#### **Systematic Fusion of CUDA Kernels for Iterative Sparse Linear System Solvers**

José I. ALIAGA Joaquín PÉREZ

Enrique S. QUINTANA-ORTÍ



Euro-Par 2015, Vienna

Motivation Sparse Linear Systems

#### Ubiquitous problem:

- Boundary value problems and FEM for PDEs
- Quantum chemistry
- Economic modeling
- Web search
- Information retrieval
- Pattern recognition





# → Intel MKL, IBM ESSL, NVIDIA cuSPARSE,...



#### Motivation Power/Energy/Utilization walls

#### Accelerator computing

- Moore Law in place
- End of Dennard's scaling
- Dark silicon



Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten Dotted line extrapolations by C. Moore



# Motivation GPU Computing

#### Why?

- Moderate Price
- High performance
- Favorable throughput-per-Watt
- Powerful and simple APIs (remember Cell B.E.)
  - OpenACC
  - CUDA
  - OpenCL



June 2015: 90 systems equipped with accelerators

- 52 with NVIDIA GPUs
- 4 with ATI Radeon GPUs
- 4 with NVIDIA GPUs + Intel Xeon Phi



#### Motivation Efficiency of sparse solvers on GPUs

#### Use of libraries (CUBLAS, cuSPARSE) to ease development:

while 
$$(k < maxiter) \& (res > \varepsilon)$$
  
 $z_k := Ad_k$   
 $\rho_k := \frac{\beta_k}{d_k^T z_k}$   
 $\gamma_k := \beta_k$   
 $x_{k+1} := x_k + \rho_k d_k$   
 $r_{k+1} := r_k - \rho_k z_k$   
 $\beta_{k+1} := r_{k+1}^T r_{k+1}$   
 $\alpha_k := \frac{\beta_{k+1}}{\gamma_k}$   
 $d_{k+1} := r_{k+1} + \alpha_k d_k$   
 $res := || r_{k+1} ||_2$   
 $k := k + 1$   
end

```
while(( k < maxiter ) && ( res > epsilon ) ){
   SSpMV <<<Gs,Bs>>> ( n, rowA, colA, valA, d, z );
   tmp = cublasSdot ( n, d, 1, z, 1 );
   rho = beta / tmp;
   gamma = beta;
   cublasSaxpy (n, rho, d, 1, x, 1 );
   cublasSaxpy (n, -rho, z, 1, r, 1 );
   beta = cublasSdot( n, r, 1, r, 1 );
   alpha = beta / gamma;
   cublasSscal (n, alpha, d, 1 );
   res = sqrt( beta );
   k++;
} // end-while
```



### Motivation Efficiency of sparse solvers on GPUs

Kernel fusion:

- Improve memory locality
- Reduce overhead of procedure calls
- Richer space for compiler optimizations...

```
while(( k < maxiter ) && ( res > epsilon ) ){
   SSpMV <<<Gs,Bs>>> ( n, rowA, colA, valA, d, z );
   tmp = cublasSdot ( n, d, 1, z, 1 );
   rho = beta / tmp;
   gamma = beta;
   cublasSaxpy (n, rho, d, 1, x, 1 );
   cublasSaxpy (n, -rho, z, 1, r, 1 );
   beta = cublasSdot( n, r, 1, r, 1 );
   alpha = beta / gamma;
   cublasSscal (n, alpha, d, 1 );
   res = sqrt( beta );
   k++;
} // end-while
```



#### Motivation Efficiency of sparse solvers on GPUs

#### Waste of energy for CPU: Hard to leverage energy-saving C-states



# Outline

#### Systematic Fusion of CUDA Kernels for Iterative Solvers

- Iterative solvers for sparse linear systems
  - Identification of kernels
- Fusion
  - Characterization of kernels
- Case studies: BiCG and others
- Experiments
- Concluding remarks



