



OSKI: A Library of Automatically Tuned Sparse Matrix Kernels



Richard Vuduc (LLNL), James Demmel, Katherine Yelick

Berkeley Benchmarking and OPTimization (BeBOP) Project

bebop.cs.berkeley.edu

EECS Department, University of California, Berkeley

SIAM CSE

February 12, 2005



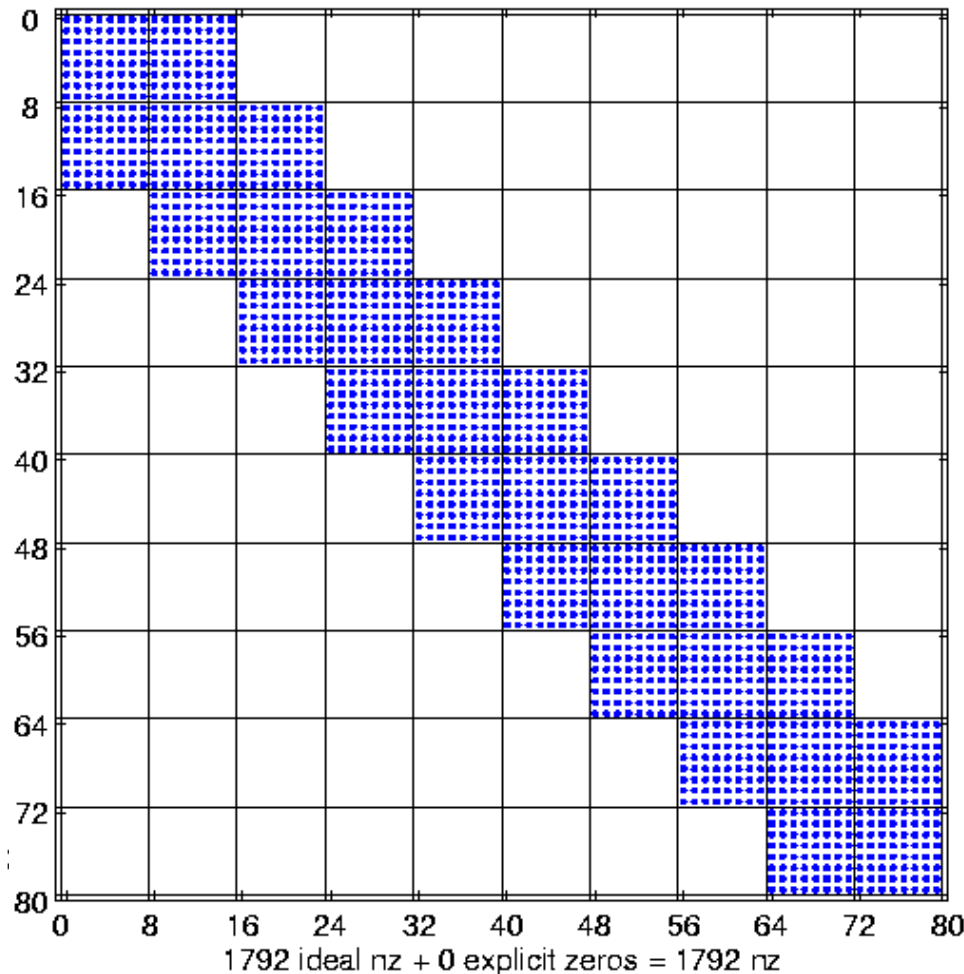
OSKI: Optimized Sparse Kernel Interface

- Sparse kernels tuned for user's matrix & machine
 - Hides complexity of run-time tuning
 - Low-level BLAS-style functionality
 - Sparse matrix-vector multiply (SpMV), triangular solve (TrSV), ...
 - Includes fast locality-aware kernels: $A^T A * x$, ...
 - Initial target: cache-based superscalar uniprocessors
- Faster than standard implementations
 - Up to **4x** faster SpMV, **1.8x** TrSV, **4x** $A^T A * x$
- For “advanced” users & solver library writers
 - Available as stand-alone open-source library (pre-release)
 - PETSc extension in progress
- Written in C (can call from Fortran)



