



Java on ARM Theory, Applications and Workloads

Dmitry Chuyko
JVM Team

Who we are

Dmitry Chuyko

 [@dchuyko](https://twitter.com/dchuyko)

 BELL^{SOFT}

Liberica JDK – verified OpenJDK binary

<http://bell-sw.com>

Ex-employers

 ORACLE®





Committed to freedom

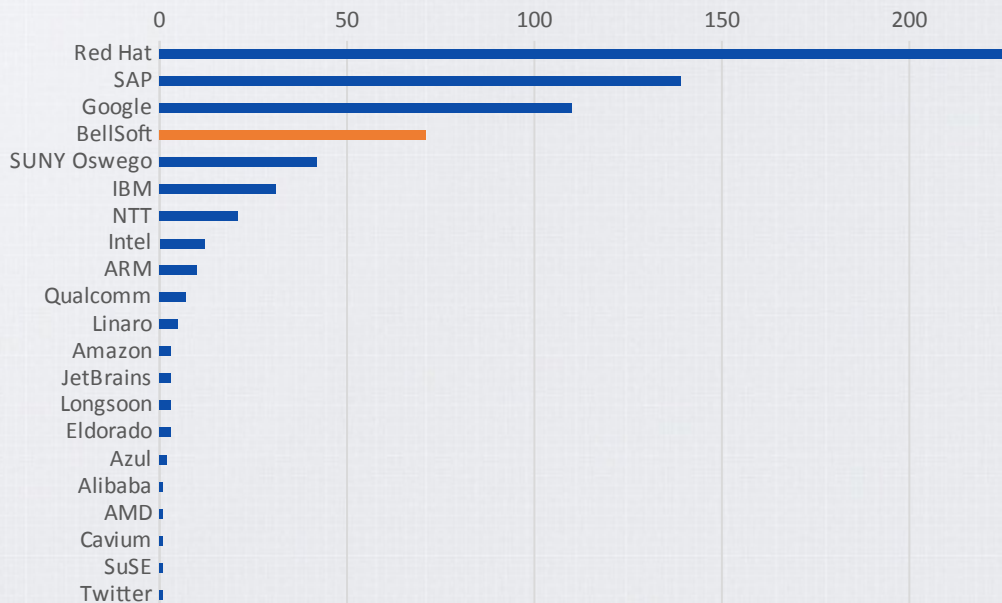


Liberica

– supported OpenJDK binaries

<http://bell-sw.com>

External contributions to OpenJDK jdk/jdk Aug '17 - Aug '18
*Note: Oracle contributed ~3965 patches in the same period



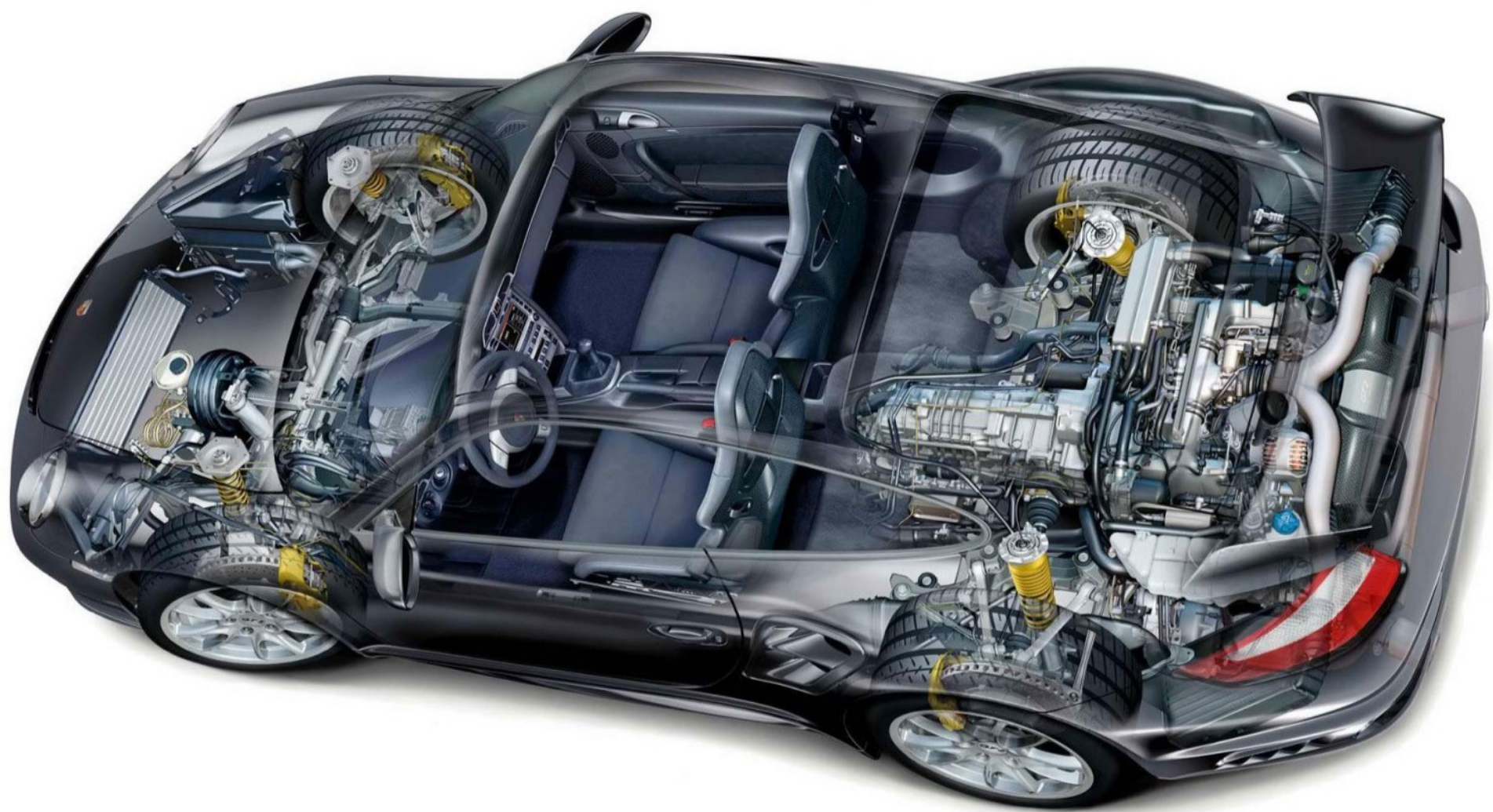


What do we know about Arm?

- Arm = Advanced RISC Machine/Acorn RISC Machine
- Founded in 1985
- UK, Cambridge
- ARM is a RISC architecture

- 30 billion processors shipped in 2013
- Plans to ship 100 billion processors by 2020







IoT Gateways

Liberica JDK



SuperMicro



Dell



Eurotech



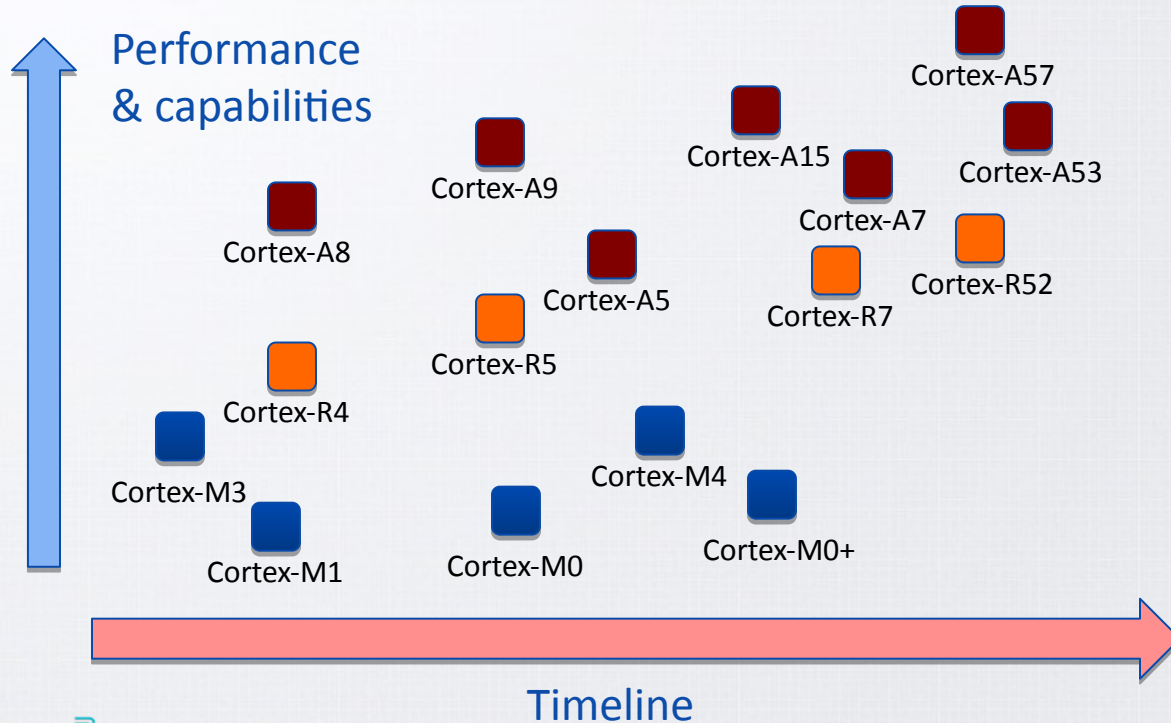
Advantech

....

