
Understanding the Performance of Ultrascale Systems with Performance Counters

Jeffrey Vetter

Future Technologies Team Leader
Oak Ridge National Laboratory

Presented to

SIAM PP MS32

26 February 2004

San Francisco, CA, USA



OAK RIDGE NATIONAL LABORATORY

Highlights

Massively parallel computing is here

Traditional techniques for performance analysis of hardware counter activity do not scale very well

- Many levels of concurrency

Microprocessor performance counters provide rich information about application performance

- How do we interpret all that performance counter data?

We propose a new solution

- Apply multivariate statistical methods

Our experiments reveal evidence that these techniques allow

- Easier understanding of counter metrics
- Reduce data management problems



Ultrascale Platforms will Rely on Extreme Concurrency

Over the past 15 years, supercomputing has moved toward massively parallel systems

- Scalar processors
- Vector processors
- SoC

This trend will continue

- Red Storm
- BlueGene/L
- Cray X2
- Commodity clusters



**Cray X1/X2:
Leadership class
computer for science**



**IBM Power4:
8th in the world**



**IBM Power3:
DOE-SC's first
terascale system**



**Intel Paragon:
World's fastest computer**



With this Extreme Concurrency Comes the Challenge of Efficient Application Execution

Performance analysis will become increasingly important as processor counts grow

- Even small inefficiencies can be amplified

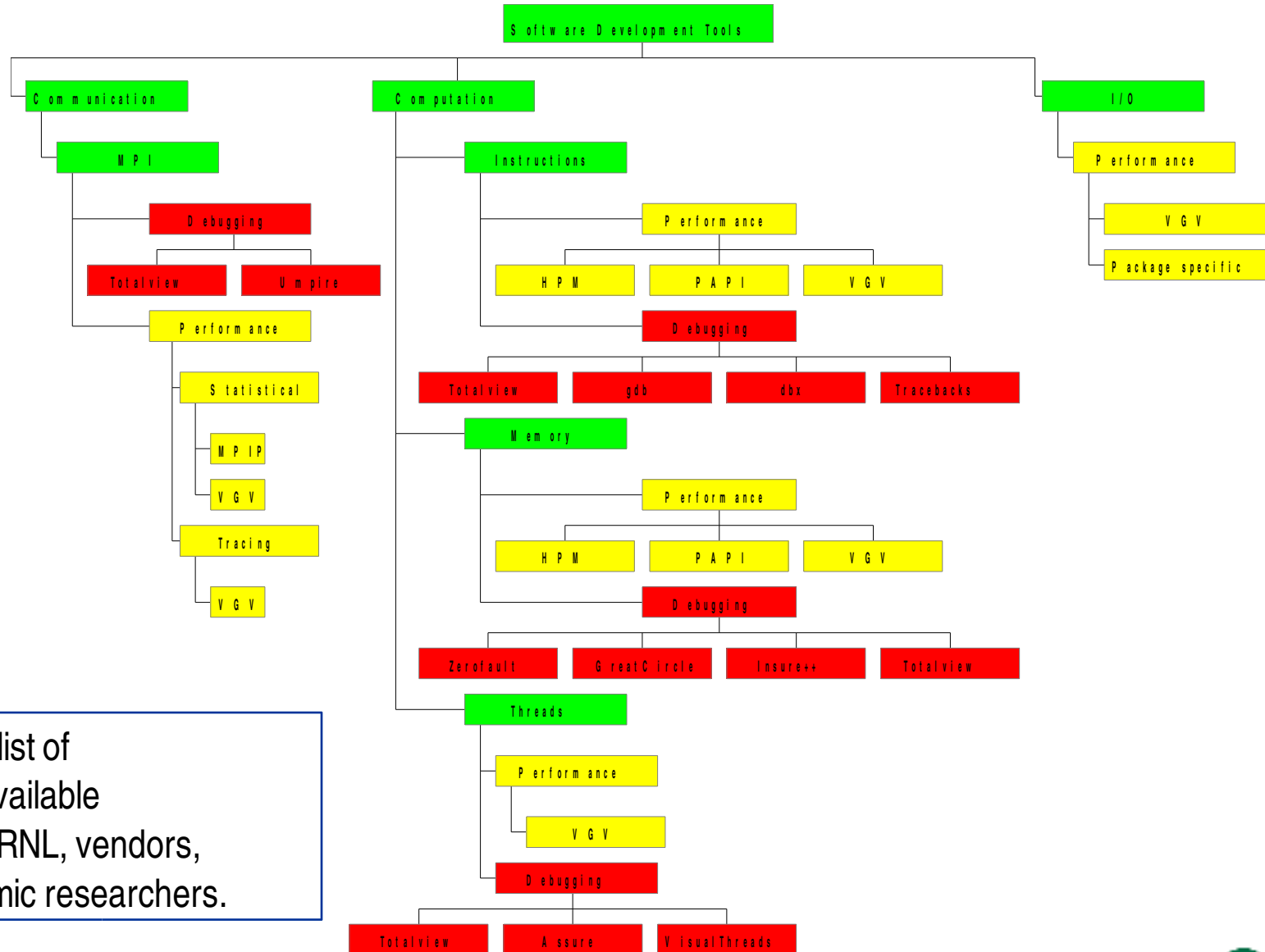
Understanding performance at this level of concurrency can be impractical and burdensome

Several areas of performance analysis technology don't scale well

Concurrency	O(100)	O(1,000)	O(10,000)	O(100,000)
Instrumentation	OK	OK	OK	OK
Instrumentation management	OK	Hurdle	Barrier	Barrier
Data management	OK	Hurdle	Barrier	Barrier
Data interpretation	Hurdle	Barrier	Barrier	Barrier