Motivation: The Difficulty of Tuning

Matrix 02-raefsky3

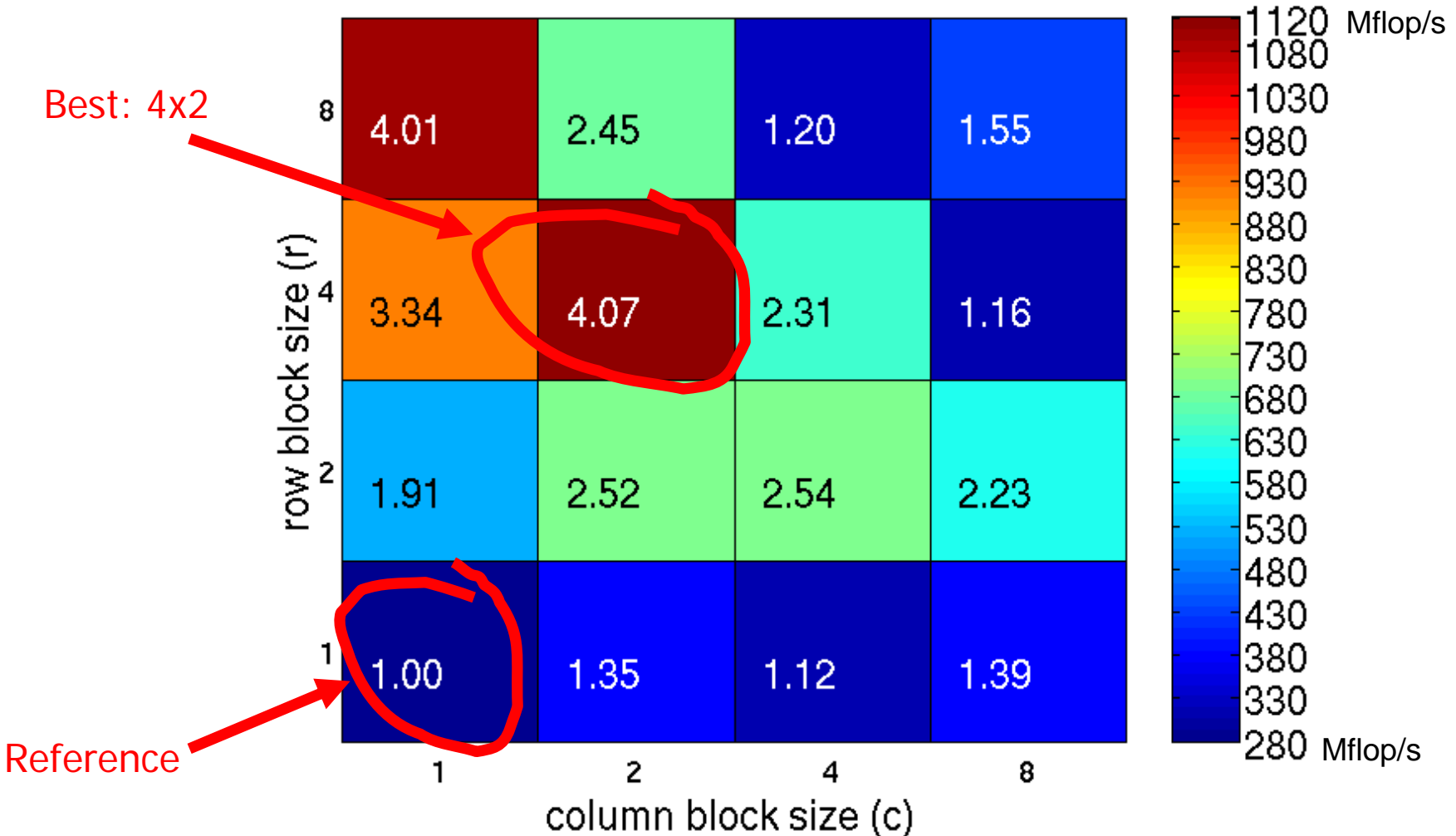


- $n = 21216$
- $nnz = 1.5 \text{ M}$
- kernel: SpMV
- Source: NASA structural analysis problem
- **8x8** dense substructure

