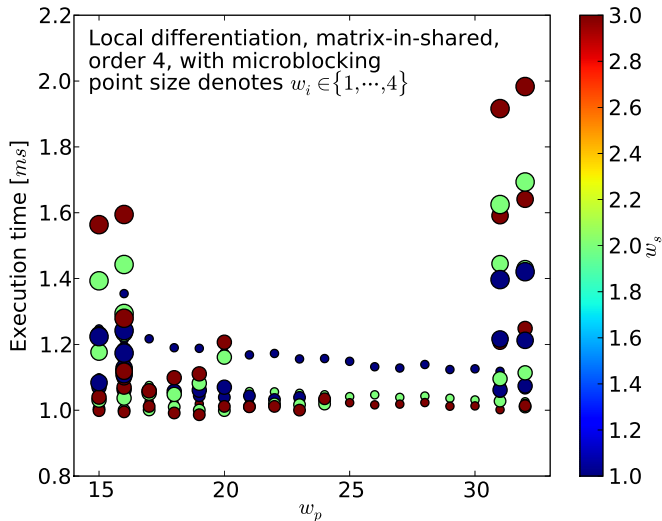


Metaprogramming for GPU-DG

- Specialize code for user-given problem:
 - Flux Terms (*)
- Automated Tuning:
 - Memory layout
 - Loop slicing (*)
 - Gather granularity
- Constants instead of variables:
 - Dimensionality
 - Polynomial degree
 - Element properties
 - Matrix sizes
- Loop Unrolling



Loop Slicing for Differentiation



Metaprogramming DG: Flux Terms

$$0 = \int_{D_k} u_t \varphi + [\nabla \cdot F(u)] \varphi \, dx - \underbrace{\int_{\partial D_k} [\hat{n} \cdot F - (\hat{n} \cdot F)^*] \varphi \, dS_x}_{\text{Flux term}}$$

Metaprogramming DG: Flux Terms

$$0 = \int_{D_k} u_t \varphi + [\nabla \cdot F(u)] \varphi \, dx - \underbrace{\int_{\partial D_k} [\hat{n} \cdot F - (\hat{n} \cdot F)^*] \varphi \, dS_x}_{\text{Flux term}}$$

Flux terms:

- vary by problem
- expression specified by user
- evaluated pointwise

Metaprogramming DG: Flux Terms Example

Example: Fluxes for Maxwell's Equations

$$\hat{n} \cdot (F - F^*)_E := \frac{1}{2} [\hat{n} \times (\llbracket H \rrbracket - \alpha \hat{n} \times \llbracket E \rrbracket)]$$

Metaprogramming DG: Flux Terms Example

Example: Fluxes for Maxwell's Equations

$$\hat{n} \cdot (F - F^*)_E := \frac{1}{2} [\hat{n} \times (\llbracket H \rrbracket - \alpha \hat{n} \times \llbracket E \rrbracket)]$$

User writes: Vectorial statement in math. notation

```
flux = 1/2*cross(normal, h.int - h.ext  
              - alpha*cross(normal, e.int - e.ext))
```

Metaprogramming DG: Flux Terms Example

Example: Fluxes for Maxwell's Equations

$$\hat{n} \cdot (F - F^*)_E := \frac{1}{2} [\hat{n} \times (\llbracket H \rrbracket - \alpha \hat{n} \times \llbracket E \rrbracket)]$$

We generate: Scalar evaluator in C (6×)

```
a_flux += (
    ((( val_a_field5  - val_b_field5 ) * fpair ->normal[2]
      - ( val_a_field4  - val_b_field4 ) * fpair ->normal[0])
    + val_a_field0  - val_b_field0 ) * fpair ->normal[0]
    - ((( val_a_field4  - val_b_field4 ) * fpair ->normal[1]
      - ( val_a_field1  - val_b_field1 ) * fpair ->normal[2])
    + val_a_field3  - val_b_field3 ) * fpair ->normal[1]
    ) * value_type (0.5);
```

Hedge DG Solver



- High-Level Operator Description
 - Maxwell's
 - Euler
 - Poisson
 - Compressible Navier-Stokes, ...
- One Code runs...
 - ...on CPU, CUDA
 - ...on {CPU,CUDA}+MPI
 - ...in 1D, 2D, 3D
 - ...at any order
- Uses CPU, GPU code generation
- Open Source (GPL3)
- Written in Python,