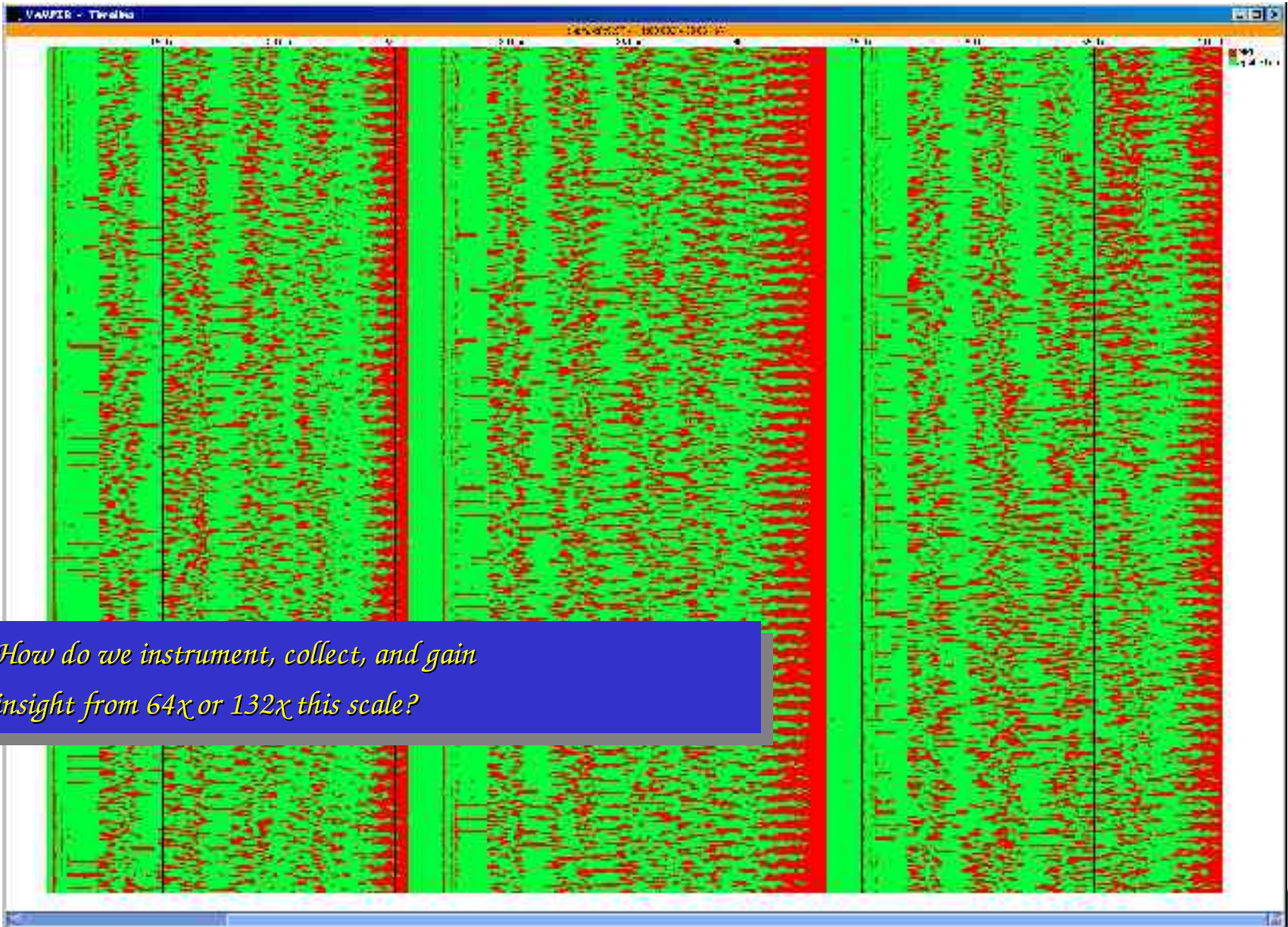
Assortment of Tools to Address Different Topics at Varying Levels of Detail



Partial list of tools available from ORNL, vendors, academic researchers.



Scalability Remains a Major Challenge on Today's Platforms: sPPM at 1024 tasks



How do we instrument, collect, and gain insight from 64x or 132x this scale?



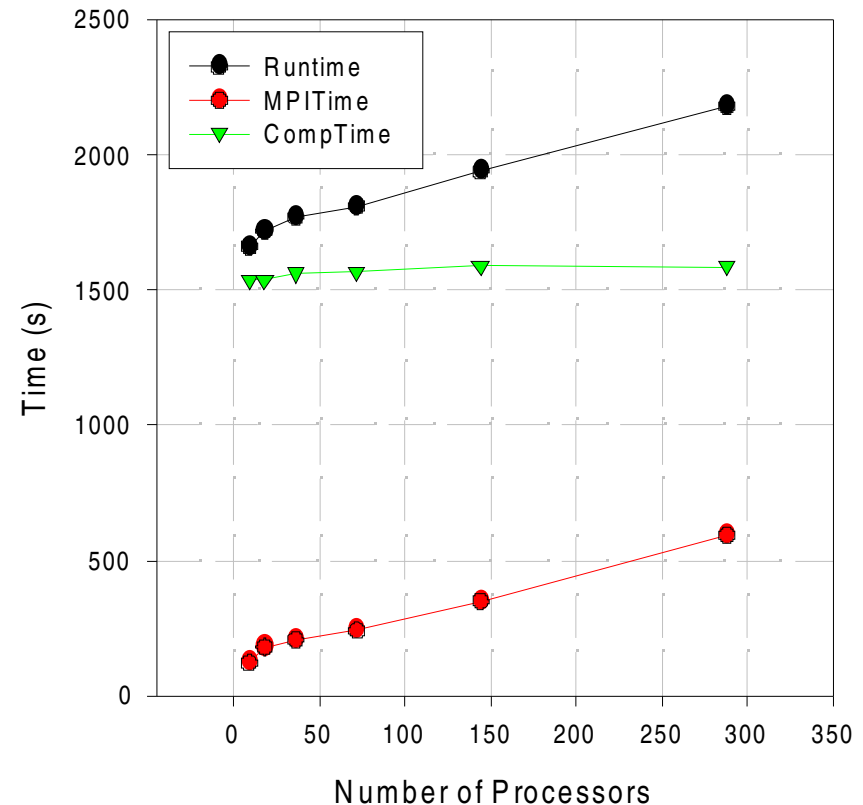
MPIP Provides Scalable MPI Performance Data

MPIP provides lightweight, scalable timing performance data

- Does NOT collect a trace

Tested up to 4096 processors

Gives statistical timing and payload information about MPI callsites



MPIP release 2.5 planned for March. See www.ccs.ornl.gov/~vetter

J.S. Vetter and M.O. McCracken, "Statistical Scalability Analysis of Communication Operations in Distributed Applications," Proc. ACM SIGPLAN Symp. on Principles and Practice of Parallel Programming (PPOPP), 2001.



Hardware Performance Counters are Invaluable for Understanding Computation Performance

Hardware counters measure empirical data that can help a user optimize and understand their code

- Low perturbation
- Real values (cf. static analysis, analytical modeling)

Despite some of their shortcomings, these counters are one of the few practical tools available

- Ill-defined
- Inconsistencies across platforms

Some systems have many counters and events

- POWER4 – over 200 events
- IA-64 – over 900 events

High dimensionality is compounded by extreme concurrency

Traditional Techniques (Sweep3d 256 mpi tasks)

Yield Huge Number of Data Points

G: Instrumentation ID	P: MPI Task	S: Instance								
1	1	1	83057605	77956513	22653498	1468948	7299392	37442673	21232357	1253921
1	1	2	82331144	75877135	22574422	8987587	7281693	36129321	21007529	1260843
1	2	1	81973603	77017507	22330703	1469595	7342599	37369569	20757228	1237291
1	2	2	81351386	75937600	22074563	9172755	7369905	35903746	20603110	1230463
2	1	1	83263293	75591985	24017955	1458386	7238297	37173268	20786530	1233867
2	1	2	82917911	74212484	23346289	8509892	7207491	35405216	20600234	1230893
2	2	1	84051067	76450556	24153969	1465589	7278521	37087982	21047398	1248838
2	2	2	83810819	75237837	23772760	8606055	7260832	35532887	20844958	1256918
			56	02	28		9	76	57	516

