

SPARSE MATRIX SOLVERS ON THE GPU: CONJUGATE GRADIENTS AND MULTIGRID

Jeffrey Bolz

Ian Farmer

Eitan Grinspun

Peter Schröder

Caltech



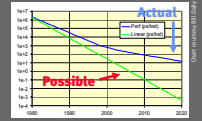
MULTI-RES MODELING GROUP

1

WHY USE THE GPU?

Semiconductor trends

- cost
- wires vs. compute
- Stanford streaming supercomputer



Parallelism

- many functional units
- graphics is prime example



MULTI-RES MODELING GROUP

2

NEW HARDWARE FEATURES

Latest generation graphics HW

- floating point throughout
- programmability
- high resource limits
 - dependent texturing
 - many registers
 - many instructions



MULTI-RES MODELING GROUP

3

HOW TO EXPLOIT?

Harvesting this power

- what application suitable?
- what abstractions useful?

History

- massively parallel SIMD machines
- media processing



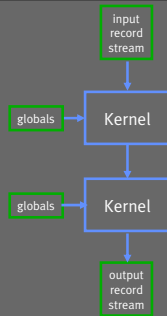
MULTI-RES MODELING GROUP

4

STREAMING MODEL

What is the right abstraction?

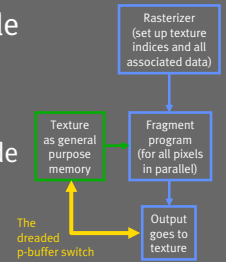
- Purcell, et al. 2002
- data structures \rightsquigarrow streams \rightsquigarrow textures
- algorithms \rightsquigarrow kernels \rightsquigarrow fragment programs



MAPPING THE PROGRAM

How does a program execute?

- “render” a rectangle
- memory as texture
- fragment program
 - many pixels provide parallelism
 - write to texture
 - close the loop



OUR PROGRAM

Kernels for scientific computing

- sparse matrix solvers
 - not high arithmetic intensity...
 - but... we need them everywhere
- unstructured: conjugate gradients
- structured: multigrid

Lessons to learn here...



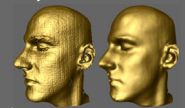
SPARSE MATRICES

Ubiquitous in numerical computing

- discretization of PDEs: animation
 - finite elements, difference, volumes
- optimization, editing, etc., etc.

Example here:

- processing of surfaces
- 2-manifold triangle meshes



GEOMETRIC FLOW

Canonical non-linear problem

- mean curvature flow

$$\partial_t x_i(t) = -\lambda_i(t) H_i(t) \tilde{n}_i(t)$$

Velocity opposite mean curvature normal

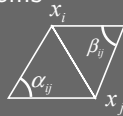
- implicit time discretization

$$Ax^{(i+1)} = x^{(i)}$$

- solve sequence of SPD systems

$$a_{ij} = -\lambda \Delta t (\cot(\alpha_{ij}) + \cot(\beta_{ij}))$$

$$a_{ii} = 4A_i - \sum_{j \in N(i)} a_{ij}$$



NON-LINEAR PROBLEMS

Basic structure

- solver for SPD systems
 - conjugate gradients
 - other variants if not SPD
- recompute matrix entries on GPU
 - minimize CPU to GPU traffic
- control structure on CPU



CONJUGATE GRADIENTS

High level code

- inner loop
- matrix-vector multiply
- sum-reduction
- scalar-vector MAD

while($\delta_{new} < \epsilon^2 \delta_0$)

