

# Performance Monitoring Impact of Intel® Transactional Synchronization Extension Memory Ordering Issue

White paper

June 2021

**Revision 1.4** 

Document Number: 604224



Notice: 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 Intel product specifications and roadmaps. Do not finalize a design with this information.

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software, or service activation. Learn more at intel.com, or from the OEM or retailer.

No computer system can be absolutely secure. Intel does not assume any liability for lost or stolen data or systems or any damages resulting from such losses.

You may not use or facilitate the use of this document in connection with any infringement or other legal analysis concerning Intel products described herein. You agree to grant Intel a non-exclusive, royalty-free license to any patent claim thereafter drafted which includes subject matter disclosed herein.

No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document. The products described may contain design defects or errors known as errata which may cause the product to deviate from published specifications. Current characterized errata are available on request.

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.

Copies of documents which have an order number and are referenced in this document may be obtained by calling 1-800-548-4725 or by visiting <a href="https://www.intel.com/design/literature.htm">www.intel.com/design/literature.htm</a>.

Intel, the Intel logo, Intel® Xeon Phi™, Intel Atom, and Xeon 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

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. All Rights Reserved.



### **Contents**

1	Introduction			
	1.1 Implications for users	5		
	1.2 Implications for PMU drivers and performance tools			
	1.3 Implications for Intel TSX library developers	7		
2	TSX Disable Update			
	2.1 Unsupported software development mode			
	2.2 Updated MSR definition and affected products			
Α	Enabling FORCE_ABORT_RTM Mode to Use All Four Counte	rs11		
В	Guidance on Specific Profilers	13		
	B.1 Linux* perf			
	B.2 Processor Counter Monitor			
<b>Tables</b>				
	Table A-1. Description of TSX_FORCE_ABORT_MSR	11		

# **Revision History**

Document Number	Revision Number	Description	Date
604224	1.0	Initial release of the document.	October 2018
604224	1.1	Update to RTM disable in SGX and SMM modes	October 2018
604224	1.2	<ul> <li>Document RTM retry bit behavior. Perf updates.</li> <li>PMU counter CPUID enumeration changes. Now number of general purpose counters</li> </ul>	January 2019
604224	1.3	Additional details on re-enabling TSX in Linux	March 2019
604224	1.4	Document client TSX changes	June 2021

#### 1 Introduction

This whitepaper describes Intel® Transactional Synchronization Extension (Intel® TSX) and Performance Monitoring Unit (PMU) behavior due to the updated microcode for Intel® Xeon® D (code name Skylake-D), Intel® Xeon® Scalable Processor, and certain Intel® Xeon® Processor E3 v5 and v6 Family (code name Skylake and Kaby Lake) and 6th, 7th, and 8th Generation Intel® Core™ i7 and i5 (code name Skylake, Kaby Lake, Coffee Lake, and Whiskey Lake).

Note that for the affected Intel® Xeon® Processor E3 v5 and v6 Family (code name Skylake and Kaby Lake) and the 6th, 7th, and 8th Generation Intel® Core™ i7 and i5 (code name Skylake, Kaby Lake, Coffee Lake, and Whiskey Lake), a newer microcode update will be released in 2021.1 IPU that will disable Intel TSX by default. For these processors, the behavior documented in this section will be superseded by the behavior documented in section 2.

Intel TSX is a technology to enable hardware transactional memory. Intel TSX provides two software interfaces – Hardware Lock Elision (HLE) and Restricted Transactional Memory (RTM). HLE is an instruction prefix-based interface designed to be backward compatible with processors without Intel TSX support. RTM is a new instruction set interface using the XBEGIN and XEND instructions. For more details on Intel TSX please see <a href="http://www.intel.com/software/tsx">http://www.intel.com/software/tsx</a>.

The PMU measures performance events using performance counters. With the microcode update described in the TSX Memory Ordering Issue disclosure, released in October 2018, general purpose (GP) counters became available to the PMU driver, but the fourth performance counter may contain unexpected values. The October 2018 microcode update also disabled the HLE instruction prefix of Intel TSX and force all RTM transactions to abort when operating in Intel SGX mode or System Management Mode (SMM).

Intel does not expect these microcode updates to affect users who do not use the PMU, or who only use updated PMU drivers and tools. However, we recommend that PMU driver developers and performance tool developers follow the guidance in this document. Some advanced users of performance monitoring (Perfmon) may need to change their collection scripts and methodologies. The purpose of this whitepaper is to enable Perfmon users and tool developers to understand and, if necessary, work around the implications of these changes.

