



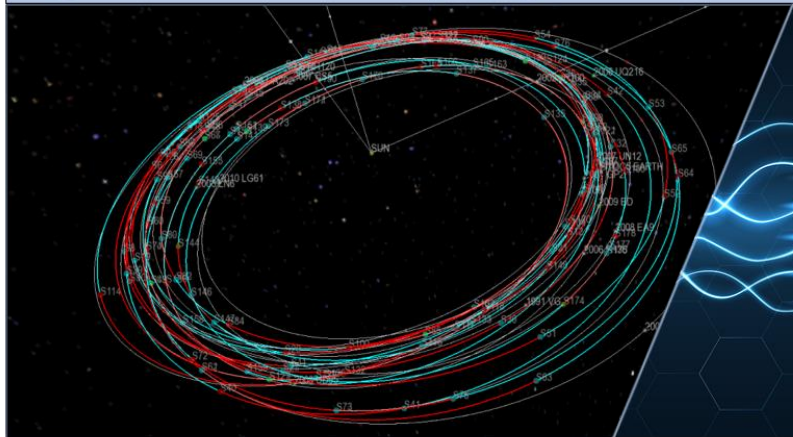
WILLIAM & MARY

CHARTERED 1693

Fault Site Pruning for Practical Reliability Analysis of GPGPU Applications

Bin Nie, **Lishan Yang**, Adwait Jog, and Evgenia Smirni
College of William & Mary

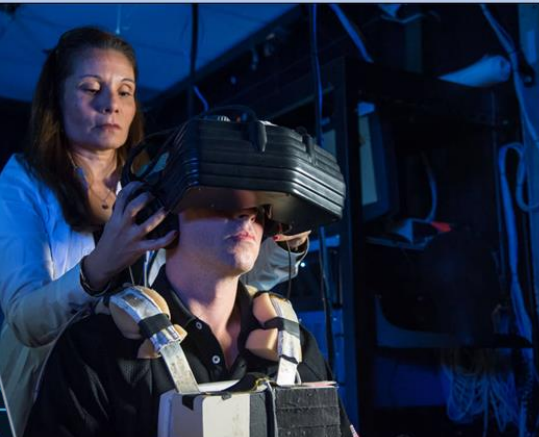
High Performance Computing



Deep Learning



Virtual Reality



Soft Errors

- Soft errors: most commonly observed errors
 - Single bit flips
- Masked output
 - Correct answer.
- Silent Data Corruption (SDC) output
 - Wrong answer.
- Other output
 - Crash, hang, ...

Reliability Research: Fault Injection

- Fault injection method
 - Injecting single-bit errors into different locations (fault sites) in applications
- Ground truth: *huge unreachable exhaustive fault sites!*

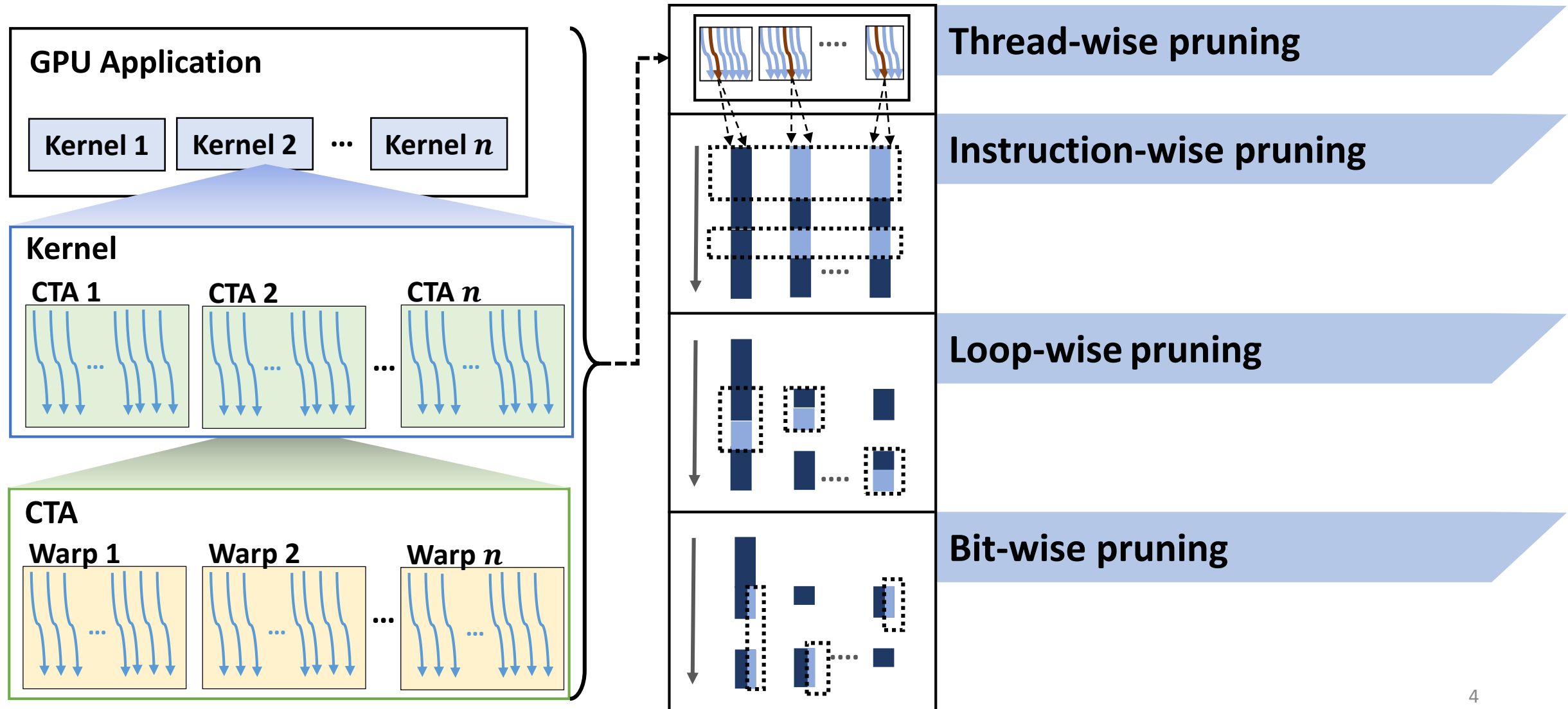
GEMM:

$16384 \text{ threads} \times 1305 \text{ insns/thd} \times 29.6 \text{ bits/insn} = 6.23 \times 10^8 \text{ fault sites!}$

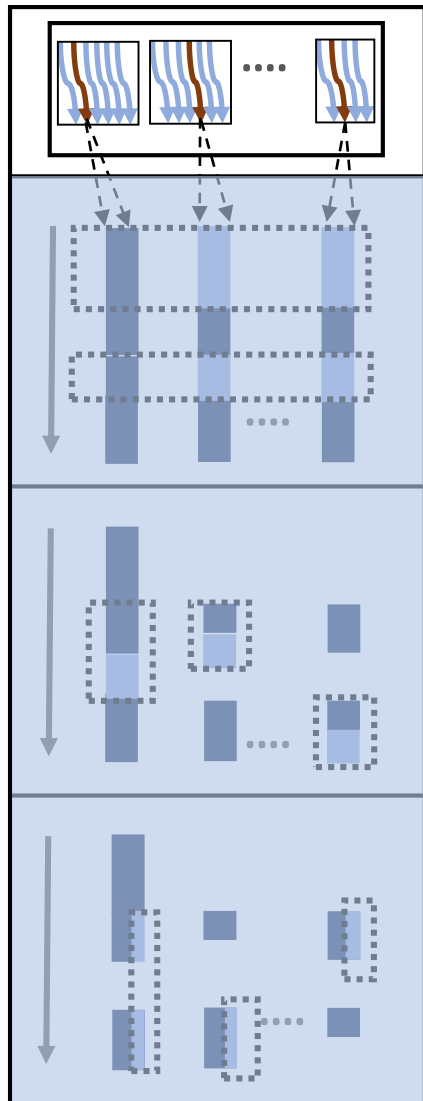
- Closest to ground truth (Baseline):
 - *Random sampling* based on statistics

