# PAPI 7.0 and beyond

15<sup>th</sup> Scalable Tools Workshop

<u>Anthony Danalis</u>, Heike Jagode, Giuseppe Congiu, Daniel Barry Tahoe City, CA June 18-22, 2023





- Library that provides a **consistent interface** (and methodology) for hardware performance counters, found across the system: i.e., CPUs, GPUs, on-/off-chip Memory, Interconnects, I/O, FS, Energy/Power.
- PAPI enables SW engineers to see, in near real time, the relation between **SW performance** and **HW events across the entire compute system**.



- Library that provides a **consistent interface** (and methodology) for hardware performance counters, found across the system: i.e., CPUs, GPUs, on-/off-chip Memory, Interconnects, I/O, FS, Energy/Power.
- PAPI enables SW engineers to see, in near real time, the relation between **SW performance** and **HW events across the entire compute system**.

#### **SUPPORTED ARCHITECTURES:**

- AMD <u>up to Zen4</u>
- ARM Cortex A8, A9, A15, ARM64, <u>ARM uncore-support</u>
- IBM Blue Gene Series
- IBM Power Series, **PCP for POWER9-nest**
- Intel KNC, KNL, KNM, up to **<u>Sapphire Rapids</u>**





1999 - 2009









- Library that provides a **consistent interface** (and methodology) for hardware performance counters, found across the system: i.e., CPUs, GPUs, on-/off-chip Memory, Interconnects, I/O, FS, Energy/Power.
- PAPI enables SW engineers to see, in near real time, the relation between **SW performance** and **HW events across the entire compute system**.

#### **SUPPORTED ARCHITECTURES:**

- AMD <u>up to Zen4</u>
- AMD <u>GPUs up to MI2x0</u>
- ARM Cortex A8, A9, A15, ARM64, <u>ARM uncore-support</u>
- CRAY: Gemini and Aries interconnects, power/energy
- IBM Blue Gene Series, Q: 5D-Torus, I/O system
- IBM Power Series, **PCP for POWER9-nest**
- Intel KNC, KNL, KNM, up to Sapphire Rapids
- Intel GPUs
- InfiniBand
- Lustre FS
- NVIDIA Tesla, Kepler, Maxwell, Pascal, Volta, <u>**Turing, Ampere**</u>: support for multiple GPUs
- NVIDIA: support for NVLink



.l.u.s.t.r.e.



INFINIBAND<sup>®</sup>

#### 2009 - 2018

Applications / 3<sup>rd</sup> Party Tools

Low-Level API High-Level API

PORTABLE LAYER

Developer API Developer API

- Library that provides a **consistent interface** (and methodology) for hardware performance counters, found across the system: i.e., CPUs, GPUs, on-/off-chip Memory, Interconnects, I/O, FS, Energy/Power.
- PAPI enables SW engineers to see, in near real time, the relation between **SW performance** and **HW events across the entire compute system**.

#### **SUPPORTED ARCHITECTURES:**

- AMD up to Zen4, power for Fam17h
- AMD GPUs up to MI2x0, power, temperature, fan
- ARM Cortex A8, A9, A15, ARM64, <u>ARM uncore-support</u>
- CRAY: Gemini and Aries interconnects, power/energy
- IBM Blue Gene Series, Q: 5D-Torus, I/O system, EMON power/energy
- IBM Power Series, PCP for POWER9-nest, power monitoring & capping on POWER9
- Intel KNC, KNL, KNM, up to **<u>Sapphire Rapids</u>**
- Intel RAPL (power/energy), **power capping**
- Intel GPUs
- InfiniBand
- Lustre FS
- NVIDIA Tesla, Kepler, Maxwell, Pascal, Volta, <u>**Turing, Ampere**</u>: support for multiple GPUs
- NVIDIA: support for NVLink
- NVIDIA NVML (power/energy); power capping
- Virtual Environments: VMware, KVM

Performance Counter Hardware
PAPI currently has >30 Components

2009 - 2018

Applications / 3<sup>rd</sup> Party Tools

Low-Level API High-Level API

PORTABLE LAYER

OS + Kernel Ext.