But Servers?



Arm: architecture, profile, implementation

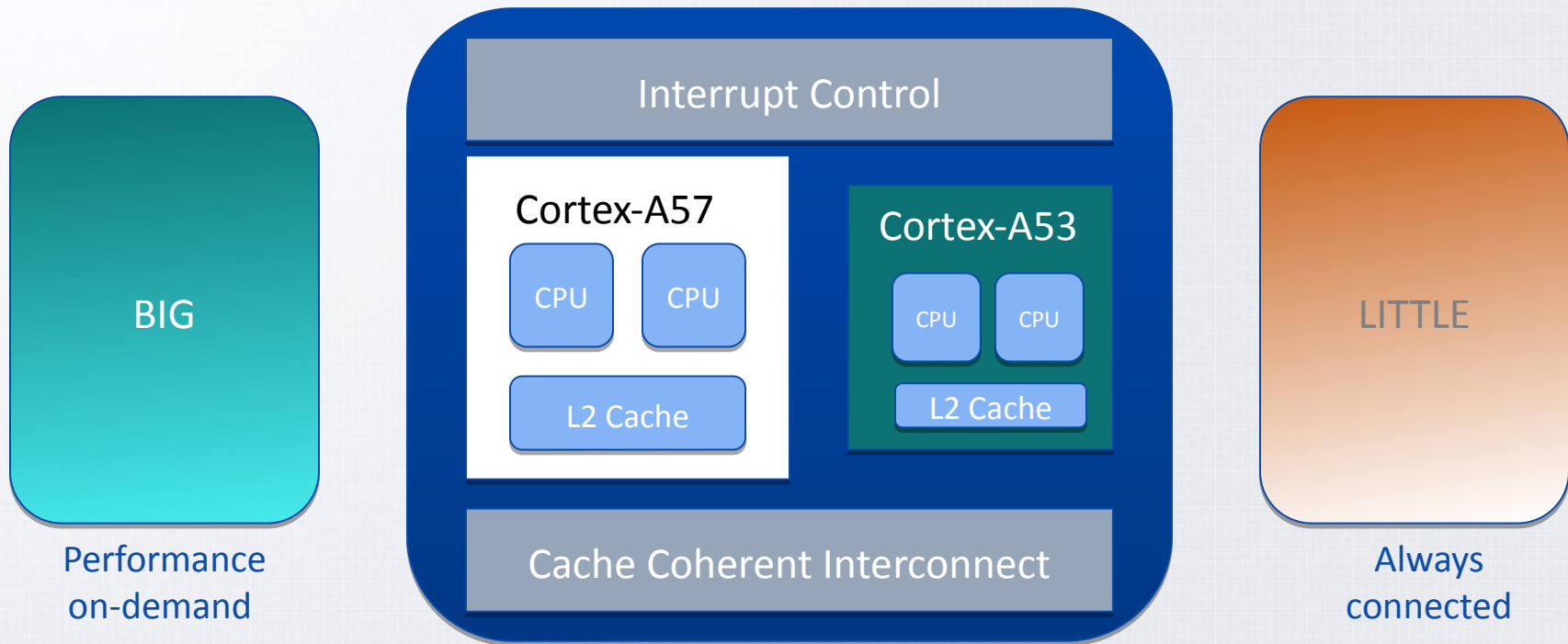


- ARM v7
- Architecture profiles
 - v7-M (Embedded)
 - V7-R (Real-Time)
 - V7-A (Application)

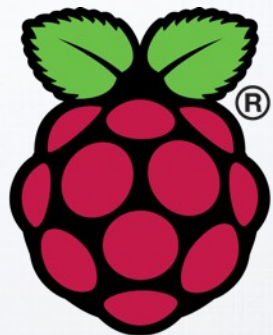
- ARM v8
- Architecture profiles
 - v8-M (Embedded)
 - V8-R (Real-Time)
 - V8-A (Application)



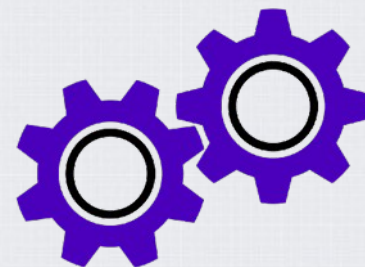
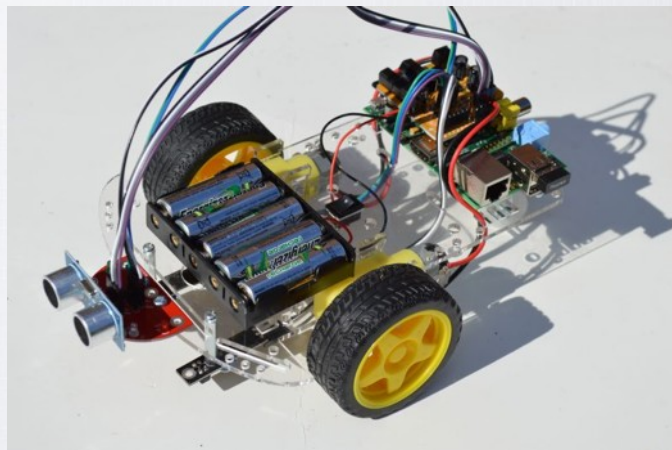
Arm: big.LITTLE



DIY



Raspberry Pi



Robo4J



OpenJDK Arm32 port

- Available since OpenJDK 9
 - Minimal VM, Client VM, Server VM
- Works on the Raspberry Pi
- jlink + jdeps
 - Allows to create a smaller runtime (as small as 16 Mb)
- Java FX Embedded
 - Allows to build fancy UI for the Raspberry Pi
 - EGL/DFB acceleration
 - Touch screen support



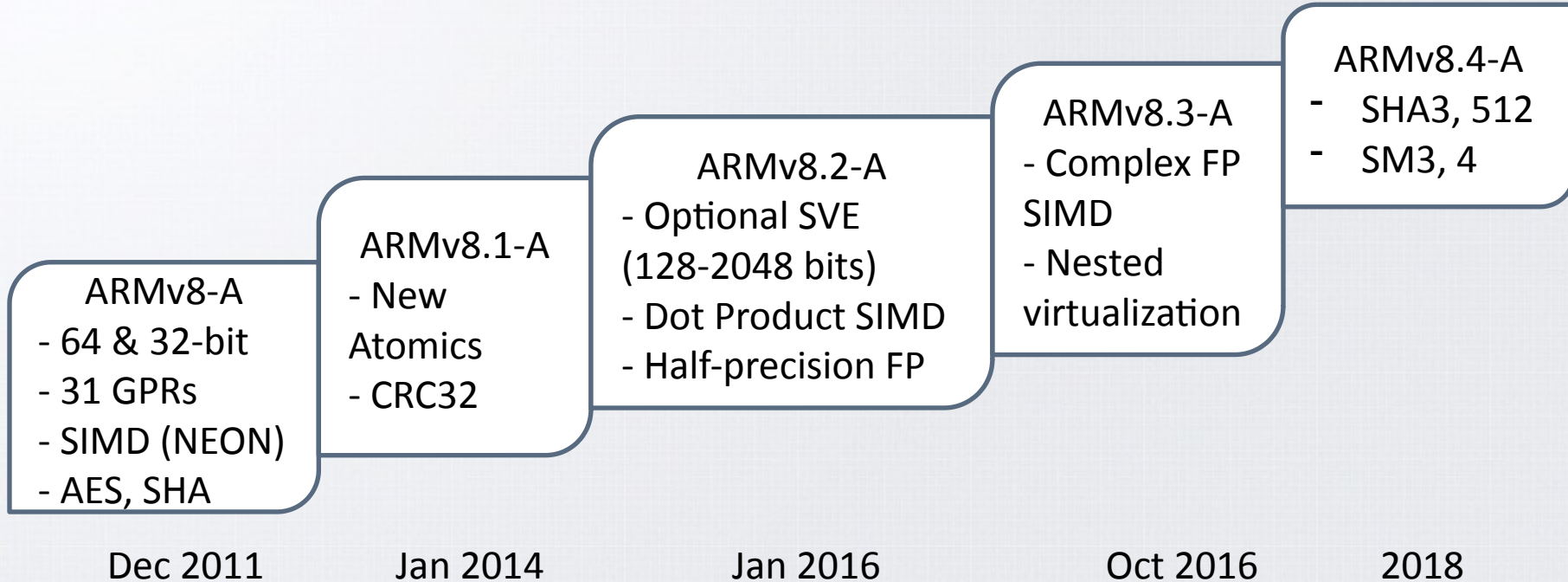
Minimal VM

- Optimized for footprint, rather than functionality
 - Serial GC
 - C1 JIT compiler
 - No JDWP support
 - No JMX support
- But... it is **< 4 Mb!**
 - Linux x86_64 Server VM: 23 Mb
- jlink @since jdk9
 - java.base with Minimal VM under 16 Mb!
 - Modules for jetty: under 32 Mb





ARMv8-A Specification



....

