

# Sparse Matrix-Vector Multiplication (Sparsity, Bebop)

María Jesús Garzarán  
CS 498: Compiler Optimizations, Fall 2006

(based on the slides from Markus Pueschel, CMU)

University of Illinois at Urbana-Champaign



## MMM versus MVM

- ◆ MMM
  - BLAS3
  - $O(n^2)$  data,  $O(n^3)$  computation, implies  $O(n)$  reuse per data.
- ◆ MVM:  $y=Ax$ 
  - BLAS2
  - $O(n^2)$  data,  $O(n^2)$  computation
  - Optimizations that are still useful:
    - Cache blocking?
    - Register blocking?
    - Unrolling?
    - Scalar replacent?
    - Add/mult interleaving, skewing?



# MMM versus MVM: Performance

---

- ◆ Performance for 2000 x 2000 matrices
- ◆ Best code out of ATLAS, vendor lib. , Goto's assembly codes

Processor and compiler	Clock (MHz)	Data cache sizes	DGEMV (MFLOPS)	DGEMM (MFLOPS)
Sun UltraSPARC III Sun C v6.0	333	L1: 16 KB L2: 2 MB	58	425
Intel Pentium III Mobile (Coppermine) Intel C v6.0	800	L1: 16 KB L2: 256 KB	147	590
IBM Power4 IBM xlc v6	1300	L1: 64 KB L2: 1.5 MB L3: 32 MB	915	3500
Intel Itanium 2 Intel C v7.0	900	L1: 16 KB L2: 256 KB L3: 3 MB	1330	3500



Source: Eun-JIN Im, Katherine A. Yelick, Richard Vuduc. SPARSITY: "An Optimization Framework for Sparse Matrix Kernels, Int'l Journal of High Performance Computing App. 18 (1), pp. 135 -158, 2004 3

## Sparse MVM

---

- ◆  $y = Ax$ ,  $A$  sparse but known
- ◆ Important routine in:
  - Finite element methods
  - PDE solvers
  - Physical/chemical simulations (e.g. fluid dynamics)
  - Linear programming
  - Scheduling
  - Signal processing (e.g. filters)
  - ...
- ◆ In these applications,  $y = Ax$  is performed many times.
  - Justifies one-time tuning effort



# Storage of Sparse Matrices

---

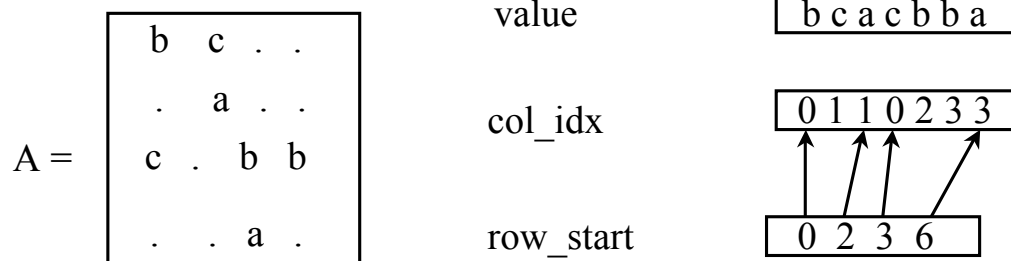
- ◆ Standard storage (as 2-D array) inefficient
  - Many zeros are stored
- ◆ Several sparse format are available.
  - Compressed Sparse Row Format
  - Blocked CSR (BCSR) Format