#### 1.1 Implications for users

The microcode update (CPUID.07H.EDX[bit 13]=1) is not expected to impact users who do not utilize Perfmon or HLE. We recommend that all users who do use Perfmon to update their PMU profiling tools to the latest version. Refer to Appendix B Guidance on Specific Profilers and the PMU tools documentation for more information on specific tools. No further action is required for PMU users who do not use groups.

Performance tools often use event multiplexing to collect data using more events than the number of available GP counters in the CPU. In a typical user scenario, such

as Microarchitecture Exploration Analysis Type in Intel® VTune™ Amplifier, there are predefined tool configurations, profiles, or scripts that can specify the event groupings. The primary impact to users in this scenario is that the GP counter collection groups would be split into three events each instead of four events.¹ Updated versions of these profiling tools (see Appendix B Guidance on Specific Profilers) automatically handle this change.

Some advanced users who choose to develop their own event groupings in collection methodologies or scripts will need to modify their input to ensure they only utilize the number of available counters for each counter grouping, or use the method described in Appendix A Enabling FORCE\_ABORT\_RTM Mode to Use All Four Counters.

# 1.2 Implications for PMU drivers and performance tools

For more details on the PMU, refer to the Software Developer's Manual (<a href="http://www.intel.com/sdm">http://www.intel.com/sdm</a>) Volume 3, Chapter 18 "Performance Monitoring".

When Restricted Transactional Memory (RTM) is supported (CPUID.07H.EBX.RTM [bit 11] = 1) and CPUID.07H.EDX[bit 13]=1 and

TSX\_FORCE\_ABORT[RTM\_FORCE\_ABORT]=0 (described later in this document), then Performance Monitor Unit (PMU) general purpose counter 3 (IA32\_PMC3, MSR C4H and IA32\_A\_PMC3, MSR 4C4H) may contain unexpected values. Specifically, IA32\_PMC3 (MSR C4H), IA32\_PERF\_GLOBAL\_CTRL[3] (MSR 38FH) and IA32\_PERFEVTSEL3 (MSR 189H) may contain unexpected values, which also affects IA32\_A\_PMC3 (MSR 4C4H) and IA32\_PERF\_GLOBAL\_INUSE[3] (MSR 392H). PMU driver should avoid using general purpose counter 3. General purpose counters beyond 3, if reported in CPUID.(EAX=0xA).EAX[15:08], can still be used. Using counter 3 will result in nondeterministic counting, especially in the presence of RTM transactions; however, this should not crash the PMU driver.

When supporting event multiplexing, the PMU driver needs to split the event list into the correct configuration and groups based on the number of available GP counters. Also, any tool configurations or scripts which have hard-coded specific groups of counters must be changed to support the possibility of having fewer counters available.

New versions of the PMU driver tools can add an option to gain use of all GP counters by enabling FORCE\_ABORT\_RTM mode during the measurement (see Appendix A Enabling FORCE\_ABORT\_RTM Mode to Use All Four Counters and Appendix B Guidance on Specific Profilers).

Table 1.2-1 Affected products if Intel TSX is supported

Family-Model	Stepping	Processor Families / Processor Number Series
06_55H	<=5	First generation Intel® Xeon® Scalable Processor Family and Intel® Xeon® Processor D Family based on Skylake microarchitecture

Family-Model	Stepping	Processor Families / Processor Number Series
06_4EH, 06_5EH	All	6th generation Intel <sup>®</sup> Core <sup>™</sup> processors and Intel <sup>®</sup> Xeon <sup>®</sup> processor E3-1500m v5 product family and E3- 1200 v5 product family based on Skylake microarchitecture
06_8EH	<=0xB	7th/8th generation Intel <sup>®</sup> Core <sup>™</sup> processors and Intel <sup>®</sup> Pentium <sup>™</sup> processors based on Kaby Lake/Coffee Lake/Whiskey Lake microarchitecture
06_9EH	<=0xC	8th/9th generation Intel <sup>®</sup> Core <sup>™</sup> processors and Intel <sup>®</sup> Pentium <sup>™</sup> processors based on Coffee Lake microarchitecture

#### 1.3 Implications for Intel TSX library developers

Libraries using RTM transactions often check the return (or abort) value of \_xbegin() to decide when and how often to retry transactions. Normally it is beneficial to retry transactions on a conflict abort. For a normal conflict abort the \_XBEGIN\_CONFLICT and \_XBEGIN\_RETRY bits are in the abort value. With CPUID.07H.EDX[bit 13]=1, it is possible to see conflict aborts that only have the \_XBEGIN\_CONFLICT bit set. These should be handled like normal conflicts for best performance.

With CPUID.07H.EDX[bit 13] =1, the correct implementation of lock elision for aborts that only have the XBEGIN CONFLICT bit set:

For other non-XBEGIN aborts the retry bit should still be taken into account.