Arm architecture licensees

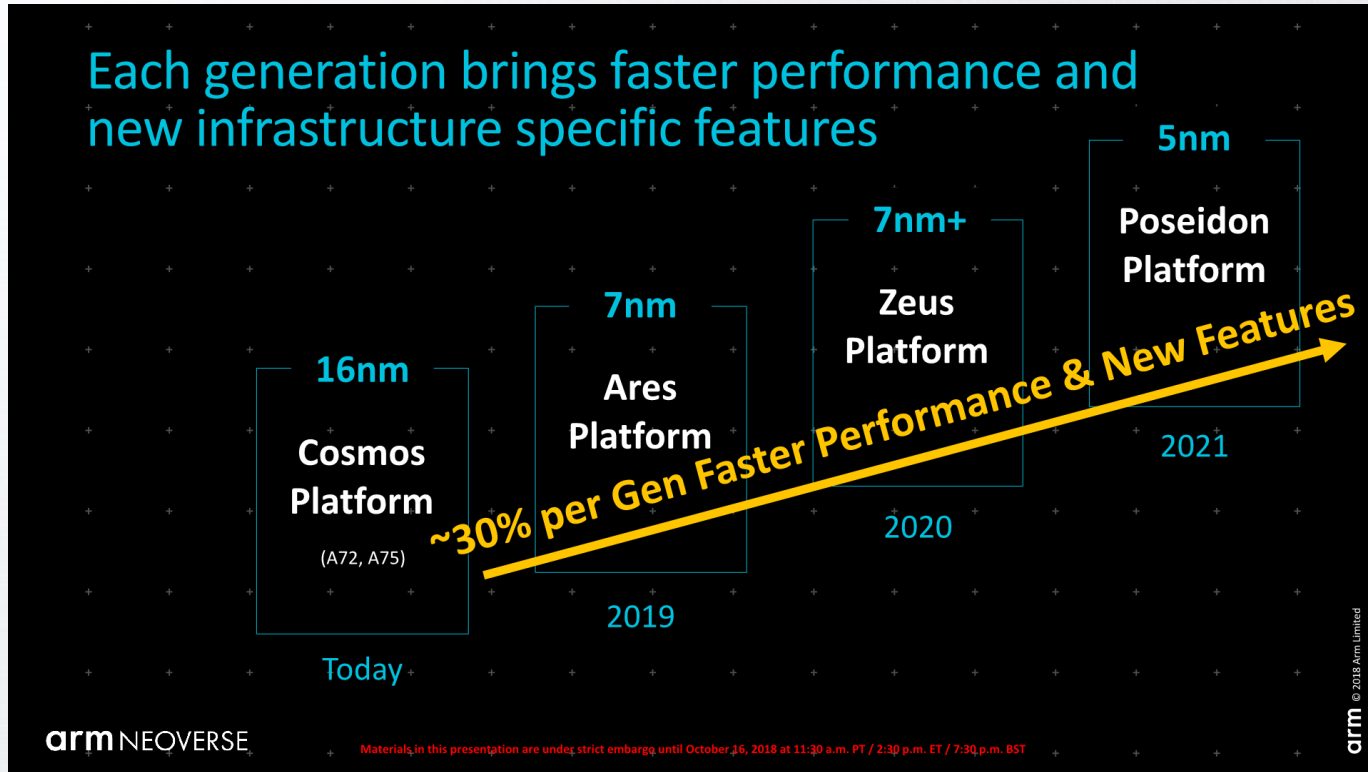


Apple



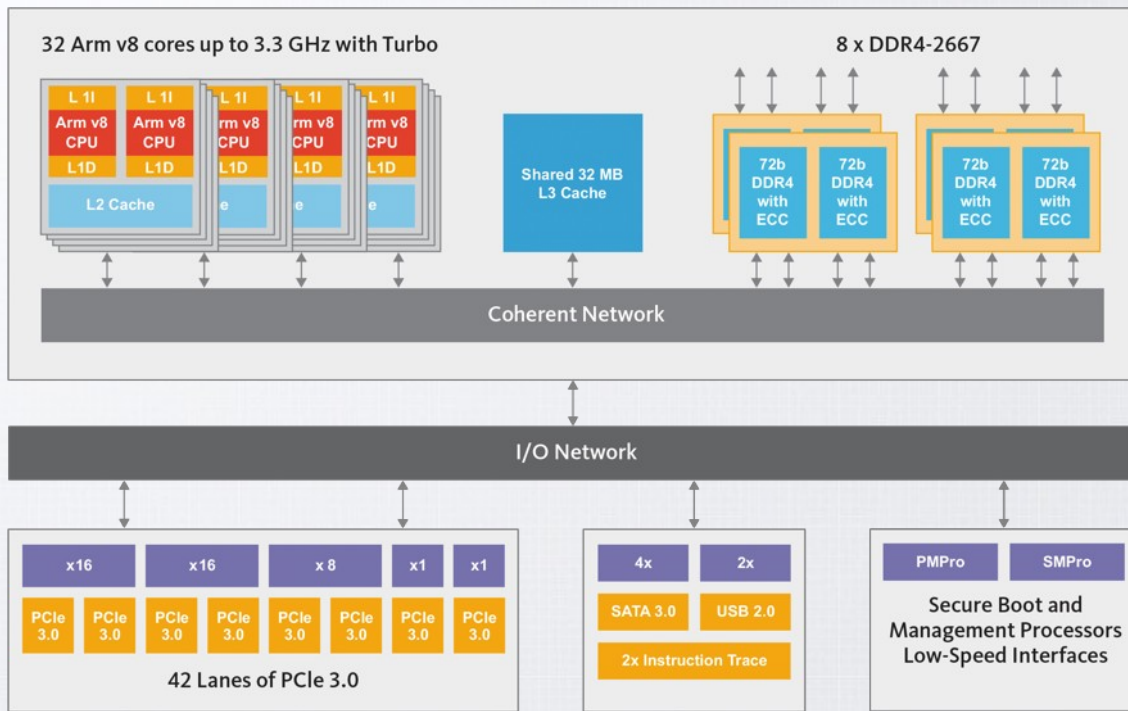


Neoverse Infrastructure IP Roadmap





Ampere Computing (ex APM)

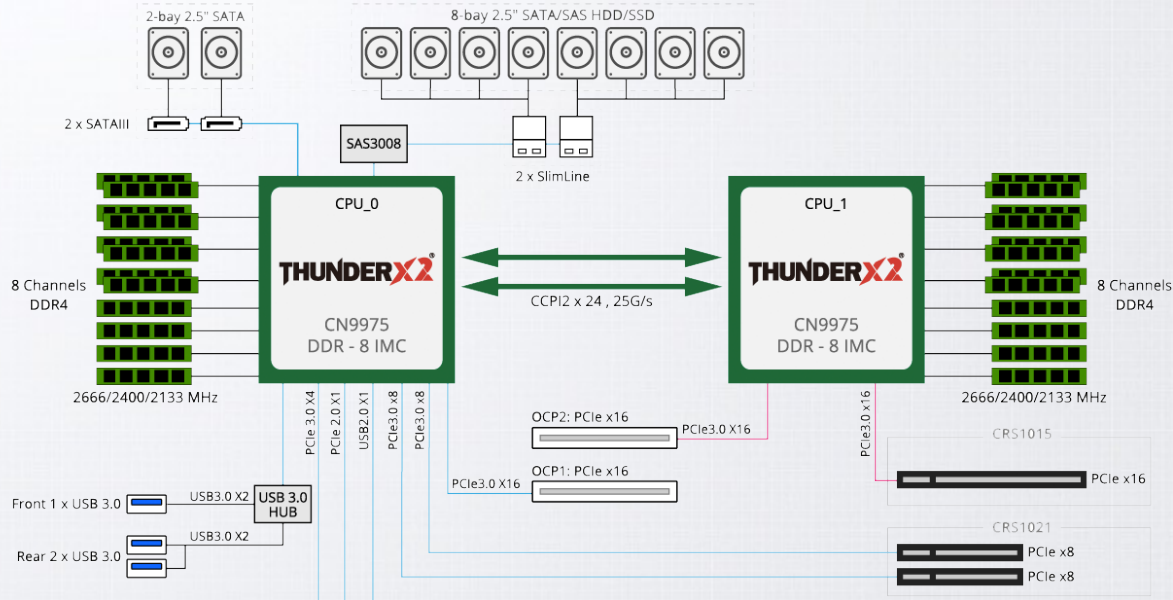


Up to 32 cores
Up to 32 threads
8 DDR Channels
32 Mb L3





Cavium/Marvell ThunderX2



32 cores/128 threads
32 Mb L3
8 DDR Channels/socket
Multi-socket
Up to 4 TB RAM



That thing
is real!