5

## Compressed Sparse Row (CSR)

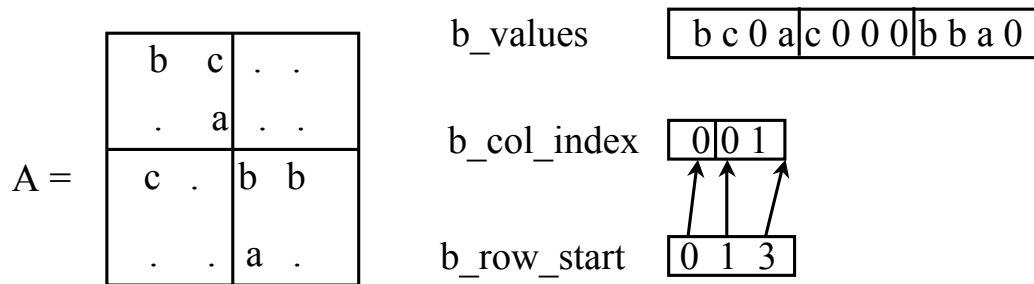
---



6

# Blocked CSR (BCSR) Format

- ◆ Block A into  $r \times c$  blocks
- ◆ Assume  $r=c=2$



7

# Direct Implementation $y=Ax$ , A in CSR

```

void smvm_1x1( int m, const double* value, const int* col_idx, const inst*
row_start, const double* x, double y)
{
  int i, jj;
  /* loop over rows */
  for (i=0; i<m; i++) {
    double y_i= y[i];
    /*loop over non-zero elements */
    for (jj = row_start[i]; jj<row_start[i+1];
        jj++, col_idx++, value++) {
      y_i += value[0]*x[col_idx[0]];
    }
    y[i] = y_i;
  }
}

```

scalar replacement

indirect array accessing



8

# Impact of Matrix-Sparsity on Performance

---

- ◆ Addressing overhead (dense MVM vs. dense MVM in CSR):
  - ~2x slower (mflops, example only)
- ◆ Irregular structure
  - ~5x slower (mflops, example only) for “random” sparse matrices
- ◆ Fundamental difference between MVM and sparse MVM (SMVM)
  - Sparse MVM is input dependent (sparsity pattern of A)
  - Changing the order of computation (blocking) requires change of data structure (CSR)



9

## Bebop/Sparsity: SMVM Optimizations

---

- ◆ Register blocking
- ◆ Cache blocking



10

# Register Blocking

---

- ◆ Idea: divide SMVM  $y = Ax$  into micro (dense) MVMs of matrix size  $r \times x$ 
  - Store  $A$  in  $r \times c$  block CSR ( $r \times c$  BCSR)
- ◆ Advantages:
  - Reuse of  $x$  and  $y$  (as for dense MVM)
  - Reduces index overhead
- ◆ Disadvantages:
  - Computational overhead (zeros added)
  - Storage overhead (for  $A$ )



11

## Example: $y = Ax$ in $2 \times 2$ BCSR

---

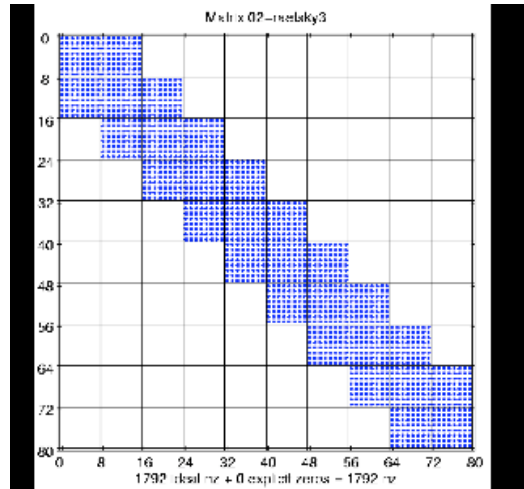
```
void smvm_2x2( int bm, const int *b_row_start, const int *b_col_idx, const double
    *b_value, const double *x, double *y)
    int I, jj;
    /* loop over block rows */
    for (i=0; i<bm; i+=2, y+=2){
        register double d0 = y[0];
        register double d1 = y[1];
        /*dense micro MVM */
        for (jj = b_row_start[i]; jj<b_row_start[i+1];
            jj++, b_col_idx++, b_value += 2*2) {
            d0 += b_value[0] * x[b_col_idx[0]+0];
            d1 += b_value[2] * x[b_col_idx[0]+0];
            d0 += b_value[1] * x[b_col_idx[0]+1];
            d1 += b_value[3] * x[b_col_idx[0]+1];
        }
        y[0]=d0;
        y[1]=d1;
    }
```



12

# Which Block Size (r x c) is Optimal?

- ◆ Example: ~20,000 x 20,000 matrix with perfect 8x8 block structure, 0,33% non-zero entries
- ◆ In this case, no overhead when blocked r x c, with r, c divides 8.

