



DPDK

USERSPACE, October 2016

Improving Driver Performance – A Worked Example

Bruce Richardson



Legal Disclaimers

No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document.

Intel disclaims all express and implied warranties, including without limitation, the implied warranties of merchantability, fitness for a particular purpose, and non-infringement, as well as any warranty arising from course of performance, course of dealing, or usage in trade.

This document contains information on products, services and/or processes in development. All information provided here is subject to change without notice. Contact your Intel representative to obtain the latest forecast, schedule, specifications and roadmaps.

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. **No computer system can be absolutely secure.** Check with your system manufacturer or retailer or learn more at intel.com.

© 2016 Intel Corporation. Intel, the Intel logo, Intel. Experience What's Inside, and the Intel. Experience What's Inside logo are trademarks of Intel. Corporation in the U.S. and/or other countries.

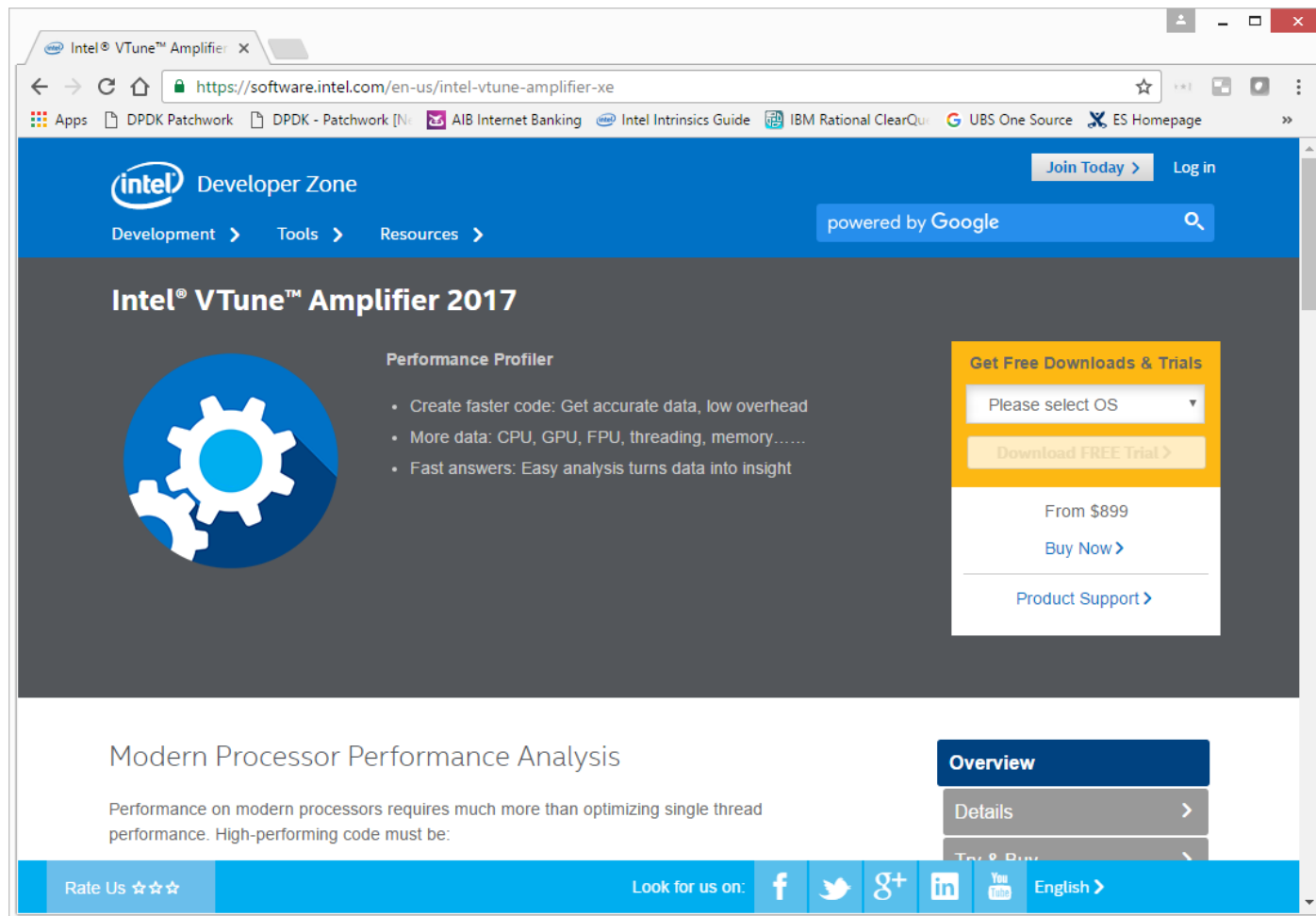
*Other names and brands may be claimed as the property of others.

