

Abstract Visualization of Runtime Memory Behavior

A.N.M. Imroz Choudhury and Paul Rosen

Scientific Computing and Imaging Institute
School of Computing
University of Utah

September 29, 2011

Why Visualize Memory Behavior?

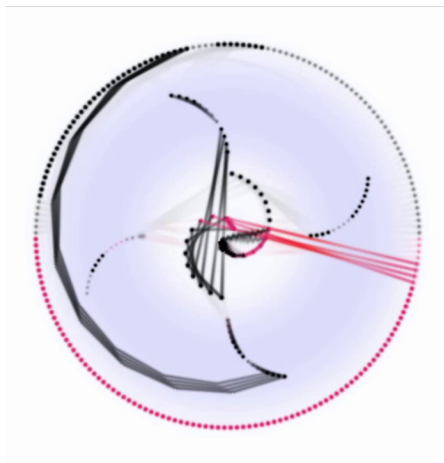
- Modern hardware + software → complex interactions
- e.g. caches: simple rules of operation, but complex behavior
- Sometimes performance is a design goal; cache performance is key to overall performance
 - Performance counters → coarse-grained information
 - Memory reference trace + cache simulation instead

Why Visualize Memory Behavior?

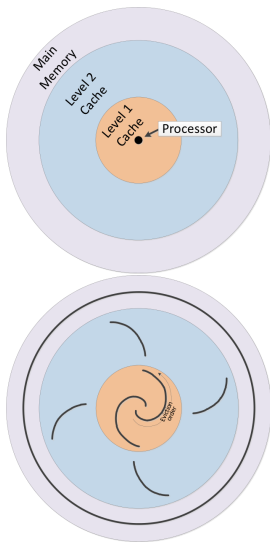
- Modern hardware + software → complex interactions
- e.g. caches: simple rules of operation, but complex behavior
- Sometimes performance is a design goal; cache performance is key to overall performance
 - Performance counters → coarse-grained information
 - Memory reference trace + cache simulation instead

Our approach: simulate cache, visualize what happens inside, highlight performance issues, **move towards understanding memory performance**

Example



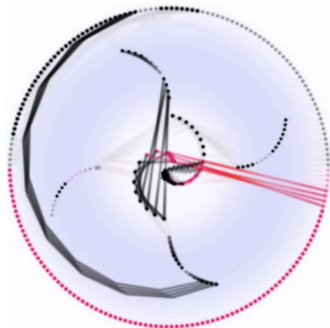
Overview of Approach



- Collect memory reference trace
 - Modify read/write instructions at run time
 - Trace includes source code information
- Cache simulation
 - Trace records (e.g. R 0x400500f) fed through simulator → hit or miss? eviction? etc.
- Visualization
 - Structural layout reflects cache architecture
 - Simulation results play out over layout

Visualization

Visual Channels and Frequencies

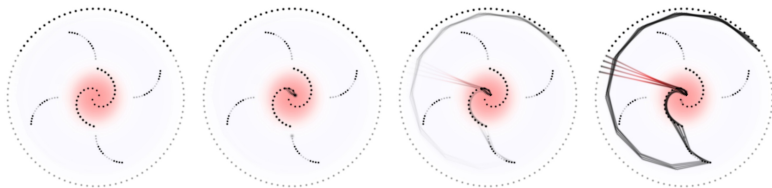


	Low frequency	High frequency
Color	home memory region	cache miss status
Motion	change in eviction order	change in cache level
Structure	associativity set	change in eviction order
Size	N/A	access