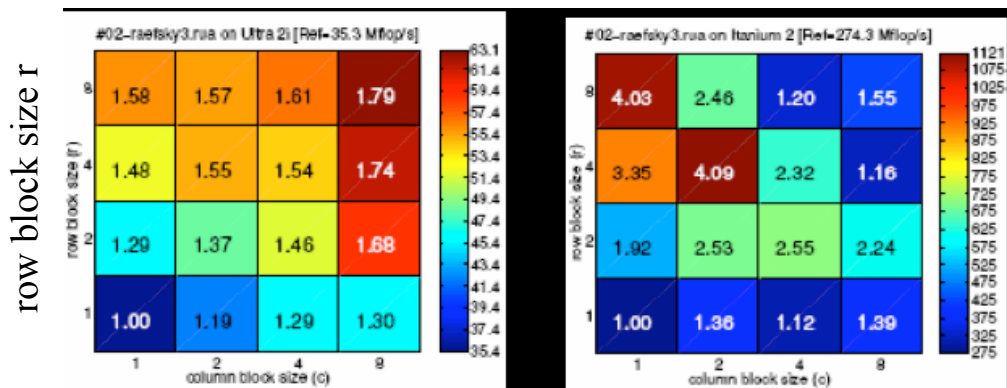


Source: Vuduc, LLNL

# Speed-up through r x c Blocking

Ultra 2i

Itanium 2



col. Block size c

col. Block size c

- Machine dependent
- Hard to predict



Source: Eun-JIN Im, Katherine A. Yelick, Richard Vuduc. SPARSITY: "An Optimization Framework for Sparse Matrix Kernels, Int'l Journal of High Performance Computing App. 18 (1), pp. 135 -158, 2004

## How to Find the Best Register Blocking for given A

---

- ◆ Best blocksize is hard to predict (see the previous slide)
- ◆ Searching over all  $r \times c$  (within a range 1..12) is expensive
  - Conversion of A in CSR to BCSR roughly as expensive as 10 SMVMs
- ◆ Solution: Performance model for given A



15

## Performance Model for given A

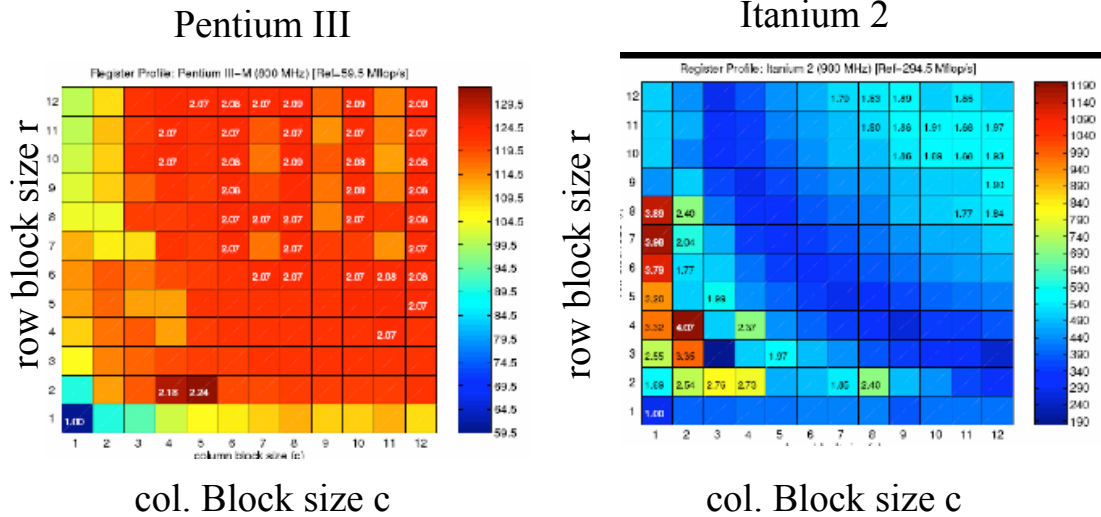
---

- ◆ Model for given A built from
  - Gain of blocking:  
 $Gr,c = \text{Performance } r \times c \text{ BCSR} / \text{performance CSR for dense MVM}$   
Machine depending, independent of matrix A
  - Computational overhead:  
 $Or,c = \text{size of A in } r \times c \text{ BCSR} / \text{size of A in CSR}$   
machine independent, dependent of A  
computed by scanning only a fraction of the matrix
- ◆ Model: Performance gain from  $r \times c$  blocking of A:  
 $Pr,c = Gr,c / Or,c$
- ◆ For given A, use this model to search over all  $r, c$  in  $\{1..12\}$