Introduction

- This is a real-world use-case of performance improvement using performance analysis tools – in this case “Intel® VTune™ Amplifier”
- The issues discovered here were largely discovered by accident when investigating other DPDK behaviour
- Hope the walk-through of the issue debug may be of use to others when dealing with performance issues

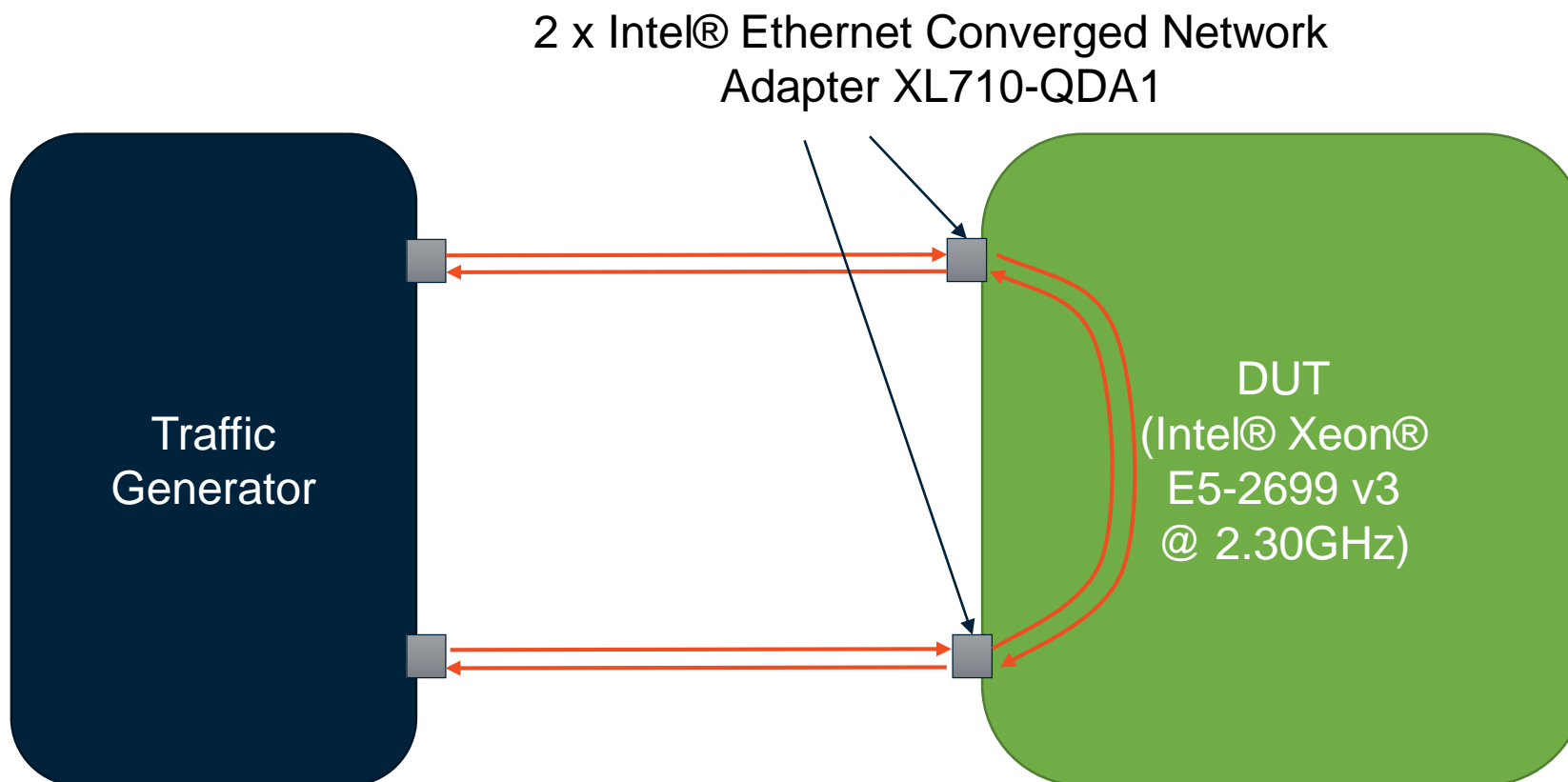
Intel® Vtune Amplifier

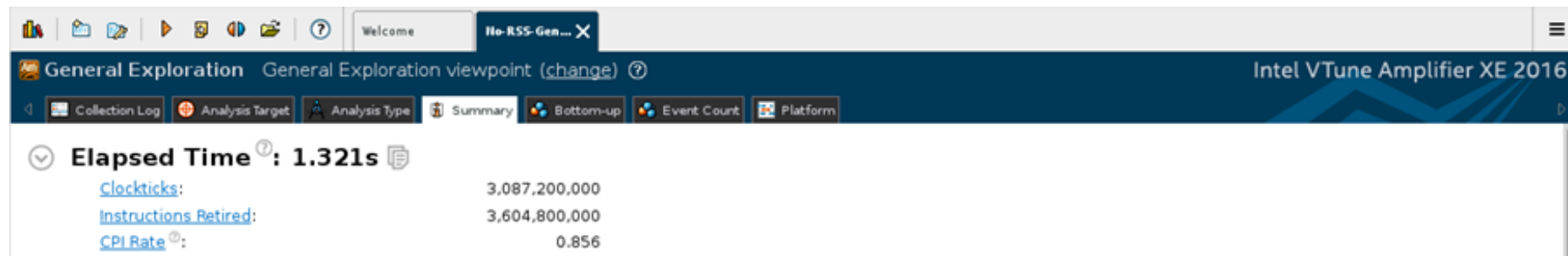
- Tool for performance analysis and debugging
- Uses hardware event counters to report on issues affecting program execution, e.g. CPU stalls due to memory access



The screenshot shows the Intel Developer Zone website for the Intel® Vtune™ Amplifier 2017. The page features a blue header with the Intel logo, navigation links for Development, Tools, and Resources, and a search bar. The main content area is dark grey and includes a large blue gear icon, a 'Performance Profiler' section with bullet points, and a yellow 'Get Free Downloads & Trials' box with a 'Download FREE Trial' button. Below this, there is a 'Modern Processor Performance Analysis' section and a sidebar with 'Overview' and 'Details' tabs. The footer contains social media links and a language selector.

Test Setup





⊖ **L1 Bound** [?]: **0.288**

This metric shows how often machine was stalled without missing the L1 data cache. The L1 cache typically has the shortest latency. However, in certain cases like loads blocked on older stores, a load might suffer a high latency even though it is being satisfied by the L1. Note that this metric value may be highlighted due to DTLB Overhead or Cycles of 1 Port Utilized issues.

[DTLB Overhead](#) [?]: 0.000