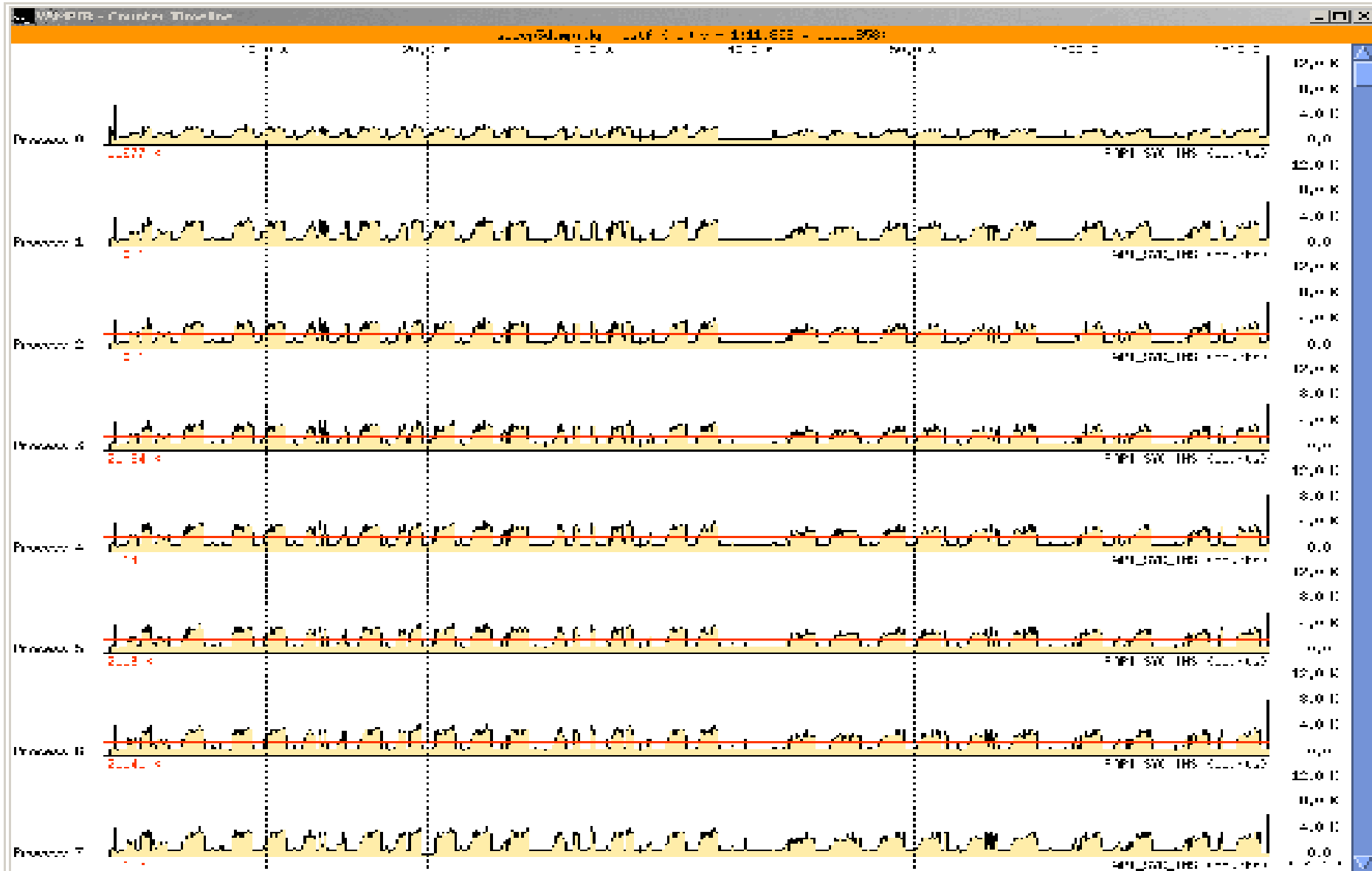
<num events> X <num sections> X <num instances> X <num tasks>

