

Intel[®] TDX Module Architecture Specification: TD Migration

DRAFT – WORK IN PROGRESS Updated Text is Highlighted in Red

Draft June 2025 5

Notices and Disclaimers

Intel Corporation ("Intel") provides these materials as-is, with no express or implied warranties.

All products, dates, and figures specified are preliminary, based on current expectations, and are subject to change without notice. Intel does not guarantee the availability of these interfaces in any future product. Contact your Intel representative to obtain the latest Intel product specifications and roadmaps.

The products described might contain design defects or errors known as errata, which might cause the product to deviate from published specifications. Current, characterized errata are available on request.

Intel technologies might require enabled hardware, software, or service activation. Some results have been estimated or simulated. Your costs and results might vary.

10 No product or component can be absolutely secure.

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 that includes the subject matter disclosed herein.

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

15 This document contains information on products, services and/or processes in development. All information provided here is subject to change without notice.

Copies of documents that have an order number and are referenced in this document or other Intel literature may be obtained by calling 1-800-548-4725 or by visiting <u>http://www.intel.com/design/literature.htm.</u>

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands might be claimed as the property of others.

Table of Contents

	Table of	f Contents	3
	SECTION	N 1: TO MIGRATION INTRODUCTION AND OVERVIEW	7
	0201101		
	1. /	About this Document	8
5	1.1.	Scope of this Document	8
	1.2.	Document Organization	8
	1.3.	Glossary	9
	1.4.	Notation	9
	15	References	10
10	1.5.	1 Intel Public Documents	10
10	1.5	2. Intel TDX Public Documents	10
	1.5	i.3. Non-Intel Public Documents	10
	2 1	TD Migration Overview	11
			11
	2.1.		11
15	2.2.	Updated: TD Migration Scenarios	11
	2.2	2.1. Cold migration	11
	2.2	2.2. Live Migration	12
	2.2	2.3. Image Snapshot and Jumpstart	12
	2.3.	Components Involved in TD Migration	12
20	2.4.	Migrated Assets	13
	2.5.	Guest TD Migration Life Cycle Overview	14
	2.5	5.1. Updated: Reservation and Session Setup	14
	2	2.5.1.1. Guest TD Build, Migration TD Binding and TD Execution on the Source Platform	14
	2	2.5.1.2. Guest TD Initial Build on the Destination Platform	14
25	2	2.5.1.3. Migration TDs Session Establishment	14
	2	2.5.1.4. Migration Session Key and Protocol Version Exchange	14
	2.5	5.2. Updated: In-Order Memory Migration Phase	15
	4	2.5.2.1. ID-Scope Infinitiable Metadata (Non-Memory State) Migration	15
30	2	2.5.2.2. Undated: Source TD Pause and Final Non-Memory State Migration	10
50	2	2.5.2.4. TD-Scope and VCPU-Scope Mutable Non-Memory State Migration	17
	2.5	5.3. Out-Of-Order Memory Migration Phase	17
	2	2.5.3.1. Migration of Memory State and Commitment of Import	17
	2	2.5.3.2. Updated: Post-Copy of Memory State	17
35	2.5	6.4. Migration Commitment	18
	2.5	5. Updated: Migration Abort	18
	2.6.	Impact of Migration on Measurement and Attestation	18
	2.7.	Intel TDX Module TD Migration Interface Functions Overview	18
	3. l	Updated: TD Migration Software Flows	19
40	3.1.	Typical TD Migration Flow Overview (Write-Blocking Based Export)	19
	3.2.	New: Typical TD Migration Flow Overview (Non-Blocking Export)	20
	3.3.	Successful Write-Blocking Based Export	20
	3.4.	New: Successful Non-Blocking Export	22
	3.5.	Successful Import	23
45	3.6.	TD Import Abort	24

	3.6.1.	TD Import Abort During the In-Order Import Phase	
	3.6.2.	TD Import Abort During the Out-Of-Order Import Phase	
	37 TD F	xnort Abort	25
	3.7.1.	Export Abort During the In-Order Export Phase	
5	3.7.2.	Export Abort During the Out-Of-Order Export Phase	
			27
	4. Migrati	ion TD, Migration Policy and the Extended TCB	
	4.1. Exter	nded TCB and the Migration Policy	
	4.2. Attes	station of the Migration TD and its Migration Policy	
10	4.3. Input	ts to the Migration TD's Migration Policy Evaluation	
	4.4. Migro	ated TD Information Provided by TDG.SERVD.RD	
	4.5. Migro	ation Protocol Version Setup	
	4.6. Migro	ation Session Keys (MSKs) Exchange	
	4.7. Exam	nple Migration Session Establishment	
15	5. Update	ed: Common TD Migration Mechanisms	31
	5.1. Migro	ation Bundles	
	5.1.1.	Overview	
	5.1.2.	Migration Data	
	5.1.3.	Migration Bundle Metadata (MBMD)	
20	5.1.4.	Untrusted Metadata	
	5.2. Upda	ated: Export and Import Functions Interface	
	5.2.1.	Migrating a Multi-Page Migration Bundle	
	5.2.2.	Updated: Migration Functions Interruptibility	
	5.3. Crypt	tographic Protection of Migration Data	
25	5.3.1.	Encryption Algorithm	
	5.3.2.	Migration Session Keys	
	5.4. Migro	ation Streams and Migration Queues	
	5.5. Meas	surement and Attestation	
	5.5.1.	TD Measurement Registers Migration	
30	5.5.2.	TD Measurement Reporting Changes	
	5.5.3.	TD Measurement Quoting Changes	
	5.5.4.	TCB Recovery and Migration	
	5.6. Upda	ated: TDX Control Structure Updates	
	5.6.1.	MIGSC: Migration Stream Context	
35	6. Update	ed: Migration Session Control and State Machines	
	6.1. Upda	ated: Overview	
	6.1.1.	Updated: Pre-Migration	
	6.1.2.	Updated: Successful Migration Session	
	.6.1.3	Aborted Migration Session	
40	6.1.3.1.	Abort during the In-Order Phase	
	6.1.3.2.	Abort during the Out-Of-Order Phase	
	0.1.4.		
	6.2. Upda	ated: Migration Session Control	
45	6.2.1.	Ivigration TD Binding and Migration Key Assignment	
45	6.2.2.	Updated: Export Side (Source Platform)	
	0.2.3. 6 7 <i>1</i>	Details' Migration State Machine	
	6.2.4.1	Details: Reminder: TD Lifecycle State Machine	

	6.2.4.2. 6.2.4.3.	Details: OP_STATE: TD Operation State Machine Details: OP_STATE Summary	
	6.3. Migrati	on Tokens	
	6.4. Miarati	on Protocol Versioning	
5	6.4.1. Int	roduction	
	6.4.2. En	umeration of Supported Migration Versions	
	6.4.3. Se	tting the Migration Protocol Version for a Migration Session	
	6.5. Update	d: Migration Session Control Functions Summary	
	7. TD Non-M	Jemory State Migration	
10	7.1 TD Non	Mamon State Migration Operation	50
10	7.1. ID NON 7.1.1 No	-Memory State Migration Data	50
	7.1.1. NC	n-Memory State MRMD	50 S
	7.1.2. No 7.1.3 Im	mutable vs. Mutable TD State	50
	7.2 Now: E	innected Configuration by the Host VMM	50
45	7.2. New. L	amony State Migration Synthesis Summary	
15	7.3. NOTI-IVI	emory state migration Functions Summary	
	8. Updated:	TD Private Memory Migration	52
	8.1. Update	d: Overview	
	8.1.1. <mark>U</mark> p	odated: In-Order and Out-of-Order Migration	52
	8.1.2. Ne	ew: Write-Blocking Export vs. Non-Blocking Live Export	
20	8.2. New: C	Conventions: SEPT Entry State Diagrams Color Coding	52
	8.3. Update	d: GPA Lists and Private Memory Migration Bundles	53
	8.3.1. Ov	/erview	53
	8.3.2. Up	odated: GPA List	
	8.3.3. Pa	ge Attributes List (Required for Partitioned TDs)	
25	8.3.4. Pri	ivate Memory Migration Buffer	
	8.4. Update	d: Write-Blocking Based Memory Export	
	8.4.1. Ed	ited: Host VMM Perspective	54
	8.4.1.1.	Typical Export Session	55
	8.4.1.2.	Live Export: Blocking for Writing, TLB Tracking and Exporting a Page	55
30	8.4.1.3.	Exporting a Page after the Source TD is Paused	
	8.4.1.4.	Unblocking for Write, Tracking Dirty Pages and Re-Exporting	
	8.4.1.5.	Using the same GPA List for TDH.EXPORT.BLOCKW and TDH.EXPORT.MEM	
	8.4.1.6.	Prohibited Operations on Exported Pages and Export Cancellation	
	8.4.1.7.	Updated: Exporting Pending Pages	
35	8.4.1.8.	Re-Exporting a Non-Dirty Page	
	8.4.1.9. 8.4.2 Ed	SEPT Cleanup after Export Abort	
	0.4.2. LU 9/101	Details: 11 SEDT Loaf Entry Partial State Diagram for Manned Page Expert	
	8422	Details: LI SEPT Leaf Entry Partial State Diagram for Pending Page Export	
40	8.4.2.3.	Details: TDCS.DIRTY_COUNT: TD-Scope Dirty Page Counter	
	8.5. New: N	Ion-Blocking Memory Export	
	8.5.1. Ho	ost VMM Perspective	61
	8.5.1.1.	EPT Access and Dirty Bits Background	62
	8.5.1.2.	Memory Export Concept: Scan and Export	
45	8.5.1.3.	Conceptual, Simplified Page State Diagram	
	8.5.1.4.	Scanning for Candidate Pages to Export or Re-export	64
	8.5.1.5.	Typical Non-Blocking Export Session	65
	8.5.1.6.	Interaction with Memory Management Operations	67
	8.5.1.7.	Exporting Pending Pages	67
50	8.5.1.8.	SEPT Cleanup after Export Abort	67
	8.5.2. De	tails of Non-Blocking Export	
	8.5.2.1.	Details: L1 SEPT Leaf Entry Partial State Diagram for Mapped Page Export	

	8.5.2.2.	Details: L1 SEPT Leaf Entry Partial State Diagram for Pending Page Export	
	8.5.2.3.	Details: TD Partitioning Considerations for Dirty Bit Operations	
	8.5.2.4.	Details: Pending Pages Considerations	71
	8.5.2.5.	Details: Blocked Pages Considerations	71
5	8.5.2.6.	Details: Memory Management Considerations	71
	8.5.2.7.	Details: Export Completeness Tracking	73
	8.6. Mem	ory Import	74
	8.6.1.	Host VMM Perspective	74
	8.6.1.1	In-Order Import Phase	74
10	8.6.1.2	Out-of-Order import Phase	74
	8.6.1.3	In-Place Import	75
	8.6.2.	Updated: Details of Memory Import	75
	8.6.2.1	Updated: Details: In-Order Import Phase	75
	8.6.2.2.	Details: Out-of-Order import Phase	77
15	8.7. Secu	re EPT Concurrency Considerations	
	8.7.1.	Overview	
	8.7.2.	New: GPA List Processing Implications	79
	8.8. Move	ed & Updated: Security Analysis: Achieving Memory Migration Security Objectives	
	8.8.1.	Updated: General	79
20	8.8.2.	Updated: Preventing Usage of Stale Memory Copies due to Mis-Ordering	
	8.8.3.	New: Enforcing Export of the Entire Memory Image	80
	8.8.4.	New: Non-Blocking Export: Detecting Memory State Change	
	8.8.5.	Updated: Preventing Usage of Stale Memory Copies due to Failure to Re-export	
	8.8.6.	Updated: Preventing Usage of Missing or Stale Memory Copies due to Failure to Impor	t81
25	8.8.7.	Preventing Usage of Stale Memory GPA Mapping and Attributes	
	8.8.8.	Out-Of-Order Phase and Its Usage for Post Copy	
	8.9. Upda	nted: Memory Migration Interface Functions Summary	81

5

SECTION 1: TD MIGRATION INTRODUCTION AND OVERVIEW

1. About this Document

1.1. Scope of this Document

5

10

20

This document describes the architecture and the external Application Binary Interface (ABI) of the Intel® Trust Domain Extensions (Intel® TDX) module's Live Migration feature, implemented using the Intel TDX Instruction Set Architecture (ISA) extensions, for cold or live migration of Trust Domains in an untrusted hosted cloud environment.

This document is part of the **TDX Module Architecture Specification Set**, which includes the following documents:

 Table 1.1: TDX Module Architecture Specification Set

Document Name	Reference	Description
TDX Module Base Architecture Specification	[TDX Module Base Spec]	Base TDX module architecture overview and specification, covering key management, TD lifecycle management, memory management, virtualization, measurement and attestation, service TDs, debug aspects etc.
TDX Module TD Migration Architecture Specification	[TD Migration Spec]	Architecture overview and specification for TD migration
TDX Module TD Partitioning Architecture Specification	[TD Partitioning Spec]	Architecture overview and specification for TD Partitioning
TDX Module TDX Connect Specification	[TDX Connect Spec]	Architecture overview and specification for TDX Connect
TDX Module ABI Reference Tables	[TDX Module ABI Tables]	A set of JSON format files detailing TDX module Application Binary Interface (ABI)
TDX Module TDX Connect ABI Reference Specification	[TDX Connect ABI Spec]	Detailed TDX module Application Binary Interface (ABI) reference specification, covering the TDX connect architecture
TDX Module ABI Reference Specification	[TDX Module ABI Spec]	Detailed TDX module Application Binary Interface (ABI) reference specification, covering the entire TDX module architecture

This document is a work in progress and is subject to change based on customer feedback and internal analysis. This document does not imply any product commitment from Intel to anything in terms of features and/or behaviors.

Note: The contents of this document are accurate to the best of Intel's knowledge as of the date of publication, though Intel does not represent that such information will remain as described indefinitely in light of future research and design implementations. Intel does not commit to update this document in real time when such changes occur.

15 **1.2.** Document Organization

The document has two main sections:

- Section 1 contains an introduction to the document, overview of TD Migration, scenarios and requirements.
- Section 2 contains the Intel TDX Module Migration architecture

The detailed reference specification of TD Migration data structures and interface functions is provided in the [TDX Module ABI Spec].

1.3. Glossary

For a complete TDX module glossary, see the [TDX Module Base Spec].

Table 1.2: Intel TDX Module Glossary for TD Migration

Acronym	Full Name	New for TDX	Description	
	Blackout Period	No	The period where that TD does not run anymore on the source platform and does not run yet on the destination platform.	
	Cold Migration	No	Migration usage mode where the TD does not run during the migration session.	
	In-Order Phase	Yes	The first phase of the TD migration session, where a strict order of memory export vs. memory import is maintained.	
	Live Migration	No	Migration usage mode where the TD runs during most of the migration session.	
	Migration Bundle	Yes	The basic unit of migrated information, composed of headers, metadata and/or migrated data.	
	Migration Commitment	Yes	The act of committing the migration, disallowing the TD from running on the source platform and allowing it to run of the destination platform.	
	Migration Epoch	Yes	A mechanism used to enforce ordering across multiple concurrent streams. Any specific page can be migrated only once per epoch.	
	Migration Policy	No	Policy enforced by the Migration TDs, e.g., based on the source and destination platform properties.	
	Migration Stream	Yes	A sequential stream of migration bundles, where order is enforced.	
MigTD	Migration TD	Yes	A specific type of Service TD used to provide Live Migration capability for TD VMs. A Migration TD extends the TCB of the serviced tenant TD.	
MSK	Migration Session Key	Yes	AES-GCM-256 key generated by the source MigTD and shared with the destination MigTD (protected by the Migration Transport key). This key helps protect the TD private data and is used for export and import of the TD confidential assets.	
МТК	MTK Migration Yes Transport Key		Authenticated Diffie-Helman negotiated symmetric key generated after mutual attestation of the MigTDs and is used to help protect the transport of the Migration Session Key from the source to the destination platform.	
	Out-Of-Order Phase	Yes	The last phase of the TD migration session, where a strict order of memory export vs. memory import is not maintained.	
	Post-Copy	No	Migration usage mode where some of the TD memory image is migrated after the TD is allowed to run on the destination platform.	
	Pre-Copy	No	Migration usage mode where (most of) the TD memory image is migrated before the TD is allowed to run on the destination platform.	
ТСР	Transmission Control Protocol	No	Transmission Control Protocol (TCP) is a communications standard that enables application programs and computing devices to exchange messages over a network. It is designed to send packets across the internet and ensure the successful delivery of data and messages over networks.	