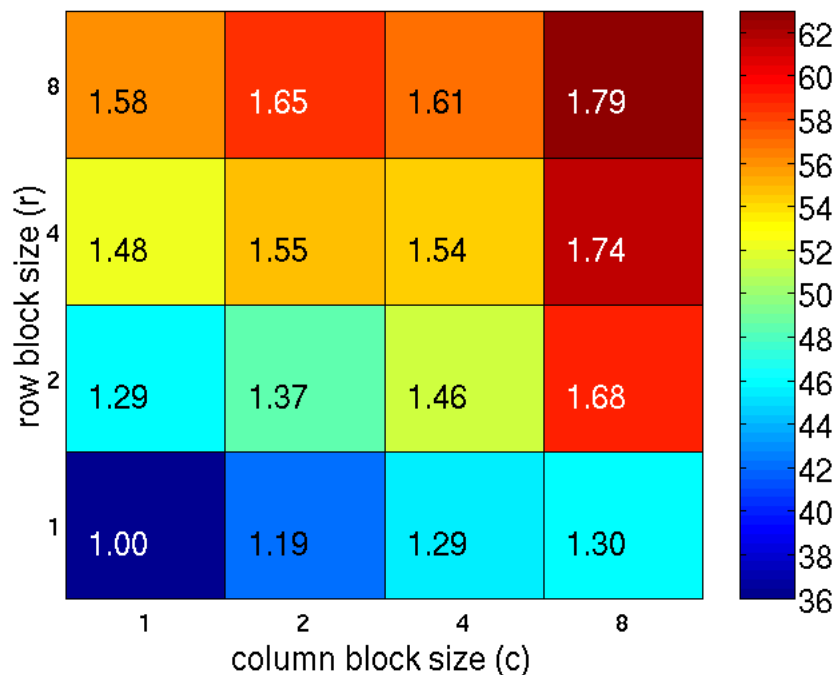


Speedups on Itanium 2: The Need for Search

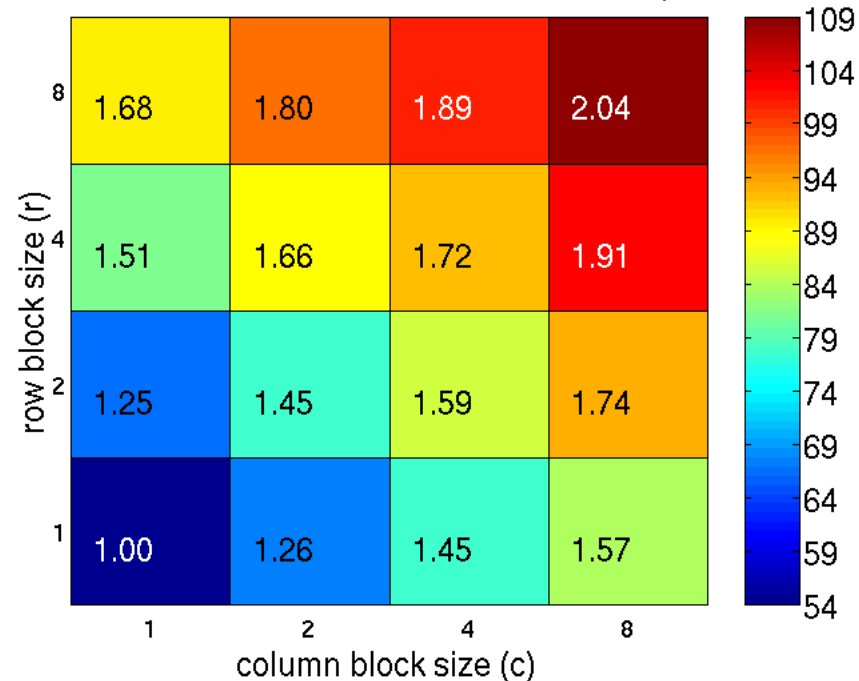
900 MHz Itanium 2, Intel C v8: ref=275 Mflop/s



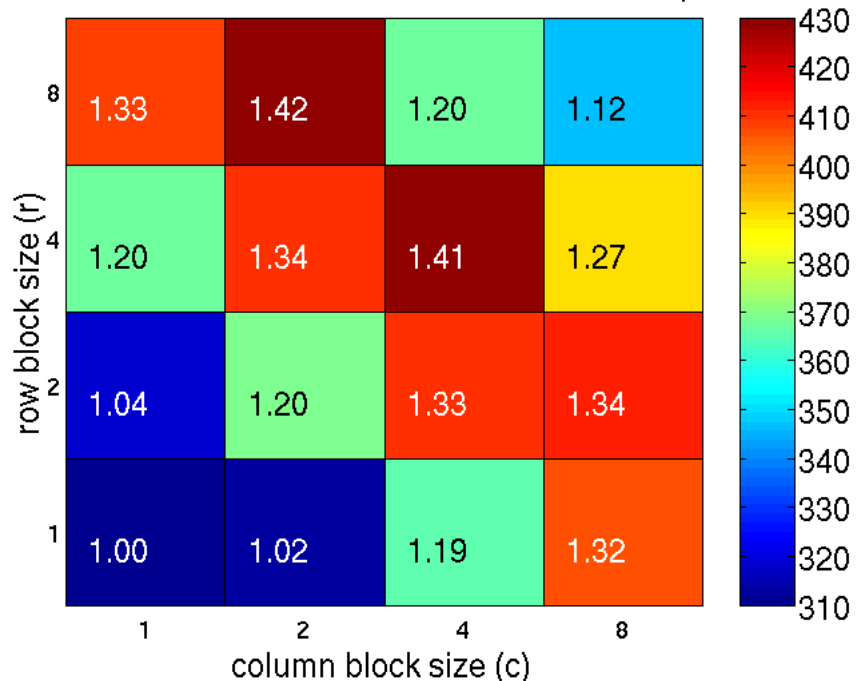
333 MHz Sun Ultra 2i, Sun C v6.0: ref=35 Mflop/s