Confidence Interval: 99.8%, Error Margin: 1.26%
→ 60K fault sites

GPU Architecture → Progressive Fault Sites Pruning



Progressive Fault Sites Pruning



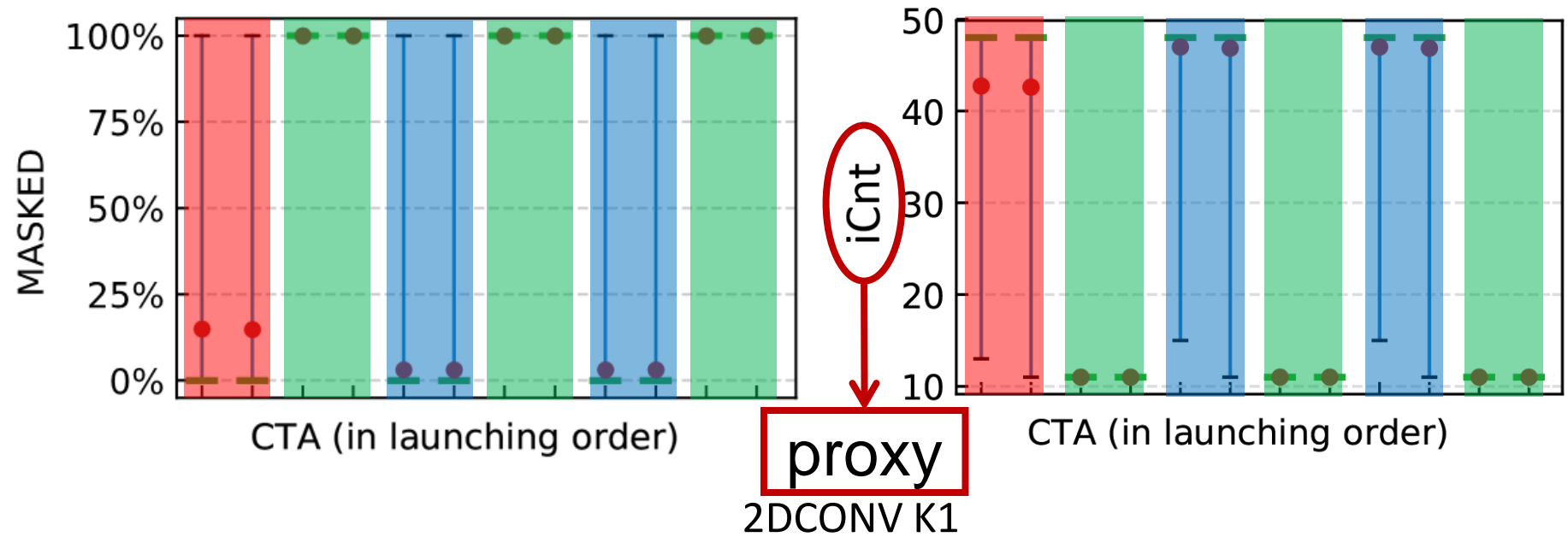
Thread-wise pruning

• CTA-level

Instruction-wise pruning

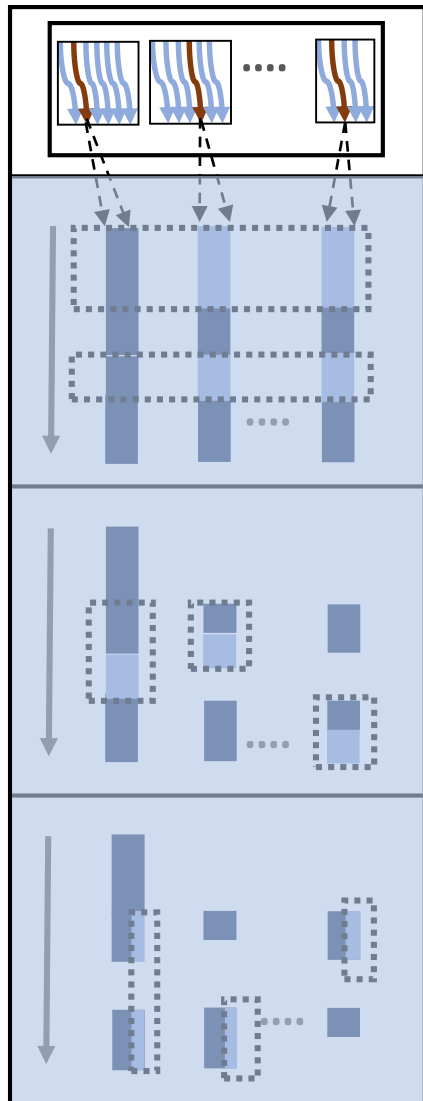
Loop-wise pruning