Wait, how many threads?

1	[91.1%]	57	[90.0%]	113	[96.2%]	169	[96.2%]
2	[90.7%]	58	[90.8%]	114	[96.2%]	170	[96.2%]
3	[91.2%]	59	[90.3%]	115	[96.6%]	171	[96.2%]
4	[90.3%]	60	[89.9%]	116	[94.6%]	172	[95.8%]
5	[90.8%]	61	[89.5%]	117	[95.8%]	173	[95.4%]
6	[90.8%]	62	[89.5%]	118	[96.2%]	174	[95.8%]
7	[89.9%]	63	[90.3%]	119	[96.2%]	175	[96.7%]
8	[89.9%]	64	[89.9%]	120	[95.8%]	176	[96.2%]
9	[92.1%]	65	[91.7%]	121	[96.2%]	177	[96.2%]
10	[91.6%]	66	[89.9%]	122	[96.2%]	178	[95.8%]
11	[91.6%]	67	[91.6%]	123	[96.2%]	179	[96.2%]
12	[91.2%]	68	[89.5%]	124	[95.8%]	180	[96.2%]
13	[91.2%]	69	[91.6%]	125	[96.2%]	181	[95.8%]
14	[89.1%]	70	[91.2%]	126	[96.2%]	182	[94.6%]
15	[89.1%]	71	[88.7%]	127	[96.2%]	183	[95.8%]
16	[88.6%]	72	[89.2%]	128	[96.2%]	184	[95.8%]
17	[88.7%]	73	[89.5%]	129	[95.8%]	185	[95.8%]
18	[89.9%]	74	[89.5%]	130	[95.8%]	186	[94.6%]
19	[89.9%]	75	[89.9%]	131	[96.2%]	187	[95.4%]
20	[88.8%]	76	[88.3%]	132	[95.8%]	188	[95.8%]
21	[89.7%]	77	[89.7%]	133	[95.8%]	189	[95.4%]
22	[88.7%]	78	[90.0%]	134	[95.4%]	190	[95.8%]
23	[91.6%]	79	[89.5%]	135	[95.8%]	191	[95.8%]
24	[87.8%]	80	[88.7%]	136	[96.2%]	192	[94.6%]
25	[89.0%]	81	[90.3%]	137	[96.2%]	193	[95.8%]
26	[90.0%]	82	[88.3%]	138	[95.8%]	194	[94.6%]
27	[90.0%]	83	[89.1%]	139	[96.2%]	195	[96.2%]
28	[87.9%]	84	[88.3%]	140	[94.6%]	196	[95.8%]
29	[89.5%]	85	[91.6%]	141	[96.2%]	197	[96.2%]
30	[91.6%]	86	[90.4%]	142	[96.2%]	198	[96.2%]
31	[90.8%]	87	[90.8%]	143	[96.2%]	199	[96.6%]
32	[90.3%]	88	[89.5%]	144	[96.2%]	200	[95.9%]
33	[90.8%]	89	[90.3%]	145	[95.8%]	201	[95.8%]
34	[91.2%]	90	[90.8%]	146	[96.2%]	202	[95.8%]
35	[89.8%]	91	[89.5%]	147	[96.2%]	203	[96.2%]
36	[90.3%]	92	[89.5%]	148	[97.1%]	204	[96.2%]
37	[90.0%]	93	[89.9%]	149	[95.8%]	205	[95.8%]
38	[89.9%]	94	[91.2%]	150	[96.2%]	206	[96.2%]
39	[89.5%]	95	[89.9%]	151	[96.2%]	207	[96.2%]
40	[91.2%]	96	[89.6%]	152	[95.8%]	208	[96.2%]
41	[89.2%]	97	[90.0%]	153	[96.6%]	209	[96.2%]
42	[89.5%]	98	[90.4%]	154	[95.8%]	210	[96.2%]
43	[90.0%]	99	[89.1%]	155	[95.8%]	211	[96.2%]
44	[88.7%]	100	[88.2%]	156	[95.8%]	212	[95.8%]
45	[87.8%]	101	[89.5%]	157	[95.8%]	213	[96.2%]
46	[88.3%]	102	[87.9%]	158	[96.2%]	214	[96.2%]
47	[89.0%]	103	[89.1%]	159	[96.2%]	215	[95.8%]
48	[88.7%]	104	[88.2%]	160	[96.2%]	216	[94.9%]
49	[89.1%]	105	[89.1%]	161	[96.2%]	217	[96.2%]
50	[87.9%]	106	[87.8%]	162	[95.8%]	218	[95.8%]
51	[89.1%]	107	[90.0%]	163	[95.8%]	219	[95.8%]
52	[87.0%]	108	[88.2%]	164	[95.8%]	220	[95.8%]
53	[89.5%]	109	[88.7%]	165	[95.4%]	221	[95.8%]
54	[89.1%]	110	[88.3%]	166	[96.2%]	222	[95.8%]
55	[88.3%]	111	[89.9%]	167	[94.6%]	223	[96.2%]
56	[87.8%]	112	[88.2%]	168	[96.2%]	224	[95.8%]

Mem [|||||] 13.1G/511G Tasks: 50, 660 thr: 244 running
Swp [|||||] 0K/0K Load average: 205.40 100.67 39.63
Uptime: 30 days, 17:53:24

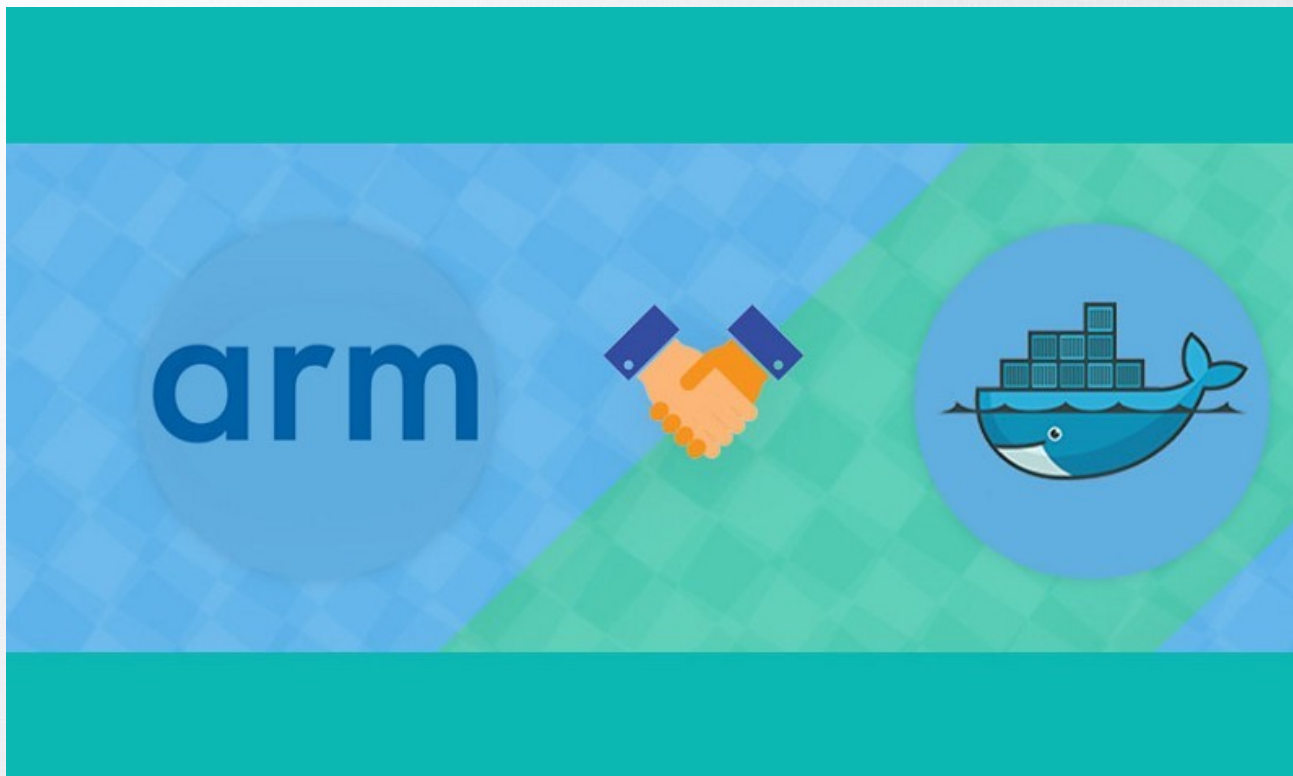
Arm Software ecosystem

<p>Key Applications Middleware</p>	 	 	 	 	 	
<p>Operating System, Virtualization & Firmware</p>		<p>ACPI</p>	 	 	  <small>Supported by Canonical</small>	