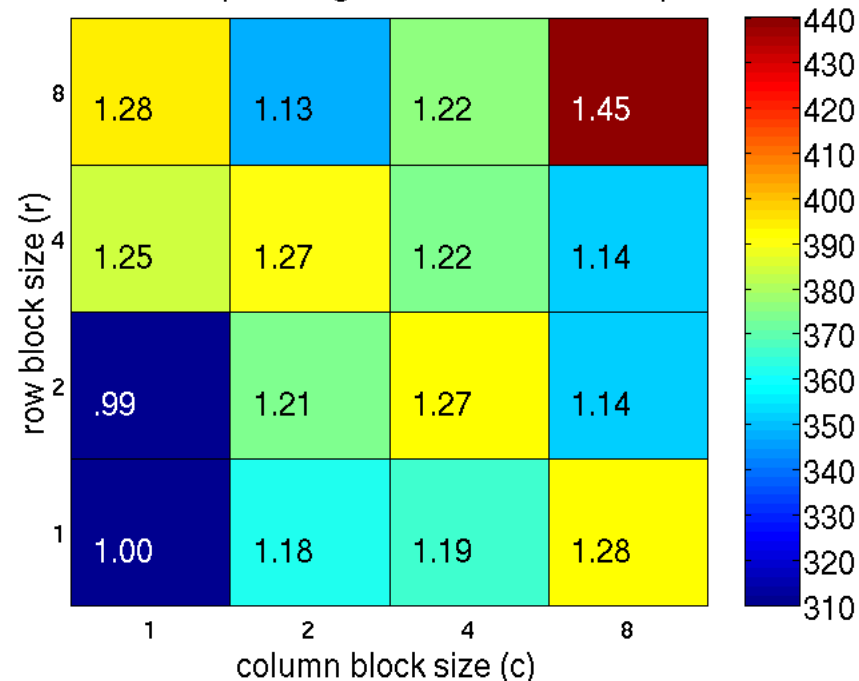
900 MHz Ultra 3, Sun CC v6: ref=54 Mflop/s



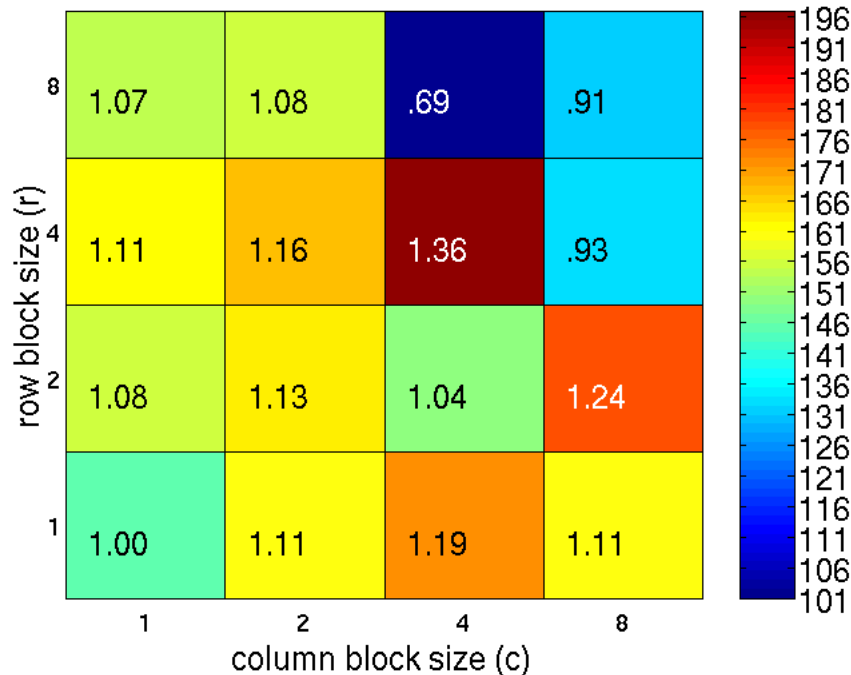
2 GHz Pentium M, Intel C v8.1: ref=308 Mflop/s