[Loads Blocked by Store Forwarding](#) [?]: **0.189**

Loads are blocked during store forwarding. To streamline memory operations in the pipeline, a load can avoid waiting for memory if a prior store, still in flight, is writing the data that the load wants to read (a 'store forwarding' process). However, in some cases, generally when the prior store is writing a smaller region than the load is reading, the load is blocked for a significant time pending the store forward. This metric measures the performance penalty of such blocked loads.

[Split Loads](#) [?]:

[4K Aliasing](#) [?]:

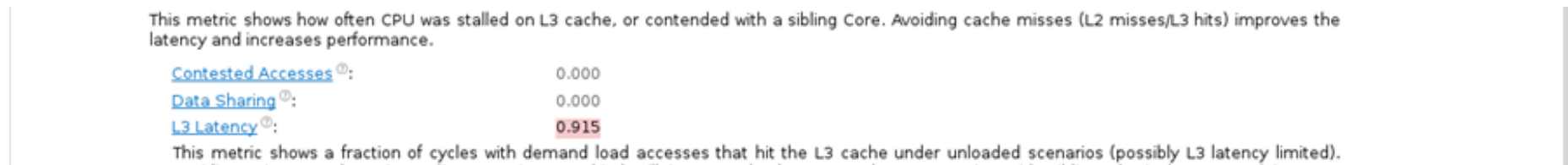
[FB Full](#) [?]:

This metric does a rough estimation of the number of additional demand L1D demand requests to proceed.

[L2 Bound](#) [?]: 0.000

⊖ **L3 Bound** [?]: **0.179**

This metric shows how often CPU was stalled on L3 cache, or contended with a sibling Core. Avoiding cache misses (L2 misses/L3 hits) improves the latency and increases performance.



What does this mean?

- We are doing a write to a memory location
- We are subsequently doing a read from that memory location
- The read is getting blocked by the write because the read is for more data than just the write
 - if the write is for the same amount of data, then we can do “store forwarding” to return the data write to the read without hitting cache/mem
 - if the read can be delayed till later, the write can complete and the read can come from cache
- This should not generally show up as a problem in good code

Intel VTune Amplifier XE 2016

General Exploration General Exploration viewpoint (change) ?

Analysis Target Analysis Type Collection Log Summary Bottom-up Event Count Platform i40a_rxtx_... i40a_rxtx_...