Developer API Developer API



·l·u·s·t·r·e·

ARM

THE SUPERCOMPU







# The PAPI repo has moved:

# https://github.com/icl-utk-edu/papi





# Sysdetect component



#### **System Detection Capabilities**



Applications / Tools PAPI High-level APIs PAPI Low-level APIs Components Substrate sysdetect rocm nic rocm сри 000 probe probe probe ROC Profiler ROCr сри Verbs 000 (hsa)



TENNESSEE

#### **System Detection Component**

#### \$ papi\_hardware\_avail

Device Summary Vendor DevCount		• In • Out
AuthenticAMD (1)		• PAPT GE
<pre>\-&gt; Status: Device Initialized</pre>		• In
(I)		• In
(-> Status. MP1 Not configured		• Out
Device Information		• PAPI_qe
Vendor	: AuthenticAMD	• In
Id	: 0	• In
Name	: AMD EPYC 7413 24-Core Processor	• Tn
CPUID	: Family/Model/Stepping 25/1/1	• Out
0x19/0x01/0x01		040
Sockets	: 2	
Numa regions	: 2	Wavefront
Cores per socket	: 24	SIMD per o
Cores per NUMA region	: 48	Max threa
SMT threads per core	: 2	
L1i Cache	: Size/LineSize/Lines/Assoc 32KB/64B/512/8	Max waves
L1d Cache	: Size/LineSize/Lines/Assoc 32KB/64B/512/8	Max shared
		Max workg
Numa Node 0 Memory	: 264023436 KB	Max workg
Numa Node 1 Memory	: 264199004 KB	Max workg
Numa Node 0 Threads	: 0 1 2 3 4 5 6 7 8 9 10 11	Max grid o
Numa Node 1 Threads	: 24 25 26 27 28 29 30 31 32	Max grid (
Vender		Max grid d
Name	<u>:                                    </u>	compute ca

#### Sysdetect low-level calls

- PAPI\_enum\_dev\_type
  - In modifier (e.g. PAPI\_DEV\_TYPE\_ENUM\_ROCM)
  - Out handle (opaque)
- PAPI\_get\_dev\_type\_attr
  - In handle
  - In attribute (e.g. PAPI\_DEV\_TYPE\_ATTR\_\_INT\_PAPI\_ID)
  - Out value
- PAPI\_get\_dev\_attr
  - In handle
  - In identifier (e.g. numa/gpu/... number)
  - In attribute (e.g. PAPI\_DEV\_ATTR\_\_ROCM\_<u>ULONG\_UID)</u>
  - Out value

device type attribute

Wavefront size	64
SIMD per compute unit	4
Max threads per workgroup	1024
Max waves per compute unit	40
Max shared memory per workgroup	65536
Max workgroup dim x	1024
Max workgroup dim y	1024
Max workgroup dim z	1024
Max grid dim x	429496729
Max grid dim y	429496729
Max grid dim z	429496729
Compute unit count	120
Compute capability	1.1





TENNESSEE

KNOXVILLE

THITTE

Ann

# **Counter Analysis Toolkit (CAT)**

## **Key Concepts**

#### • Goal:

Create a set of micro-benchmarks for illustrating details in hardware events and how they relate to the behavior of the micro-architecture

#### • Target audience:

- Performance conscious application developers
- PAPI developers working on new architectures (think preset events)
- Developers interested in validating hardware event counters



#### **CAT kernel example for Branch Events** 50% Mispredicted 0% Mispredicted do{ do{ iter\_count++; if ( iter\_count < (size/2) ){</pre> BUSY\_WORK(); global\_var2 += 2; BRNG(); if ( (result % 2) == 0 ){ BRNG(); iter count++; global var1+=2; 150% Taken }while(iter\_count<size);</pre> 100% Direct }while(iter\_count<size);</pre> do{ iter\_count++; do{ BUSY\_WORK(); BRNG(); BRNG(); global\_var2+=2; if ( (result % 2) == 0 ){ if ( iter\_so it < global\_var2 ){</pre> if( (global\_var1 % 2) != 0 ){ glob var1+=2; global var2++; goto zzz; global\_var1+=2; BRNG(): BUSY WORK(); zzz: iter\_count++; BRNG(); }while(iter\_count<size);</pre> }while(iter\_count<size);</pre> **¢iCl**