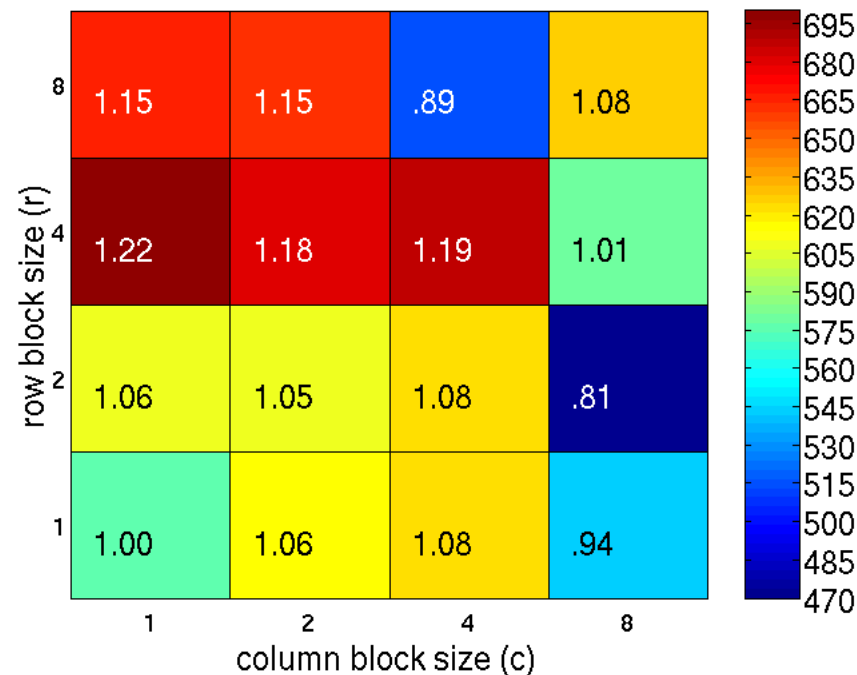
1.4 GHz Opteron, gcc 3.4.2: ref=308 Mflop/s



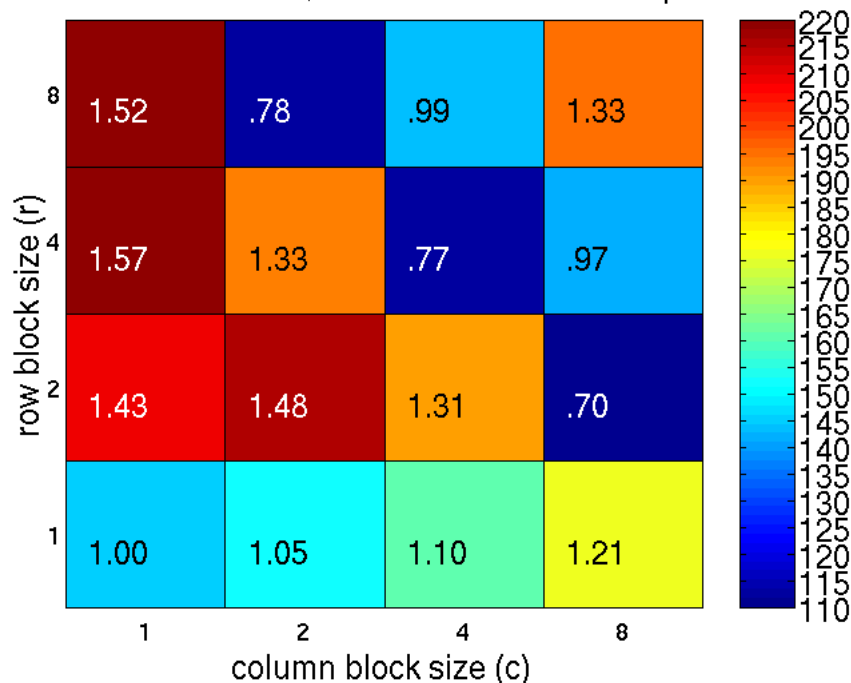
375 MHz Power3, IBM xlc v6: ref=145 Mflop/s



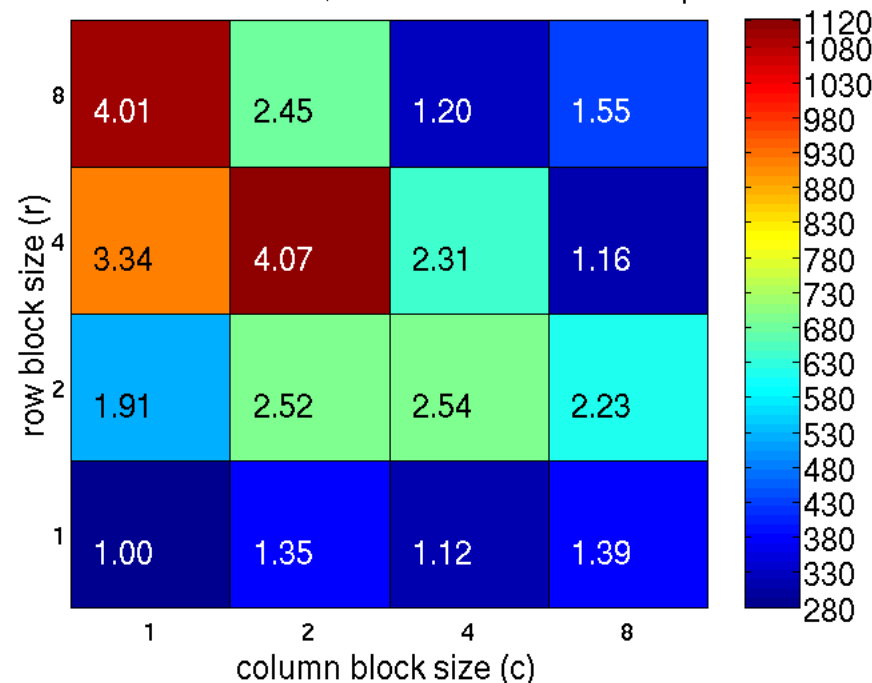
1.3 GHz Power4, IBM xlc v6: ref=577 Mflop/s