Check out if it works on Arm: <https://worksonarm.com>



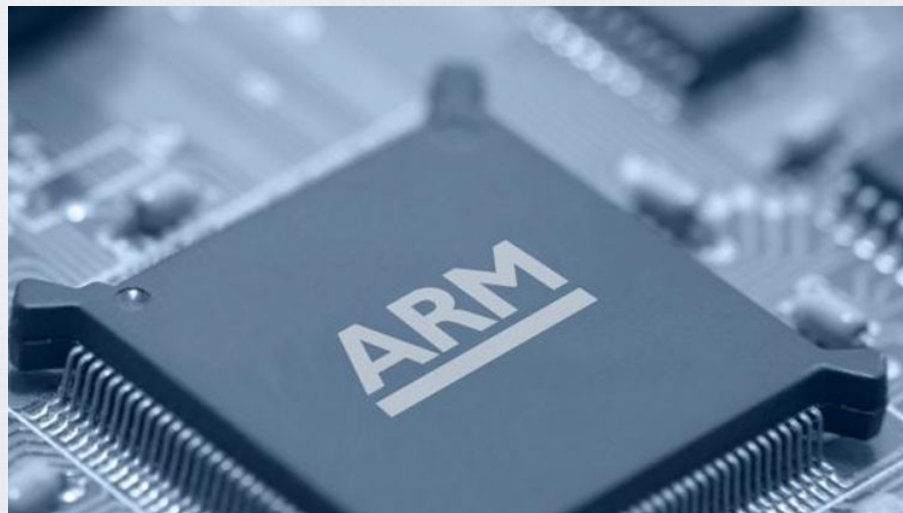
Cross-build Docker images





OpenJDK ARM ports

- ARM (32 bit & 64 bit)
 - Full Java SE Spec
 - ARM v6/v7/v8
 - C1 & C2
- AARCH64 (64 bit only)
 - Full Java SE Spec
 - C1 & C2
 - G1 / Parallel GC / Shenandoah
(and ZGC is coming)
 - AppCDS, JFR, NMT, AOT





Intrinsics

Intrinsic:

“function (subroutine) available for use in a given programming language which implementation is handled specially by the compiler.”



Intrinsics

- GCC/LLVM
 - Specialized instructions not expressible through regular language constructs
 - Wrappers around libc calls
- HotSpot
 - Manipulation on C2 IR
 - Typically a specialized assembly instruction call for a given architecture
 - Stub – assembler or native routines

What will C2 do with math Java code?

java.lang.Math:

```
/**
 * Returns as a {@code long} the most significant 64 bits of the
 * 128-bit product of two 64-bit factors.
 * @since 9
 */
public static long multiplyHigh(long x, long y) {
    if (x < 0 || y < 0) {
        long x1 = x >> 32;
        long x2 = x & 0xFFFFFFFFL;
        long y1 = y >> 32;
        long y2 = y & 0xFFFFFFFFL;
        long z2 = x2 * y2;
        long t = x1 * y2 + (z2 >>> 32);
        long z1 = t & 0xFFFFFFFFL;
        long z0 = t >> 32;
        z1 += x2 * y1;
        return x1 * y1 + z0 + (z1 >> 32);
    } else { ...
```

What will C2 do with math Java code?

java.lang.Math:

```
/**  
 * Returns as a {@code long} the most significant 64 bits of the 128-bit  
 * product of two 64-bit factors.  
 * @since 9  
 */  
public static long multiplyHigh(long x, long y) {  
    // Use technique from section 8-2 of Henry S. Warren, Jr.,  
    // Hacker's Delight (2nd ed.) (Addison Wesley, 2013), 173-174.  
    ...  
    // Use Karatsuba technique with two base 232 digits.  
    ...  
    return ...;  
}
```

Math code in assembly

```
    mul x12, x10, x11
    mul x13, x12, x12, #32
    lsr x13, x13, #32
    mul x14, x13, x13
    mul x15, x14, x14
    mul x11, x14, x11
    add x10, x15, x10
    mul x14, x15, x14
    mul x10, x10, x11
```

14 operations with latency

1

```
sub x10, x10, x14
```

Can we make it faster?

- Rewrite as a C + JNI call
 - Well, it will be slower
- Tune HotSpot to optimize IR for this code better*
 - Even if this is possible, this might lead to regressions
- Tune HotSpot to detect this method and substitute optimal code instead

SMULH Xd, Xn, Xm (cost: 4)
“Signed multiply high”

C2 Intrinsic How-to

- 1) Add **SMULH** instruction into `${arch}/assembler_${arch}.hpp`
- 2) Describe a node with this instruction and its cost in `${arch}.ad`
- 3) Mark this method as intrinsic in `share/classfile/vmSymbols.hpp`
- 4) Substitute the method with the node

```
bool LibraryCallKit::inline_math_multiplyHigh() {  
    set_result(_gvn.transform(new MulHiLNode(arg (0), arg (2))));  
    return true;  
}
```

- 5) Annotate `j.l.Math.multiplyHigh()` **@HotSpotIntrinsicCandidate**
- 6) Measure performance

Benchmarking (throughput)

```
public class MultiplyHighJMHBench {  
  
    @Benchmark  
    @OperationsPerInvocation(10000)  
    public long bench() {  
        long op = System.currentTimeMillis();  
        long accum = 0;  
        for (int i = 0; i < 10000; i++) {  
            accum += Math.multiplyHigh(op + i, op + i);  
        }  
        return accum;  
    }  
}
```

3.5x better

SMULH cost: 4
Good for JDK 11!



Let's do something useful for enterprise apps

- What does a JVM do when executing a typical enterprise program?
 - Creates, copies objects, strings, arrays, frees memory
 - Searches and compares objects, strings, arrays
 - Checks that the right information is received





```
String s = new String("Can this work faster?");
```

- Compact Strings @since JDK 9
 - Most strings do not require UTF-16 as inner representation
 - Inner representation of strings:
 - char[] -> byte[], coder
 - Either ISO-8859-1/Latin-1
 - Either UTF-16 if required