# Iterative Solvers for Sparse Linear Systems Identification of kernels

#### CG, BiCG, BiCGStab,... with preconditioner

 $A \to M$ Initialize  $r_0, r_0^*, p_0, p_0^*, x_0, \sigma_0, \tau_0; j := 0$ while  $(\tau_i > \tau_{\max})$ 1.  $v_j := A p_j$ 2.  $\alpha_j := \sigma_j / (v_j, p_j^*)$ 3.  $x_{j+1} := x_j + \alpha_j p_j$ 4.  $r_{j+1} := r_j - \alpha_j v_j$ 5.  $z_i := M^{-1} r_{i+1}$ 6.  $v_i^* := A^T p_i^*$ 7.  $r_{j+1}^* := r_j^* - \alpha_j v_j^*$ 8.  $z_j^* := M^{-1} r_{j+1}^*$ 9.  $\zeta_j := (z_j, r_{j+1}^*)$ 10.  $\beta_i := \zeta_j / \sigma_j$ 11.  $\sigma_i = \zeta_i$ 12.  $p_{j+1} := z_j + \beta_j p_j$ 13.  $p_{i+1}^* := z_i^* + \beta_i p_i^*$ 14.  $\tau_{j+1} := ||r_{j+1}||_2$ i := i + 1endwhile

Compute Jacobi preconditioner Loop for iterative solver SpMV 1. 2.DOT 3. AXPY 4. AXPY 5. JPRED (Jacobi preconditioner) 6. SpMV 7. AXPY JPRED (Jacobi preconditioner) 8. 9. DOT 10. Scalar op 11. Scalar op 12. XPAY (AXPY-like) 13. XPAY (AXPY-like) 14. Vector 2-norm (DOT + sqrt)



SpMV (different formats) Preconditioner appl. BLAS-1: AXPY-like, DOT



#### Fusion Characterization of kernels

#### Mapped vs unmapped:

• A kernel *K* performs a *mapped access* to a vector *v* if each thread of *K* accesses one element of *v*, independently of other threads, and the global access is coalesced

Operation		Input $vector(s)$		Output vector
		x	$y/M^{-1}$	y
AXPY	$y := \alpha x + y$	mapped	mapped	mapped
XPAY	$y := \alpha y + x$	mapped	mapped	$\operatorname{mapped}$
JPred	$y := M^{-1}x$	mapped	mapped	$\operatorname{mapped}$
DOT	$\alpha := x^T y = (x, y)$	mapped	mapped	unmapped
SPMV vector CSR	y := Ax	unmapped	—	unmapped
SPMV scalar CSR	y := Ax	unmapped	_	$\operatorname{mapped}$
SPMV ELL	y := Ax	unmapped	—	$\mathbf{mapped}$



### Fusion Conditions

# Merge of $K_1$ and $K_2$ possible?

- $K_1$  and  $K_2$  are independent
- If  $K_1$  produces an output v and  $K_2$  takes this vector as an input (RAW dependency):

 $K_1$  and  $K_2$  perform mapped accesses to the output/input vector v



- 1. Both kernels apply the same mapping of threads to the vector elements shared (exchanged) via registers
- 2. Both kernels apply the same mapping of thread blocks to the vector elements shared (exchanged) via shared memory
- 3. A global barrier is not necessary between the two kernels



 $K_1$  and  $K_2$  perform mapped accesses to the input/output vector v

Operat	ion	Input ve	$\operatorname{ctor}(\mathrm{s})$	Output vector
		x	$y/M^{-1}$	y
AXPY	$y := \alpha x + y$	mapped	mapped	mapped
XPAY	$y := \alpha y + x$	mapped	mapped	$\operatorname{mapped}$
JPred	$y := M^{-1}x$	mapped	mapped	mapped
DOT	$\alpha := x^T y = (x, y)$	mapped	mapped	unmapped
SPMV vector CSR	y := Ax	unmapped	—	unmapped
SPMV scalar CSR	y := Ax	unmapped	—	$\mathbf{mapped}$
SPMV ELL	y := Ax	unmapped	—	mapped



