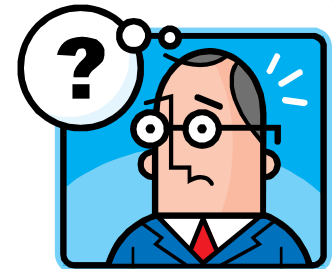




# Motivation

- Problem: HPC systems operate far below peak
  - Performance optimization complexity is growing
- Status: Most performance tools are hard to use
  - Require detailed performance and system expertise
  - HPC application developers are domain experts
- Result: HPC programmers do not use these tools
  - 75% of users haven't used performance tool on Ranger
  - Do not know how to apply information

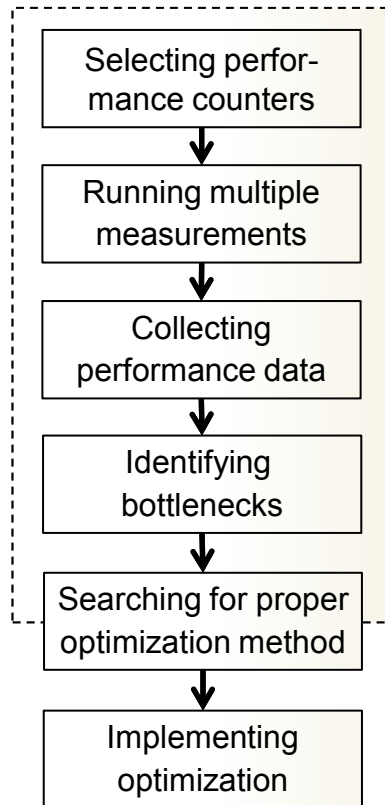


# Performance Counter Tool Workflow

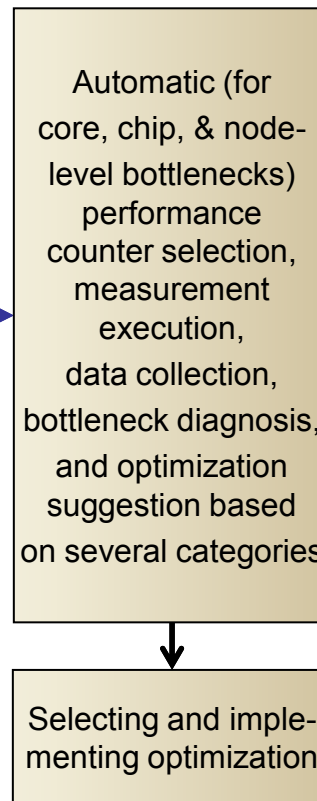
## Basic tools [mostly manual]

Provide no aid with:

- Counter selection
  - 100s of possibilities
  - cryptic descriptions
  - unclear what counted
- Result interpretation
  - Is there a problem?
  - What is the problem?
- Solution finding
  - How do I fix it?



## PerfExpert [mostly automated]



PerfExpert features:

- Automatic bottleneck detection & analysis
  - at core/chip/node level
- Recommends remedy
  - includes code examples & compiler switches
- Simple user interface
  - use provided job script
  - intuitive output



# Overview

- PerfExpert case studies on four Ranger codes
  - Mangll: mantle advection production code (C)
  - Homme: atmospheric acceptance benchmark (F95)
  - Libmesh: Navier-Stokes example code (C++)
  - Asset: astrophysics production code (F90)
- Step-by-step usage example
- Internal operation and performance metric
- Future work
- Summary





# Navier-Stokes Case Study


- Illustrates optimization-benefit tracking ability

```

total runtime in ex18.xml is 144.78 seconds
total runtime in ex18-cse.xml is 137.91 seconds

...

NavierSystem::element_time_derivative (runtimes are 33.29s and 25.24s)
-----
performance assessment   great.....good.....okay.....bad.....problematic
- overall                >>>>>>>>>>>>>>>>>>>222
upper bound by category
- data accesses          >>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>2
- instruction accesses   >>>>>>>
- floating-point instr   >>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>11111111111
- branch instructions    >
- data TLB               >
- instruction TLB        >
  
```







# Step-by-Step Usage Example

- Scenario
  - Developer's HPC code performs poorly
  - May know code section but not how to accelerate it
- Example: matrix-matrix multiplication
  - Coded inefficiently for illustration purposes
- PerfExpert reports **where** the slow code is, **why** it performs poorly, and suggests **how** to improve it



# Optimize Critical Code Section

- Loop nest around line 25

```
for (i = 0; i < n; i++)  
    for (j = 0; j < n; j++)  
        for (k = 0; k < n; k++)  
            c[i][j] += a[i][k] * b[k][j];
```

- Identified main bottleneck
  - Cause: memory accesses & data TLB
- Focus on data TLB problem first
  - No need to know what a data TLB is, just used as label to locate corresponding optimizations on web page

# Data TLB Optimization Suggestions

## 1) Improve the data locality

- a) use superpages (larger page sizes)

not yet enabled on all Ranger nodes

- b) change the order of loops

`loop i {...} loop j {...} → loop j {...} loop i {...}`

- c) employ loop blocking and interchange (change the order of the memory accesses)

`loop i {loop k {loop j {c[i][j] = c[i][j] + a[i][k] * b[k][j];}}} →  
loop k step s {loop j step s {loop i {for (kk = k; kk < k + s; kk++)  
{for (jj = j; jj < j + s; jj++) {c[i][jj] = c[i][jj] + a[i][kk] * b[kk][jj];}}}}}`

## 2) Reduce the data size

- a) use smaller types (e.g., float instead of double or short instead of int)