Visualization

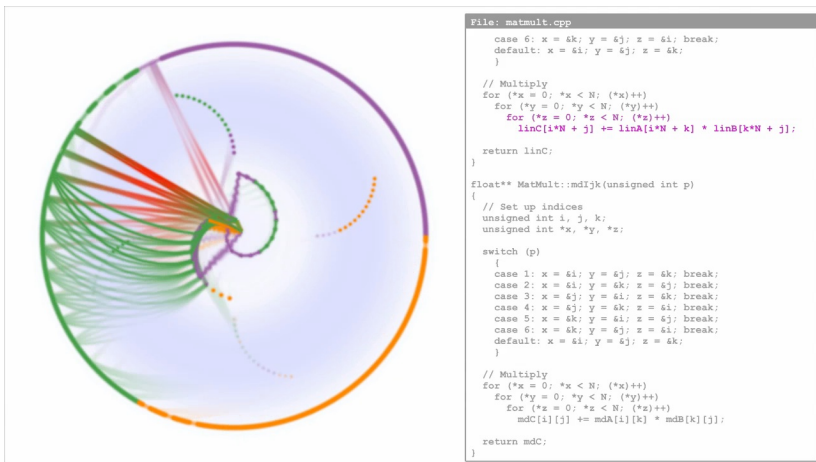
Time-Lapse Mode

- Very long traces, but can't just speed up playback
- Compress time → glyph *pathlines*
- Compute fewer such frames → time speeds up



Visualization

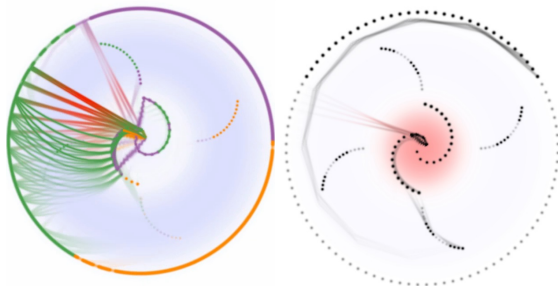
Time-Lapse Mode



Visualization

Summary View

- Expansive space behind data glyphs—display a summary statistic here
- “Cache temperature” summarizes recent cache performance, gives sense of history