= 23 events X 5 sections X (1~12 instances) X 256

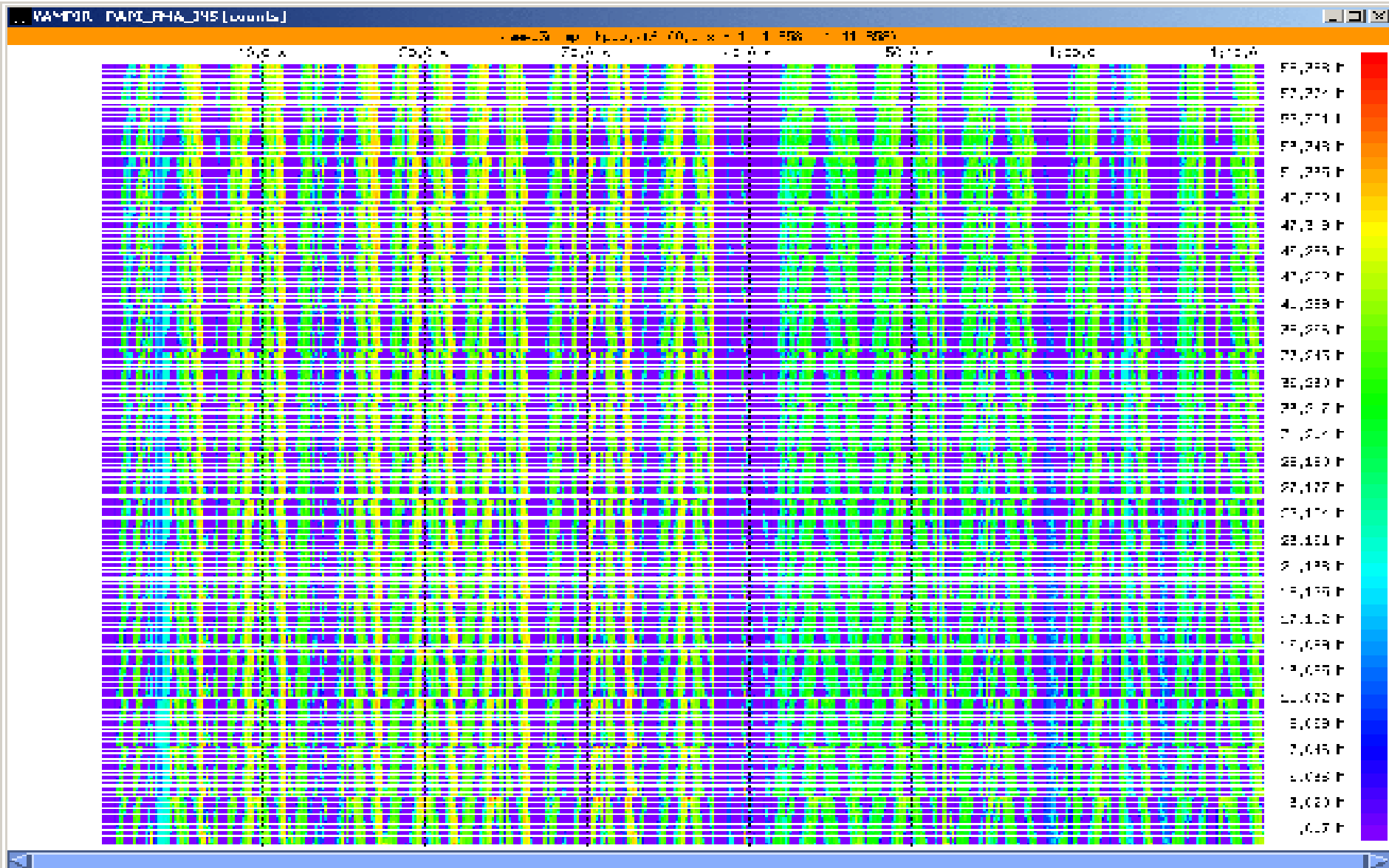
= **353,280 values !!**



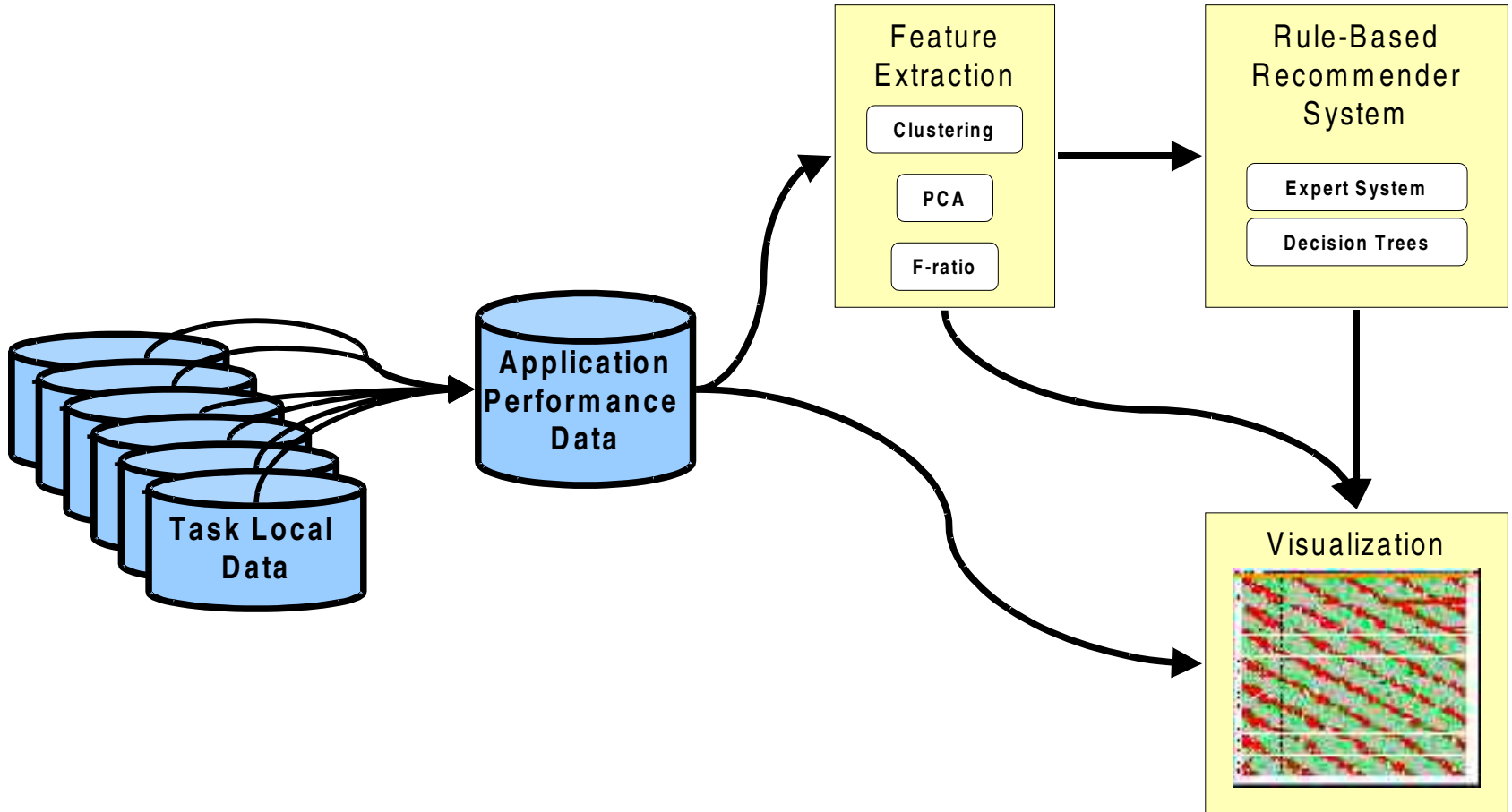
Try Mapping One or Two Dimensions to Timeline? Scalable?



Try Mapping One or Two Dimensions to Timeline w/ Color-coding? Scalable?