Outline

1 PyCUDA

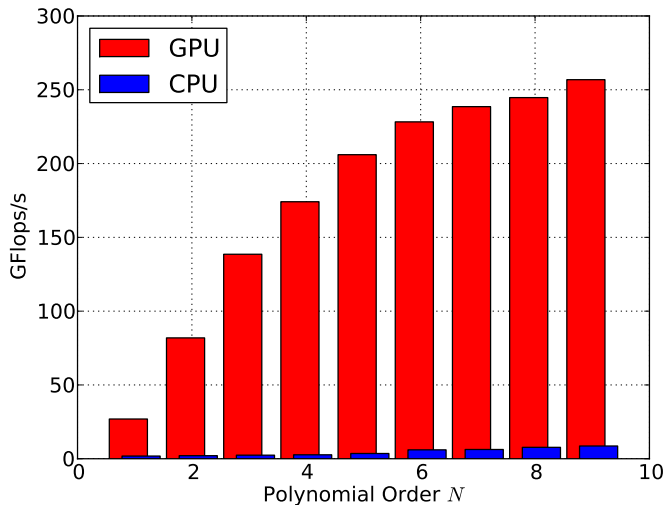
2 Automatic GPU Programming

3 GPU-DG: Challenges and Solutions

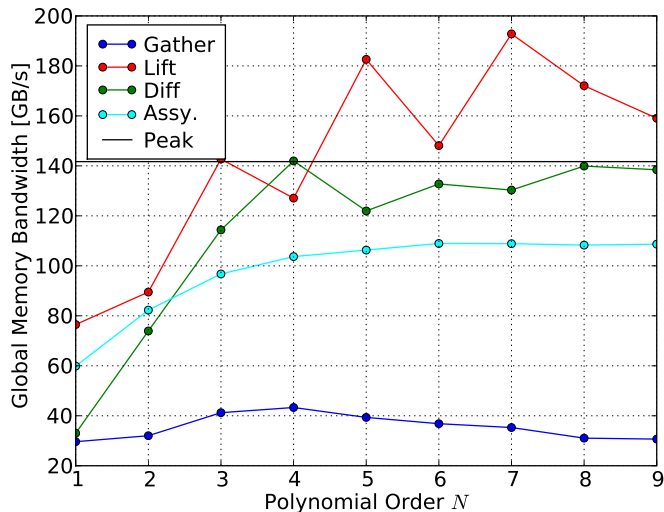
- Introduction
- Challenges
- Benefits of Metaprogramming
- GPU-DG: Performance and Generality
- Viscous Shock Capture



Nvidia GTX280 vs. single core of Intel Core 2 Duo E8400

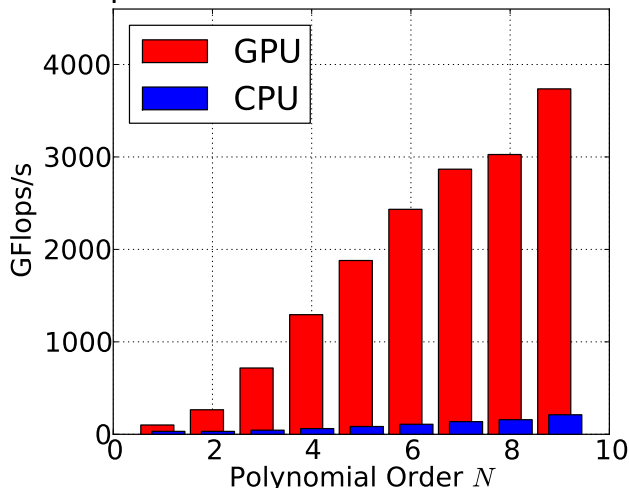


Memory Bandwidth on a GTX 280



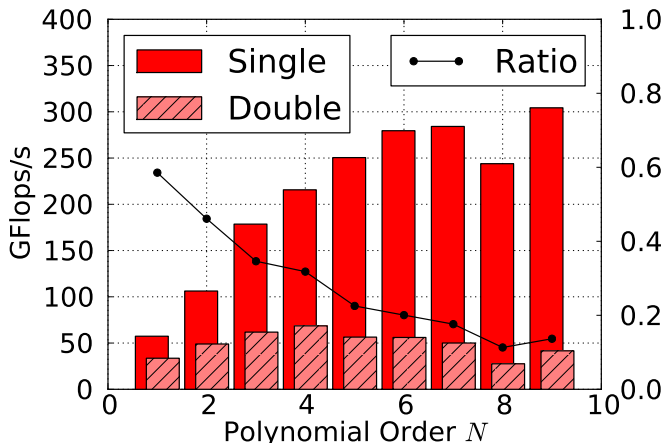
Multiple GPUs via MPI: 16 GPUs vs. 64 CPUs