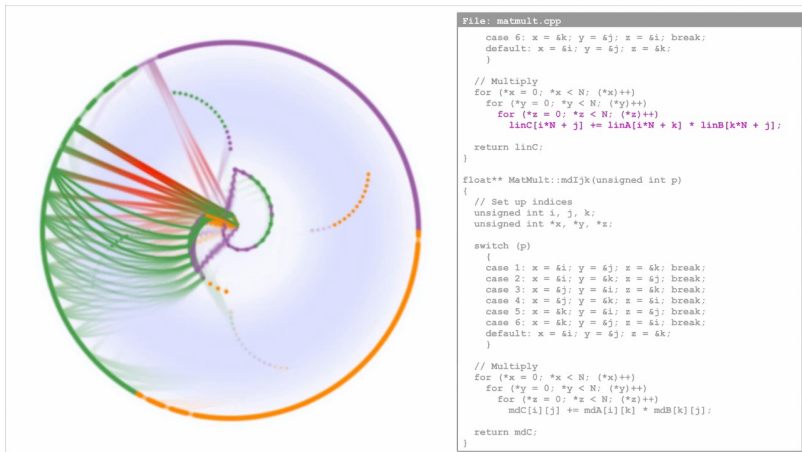
Standard Matrix Multiplication

- Left matrix: efficient access
- Right matrix: inefficient access

A _{1,1} A _{1,2} A _{1,3} A _{1,4} A _{1,5} A _{1,6} A _{1,7} A _{1,8}	X	B _{1,1} B _{1,2} B _{1,3} B _{1,4} B _{1,5} B _{1,6} B _{1,7} B _{1,8}	=	C _{1,1} C _{1,2} C _{1,3} C _{1,4} C _{1,5} C _{1,6} C _{1,7} C _{1,8}
A _{2,1} A _{2,2} A _{2,3} A _{2,4} A _{2,5} A _{2,6} A _{2,7} A _{2,8}		B _{2,1} B _{2,2} B _{2,3} B _{2,4} B _{2,5} B _{2,6} B _{2,7} B _{2,8}		C _{2,1} C _{2,2} C _{2,3} C _{2,4} C _{2,5} C _{2,6} C _{2,7} C _{2,8}
A _{3,1} A _{3,2} A _{3,3} A _{3,4} A _{3,5} A _{3,6} A _{3,7} A _{3,8}		B _{3,1} B _{3,2} B _{3,3} B _{3,4} B _{3,5} B _{3,6} B _{3,7} B _{3,8}		C _{3,1} C _{3,2} C _{3,3} C _{3,4} C _{3,5} C _{3,6} C _{3,7} C _{3,8}
A _{4,1} A _{4,2} A _{4,3} A _{4,4} A _{4,5} A _{4,6} A _{4,7} A _{4,8}		B _{4,1} B _{4,2} B _{4,3} B _{4,4} B _{4,5} B _{4,6} B _{4,7} B _{4,8}		C _{4,1} C _{4,2} C _{4,3} C _{4,4} C _{4,5} C _{4,6} C _{4,7} C _{4,8}
A _{5,1} A _{5,2} A _{5,3} A _{5,4} A _{5,5} A _{5,6} A _{5,7} A _{5,8}		B _{5,1} B _{5,2} B _{5,3} B _{5,4} B _{5,5} B _{5,6} B _{5,7} B _{5,8}		C _{5,1} C _{5,2} C _{5,3} C _{5,4} C _{5,5} C _{5,6} C _{5,7} C _{5,8}
A _{6,1} A _{6,2} A _{6,3} A _{6,4} A _{6,5} A _{6,6} A _{6,7} A _{6,8}		B _{6,1} B _{6,2} B _{6,3} B _{6,4} B _{6,5} B _{6,6} B _{6,7} B _{6,8}		C _{6,1} C _{6,2} C _{6,3} C _{6,4} C _{6,5} C _{6,6} C _{6,7} C _{6,8}
A _{7,1} A _{7,2} A _{7,3} A _{7,4} A _{7,5} A _{7,6} A _{7,7} A _{7,8}		B _{7,1} B _{7,2} B _{7,3} B _{7,4} B _{7,5} B _{7,6} B _{7,7} B _{7,8}		C _{7,1} C _{7,2} C _{7,3} C _{7,4} C _{7,5} C _{7,6} C _{7,7} C _{7,8}
A _{8,1} A _{8,2} A _{8,3} A _{8,4} A _{8,5} A _{8,6} A _{8,7} A _{8,8}		B _{8,1} B _{8,2} B _{8,3} B _{8,4} B _{8,5} B _{8,6} B _{8,7} B _{8,8}		C _{8,1} C _{8,2} C _{8,3} C _{8,4} C _{8,5} C _{8,6} C _{8,7} C _{8,8}