 $K_1$  and  $K_2$  perform mapped accesses to the input/output vector v

Operat	ion		Input ve	$\operatorname{ctor}(\mathbf{s})$	Output vector
			x	$y/M^{-1}$	y
AXPY	$y := \alpha x + y$		mapped	mapped	mapped
XPAY	$y := \alpha y + x$		$\mathbf{mapped}$	mapped	$\mathbf{mapped}$
JPred	$y := M^{-1}x$		mapped	mapped	$\mathbf{mapped}$
DOT	$\alpha := x^T y = ($	[x,y)	mapped	mapped	unmapped
SPMV vector CSR	y := Ax		unmapped	—	unmapped
SPMV scalar CSR	y := Ax		unmapped	—	mapped
${ m SPMV}$ ELL	y := Ax		unmapped	—	$\operatorname{mapped}$



 $K_1$  and  $K_2$  perform mapped accesses to the input/output vector v

Operat	ion		Input ve	$\operatorname{ctor}(s)$	Output vector
			x	$y/M^{-1}$	y
AXPY	$y := \alpha x + y$		mapped	mapped	mapped
XPAY	$y := \alpha y + x$		mapped	mapped	$\mathbf{mapped}$
JPred	$y := M^{-1}x$		mapped	mapped	mapped
DOT	$\alpha := x^T y = 0$	(x, y)	mapped	mapped	unmapped
SPMV vector CSR	y := Ax		unmapped	—	unmapped
SPMV scalar CSR	y := Ax		unmapped	—	$\operatorname{mapped}$
${ m SpMV}$ ELL	y := Ax		unmapped	—	$\operatorname{mapped}$



 $K_1$  and  $K_2$  perform mapped accesses to the input/output vector v

Operation		Input vector(s)		Output vector
		x	$y/M^{-1}$	y
AXPY	$y := \alpha x + y$	mapped	mapped	mapped
XPAY	$y := \alpha y + x$	mapped	mapped	mapped
JPRED	$y := M^{-1}x$	mapped	mapped	mapped
DOT	$\alpha := x^T y = (x, y)$	mapped	mapped	unmapped
SPMV vector CSR	y := Ax	unmapped	—	unmapped
SPMV scalar CSR	y := Ax	unmapped	—	mapped
SPMV ELL	y := Ax	unmapped	—	mapped



 $K_1$  and  $K_2$  perform mapped accesses to the input/output vector v

Operation		Input vector(s)		Output vector
		x	$y/M^{-1}$	y
AXPY	$y := \alpha x + y$	mapped	mapped	mapped
XPAY	$y := \alpha y + x$	mapped	mapped	mapped
JPred	$y := M^{-1}x$	mapped	mapped	mapped
DOT	$\alpha := x^T y = (x, y)$	mapped	mapped	unmapped
SPMV vector CSR	y := Ax	unmapped	-	unmapped
SPMV scalar CSR	y := Ax	unmapped	_	$\operatorname{mapped}$
SPMV ELL	y := Ax	unmapped	_	mapped



#### Two-stage DOT:

- DOT<sub>ini</sub>
  - CUDA kernel
  - Costly element-wise products and reduction within a single thread block
  - Mapped input  $\rightarrow$  Fusible with AXPY, XPAY, JPRED, SpMV (scalar CSR, ELL)
  - Unmapped output
- DOT<sub>fin</sub>
  - Final reduction into scalar result. Very small cost
  - CUDA routine, implemented as a collection of CUDA kernels
  - Each kernel with mapped input and unmapped output



#### When?

• Reduce as much as possible the number of macro-kernels

#### How?

- Systematic analysis of pairs of kernels
- In general, merge kernels by "stacking" their codes one after another respecting dependencies
- When merging two DOT<sub>ini</sub> or two DOT<sub>fin</sub>, interleave their codes