800 MHz Itanium, Intel C v7: ref=146 Mflop/s

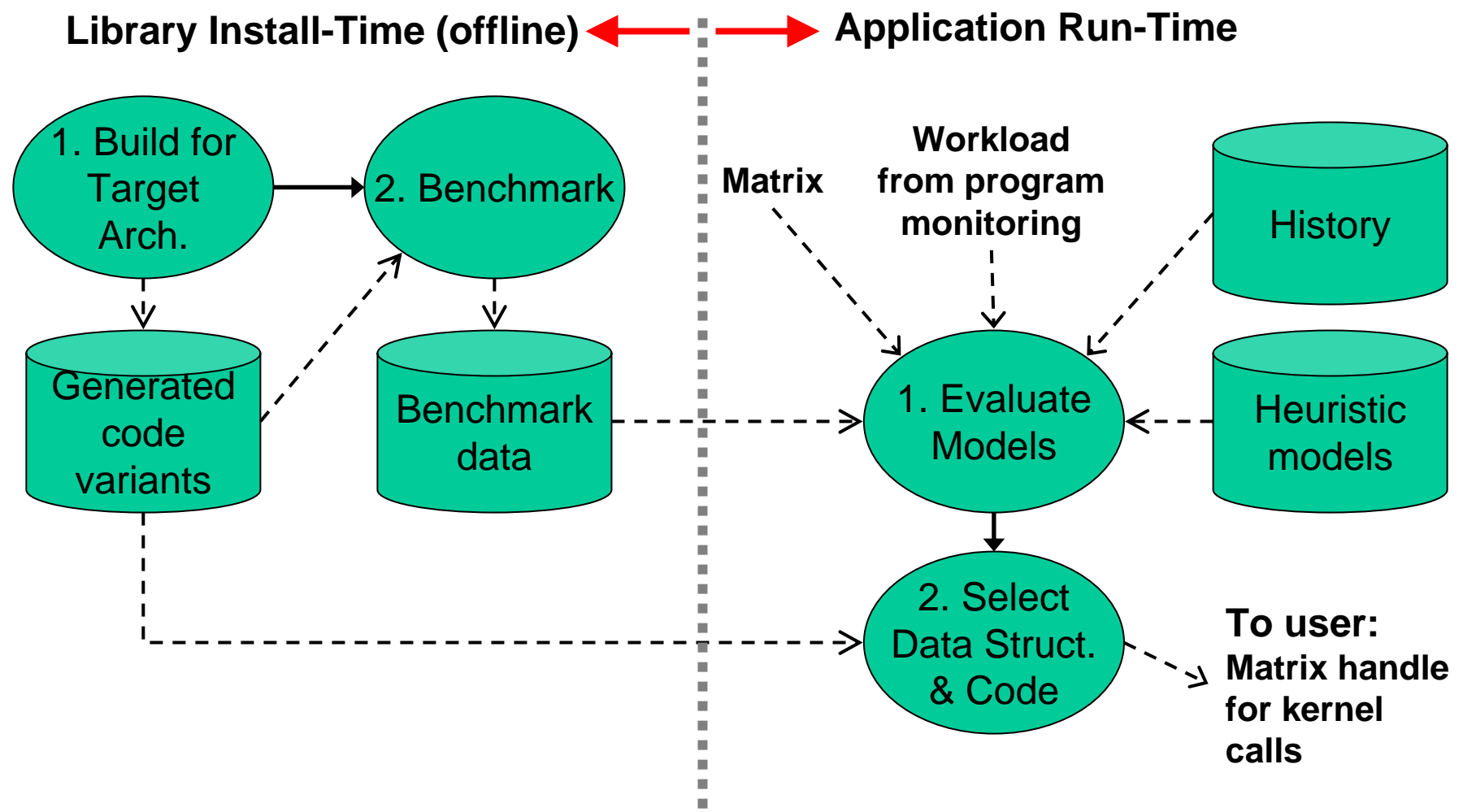


900 MHz Itanium 2, Intel C v8: ref=275 Mflop/s





How OSKI Tunes (Overview)



Extensibility: Advanced users may write & dynamically add “Code variants” and “Heuristic models” to system.