Source Assembly

So. Ln.	Source	Clo.	Instr. Retr.	CPI Rate	Fr. Bo.	Ba. Sp.	DTLB Over.	Address	Sour. Line	Assembly	Pr. Bte.	Fr. Bo.	Bad Spe.	DTLB Over.	L1 Bound	Loads Bloc. by 5...	Split Loads	4...	
326	/* B.1 load 2 mbuf point */							0x52d49f	Block 6:										
327	descs[1] = _mm_loadu_si128((__m128i *) (rxdp + 1));	1.20	200	6.000	0.0%	0.0%	0.000	0x52d4a3	341	vpunpckhqdq %xmm1, %xmm0, %xmm1									
328	descs[0] = _mm_loadu_si128((__m128i *) (rxdp));	1.00	200	5.000	0.0%	50.0%	0.000	0x52d4a7	341	vpunpckhqdq %xmm3, %xmm2, %xmm0									
329								0x52d4ab	341	vpunpckhqdq %xmm0, %xmm1, %xmm0									
330	/* B.2 copy 2 mbuf point into rx_pkts */							0x52d4ad	341	vpsrlq \$0x6, %xmm0, %xmm0									
331	_mm_storeu_si128((__m128i *) &rx_pkts[pos+2], mbp2);							0x52d4b0	341	vpand %xmm0, %xmm15, %xmm0									
332								0x52d4b4	341	vpackssdq %xmm8, %xmm0, %xmm0									
333	if (split_packet) {	1.20	1.00	1.200	0.0%	0.0%	0.000	0x52d4b9	341	vsovbq %xmm0, %rdx									
334	rte_prefetch0(&rx_pkts[pos]->cacheline1);							0x52d4be	341	mov %rdx, %r10									
335	rte_prefetch0(&rx_pkts[pos + 1]->cacheline1);							0x52d4c1	341	movb %dx, -0x3a(%rsp)									
336	rte_prefetch0(&rx_pkts[pos + 2]->cacheline1);							0x52d4c6	341	shr \$0x10, %r10									
337	rte_prefetch0(&rx_pkts[pos + 3]->cacheline1);							0x52d4ca	341	movb %r10v, -0x2a(%rsp)									
338	}							0x52d4d0	341	mov %rdx, %r10									
339	}							0x52d4d3	341	shr \$0x30, %rdx									
340	/*shift the pktlen field*/							0x52d4d7	341	shr \$0x20, %r10									
341	desc_pktlen_align(descs);							0x52d4db	341	movb %dx, -0xa(%rsp)									
342								0x52d4e0	341	movb %r10v, 0x1a(%rsp)									
343	/* avoid compiler reorder optimization */							0x52d4e6	347	vmovdqax -0x18(%rsp), %xmm11	189	0.0%	15.2%	0.000	0.000	0.000	0.000	0.000	
344	rte_compiler_barrier();							0x52d4ec	348	vmovdqax -0x28(%rsp), %xmm1	208	0.4%	0.0%	0.000	0.389	0.000	0.000	0.000	0.000
345								0x52d4f2	35	vmovdqax -0x48(%rsp), %xmm3	348	0.0%	9.3%	0.000	0.000	0.000	0.000	0.000	0.000
346	/* D.1 pkt 3,4 convert format from desc to pktmbuf */							0x52d4f8	37										
347	pkt_mb4 = _mm_shuffle_epi8(descs[3], shuf_msk);	571	178	1.208	0.4%	0.0%	0.000	0x52d4fd	370	vpaddq %xmm6, %xmm7, %xmm6									
348	pkt_mb3 = _mm_shuffle_epi8(descs[2], shuf_msk);							0x52d501	353	vmovdqax -0x38(%rsp), %xmm4	300	0.0%	0.0%	0.000	1.000	0.000	0.000	0.000	0.000
349								0x52d507	341	vpslufb %xmm9, %xmm1, %xmm5									
350	/* C.1 4->2 filter staterr info only */							0x52d50a	35	vpunpckhqdq %xmm1, %xmm11, %xmm0									
351	sterr_tmp2 = _mm_unpackhi_epi32(descs[3], descs[2]);							0x52d510	370	vpunpckhqdq %xmm11, %xmm1, %xmm1									
352	/* C.1 4->2 filter staterr info only */							0x52d515	370	movb -0x18(%r10), %r10									
353								0x52d51a	370										
354								0x52d51d	370										
355	/* C.2 get 4 pkts staterr value */							0x52d520	370										
356	zero = _mm_xor_si128(dd_check, dd_check);							0x52d525	355	vpsrlub %xmm2, %xmm12, %xmm2									
357	staterr = _mm_unpacklo_epi32(sterr_tmp1, sterr_tmp2);							0x52d52a	355	vpor %xmm2, %xmm1, %xmm2									
358								0x52d52f	355	vsovbq %xmm2, %rdx									
359	/* D.3 copy final 3,4 data to rx_pkts */							0x52d534	363	vpslufb %xmm9, %xmm3, %xmm3									
360	_mm_storeu_si128((void *) &rx_pkts[pos+3]->rx_descriptor_fields1,							0x52d539	376	vpaddq %xmm4, %xmm7, %xmm4									
361	pkt_mb4);							0x52d53e	355	movzx %dx, %ebp									
362								0x52d543	377	vpaddq %xmm3, %xmm7, %xmm3									
363								0x52d548	355	movq %rbp, 0x18(%r10)									
364								0x52d54d	355	movq %rsi, %rbp									
365	/* C.2 get 4 pkts staterr value */							0x52d552	355	mov %rdx, %r10									
366	zero = _mm_xor_si128(dd_check, dd_check);							0x52d557	376										
367	staterr = _mm_unpacklo_epi32(sterr_tmp1, sterr_tmp2);							0x52d55c	355										
368								0x52d561	355										
369	/* D.3 copy final 3,4 data to rx_pkts */							0x52d566	355										
370	_mm_storeu_si128((void *) &rx_pkts[pos+3]->rx_descriptor_fields1,							0x52d56b	355										
371	pkt_mb4);							0x52d570	355										