1.4. Notation

5 See the [TDX Module Base Spec].

Draft

1.5. References

1.5.1. Intel Public Documents

See the [TDX Module Base Spec].

- 1.5.2. Intel TDX Public Documents
- 5 See the [TDX Module Base Spec].

1.5.3. Non-Intel Public Documents

Table 1.3: Non-Intel Public Documents

Reference	Document	Version & Date
AES-256-GCM	NIST Special Publication 800-38D: Recommendation for Block Cipher Modes of Operation: Galois/Counter Mode (GCM) and GMAC	November 2007

2. TD Migration Overview

For an overview of TDX, refer to the [TDX Module Base Spec].

2.1. Introduction



5

25

Figure 2.1: TD Migration

Analogous to legacy VM migration, a cloud-service provider (CSP) may want to relocate/migrate an executing Trust Domain from a **source TDX platform** to a **destination TDX platform** in the cloud environment. A cloud provider may use TD migration to meet customer SLA, while balancing cloud platform upgradability, patching and other serviceability requirements. Since a TD runs in a CPU mode which helps protect the confidentiality of its memory contents and its CPU state from any other platform software, including the bosting Virtual Machine Monitor (VMM), this primary security

- state from any other platform software, including the hosting Virtual Machine Monitor (VMM), this primary security objective must be maintained while allowing the TD resource manager, i.e., the host VMM to migrate TDs across compatible platforms. The TD typically may be assigned a different HKID (and will be always assigned a different ephemeral key) on the destination platform chosen to migrate the TD.
- In this specification, the TD being migrated is called the **source TD**, and the TD created as a result of the migration is called the **destination TD**. An extensible **TD Migration Policy** is associated with a TD that is used to maintain the TD's security posture. The TD Migration policy is enforced in a scalable and extensible manner using a specific type of **Service TD** called the **Migration TD (a.k.a. MigTD)** (introduced in the Figure 2.2 below) – which is used to provide services for migrating TDs.

The TD Live Migration process (and the Migration TD) does not depend on any interaction with the TD guest software operating inside the TD being migrated.

2.2. Updated: TD Migration Scenarios

This section describes the usage scenarios addressed by this specification (and those explicitly out of scope). This specification documents the TD Migration functionality from a Live Migration (scenario described below) perspective. Cold Migration and other scenarios described below are effectively subset scenarios that are software managed via the Intel TDX module interface functions in this specification.

- 2.2.1. Cold migration
- Both source and destination platforms must be alive during migration.
- The TD is suspended during migration. Blackout time is typically longer than TCP timeout, which may cause remote connections to the TD to break up.

Cold migration may be useful for rolling upgrades or patches and rebooting servers (non-reboot patches can be done without migrating the TD), capacity planning and load balancing.

A TD may be cold migrated more than once using multiple sessions.

2.2.2. Live Migration

5

- Both source and destination platforms must be alive during migration.
 - The TD continues executing during migration. It is paused for a short blackout time, typically shorted than TCP timeout, so remote connections to the TD should not break up.

Live migration may be useful for customer SLA support, capacity planning and load balancing.

A TD may be live migrated more than once using multiple sessions.

10 2.2.3. Image Snapshot and Jumpstart

- Destination need not be alive during export, when the migration image is prepared.
- Source need not be alive during import, when the migration image is loaded.
- TD image may be stored for an indeterminate amount of time.

This usage has additional platform security requirements that are not comprehended in this specification. Example use cases are saving checkpoints of TDs such that TD may be pre-loaded into memory. Alternate implementations to satisfy this usage are possible. E.g., the TD could un-hibernate the image itself. **This scenario/use case is out of scope for this specification** and is not supported by it.

2.3. Components Involved in TD Migration



Figure 2.2: Components Involved in TD Migration

Migrated TD

20

25

The migrated TD's role is passive. It is not directly aware of it being migrated.

Migration TD (MigTD)

The Migration TDs (MigTD) exist on the source and destination platforms. Their main role is to implement a **migration policy** and evaluate migration sources and destinations for adherence to that policy. The migration policy may enumerate TDX platform TCB requirements, platform features and acceptable destination Migration TD TCB requirements.

The MigTDs securely exchange unique per-session **Migration Session Key (MSK)** pair. The MSKs are used to migrate assets of a specific TD. The MigTD on each side reads an encryption key generated by the TDX module and securely transfers it to the MigTD on the other side, where it is written as the decryption key for its side.

Migration TDs implement and use the TDX module Service TD protocol (for details, see [TDX Module Base Spec]'s Service
 TDs chapter). The host VMM may bind a MigTD to one or migrated TDs. The MigTD is in the TCB of the migrated TD and its measurements are in included in the migrated TD's attestation information. Thus, the MigTD must be pre-bound to the migrated TD before that TD's measurement is finalized. The MigTD lifecycle does not have to be coincidental with the migrated TD – the only requirement is that the MigTD must be bound to the migrated TD before migration can begin and must be operational until the migration session keys has been successfully exchanged.

10 Host VMM

20

The host VMMs on the source and destination sides orchestrate and manage the migration session, via their respective TDX modules and MigTDs. They are responsible for export and import of the TD's memory and non-memory state, via the TDX module, and for the transport of that state between the source and the destination platforms.

TDX Module

15 The TDX module implements a set of TD migration primitives to implement and enforce the security of migration session control, TD private memory migration and TD non-memory state migration.

2.4. Migrated Assets

The table below shows the TD assets that are migrated. Metadata includes TD-scope and VCPU-scope non-memory state (such as control state, CPU register state etc.) and memory attributes (such as GPA and access permissions). Metadata is not migrated as-is; it is serialized into a migration format and re-created on the destination platform.

TD Asset	Where Held	Export Functions	Import Functions
Immutable Non-Memory State (Metadata)	TDX module global TDR TDCS	TDH.EXPORT.STATE.IMMUTABLE	TDH.IMPORT.STATE.IMMUTABLE
Mutable Non-Memory State (Metadata)	TDCS TDVPS	TDH.EXPORT.STATE.TD TDH.EXPORT.STATE.VP	TDH.IMPORT.STATE.TD TDH.IMPORT.STATE.VP
Memory State and Metadata	TD private pages Secure EPT	TDH.EXPORT.MEM	TDH.IMPORT.MEM

Table 2.1: Migrated TD Assets

2.5. Guest TD Migration Life Cycle Overview

2.5.1. Updated: Reservation and Session Setup



Figure 2.3: Updated: Migration TD Binding and Session Setup

5 2.5.1.1. Guest TD Build, Migration TD Binding and TD Execution on the Source Platform

The source TD build and execution process is described in the [TDX Module Base Spec]. To be migratable, the TD may be initialized, using the TDH.MNG.INIT function, with **ATTRIBUTES.MIGRATABLE** bit set to 1.

Before a migration session can begin, the host VMM on the source platform must use TDH.SERVTD.BIND to bind a Migration TD to the source TD. The MigTD may read selected metadata fields of the source TD, e.g., its ATTRIBUTES and XFAM configuration, to be used in evaluating a migration policy. The bound MigTD is reflected in the migratable TD's TDREPORT.

2.5.1.2. Guest TD Initial Build on the Destination Platform

Same as a legacy TD build process, the host VMM creates a new guest TD by using the TDH.MNG.CREATE interface function. This destination TD is setup as a "template" to receive the state of the Source Guest TD. As with any TD build, the host VMM configures the TD's private HKID (which is no relation to the HKID used on the source platform) using the TDH.MNG.KEY.CONFIG interface function on each package. The host VMM can then continue to build the TDCS by adding TDCS pages using the TDH.MNG.ADDCX interface function.

Once the destination TDCS is built and before TD import can begin, the VMM on the destination platform must use TDH.SERVTD.BIND to bind a Migration TD to the destination TD. Once migration succeeds, the MigTD bound at the destination will be reflected in the migratable TD's TDREPORT.

2.5.1.3. Migration TDs Session Establishment

Prior to starting any TD migration session, the migration TDs on the source and destination platforms need to create a secure connection between them. This connection may be done for supporting one or more migration sessions. Migration TDs may be written by any vendor; thus, details may vary. Typically, the migration TDs executing on the source and destination platforms use a TD-quote-based mutual authentication protocol to create a VMM-transport-agnostic session between them. The Migration TDs typically negotiate a protected transport session (using Diffie-Hellman exchange). Using this protected transport session, the migration policy can be evaluated by the Migration TDs.

2.5.1.4. Migration Session Key and Protocol Version Exchange

The migration session keys are ephemeral AES-256-GCM keys used for confidentiality and integrity protection of the migrated TD private state, and for integrity protections of the migration session control protocol. TD shared memory

30

10

15

20

25

5

15

state is migrated by the untrusted host VMM per legacy methods – the same network transport may be used for both by the host VMM. The Migration Session Keys (MSKs) are exchanged between the TDX modules on both side with the help of the Migration TDs.

The TDX module on each side generates an ephemeral migration session encryption key. The Migration TDs on each

side uses the Service TD metadata read function (TDG.SERVTD.RD) to read the encryption key and securely transfer it to the peer Migration TD, which uses the Service TD metadata write function (TDG.SERVTD.WR) to write it as the migration session decryption key.

In addition to the decryption keys, the Migration TDs on both sides write the migration protocol version to be used by the TDX modules.

After this point, the host VMM can invoke TDX Module functions such as TDH.EXPORT.* to export state at the source 10 platform and TDH.IMPORT.* to import TD state at the destination platform.

Updated: In-Order Memory Migration Phase 2.5.2.

The TD migration session has two phases: in-order and out-order. Those terms are defined in the context of TDX module enforcement of memory import ordering vs. memory export ordering. In-order enforcement is required during the live migration phase, until the source TD is paused and the blackout period begins. During this phase, as long as the source TD pages are mutable, order enforcement is essential to help ensure the migrated memory image is correct.



Figure 2.4: Updated: In-Order Migration Phase

2.5.2.1. TD-Scope Immutable Metadata (Non-Memory State) Migration

The TDX module protects the confidentiality and integrity of a guest TD non-memory state. Control structures, which 20 hold guest TD metadata, are not directly accessible to any software (besides the Intel TDX module) or devices. These structures are stored encrypted in memory with the TD private key and managed by TDX module interface functions.

Immutable metadata is the set of TD-scope state variables that are set by TDH.MNG.INIT, may be modified during TD build but are never modified after the TD's measurement is finalized using TDH.MR.FINALIZE. Some of these state variables control how the TD and its memory is migrated. Therefore, the immutable TD control state is migrated before any of the TD memory state is migrated.

25

TD immutable state is exported via the TDH.EXPORT.STATE.IMMUTABLE interface function and imported on the destination platform via the TDH.IMPORT.STATE.IMMUTABLE interface function. TD global immutable state migration is described in Ch. 7.

Draft

2.5.2.2. Updated: Iterative Pre-Copy of Memory State

During the in-order migration phase, the VMM aims to balance TD availability with progressing the migration work. This is generally done through an iterative pre-copy of TD memory pages and their in-order import on the destination platform. The details of TD private memory migration covered in Ch. 8; this section explains its key concepts.

5 2.5.2.2.1. Updated: Migration Considerations for TD Private Memory

Memory Migration vs. Memory Encryption

TD Migration does not migrate the TD's private HKID nor the keys used to encrypt its private memory. The TD's private assets are migrated using the migration session keys as described above. The host VMM on the destination platform assigns a new free HKID and the TDX module generates new memory encryption keys.

10 Memory Modification Tracking and Iterative Migration

During live migration, the running source TD and, if enabled for TDX Connect, its bound TDIs may write to memory. A memory page that has been exported and later modified must be re-migrated in order to ensure the consistency of the migrated memory image. To enforce that, the TDX module employs a page modification tracking mechanism. Two export modes may be supported:

- Write-blocking export is based on blocking memory for write, to detect TD attempts to modify memory that has been exported.
 - Non-blocking export (if supported by the TDX module) detects memory modifications based on EPT Dirty bit indication set by the hardware.

Migration Streams

To utilize multiple LPs and improve migration bandwidth, the host VMM may create multiple **migration streams** for concurrently transferring memory state. Migration streams are described in 5.4.

Migration Epochs

25

30

Each migration stream enforces import vs. export ordering within the stream. To enforce ordering across multiple concurrent streams, migration epochs are used. Any specific page can be migrated only once per epoch. Migration epochs are described in 6.1.4.

2.5.2.2.2. Migration Considerations for EPT Structures

Secure EPT trees are not migrated but rebuilt on the destination platform using memory page attributes, which are migrated as metadata along with page contents. To do this, the host VMM calls the TDX module's TDH.MEM.SEPT.ADD to build the Secure EPTs and TDH.IMPORT.MEM to import memory pages. TDH.IMPORT.MEM uses the cryptographically protected page metadata (including GPA, Read, Write and Execute attributes, and PENDING state) to ensure Secure EPT properties from the source platform are accurately recreated, preventing remap attacks during migration.

Migrating shared memory is outside the scope of TDX. Shared memory assigned to the TD can be migrated by the host VMM using legacy mechanisms.

2.5.2.3. Updated: Source TD Pause and Final Non-Memory State Migration

Following pre-copy of TD private memory, the host VMM must detach any connected TDIs and pause the source TD for a brief period (also called the blackout period). The VMM initiates this via TDH.EXPORT.PAUSE, which checks security pre-conditions and prevents TD VCPUs from executing any more.

2.5.2.3.1. Final Non-Memory State Migration

TD **mutable non-memory state** is a set of source TD state variables that might have changed since it was finalized via TDH.MR.FINALIZE. Immutable non-memory state exists for the TD scope (as part of the TDR and TDCS control structures) and the VCPU scope (as part of the TDVPS control structure).

Once the source TD is paused, the host VMM exports the final (mutable) TD non-memory state, for the TD as a whole and for each VCPU.

Mutable TD state is exported by TDH.EXPORT.STATE.TD (per TD) and TDXPORT.STATE.VP (per VCPU) and imported by TDH.IMPORT.STATE.TD and TDH.IMPORT.STATE.VP respectively. This is described in Ch. 7. 5

2.5.2.3.2. In-Order Memory State Migration Completion

Any memory state that has been migrated must be up to date when the in-order phase completes. If a memory page had been migrated and its content or attributes were later updated by the running source TD or by a TDI, it must be remigrated. If the TD is enabled for TDX Connect, then all its private memory must be migrated.

The TDX Module enforces this using the commitment protocol described in 2.5.4 below.

2.5.2.4. TD-Scope and VCPU-Scope Mutable Non-Memory State Migration

2.5.3. Out-Of-Order Memory Migration Phase



Figure 2.5: Out-Of-Order Migration Phase

10 2.5.3.1. Migration of Memory State and Commitment of Import

If TDX connect is not enabled for the migrated TD, memory pages that have not been migrated during the in-order phase may be migrated during the out-of-order phase. Since the memory state on the source platform does not change at this stage, the order of import vs. export is not enforced.