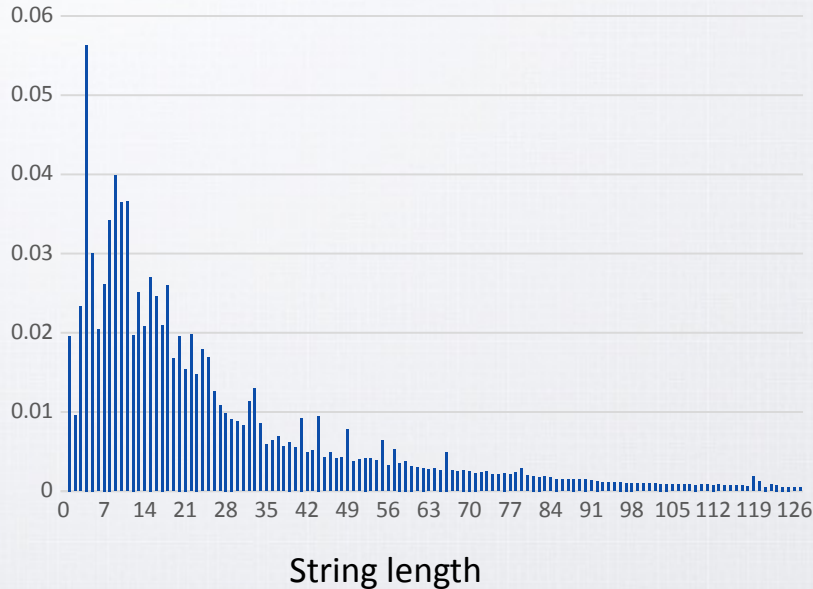
S	t	r	i	n	g
---	---	---	---	---	---

С	т	р	о	к	а
---	---	---	---	---	---



1001 Heap Dump

String length distribution

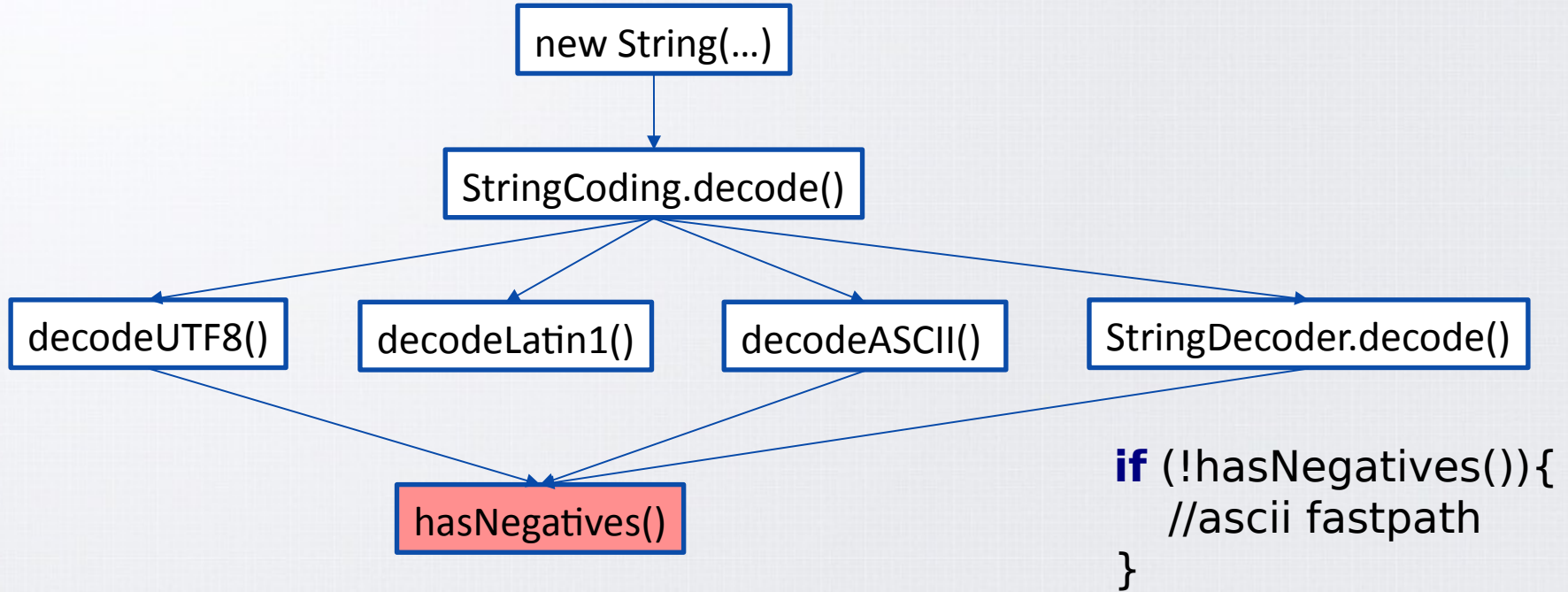


- Log-normal distribution
- < 0.3% of all strings are not Latin-1
- 18% strings < 8 symbols
- 66% strings < 32 symbols
- 95% strings < 128 symbols

Any changes to improve the current state of things should not case regressions on this dataset



String s = new String("Can this work faster?");





StringCoding.hasNegatives()

```
@HotSpotIntrinsicCandidate
public static boolean hasNegatives(byte[] ba, int off, int len) {
    for (int i = off; i < off + len; i++) {
        if (ba[i] < 0) {
            return true;
        }
    }
    return false;
}
```

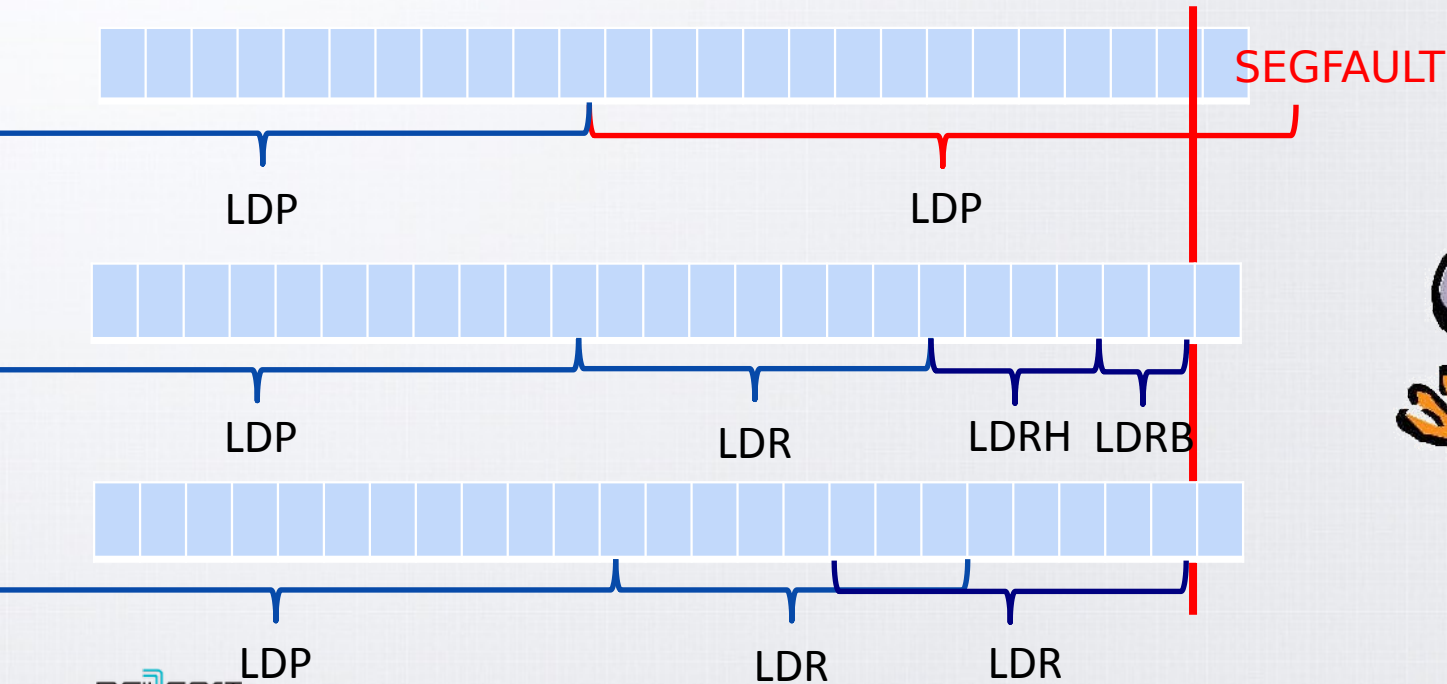
Some ARM assembly – memory reads



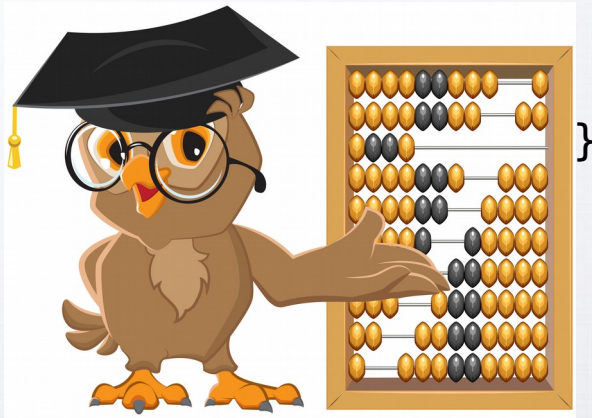
	Register	Width (bits)	Latency (cycles)
LDRB	GPR	8	4
LDRH	GPR	16	4
LDR	GPR	32 or 64	4
LDP	GPR	64+64	5

....