Bit-wise pruning



A few CTAs are enough; the rest are pruned

Progressive Fault Sites Pruning



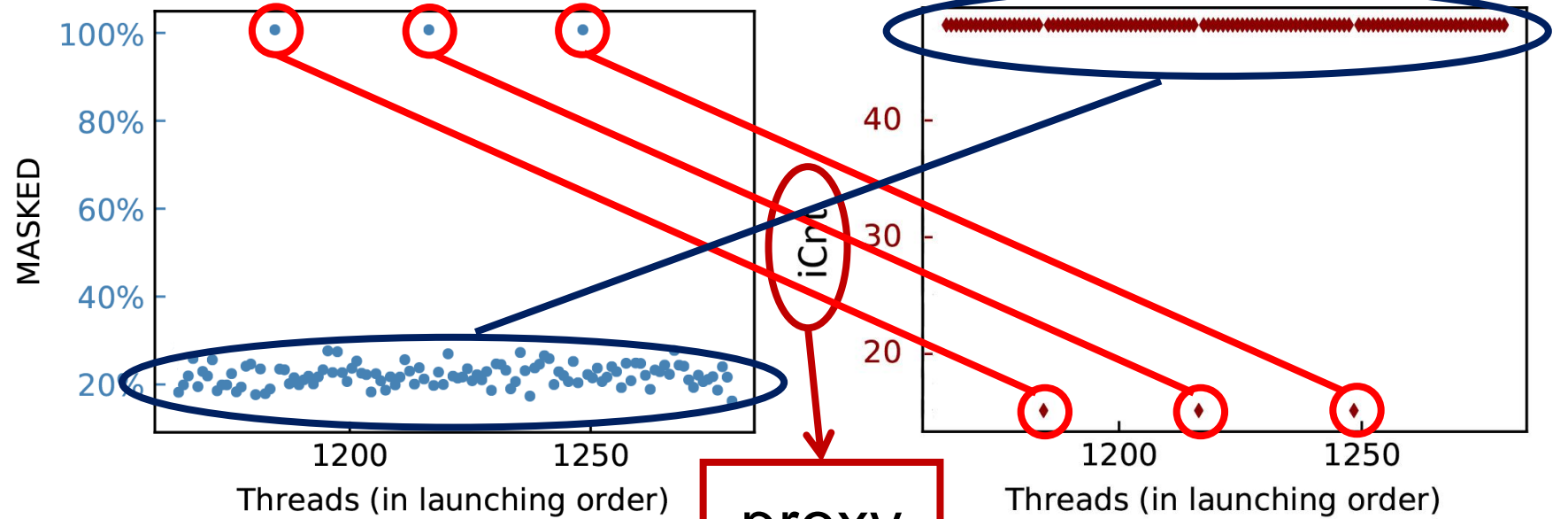
Thread-wise pruning

Instruction-wise pruning

Loop-wise pruning

Bit-wise pruning

• CTA-level → Thread-level

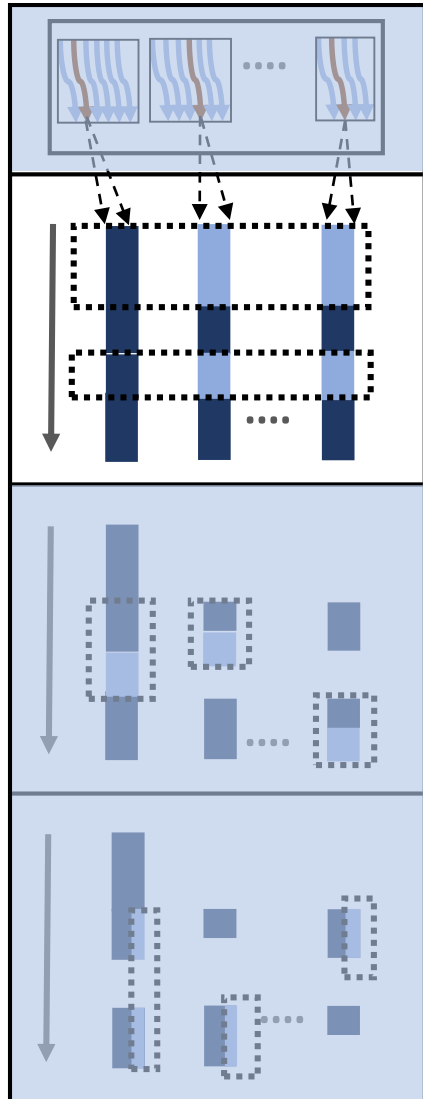


2DCONV K1, CTA(0-1-0)