### **CAT kernel example for Branch Events**



### **CAT kernel example for Branch Events**



#### **Unique Responses Reveal Mapping to Preset Events**



#### Noise in Measurements: Relevant Event



BR\_INST\_RETIRED:COND (CAT Branching Benchmark, Intel Ice Lake SP)



#### Noise in Measurements: Irrelevant Event



#### **Distance Between Two Measurements**



#### **Distance Between Two Curves**





#### **RNMSE: Root Normalized Mean Square Error**

$$RNMSE(1,2) = \sqrt{\frac{1}{N} \sum_{i} \frac{(m_{1}^{i} - m_{2}^{i})^{2}}{\overline{m_{1}} \cdot \overline{m_{2}}}}$$

$$\overline{m_1} = \frac{1}{N} \sum_i m_1^i$$
$$\overline{m_2} = \frac{1}{N} \sum_i m_2^i$$



### **RNMSE: Root Normalized Mean Square Error**





## **Events with Low Noise Measurements (Zen4)**

CYCLES\_NO\_RETIRE:EMPTY OPS\_SOURCE\_DISPATCHED\_FROM\_DECODER:DECODER RETIRED\_BRANCH\_INSTRUCTIONS RETIRED\_BRANCH\_INSTRUCTIONS\_MISPREDICTED RETIRED\_BRANCH\_MISPREDICTED\_DIRECTION\_MISMATCH RETIRED\_CONDITIONAL\_BRANCH\_INSTRUCTIONS RETIRED\_FUSED\_INSTRUCTIONS RETIRED\_INSTRUCTIONS RETIRED\_OPS RETIRED\_TAKEN\_BRANCH\_INSTRUCTIONS\_MISPREDICTED RETIRED\_TAKEN\_BRANCH\_INSTRUCTIONS\_MISPREDICTED RETIRED\_UNCONDITIONAL\_BRANCH\_INSTRUCTIONS

perf::BRANCH-INSTRUCTIONS perf::BRANCH-LOAD-MISSES perf::BRANCH-LOADS perf::BRANCH-MISSES perf::BRANCHES perf::INSTRUCTIONS perf::L1-ICACHE-LOADS perf::PERF COUNT HW BRANCH INSTRUCTIONS perf::PERF COUNT HW BRANCH MISSES perf::PERF\_COUNT\_HW\_CACHE\_BPU:ACCESS perf::PERF COUNT HW CACHE BPU:MISS perf::PERF\_COUNT\_HW\_CACHE\_BPU:READ perf::PERF\_COUNT\_HW\_CACHE\_L1I:ACCESS perf::PERF\_COUNT\_HW\_INSTRUCTIONS





THE UNIVERSITY OF TENNESSEE KNOXVILLE

THUTTE

Sin

# **Software Defined Events (SDE)**

### **PAPI Software Defined Events (SDEs)**

Support for Events that originate in Software Layers

SDEs enable **software** layers to export **arbitrary information** as if it came from hardware counters

Arguments passed to functions, residuals, tasks stolen, hash-table collisions, messages sent, memory consumption, size of internal data structures, ...



#### Pull mode: Low overhead (down to zero)





#### Pull mode: Low overhead (down to zero)







## Pull mode: Low overhead (down to zero)

The application reads whenever it deems necessary.



papi\_sde\_register\_counter(&x,...)



## **Recorders (Multi-Value SDEs)**

The library keeps track of a series of events.



#### Set data-structure





### Hash table implementation inspired by FB's F14



DIRECT STORAGE in 3-slot F14-like chunks

https://engineering.fb.com/2019/04/25/developer-tools/f14/



- Add 85 elements in hash table with capacity 5110:
- $P(collision) = 1 (5109/5110)^{(85*84/2)} = 0.502$



- Add 85 elements in hash table with capacity 5110:
- $P(collision) = 1 (5109/5110)^{(85*84/2)} = 0.502$
- HT: 365 buckets with capacity of 14, how many insertions to get 15 collisions?