Standard Matrix Multiplication

$$A \times B = C$$



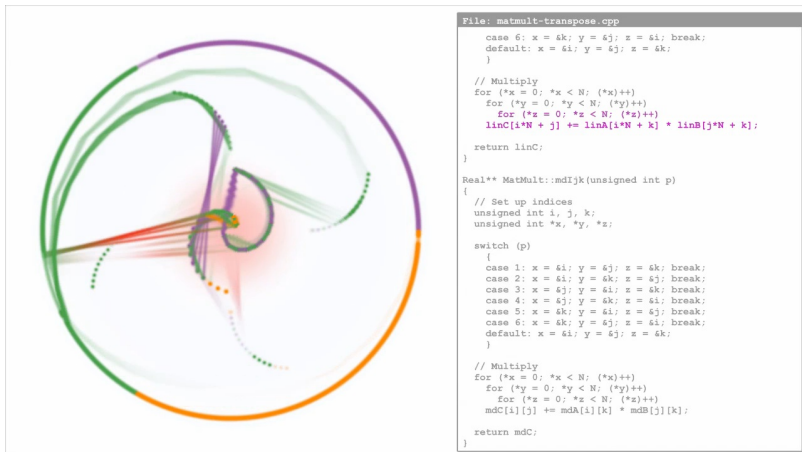
Transposed Matrix Multiplication

- Transpose right-hand matrix?
- Problem: loses generality

A _{1,1} A _{1,2} A _{1,3} A _{1,4} A _{1,5} A _{1,6} A _{1,7} A _{1,8}	X	B _{1,1} B _{1,2} B _{1,3} B _{1,4} B _{1,5} B _{1,6} B _{1,7} B _{1,8}	=	C _{1,1} C _{1,2} C _{1,3} C _{1,4} C _{1,5} C _{1,6} C _{1,7} C _{1,8}
A _{2,1} A _{2,2} A _{2,3} A _{2,4} A _{2,5} A _{2,6} A _{2,7} A _{2,8}		B _{2,1} B _{2,2} B _{2,3} B _{2,4} B _{2,5} B _{2,6} B _{2,7} B _{2,8}		C _{2,1} C _{2,2} C _{2,3} C _{2,4} C _{2,5} C _{2,6} C _{2,7} C _{2,8}
A _{3,1} A _{3,2} A _{3,3} A _{3,4} A _{3,5} A _{3,6} A _{3,7} A _{3,8}		B _{3,1} B _{3,2} B _{3,3} B _{3,4} B _{3,5} B _{3,6} B _{3,7} B _{3,8}		C _{3,1} C _{3,2} C _{3,3} C _{3,4} C _{3,5} C _{3,6} C _{3,7} C _{3,8}
A _{4,1} A _{4,2} A _{4,3} A _{4,4} A _{4,5} A _{4,6} A _{4,7} A _{4,8}		B _{4,1} B _{4,2} B _{4,3} B _{4,4} B _{4,5} B _{4,6} B _{4,7} B _{4,8}		C _{4,1} C _{4,2} C _{4,3} C _{4,4} C _{4,5} C _{4,6} C _{4,7} C _{4,8}
A _{5,1} A _{5,2} A _{5,3} A _{5,4} A _{5,5} A _{5,6} A _{5,7} A _{5,8}		B _{5,1} B _{5,2} B _{5,3} B _{5,4} B _{5,5} B _{5,6} B _{5,7} B _{5,8}		C _{5,1} C _{5,2} C _{5,3} C _{5,4} C _{5,5} C _{5,6} C _{5,7} C _{5,8}
A _{6,1} A _{6,2} A _{6,3} A _{6,4} A _{6,5} A _{6,6} A _{6,7} A _{6,8}		B _{6,1} B _{6,2} B _{6,3} B _{6,4} B _{6,5} B _{6,6} B _{6,7} B _{6,8}		C _{6,1} C _{6,2} C _{6,3} C _{6,4} C _{6,5} C _{6,6} C _{6,7} C _{6,8}
A _{7,1} A _{7,2} A _{7,3} A _{7,4} A _{7,5} A _{7,6} A _{7,7} A _{7,8}		B _{7,1} B _{7,2} B _{7,3} B _{7,4} B _{7,5} B _{7,6} B _{7,7} B _{7,8}		C _{7,1} C _{7,2} C _{7,3} C _{7,4} C _{7,5} C _{7,6} C _{7,7} C _{7,8}
A _{8,1} A _{8,2} A _{8,3} A _{8,4} A _{8,5} A _{8,6} A _{8,7} A _{8,8}		B _{8,1} B _{8,2} B _{8,3} B _{8,4} B _{8,5} B _{8,6} B _{8,7} B _{8,8}		C _{8,1} C _{8,2} C _{8,3} C _{8,4} C _{8,5} C _{8,6} C _{8,7} C _{8,8}