Cost of Tuning

- Non-trivial run-time tuning cost: up to ~ 40 mat-vecs
 - Dominated by conversion time
- Design point: user calls “tune” routine explicitly
 - Exposes cost
 - Tuning time limited using estimated workload
 - Provided by user or inferred by library
- User may save tuning results
 - To apply on future runs with similar matrix
 - Stored in “human-readable” format



How to Call OSKI: Basic Usage

- May gradually migrate existing apps
 - Step 1: “Wrap” existing data structures
 - Step 2: Make BLAS-like kernel calls

```
int* ptr = ..., *ind = ...; double* val = ...; /* Matrix, in CSR format */  
double* x = ..., *y = ...; /* Let x and y be two dense vectors */
```

```
/* Compute  $y = \beta \cdot y + \alpha \cdot A \cdot x$ , 500 times */  
for( i = 0; i < 500; i++ )  
    my_matmult( ptr, ind, val,  $\alpha$ , x,  $\beta$ , y );
```



How to Call OSKI: Basic Usage

- May gradually migrate existing apps
 - Step 1: “Wrap” existing data structures
 - Step 2: Make BLAS-like kernel calls

```
int* ptr = ..., *ind = ...; double* val = ...; /* Matrix, in CSR format */
double* x = ..., *y = ...; /* Let x and y be two dense vectors */
/* Step 1: Create OSKI wrappers around this data */
oski_matrix_t A_tunable = oski_CreateMatCSR(ptr, ind, val, num_rows,
    num_cols, SHARE_INPUTMAT, ...);
oski_vecview_t x_view = oski_CreateVecView(x, num_cols, UNIT_STRIDE);
oski_vecview_t y_view = oski_CreateVecView(y, num_rows, UNIT_STRIDE);

/* Compute  $y = \beta \cdot y + \alpha \cdot A \cdot x$ , 500 times */
for( i = 0; i < 500; i++ )
    my_matmult( ptr, ind, val,  $\alpha$ , x,  $\beta$ , y );
```



How to Call OSKI: Basic Usage

- May gradually migrate existing apps
 - Step 1: “Wrap” existing data structures
 - Step 2: Make BLAS-like kernel calls

```
int* ptr = ..., *ind = ...; double* val = ...; /* Matrix, in CSR format */
double* x = ..., *y = ...; /* Let x and y be two dense vectors */
/* Step 1: Create OSKI wrappers around this data */
oski_matrix_t A_tunable = oski_CreateMatCSR(ptr, ind, val, num_rows,
    num_cols, SHARE_INPUTMAT, ...);
oski_vecview_t x_view = oski_CreateVecView(x, num_cols, UNIT_STRIDE);
oski_vecview_t y_view = oski_CreateVecView(y, num_rows, UNIT_STRIDE);

/* Compute  $y = \beta \cdot y + \alpha \cdot A \cdot x$ , 500 times */
for( i = 0; i < 500; i++ )
    oski_MatMult(A_tunable, OP_NORMAL,  $\alpha$ , x_view,  $\beta$ , y_view); /* Step 2 */
```



How to Call OSKI: Tune with Explicit Hints

- User calls “tune” routine
 - May provide explicit tuning hints (OPTIONAL)

```
oski_matrix_t A_tunable = oski_CreateMatCSR( ... );
    /* ... */
/* Tell OSKI we will call SpMV 500 times (workload hint) */
oski_SetHintMatMult(A_tunable, OP_NORMAL,  $\alpha$ , x_view,  $\beta$ , y_view, 500);
/* Tell OSKI we think the matrix has 8x8 blocks (structural hint) */
oski_SetHint(A_tunable, HINT_SINGLE_BLOCKSIZE, 8, 8);

oski_TuneMat(A_tunable); /* Ask OSKI to tune */

for( i = 0; i < 500; i++ )
    oski_MatMult(A_tunable, OP_NORMAL,  $\alpha$ , x_view,  $\beta$ , y_view);
```



How the User Calls OSKI: Implicit Tuning

- Ask library to infer workload
 - Library profiles all kernel calls
 - May periodically re-tune

```
oski_matrix_t A_tunable = oski_CreateMatCSR( ... );  
/* ... */  
  
for( i = 0; i < 500; i++ ) {  
    oski_MatMult(A_tunable, OP_NORMAL,  $\alpha$ , x_view,  $\beta$ , y_view);  
    oski_TuneMat(A_tunable); /* Ask OSKI to tune */  
}
```