A few threads are enough; the rest are pruned

Progressive Fault Sites Pruning

PathFinder K1



Thread-wise pruning

Instruction-wise pruning

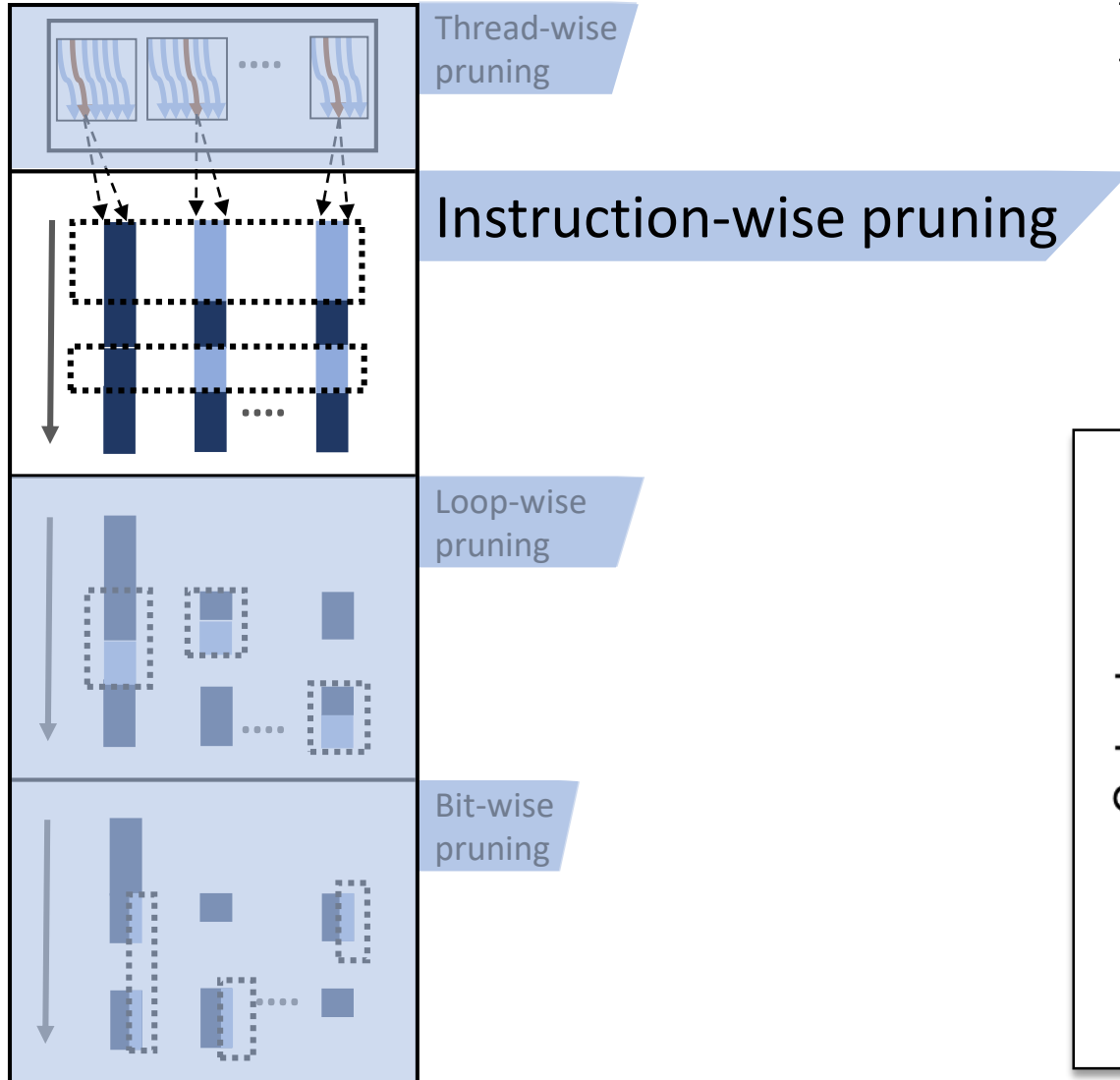
Loop-wise pruning

Bit-wise pruning

Thread "a" (<i>iCnt</i> = 533)		Thread "b" (<i>iCnt</i> = 516)	
1	<code>shl.u32 \$r3, s[0x0010], 0x00000001</code>	1	<code>shl.u32 \$r3, s[0x0010], 0x00000001</code>
2	<code>cvt.u32.u16 \$r1, %ctaid.x</code>	2	<code>cvt.u32.u16 \$r1, %ctaid.x</code>
3	<code>add.u32 \$r3, -\$r3, 0x00000100</code>	3	<code>add.u32 \$r3, -\$r3, 0x00000100</code>
4	<code>mul.wide.u16 \$r4, \$r1.lo, \$r3.hi</code>	4	<code>mul.wide.u16 \$r4, \$r1.lo, \$r3.hi</code>
5	<code>mad.wide.u16 \$r4, \$r1.hi, \$r3.lo, \$r4</code>	5	<code>mad.wide.u16 \$r4, \$r1.hi, \$r3.lo, \$r4</code>
.....		
49	<code>cvt.s32.s32 \$r2, -\$r2</code>	49	<code>cvt.s32.s32 \$r2, -\$r2</code>
50	<code>and.b32 \$p0 \$o127, \$r5, \$r2</code>	50	<code>and.b32 \$p0 \$o127, \$r5, \$r2</code>
51	<code>ssy 0x00000228</code>	51	<code>ssy 0x00000228</code>
52	<code>mov.u32 \$r2, \$r124</code>	52	<code>mov.u32 \$r2, \$r124</code>
53	<code>@\$p0.eq bra l0x00000228</code>	53	<code>@\$p0.eq bra l0x00000228</code>
54	<code>add.half.u32 \$r7, s[0x0038], \$r1</code>		
55	<code>mov.half.u32 \$r2, s[0x0030]</code>		
56	<code>mul.wide.u16 \$r8, \$r2.lo, \$r7.hi</code>		
57	<code>mad.wide.u16 \$r8, \$r2.hi, \$r7.lo, \$r8</code>		
58	<code>shl.u32 \$r8, \$r8, 0x00000010</code>		
.....			
66	<code>min.s32 \$r7, s[\$ofs2+0x0040], \$r8</code>		
67	<code>ld.global.u32 \$r2, [\$r2]</code>		
68	<code>add.u32 \$r2, \$r2, \$r7</code>		
69	<code>mov.u32 s[\$ofs3+0x0440], \$r2</code>		
70	<code>mov.u32 \$r2, 0x00000001</code>		
71	<code>l0x00000228: nop</code>	54	<code>l0x00000228: nop</code>
72	<code>bar.sync 0x00000000</code>	55	<code>bar.sync 0x00000000</code>
73	<code>set.eq.s32.s32 \$p0 \$o127, \$r6, \$r1</code>	56	<code>set.eq.s32.s32 \$p0 \$o127, \$r6, \$r1</code>
74	<code>@\$p0.ne bra l0x000002b8</code>	57	<code>@\$p0.ne bra l0x000002b8</code>
75	<code>set.ne.s32.s32 \$p1 \$r1, \$r2, \$r124</code>	58	<code>set.ne.s32.s32 \$p1 \$r1, \$r2, \$r124</code>
.....		
529	<code>set.eq.s32.s32 \$p0 \$o127, \$r6, \$r1</code>	512	<code>set.eq.s32.s32 \$p0 \$o127, \$r6, \$r1</code>
530	<code>@\$p0.ne bra l0x000002b8</code>	513	<code>@\$p0.ne bra l0x000002b8</code>
531	<code>l0x000002b8: set.ne.s32.s32 \$p0 \$o127, \$r2, \$r124</code>	514	<code>l0x000002b8: set.ne.s32.s32 \$p0 \$o127, \$r2, \$r124</code>
532	<code>bra l0x000002c8</code>	515	<code>bra l0x000002c8</code>
533	<code>l0x000002c8: @\$p0.eq retpz</code>	516	<code>l0x000002c8: @\$p0.eq retpz</code>