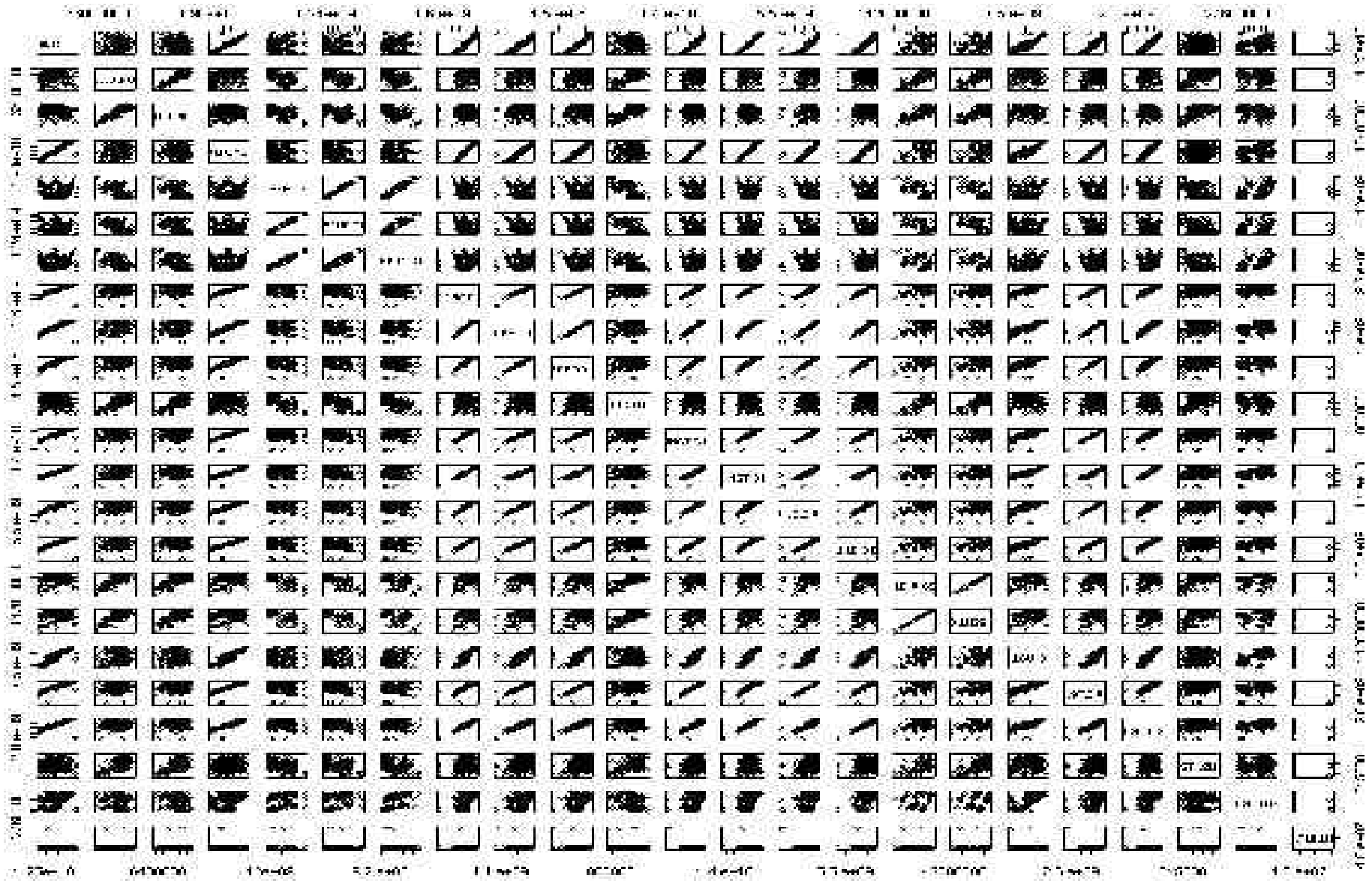


New Solution: Use Multivariate Statistics to Augment Analysis



Joint work with D.H. Ahn.

Sweep3D Scatterplot Matrix Illustrates Correlations between Counter Values



Many Facets of Multivariate Statistical Analysis are Useful

MSA provides a means to find relationships among variables (hardware metrics dimension) and individuals (i.e. mpi tasks or omp threads dimension)

Cluster Analysis groups together the processing elements (i.e. mpi tasks) that behave statistically similarly in hardware counter data space

Factor Analysis allows us to group together those metrics that explains the same underlying factor

Combining CA/FA and using their parameters (F-ratio, factor loadings) provides a variety of viewpoints on hardware counter dataset in more meaning ways

One Example of MSA: Cluster Analysis

Techniques

- Hierarchical
- K-means

Hierarchical method provides a general idea about cluster structure on dataset by building a hierarchical tree

Both methods work well to group together similar processing elements, based on dissimilarity matrix

Once datapoints are clustered, F-ratio (between-cluster variability/within-cluster variability) can identify important events that yield the particular cluster structure

Sweep3D Performance Experiment

Description

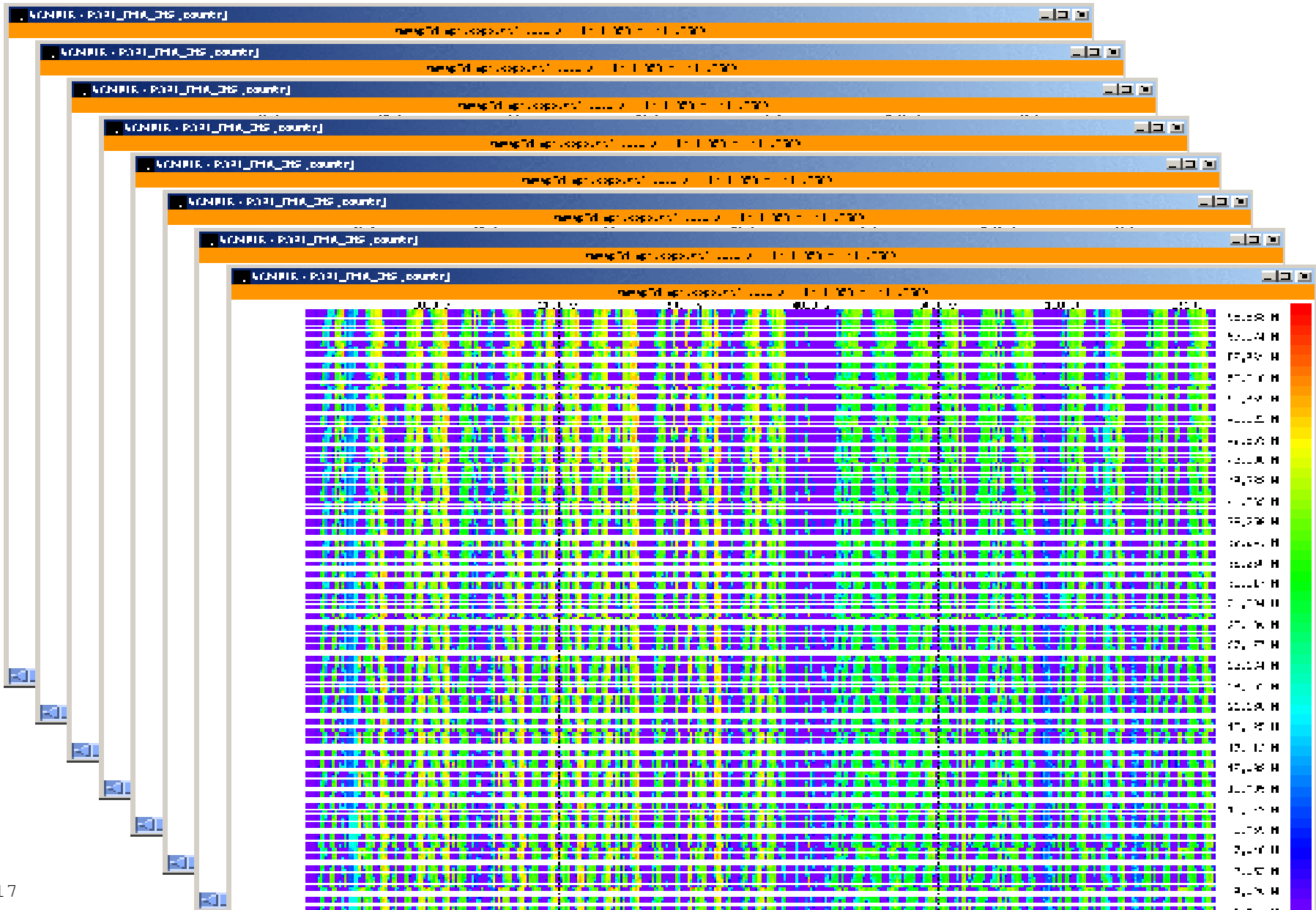
- A solver for the 3-D, time-independent, particle transport equation on an orthogonal mesh. Sweep3D uses wavefront algorithm for discrete ordinates deterministic particle transport simulation.
- Sweep3D exchanges messages between processors as wavefronts propagate diagonally across the 3-D space. Sweep section is the core of this algorithm. About 63% of its aggregate time spent in this section.