Flop Rates: 16 GPUs vs 64 CPU cores



GPU-DG in Double Precision

GPU-DG: Double vs. Single Precision



Outline

1 PyCUDA

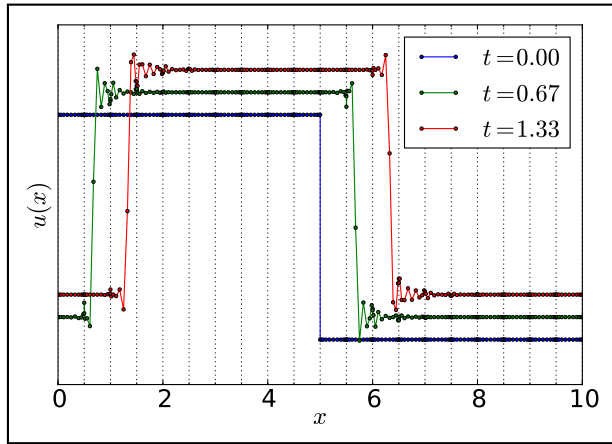
2 Automatic GPU Programming

3 GPU-DG: Challenges and Solutions

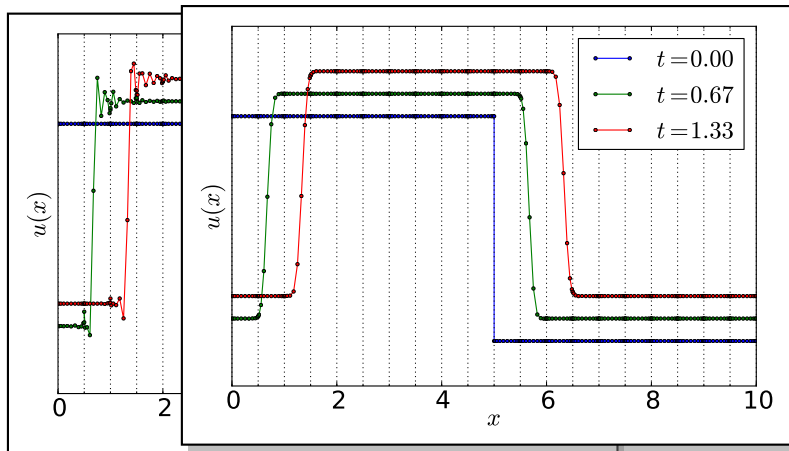
- Introduction
- Challenges
- Benefits of Metaprogramming
- GPU-DG: Performance and Generality
- Viscous Shock Capture



Nonlinear conservation laws \rightarrow shocks?



Nonlinear conservation laws \rightarrow shocks?



Nonlinear conservation laws \rightarrow shocks?

1D advection:

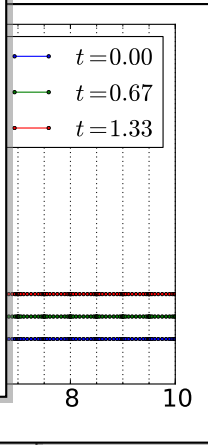
$$\partial_t u + \partial_x u = 0$$

1D advection with viscosity:

$$\partial_t u + v \cdot \nabla_x u = \nabla_x \cdot (\nu \nabla_x u).$$

Important: Conservation form.

Upwind fluxes for advection, IPDG for second-order



Nonlinear conservation laws \rightarrow shocks?