Learning to read (again)



And compare 8 bits at a time with 0

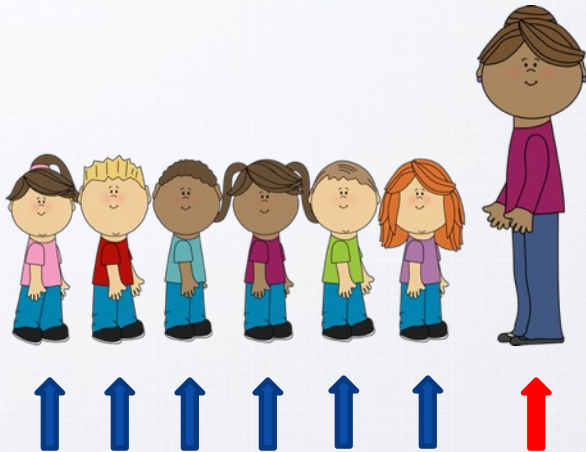


```
for(int i = off; i < off + len; i++) {
    if (ba[i] < 0) {
        return true;
    }
}
```

```
const uint64_t UPPER_BIT_MASK=0x8080808080808080;
...
__tst(rscratch2, UPPER_BIT_MASK);
```



Aligned memory access



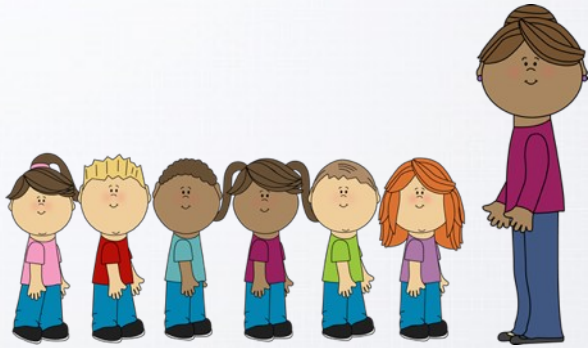
x86:

- in most cases modern processors do not have a penalty for unaligned memory access

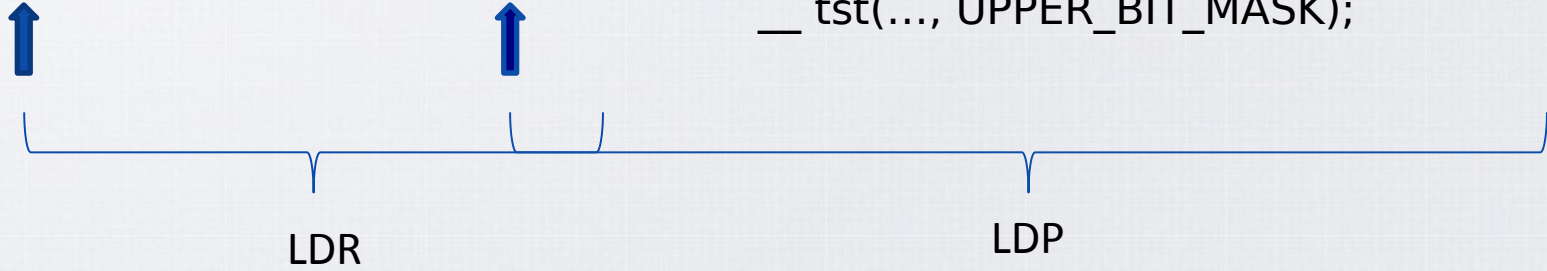
ARM is a spec:

- some CPU manufacturers do not have a penalty
- others do have (20%, 50%, 100%)

How to align memory access



```
// pre-loop
__ldr();
...
__tst(..., UPPER_BIT_MASK);
// main loop
__ldp(); //aligned
...
__tst(..., UPPER_BIT_MASK);
```





The plan for hasNegatives() intrinsic

- Read as much bytes at a time as possible, without crossing the page boundaries
 - If the page border is close
 - Read less bytes
 - Shift to the left
- Compare as many bytes with 0 as possible at a time
- Align memory access
- Reality
 - The code gets too big – 200 instructions
 - This interferes with inlining: C2 inlines up to 1500 instructions



Code is too big – what do we do?

- ARM ASM pseudo-code in Java that is short (27 instructions)
 - Not optimal, unaligned, but short

```
if (len > 32)
    return stubHasNegatives(ba, 0, len);
for (int i = 0; i < 32; i++) {
    if (ba[i] < 0) {           // ldr, tst
        return true;
    }
}
return stubHasNegatives(ba, 32, len); // ldp, tst
```

- The rest of the code goes to stub



What is a stub?

- A type of assembly inline in HotSpot
- Close analogy is a function
 - Can be called from macroAssembler
 - Code gets loaded during JVM startup once
 - Does not get inlined
- Several entry points are possible
- Some performance penalty calling stub



What should we place in stub?

```
// align memory access
__ bind(LARGE_LOOP); // 64 byte at a time
4x __ ldp(); //ary1, ary1+16, ary1+32, ary1+48
__ add(ary1, ary1, large_loop_size);
__ sub(len, len, large_loop_size);
7x __ orr(...);
__ tst(tmp2, UPPER_BIT_MASK);
__ br(Assembler::NE, RET_TRUE);
__ cmp(len, large_loop_size);
__ br(Assembler::GE, LARGE_LOOP);
```

OK, we helped C2. Can we help the hardware?



Software Prefetching

Let's give a processor a hint where we are going to read from memory next time:

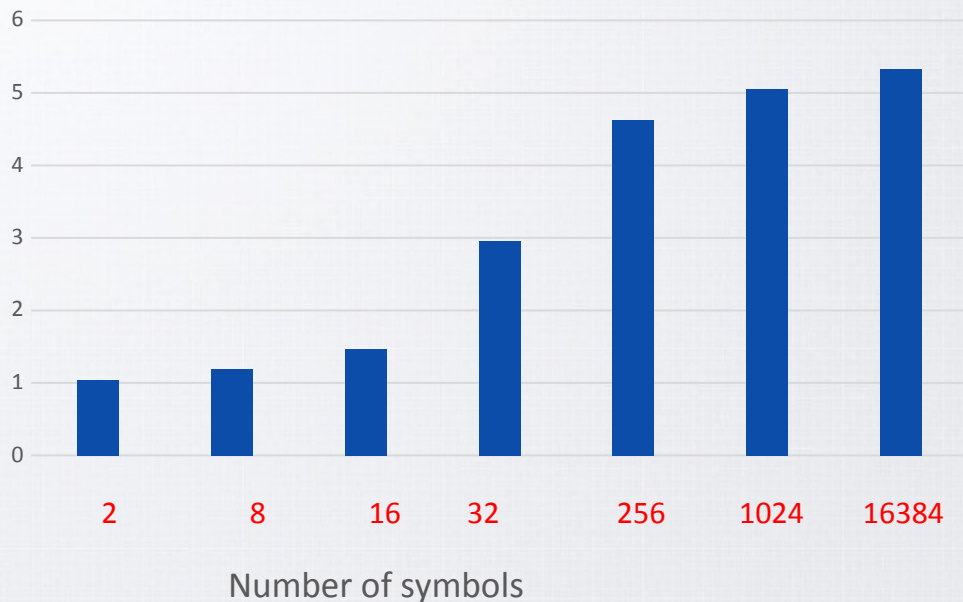
```
__prfm(Address(ary1, SoftwarePrefetchHintDistance));  
// do local register or operations on data in cache  
__ldp();
```

- Can be a major performance gain if
 - Processor has enough data to process between prfm and memory load
 - `SoftwarePrefetchHintDistance` is correctly defined:
> `d_cache_line_size`



Benchmark for new String() – long strings

Speedup compared to C2, times



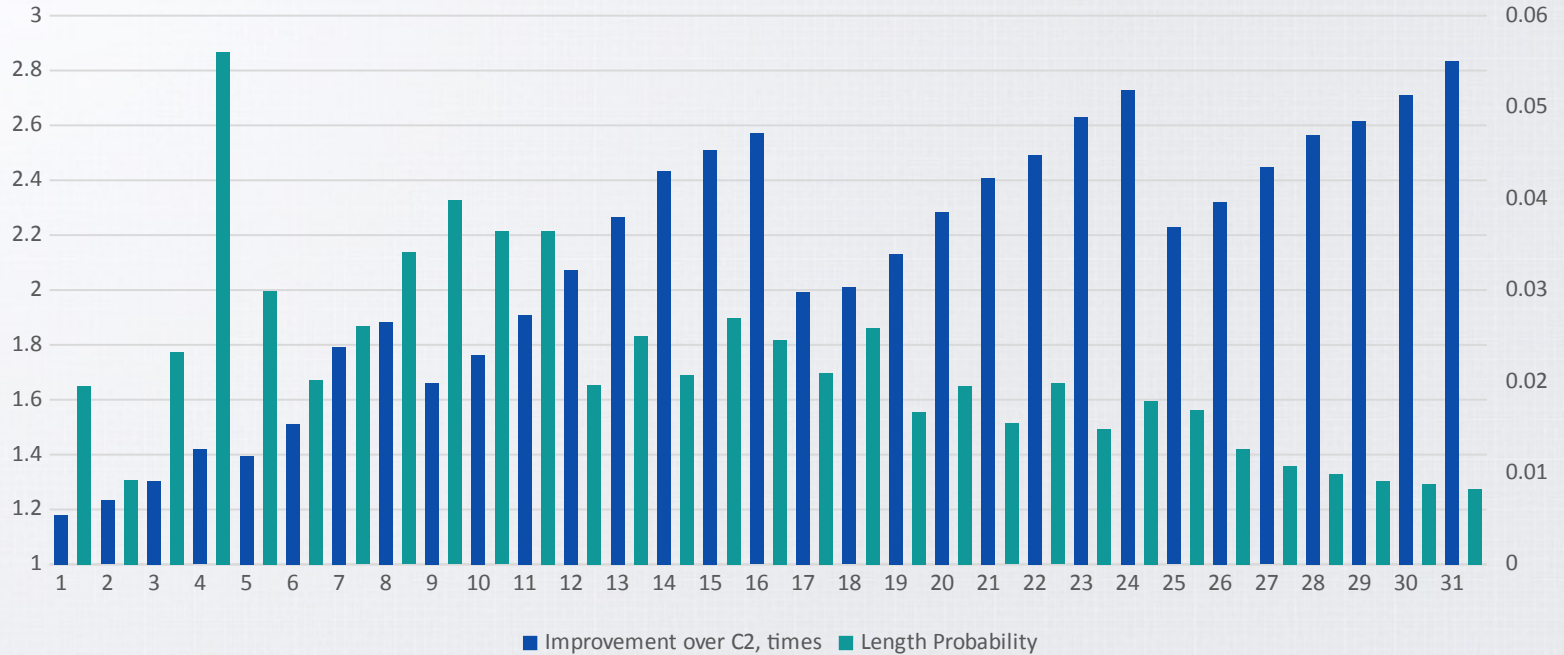
Longer string sizes experience more performance gain from optimization due to

- Optimal ldp & tst use
- Prefetching

Speedup up to 5x!