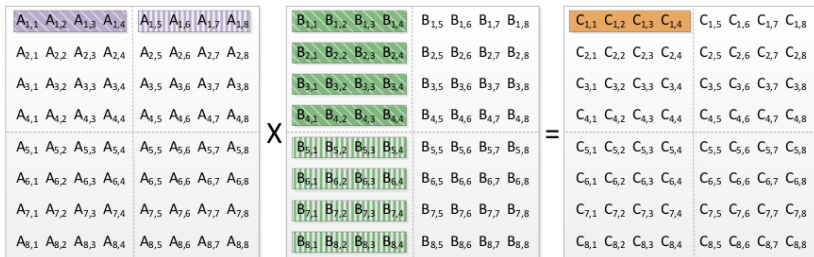
Transposed Matrix Multiplication

$$A \times B = C$$



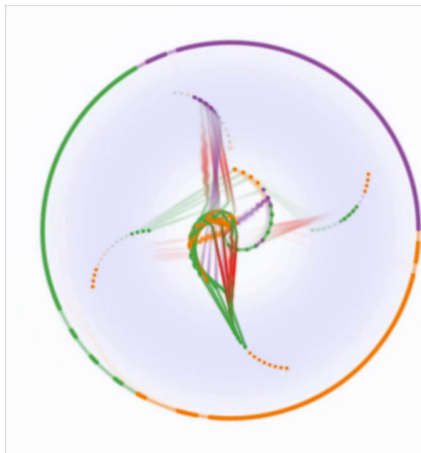
Blocked Matrix Multiplication

- Multiply small, cache-sized submatrices, accumulate into result
- Work on each cache line for longer → better data reuse



Blocked Matrix Multiplication

$$A \times B = C$$



File: matmult.cpp

```

for (unsigned int k0 = 0; k0 < N; k0 += b)
  for (unsigned int j0 = 0; j0 < N; j0 += b)
    for (unsigned int i = 0; i < N; i++)
      for (unsigned int k = k0; k < min(k0 + b, N); k++)
        {
          r = linA[i*N + k];
          for (unsigned int j = j0; j < min(j0 + b, N); j++)
            linC[i*N + j] += r*linB[k*N + j];
        }

return linC;
}

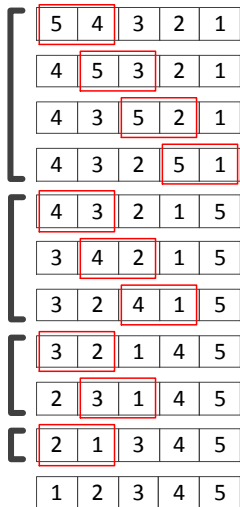
void MatMult::linRandM(float *A)
{
  for (unsigned int i = 0; i < N*N; i++)
    A[i] = ((float) rand()) / RAND_MAX;
}

void MatMult::mdRandM(float **A)
{
  for (unsigned int i = 0; i < N; i++)
    for (unsigned int j = 0; j < N; j++)
      A[i][j] = ((float) rand()) / RAND_MAX;
}

void MatMult::linPrintResults(float *A)
{
  printf("[\n");
  for (unsigned int i = 0; i < N*N; i++)
    {
      printf(" %f", A[i]);
      if ((i + 1) % N == 0)
        printf(";\n");
    }
}

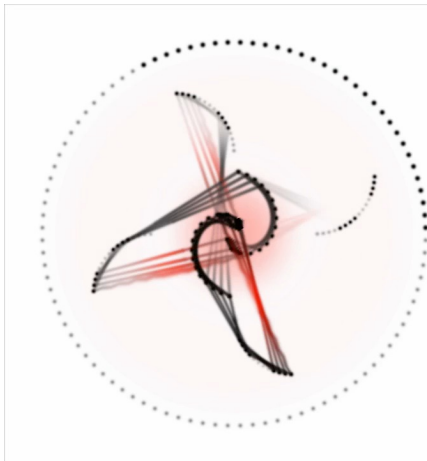
```

Bubble Sort



- Sorts list by repeated, shrinking sweeps that swap one element to correct place
- Known for slow algorithmic complexity, $O(n^2)$
- Interesting cache performance—dramatic increase when working set fits into cache