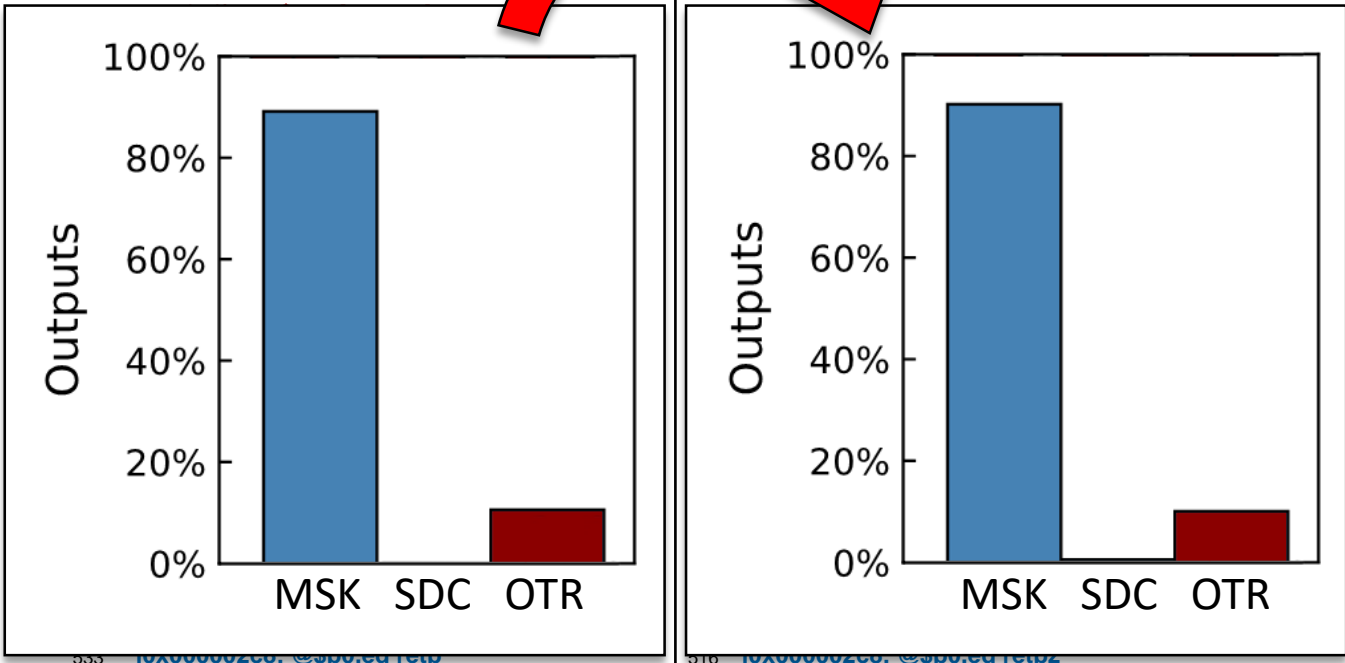
Progressive Fault Sites Pruning

PathFinder K1

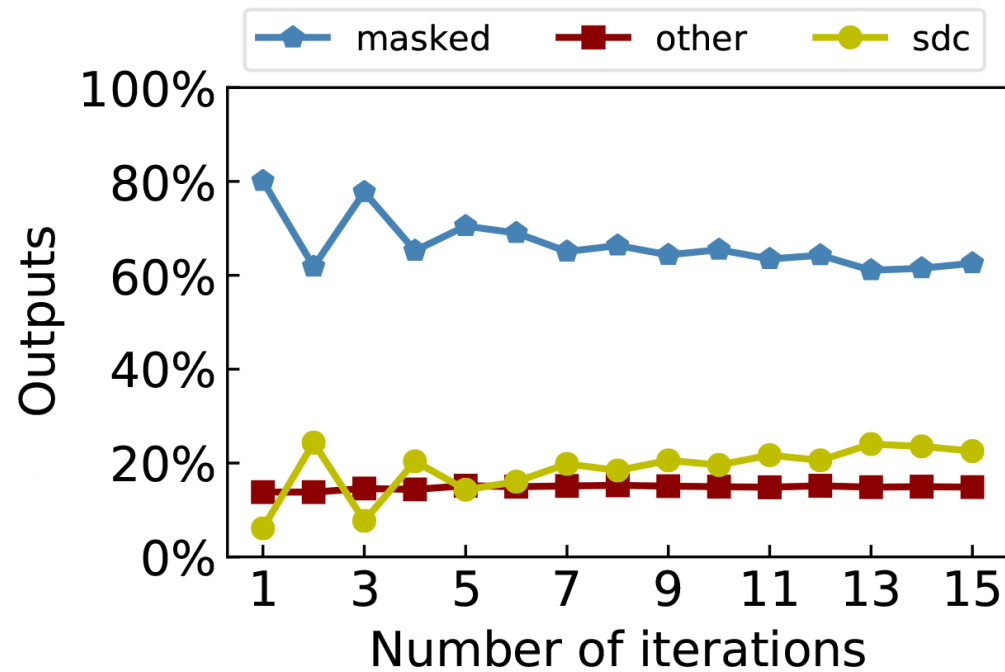
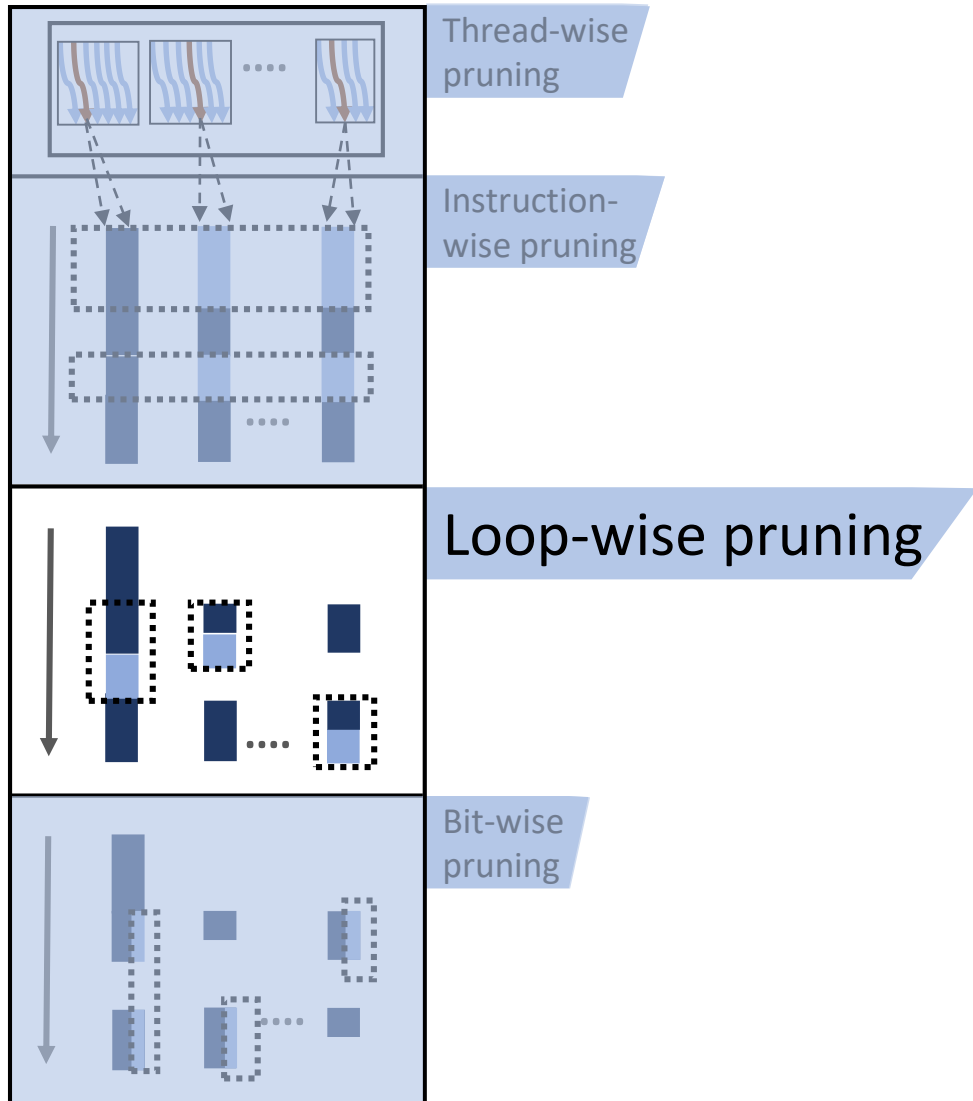