Selected 1 row(s): 571, 473, 1.208, 0.4%, 0.0%, 0.000

Highlighted 2 row(s): 1, 0.4%, 0.0%, 0.000, 0.389, 0.000, 0.000

```

/* D.1 pkt 3,4 convert format from desc to pktmbuf *
pkt_mb4 = _mm_shuffle_epi8(descs[3], shuf_msk);
pkt_mb3 = _mm_shuffle_epi8(descs[2], shuf_msk);

```

```

vmovdqax -0x18(%rsp), %xmm11
vmovdqax -0x28(%rsp), %xmm1
vmovdqax -0x48(%rsp), %xmm3

```


Reading this case

- In this case, both the instruction highlighted and the previous one are 128-bit loads.
- Therefore either one is a potential candidate for the source of the delay
- If we assume that this is the read, then we need to find the offending write:
 - occurs previous to these [nice and easy to find in C, as there is a compiler barrier]
 - is of a size smaller than 128-bits

Other Weirdness...

- Why does the code:

```
pkt_mb3 = _mm_shuffle_epi8(descs[2], shuf_msk);
```

cause a read from memory at all?

- Shouldn't descs[2] have already been loaded to xmm register previously at the line?

```
descs[2] = _mm_loadu_si128((__m128i *) (rxdp + 2));
```

- Let's look at that previous load lines in vtune...

```
311 __m128i mbp1, mbp2, /* two mbuf pointer in one XMM reg. */
312
313 /* B.1 load 1 mbuf point */
314 mbp1 = _mm_loadu_si128((__m128i *)&sv_ring[pos]);
315 /* Read desc statuses backwards to avoid race condition */
316 /* A.1 load 4 pkts desc */
317 desc[3] = _mm_loadu_si128((__m128i *)(rxdp + 3));
318
319 /* B.2 copy 2 mbuf point into rx_pkts */
320 _mm_storeu_si128((__m128i *)&rx_pkts[pos], mbp1);
321
322 /* B.1 load 1 mbuf point */
323 mbp2 = _mm_loadu_si128((__m128i *)&sv_ring[pos+2]);
324
325 desc[2] = _mm_loadu_si128((__m128i *)(rxdp + 2));
326 /* B.1 load 2 mbuf point */
```

```
0x52d41c 305 xor %eax, %eax
0x52d41e 305 vmovdqax 0xd582a(%rip), %xmm12
0x52d426 240 vpinsrw $0x1, %eax, %xmm0, %xmm0
0x52d42b 240 vpunpckldq %xmm0, %xmm8, %xmm7
0x52d42f 240 vpunpckldq %xmm8, %xmm0, %xmm0
0x52d434 240 vpunpckldq %xmm0, %xmm7, %xmm7
0x52d438      Block 4:
0x52d43d 314 vmovdqax (%r11), %xmm0
0x52d43d 317 vmovdqax 0x60(%r9), %xmm3
0x52d443 317 vmovapsx %xmm3, -0x18(%rsp)
0x52d449 320 vmovupsx %xmm0, -0x8(%r11)
0x52d44e 325 vmovdqax 0x40(%r9), %xmm1
0x52d454 323 vmovdqax 0x10(%r11), %xmm4
0x52d45a 325 vmovapsx %xmm1, -0x28(%rsp)
0x52d460 32 vmovdqax 0x20(%r9), %xmm2
0x52d466 327 vmovapsx %xmm2, -0x38(%rsp)
```

`descs[3] = _mm_loadu_si128((__m128i *)(rxdp + 3))`

```
vmovdqax 0x60(%r9), %xmm3
vmovapsx %xmm3, -0x18(%rsp)
```

What is happening?