Bubble Sort



File: sort.cpp

```
int floatCompare(const void *a, const void *b){
    return *(const float *)a - *(const float *)b;
}

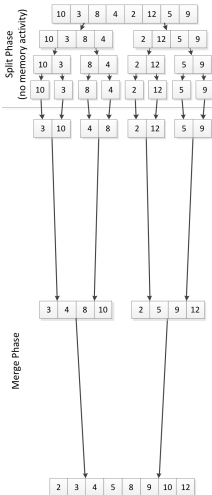
void bubblesort(std::vector<float>& v){
    for(unsigned end=v.size()-1; end >= 0; end--){
        bool swapped = false;
        for(unsigned i=0; i<end; i++){
            if(v[i] > v[i+1]){
                std::swap(v[i], v[i+1]);
                swapped = true;
            }
        }
        if(!swapped){
            break;
        }
    }
}

int main(int argc, char *argv[]){
    // Get command line arguments.
    std::string inputfile, sortalgo;
    unsigned random;
    bool print, printbefore;

    try{
        // Create command line parser.
        TCLAP::CmdLine cmd("Sort a file's worth of numbers acco
                                "input-file"
                                "Newline-sep"
                                false,
                                "",
                                "filename",
                                cmd);

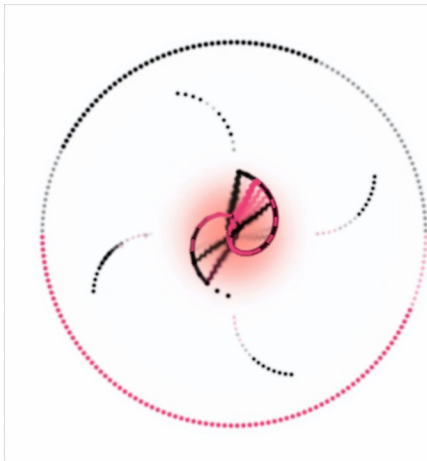
        // The file containing the numbers to sort.
        TCLAP::ValueArg<std::string> inputfileArg("i",
```

Merge Sort



- Better run time: $O(n \log n)$
- Recursive execution structure → more complex cache behavior

Merge Sort



File: mergesort.cpp

```

// Merge sort sublists
MergeSort( data, tmp, i0, i1 );
MergeSort( data, tmp, i1, i2 );

// Merge Lists
int _i0 = i0;
int _i1 = i1;
for(int i = i0; i < i2; i++){
    if( _i1 >= i2 ){
        tmp[i] = data[_i0++];
    }
    else if( _i0 >= i1 ){
        tmp[i] = data[_i1++];
    }
    else if( data[_i0] < data[_i1] ){
        tmp[i] = data[_i0++];
    }
    else if( data[_i1] < data[_i0] ){
        tmp[i] = data[_i1++];
    }
}

// Copy data from temporary location
for(int i = i0; i < i2; i++){
    data[i] = tmp[i];
}
}

int main(int argc, char ** argv){
    srand(0);
    const int size = 96;
    vector<double> data(size);
    vector<double> tmp(size);

    MTR::Registrar r;
    r.array(reinterpret_cast<MTR::addr_t>(&data[0]), d
    r.array(reinterpret_cast<MTR::addr_t>(&tmp[0]), tmp

```

Material Point Method (MPM)

- Particle-based simulation method for mechanical engineering problems
- Particles carry several attributes, stored in several arrays
- Tends to show hallmark “streaming” access pattern

Material Point Method (MPM)



Conclusions

- Memory reference trace → Cache simulation → **Structured visualization of cache internals**
- Visual design assigns meaning to several visual channels
- Resulting visuals and patterns convey memory behavior

Future Work

- Exploration of unused or underused visual channels, e.g. low-frequency motion
- Expansion of techniques to other systems where performance is important

Thank you!