`double a[n]; → float a[n];` code example

use the "-fpack-struct" compiler flag compiler flag example

- b) allocate an array of elements instead of each element individually

`loop {... c = malloc(1); ...} →  
top = n; loop {if (top == n) {tmp = malloc(n); top = 0;} ... c = &tmp[top++]; ...}`

suggested remedy

# Eliminate Inapplicable Suggestions

## 1) Improve the data locality

~~a) use superpages (larger page sizes)~~

~~— not yet enabled on all Ranger nodes~~

b) change the order of loops

~~loop i {...} loop j {...} → loop j {...} loop i {...}~~

c) employ loop blocking and interchange (change the order of the memory accesses)

~~loop i {loop k {loop j {c[i][j] = c[i][j] + a[i][k] \* b[k][j];}}} →  
loop k step s {loop j step s {loop i {for (kk = k; kk < k + s; kk++)  
{for (jj = j; jj < j + s; jj++) {c[i][jj] = c[i][jj] + a[i][kk] \* b[kk][jj];}}}}}~~

## 2) Reduce the data size

~~a) use smaller types (e.g., float instead of double or short instead of int)~~

~~— double a[n]; → float a[n];~~

~~— use the "fpack-struct" compiler flag~~

~~b) allocate an array of elements instead of each element individually~~

~~— loop {... c = malloc(1); ...} →~~

~~top = n; loop {if (top == n) {tmp = malloc(n); top = 0;} ... c = &tmp[top++]; ...}~~

# Try Remaining Suggestions

- Start with suggestion 1b because it is simpler

## 1) Improve the data locality

- b) change the order of loops

`loop i {...} loop j {...} → loop j {...} loop i {...}`

- Exchange the  $j$  and  $k$  loops of the loop nest

```
for (i = 0; i < n; i++)
  for (k = 0; k < n; k++)
    for (j = 0; j < n; j++)
      c[i][j] += a[i][k] * b[k][j];
```

- Assess transformed code with PerfExpert



# Data Access Optimization Suggestions

## 1) Reduce the number of memory accesses

- a) move loop invariant memory accesses out of loop

```
loop i {a[i] = b[i] * c[j]} → temp = c[j]; loop i {a[i] = b[i] * temp;}
```

- b) ...

## 2) Improve the data locality

- a) componentize important loops by factoring them into their own subroutines

```
... loop i {...} ... loop j {...} ... → void li() {...}; void lj() {...}; ... li(); ... lj(); ...
```

- b) employ loop blocking and interchange (change the order of the memory accesses)

```
loop i {loop k {loop j {c[i][j] = c[i][j] + a[i][k] * b[k][j];}}} →  
loop k step s {loop j step s {loop i {for (kk = k; kk < k + s; kk++)  
{for (jj = j; jj < j + s; jj++) {c[i][jj] = c[i][jj] + a[i][kk] * b[kk][jj];}}}}}
```

- c) ...

## 3) Reduce the data size

...

we will pick this one as it was already suggested before



# Try Loop Blocking Suggestion

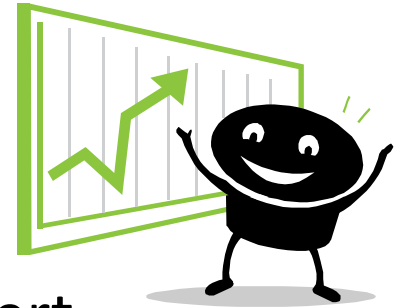
- Blocked loop code (blocking factor  $s = 70$ )

```
for (k = 0; k < n; k += s) {
    for (j = 0; j < n; j += s) {
        for (i = 0; i < n; i++) {
            for (kk = k; kk < k + s; kk++) {
                for (jj = j; jj < j + s; jj++) {
                    c[i][jj] += a[i][kk] * b[kk][jj];
                }
            }
        }
    }
}
```



# Usage Example Summary

- Performance is greatly improved
  - Optimization process guided by PerfExpert
  - Runtime dropped by 13x
- Memory access and data TLB problems fixed
  - PerfExpert correctly identified these bottlenecks
  - Suggested useful code optimizations
  - Helped verify the resolution of the problem



# Internal PerfExpert Operation

- Gather performance counter measurements
  - Multiple runs with HPCToolkit (PAPI & native counters)
  - Sampling-based results for procedures and loops
- Combine and check results
  - Check variability, runtime, consistency, and integrity
- Compute metrics and output assessment
  - Only for most important code sections
  - Correlate results from different runs

# PerfExpert's Performance Metric

- Local Cycles Per Instruction (LCPI)
  - Compute upper bounds on CPI contribution for various categories (e.g., branches, memory) and code sections
    - $(BR\_INS * BR\_lat + BR\_MSP * BR\_miss\_lat) / TOT\_INS$
    - $(L1\_DCA * L1\_dlat + L2\_DCA * L2\_lat + L2\_DCM * Mem\_lat) / TOT\_INS$
    - green = performance counter results, blue = system parameters
- Benefits
  - Highlights key aspects and hides misleading details
  - Relative metric (less susceptible to non-determinism)
  - Easily extensible (to refine or add more categories)

# Related Work

- Automatic bottleneck analysis and remediation
  - PERCS project at IBM Research
    - Less automation for bottleneck identification and analysis
    - Not open source
  - PERI Autotuning project
  - Parallel Performance Wizard
    - Event trace analysis, program instrumentation
- Analysis tools with automated diagnosis
- Projects that target multicore optimizations

# Future Work



- More case studies
  - Applications with various bottlenecks to harden tool
- Port to other systems: AMD, Intel, Power & GPU
  - Make PerfExpert available for general download
- Improve and expand capabilities
  - Finer-grained recommendations
  - Add data structure based analyses and optimizations
  - **Automatic implementation of solutions to common core, chip and node-level performance bottlenecks**