Experiment

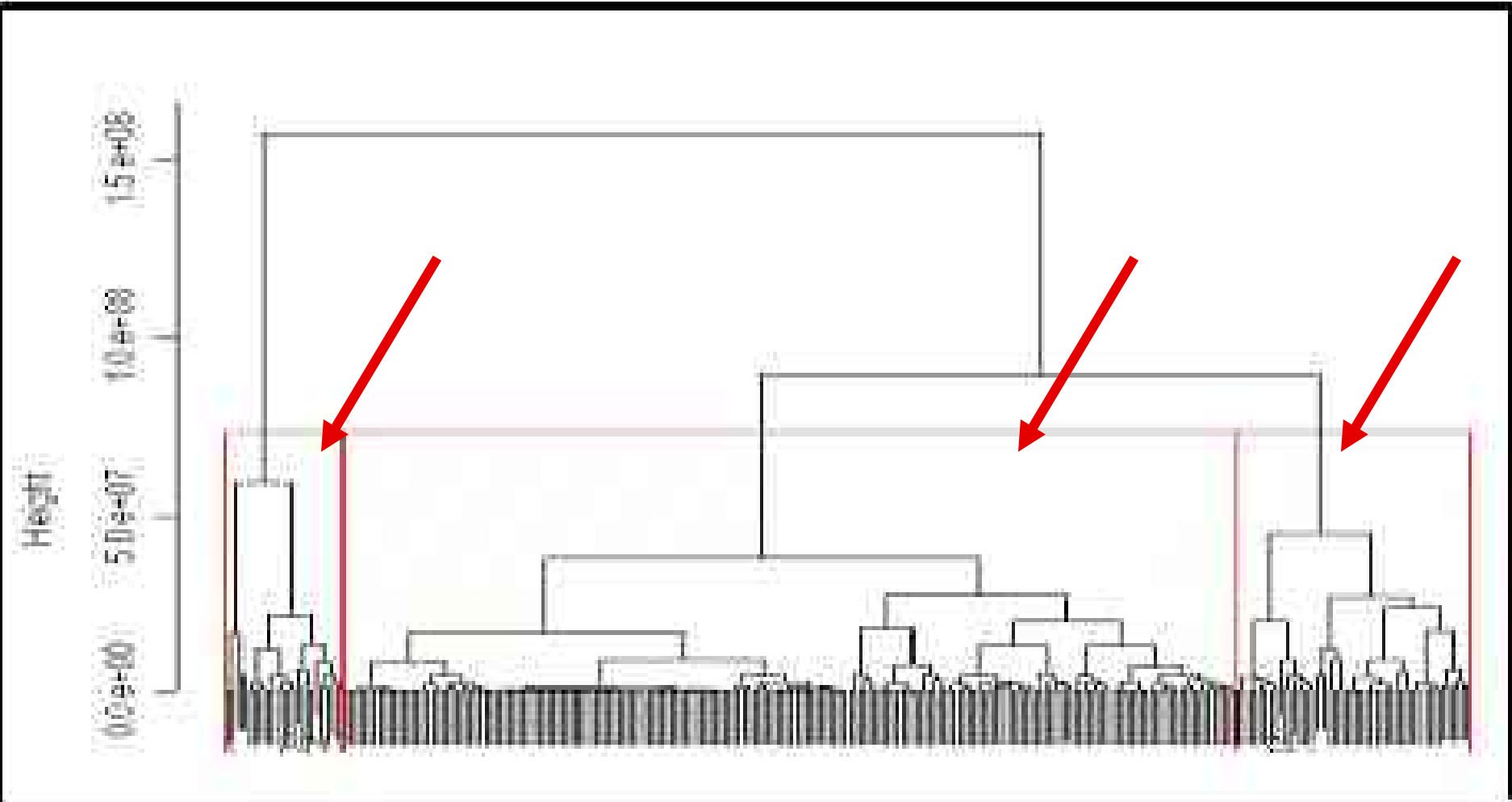
- Problem size : 512 X 512 X 150
- Num of tasks : 256 MPI tasks (16 X 16)

Capturing 8 counter events for one code section for each MPI task

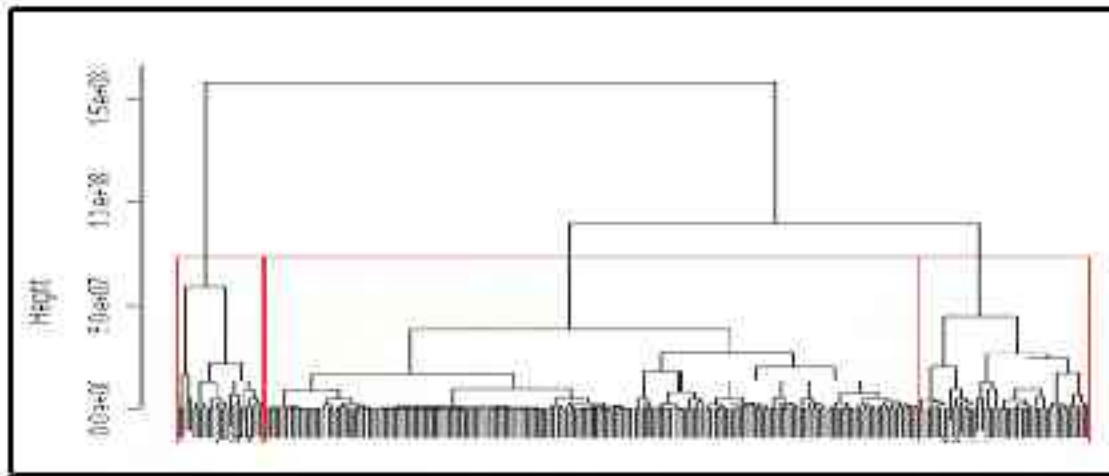
Curse of Dimensionality: Each Counter Generates a Separate Dimension of Values over Time