- Add 85 elements in hash table with capacity 5110:
- $P(collision) = 1 (5109/5110)^{(85*84/2)} = 0.502$
- HT: 365 buckets with capacity of 14, how many insertions to get 15 collisions?
- N = 2263 for P>50%



- Add 85 elements in hash table with capacity 5110:
- $P(collision) = 1 (5109/5110)^{(85*84/2)} = 0.502$
- HT: 365 buckets with capacity of 14, how many insertions to get 15 collisions?
- N = 2263 for P>50%



#### Actual bucket (sde\_hash\_bucket\_t)

	ev	ev	ev	ey	count											
	bit k	size														
0.00	64	64	64	64	64	64	64	64	64	64	64	64	64	64	ptr 🔨	ptr



#### Actual bucket (sde\_hash\_bucket\_t)

#### L1 cache line

100	ey	ev	ev	ν	ey	count										
	bit k	hit k	bit k	size												
000	64	64	64	64	64	64	64	64	64	64	64	64	64	64	ptr 🔨	ptr



#### Actual bucket (sde\_hash\_bucket\_t)

#### L1 cache line pair

ſ	ey	ev	Ň	λ i	d d	ey	ey	count									
	bit k	hit k	1 1 1 1 1 1		bit k	bit k	size										
	64	64	64	64	64	64	64	64	64	64	64	5 3	5	64	64	ptr 🔨	ptr



## The proof is in the pudding.







THE UNIVERSITY OF TENNESSEE KNOXVILLE

THUTT

SIL

# **Event Description Ambiguities**

## What does this event really measure?

#### LONGEST\_LAT\_CACHE:REFERENCES

Core-originated cacheable requests that refer to L3 (Except hardware prefetches to the L3).



## What does this event really measure?

#### LONGEST\_LAT\_CACHE:REFERENCES

Core-originated cacheable requests that refer to L3 (Except hardware prefetches to the L3).

Is "Except hardware prefetches" the same as "Demand References"?



#### **Demand References on Cascade Lake**





#### LLC References on Cascade Lake





#### **Summary**

- PAPI 7.0.1 is out (https://github.com/icl-utk-edu/papi).
- Support for GPU counters/metrics across vendors.
- Support for power management on CPUs & GPUs.
- API & utility for detecting available hardware.
- Software Defined Events as a standalone library (libsde.so).
- Counter Analysis Toolkit provides hardware insights.
- Noise is not due to hardware counters, but system operations.







- IIIIIIIIIIIIII

Summer of the local division of the local di

# **Backup Slides**

#### **Simple Buffer Traversal (C code)**

```
int main(int argc, char **argv){
    int M, N;
#define BUFFER_SIZE 256
    int buffer[BUFFER_SIZE];
```

```
M = 1024;
N = 1024;
for(int i=0; i<BUFFER SIZE; i++){</pre>
    buffer[i] = (argc*1222)/(1223+i);
int index = buffer[0];
for(int i=0; i<M; i++){</pre>
    for(int j=0; j<N; j++){</pre>
        index = buffer[index%BUFFER_SIZE];
printf("index=%d\n",index);
return 0;
```



```
Simple Buffer Traversal (C code)
int main(int argc, char **argv){
                                                     int main(int argc, char **argv){
    int M, N;
                                                         int M, N;
#define BUFFER SIZE 256
                                                     #define BUFFER SIZE 256
    int buffer[BUFFER_SIZE];
                                                         int buffer[BUFFER SIZE];
    M = 1024:
                                                         M = 1024;
    N = 1024;
                                                         N = 1024;
    for(int i=0; i<BUFFER SIZE; i++){</pre>
                                                         for(int i=0; i<BUFFER SIZE; i++){</pre>
        buffer[i] = (argc*1222)/(1223+i);
                                                             buffer[i] = (argc*1222)/(1223+i);
    int index = buffer[0];
                                                           int index = buffer[0];
    for(int i=0; i<M; i++){</pre>
                                                         uintptr_t index = buffer[0];
        for(int j=0; j<N; j++){</pre>
                                                         for(int i=0; i<M; i++){</pre>
            index = buffer[index%BUFFER SIZE];
                                                             for(int j=0; j<N; j++){</pre>
                                                                 index = buffer[index%BUFFER SIZE];
    printf("index=%d\n",index);
                                                         printf("index=%d\n",index);
    return 0;
                                                         return 0;
```