$$q = Ad$$

$$\alpha = \delta_{new} / d^T q$$

$$x = x + \alpha d$$

$$r = r - \alpha q$$

$$\delta_{old} = \delta_{new}$$

$$\delta_{new} = r^T r$$

$$\beta = \delta_{new} / \delta_{old}$$

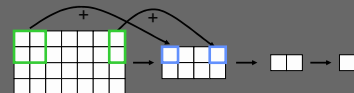
$$d = r + \beta d$$



VECTOR INNER PRODUCTS

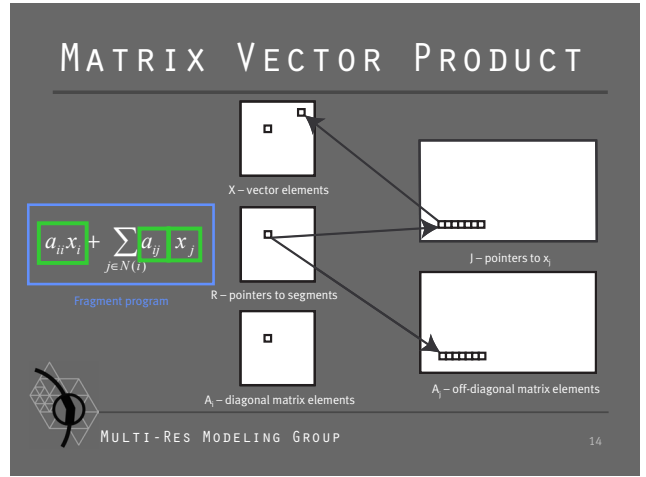
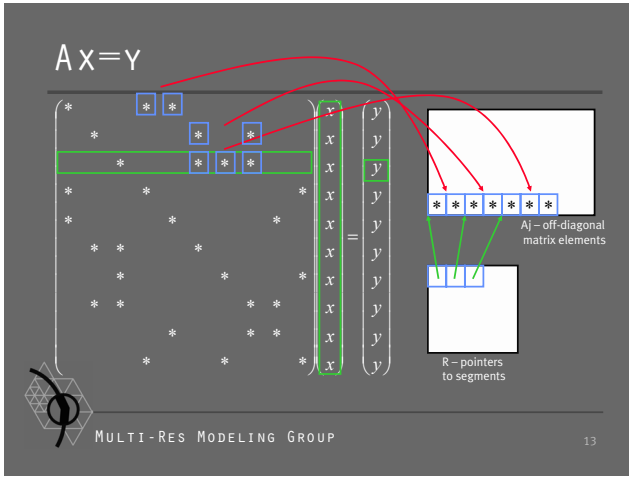
Decompose

- fragment-wise multiply
- followed by sum-reduction



- odd dimensions can be handled
- useful for other operations as well





APPLY TO ALL PIXELS

Two extremes

- one row at a time
 - setup overhead
- all rows at once
 - limited by worst row
- middle ground
 - organize "batches" of work
 - what size? how to organize?

MULTI-RES MODELING GROUP 15

WHAT SIZE BATCHES?

We choose fixed size rectangles

- fragment pipe is quantized

- simple experiments reveal best size
 - performance model
 - area, diagonal, wasted frags

MULTI-RES MODELING GROUP 16

PACKING (GREEDY)

15 13 13 12 12 11 10 9 9 9 8 8 8 8 8 7 7 7 7 7 7 7 7 6 5 4 ...

non-zero entries per row

15	13	13	9	9	8	7	7
12	12	11	8	8	7	7	7
10	9	9	8	7	7	7	6

still some zero padding required

each batch bound to an appropriate fragment program

All this setup done once only at the beginning of time. Depends only on mesh connectivity



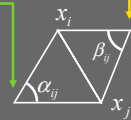
RECOMPUTING MATRIX

Matrix entries depend on surface

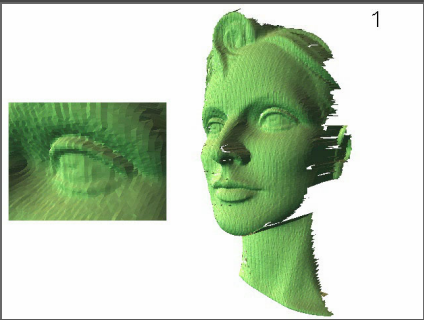
- must "render" into matrix
- two additional indirection textures
 - previous and next

$$a_{ij} = -\lambda \Delta t (\cot(\alpha_{ij}) + \cot(\beta_{ij}))$$

$$a_{ii} = 4A_i - \sum_{j \in N(i)} a_{ij}$$



SURFACE SMOOTHING



RESULTS (NV30@500MHZ)

37k elements

- matrix multiply
 - 33 instructions
 - 1/100th of a second
- reduction
 - 7 instructions/fragment/pass
 - 1/1900th of a second
- CG solve in 1/20th of a second



RESULTS

How efficient?

- lots of indirection
 - 33 instructions for 13 (average) flops
- bandwidth limited
 - fat texture fetches are slow
- not the ideal example for GPU...
 - small op/fetch ratio



SO FAR...