The host VMM on the destination platform commits the import by calling TDH.IMPORT.COMMIT (if post-copy is to be used) or TDH.IMPORT.END. At that point, the TD may run on the destination platform and may not run on the source platform.

2.5.3.2. Updated: Post-Copy of Memory State

In some live migration scenarios, the host VMM may stage some memory state transfer to occur lazily after the destination TD has started execution. In this case, the host VMM will be required to fetch the required pages as accesses occur by the destination TD – this order of access is indeterminate and will likely differ from the order in which the host VMM has queued memory state to be transferred.

20

In order to support that on-demand model, the order of memory migration during this **post-copy stage** is not enforced by TDX. The host VMM may implement multiple migration queues with multiple priorities for memory state transfer. For example, the host VMM on the source platform may keep a copy of each encrypted migrated page until it receives a

25 confirmation from the destination that the page has been successfully imported. If needed, that copy can be re-sent on a high priority queue. Another option is, instead of holding a copy of exported pages, to call TDH.EXPORT.MEM again on demand.

Also, to simplify host VMM software for this model, the TDX module interface functions used for memory import in this post-copy stage return additional informational error codes to indicate that a stale import was attempted by the host

VMM to account for the case where the low latency import operation for a GPA superseded the import from the higher latency import queue.

Post-copy migration is not supported if the TD is enabled for TDX Connect.

2.5.4. Migration Commitment

5 The commitment protocol is enforced by the Intel TDX Module to help ensure that a host VMM cannot violate the security objectives of TD Live migration – for example, either the destination or source TD, but not both, may execute after live migration of the source TD to a destination TD, even if an error causes the TD migration to be aborted.

On the source platform, TDH.EXPORT.PAUSE starts the blackout phase of TD live migration and TDH.EXPORT.TRACK(DONE) ends the blackout phase of live migration (and marks the end of the transfer of TD memory pre-copy, mutable TD VP and mutable TD global control state). TDH.EXPORT.TRACK(DONE) generates a secure **start token**, which indicates that the TD will not run on the source platform and allows the destination TD to become runnable. On the destination platform, TDH.IMPORT.TRACK consumes the start token. TDH.IMPORT.COMMIT (in case post-copy is used) or TDH.IMPORT.END commits the migration and allows the TD to run on the destination.

2.5.5. Updated: Migration Abort

15 There are two abort scenarios:

10

Source Initiated Abort:	In a live migration scenario, an error may cause migration orchestration to abort the
	migration session while pre-copy is in progress, i.e., the source TD may still run or export
	blackout started, but no start token has yet been generated. In such a scenario, the host
	VMM on the source platform may initiate an abort via TDH.EXPORT.ABORT.

20 **Destination Initiated Abort:** If the pre-copy is complete and a start token has been generated, abort must be initiated by the host VMM on the destination platform. The TD may be re-allowed to run on the source platform unless the import session has been committed by TDH.IMPORT.COMMIT or TDH.IMPORT.END. To do that, the host VMM can generate an **abort token** using TDH.IMPORT.ABORT. The abort token, which indicates that the TD will not run on the destination platform, is sent to the source platform where it may be consumed by TDH.EXPORT.ABORT.

In both scenarios, the host VMM on the source platform must restore the SEPT state of exported pages on the source platform, using TDH.EXPORT.RESTORE or TDH.MEM.SCAN.RANGE(EXPORT_RESTORE) (if supported).

The detailed operations are described in Ch. 26.

2.6. Impact of Migration on Measurement and Attestation

TD measurement is extended for the MigTD bound to the TD being migrated, and the ATTRIBUTES.MIGRATABLE bit is part of the TD attestation. For details, see the [TDX Module Base Spec].

2.7. Intel TDX Module TD Migration Interface Functions Overview

See the [TDX Module Base Spec]'s overview chapter.

3. Updated: TD Migration Software Flows

This chapter summarizes the software flows used for TD migration using Intel TDX Module interface functions.

3.1. Typical TD Migration Flow Overview (Write-Blocking Based Export)

Source Platform				Destination Platform		
1.	Create Worker Threads, disable device optimizations		-		 Receive memory layout and device resource information 	
2.	Save memory layout and device resource information		1		5. Allocate guest memory and reserve resources on devices	
3.	Send memory layout and resource information to the destination host	Send Migration Information	<i>></i>	Receive Migration Information		
			_		6. Setup worker process and signal the completion to source host	
		TDH.SERVTD.BIND		TDH.SERVTD.BIND		
		Establish DH session between MigTDs	\leftrightarrow	Establish DH session between MigTDs		
7.	MigTDs exchange session keys, version information etc.	TDG.SERVTD.RD, TDG.SERVTD.WR	\leftrightarrow	TDG.SERVTD.RD TDG.SERVTD.WR	 MigTDs exchange session keys, version information etc. 	
		TDH.EXPORT.STATE. IMMUTABLE	\rightarrow	TDH.IMPORT.STATE. IMMUTABLE		
		Memory Pre-Export		Memory Pre-Import		
		TDH.EXPORT.BLOCKW				
9.	Run multiple passes of memory	TDH.MEM.TRACK				
	walker to transfer dirty guest-	TDH.EXPORT.MEM	\rightarrow	TDH.IMPORT.MEM	10. Import TD memory contents	
	memory pages	TDH.EXPORT.UNBLOCKW			-	
		TDH.EXPORT.TRACK	\rightarrow	TDH.IMPORT.TRACK		
		TDH.EXPORT.PAUSE			-	
11	Fundant mutable state	TDH.EXPORT.STATE.TD	\rightarrow	TDH.IMPORT.STATE.TD		
11	. Export mutable state	TDH.EXPORT.STATE.VP	\rightarrow	TDH.IMPORT.STATE.VP		
		TDH.EXPORT.TRACK(start token)	\rightarrow	TDH.IMPORT.TRACK(start token)		
13	. Export memory	TDH.EXPORT.MEM	\rightarrow	TDH.IMPORT.MEM	14. Import memory	
		15. (optional) Send end of export indication	\rightarrow	16. (optional) Receive end of export indication		
				TDH.IMPORT.COMMIT		
					17. Resume VPs	
			÷	18. (optional) Request post-copy of pages		
19	. (optional) Export memory (post- copy)	TDH.EXPORT.MEM	\rightarrow	TDH.IMPORT.MEM		
			_	TDH.IMPORT.END		
		Receive migration success indication	←	Send migration success indication		
		Indicate success to VM Management		Indicate success to VM Management		

20. Tear-down source TD

Figure 3.1: Typical TD Migration Flow

3.2. New: Typical TD Migration Flow Overview (Non-Blocking Export)

	Source	Platform		Destinati	ion I	Platform
1. 2.	Create Worker Threads, disable device optimizations Save memory layout and device		_		4. 5.	Receive memory layout and device resource information Allocate guest memory and
	resource information		_			reserve resources on devices
3.	Send memory layout and resource information to the destination host	Send Migration Information	\rightarrow	Receive Migration Information		
			_		6.	Setup worker process and signal
		TDH.SERVTD.BIND	1	TDH SERVED BIND		the completion to source host
		Establish DH session between MigTDs	\leftrightarrow	Establish DH session between MigTDs		
7.	MigTDs exchange session keys, version information etc.	TDG.SERVTD.RD, TDG.SERVTD.WR	\leftrightarrow	TDG.SERVTD.RD TDG.SERVTD.WR	8.	MigTDs exchange session keys, version information etc.
		TDH.EXPORT.STATE. IMMUTABLE	→	TDH.IMPORT.STATE. IMMUTABLE		
		Memory Pre-Export		Memory Pre-Import		
		TDH.MEM.SCAN.RANGE(DSCAN)				
9.	Initial In-Order Memory Export	TDH.MEM.TRACK			10	Initial In-Order Memory Import
	Round: Run memory walker to	IPIs			10	Round: Import TD memory
	pages	TDH.EXPORT.MEM	\rightarrow	TDH.IMPORT.MEM		contents
		TDH.EXPORT.TRACK	\rightarrow	TDH.IMPORT.TRACK		
		TDH.MEM.SCAN.RANGE(DSCAN)				
11	. Additional In-Order Memory	TDH.MEM.TRACK			12	Additional In-Order Memory
	Export Rounds: Run memory walker to transfer dirty guest-	IPIs			12	Import Rounds: Import TD
	memory pages	TDH.EXPORT.MEM	\rightarrow	TDH.IMPORT.MEM		memory contents
		TDH.EXPORT.TRACK	\rightarrow	TDH.IMPORT.TRACK		
		TDH.EXPORT.PAUSE				
13	B. Final In-Order Memory Export	TDH.MEM.SCAN.COMP(DCHECK)			14	. Final In-Order Memory Import
	Round	TDH.EXPORT.MEM	\rightarrow	TDH.IMPORT.MEM		Round
15	Export mutable state	TDH.EXPORT.STATE.TD	\rightarrow	TDH.IMPORT.STATE.TD	16	Import mutable state
1.	. Export mutable state	TDH.EXPORT.STATE.VP	\rightarrow	TDH.IMPORT.STATE.VP	10	
		TDH.EXPORT.TRACK(DONE)	\rightarrow	TDH.IMPORT.TRACK		
		(optional) Send end of export indication	\rightarrow	(optional) Receive end of export indication		
				TDH.IMPORT.COMMIT		
					16	. Resume VPs
			+	(optional) Request post-copy of pages		
17	7. (optional – non TDX Connect) Export memory (post-copy)	TDH.EXPORT.MEM	\rightarrow	TDH.IMPORT.MEM		
			-	TDH.IMPORT.END		
		Receive migration success indication	÷	Send migration success indication		
		Indicate success to VM Management		Indicate success to VM Management		
18	3. Tear-down source TD					

3.3. Successful Write-Blocking Based Export

The following sequence is typically used to export a TD from a source platform using write blocking.

Migration Step Des Phase Des		Description	Plurality	TDX Module Interface Functions
In-Order Export	Start of Export Session	VMM initializes MigTD (as a non- migratable TD) and binds it to the source TD.	Once	TDH.SERVTD.BIND
		VMM/orchestration sets up transport session between source and destination MigTD. MigTDs setup their own protected channel.	Once	
		MigTDs reads the session encryption key, version information and other metadata from the source TD and sends it to the MigTD on the destination.	Once	TDG.SERVTD.RD
		MigTD receives the session encryption key from the MigTD on the destination. MigTD writes it as the session decryption key to the source TD. It may also write the migration version.	Once	TDG.SERVTD.WR
		VMM starts the export session and exports immutable state creating a state migration bundle.	Once	TDH.EXPORT.STATE.IMMUTABLE
	Live Memory Export	Host VMM blocks a set of pages for writing.	Multiple	TDH.EXPORT.BLOCKW
		Host VMM increments the TD's TLB epoch	Once per migration epoch	TDH.MEM.TRACK
		Host VMM starts migration epoch and creates epoch token migration bundle; a page can be exported once per epoch.	Once per migration epoch	TDH.EXPORT.TRACK(epoch token)
		Host VMM exports, re-exports or cancels the export of TD private pages and creates a memory migration bundle.	Multiple	TDH.EXPORT.MEM
		TD write attempt to write to page blocked for writing results in an EPT violation. The host VMM unblocks the page; if already exported, it will need to be re-blocked and re-exported.	Multiple	TDH.EXPORT.UNBLOCKW
	Mutable Non-	VMM pauses the source TD	Once	TDH.EXPORT.PAUSE
	Memory State Export	VMM exports mutable TD-scope state and creates a state migration bundle.	Once	TDH.EXPORT.STATE.TD
		VMM exports mutable VCPU-scope state and creates a state migration bundle.	Per VCPU	TDH.EXPORT.STATE.VP
Out-Of- Order Export	Cold Memory Export	Host VMM starts the out-of-order export phase and creates a start token migration bundle.	Once	TDH.EXPORT.TRACK(start token)
		Host VMM exports TD private pages and creates a memory migration bundle.	Multiple	TDH.EXPORT.MEM

Table 3.1: Typical Write-Blocking Based TD Export Sequence

Migration Phase	Step	Description	Plurality	TDX Module Interface Functions
TD Teardown	End	Host VMM gets success notification from the destination platform, terminates the export session and tears down the TD on the source platform	Once	TDH.MNG.VPFLUSHDONE TDH.PHYMEM.CACHE.WB TDH.MNG.KEY.FREEID TDH.PHYMEM.PAGE.RECLAIM TDH.PHYMEM.PAGE.WBINVD

3.4. New: Successful Non-Blocking Export

The following sequence is typically used to export a TD from a source platform using the non-blocking method.

Table 3.2: Typical Non-Blocking TD Export Sequence				
Migration Phase	Step	Description	Plurality	TDX Module Interface Functions
In-Order Export	Start of Export Session	Same as above		
	Live Memory Export (Initial	Host VMM scans memory for page export candidates.	Multiple	TDH.MEM.SCAN.RANGE(DSCAN)
	Round)	Host VMM increments the TD's TLB epoch.	Once per migration epoch	TDH.MEM.TRACK
		Host VMM exports TD private pages reported by TDH.MEM.SCAN and creates memory migration bundles.	Multiple	TDH.EXPORT.MEM
	Live Memory Export (Optional Additional Rounds)	Host VMM scans memory for page export candidates.	Multiple	TDH.MEM.SCAN.RANGE(DSCAN)
		Host VMM increments the TD's TLB epoch	Once per migration epoch	TDH.MEM.TRACK
		Host VMM starts migration epoch and creates epoch token migration bundle; a page can be exported once per epoch.	Once per migration epoch	TDH.EXPORT.TRACK(epoch token)
		Host VMM exports, re-exports or cancels the export of TD private pages reported by TDH.MEM.SCAN and creates a memory migration bundle.	Multiple	TDH.EXPORT.MEM
	Final Memory	VMM pauses the source TD	Once	TDH.EXPORT.PAUSE
	Export and Mutable Non- Memory State Export	VMM exports mutable TD-scope state and creates a state migration bundle.	Once	TDH.EXPORT.STATE.TD
		VMM exports mutable VCPU-scope state and creates a state migration bundle.	Per VCPU	TDH.EXPORT.STATE.VP
		Host VMM scans memory for page export candidates.	Multiple	TDH.MEM.SCAN.COMP(DCHECK)
		Host VMM exports, re-exports or cancels the export of TD private pages reported by TDH.EXPORT.SCAN and creates a	Multiple	TDH.EXPORT.MEM

Table 2.2. Typical Nr Plackin rt S

memory migration bundle.

Migration Phase	Step	Description	Plurality	TDX Module Interface Functions
Out-Of- Order Export	f- Cold Memory Export t	Host VMM starts the out-of-order export phase and creates a start token migration bundle.	Once	TDH.EXPORT.TRACK(start token)
		Host VMM exports TD private pages and creates a memory migration bundle.	Multiple	TDH.EXPORT.MEM
TD Teardown	End	Host VMM gets success notification from the destination platform, terminates the export session and tears down the TD on the source platform	Once	TDH.MNG.VPFLUSHDONE TDH.PHYMEM.CACHE.WB TDH.MNG.KEY.FREEID TDH.PHYMEM.PAGE.RECLAIM TDH.PHYMEM.PAGE.WBINVD

3.5. Successful Import

The following sequence is typically used to import a TD to a destination platform.