Benchmark for new String() – results



JEP 315: Improve Aarch64 Intrinsics

Workshop

[OpenJDK FAQ](#)
[Installing](#)
[Contributing](#)
[Sponsoring](#)
[Developers' Guide](#)

[Mailing lists](#)
[IRC](#) · [Wiki](#)

[Bylaws](#) · [Census](#)
[Legal](#)

JEP Process

Source code

[Mercurial](#)
[Bundles \(6\)](#)

Groups

[\(overview\)](#)
[2D Graphics](#)
[Adoption](#)
[AWT](#)
[Build](#)
[Compatibility & Specification](#)
[Review](#)
[Compiler](#)

Owner [Dmitrij Pochevko](#)

Type [Feature](#)

Scope [Implementation](#)

Status [Closed / Delivered](#)

Release [11](#)

Component [hotspot / compiler](#)

Discussion [hotspot dash compiler dash dev at openjdk dot java dot net](#)

Effort [L](#)

Duration [L](#)

Reviewed by [Mikael Vidstedt, Vladimir Kozlov](#)

Endorsed by [Vladimir Kozlov](#)

Created [2017/10/10 12:40](#)

Updated [2018/09/10 14:45](#)

Issue [8189104](#)

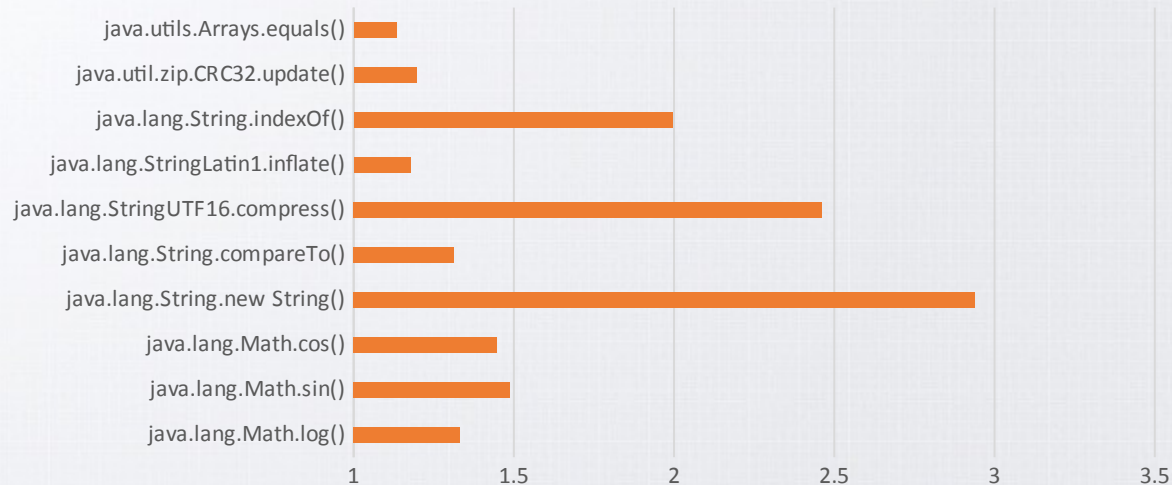
Summary

Improve the existing string and array intrinsics, and implement new intrinsics for the `java.lang.Math` `sin`, `cos` and `log` functions, on AArch64 processors.



Performance improvement

Average performance improvement*, times

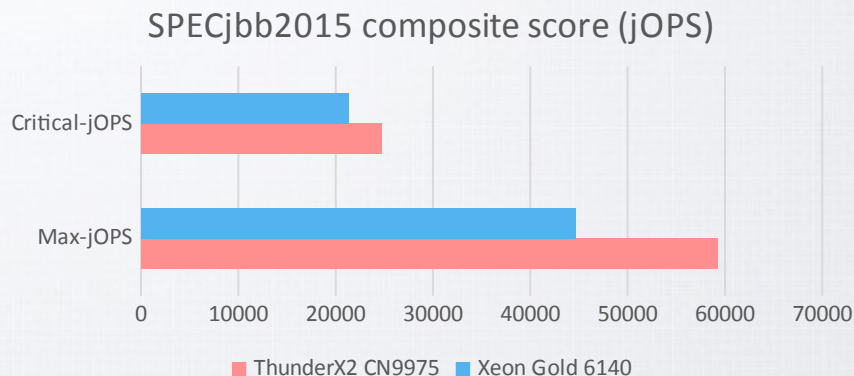


- Speedup up to 78x in microbenchmarks

* mean improvement over different size, length, encodings



JVM Benchmark #1 results

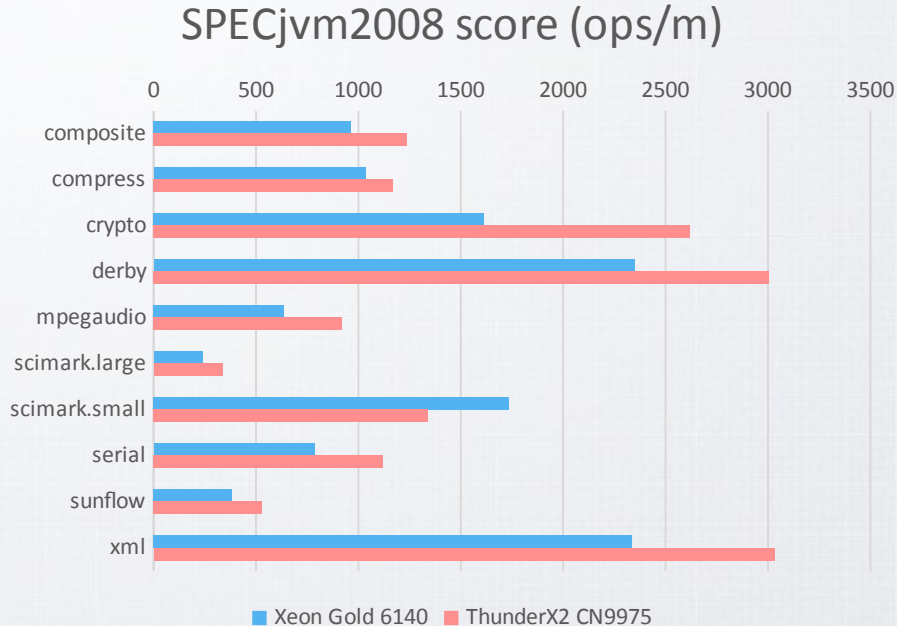


- Liberica JDK 11
- Average over 20 runs
- JEP 315 in JDK 11
- Cavium Thunder X2 outperforms Xeon 6140
 - by 33% in Max-jOPS score
 - by 16% in Critical-jOPS score

ARMv8: *-Xmx24G -Xms24G -Xmn16G -XX:+AlwaysPreTouch -XX:+UseParallelGC -XX:+UseTransparentHugePages -XX:-UseBiasedLocking*
X86: *-Xmx24G -Xms24G -Xmn16G -XX:+AlwaysPreTouch -XX:+UseParallelGC -XX:+UseTransparentHugePages -XX:+UseBiasedLocking*



JVM Benchmark #2 results



- Liberica JDK 11
- Default JVM settings
- Average over 20 runs
- Thunder X2 outperforms Xeon 6140
 - by 62% in Crypto
 - by 42% in MpegAudio
 - By 29% in XML
 - by 12% in Compress
- Xeon 6140 outperforms Thunder X2
 - By 29% in scimark.small



Where to try ARM servers?

packet

Bare Metal



scaleway

VPS



VPS



Conclusions

- Arm server vendors did a great job
- Cloud providers provide access to Arm servers right now
- Ubuntu, Red Hat, Oracle Linux, SuSE have ARMv8 support
- The software ecosystem just works as expected on ARMv8
- OpenJDK 11 is optimized for ARMv8

Download and install Liberica for ARMv8