Unstructured matrices

- irregular triangle meshes
 - also tet meshes
- main issue is layout of matrix

Structured matrices

- much nicer layout
- example: fluids, image processing



REGULAR GRIDS

PDEs again

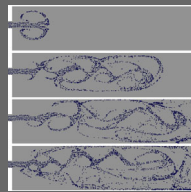
- this time variables on “pixel grid”
 - e.g.: Navier-Stokes

$$\nabla \cdot \mathbf{u} = 0$$

$$\rho \frac{\partial \mathbf{u}}{\partial t} = -(\mathbf{u} \cdot \nabla) \mathbf{u} + \nu \nabla^2 \mathbf{u} + \rho \mathbf{b}$$

after discretization:
solve Poisson eq.
at each time step

$$\nabla^2 p = \nabla \cdot \mathbf{u}$$



POISSON EQUATION

Appears all over the place

- easy to discretize on regular grid
- matrix multiply is stencil application
- FD Laplace stencil:

0	1	0
1	-4	1
0	1	0

(i,j)

$$\nabla^2 X_{i,j} = X_{i-1,j} + X_{i+1,j} + X_{i,j-1} + X_{i,j+1} - 4X_{i,j}$$



SOLVER

Use iterative matrix solver

- just need application of stencil
- easy: just like filtering
- incorporate geometry (Jacobian)
- variable coefficients

But..., very ill-conditioned

- use multigrid to counteract



MULTIGRID (BASIC)

Principles

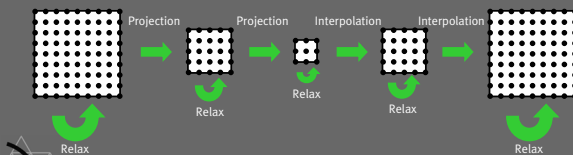
- lower frequency error resolved on coarser grid
- implementation needs:
 - interpolation (coarse → fine)
 - projection (fine → coarse)
 - relaxation
 - Jacobi iterations



V-CYCLE

Fine to coarse to fine cycle

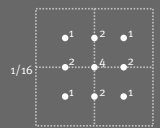
- project residual, relax, interpolate correction vector



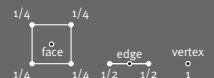
COMPUTATIONS

Lots of stencil applications

- matrix multiply: 3x3 stencil
- projection: 3x3 stencil
- interpolation: 2x2(!)
- floor op in indexing...



$$\mathbf{v}_h[\hat{i}] = 1/4 \sum_{d \in \{0,1\}^2} \mathbf{v}_{2h}[\lfloor (i+d)/2 \rfloor]$$



DETAILS

Boundaries

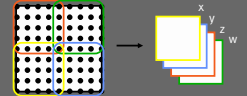
- Dirichlet boundaries
 - boundary variables are fixed
- Neumann boundaries
 - out of bounds access clamped to 0
- gradient and divergence operators
 - one sided differences at boundary



TEXTURE LAYOUT

Storage for matrices and DOFs

- variables in one texture
- matrices in 9(=3x3) textures
- all textures packed
 - exploit 4 channels
 - domain decomp.
 - padded boundary



COARSER MATRICES

Operator at coarser level

- needed for relaxation at all levels
- triple matrix product

$$A^c = P := 1/4 S^T \quad A^f \quad S$$

Effectively:
Stencil composition



STENCIL COMPOSITION

Triple matrix product...

- work out terms and map to stencils
- exploit local support of stencils
- straightforward but t-e-d-i-o-u-s

$$A_{2,h}^d[i] = 1/4 \sum_{e,g \in \{-1,0,1\}^2} S^e S^{e+g-2d} A_h^g[2i+e]$$

$$= 1/4 \sum_{e,g \in \{-1,0,1\}^2} S^e S^g [e+g-2d] A_h^g[2i+e]$$

- store S in texture S' with a 0 boundary



MORE DETAILS

What is variable?

- only matrix entries (stencils) vary
- all other operators hard wired
 - still general solution!
- debugging
 - oh, joy...
- obstacles resolved at coarsest level



FLOW EXAMPLE

Putting it all together

- here: fixed velocity in and out
- tracer particles for visualization
 - simple advection in evolving flow



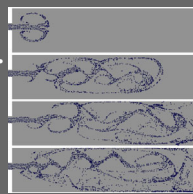
RESULTS (NV30@500MHZ)

257x257 grid

- matrix multiply - 27 instructions
 - 1/800th of a second
- interpolation 10 instr.
- projection 19 instr.

Overall performance

- 257x257 at 90 fps!



PROBLEMS

Bleeding edge...

- PBuffer overhead is killing us
- managing layout in a buffer by hand... OUCH
- scalar versus vector problems
- give us rectangular texture border



ENHANCEMENTS

Small/large changes?

- global registers for reductions
- texture fetch with offset
- scatter (displacement mapping?!)
 - rasterization order undefined
 - scatter w/ undefined order still useful
- scientific computing compiler



CONCLUSION

Where are we now?

- performance not up to expectations
- still a good streaming processor
- most kernels run at 1-2GFlops/s
 - SSE on P4 still very competitive
- should beat CPUs in about a year
- better languages! Brook? C*?