Table 3.3: Typical TD Import Sequence

Migration Phase	Step	Description	Plurality	TDX Module Interface Functions
In-Order Import	Start of Import Session	VMM creates the destination TD's skeleton	Once	TDH.MNG.CREATE TDH.MNG.KEY.CONFIG TDH.MNG.ADDCX
		VMM initializes MigTD (as a non- migratable TD) and binds it to the destination TD.	Once	TDH.SERVTD.BIND
		VMM/orchestration sets up transport session between source and destination MigTD. MigTDs setup their own protected channel.	Once	
		MigTDs reads the session encryption key, version information and other metadata from the destination TD and sends it to the MigTD on the source.	Once	TDG.SERVTD.RD
		MigTD receives the session encryption key from the MigTD on the source. MigTD writes it as the session decryption key to the destination TD. It may also write the migration version.	Once	TDG.SERVTD.WR
		VMM starts the import session and imports immutable state with a state migration bundle received from the source platform.	Once	TDH.IMPORT.STATE.IMMUTABLE
	Pre-Copy Memory	Host VMM builds the Secure EPT by allocating physical pages.	Multiple	TDH.MEM.SEPT.ADD
	Import	Host VMM imports TD private pages with a memory migration bundle received from the source platform.	Multiple	TDH.IMPORT.MEM

Section 1: Introduction and Overview

Migration Phase	Step	Description	Plurality	TDX Module Interface Functions
		Host VMM starts migration epoch with an epoch token migration bundle received from the source platform; a page can be imported once per epoch.	Once per migration epoch	TDH.IMPORT.TRACK(epoch token)
	Mutable TD- scope and VCPU-scope	VMM imports mutable TD-scope state with a state migration bundle received from the source platform.	Once	TDH.IMPORT.STATE.TD
	non-memory state import	VMM creates VCPU.	Per VCPU	TDH.VP.CREATE
		VMM allocates physical pages for the VCPU's TDVPS.	Multiple per VCPU	TDH.VP.ADDCX
		VMM imports mutable VCPU-scope state with a state migration bundle received from the source platform.	Per VCPU	TDH.IMPORT.STATE.VP
Out-Of- Order Import	Pre-Copy Memory Import	Host VMM starts the out-of-order import phase with a start token migration bundle received from the source platform.	Once	TDH.IMPORT.TRACK(start token)
		Host VMM builds the Secure EPT by allocating physical pages.	Multiple	TDH.MEM.SEPT.ADD
		Host VMM imports TD private pages with a memory migration bundle received from the source platform.	Multiple	TDH.IMPORT.MEM
	Post-Copy Memory Import	Host VMM commits the import session, allowing the TD to run on the destination platform.	Once	TDH.IMPORT.COMMIT
		Host VMM can execute the TD as usual. Memory can be imported on demand.	Per VCPU	TDH.VP.ENTER
		On EPT violation, host VMM requests a page import from the source platform.	Multiple	N/A
		Host VMM builds the Secure EPT by allocating physical pages.	Multiple	TDH.MEM.SEPT.ADD
		Host VMM imports TD private pages with a memory migration bundle received from the source platform.	Multiple	TDH.IMPORT.MEM
	End	Host VMM terminates the import session	Once	TDH.IMPORT.END

3.6. TD Import Abort

The following sequences are typically used to abort an import of TD to a destination platform, if an error is detected.

3.6.1. TD Import Abort During the In-Order Import Phase

5

Table 3.4: Typical TD Import Sequence Abort During In-Order Input

Migration Phase	Step	Description	Plurality	TDX Module Interface Functions

Migration Phase	Step	Description	Plurality	TDX Module Interface Functions
In-Order Import Pre-Copy Memory Import (Failed) Abort Token Transmission (Optional)	Pre-Copy Memory	Host VMM builds the Secure EPT by allocating physical pages.	Multiple	TDH.MEM.SEPT.ADD
	Import (Failed)	Host VMM imports TD private pages with a memory migration bundle received from the source platform. TDH.IMPORT.MEM returns an error status indicating a failed import session.	Multiple	TDH.IMPORT.MEM
	Abort Token Transmission (Optional)	VMM creates an abort token and transmits it to the source platform.	Once	TDH.IMPORT.ABORT
TD Teardown	End	Host VMM terminates the import session and tears down the TD on the destination platform	Once	TDH.MNG.VPFLUSHDONE TDH.PHYMEM.CACHE.WB TDH.MNG.KEY.FREEID TDH.PHYMEM.PAGE.RECLAIM TDH.PHYMEM.PAGE.WBINVD

3.6.2. TD Import Abort During the Out-Of-Order Import Phase

Table 3.5: Typical TD Import Sequence Abort During Out-of-Order Input

Migration Phase	Step	Description	Plurality	TDX Module Interface Functions
In-Order Import				
Out-Of-	Memory			
Order Import	(Failed)	Host VMM imports TD private pages with a memory migration bundle received from the source platform.	Multiple	TDH.IMPORT.MEM
		TDH.IMPORT.MEM returns an error status indicating a failed import session.		
	Abort Token Transmission	VMM creates an abort token and transmits it to the source platform.	Once	TDH.IMPORT.ABORT
TD	End	Host VMM terminates the import session	Once	TDH.MNG.VPFLUSHDONE
Teardown		and tears down the TD on the destination		TDH.PHYMEM.CACHE.WB
		platform		TDH.MNG.KEY.FREEID
				TDH.PHYMEM.PAGE.RECLAIM
				TDH.PHYMEM.PAGE.WBINVD

5 **3.7. TD Export Abort**

The following sequence is typically used to export a TD from a source platform, if the export is aborted.

3.7.1. Export Abort During the In-Order Export Phase

Migration Phase	Step	Description	Plurality	TDX Module Interface Functions
In-Order	Start of			
Export Export Session Export Abo	Export Session	VMM starts the export session and exports immutable state creating a state migration bundle.	Once	TDH.EXPORT.STATE.IMMUTABLE
	Export Abort	Host VMM aborts the export session.	Once	TDH.EXPORT.ABORT
Source TD Run and		Host VMM may run and manage the source TD	Multiple	TDH.VP.ENTER etc.
Restore		Host VMM restores SEPT entries to their normal non-export state	Multiple	TDH.EXPORT.UNBLOCKW TDH.EXPORT.RESTORE

Table 3.6: Typical TD Export Sequence Abort During In-Order Export

3.7.2. Export Abort During the Out-Of-Order Export Phase

Table 3.7: Typical TD Export Sequence Abort During Out-Of-Order Export

Migration Phase	Step	Description	Plurality	TDX Module Interface Functions
In-Order Export				
Out-Of-				
Order Export	Export Abort	Host VMM receives an abort token from the destination platform and abort the export session.	Once	TDH.EXPORT.ABORT
Source TD Run and		Host VMM may run and manage the source TD	Multiple	TDH.VP.ENTER etc.
Restore		Host VMM restores SEPT entries to their normal non-export state	Multiple	TDH.EXPORT.UNBLOCKW TDH.EXPORT.RESTORE

5

SECTION 2: TD MIGRATION ARCHITECTURE SPECIFICATION

4. Migration TD, Migration Policy and the Extended TCB

This chapter describes the role of the Migration TD, the migration policy it implements and the extended TCB for migrated TDs. For details of the Intel reference MigTD, refer to the [MigTD Spec].

4.1. Extended TCB and the Migration Policy

5 The TCB of a migrated TD is extended to include the Migration TD bound to it. The Migration TD builds trust relationships with to other Migration TDs, and through them, with other platforms and their component. Thus, it can be said that **a migrated TD's effective TCB extends to the whole migration pool of platforms**, each with its own Migration TD, TDX modules etc.



10

Figure 4.1: Trust Relationship with Migration TDs

This extended TCB is controlled by the **migration policy** implemented by the Migration TD. TDX does not enforce any specific migration policy; this is up to the migration TD's writer.

4.2. Attestation of the Migration TD and its Migration Policy

- A migrated TD's attestation reflects the current platform on which it is running. I.e., calling TDG.MR.REPORT generates a TDREPORT_STRUCT with CPUSVN and TEE_TCB_INFO for the current platform. However, there is no requirement to re-attest a TD after migration. The initial attestation, as well as attestation done at any time during the migrated TD's lifetime, contains the information about the TD's extended TCB, provided as a measurement of the Migration TD and its migration policy.
- A Migration TD is a private case of a **Service TD**, which is described in the [Base Spec]. As such, a hash of the bound Migration TD's TDREPORT_STRUCT is included in the migrated TD's TDREPORT_STRUCT (as the SERVTD_HASH field). The service TD protocol allows, at binding time, to select which field of the TDREPORT_STRUCT is included in the calculation of SERVTD_HASH.

The static components of the Migration TD, i.e., code and data added and measured during its build, are measured by its MRTD. This may include, for example, a baseline migration policy. Migration TD writers are expected to measure the

25 configurable part of the migration policy and its parameters, i.e., any change that can be made after the Migration TD build was finalized, using one or more of the RTMRs (run time measurement registers) provided by the TDX architecture via the TDG.MR.RTMR.EXTEND interface function. Those RTMRs should be included in the SERVTD_HASH calculation. Thus, the migration policy is attested as part of the Migration TD, and as such, as part of the migrated TD.

4.3. Inputs to the Migration TD's Migration Policy Evaluation

The following data can be read and evaluated by the Migration TD, as part of its migration policy evaluation. The migration TD may reflect such inputs in an RTMR (using TDG.MR.RTMR.EXTEND) so that it becomes part of any migrated TD's attestation.

System-Scope Information Provided by TDG.SYS.RD*

System-scope (platform and TDX module) information can be read using TDG.SYS.RD or TDG.SYS.RDALL. Examples of such information are:

- TDX features supported on this platform (e.g., TDX_FEATURESO field)
 - Supported TD features (e.g., ATTRIBUTES_FIXED0/1, XFAM_FIXED0/1, CPUID_CONFIG_VALUES and IA32_ARCH_CAPABILITIED_CONFIG_MASK fields)
 - TD Migration protocol features (e.g., MIN/MAX_EXPORT/IMPORT_VERSION fields)

The complete list of fields enumerated by TDG.SYS.RD* is provided in the [ABI Spec].

15 System-Scope Information Provided by TDG.MR.REPORT

The Migration TD may call TDG.MR.REPORT to get the TDREPORT_STRUCT, which contains information about the platform (CPUSVN) and the TDX module (TEE_TCB_INFO).

Migration Policy Configuration

The Migration TD's migration policy may be configured by, e.g., the host VMM.

20 4.4. Migrated TD Information Provided by TDG.SERVD.RD

Migration TDs on both sides may read migrated TD information using TDG.SERVD.RD, in order to decide whether that TD can be migrated to a specific destination platform. Note that similar checks are also done by the TDX module on import, so there's no strict requirement for the Migration TD to do them.

Examples of such information are:

- 25 ATTRIBUTES_FIXED0/1
 - XFAM_FIXED0/1
 - GPAW

35

40

5

CPUID_VALUES

4.5. Migration Protocol Version Setup

30 The migration protocol supports versioning, to allow for future updates.

Before starting a migration session, the MigTDs on the source and destination should agree on migration protocol version that is supported by both sides. To do so, each MigTD can read the following fields using TDG.SYS.RD:

- MIN_EXPORT_VERSION
- MAX_EXPORT_VERSION
- MIN_IMPORT_VERSION
- MAX_IMPORT_VERSION

The MigTD on each side then should write the MIG_VERSION to the migrated TD using TDG.SERVD.WR.

4.6. Migration Session Keys (MSKs) Exchange

- 1. The MigTD on each side reads the migration encryption key (MIG_ENC_KEY using TDG.SERVTD.RD. The TDX module generates a new key on each such read operation.
- 2. The MigTD on each side sends the key value over a secure channel to the peer MigTD on the other side.

- **Note:** A MigTD must never send the same key value to more than one peer MigTD. Failing to do so may expose TDX to an attack where the same TD is migrated to more than one destination. The MigTD should always read a new MIG_ENC_KEY value using TDG.SERVTD.RD if it needs to resend it.
- 3. The MigTD on each side writes the received key value as the migration decryption key (MIG_DEC_KEY), using TDG.SERVTD.WR.

4.7. Example Migration Session Establishment

Note: The example below illustrates migration session establishment and doesn't necessarily imply actual implementation.

The goal is to stablish secure transport channel between Intel TDX Modules and MigTDs on both sides across compatible platforms and reserve resources for the migration session.

- 1. On the source platform, TD-s to be migrated is created with MIGRATABLE attribute.
- 2. On the source platform, MigTD-s is created and built.
- 3. On the source platform, MigTD-s evaluates the platform and TDX module information using TDG.SYS.RD* and TDG.MR.REPORT and reflects its updated migration policy in the MigTD's RTMR[3], using TDG.MR.RTMR.EXTEND.
- The host VMM on the source platform binds MigTDs to TD-s via an unsolicited Service TD binding using TDH.SERVTD.BIND:
 - 4.1. MigTD-s requests the host VMM to be bound to TD-s, using TDG.VP.VMCALL.
 - 4.2. The VMM invokes TDH.SERVTD.BIND to bind TD0 to MigTD-s. As a result, the MigTD's TDREPORT fields' hash is included in the migrated TD's TDREPORT (as SERVTD_HASH).
 - 4.3. TDH.SERVTD.BIND returns the migrated TD's binding handle and TD_UUID. The host VMM communicates them to MigTD-s.
 - 5. The host VMM on the source platform finalizes the build of TD-s using TDH.MR.FINALIZE.
 - 6. TD-s executes on the source platform.
 - 7. Similar to MigTD-s above, on the destination platform, MigTD-d is created and built. It evaluates its platform and TDX module and reflects the updated migration policy in RTMR[3].
 - 8. The migration orchestrator triggers a migration of TD-s from the source platform to the destination platform.
 - 9. MigTD-s may read selected metadata of TD-s (e.g., its ATTRIBUTES and XFAM) using TDG.SERVTD.RD, to use as an input to the migration policy evaluation.
 - 10. The migration orchestrator requests the source and destination Mig TDs to establish a secure session and negotiate the migration policy agreement.
 - 11. MigTD-s verifies the MigTD-d quote and establishes a secure transport key for the session with the destination platform (and vice-versa).
 - 12. MigTD-s authenticates the Migration Policy and evaluates it per the capabilities (SVN etc.) of the destination platform (learnt via the quote) for the specified live migration session.
 - 13. On the destination platform, a destination TD-d skeleton is created via legacy process.
 - 14. On the destination platform, the host VMM may bind MigTD-d to TD-d using TDH.SERVTD.BIND.
 - 15. Migration Keys Exchange:
 - 15.1. MigTD-s reads the source migration encryption key using TDG.SERVTD.RD. This key is used as the forward migration session key, to encrypt information exported by the TDX module.
 - 15.2. MigTD-s sends the forward migration session key to MigTD-d.
 - 15.3. On the destination platform, MigTD-d writes the forward migration session key, as the migration decryption key, using TDG.SERVTD.WR.
 - 15.4. MigTD-d reads the destination migration encryption key using TDG.SERVTD.RD. This key is used as the backward migration session key, to encrypt information sent by the TDX module on the destination.
 - 15.5. MigTD-d sends the backward migration session key to MigTD-s.
 - 15.6. On the destination platform, MigTD-d writes the backward migration session key, as the migration decryption key, using TDG.SERVTD.WR.
 - 16. The host VMM on the source platform can now initiate the state export via TDH.EXPORT* SEAMCALLs and import state via TDH.IMPORT* SEAMCALLs

5

15

20

25

30

35

40

45

5. Updated: Common TD Migration Mechanisms

This chapter describes the infrastructure used by all Import/Export APIs to migrate TD private memory and metadata.

5.1. Migration Bundles

A migration bundle is the basic unit of information transported between the source and destination platforms. This section describes the generic migration bundle structure. Private memory migration uses an enhanced format, described in 8.3.