16

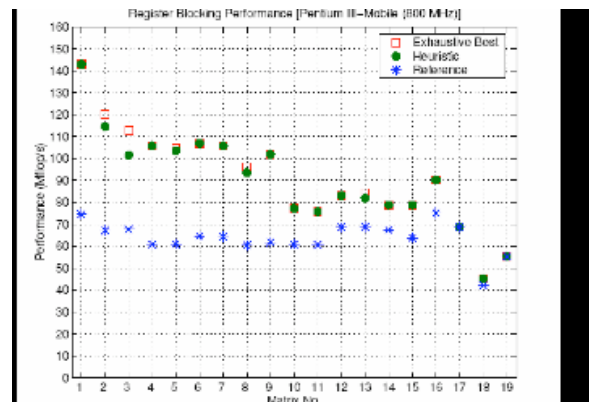
# Gain from Blocking (Dense Matrix in BCSR)



Source: Eun-JIN Im, Katherine A. Yelick, Richard Vuduc. SPARSITY: "An Optimization Framework for Sparse Matrix Kernels, Int'l Journal of High Performance Computing App. 18 (1), pp. 135 -158, 2004 17

## Register Blocking: Experimental Results

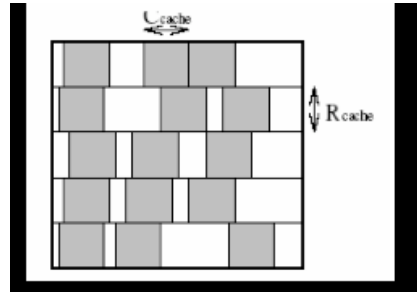
- ◆ Paper applies method to a large set of sparse matrices
- ◆ Performance gains between 1x (no gain) for very unstructured matrices and 4x



Source: Eun-JIN Im, Katherine A. Yelick, Richard Vuduc. SPARSITY: "An Optimization Framework for Sparse Matrix Kernels, Int'l Journal of High Performance Computing App. 18 (1), pp. 135 -158, 2004 18

# Cache Blocking

- ◆ Idea: divide sparse matrix into blocks of sparse matrices

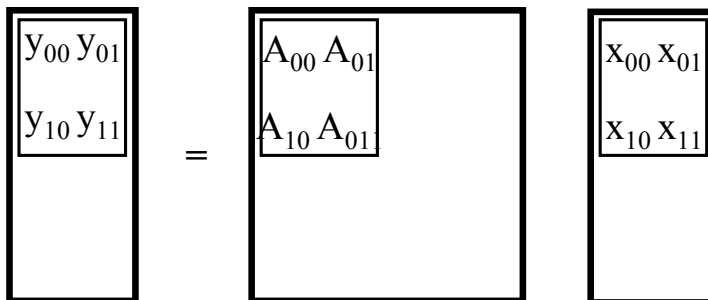


- ◆ Experiments:
  - Requires very large matrices (x and y do not fit in cache)
  - Speedup up to 80% only for few matrices, with 1x1 BSCR



Source: Eun-JIN Im, Katherine A. Yelick, Richard Vuduc. SPARSITY: “An Optimization Framework for Sparse Matrix Kernels, Int’l Journal of High Performance Computing App. 18 (1), pp. 135 -158, 2004 19

# Multiple Vector Optimization



OPTION 1

OPTION 2

$$\begin{aligned} (1) \quad y_{00} &= A_{00} x_{00} + A_{01} x_{10} \\ (2) \quad y_{10} &= A_{10} x_{00} + A_{11} x_{10} \end{aligned}$$

$$\begin{aligned} (1) \quad y_{00} &= A_{00} x_{00} + A_{01} x_{10} \\ (2) \quad y_{10} &= A_{10} x_{00} + A_{11} x_{10} \\ (3) \quad y_{01} &= A_{00} x_{01} + A_{01} x_{11} \\ (4) \quad y_{11} &= A_{10} x_{01} + A_{11} x_{11} \end{aligned}$$