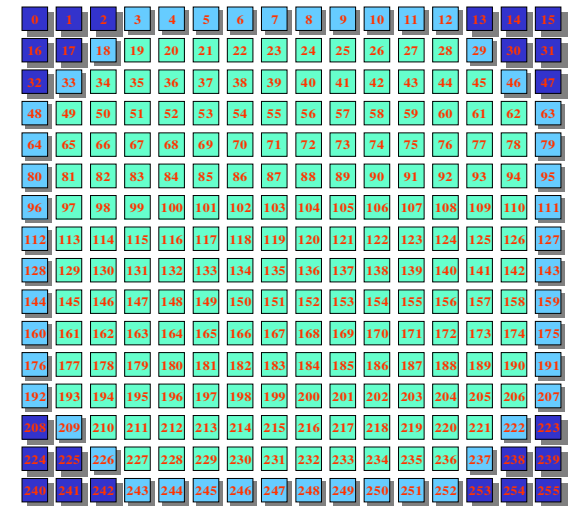
Dendrogram Illustrates the Clusters for the Counter Values across All 256 MPI tasks



Mapping the Clusters to the Application Topology Reveals Interesting Results



Task mapping



Performance optimizations to a representative translate to entire cluster

Easily identify metrics that separate performance clusters

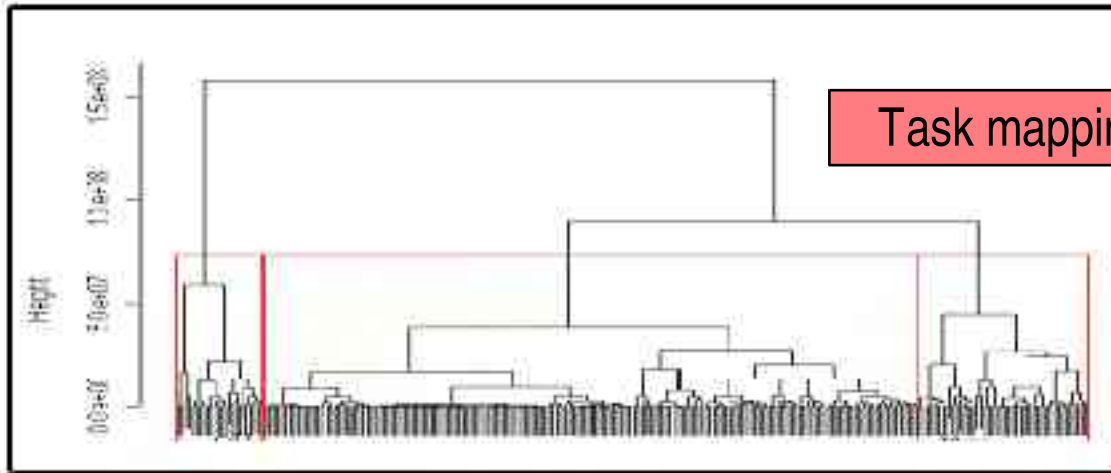
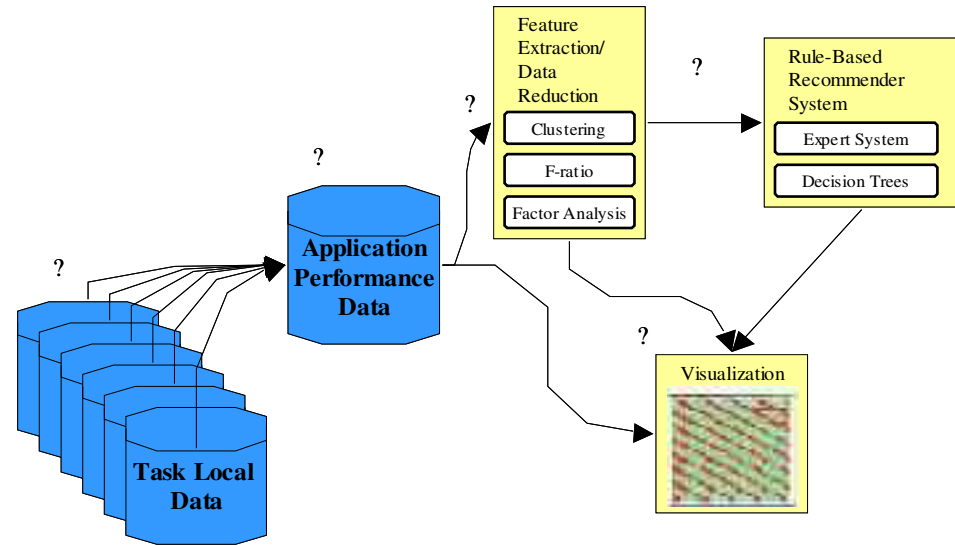
Collapse 256 tasks into 3!

Summary: Multivariate Statistical Analysis of Hardware Counter Data

Hardware counters produce huge amounts of data on large systems

Multivariate statistical techniques help distill important features

Clustering, Factor analysis, PCA



Task mapping

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
16	17	18	19	20	21	22	23	24	25	26	27	28	29	30
31	32	33	34	35	36	37	38	39	40	41	42	43	44	45
46	47	48	49	50	51	52	53	54	55	56	57	58	59	60
61	62	63	64	65	66	67	68	69	70	71	72	73	74	75
76	77	78	79	80	81	82	83	84	85	86	87	88	89	90
91	92	93	94	95	96	97	98	99	100	101	102	103	104	105
106	107	108	109	110	111	112	113	114	115	116	117	118	119	120
121	122	123	124	125	126	127	128	129	130	131	132	133	134	135
136	137	138	139	140	141	142	143	144	145	146	147	148	149	150
151	152	153	154	155	156	157	158	159	160	161	162	163	164	165
166	167	168	169	170	171	172	173	174	175	176	177	178	179	180
181	182	183	184	185	186	187	188	189	190	191	192	193	194	195
196	197	198	199	200	201	202	203	204	205	206	207	208	209	210
211	212	213	214	215	216	217	218	219	220	221	222	223	224	225
226	227	228	229	230	231	232	233	234	235	236	237	238	239	240
241	242	243	244	245	246	247	248	249	250	251	252	253	254	255

D.H. Ahn and J.S. Vetter, "Scalable Analysis Techniques for Microprocessor Performance Counter Metrics," Proc. SC 2002, 2002.