Thread "a" (iCnt = 533)		Thread "b" (iCnt = 516)	
1	shl.u32 \$r3, s[0x0010], 0x00000001	1	shl.u32 \$r3, s[0x0010], 0x00000001
2	cvt.u32.u16 \$r1, %ctaid.x	2	cvt.u32.u16 \$r1, %ctaid.x
3	add.u32 \$r3, -\$r3, 0x00000100	3	add.u32 \$r3, -\$r3, 0x00000100
4	mul.wide.u16 \$r4, \$r1.lo, \$r3.hi	4	mul.wide.u16 \$r4, \$r1.lo, \$r3.hi
5	mad.wide.u16 \$r4, \$r1.hi, \$r3.lo, \$r4	5	mad.wide.u16 \$r4, \$r1.hi, \$r3.lo, \$r4
...
49	cvt.s32.s32 \$r2, -\$r2	49	cvt.s32.s32 \$r2, -\$r2
50	and.b32 \$p0 so127, \$r5, \$r2	50	and.b32 \$p0 so127, \$r5, \$r2
51	ssy 0x00000228	51	ssy 0x00000228
52	mov.u32 \$r2, \$r124	52	mov.u32 \$r2, \$r124
53	@\$p0.eq bra l0x00000228	53	@\$p0.eq bra l0x00000228
54	add.half.u32 \$r7, s[0x0038], \$r1	54	add.half.u32 \$r7, s[0x0038], \$r1

Extrapolate

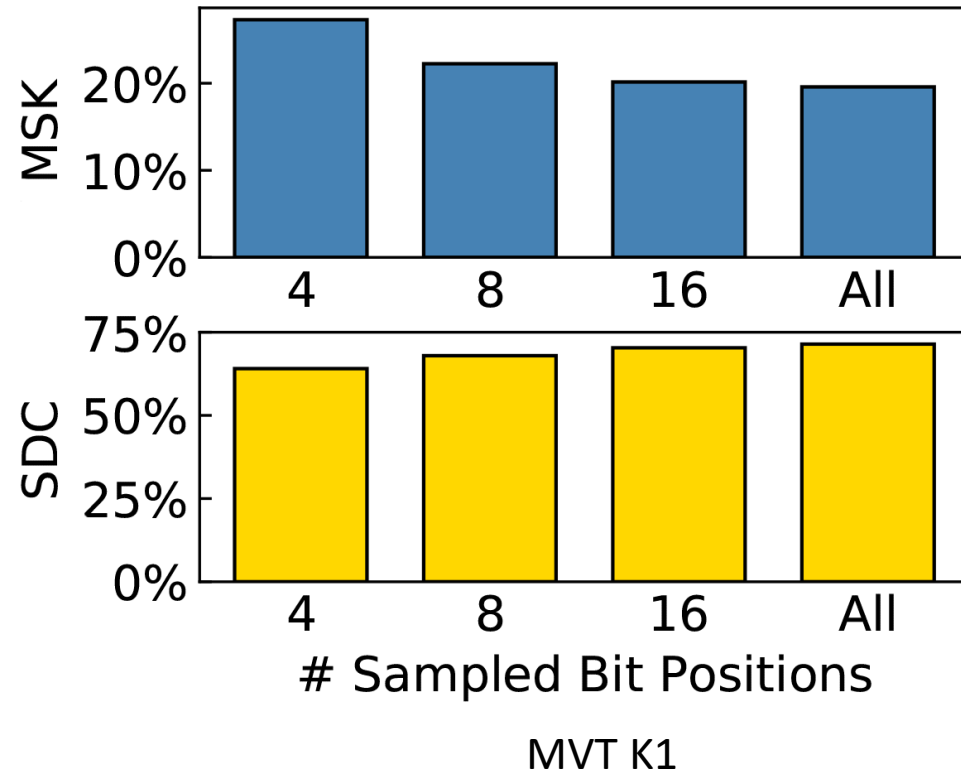
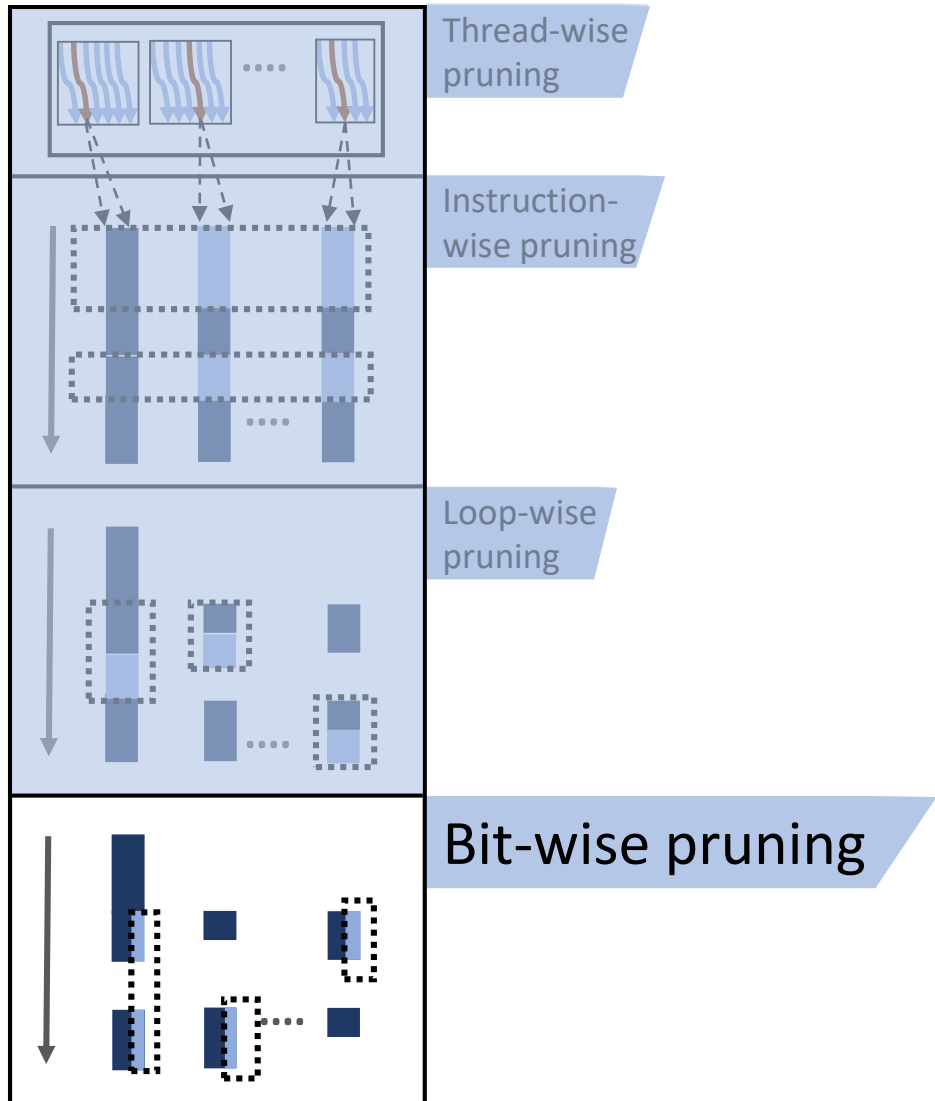


