



Using GPUs in Pre-stack AVO Inversion for Performance Improvement

Arthur Lee

CREWES Weekly Friday Morning Meeting Feb 8, 2019





- 1. Objective and Motivation
- 2. Pre-stack AVO Inversion (Hampson et al., 2005) using Conjugate Gradient
- 3. Software Engineering Solutions to a Large Data Problem
- 4. Using GPUs in Conjugate Gradient
- 5. What is GPU Programming
- 6. Results
- 7. Conclusion
- 8. Discussions

Dan Hampson, CGG

CREWES Kevin Hall, University of Calgary

Jorge Monsegny Parra, University of Calgary Daniel Trad, University of Calgary Larry Lines, University of Calgary **Reflection Seismic Method**

Acquisition \rightarrow Processing \rightarrow Interpretation and Reservoir Characterization

Acquisition \rightarrow Processing \rightarrow Interpretation:

Acquisition \rightarrow Processing \rightarrow Pre-stack AVO Inversion \rightarrow Elastic Properties (Zp, Zs, and Density)

Purposes:	(1) Reconnaissance Analysis and
	(2) Detailed Rock/Fluid Property Estimation
Requirements:	Interactive and Fast, and
	Faster

Would GPU help?

Pre-stack AVO Inversion

Seismic Inversion is "to extract geological model information from data". (Lines, Chapter 16)



<u>INPUT:</u>

Pre-stack migrated, or NMO-corrected, and AVO friendly processed Angle gathers (AVO -> AVA)

<u>OUTPUT</u>:

Simultaneously solves for elastic properties: P-Impedance, S-Impedance and Density.

Available tools (linear/non-linear):

- 1. Artificial Intelligent
- 2. Statistical
- 3. Deterministic

Pre-stack AVO Inversion – Assumptions

The Forward Model Development ...

Equation:

 $T = \mathbf{W} R$ -- Convolutional Model

 $T = \mathbf{W} \ (c_1 R_P + c_2 R_S + c_3 R_D)$ -- Linearization

$$T = \mathbf{W} \left(\frac{1}{2} c_1 \mathbf{D} L_P + \frac{1}{2} c_2 \mathbf{D} L_S + c_3 \mathbf{D} L_D \right)$$
 -- Reflectivity as small changes in Elastic Parameters

Assumptions:

Additional assumption by Hampson et al., 2005

Model parameters Z_S and ρ are related to Z_P .



Assume: $L_S = k L_P + k_c + \Delta L_S$ $L_D = m L_P + m_c + \Delta L_D$

Then: $\mathbf{D}L_S = \mathbf{k} \mathbf{D}L_P + \mathbf{D}\Delta L_S$ $\mathbf{D}L_D = \mathbf{m} \mathbf{D}L_P + \mathbf{D}\Delta L_D$

The Forward Model

Equation:

 $T = \mathbf{W} R$

- $T = \mathbf{W} \ \left(c_1 R_P + c_2 R_S + c_3 R_D \right)$
- $T = \mathbf{W} \left(\frac{1}{2} c_1 \mathbf{D} L_P + \frac{1}{2} c_2 \mathbf{D} L_S + c_3 \mathbf{D} L_D \right)$
- $T = \mathbf{W} \left(\widetilde{c_1} \mathbf{D} L_P + \widetilde{c_2} \mathbf{D} \Delta L_S + \widetilde{c_3} \mathbf{D} \Delta L_D \right)$

Assumptions:

-- Convolutional Model

-- Linearization

-- Reflectivity as small changes in Elastic Parameters

-- Elastic Parameters are related.

-- Invert for elastic parameter changes from a trend.

The equation in Matrix Notation for one angle gather: (with *M* angle traces per gather, and *N* samples per trace)

$$\begin{bmatrix} T(\theta_{1}) \\ \vdots \\ T(\theta_{2}) \\ \vdots \\ T(\theta_{M}) \end{bmatrix} = \begin{bmatrix} W_{1}(\theta_{1}) \dots & W_{2}(\theta_{1}) \dots & W_{3}(\theta_{1}) \dots \\ \vdots & \vdots & \vdots \\ W_{1}(\theta_{2}) \dots & W_{2}(\theta_{2}) \dots & W_{3}(\theta_{2}) \dots \\ \vdots & \vdots & \ddots & \ddots \\ W_{1}(\theta_{M}) \dots & W_{2}(\theta_{M}) \dots & W_{3}(\theta_{M}) \dots \end{bmatrix} \begin{bmatrix} L_{P} \\ \vdots \\ \Delta L_{S} \\ \vdots \\ \Delta L_{D} \end{bmatrix} \longrightarrow T = \mathbf{A}m$$

$$\begin{bmatrix} Wavelets and linearized AVO \\ coefficients combined at \\ different angles. \end{bmatrix}$$

The basic equation.