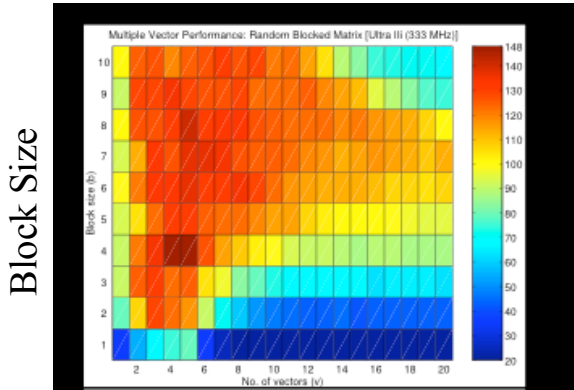


$$\begin{aligned} &\dots \\ (nz + 1) \quad y_{01} &= A_{00} x_{01} + A_{01} x_{11} \\ (nz + 2) \quad y_{11} &= A_{10} x_{01} + A_{11} x_{11} \end{aligned}$$

# Multiple Vector Optimization

Ultra Ili

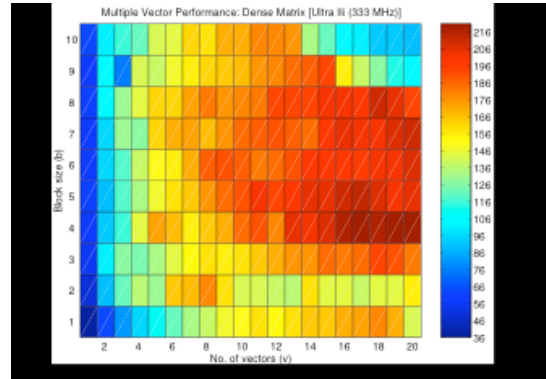
Random Blocked Matrix



Number of Vectors

Ultra Ili

Dense Matrix



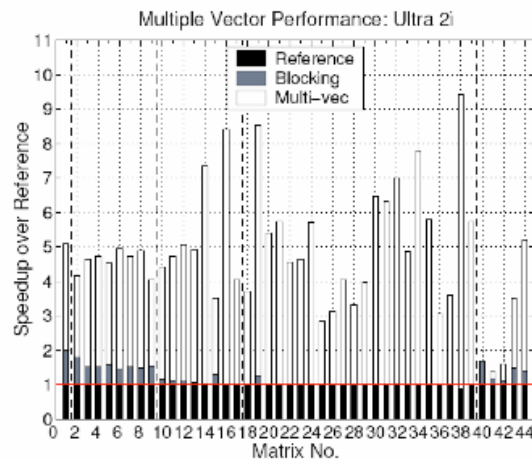
Number of Vectors



Source: Eun-JIN Im, Katherine A. Yelick, Richard Vuduc. SPARSITY: "An Optimization Framework for Sparse Matrix Kernels, Int'l Journal of High Performance Computing App. 18 (1), pp. 135 -158, 2004 21

# Multiple Vector Optimization

- ◆ Experiments: Up to 9x speedup for 9 vectors



Source: Eun-JIN Im, Katherine A. Yelick, Richard Vuduc. SPARSITY: "An Optimization Framework for Sparse Matrix Kernels, Int'l Journal of High Performance Computing App. 18 (1), pp. 135 -158, 2004 22

# Principles in Bebop/Sparsity

---

- ♦ Optimization for memory hierarchy = increasing locality
  - Blocking for registers (micro-MMMs) + **change of data structure for A**
  - **Less important: blocking for cache**
  - Optimizations are **input dependent** (on sparse structure of A)
- ♦ Fast basick blocks for small sizes (micro-MMM):
  - Loop unrolling (reduced loop overhead)
  - Some scalar replacement (enables better compiler optimizations)
- ♦ Search for the fastest over a relevant set of algorithm/implementation alternatives ( $=r,c$ )
  - **Use of performance model** (versus measuring runtime) to evaluate expected gain



**red = different from ATLAS**