Additional Features

- Embedded scripting language for selecting customized, complex transformations
 - Mechanism to save/restore transformations

```
/* In "my_app.c" */  
fp = fopen("my_xform.txt", "rt");  
fgets(buffer, BUFSIZE, fp);  
  
oski_ApplyMatTransform(A_tunable,  
    buffer);  
oski_MatMult(A_tunable, ...);
```

```
# In file, "my_xform.txt"  
# Compute  $A_{fast} = P * A * P^T$  using  
  Pinar's reordering algorithm  
A_fast, P =  
    reorder_TSP(InputMat);  
# Split  $A_{fast} = A_1 + A_2$ , where  $A_1$  in 2x2  
  block format,  $A_2$  in CSR  
A1, A2 =  
    A_fast.extract_blocks(2, 2);  
  
return transpose(P) * (A1+A2) * P;
```



Additional Features

- GNU AutoTools (autoconf) based install process
- Support for several scalar type combinations
 - {32-bit, 64-bit int} x {single, double, complex, double_complex}
- “Plug-in” extensibility
 - Very advanced users may customize library (at run-time)
 - New heuristics
 - Alternative data structures & code variants



Optimizations Available in the Initial Release

- Optimizations for SpMV (**bold** → **heuristics**)
 - **Register blocking (RB)**: up to **4x** over CSR
 - Variable block splitting: 2.1x over CSR, 1.8x over RB
 - Diagonals: 2x over CSR
 - Reordering to create dense structure + splitting: 2x over CSR
 - **Symmetry**: **2.8x** over CSR, 2.6x over RB
 - **Cache blocking**: **3x** over CSR
 - **Multiple vectors (SpMM)**: **7x** over CSR
 - And combinations...
- Sparse triangular solve
 - Hybrid sparse/dense data structure: **1.8x** over CSR
- Higher-level kernels
 - **$AA^T * x$** , **$A^T A * x$** : **4x** over CSR, 1.8x over RB
 - **$A^2 * x$** : 2x over CSR, 1.5x over RB



Current and Future Work

- Pre-release and docs available at bebop.cs.berkeley.edu/oski
 - Fortran wrappers in progress
 - Comments on interface welcome!
- Future work
 - PETSc integration
 - Port to additional architectures
 - Vectors
 - SMPs
 - Additional heuristics
 - Buttari, et al. (2005)



The End

(Extra slides follow)



Installation

- `./configure [options]` [detect system info]
 - `--with-blas={<lib>, no, yes}`
 - `--with-papi={<lib>, no, yes}`
 - `--with-index-type={int, long, <C-type>}`
 - `--with-value-type={single, double, complex, doublecomplex}`
 - ...
- `make` [build lib & run off-line benchmarks]
- `make install`
- `make check` [optional testing]



Implementation

- Uses preprocessor to generate different integer/value type combinations from single set of sources
- Matrix type modules
 - Each matrix type is its own dynamically loaded module
 - Parameterized by scalar type, e.g., CSR<int, double>
 - Types “registered” at run-time
 - Module interface includes kernels, conversion, ...
- Kernels
 - Dispatch based on matrix type
 - Each type implements SpMV + any subset of other 4 kernels
 - Default implementations if matrix type does not implement particular kernel
 - Self-profiling: time, number of calls



What Happens at Tuning Time?

- Available information
 - Hints about matrix structure
 - Workload hints from user (# of calls to each kernel w/ particular options)
 - Trace: Observed calls and execution time
 - Time to stream through matrix (at matrix creation time)
- Tuning procedure
 - Estimate a “tuning budget” from trace & workload hints
 - (fraction) * MAX(workload “time”, trace time)
 - WHILE (time left for tuning) & (not tuned) DO
 - Get and try a heuristic
- Currently does not re-tune



Heuristic Models

- Each in its own dynamically loadable module
- Module interface to a heuristic
 - `IsApplicable(<tuning info, e.g., matrix, trace, hints>);`
 - `<estimated time> = GetEstimatedCost(...);`
 - `<results> = Evaluate(...);`
 - `Transform(matrix, results);`

Requires High-Resolution Timers



- Inline assembly cycle counter readers for most platforms
 - Adapted from FFTW-3.0 (MIT license)
 - Includes x86-32, x86-64, IA-64, Sun, PowerPC, PA-RISC
 - Also wraps around PAPI if available (configure-time)

Documentation & Testing



- Doxygen for API
- Large test suite
 - ~30k line matrix multiply test program tries {precision} x {pattern} x {0,1-based inds} x {op(A)} x {x-orient} x {y-orient}



Optimizations in the Initial OSKI Release

- Fully automatic heuristics for
 - Sparse matrix-vector multiply
 - Register-level blocking
 - Register-level blocking + symmetry + multiple vectors
 - Cache-level blocking
 - Sparse triangular solve with register-level blocking and “switch-to-dense” optimization
 - Sparse $A^T A * x$ with register-level blocking
- User may select other optimizations manually
 - Diagonal storage optimizations, reordering, splitting; tiled matrix powers kernel ($A^k * x$)
 - All available in dynamic libraries
 - Accessible via high-level embedded script language