Basic Equation:

$$T = Am$$
Normal Equation:

$$A^{T}T = A^{T}Am$$
RHS:

$$A^{T}Am = \begin{bmatrix} D^{T}G_{11}D & D^{T}G_{12}D & D^{T}G_{13}D \\ \vdots & \vdots & \vdots \\ D^{T}G_{21}D & D^{T}G_{22}D & D^{T}G_{23}D \\ \vdots & \vdots & \vdots \\ D^{T}G_{31}D & D^{T}G_{32}D & D^{T}G_{33}D \end{bmatrix} \begin{bmatrix} L_{P} \\ \vdots \\ \Delta L_{S} \\ \vdots \\ \Delta L_{D} \\ \vdots$$

Basic equation: Normal equation:

 $T = \mathbf{A}m$ $\mathbf{A}^{\mathrm{T}}T = \mathbf{A}^{\mathrm{T}}\mathbf{A}m$

Conjugate Gradient Algorithm by Hestenes and Stiefel (1952) as mentioned by Scales (1989):



Hampson's conjugate gradient implementation concentrates on the combined **A^TA** operator.

$$\mathbf{A}^{\mathrm{T}}\mathbf{A}m = \begin{bmatrix} \mathbf{D}^{\mathrm{T}}\mathbf{G}_{11}\mathbf{D} & \mathbf{D}^{\mathrm{T}}\mathbf{G}_{12}\mathbf{D} & \mathbf{D}^{\mathrm{T}}\mathbf{G}_{13}\mathbf{D} \\ \mathbf{D}^{\mathrm{T}}\mathbf{G}_{21}\mathbf{D} & \mathbf{D}^{\mathrm{T}}\mathbf{G}_{22}\mathbf{D} & \mathbf{D}^{\mathrm{T}}\mathbf{G}_{23}\mathbf{D} \\ \mathbf{D}^{\mathrm{T}}\mathbf{G}_{31}\mathbf{D} & \mathbf{D}^{\mathrm{T}}\mathbf{G}_{32}\mathbf{D} & \mathbf{D}^{\mathrm{T}}\mathbf{G}_{33}\mathbf{D} \end{bmatrix} \begin{bmatrix} L_{P} \\ \Delta L_{S} \\ \Delta L_{D} \end{bmatrix}$$

- It employs 9 convolutions which can be running in *parallel*.
- The implementation should fit well into a typical multi-threading optimization style.

Software Engineering Solutions to a Large Data Problem

Typical software engineering solutions (using existing science):

<u>Coarse-grained Parallelization</u> Multiple Nodes (CPU's / Machines / Clusters) Network connection. Non-shared memory.	E.g. MPI					
Medium-grained Parallelization						
Multiple threads in one CPU						
Threads ~= cores. Shared memory.	E.g. pthread, OpenMP or MPI					
<u>Fine-grained Parallelization</u> Multiple threads (in one CPU) Multiple GPU cards connected by PCI-e	E.g. pthread, OpenMP	heterogenous				
<u>Ultra-fine-grained Parallelization</u> Multiple threads in a single GPU card Array or Vector Processors	E.g. "OpenACC", " CUDA ", SIMD (Not seen anymore in geophysical industry.)					

Using GPUs in *the* Pre-stack AVO Inversion



Using GPUs in AVO Inversion







As described in www.nvidia.com,

- parallel computing architecture
- many computing cores \rightarrow mathematical calculations
- supports OpenCL and DirectX
- supports C and Fortran.

In other words, it is a C/C++ pre-processor and compiler which comes with some libraries.

A Hard way to explain the Software operations

Block Diagram of the GP104 chip of GeForce GTX 1090, a 4 cluster card.



GP104 Multiprocessor Diagram		SM		_		Instructi	ion Cache I	_	_	_	_	-
			Inst	uction B	uffer				nstructio	on Buffe	r	
Available Cores			Wa	rp Schedu	ller				Warp Sc	heduler		
Available Colles.			ispatch Unit		Dispatch	Unit	Di	spatch Uni	it	D	ispatch Unit	
In one Stream Multiprocessor (SM or MP):			Register F	ile (16,38	4 x 32-bit))		Regist	er File (1	6,384 x	32-bit)	
		Core	Core C	ore Co	re LDIST	SFU	Core	Core	Core	Core	LD/ST	SFU
• 4 warps		Core	Core C	ore Co	re LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU
• Each warp has 4 x 8 = 32 Cores	_	Core	Core C	ore Co	re LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU
		Core	Core C	ore Co	re LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU
in one Cluster :		Core	Core C	ore Co	re LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU
 5 SM x 4 Warps/SM x 32 Cores/Warps 		Core	Core C	ore Co	re LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU
-640 Cores		Core	Core C	ore Co	re LDIST	SFU	Core	Core	Core	Core	LD/ST	SFU
= 040 Cores		Core	Core C	ore Co	re LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU
In one Device :						Texture /	L Cache					
 4 Clusters x 640 Cores/Cluster 			Tex		Тех			Төх			Төх	
$2 \sum (0, C_{a}) = (a_{a} \sum b_{a})$		Instruction Buffer Instruction Buffer Warp Scheduler Warp Scheduler			r							
 = 2560 Cores (or Threads) 		Dispatch Unit Dispatch Unit Dispatch Unit Dispatch Unit		ispatch Unit								
			Register	gister File (16,384 x 32-bit) Register File (16,384 x 32-bit)								
Available Memory:		Core	Core C	ore Co	re LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU
		Core	Core C	ore Co	re LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU
register file capacity		Core	Core C	ore Co	re LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU
$(16.384x32bit) \times 4 = (16Kx4)x4 = 256KB$		Core	Core	ore Co	re LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU
		Core	Core	ore Co	re LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU
		Core	Core	ore Co	re LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU
• 48 KB L1 cache		Core	Core C	ore Co	re LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU
		Core	Core 0	ore Co	re LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU
						Texture	/ L1 Cache					
 96 KB shared memory 			Tex		Tex	t.		Tex			Tex	
,						96KB Sha	red Memor	у				