Other Performance Analysis Techniques to Address Scalability

Automatic Classification for MPI Trace Analysis

Use decision tree classification (a supervised learning technique) to classify application's messages automatically

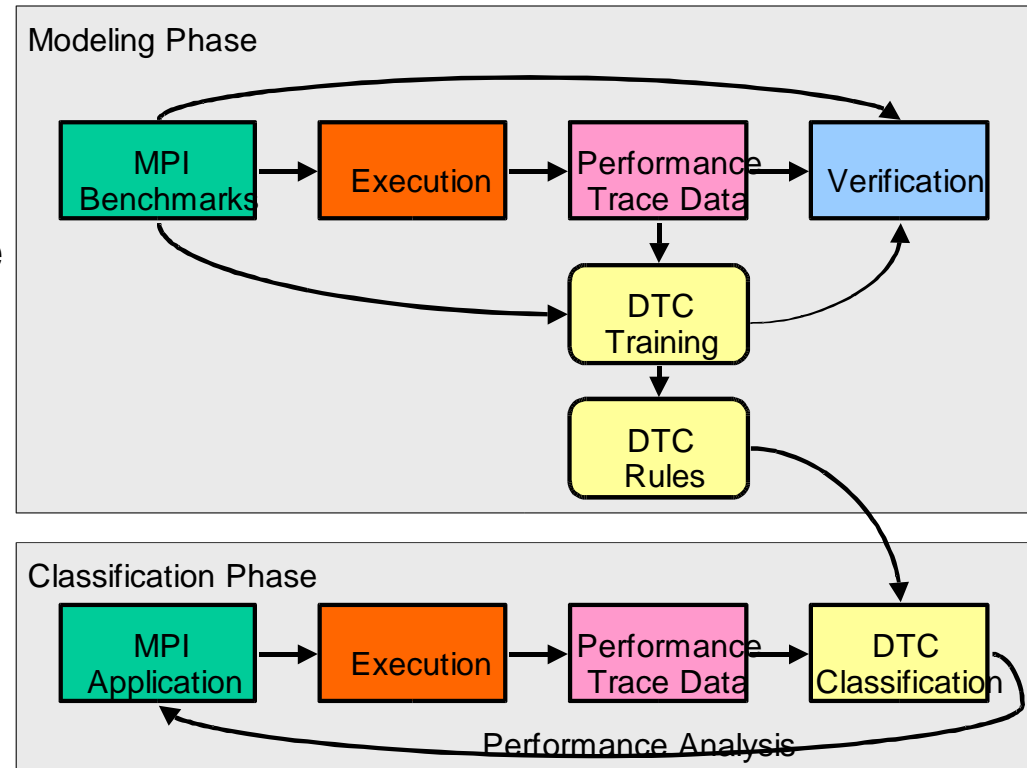
Compare an application's message operations to 'normal' communication for a particular MPI configuration

Modeling Phase (once)

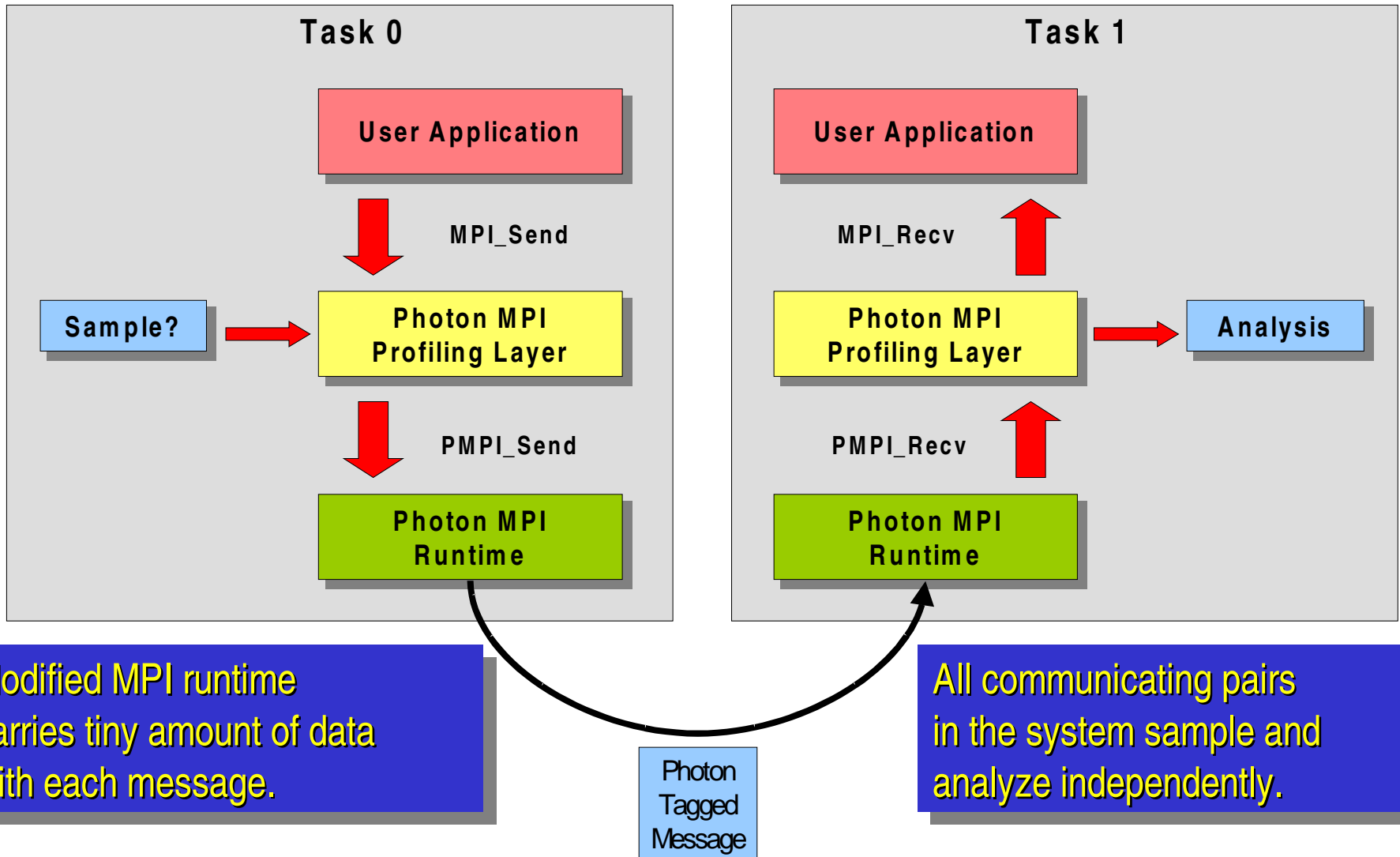
- Use benchmarks to generate decision tree
- Both efficient and inefficient

Classification Phase (many)

- Execute application
- Analyze application trace with classifier based on decision tree



Use Message Sampling and Runtime Analysis to Reduce Data and Perturbation



Conclusions

Massively parallel computing is here

Microprocessor performance counters provide rich information about application performance

Traditional techniques for performance analysis of hardware counter activity do not scale very well

- How do we interpret all the counter data?

We propose a new solution

- Apply multivariate statistical methods

Experiments reveal evidence that these techniques allow

- Easier understanding of counter metrics
- Reduce data management problems
- Improved productivity

More information at www.ccs.ornl.gov/~vetter



Bonus Slides