For more details on Intel TSX please see the Intel Software Developer's manual volume 1 chapter 16 (Programming with Intel® Transactional Synchronization Extensions) and the Intel Optimization Manual Chapter 16 (Intel® TSX

Recommendations). Both available from  $\underline{\text{http://www.intel.com/sdm}}$ . For generic Intel TSX resources, refer to  $\underline{\text{http://www.intel.com/software/tsx}}$ .

§

Document Number: 604224, Revision 1.4

# 2 TSX Disable Update

This section describes behavior of a newer microcode update in 2021.1 IPU for a subset of the updated processors. This update for affected Intel® Xeon® Processor E3 v5 and v6 Family (code name Skylake and Kaby Lake) and the 6th, 7th, and 8th Generation Intel® Core™ i7 and i5 (code name Skylake, Kaby Lake, Coffee Lake, and Whiskey Lake) will disable Intel TSX by default.

By default, the processor will force abort all RTM transactions. CPUID bit CPUID.07H.0H.EDX[11](RTM\_ALWAYS\_ABORT) is set to indicate to updated software that the loaded microcode is forcing RTM abort. This bit can also be used to determine that the microcode update has been loaded with default settings that force aborts (see below section).

On processors that enumerate support for RTM, the CPUID enumeration bits for Intel TSX (CPUID.07H.0H.EBX[11] and CPUID.07H.0H.EBX[4]) continue to be set by default after the microcode update. System software may use new TSX\_FORCE\_ABORT[TSX\_CPUID\_CLEAR] functionality to clear those bits to indicate to software that RTM is disabled.

This microcode update eliminates the PMU interactions described in the previous section. PMU general purpose counter 3 can be considered reliable regardless of the value of the TSX\_FORCE\_ABORT MSR. Reads of the TSX\_FORCE\_ABORT[RTM\_FORCE\_ABORT] bit return value 1 by default to indicate the force abort behavior. Writes to the TSX\_FORCE\_ABORT[RTM\_FORCE\_ABORT] bit are ignored.

Consistent with prior updates, this microcode update will continue to unconditionally disable the HLE instruction prefix of Intel TSX.

#### 2.1 Unsupported software development mode

The default RTM force-abort behavior can be optionally disabled by setting MSR bit TSX\_FORCE\_ABORT.SDV\_ENABLE\_RTM=1. However, when RTM force abort is disabled in this way, RTM usage may be subject to memory-ordering correctness issues. Due to these issues, this unsupported mode **should not be enabled for production use**. System software might typically choose not to directly expose this functionality to users.

When TSX\_FORCE\_ABORT.SDV\_ENABLE\_RTM=1. CPUID bit CPUID.07H.0H.EDX[11](RTM\_ALWAYS\_ABORT) is cleared.

# 2.2 Updated MSR definition and affected products

The updated definition of the thread-scope TSX\_FORCE\_ABORT MSR is described in the following table. Support of this updated MSR definition can be determined by checking for the combination of the following conditions:

• CPUID.07H.0H.EDX[13] = 1

• CPUID.07H.0H.EDX[11](RTM\_ALWAYS\_ABORT) = 1 or TSX\_FORCE\_ABORT[SDV\_ENABLE\_RTM] = 1

Table 2.2-1 Description of updated TSX\_FORCE\_ABORT\_MSR

Register address		Register Name / Bit fields	Bit Description	Comment
Hex	Dec			
10f	271	TSX_FORCE_ABORT		
		0	RTM_FORCE_ABORT: Reads as 1, unless bit 2 is set. No implication on Counter 3.	Writes ignored, Default: 1
		1	TSX_CPUID_CLEAR: When set, CPUID.07H.0H.EBX[11]=0 and CPUID.07H.0H.EBX[4]=0.	R/W, Default: 0
		2	SDV_ENABLE_RTM: When set, processor may not force abort RTM. This unsupported mode should only be used for software development and not for production usage.	R/W, Default: 0
		3:63	Reserved	

Table 2.2-2 Affected products with Intel TSX disable microcode update

Family-Model	Stepping	Processor Families / Processor Number Series
06_55H	<=5	First generation Intel® Xeon® Scalable Processor Family and Intel® Xeon® Processor D Family based on Skylake microarchitecture
06_4EH, 06_5EH	All	6th generation Intel <sup>®</sup> Core <sup>™</sup> processors and Intel <sup>®</sup> Xeon <sup>®</sup> processor E3-1500m v5 product family and E3- 1200 v5 product family based on Skylake microarchitecture
06_8EH	<=0xB	7th/8th generation Intel <sup>®</sup> Core <sup>™</sup> processors and Intel <sup>®</sup> Pentium <sup>™</sup> processors based on Kaby Lake/Coffee Lake/Whiskey Lake microarchitecture
06_9EH	<=0xC	8th/9th generation Intel <sup>®</sup> Core <sup>™</sup> processors and Intel <sup>®</sup> Pentium <sup>™</sup> processors based on Coffee Lake microarchitecture