CUDA Basics – Write your own kernel function

CPU: C/C++ simple example:

```
for (int i =0; i < N; i++) { function(i, a, b, c ...); }
```

GPU: corresponding CUDA kernel function:

```
function<<< Block Dimension, Grid Dimension >>>(a, b, c)
{
    //Code for one thread, asynchronous:
    int i = (blockIdx.x * blockDim.x) + threadIdx.x;
    if (i < N) { .... function(i, a,b,c); }
}</pre>
```

← Row major linear arrangement of blocks of threads.

CPU: C/C++ simple example:

```
fftwf_plan d_plan_forward =
fftwf_plan_dft_r2c_1d (myNFFT, d_real, d_complex, FFTW_ESTIMATE);
fftwf_execute (d_plan_forward);
```

Corresponding CUDA GPU kernel launches:

```
cufftHandle d_plan_forward = 0;
cufftPlan1d (&d_plan_forward, myNFFT, CUFFT_R2C, N);
cufftExecR2C ( d_plan_forward, d_real, d_complex );
```

GPU Programming – Only two main ideas are involved

1) Putting serial workloads to the CPU (host) and parallel workloads to the GPU (device).

Conjugate gradient iterative solver:

Loop (a serial process) in CPU.

FFT, additions and subtractions (parallel data) in GPU.

```
For all iterations {

If not the first iteration {

\beta = r^T r / r^T r_{old}

p = r + \beta p

}

\alpha = r^T r / (p^T (A^T A)p)

r = r - \alpha (A^T A)p

m = m + \alpha p

}
```

2) Coalescing data access to avoid memory bandwidth limitations. (To be discussed in the future.)



Current Status and Expectation:

- Similar performance after changing the conjugate gradient from CPU • to GPU. The CPU thread can be freed up for other activities.
- Amazon's Best CAD\$. •
- In our tests, the CPU is running in a single core. Two times will roughly balance the cost only. We hope that the GPU can run ~10 times faster.

Future Work:

- Optimize the conjugate gradient implementation.
- Invert multiple CDPs in parallel.



C\$1817	C\$581
Intel	Quadro
Xeon	P2000
W3680	1024 cores
6 cores	
	C\$388
	GeForce
	GTX 1060
	1280 cores



Converted from CPU to GPU a linearized pre-stack AVO least-squares inversion.

Pros:

Highly parallel and scalable. Suitable for conjugate gradient methods. Work well for a large number of data samples and number of iterations.

Cons:

Steep learning curve.

Sensitive to hardware setup.

Code changes may easily run --- slower!

Extra attention is needed for performance profiling and tuning.

Three steps to deliver:

- 1. Profile on CPU, identify and evaluate the performance bottleneck carefully.
- 2. Code and debug in CUDA environment. Use existing libraries.
- 3. Analyze and improve core occupancy, memory bandwidth and other GPU resource usages with tools provided.



The following may interest you:

- 1. Linear Algebra and Sparse Matrix handling. (cuBLAS and cuSPARSE)
- 2. Linear Solvers like Cholesky, SVD. Spare and dense matrix. (cuSOLVER)
- 3. Artificial Intelligent (cuDNN)
- 4. Shared memory (low latency memory) →
 Laplace operator stencils in Finite Difference methods.
- 5. Texture memory \rightarrow trace data interpolation ?



Hardware resources at University of Calgary: barracuda

Two CUDA enabled graphic cards.

Tesla K10 GK104GL dual Kepler (v3) processors.



Software resources: CUDA API,

e lesources. CODA API,	CUDA API References
	CUDA Runtime API
CUDA TOOIKIT V10.0.130	CUDA Driver API
	CUDA Math API
	cuBLAS
	NVBLAS
	nvJPEG
	cuFFT
	nvGRAPH
	cuRAND
	cuSPARSE
	NPP
	NVRTC (Runtime Compilation)
	Thrust
	cuSOLVER

and YOU to build up the CREWES library!