5.1.1. Overview

10

TD information is transported from the source platform to the destination platform in **migration bundles**. A migration bundle consists of **migration data**, which may span one or more 4KB pages, and **migration bundle metadata (MBMD)**. The host VMM may add its own untrusted metadata to the migration bundle, for managing the migration. Migration

bundle transport is the responsibility of untrusted software and is out of the scope of this specification.





5.1.2. Migration Data

- 15 Migration data contains either TD private memory contents or TD non-memory state. It is confidentiality-protected using AES-GCM with the TD migration key and a running migration session counter. Migration data is integrity-protected by its associated MBMD. For encryption details, see 5.3.
 - **Note:** Migration of shared memory pages is the responsibility of untrusted software and is out of the scope of this specification.
- In memory, migration data occupies one or more 4KB shared memory pages, managed by the host VMM.

5.1.3. Migration Bundle Metadata (MBMD)

A migration bundle metadata (MBMD) structure provides metadata for the migrated data in the migration bundle. In memory, MBMD resides in a shared page, managed by the host VMM, and must be naturally aligned.

An MBMD is not confidentiality-protected. The host VMM can read the MBMD; this is required in order for the VMM to perform the required operations. E.g., the host VMM on the destination platform reads the MBMD to decide which import function to call.

The MBMD provides integrity protection for itself and for its associated migration data.

The MBMD structure consists of a fixed header and a per-type variable part. The header contains the following fields:

SIZE: Overall size of the MBMD structure, in bytes

VIG_VERSION:	Migration protocol version
_	

MB_TYPE: The type of information being migrated

MB_COUNTER: Per-stream migration bundle counter

MIG_EPOCH: Migration epoch number

MIGS_INDEX: Index of the migration stream

IV_COUNTER: Monotonically increasing counter, used as a component in the AES-GCM IV

The last field of each MBMD is an AES-256-GCM MAC over other MBMD fields and other associated migration data (migration pages).

The detailed MBMD definition is provided in [TDX Module ABI Spec].

10 5.1.4. Untrusted Metadata

5

The host VMM may add its own metadata to migration bundle, e.g., to support its implementation of the migration protocol. This metadata is untrusted and is not used by the TDX module.

5.2. Updated: Export and Import Functions Interface

On each invocation, export and import functions operate on a single migration bundle and a specific migration stream.

15 **5.2.1.** Migrating a Multi-Page Migration Bundle

Migration bundles may consist of multiple pages of migrated information. To export a multi-page migration bundle, the host VMM on the source platform prepares a set of buffers in shared memory. The host VMM provides the MBMD's HPA and a list of HPA pointers to the migration pages as an input to the TDH.EXPORT* function. The required number of migration pages per TDH.EXPORT.* function is enumerated by the TDX module. See the [ABI Spec] for details.

To import a multi-page migration bundle, the host VMM on the destination platform prepares the set of migration pages and the MBMD, as received from the source platform, in shared memory. The host VMM provides the MBMD's HPA and a list of HPA pointers to the migration pages as an input to the TDH.IMPORT* function.



Figure 5.2: Migrating a Multi-Page Migration Bundle

5.2.2. Updated: Migration Functions Interruptibility

TDH.EXPORT.* and TDH.IMPORT.* functions may take relatively long time to execute. This is especially true for memory migration, which can process up to 512 4KB pages. To avoid latency issues, such functions may be **interruptible** and **resumable**. This is supported as follows:

- TDH.EXPORT.* and TDH.IMPORT.* functions are designed to synchronously check for a pending external event by reading MSR_INTR_PENDING (once after every pre-determined number of cycles, chosen to be smaller than the maximum allowed cycle latency).
 - As described later, migration functions that work with a Migration Stream use Migration Stream Context (MIGSC). If an external event is pending, the functions store their context in the proper MIGSC and return with a TDX_INTERRUPTED_RESUMABLE completion status.
 - The host VMM is expected to call the TDH.EXPORT.* or TDH.IMPORT.* function again with the same set of inputs until the operation is completed successfully (completion status is TDX_SUCCESS) or some error occurs (completion status indicates an error).

15

10

5

• An input flag indicates whether the invocation of a TDH.EXPORT.* or TDH.IMPORT.* function starts a new operation (and possibly aborts an interrupted one) or resumes an interrupted operation. A migration function which is called as a resumption of an interrupted operation checks to see if an intermediate state has been saved, and if so, it checks that it is being invoked with the same input arguments as last time when it was interrupted.

5.3. Cryptographic Protection of Migration Data

5.3.1. Encryption Algorithm

- 20 TD migration uses AES in Galois/Counter Mode (GCM) to transfer state between the source and destination platform platforms. Per [AES-256-GCM] definitions, the TD data private memory or non-memory state temporarily held in the CPU cache during TDH.EXPORT.* forms the "Plaintext", and some of the MBMD fields form the "Additional Authenticated Data". The "Plaintext" is encrypted using a Migration key (described below). The MAC size, also known as t, as defined in [AES-256-GCM], must be 128 bits.
- The Initialization Vector (IV) is 96 bits. It is composed as described below. Since 64 bits will never wrap around in practice, this helps ensure a unique counter for each stream.

Bits	Size	Name	Description
63:0	64	IV_COUNTER	Starts from 1, incremented by 1 every time AES-GCM is used to encrypt data and/or generate a MAC for a migration bundle. The counter is incremented even if the data is discarded and not used for migration.
79:64	16	MIGS_INDEX	Stream index (see 5.3)
95:80	16	RESERVED	Set to 0

Table 5.1: Components of the 96-bits IV

5.3.2. Migration Session Keys

- 30 Two migration session keys are used, one in each direction:
 - The TDX module on the source platform generates a **migration session forward key** for encrypting migration bundles by the source TDX module and decrypting them by the destination TDX module.
 - The TDX module on the destination platform generates a **migration session backward key** for encrypting migration bundles by the destination TDX module and decrypting them by the source TDX module.
- ³⁵ Each of the MigTDs on the source and the destination platforms reads the key generated on its side, known as the **migration encryption key**, from the TDX module's migrated TD's metadata, and transfers it using a secure connection to its peer MigTD on the other side of the migration session. The peer MigTD then writes the key, to be used as the **migration decryption key**, to the TDX module's migrated TD's metadata. TD metadata read and write use the Service TD protocol, as described in the [TDX Module Base Spec].



Figure 5.3: Migration Session Keys Exchange

The migration keys properties are as follows:

- The key strength is 256 bits.
- A new encryption key is generated by the TDX module on TD creation (TDH.MNG.CREATE) and at the beginning of each migration session (TDH.*PORT.STATE.IMMUTABLE).
 - The encryption key is read by the MigTD from the migrated TD's metadata using the Service TD metadata protocol, as described in the [TDX Module Base Spec].
- The MigTD transfer the encryption key to its peer MigTD, over a secure channel both MigTDs have created.
- The decryption key is written by the MigTD to the migrated TD's metadata using the Service TD metadata protocol, as described in the [TDX Module Base Spec].
 - The keys are accessible only by the MigTD and the Intel TDX Module.
 - On migration session start by TDH.*PORT.STATE.IMMUTABLE, the Intel TDX module copies the encryption and decryption keys into working keys that are used throughout the session.
 - The keys are used by TDH.EXPORT.*/TDH.IMPORT.* to control the migration session and migrate TD memory and non-memory. The migration stream AES-GCM protocol requires that state is migrated in-order between the source and destination platform. This helps guarantee the order within each migration stream.
 - The keys are **destroyed** when a TD holding them is torn down, or when new keys are generated or programmed.

5.4. Migration Streams and Migration Queues

20 **Migration stream** is a **TDX concept**. Multiple streams allow multi-threaded, concurrent export and import, and enable the Intel TDX Module to enforce proper ordering of migration bundles during the in-order phase where this is essential.

Migration queue is a **host VMM concept**. Multiple queues allow QoS and prioritization. E.g., Post-copy of pages on demand (triggered by an EPT violation on the destination platform) may have a higher priority than other post-copy of pages. To avoid head-of-line blocking by waiting in the same queue as lower priority pages, a separate high priority queue can be used by the host VMM.

25

5

15

From the migration streams and migration queues perspective, a migration session is divided into two main phases:

- **In-order**, where the source TD may run, and its memory and non-memory state may change. During the in-order phase, the order of memory migration is critical. A newer export of the same memory page must be imported after an older export of the same page. Furthermore, for any memory page that has been migrated during the in-order phase, the most up-to-date version of that page must be migrated before the in-order phase ends. In the in-order
- 30
- phase, one or more migration streams are mapped to each migration queue.
- **Out-of-order**, where the source TD does not run, and its memory and non-memory state may not change. During out-of-order, the order of memory migration is not important except that migration bundles exported during the

in-order phase can't be imported during the out-of-order phase. Furthermore, the host VMM may assign exported pages (even multiple copies of the same exported page) to different priority queue. This is used, e.g., for on-demand migration after the destination TD starts running.

The **start tokens**, generated by TDH.EXPORT.TRACK(DONE) and verified by TDH.IMPORT.TRACK, serve as markers to indicate the end of the in-order phase and start of the out-of-order phase. They are used to implement a rendezvous point, enforcing all the in-order state (across all streams) to have been imported before the out-of-order phase starts and the destination TD may execute.

5.3 describes how the stream context held by the source and the destination platforms and the MBMD fields included in each migration bundle are used to construct the non-repeated AES-GCM IV. Note also that the same stream queues can be used for both in-order and out-of-order. The semantic use of the queues is up to the host VMM.



Figure 5.4: Migration Streams

Migration streams have the following characteristics:

- Within each stream, state is migrated in-order. This is enforced by the MB_COUNTER field of MBMD.
- Export or import operations using a specific migration stream must be serialized. Concurrency is supported only between streams.
 - The host VMM should use the same stream index to import memory on the destination TD (which should be in MEMORY_IMPORT, STATE_IMPORT or RUNNABLE state). This is enforced by TDH.IMPORT.MEM.
- Non-memory state can only be migrated once; there is no override of older migrated non-memory state with a newer one. Ordering requirements (e.g., TD-scope non-memory state must be imported before VCPU non-memory state) are enforced by the lifecycle state machine, as described in 6.2.
 - The maximum number of forward streams is implementation dependent:
 - \circ \quad Each stream requires context space allocation.
 - Stream ID requires a field in the MBMD header.
- There is one backward stream.

15

20

10

Draft

- Draft
- The maximum total number of forward and backward migration streams is enumerated by MAX_MIGS, readable TDH.SYS.RD*.

5.5. Measurement and Attestation

5.5.1. TD Measurement Registers Migration

- 5 TDs have two types of measurement registers:
 - MRTD: Static measurement of the TD build process and the initial contents of the TD. This state is migrated as part of the global immutable state of the TD (via TDH.EXPORT.STATE.IMMUTABLE and TDH.IMPORT.STATE.IMMUTABLE).
- RTMR: An array of general-purpose measurement registers, available to the TD software for measuring additional logic and data loaded into the TD at runtime. Since this measurement covers dynamic state beyond the static state and can be extended by TD software via TDG.MR.RTMR.EXTEND, this state is migrated only during the blackout period, as part of the TD's mutable state (via TDH.EXPORT.STATE.TD and TDH.IMPORT.STATE.TD).

All TD measurements are reflected in TD attestations.

15 **5.5.2. TD Measurement Reporting Changes**

The TDINFO structure is enhanced to include hashes of Service TDs' TDINFO; for TD migration, the applicable Service TD is the Migration TD. Refer to the [TDX Module Base Spec] for a discussion of Service TDs.

5.5.3. TD Measurement Quoting Changes

To create a remotely verifiable attestation, the TDREPORT_STRUCT must be converted into a Quote signed by a certified Quote signing key, as described in the [TDX Module Base Spec].

TDREPORT_STRUCT is HMAC'ed using an HMAC key unique to each platform and accessible only to the CPU. This protects the integrity of the structure and can only be verified on the <u>local platform</u> via the SGX ENCLU(EVERIFYREPORT2) instruction. TDREPORT_STRUCT cannot be sent off platforms for verification; it first must be converted into signed Quotes.

²⁵ If a report is generated by TDG.MR.REPORT on the source platform, but the TD is migrated to a destination platform, the local HMAC key is different and hence the EVERIFYREPORT2 on the migrated TDG.MR.REPORT is expected to fail. The TD software is typically unaware of being migrated. It is expected to retry the TDG.MR.REPORT operation if it fails.

5.5.4. TCB Recovery and Migration

The TDX architecture has several levels of TCB:

- CPU HW level, which includes microcode patch, ACMs etc.
- Intel TDX Module

30

35

• Attestation Enclaves, which include the TD Quoting Enclave and Provisioning Certification Enclave

The TCB Recovery story is different for each level. The existing TCB Recovery model for CPU SW level items applies similarly to TDX and SGX and requires a restart of the platform to take effect. The Intel TDX Module can be unloaded and reloaded to reflect an upgraded Intel TDX Module. The enclaves can be upgraded at runtime, but if the PCE is upgraded, a new certificate must be downloaded.

5.6. Updated: TDX Control Structure Updates

This section discusses updates and additions to the global and TD-scope control structures.

5.6.1. MIGSC: Migration Stream Context

40 Migration streams are defined in 5.3.

MIGSC (Migration Stream Context) is an opaque control structure that holds migration stream context. MIGSC occupies a single 4KB physical page; it is created using the TDH.MIG.STREAM.CREATE function. MIGSC can only be created if a migration session is not in progress.


Figure 5.5: TD-Scope Control Structures Overview

6. Updated: Migration Session Control and State Machines

This chapter discusses the TD migration session control, state machine and messaging protocol.

6.1. Updated: Overview

6.1.1. Updated: Pre-Migration

- 5 Prior to starting a migration session, the following should have happened:
 - On the destination platform, the TD has been created as a skeleton (control structure pages only).
 - Migration TDs should be bound as service TDs on both source and destination platforms.
 - Migration TDs must have exchanged the migration session keys and decided on a migration protocol version.



10

15

Figure 6.1: Updated: Pre-Migration

6.1.2. Updated: Successful Migration Session

Figure 6.3 below shows an overview of a successful migration session. This figure shows the following:

- Migration session control interface functions (TDH.EXPORT.PAUSE, TDX.EXPORT.TRACK etc.)
- States of the TD Operation State Machine (RUNNABLE, LIVE_EXPORT etc.) are also shown. The state machine itself is discussed in 6.2.4.2 below.
- Phases of the export and import sessions



Figure 6.2: Updated: Migration Session In-Order Phase (Success Case)

On the source platform, an export session's in-order export phase starts with the host VMM invoking the TDH.EXPORT.STATE.IMMUTABLE function. This function creates a migration bundle that is transmitted by the host VMM

5 to the destination platform, where TDH.IMPORT.STATE.IMMUTABLE is invoked to start the import session's in-order import phase.

TDH.EXPORT.PAUSE pauses the TD on the source platform and starts the TDX-imposed blackout period. If the TD is configured with TDX Connect enabled, the TDX module checks that all non-migratable assets have been released (e.g., MMIO mappings, TDIs).

- 10 TDH.EXPORT.TRACK(DONE), when invoked on the source platform, verifies proper in-order export:
 - The TD's non-memory state must have been exported.
 - If any TD private page has been exported, the latest version of that page must have been exported.
 - If the TD is configured with TDX Connect enabled, all TD private pages must have been exported.

TDH.EXPORT.TRACK(DONE) ends the in-order memory export phase and creates a **start token** migration bundle that is transmitted by the host VMM to the destination platform, where TDH.IMPORT.TRACK(DONE) is invoked to end the inorder import phase. See also the discussion of migration epochs in 6.1.4 below.



Figure 6.3: Migration Session Out-Of-Order Phase (Success Case)

TDH.IMPORT.COMMIT, invoked on the destination platform, commits the migration session and enables the TD to run on it, while memory import may continue. This also helps ensure that the TD will not run on the source platform, since an abort token can no longer be generated.