Document Number: 604224, Revision 1.4

# A Enabling FORCE\_ABORT\_RTM Mode to Use All Four Counters

It is possible for the PMU driver to opt-in to use all GP counters by enabling FORCE\_ABORT\_RTM mode. This requires setting bit 0 (FORCE\_ABORT\_RTM) in the TSX\_FORCE\_ABORT (0x10f) MSR for each logical CPU that is affected. The driver should only access this MSR when CPUID 7.EDX[13] is set.

When FORCE\_ABORT\_RTM is enabled, all RTM transactions on the logical CPU will forcefully abort, which can potentially impact performance of Intel TSX-enabled software, but the general purpose counter 3 will report correct values.

Application functionality should not be impacted because software that uses RTM is required to implement valid, non-transactional fallback paths for potential aborts, which are already exercised. When FORCE\_ABORT\_RTM mode is disabled, the RTM transactions will be allowed to commit again.

FORCE\_ABORT\_RTM mode does not change the CPUID feature enumeration for RTM or HLE.

FORCE\_ABORT\_RTM mode should always be disabled when the measurement session is finished to prevent applications that use RTM from experiencing performance impacts.

Table A-1. Description of TSX\_FORCE\_ABORT\_MSR

Register address		Register Name / Bit fields	Bit Description	Comment
Hex	Dec			
10f	271	TSX_FORCE_ABORT		MSR existence enumerated by CPUID 7:0 EDX[13]
		0	RTM_FORCE_ABORT: When set to 1 all RTM transactions abort with EAX code 0 while the bit it set. Counter 3 becomes usable.	
		1:63	Reserved	

# **B** Guidance on Specific Profilers

#### B.1 Linux\* perf

Linux\* perf is a PMU profiler integrated into the Linux kernel. With old or unpatched kernel versions, when Intel TSX is used, measurements using general purpose counter 4 (PMC3) may report incorrect values. This can be resolved with a kernel update.

The updated kernel exposes a new API as a <code>/sys/devices/cpu/allow\_tsx\_force\_abort</code> sys file that takes the values of 0 or 1. If the value is set by the administrator to 0, then the system will only support three general purpose counters for use with perf on affected systems with Intel TSX. In this case, RTM transactions will not be forced to abort.

However, when the value of /sys/devices/cpu/allow\_tsx\_force\_abort is 1 then RTM transactions will force abort only while anyone on the system is running a perf session under the following two scenarios (exclusively):

- 1) Using PMC3 while profiling the process using Intel TSX instructions. Only root or the user running the program can do this.
- 2) Using PMC3 while profiling the global system that includes the kernel and all processes running. This is done by using perf with the -a option. With default perf permission settings, a user needs to be root to use the -a option.

The default setting of the upstream kernel is to set allow\_tsx\_force\_abort to 1. However, the Linux distribution may choose to set allow\_tsx\_force\_abort to 0. When set to 1, Intel TSX users need to be aware of the possibility of performance variation due to Intel TSX instructions aborting associated with the use of perf under the two scenarios outlined above. To avoid this, system administrator can set allow tsx force abort to 0, or avoid using perf in the scenarios described above.

When allow\_tsx\_force\_abort is set to 0, then any perf command lines defining groups with four generic counter events will need to be updated to use at most three generic events per group.

Tools that generate Linux perf groups (such as pmu-tools or Intel® VTune $^{\text{TM}}$ ) will also need to be updated.

#### **B.2** Processor Counter Monitor

Processor Counter Monitor (PCM) is an application programming interface (API) and a set of tools based on the API used to monitor performance and energy metrics of Intel® Core $^{\text{\tiny TM}}$ , Intel® Xeon®, Intel® Atom $^{\text{\tiny TM}}$  and Intel® Xeon Phi $^{\text{\tiny TM}}$  processors.

When ALLCTR (force\_RTM\_abort) mode is disabled, these PCM versions automatically limit the number of metrics and events being collected simultaneously. These versions also support enabling ALLCTR mode (force\_RTM\_abort) using a command line switch (see the output of pcm.x --help). Users can also control the modes programmatically with the enableForceRTMAbortMode() and disableForceRTMAbortMode() API calls.

§

<sup>&</sup>lt;sup>1</sup> Eight counters when hyperthreading is disabled