- A load intrinsic is resulting in an xmm load followed by a store?
- We have a second mystery.
- Let's trace back through what happens to the descriptors through the code between the two points...

desc_pktlen_align

- Only work done between desc[2] load and the offending line is function “desc_pktlen_align”
- Again look at assembler listing

Intel VTune Amplifier XE 2016

General Exploration viewpoint (change)

Analysis Target: Analysis Type: Collection Log: Summary: Bottom-up: Event Count: Platform: i40e_rxtx_v... i40e_rxtx_v...

Source Assembly

So. Lin.	Source	Cl.	Instr. Retr...	CPI Rate	Fr. Bo.	Ba. Sp.	DTLB Over	Address	Sour. Line	Assembly	PI Re	Fr. Bo.	Bad Spe.	DTLB Over	Loads Bloc. by S	Split Loads	L1 Bound
302	* D. fill info. from desc to mbuf							0x52d489	335	movq (%rsi), %rdx							
303	*/							0x52d48c	335	prefetch0z 0x40(%rdx)							
304								0x52d490	336	movq 0x8(%rsi), %rdx							
305	for (pos = 0, nb_pkts_rcvd = 0; pos < RTE_I40E_VPMD_RX_BURST;	2.80	200	14.000	0.0%	0.0%	0.000	0x52d494	336	prefetch0z 0x40(%rdx)							
306	pos += RTE_I40E_DESCS_PER_LOOP,							0x52d498	337	movq (%r8), %rdx							
307	rxdp += RTE_I40E_DESCS_PER_LOOP) {	7.60	0	0.0%	0.0%	0.000		0x52d49b	337	prefetch0z 0x40(%rdx)							
308	__m128i desc[RTE_I40E_DESCS_PER_LOOP];							0x52d49f	341	vpurpckhdq %xmm1, %xmm0, %xmm1							
309	__m128i pkt_mb1, pkt_mb2, pkt_mb3, pkt_mb4;							0x52d4a3	341	vpurpckhdq %xmm3, %xmm2, %xmm0							
310	__m128i zero, staterr, sterr_tmp1, sterr_tmp2;							0x52d4a7	341	vpurpckhdq %xmm0, %xmm1, %xmm0							
311	__m128i mbp1, mbp2; /* two mbuf pointer in one XMM reg. */							0x52d4ab	341	vpsrlq \$0x6, %xmm0, %xmm0							
312								0x52d4b0	341	vpsrad %xmm0, %xmm15, %xmm0							
313	/* B.1 load 1 mbuf point */							0x52d4b4	341	vpackssdq %xmm0, %xmm0, %xmm0							
314	mbp1 = __m_loadu_si128((__m128i *)&sw_ring[pos]);							0x52d4b9	341								
315	/* Read desc statuses backwards to avoid race condition */							0x52d4be	341								
316	/* A.1 load 4 pkts desc */							0x52d4c1	341	movq %dx, -0x2a(%rsp)							
317	descs[3] = __m_loadu_si128((__m128i *)&(rxdp + 3));	3.00	1.00	1.875	0.0%	0.0%	0.000	0x52d4c6	341	shr \$0x10, %r10							
318								0x52d4ca	341	movq %r10v, -0x2a(%rsp)							
319	/* B.2 copy 2 mbuf point into rx_pkts */							0x52d4cd	341	mov %rdx, %r10							
320	__m_storeu_si128((__m128i *)&rx_pkts[pos], mbp1);							0x52d4d3	341	shr \$0x30, %rdx							
321								0x52d4d7	341	shr \$0x20, %r10							