#### **Simple Buffer Traversal (asm code)**

int index = Du	ITTEF[0];		
000000000010cc	<+76>:	MOV	(%rsp),%edx
000000000010cf	<+79>:	mov	\$0x400,% <mark>esi</mark>
000000000010d4	< <b>+</b> 84>:	nopl	0x0(% <b>rax</b> )
<pre>for(int i=0; i</pre>	_ <m; i++){<="" th=""><th></th><th></th></m;>		
)000000000010fa	<+122>:	sub	\$0x1,% <b>esi</b>
000000000010fd	<+125>:	jne	0x10d8 <b><main< b="">+88&gt;</main<></b>
for(int_j=	:0; j <n; j+<="" th=""><th>++){</th><th></th></n;>	++){	
)000000000010f5	<+117>:	sub	\$0x1,%ecx
)000000000010f8	<+120>:	jne	0x10e0 <b><main< b="">+96&gt;</main<></b>
ındex	= buffer[1	Index%Bl	JFFER_SIZE];
000000000010e0	<+96>:	MOV	%edx,%eax
)000000000010e2	<+98>:	sar	\$0x1f,%eax
000000000010e5	<+101>:	shr	\$0×18,%eax
000000000010e8	< <b>+1</b> 04>:	add	%eax,%edx
000000000010ea	<+106>:	movzbl	%dl,%edx
00000000010ed	<+109>:	sub	%eax,%edx
00000000010ef	<+111>:	movslq	%edx,%rdx
)00000000010f2	<+114>:	MOV	(%rsp,%rdx,4),%edx



#### Simple Buffer Traversal (asm code)

int index = bi 0000000000010cc 000000000010cf 0000000000	<pre> iffer[0]; &lt;+76&gt;: &lt;+79&gt;: &lt;+84&gt;: </pre>	mov mov nopl	(% <mark>rsp),%edx</mark> \$0x400,% <mark>esi</mark> 0x0(%rax)
<b>for(int i=0;</b> i	<pre><m; <+122="" i++){="">: &lt;+125&gt;; </m;></pre>	sub	\$0x1,%esi
<b>for(int j</b> =	<+1232. =0; j <n; j+<br="">&lt;+117&gt;:</n;>	++){ sub	\$0x1.%ecx
00000000010f8	<+120> <b>:</b>	jne	0x10e0 <main+96></main+96>
index	= buffer[i	index%Bl	JFFER_SIZE];
000000000010e0	<+96>:	MOV	%edx,%eax
000000000010e2	<+98>:	sar	\$0x1f,% <b>eax</b>
000000000010e5	<+101>:	shr	\$0x18,% <b>eax</b>
00000000010e8	<+104>:	add	%eax,%edx
00000000010ea	<+106>:	movzbl	%dl,%edx
00000000010ed	<+109>:	sub	%eax,%edx
000000000010ef	<+111>:	movslq	%edx,%rdx
)00000000010f2	<+114>:	MOV	(% <b>rsp</b> ,% <b>rd</b> x,4),% <b>ed</b> x

uintptr\_t index = buffer[0]; )00000000010cc <+76>: movslq (%rsp),%rdx 000000000010d0 <+80>: \$0x400,%ecx MOV )00000000010d5 <+85>: nopl (%**rax**)

for(int i=0; i<M; i++){</pre> 000000000010ec <+108>: sub )000000000010ef <+111>: jne

\$0x1,%ecx 0x10d8 <main+88>

for(int j=0; j<N; j++){</pre> )00000000010e7 <+103>: sub )00000000010ea <+106>: jne

\$0x1,%eax 0x10e0 <main+96>

index = buffer[index%BUFFER\_SIZE]; )00000000010e0 <+96>: movzbl %dl,%edx

)00000000010e3 <+99>: movslq (%rsp,%rdx,4),%rdx 🥌