1D advection:

$$\partial_t u + \partial_x u = 0$$

1D advection

$$\partial_t u +$$

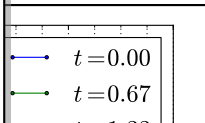
Important: Co

Upwind fluxes
order

0

2

- Detector $\rightarrow \nu$?
GPU-suitability? Data locality?
Properties? [Build on work by
Persson/Peraire '06]
- Time integration
Implicit/explicit? Adaptivity? RKC for
bigger Δt with viscosity?
- Accuracy?
Near shocks? Away from them?



Results: Euler's Equations of Gas Dynamics

Euler's equations with viscosity:

$$\begin{aligned}\partial_t \rho + \nabla_x \cdot (\rho \mathbf{u}) &= \nabla_x \cdot (\nu \nabla_x \rho), \\ \partial_t (\rho \mathbf{u}) + \nabla_x \cdot (\mathbf{u} \otimes (\rho \mathbf{u})) + \nabla_x p &= \nabla_x \cdot (\nu \nabla_x (\rho \mathbf{u})), \\ \partial_t E + \nabla_x \cdot (\mathbf{u}(E + p)) &= \nabla_x \cdot (\nu \nabla_x E).\end{aligned}$$

Again: Single ν , sensed on ρ . \rightarrow Undue pollution of the other field?

[Persson/Peraire '06] suggest Navier-Stokes-like viscosity. No good: can't control jumps in ρ .

Rusanov fluxes for Euler, IPDG for viscosity.

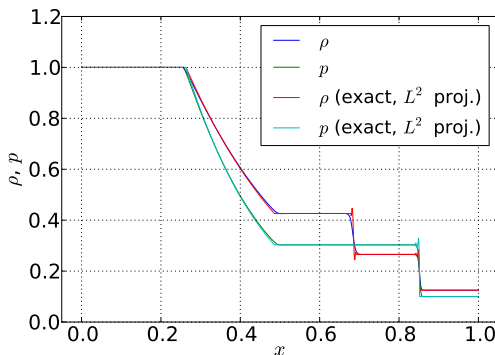
Results: Euler's Equations of Gas Dynamics

Euler's equations with viscosity:

$$\partial_t \rho + \nabla_x \cdot (\rho \mathbf{u}) = \nabla_x \cdot (\nu \nabla_x \rho),$$

$\partial_t(\rho \mathbf{u}) + \nabla_x \cdot (\rho \mathbf{u} \mathbf{u}) = -\nabla_x p + \nabla_x \cdot (\nu \nabla_x \mathbf{u})$

Sod's Problem with $N=5$ and $K=80$



Against
the ot
[Perss
cosity.
Rusan

of
is-

Results: Euler's Equations of Gas Dynamics

Euler's equations with viscosity:

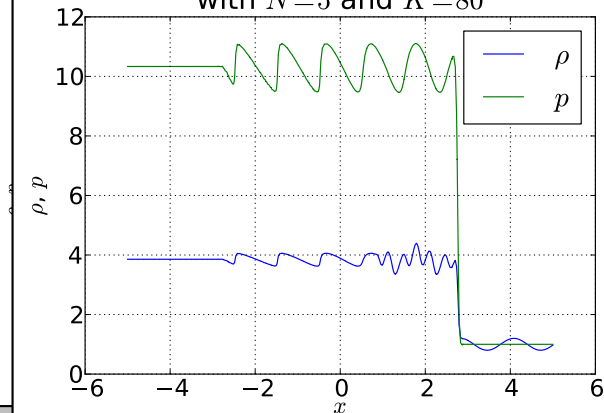
$\partial_t(\rho u)$

Again
the ot

[Perss
cosity.

Rusan

Shock-Wave Interaction Problem
with $N=5$ and $K=80$



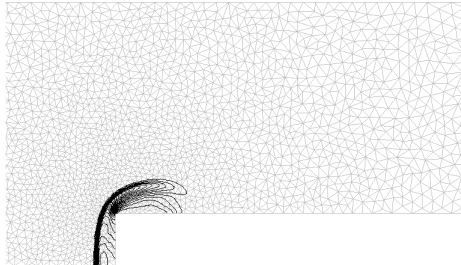
Results: Euler's Equations of Gas Dynamics

Euler's equations with viscosity:

Shock-Wave Interaction Problem

$\partial_t(\rho u)$

Again
the ot
[Perss
cosity.
Rusan



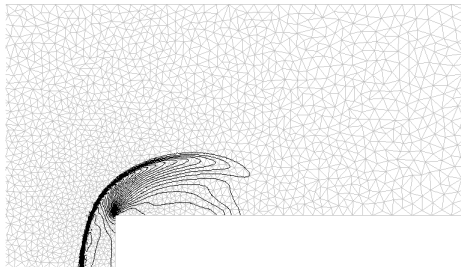
Results: Euler's Equations of Gas Dynamics

Euler's equations with viscosity:

Shock-Wave Interaction Problem

$\partial_t(\rho u)$

Again
the ot
[Perss
cosity.
Rusan



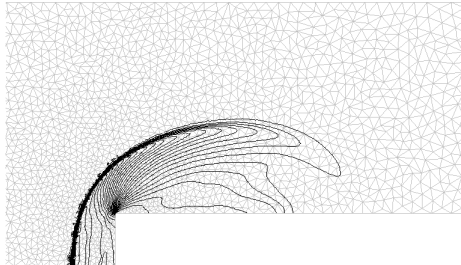
Results: Euler's Equations of Gas Dynamics

Euler's equations with viscosity:

Shock-Wave Interaction Problem

$\partial_t(\rho u)$

Again
the ot
[Perss
cosity.
Rusan



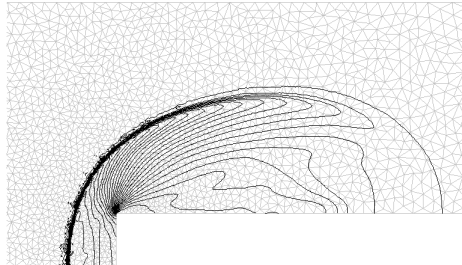
Results: Euler's Equations of Gas Dynamics

Euler's equations with viscosity:

Shock-Wave Interaction Problem

$\partial_t(\rho u)$

Again
the ot
[Perss
cosity.
Rusan



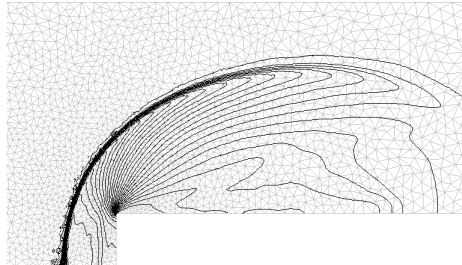
Results: Euler's Equations of Gas Dynamics

Euler's equations with viscosity:

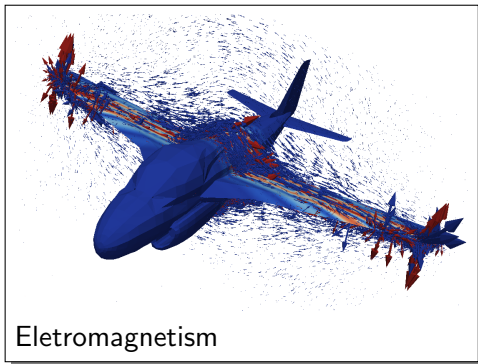
Shock-Wave Interaction Problem

$\partial_t(\rho u)$

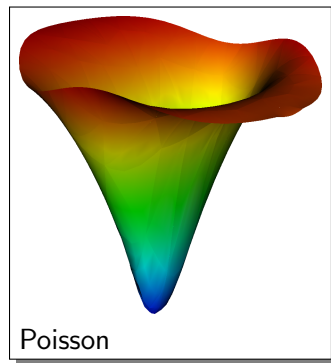
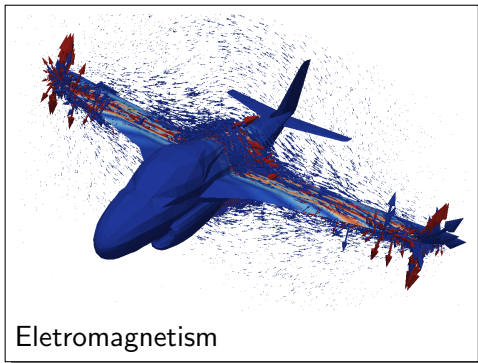
Again
the ot
[Perss
cosity.
Rusan