Optionally, TDH.IMPORT.END, invoked on the destination platform, commits the migration session and enables the TD to run on it if not already done by TDH.IMPORT.COMMIT. TDH.IMPORT.END ends the migration session; memory import

6.1.3. Aborted Migration Session

is no longer allowed.

5

10 6.1.3.1. Abort during the In-Order Phase

Figure 6.4 below shows a case where a migration session is aborted during the in-order migration phase. TDH.EXPORT.ABORT, invoked by the host VMM, terminates the export session and enables the TD to resume running on the source platform. By design, the TD should not be able to run on the destination platform – it is up to the host VMM to free up any resource allocated there.





6.1.3.2. Abort during the Out-Of-Order Phase

Figure 6.5 below shows a case where a migration session is aborted during the out-of-order migration phase. TDH.IMPORT.ABORT is invoked by the host VMM on the destination platform. This function terminates the import session and puts the TD in a state where, by design, it is able to run – it is up to the host VMM to free up any resource allocated there. TDH.IMPORT.ABORT also creates an abort token, which is transmitted by the host VMM back to the source platform.

On the source platform, the host VMM invokes TDH.EXPORT.ABORT, which checks the validity of the abort token and enables the TD to resume running.



Figure 6.5: Migration Session Control Overview (Abort During the Out-Of-Order Phase)

Draft

6.1.4. Updated: Migration Epochs

As described in 5.4, within each migration stream, proper ordering is maintained by the migration bundle counter (MB_COUNTER) of each MBMD. However, there is no intrinsic guarantee of ordering across migration streams.

To help ensure overall ordering, the migration session is divided to **migration epochs**. A given page can only be imported, or its import can be cancelled, once per migration epoch. An **epoch token**, generated on the source platform by TDH.EXPORT.TRACK, serves as an epoch separator. It provided the total number of migration bundles exported so far. This helps TDH.IMPORT.TRACK, which imports the epoch token, check that all migration bundles of the previous epoch have been received. No migration bundle of an older epoch may be imported.

The **start token**, which starts the out-of-order phase, is a special version of the epoch token. Epoch number 0xFFFFFFFF indicates the out-of-order phase.

Notes

- Do not confuse TDH.MEM.TRACK (which is used for TLB tracking) with TDH.EXPORT.TRACK (which rendezvous all migration streams).
- Migration epoch is a TDX concept. It roughly corresponds to **migration round** (or **migration pass**) which is a usage concept.

15

5



Figure 6.6: Migration Epochs Overview

6.2. Updated: Migration Session Control

6.2.1. Migration TD Binding and Migration Key Assignment

- 20 Migration TD binding (using TDH.SERVTD.BIND) must happen before a migration session can start. This may happen during TD build, before the measurement has been finalized (by TDH.MR.FINALIZE). Alternatively, pre-binding (using TDH.SERVTD.PREBIND) can be done during TD build, and actual binding can happen later. On the destination platform migration TD binding and TD import must happen before the TD is initialized (by TDH.MNG.INIT).
- Migration key assignment, done by TDG.SERVTD.WR, may happen at any time after migration TD binding, except during the PAUSED_EXPORT and POST_EXPORT states. A new migration key must be written for any migration session.

6.2.2. Updated: Export Side (Source Platform)

To begin an export session, the TD's OP_STATE must either be RUNNABLE, indicating that its measurement has been finalized (by TDH.MR.FINALIZE), or LIVE_IMPORT, indicating that this TD has been previously imported.

An export session begins with immutable TD state export (using TDH.EXPORT.STATE.IMMUTABLE). This function copies the migration key to a working migration key. It then starts the **in-order export phase**. It transitions the OP_STATE to LIVE_EXPORT, allowing the source TD to continue running normally while private memory is being exported.

TDH.EXPORT.PAUSE transitions the source TD's OP_STATE into the PAUSED_EXPORT state. In this state, TD private memory and TD non-memory state modification are prevented. None of the TD VCPUs may be running (i.e., in TDX non-

root mode), and no host-side (SEAMCALL) function is allowed to change any TD non-memory state that is to be exported. For TDX Connect, all non-migratable assets have been released (e.g., MMIO mappings, TDIs). Memory export (via TDH.EXPORT.MEM etc.) may still continue. Per-TD and per-VCPU mutable control state are exported using TDH.EXPORT.STATE.TD and TDH.EXPORT.STATE.VP respectively.

- 5 At any time, the export may be aborted by the host VMM using TDH.EXPORT.ABORT, which returns the source TD to the RUNNABLE state, where it can continue to run normally. No abort token is required at this phase since no start token has been generated and the destination TD, by design, should not be able to run.
 - **Note:** TDH.EXPORT.STATE.TD is expected to be called by the exporting host VMM prior to TDH.EXPORT.STATE.VP, but this is only enforced on the import side.
- 10 TDH.EXPORT.TRACK(DONE) generates a **start token** which the host VMM transmits to the destination VMM. It transitions the source TD OP_STATE into the POST_EXPORT state, starting the **out-of-order export phase**. If the TD has not been configured with TDX Connect enabled, memory export (TDH.EXPORT.MEM) may continue, to support the out-of-order stage of the TD live migration.

In the TD Migration Session POST_EXPORT state, a TDH.EXPORT.ABORT with a valid **abort token**, received from the destination VMM, indicates that the TD, by design, should not be able to run on the destination platform. It terminates the export session and returns the source TD to the RUNNABLE state, where it can continue to run normally.

The host VMM can start tearing down the source TD at any time, by ensuring that no VCPU is associated with an LP (i.e., by executing TDH.VP.FLUSH for all VCPUs) and issuing TDH.MNG.VPFLUSHDONE. Typically, it will do so in the after it gets a notification from the destination platform that import has been successful.

20 6.2.3. Import Side (Destination Platform)

Migration TD binding (using TDH.SERVTD.BIND) and migration key assignment (using TDG.SERVTD.WR) must happen in the UNINITIALIZED state, where TDCS memory has already been allocated but the destination TD has not been initialized yet. This is required since the destination TD is going to be initialized by importing immutable state from the source TD.

TDH.IMPORT.STATE.IMMUTABLE starts the **in-order import phase**. It initializes the destination TD's TDCS with imported immutable state and transitions the destination TD's OP_STATE into MEMORY_IMPORT. In this state, TD private memory can be imported using TDH.IMPORT.MEM etc.

TDH.IMPORT.STATE.TD imports the per-TD mutable state and transitions the destination TD's OP_STATE into STATE_IMPORT. In this state, mutable VCPU state can be imported using TDH.IMPORT.STATE.VP. TD private memory import also continues.

30 Upon executing TDH.IMPORT.TRACK with a valid **start token** as operand, the destination TD's OP_STATE transitions into the POST_IMPORT state, starting the **out-of-order import phase**. Memory import (a.k.a. **post-copy**) may continue, but pages can only be imported if their GPA is free (i.e., the Secure EPT state is FREE).

An import failure up to this point, e.g., improper sequence of page import vs. alias import, or executing TDH.IMPORT.TRACK with a bad start token received from the source platform, transitions the TD's OP_STATE to the

³⁵ FAILED_IMPORT state. In addition, the host VMM can explicitly abort the import by using TDH.IMPORT.ABORT. In the FAILED_IMPORT state, the TD is designed not to run; it can only be torn down. TDH.IMPORT.ABORT generates an **abort token**, which can be transmitted to the source platform.

TDH.IMPORT.COMMIT transitions the destination TD's OP_STATE transitions into the LIVE_IMPORT state. In this state, the destination TD may run normally. Out-of-order memory import may continue as long as the destination TD is in the LIVE_IMPORT state. An input failure in the LIVE_IMPORT terminates the import session; it transitions the TD's OP_STATE to the RUNNABLE state, where the TD can continue running normally. An abort token can no longer be generated.

TDH.IMPORT.END ends the import session and transitions the destination TD's OP_STATE into the RUNNABLE state. This transition is optional if TDH.IMPORT.COMMIT has already been executed; it removes any limitations on TD memory management that exist during the out-of-order import phase.

45 A new export session (TDH.EXPORT.STATE.IMMUTABLE) terminates a previous out-of-order import.

6.2.4. Details: Migration State Machine

This section provides a detailed view of migration session control state machine. Details may be of interest to host VMM programmers who require deeper understanding of TD Migration.

6.2.4.1. Details: Reminder: TD Lifecycle State Machine

The whole TD migration process happens within the TD_KEYS_CONFIGURED state of the TD life cycle state machine, where an HKID has been assigned to the TD and the keys have been configured on the hardware. As a reminder, the TD life cycle state diagram is shown in Figure 6.7 below. For details, see the [TDX Module Base Spec].



5

Figure 6.7: TD Life Cycle State Diagram

Within the TD_KEYS_CONFIGURED state, a secondary-level **TD Operation state machine** controls the overall TD operation, including migration.

6.2.4.2. Details: OP_STATE: TD Operation State Machine

10 The TD Operation state machine is shown in Figure 6.8 below. The baseline state machine is extended with new migration-related states and transitions, highlighted in red text and lines. The export states are highlighted in purple, and the import states are highlighted in blue.



Figure 6.8: TD Operation State Machine (Sub-States of TD_KEYS_CONFIGURED)

6.2.4.3. Details: OP_STATE Summary

Table 6.1: OP	_STATE Sub-	States of TD	_KEYS_	CONFIGURED
---------------	-------------	--------------	--------	------------

Sub-State	Source / Destination	Description
UNALLOCATED	Both	TDCS memory is being allocated.
UNINITIALIZED	Both	 TDCS is pending initialization. On the destination platform, migration TD binding and migration key assignment must happen in this state.

Sub-State	Source / Destination	Description
INITIALIZED	Source	 TD is being built. Memory is added and measured. VCPUs are created. Migration TD binding may happen in this state.
RUNNABLE	Both	• TD can run.
LIVE_EXPORT	Both	 TD can run. Immutable non-memory state (TDH.EXPORT.STATE.IMMUTABLE) has been exported. Live memory can be exported (TDH.EXPORT.MEM etc.).
PAUSED_EXPORT	Source	 No TD VCPU may run. TD memory can't be written. Memory can be exported. Mutable non-memory state is being exported.
POST_EXPORT	Source	 Start token has been generated. Mutable control state has been exported. No TD VCPU may run. TD memory can't be written. Memory can be exported (post-copy).
MEMORY_IMPORT	Destination	• TD memory can be imported.
STATE_IMPORT	Destination	 TD memory can be imported. TD VCPU non-memory state can be imported
POST_IMPORT	Destination	TD memory can be imported (post-copy).
LIVE_IMPORT	Destination	TD VCPUs may run.TD memory can be imported (post-copy).
FAILED_IMPORT	Destination	Destination TD will not run.

6.3. Migration Tokens

Migration tokens are transmitted from the source platform to the destination platform and vice versa as part of the migration session control.

5 An **epoch token** is generated by TDH.EXPORT.TRACK. it serves as a separator between migration epochs. A **start token** is a special version of an epoch token which starts the out-of-order phase. The start token helps ensure that no newer version of any memory page exported prior to the start token exists on the source platform.

The **abort token** is generated by TDH.IMPORT.ABORT on the destination platform if import fails for any reason. It helps ensure that the TD will not run on the destination platform, and therefore may be restored on the source platform.

10 Migration tokens are formatted as Migration Bundles, with only an MBMD. Its format is defined in the [TDX Module ABI Spec].



Figure 6.9: Migration Tokens

6.4. Migration Protocol Versioning

6.4.1. Introduction

10

- 5 Migration protocol version number is provided as part of the MBMD header. Migration protocol changes may require migration version increment and may impact source and destination compatibility. For example, this may happen due to:
 - Incompatible MBMD format changes
 - New values of MBMD fields
 - New memory migration variants (e.g., support of aliases for VM nesting)
 - Incompatible migration session state machine changes

Non-memory state (metadata) migration changes may also require migration version increment. For example, this may happen due to:

- New exported metadata fields
- New values or format of exported metadata fields

6.4.2. Enumeration of Supported Migration Versions

TDX module enumerates supported migration versions using global metadata fields that can be read by the host VMM (TDH.SYS.RD) and MigTD (TDG.SYS.RD).

On export, the TDX module can work with MIG_VERSION in the range specified by the following metadata fields:

20 MAX_EXPORT_VERSION: Maximum value of migration version supported for export

MIN_EXPORT_VERSION: Minimum value of migration version supported for export

For example, a module may be updated to support version X for new memory migration formats for VM Nesting. But it may be written to export using version X-1 if a non-VM-Nesting TD is exported to an older module.

On import, the TDX module understands MIG_VERSION in the range specified by the following:

25 MAX_IMPORT_VERSION: Maximum value of migration version supported for import

MIN_IMPORT_VERSION: Minimum value of migration version supported for import

For example, if the module supports version X that was created to support new memory migration formats for VM Nesting, it can still understand version X-1 that doesn't use the new formats.

The migration protocol version to be used for a migration session is set the MigTD before the session starts. MigTDs at the source and destination platform enumerate export and import versions support, respectively. They decide on the version that is compatible between the platforms, to be used for the migration session. TD-scope metadata field MIG_VERSION is writable by the MigTD using TDH.SERVTD.WR. At the start of the migration session, the TDX module copies MIG_VERSION to an internal WORKING_MIG_VERSION that is used throughout the session.

6.5. Updated: Migration Session Control Functions Summary

This section provides an overview of the export session control functions. A detailed description is provided in [TDX Module ABI Spec].

10

5

Name	Description	Preconditions
TDH.EXPORT.STATE.IMMUTABLE	Start an export session. This function exports the TD's immutable state – that functionality is discussed in Ch. 7.	 TD is runnable. A new migration key has been set. Enough migration streams contexts have been created using TDH.MIG.STREAM.CREATE.
TDH.EXPORT.PAUSE	Pauses the source TD and starts the TDX-enforced blackout period. This operation is local to the source platform and is not communicated to the destination platform.	 TD immutable state has been exported. All TD VCPUs have stopped executing and no other TD-specific SEAMCALL is running.
TDH.EXPORT.TRACK	Starts a new migration epoch and generate an epoch token. If so requested, starts the out-of-order phase and generates a start token.	 The export session is in progress, but its out-of-order phase has not begun yet. See the discussion of export completeness checks in Ch. 8.
TDH.EXPORT.ABORT	Aborts the export session.	 An export session is in progress. If invoked during the out-of-order phase, the abort token received from the destination platform must be correct.
TDH.IMPORT.STATE.IMMUTABLE	Start an export session. This function exports the TD's immutable state – that functionality is discussed in Ch. 7.	 The TDCS been allocation but not initialized. A new migration key has been set. Enough migration streams contexts have been created using TDH.MIG.STREAM.CREATE.
TDH.IMPORT.TRACK	Starts a new migration epoch based on an imported epoch token. If the token is a start token, starts the out- of-order phase.	 The import session is in progress, but its out-of-order phase has not begun yet. The epoch token migration bundle is imported successfully. If a start token is received, all mutable state must have been imported.
TDH.IMPORT.COMMIT	Enable the TD to run.	The import session out-of-order phase is in progress.

Section 2: TD Migration Architecture Specification

Name	Description	Preconditions		
TDH.IMPORT.END	Ends the import session.	 The import session out-of-order phase is in progress. 		
TDH.IMPORT.ABORT	Aborts the import session (if not already aborted) and generates an abort token.	 The import session is in progress. The TD is not allowed to run yet (OP_STATE is not LIVE_IMPORT). 		

7. TD Non-Memory State Migration

This chapter discusses all non-memory state migration, immutable and mutable.