# Case Studies BiCG

#### Dependencies (two-stage DOT):







# Case Studies BiCG

#### **Fusions:**

 $\begin{array}{l} A \to M \\ \text{Initialize } r_0, r_0^*, p_0, p_0^*, x_0, \sigma_0, \tau_0; j := 0 \\ \textbf{while } (\tau_j > \tau_{\max}) \\ 1. \quad v_j := Ap_j \\ 2. \quad \alpha_j := \sigma_j / (v_j, p_j^*) \\ 3. \quad x_{j+1} := x_j + \alpha_j p_j \\ 4. \quad r_{j+1} := r_j - \alpha_j v_j \\ 5. \quad z_j := M^{-1} r_{j+1} \\ 6. \quad v_j^* := A^T p_j^* \\ 7. \quad r_{j+1}^* := r_j^* - \alpha_j v_j^* \\ 8. \quad z_j^* := M^{-1} r_{j+1}^* \\ 9. \quad \zeta_j := (z_j, r_{j+1}^*) \\ 10. \quad \beta_j := \zeta_j / \sigma_j \\ 11. \quad \sigma_j = \zeta_j \\ 12. \quad p_{j+1} := z_j^* + \beta_j p_j \\ 13. \quad p_{j+1}^* := z_j^* + \beta_j p_j^* \\ 14. \quad \tau_{j+1} := \| r_{j+1} \|_2 \\ \quad j := j+1 \\ \textbf{endwhile} \end{array}$ 





# **Case Studies**





#### UFMC and Laplace problem

Matrix	$n_z$	n	$n_z/n$
BMWCRA1_1	$10,\!641,\!602$	148,770	71.53
CRANKSEG_2	14, 148, 858	63,838	221.63
F1	26,837,113	343,791	78.06
INLINE_1	38, 816, 170	503,712	77.06
LDOOR	42,493,817	952,203	44.62
AUDIKW_1	77,651,847	$943,\!645$	82.28
FEM_3DTH2	3,489,300	$147,\!900$	23.59

Matrix	$n_z$	n	$n_z/n$
A100	6,940,000	1,000,000	6.94
A126	13,907,370	2,000,376	6.94
A159	27,986,067	4,019,679	6.94
A200	55,760,000	8,000,000	6.94
A252	$111,\!640,\!032$	16,003,001	6.94

#### Server:

- Intel core i7-3770K (4 cores, 3.5 GHz) with 16 GB DDR3 RAM
- NVIDIA Kepler K20c GPU (2,496 CUDA cores, 706 MHz) with 5 GB GDDR5 RAM
- National Instruments DAS (NI9205 module and the NIcDAQ-9178)



#### Software

- Single-precision arithmetic
- Solvers: CG, BiCG, BiCGStab with/out preconditioning
- SpMV: scalar CSR, vector CSR, ELL
- Implementations of solvers:
  - CUBLASL: Plain implementation + ad-hoc SpMV + CUBLAS
  - CUBLASN: Maintain scalars in GPU (new CUBLAS interface)
  - CUDA: Replace CUBLAS by ad-hoc kernels (including two-stage DOT)
  - MERGE: Applies fusion
  - MERGE\_10: Check convergece every 10 iterations
- Execution modes: CUDA polling/blocking



3 solvers x with/out preconditioning x 3 SpMV x 5 implementations x 2 exec. modes x 12 matrices...



- Variations of time/energy w.r.t. CUBLASL in polling mode
- Average results for all matrices giving the same weight to all cases
- Vector CSR for UFMC and ELL for Laplace (best options)



- In polling mode, CUBLASL and CUDA incur small overhead in time and energy
- Blocking mode trades off time for energy



- MERGE and MERGE 10 reduce both time and energy
- Polling mode prioritizes reduction of time vs energy
- Blocking mode prioritizes energy vs time





#### Similar trends for BiCG, BiCGStab











# Concluding Remarks

Systematic fusion to derive efficient versions of three popular iterative solvers (with/out preconditioning) for sparse linear systems

- Analysis of type of thread accesses on kernel input/outputs
- Identification of data dependencies
- Eager fusion
- Divide DOT into two stages
- Reduce #kernels to improve performance/energy consumption
  - $10 \rightarrow 5$  for PCG,  $13 \rightarrow 5$  for PBiCG and  $14 \rightarrow 8$  for PBiCGStab
  - Remarkable energy savings when executed in blocking mode
  - Match performance of original versions when executed in polling mode



# THANKS... and QUESTIONS?