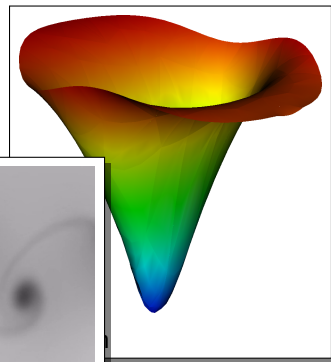
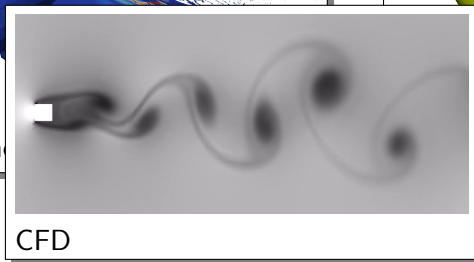
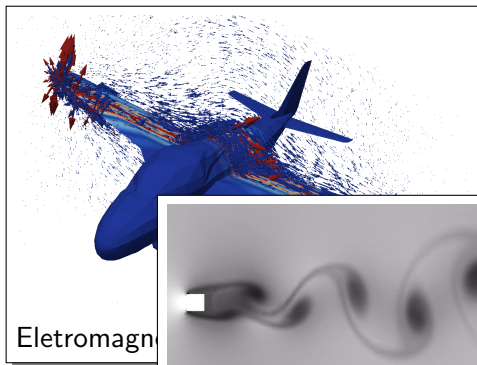
GPU DG Showcase



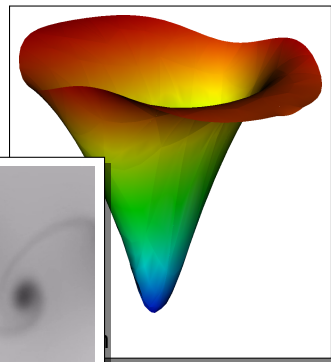
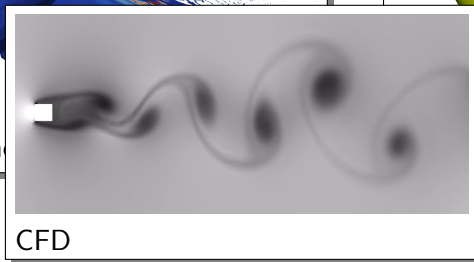
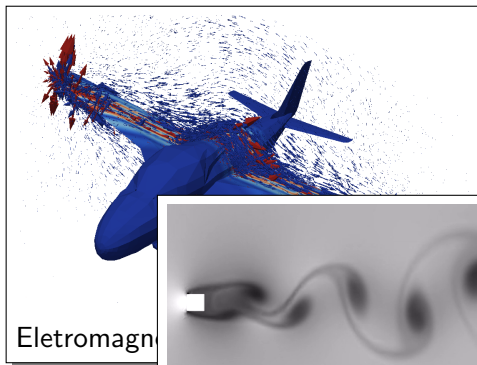
GPU DG Showcase



GPU DG Showcase



GPU DG Showcase



Where to from here?

PyCUDA, PyOpenCL, hedge

→ <http://www.cims.nyu.edu/~kloeckner/>

GPU-DG Article

AK, T. Warburton, J. Bridge, J.S. Hesthaven, “*Nodal Discontinuous Galerkin Methods on Graphics Processors*”, J. Comp. Phys., 228 (21), 7863–7882.

GPU RTCG

AK, N. Pinto et al. *PyCUDA: GPU Run-Time Code Generation for High-Performance Computing*, submitted.

Conclusions

- GPUs and scripting work surprisingly well together
 - Enable Run-Time Code Generation
- GPU-DG is significantly faster than CPU-DG
 - Method well-suited a priori
 - Numerous tricks enable good performance
- Further work in GPU-DG:
 - Curvilinear Elements (T. Warburton)
 - Local Time Stepping
 - Shock Capturing for Nonlinear Equations



Questions?



?

Thank you for your attention!

<http://www.cims.nyu.edu/~kloeckner/>

► image credits

Image Credits

- Exclamation mark: sxc.hu/cobrasoft
- Adding Machine: flickr.com/thomashawk 
- Floppy disk: flickr.com/ethanhein 
- Carrot: OpenClipart.org
- Dart board: sxc.hu/195617
- Question Mark: sxc.hu/svilen001
- Question Mark: sxc.hu/svilen001
- ?/! Marks: sxc.hu/svilen001