TD-scope non-memory state resides in control structures TDR and TDCS. TD VCPU state resides (while the VCPU is not running) in TDVPS, which includes the TD VMCS. This chapter discusses how non-memory state is migrated.

5 7.1. TD Non-Memory State Migration Operation

7.1.1. Non-Memory State Migration Data

Non-memory state migration data is used for migrating immutable state, at the beginning of the migration process, by TDH.EXPORT.STATE.IMMUTABLE and TDH.IMPORT.STATE.IMMUTABLE, and for migrating mutable state, at the end of the migration process, by TDH.EXPORT.STATE.TD, TDH.IMPORT.STATE.TD, TDH.EXPORT.STATE.VP and TDH.IMPORT.STATE.VP.

10

The non-memory state is migrated in a way that abstracts the actual TD control structure format, allowing that format to remain implementation-dependent and vary between the source and destination platforms.

7.1.2. Non-Memory State MBMD

The MBMD for each non-memory state migration bundle contains the following type-specific fields:

- Metadata type: Immutable TD-scope metadata or mutable L1 VCPU-scope metadata
 - VM index and VCPU index (if applicable).

Details of the non-memory state MBMD are defined in the [TDX Module ABI Spec].

7.1.3. Immutable vs. Mutable TD State

In the scope of TD migration, **immutable** state is defined as any TD state that may not change after TD build, i.e., after TD measurement has been finalized (by TDH.MR.FINALIZE).

Migrated immutable state includes the following:

- Platform-scope immutable state required so that the TDX module on the destination platform can verify compatibility. Namely, it includes the **source TDX module's version information**.
- TD-scope immutable state of the source TD
- 25 **Immutable TD state export and import** functions (TDH.EXPORT.STATE.IMMUTABLE, TDH.IMPORT.STATE.IMMUTABLE) start the migration session. Migration session control is discussed in Ch. 6.

Migrated mutable state includes the following:

- TD-scope mutable state
- VCPU-scope mutable state
- 30 **Mutable TD state export** is done after the TD has been paused (by TDH.EXPORT.PAUSE); it helps ensure that the state will not change anymore until TD export completes. TDH.EXPORT.STATE.TD exports TD-scope mutable state, followed by multiple, per-VCPU TDH.EXPORT.STATE.VP calls which export VCPUs mutable state.

Mutable TD state import must begin with TD-scope state import (by TDH.IMPORT.STATE.TD), followed by multiple, per-VCPU TDH.IMPORT.STATE.VP calls which import VCPUs state.

7.2. *New: Expected Configuration by the Host VMM*

The host VMM is expected to configure a migratable TD in a way that will be compatible with the set of possible destination platforms and their Intel TDX module configurations. For example:

- The configured TD ATTRIBUTES and XFAM bits should be supported by all destination platforms and their Intel TDX module configurations.
- The configured virtual CPUID values should be supported by all destination platforms and their Intel TDX module configurations.

7.3. Non-Memory State Migration Functions Summary

This section provides a short summary of the non-memory state migration interface functions. A detailed specification is provided in [TDX Module ABI Spec].

Name	Description	Preconditions	
TDH.EXPORT.STATE.IMMUTABLE	Export the TD's immutable state as a multi-page migration bundle. This function starts the export session; that functionality is described in Ch. 6.	 TD is runnable. A new migration key has been set. 	
TDH.EXPORT.STATE.TD	Export the TD-scope mutable state as a multi-page migration bundle.	• The export session is in the in-orde phase and the TD has been paused	
TDH.EXPORT.STATE.VP	Export the VCPU-scope mutable state as a multi-page migration bundle.	• The export session is in the in-order phase and the TD has been paused.	
TDH.IMPORT.STATE.IMMUTABLE	Import the TD's immutable state as a multi-page migration bundle. This function starts the import session; that functionality is described in Ch. 6.	 TD has not been initialized. A new migration key has been set. 	
TDH.IMPORT.STATE.TD	Import the TD's mutable state as a multi-page migration bundle.	• TD immutable state has been imported.	
TDH.IMPORT.STATE.VP	Imports a VCPU mutable state as a multi-page migration bundle.	 TDVPS pages have been allocated by the host VMM, but the VCPU has not been initialized. TD-scope state has been imported. 	

Table 7.4.			Minuchieus	Interferer	From estimate
	Non-Wemory	/ State	iviigration	Interface	Functions

8. Updated: TD Private Memory Migration

This chapter described how Intel TDX Module manages TD private memory and guest-physical (GPA) address translation metadata migration.

8.1. Updated: Overview

8.1.1. Updated: In-Order and Out-of-Order Migration 5

TD private memory migration can happen in the in-order migration phase and out-of-order migration phase.

During the in-order phase, the host VMM may implement live migration pre-copy, by exporting memory content (using TDH.EXPORT.MEM etc.) while the TD is running (TDCS.OP_STATE is LIVE_EXPORT). This is not enforced by TDX; the host VMM may implement cold migration by avoiding memory export until the TD is paused.

- During the out-of-order phase, the host VMM may implement **post-copy** by allowing the TD to run on the destination 10 platform (using TDH.IMPORT.COMMIT). This is not enforced by TDX; the host VMM can first complete all memory migration before allowing the TD to run, yet benefit from the simpler and potentially higher performance operation supported during the out-of-order phase.
- Note: Even if the TDX module supports post-copy with non-blocking export, post-copy is only allowed if TDX Connect is not enabled for the TD. 15

8.1.2. New: Write-Blocking Export vs. Non-Blocking Live Export

During live migration, the TDX module tracks exported page modification to enforce consistent memory image migration. The TDX module supports two modes of memory live export:

- Write Blocking: Memory pages are blocked for writing before being exported. A guest TD attempt to modify exported memory results in an EPT violation TD exit. The host VMM is expected to unblock the page and later block it again and re-export it. TDX module support of write-blocking based export is enumerated by TDX_FEATURES0.TD_MIGRATION (bit 0), readable by TDH.SYS.RD*.
- **Non-Blocking:** Memory pages are not blocked for writing. If a memory page is modified by the guest TD, DMA or by the TDX module, its SEPT entry's Dirty bit is set. The host VMM is expected to call a TDX module interface function that scans the TD's GPA space for dirty pages, and re-export those pages. TDX module support of non-blocking export is enumerated by TDX_FEATURES0.NON_BLOCKING_EXPORT (bit 41), readable by TDH.SYS.RD*.

The export mode is a TDX module configuration parameter, selected by the host VMM via a setting as part of the TDX module initialization sequence.

30

20

25

- By default, write blocking export is used.
 - Non-blocking export can be configured during first-time module initialization (by TDH.SYS.CONFIG). •
 - On TD-preserving TDX module update, the export mode may be updated from the default write blocking mode to non-blocking mode, if no export session has been used before.

8.2. New: Conventions: SEPT Entry State Diagrams Color Coding

35 This chapter contains multiple SEPT entry state diagrams. For ..., which use the following color-coding conventions.

\subset	
\bigcap	

GPA is guest accessible

GPA is blocked, TLB entries may still be cached

GPA is not guest accessible, GPA may not be reused, no physical page is associated



GPA is guest accessible for read and execute only



GPA is not guest accessible, GPA may not be reused, physical page still associated

GPA is not guest accessible, no physical page is associated

Figure 8.1: Secure EPT Entry State Diagrams Color Conventions

8.3. Updated: GPA Lists and Private Memory Migration Bundles

8.3.1. Overview



Figure 8.2: Private Memory Migration

- 5 Unlike the generic migration bundle structure described in 5.1, private memory migration bundle is composed of multiple MAC-protected components:
 - MAC-protected MBMD
 - For each 4KB page:
 - Encrypted and MAC-protected 4KB migration buffer
 - MAC-protected page GPA and additional metadata
 - Optionally, for partitioned TDs, encrypted and MAC-protected page attributes

This structure allows the export and import functions to process the MBMD and each page and its metadata separately, avoiding the need to perform SEPT walks twice and to hold intermediate SEPT entry states. The separate parts of the migration bundle are cryptographically bound together as follows:

- A per-stream monotonically incrementing IV_COUNTER and the migration steam index are used for calculating the AES-GCM IV value, as described in 5.3.
 - This is first done for the migration bundle's MBMD MAC.
 - For each page, the IV_COUNTER is incremented by 1 and a new IV value is calculated and used for the page metadata MAC.
- The MBMD specifies the number of pages migrated by the migration bundle. This helps check that the whole migration bundle is imported on the destination platform.

8.3.2. Updated: GPA List

As shown in the example in the diagram above, a GPA list is used as part of the private memory migration bundle. It is also used as an input and output of multiple memory migration interface functions: TDH.EXPORT.BLOCKW, TDH.EXPORT.MEM, TDH.EXPORT.RESTORE, TDH.IMPORT.MEM and TDH.MEM.SCAN.

25

A GPA list contains up to 512 entries, each containing the following information:

Table 8.1: Updated: GPA List Entry Abstract Definition

Field	Usage for Page Migration	Encrypted?	Details
Page Size	4КВ	No	
GPA	GPA bits 51:12	No	
State	MAPPED or PENDING	No	
Operation	NOP, MIGRATE, REMIGRATE or CANCEL	No	
Attributes	Page attributes for each L2 VM	Yes	Optional, provided in a separate list (see below)
MAC	Integrity protection for the above fields and for the migrated page content	N/A	Provided in a separate list

A single GPA list entry, a separate page MAC list entry and an optional separate page attributes list entry compose the page metadata.

The GPA list is MAC-protected but is not encrypted. This allows the host VMM on the destination to parse the GPA list in order to prepare for calling TDH.IMPORT.MEM; e.g., build the Secure EPT to map the imported pages.

A detailed definition of the GPA list is provided in the [TDX Module ABI Spec].

8.3.3. Page Attributes List (Required for Partitioned TDs)

10 If the migrated TD is partitioned, the GPA list entry is extended with L2 page attributes. These are provided in a separate page attributes list, containing the same number of entries as the GPA list. Each page attributes list entry contains the migratable page attributes for each L2 VM: R, W, Xs, Xu, SSS, VGP, PWA and SVE. The page attributes list is encrypted and MAC-protected.

A detailed definition of the page attributes list is provided in the [TDX Module ABI Spec].

15 8.3.4. Private Memory Migration Buffer

The migration buffer holds the encrypted migrated page content; thus, it is included only if the page metadata indicates a MIGRATE or a REMIGRATE request, and the page is MAPPED. The migration buffer is allocated by the host VMM and resides in shared memory. The migration buffer is encrypted and MAC-protected.

On first time import of a page (MIGRATE request), the host VMM can select in-place import: the migration buffer becomes the TD private page which holds the imported and decrypted content.

8.4. Updated: Write-Blocking Based Memory Export

This section describes private memory export using the write-blocking method.

8.4.1. Edited: Host VMM Perspective

This section describes write-blocking based export from the host VMM perspective. This is a simplified view; a more detailed view is discussed later.

Section 2: TD Migration Architecture Specification

8.4.1.1. Typical Export Session



Figure 8.3: Typical Write-Blocking Based Export Session

Typical expected usage divides the export session to **export rounds** (or **passes**). An export round may have the following steps:

- 1. If the TD has not been paused by TDH.EXPORT.PAUSE, ensure TLB shootdown:
 - 1.1. Invoke TDH.EXPORT.BLOCKW with a list of pages to be exported.
 - 1.2. Invoke TDH.MEM.TRACK.
 - 1.3. Issue IPIs to ensure TD re-entry on all VCPUs and TLB invalidation.
- 2. Start a new migration epoch by invoking TDH.EXPORT.TRACK.
- 3. Invoke TDH.EXPORT.MEM with a list of pages.
 - 3.1. If a page is being exported, mark its list entry as MIGRATE.
 - 3.2. If a page has been exported before but needs to be removed, promoted or demoted, cancel its migration by marking its list entry as CANCEL.
- 15 Note: In the example above, steps 1 and 2 need to be performed before step 3, but there is no strict requirement for the order of step 2 vs. step 1.

8.4.1.2. Live Export: Blocking for Writing, TLB Tracking and Exporting a Page

During the live export phase (when TDCS.OP_STATE is LIVE_EXPORT), exporting a private memory page requires that page modification must be tracked by the TDX module. This includes:

20 • Page content

5

10

25

30

- Page attributes
- For partitioned TDs, L2 page attributes

To achieve this, the page's L1 SEPT entry and any L2 SEPT entries must be blocked for writing by TDH.EXPORT.BLOCKW.

If the TD has not been paused, the host VMM must execute the TLB tracking sequence below, which together with the checks done by TDH.EXPORT.MEM helps ensure that no cached TLB entries that have been created before blocking for writing are left.

- 1. Execute TDH.EXPORT.BLOCKW on each page to be exported, blocking subsequent creation of writable TLB translations to that page. Note that cached translations may still exits at this stage.
- 2. Execute TDH.MEM.TRACK, advancing the TD's epoch counter.
- 3. Send an IPI (Inter-Processor Interrupt) to each RLP (Remote Logical Processor) on which any of the TD's VCPUs is currently scheduled.

4. Upon receiving the IPI, each RLP will TD exit to the host VMM.

At this point the blocked pages are considered tracked for export. Even though some LPs may still hold writable TLB entries to the target GPA ranges, those are designed to be flushed on the next TD entry. Normally, the host VMM on each RLP will treat the TD exit as spurious and will immediately re-enter the TD.

5 5. Export each page using TDH.EXPORT.MEM.

8.4.1.3. Exporting a Page after the Source TD is Paused

After the source TD is paused, no blocking is required since the TD is not running. This reduces the amount of work that needs to be done by the host VMM during the TD's blackout period. This is shown in the dashed transitions in Figure 8.8 below.

10 8.4.1.4. Unblocking for Write, Tracking Dirty Pages and Re-Exporting

8.4.1.4.1. Overview

20

During the live export phase (when TDCS.OP_STATE is LIVE_EXPORT), the source TD may attempt to write a page that has been blocked for writing, or to modify the page attributes (for partitioned TDs, this includes L2 attributes). The TDX migration architecture allows the host VMM to unblock the page. The Intel TDX module tracks such pages as "dirty". All

dirty pages must be re-exported by the host VMM for the in-order migration phase to be completed. This assures that either the latest version of a page has been exported by the time the source TD is paused, or that page has not been exported at all.

8.4.1.4.2. Unblocking for Write and Re-Exporting a Page

If the source TD attempts to write to a page that has been blocked for writing, a TD exit will occur, indicating an EPT violation due to a write attempt to a non-writable page.

Note: No indication is directly provided to the host VMM whether this page is blocked for writing by TDH.EXPORT.BLOCKW or whether writing is disabled due to some other reason (e.g., the page is BLOCKED).



Figure 8.4: Typical Sequence for Unblocking a Page on Guest TD Write Attempt

- To enable access to the page, the host VMM is expected to execute TDH.EXPORT.UNBLOCKW and then resume the TD VCPU by TDH.VP.ENTER.
 - If the page has not yet been exported, TDH.EXPORT.UNBLOCKW restores its SEPT entry's original MAPPED state.
 - If the page has been exported, TDH.EXPORT.UNBLOCKW updates its SEPT state to EXPORTED_DIRTY. This state is similar to MAPPED from the guest TD's memory access perspective, but it indicates that the page is dirty and needs to be re-exported.
 - For partitioned TDs, if the page has any L2 mappings, TDH.EXPORT.UNBLOCKW unblocks their L2 SEPT entries by restoring their W bit value.

The host VMM re-exports the page by TDH.EXPORT.BLOCKW, TLB tracking and TDH.EXPORT.MEM as described below.