322	/* B.1 load 1 mbuf point */							0x52d4db	341	movq %dx, -0xa(%rsp)							
323	mbp2 = __m_loadu_si128((__m128i *)&sw_ring[pos+2]);							0x52d4de	341	movq %r10v, -0x1a(%rsp)							
324								0x52d4e6	347	vmsvdqax -0x10(%rsp), %xmm11	208	0.4%	0.0%	0.000	0.389	0.000	0
325	descs[2] = __m_loadu_si128((__m128i *)&(rxdp + 2));	1.40	600	2.883	0.0%	0.0%	0.000	0x52d4ec	348	vmsvdqax -0x20(%rsp), %xmm1	348	0.0%	3.3%	0.000	0.000	0.000	0
326	/* B.1 load 2 mbuf point */							0x52d4f2	353	vmsvdqax -0x40(%rsp), %xmm3							
327	descs[1] = __m_loadu_si128((__m128i *)&(rxdp + 1));	1.20	200	6.000	0.0%	0.0%	0.000	0x52d4f8	347	vpshufb %xmm0, %xmm11, %xmm6							
328	descs[0] = __m_loadu_si128((__m128i *)&(rxdp));	1.00	200	5.000	0.0%	30.0%	0.000	0x52d4fd	370	vpaddq %xmm6, %xmm7, %xmm6							
329								0x52d501	353	vmsvdqax -0x30(%rsp), %xmm4	200	0.0%	0.0%	0.000	1.000	0.000	0
330	/* B.2 copy 2 mbuf point into rx_pkts */							0x52d507	348	vpshufb %xmm0, %xmm1, %xmm5							
331	__m_storeu_si128((__m128i *)&rx_pkts[pos+2], mbp2);							0x52d50c	351	vpurpckhdq %xmm1, %xmm11, %xmm0							
332								0x52d510	355	vpurpckhdq %xmm1, %xmm1, %xmm1							
333	if (!split_packet) {	1.20	1.00	1.300	0.0%	0.0%	0.000	0x52d515	355	movq -0x8(%rsi), %r10							
334	rte_prefetch0(&rx_pkts[pos->cacheline]);							0x52d519	372	vpaddq %xmm5, %xmm7, %xmm5							
335	rte_prefetch0(&rx_pkts[pos + 1->cacheline]);							0x52d51d	355	vpurpckhdq %xmm4, %xmm3, %xmm2							
336	rte_prefetch0(&rx_pkts[pos + 2->cacheline]);							0x52d521	355	vpurpckldq %xmm1, %xmm2, %xmm2							
337	rte_prefetch0(&rx_pkts[pos + 3->cacheline]);							0x52d525	355	vpsrad %xmm2, %xmm14, %xmm2							
338	}							0x52d529	353	vpurpckhdq %xmm3, %xmm4, %xmm10							
339								0x52d52d	355	vpshufb %xmm2, %xmm13, %xmm1							
340	/*shift the pktlen field*/							0x52d532	355	vpsrlq \$0xc, %xmm2, %xmm2							
341	desc_pktlen_align(descs);							0x52d537	362	vpshufb %xmm0, %xmm4, %xmm4							
342								0x52d53c	367	vpurpckldq %xmm0, %xmm10, %xmm0							
343	/* avoid compiler reorder optimization */							0x52d540	355	vpshufb %xmm2, %xmm12, %xmm2							
344	rte_compiler_barrier();							0x52d545	355	vpor %xmm2, %xmm1, %xmm2							
345								0x52d549	355	vmsvdqax %xmm2, %rdx							
346	/* D.1 pkt 3,4 convert format from desc to pktsbuf */							0x52d54e	363	vpshufb %xmm0, %xmm3, %xmm3							
347	pkt_mb4 = __m_shuffle_epi8(descs[3], shuf_mask);	5.20	1.80	2.883	0.0%	15.2%	0.000										