Progressive Fault Sites Pruning



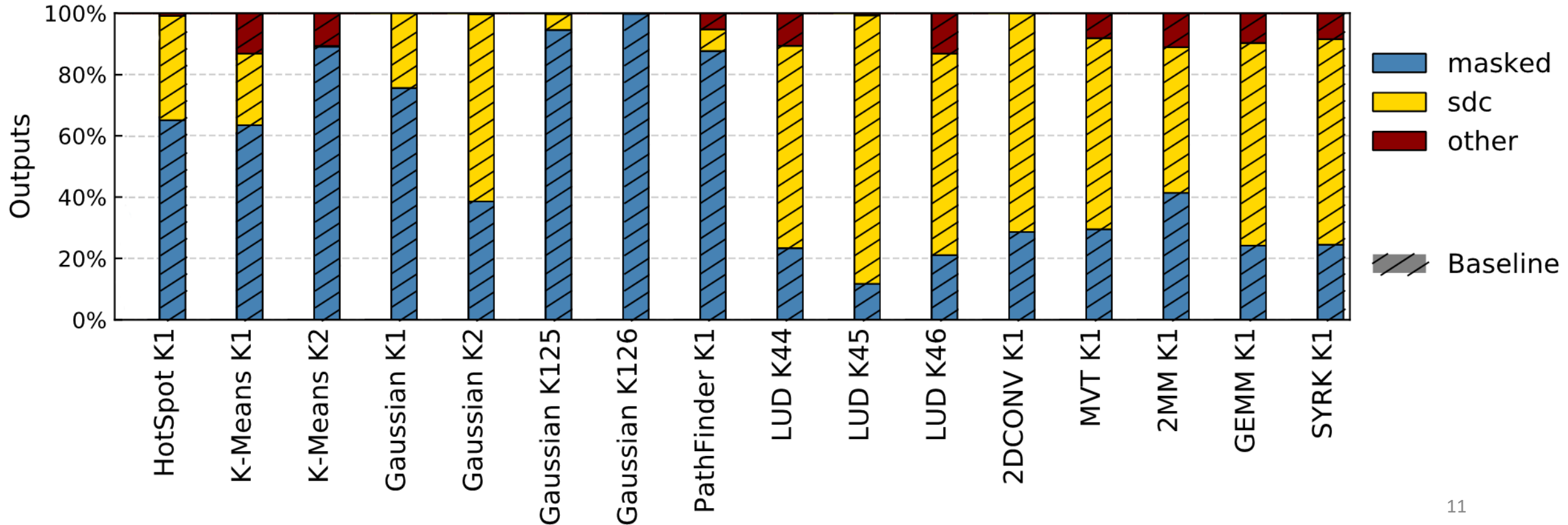
K-Means K1
34 loop iterations

Progressive Fault Sites Pruning



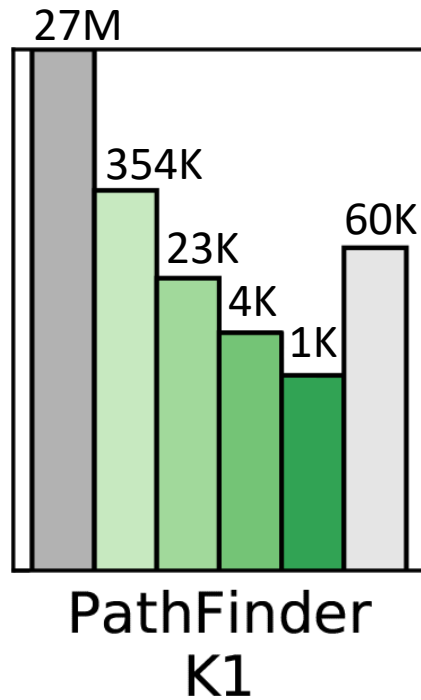
Evaluation

- Fault Site Pruning vs. Baseline
 - **Baseline: 60K** fault sites (confidence interval=99.8%, error margin=1.26%)
 - Closest to ground truth



Evaluation

- Effectiveness



	Exhaustive Fault Sites	Pruned Fault Sites
Average	1.29×10^8	1190
Min	1.09×10^5	318
Max	6.23×10^8	4678

108,403x Reduction!



Conclusion

- GPGPU applications: huge unreachable exhaustive fault sites
- **Progressive Fault Site Pruning**
 - GPGPU-specific features
- **Accurate** GPU reliability assessment
- Significant reduction:
 - Up to **7** orders of magnitude
 - **108,403x** on average



WILLIAM & MARY

CHARTERED 1693

Thank you :)

Fault Site Pruning for Practical Reliability Analysis of GPGPU Applications

Bin Nie, **Lishan Yang**, Adwait Jog, and Evgenia Smirni
College of William & Mary

This work is supported by NSF grants CCF-17175332 (CRII) and CCF-1750667 (CORE Small)