Selected 1 row(s):

Highlighted 17 row(s):

Vector Code

Scalar Code [Including writes]

The Source of the problem

- When assigning the lengths, we use 16-bit writes:
 - have to go to memory (can't assign to an xmm register)
 - causes the xmm load to immediately store to stack
 - causes the shuffle op to trigger a second load
- That second load (128b) blocks on 16b write

```
pktlen0 = _mm_srli_epi32(pktlen0, PKTLEN_SHIFT);
pktlen0 = _mm_and_si128(pktlen0, pktlen_msk);

pktlen0 = _mm_packs_epi32(pktlen0, zero);
vol.dword = _mm_cvtsi128_si64(pktlen0);

/* let the descriptor byte 15-14 store the pkt len */
*((uint16_t *)&descs[0]+7) = vol.e[0];
*((uint16_t *)&descs[1]+7) = vol.e[1];
*((uint16_t *)&descs[2]+7) = vol.e[2];
*((uint16_t *)&descs[3]+7) = vol.e[3];
```


The Fix

Rewrite to use entirely vector operations:

- keeps things in xmm registers
- saves unnecessary loads and stores
- prevents store forward errors
- improves performance.
- makes people happy*

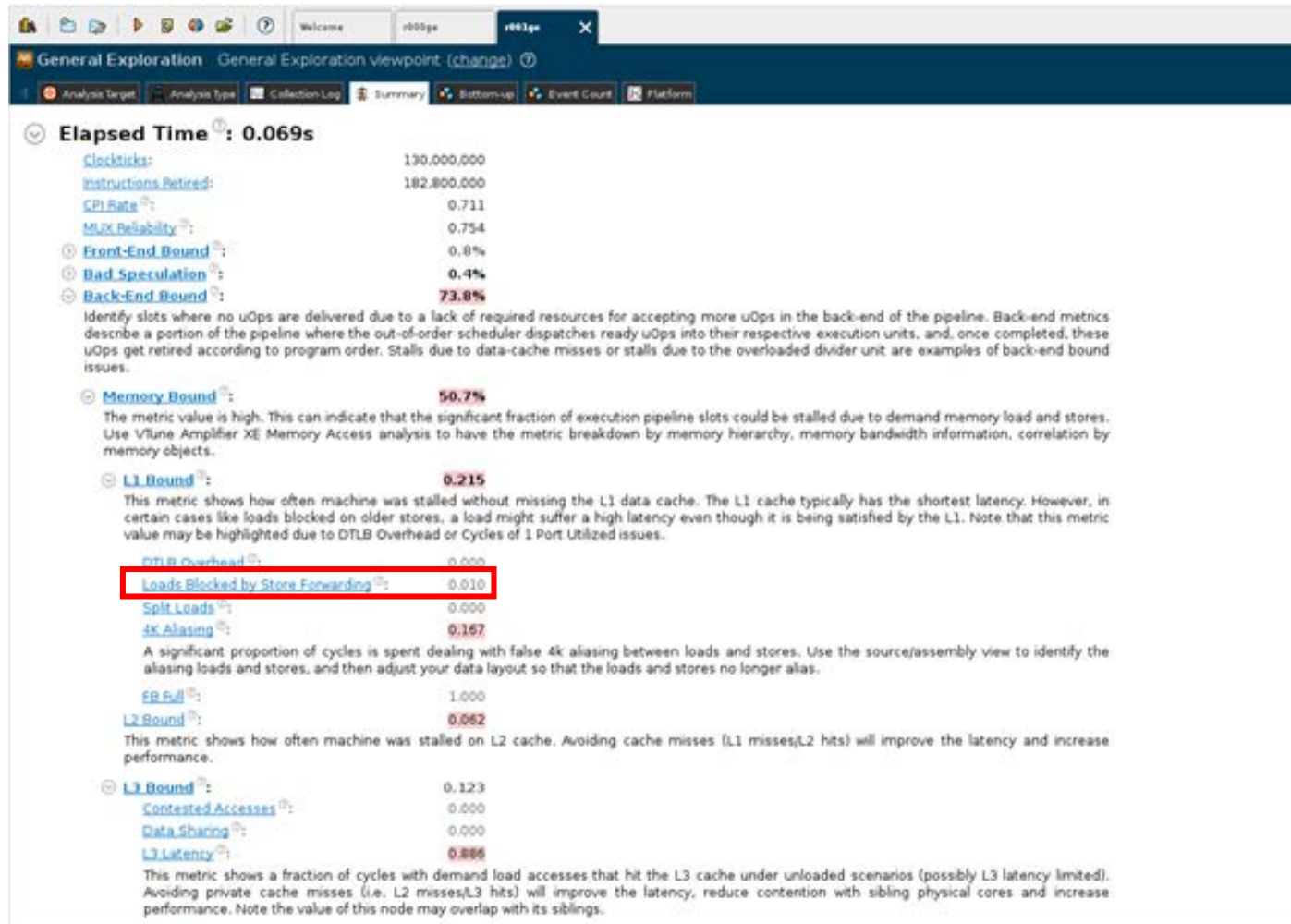
*NOTE: happiness not guaranteed

```
pktlen0 = _mm_srli_epi32(pktlen0, PKTLEN_SHIFT);
pktlen0 = _mm_and_si128(pktlen0, pktlen_msk);

pktlen0 = _mm_packs_epi32(pktlen0, pktlen0);

descs[3] = _mm_blend_epi16(descs[3], pktlen0, 0x80);
pktlen0 = _mm_slli_epi64(pktlen0, 16);
descs[2] = _mm_blend_epi16(descs[2], pktlen0, 0x80);
pktlen0 = _mm_slli_epi64(pktlen0, 16);
descs[1] = _mm_blend_epi16(descs[1], pktlen0, 0x80);
pktlen0 = _mm_slli_epi64(pktlen0, 16);
descs[0] = _mm_blend_epi16(descs[0], pktlen0, 0x80);
```

Result



- Stalls due to Loads Blocked by Store Forwarding dropped about 19x:
 - Before: 0.189
 - After: 0.010
- Overall PMD performance measured by testpmd increased by over 5%.