8.4.1.5. Using the same GPA List for TDH.EXPORT.BLOCKW and TDH.EXPORT.MEM

TDH.EXPORT.BLOCKW and TDH.EXPORT.MEM use GPA lists with compatible formats. This allows the same list to be used for blocking and exporting memory, as follows:

- 1. The host VMM on the source platform may prepare a GPA list with MIGRATE and CANCEL commands (see later), and
- 5
- provide it as input to TDH.EXPORT.BLOCKW.
 2. TDH.EXPORT.BLOCKW will attempt to block pages whose command is MIGRATE, and update the GPA list depending on the success of the operation.
- 3. The same GPA list can be provided as input to TDH.EXPORT.MEM, which then updates it depending on the success of the operation and whether this is a first-time export (MIGRATE) or re-export (REMIGRATE), and adds page attributes information.

10



Figure 8.5: Typical Write-Blocking Based Memory Export Round and the GPA List

A detailed definition of the GPA list is provided in the [TDX Module ABI Spec].

8.4.1.6. Prohibited Operations on Exported Pages and Export Cancellation

Once a page has been exported during the current export session, it can't be blocked, removed, promoted, demoted or relocated. This prevents the destination platform from using a stale copy of that page.

In order to perform such memory management operations on an exported page, the host VMM must first execute TDH.EXPORT.MEM indicating a CANCEL operation for the page. No migration buffer is required for this GPA list entry. When the GPA list is processed on the destination platform by TDH.IMPORT.MEM, the previously migrated page is

removed from the destination TD. TDH.EXPORT.MEM restores the page SEPT entry to its pre-export MAPPED or PENDING state.



Figure 8.6: Typical Sequence for Cancelling a Page Export

8.4.1.7. **Updated:** Exporting Pending Pages 5

The host VMM is not directly aware if a page is in a PENDING state or not; the guest TD may accept a PENDING page by TDG.MEM.PAGE.ACCEPT at any time. If supported, the guest may release a MAPPED page by TDG.MEM.PAGE.RELEASE, converting it to a PENDING page. The page content of a PENDING page is not exported, and no migration buffer is used. The page attributes (including the optional page attributes list entry) are exported. On the destination platform, TDH.IMPORT.MEM creates the page in a PENDING state.

10

If the guest TD accepts a pending page that has been exported, TDG.MEM.PAGE.ACCEPT results in an EPT violation. The host VMM is expected to call TDH.EXPORT.UNBLOCKW, which marks the page as PENDING_EXPORT_DIRTY, and resumes the guest TD. TDH.MEM.PAGE.ACCEPT then re-executes; it initialized the page and updates the SEPT state to mark the page as EXPORTED_DIRTY (where the page is mapped and accessible to the guest TD). The host VMM can then re-export

the page, as described in 8.4.1.4 above. 15





8.4.1.8. Re-Exporting a Non-Dirty Page

In the out-of-order phase, where strict migration order is not enforced, the host VMM may re-export a previously exported page even if it has not been unblocked for writing and its contents have not been modified.

This allows a page to be re-exported and transferred to the destination platform over a high-priority stream. This helps reduce destination TD latency while waiting for a page to be imported.

Such an operation is tagged MIGRATE, not REMIGRATE, in the exported GPA list. This is because the exact same version of the page is being exported.

8.4.1.9. SEPT Cleanup after Export Abort

10

5

After an export session is aborted (by TDH.EXPORT.ABORT), the source TD is allowed to run. However, SEPT entries and, for partitioned TDs, L2 SEPT entries, that have been modified during the aborted export session keep their state. Such SEPT entries must be cleaned up by the host VMM before memory management operations are allowed on them, and/or before a new export session is attempted, as follows:

- Cleanup of SEPT entries that have been blocked for writing is done by TDH.EXPORT.UNBLOCKW (if the page is to be written) or TDH.RANGE.BLOCK (is the page is to be blocked for some memory management operation).
- Cleanup of SEPT entries that have been exported is done by TDH.EXPORT.RESTORE.

The TDX module holds a TDCS field called MIG_COUNT, which counts the number of exported pages that require cleanup. This field is readable by the host VMM, using TDH.MNG.RD. The counter is initialized to 0. To start a new migration session, its value must be 0.

8.4.2. Edited: Details of Write-Blocking Based Export

20 This section provides a detailed view of write-blocking based export. Details may be of interest to host VMM programmers who require deeper understanding of TD Migration.

8.4.2.1. Details: L1 SEPT Leaf Entry Partial State Diagram for Mapped Page Export

Figure 8.8 below shows a partial SEPT entry state diagram for exporting mapped pages. The following sections describe the details. The color-coding convention used in this diagram is described in 8.2.



Figure 8.8: Partial L1 SEPT Leaf Entry State Diagram for Mapped Page Write-Blocking Based Export

8.4.2.2. Details: L1 SEPT Leaf Entry Partial State Diagram for Pending Page Export

The figure below shown the partial state diagram for PENDING page export.



Figure 8.9: Partial L1 SEPT Leaf Entry State Diagram for Pending Page Write-Blocking Based Export

8.4.2.3. Details: TDCS.DIRTY COUNT: TD-Scope Dirty Page Counter

TDCS.DIRTY COUNT is a TD-scope dirty page counter.

- DIRTY COUNT is cleared when a new migration session begins (by TDH.EXPORT.STATE.IMMUTABLE). •
 - DIRTY COUNT is incremented when a page that has previously been exported in the current session is unblocked for writing by TDH.EXPORT.UNBLOCKW.
 - DIRTY_COUNT is decremented when a newer version of a page, which has previously been exported in the current session, is exported by TDH.EXPORT.MEM.
- For successful start token generation by TDH.EXPORT.TRACK, the value of the DIRTY COUNT must be 0, indicating that 10 all pages exported so far have their newest pages exported. At this point, since the source TD is paused, no newer versions of any page can be created, and the destination TD can start execution. Private pages which have not been exported yet in the current session may still be remaining for post copy export. Note that exported pages may not have been transported yet. The start token MBMD's TOTAL MB field verification enforces that all exported state has been imported (in-order) on the destination - see the [TDX Module ABI Spec] for details.
- 15

5

New: Non-Blocking Memory Export **8.5**.

8.5.1. **Host VMM Perspective**

This section describes non-blocking export from the host VMM perspective. This is a simplified view; a more detailed view is discussed later.

8.5.1.1. EPT Access and Dirty Bits Background

Intel SDM, Vol. 3, 30.3.5 Accessed and Dirty Flags for EPT

The Dirty Bit

Non-blocking live export relies on detecting memory changes using the Secure EPT entry's Dirty bit (9). This bit is set by
the h/w (CPU or IOMMU) in the leaf EPT entry for a certain GPA when it performs an EPT walk for translating that GPA for a write operation. The Dirty bit is sticky – the h/w may only set it to 1 but never clear it. The Dirty bit is only cleared by s/w – in case of Secure EPT, by the TDX module.

Implication of Address Translation Caching

10

The CPU and IOMMU cache page address translations and paging structures. An EPT entry's Dirty bit is only set during an EPT walk, when there is a need to do address translation, and the applicable entries are not cached. This means that even if the TDX module clears the Dirty bit of an SEPT entry, the h/w may not be aware of this since it holds a cached entry; it will not set the Dirty bit if there's a new write access. Since page modification detection relies on the correct value of the Dirty bit, there is a need to do a TLB shootdown, as described in the following sections. TLB shootdown is tracked by the TDX module as a prerequisite to page export, to help ensure secure operation.

15 8.5.1.2. Memory Export Concept: Scan and Export

Non-blocking export is based on the concept of scan and export. Export typically consists of multiple rounds, with the following steps per round:

- 1. Scan the TD's GPA space for memory export candidates.
- 2. Do a TLB shootdown.
- 20 3. Based on the scan results, prepare a list of pages and export them.



Figure 8.10: Memory Export Round Concept: Scan and Export

The host VMM is **not required** to track the state of the TD private pages. It can rely on the provided scan functions to provide the required information.

25 **8.5.1.3.** *Conceptual, Simplified Page State Diagram*

Although the host VMM is not required to track the TD private page state, it is worthwhile to understand how the TDX module tracks page state using the Secure EPT state. This section provides a conceptual view of the page state machine, which serves as a simplified overview of the export operation from the host VMM's perspective.

Note:The actual SEPT entry state machine is more complex. It is described in the following sections and in the [Base30Spec].

Draft

The following table lists the conceptual state machine's states and their handling by the memory scan functions (which are described later):

State	Description	TDH.MEM.SCAN.*
NOT_EXPORTED	Pages that either have not been exported during the current migration session, or that were exported but their export was later cancelled.	Page GPA is reported by TDH.MEM.SCAN.*, to let the host VMM know it needs to be exported.
EXPORTED	Pages that were exported. Pages may have their Dirty bit set due to content modifications by the CPU, DMA or the TDX module, or due to attributes modifications by the guest TD (e.g., TDG.MEM.PAGE.ACCEPT or TDG.MEM.PAGE.ATTR.WR).	If the page's SEPT entry's Dirty bit is set, TDH.MEM.SCAN.* clears the Dirty bit, sets the page state to EXPORTED_MODIFIED and reports the page GPA to let the host VMM know it needs to be re-exported.
EXPORTED_MODIFIED	Pages that were exported, and later either scanned and found to require re-export since their Dirty bit was set, or their attributes were changed by some TDX module memory management function.	TDH.MEM.SCAN.* reports the page GPA to let the host VMM know it needs to be re-exported.
EXPORT_CANCEL_ REQUIRED	Pages that were exported but their export must be cancelled, due to memory management operations (such as TDH.MEM.RANGE.BLOCK) by the host VMM. Note: This conceptual state represents multiple states, shown in the detailed diagram later; those are required since different memory management operations are handled differently.	TDH.MEM.SCAN.* reports the page GPA to let the host VMM know it needs to be re-exported.
BLOCKED	 When a page that was in an EXPORT_CANCEL_REQUIRED state due to blocking is exported (as a CANCEL operation), the page state becomes the normal BLOCKED state. Note: This conceptual state represents multiple states, shown in the detailed diagram later; those are required since different memory management operations are handled differently. 	Page GPA is not reported by TDH.MEM.SCAN.RANGE(DSCAN) since it is not exported.

Table 0.2.	Concentual	Circulified I	Jana Chatan	for Dist.	Dit Deced Evenent
I dule 0.2.	conceptual,	Simplified	age states	IOI DIILY	DIL Daseu Export



Figure 8.11: Conceptual Page State Diagram for Dirty Bit Based Export

8.5.1.4. Scanning for Candidate Pages to Export or Re-export

8.5.1.4.1. Overview

10

15

25

- 5 During the live export phase, memory pages may be modified asynchronously to the host VMM's operations, by the running TD or its bound TDIs. The host VMM does not need to know the SEPT entry states. It uses the export candidates list provided by following interface functions:
- **TDH.MEM.SCAN.RANGE(DSCAN):** This function is used iteratively by the host VMM during the LIVE_EXPORT phase of TD migration to identify pages in the specified GPA range needing export. Identification is based on the SEPT entry's page state and the Dirty bit. The function returns a list of pages that the host VMM can use to prepare export requests (as an input to TDH.EXPORT.MEM).
- **TDH.MEM.SCAN.COMP(DCHECK):** This function is used by the host VMM during the export blackout period, after the TD is paused and all TDIs are unbound. The function performs a comprehensive scan of the TD's GPA address space and identifies the remaining pages needing export and returns a page list. It ensures migration consistency by requiring that the whole GPA range will be scanned and maintaining a counter of pages that need to be exported. TDH.EXPORT.TRACK(DONE) then checks that those pages have indeed been exported, as a precondition for generating a start token.
- 20 Although memory export only supports 4KB page mapping, the scanning functions return information for 4KB, 2MB and 1GB page mapping sizes. The host VMM is responsible for demoting 2MB and 1GB pages before calling TDH.EXPORT.MEM.

The host VMM is expected to maintain an export GPA list:

- A GPA should be added to the list when it is reported by TDH.MEM.SCAN.*.
- A GPA should be removed from the list when it is successfully exported by TDH.EXPORT.MEM.

TDH.MEM.SCAN.* returns a list of export candidates regardless of whether pages in that list have been returned by a previous call to TDH.MEM.SCAN.*. I.e., if page X has been reported as an export candidate, and had not yet been exported, it may be reported again.

8.5.1.4.2. Using the Memory Scanning Functions

To help reduce the scan time, especially during the export blackout period, TDH.MEM.SCAN.* can be called concurrently on multiple LPs.

A comprehensive scan of the whole GPA range is done as follows:

- The host VMM needs to configure the comprehensive scan using TDH.MEM.SCAN.CONFIG. The host VMM can divide the TD's GPA address space into ranges, typically to optimize the scan for NUMA configurations. The configuration is not migrated with the TD; it needs to be redone if the TD is to be migrated again.
- The host VMM can call TDH.MEM.SCAN.COMP on multiple threads, specifying the pre-configured GPA range to scan by each instance (multiple instances can concurrently scan the same range).
 - If the scan fails (e.g., there are blocked pages), the host VMM can reset the comprehensive scan using TDH.MEM.SCAN.RESET and then retry the scan.

For further details, see the [TDX Module ABI Spec].

10 8.5.1.5. Typical Non-Blocking Export Session

Typical expected usage divides the export session to **export rounds** (or **passes**). The diagram below shows a typical export session, which is composed of several **live migration export rounds** and a **final export round** during the blackout period, when the TD is not running. It also shows post-copy (which is not applicable if the TD is configured for TDX Connect).

The diagram below shows 4 migration streams and the usage of migration epochs to synchronize them. Migration epochs are defined in 6.1.4.





8.5.1.5.1. Typical Live Export Round

Note: Do not confuse TDH.MEM.TRACK (which is used for TLB tracking) with TDH.EXPORT.TRACK (which rendezvous all migration streams).

An export round during the live migration phase of the export session typically has the following steps:

- Prepare a list of pages to be exported by calling TDH.MEM.SCAN.RANGE(DSCAN) one or more times.
 1.1. TDH.MEM.SCAN.RANGE(DSCAN) can be called concurrently on multiple LPs.
- 2. Do the TDX TLB (and IOTLB, if the TD is enabled for TDX Connect) shootdown sequence. This is required to ensure that, after the following export, the CPU or IOMMU will indeed set the Dirty bit if a page is written.
 - 2.1. Call TDH.MEM.TRACK.
 - 2.2. Issue IPIs (Inter-Processor Interrupts) to all RLP (Remote Logical Processors) on which any of the TD's VCPUs is currently scheduled. This causes TD re-entry and TLB (and IOTLB, if the TD has bound TDI) invalidation on all VCPUs.
- 30 3. Start a new **migration epoch** by invoking TDH.EXPORT.TRACK. This is required for the TDX module on the destination to ensure proper ordering on import.

20

25

5

- 4. Call TDH.EXPORT.MEM with a list of pages based of the export candidates list reported by TDH.MEM.SCAN.RANGE(DSCAN) above.
 - 4.1. The host VMM may choose not to include pages reported as blocked, since this is an interim state as a preparation to some memory management operation.
 - 4.2. TDH.EXPORT.MEM determines the required export operation (MIGRATE, REMIGRATE or CANCEL) based on each page state.
 - 4.3. TDH.EXPORT.MEM checks that the TLB shootdown above has indeed been done, using the TDX TLB tracking mechanism (see the [TDX Module Base Spec] for details).
 - 4.4. TDH.EXPORT.MEM can be called concurrently on multiple LPs, each exporting to a different migration stream.

10 Concurrency and Order during Live Export

- There is no strict ordering requirement between steps 2 and 3, as long as they precede step 4.
- TDH.MEM.SCAN.RANGE(DSCAN) can run concurrently with TDH.EXPORT.MEM, as long as different pages are processed. E.g., it is possible to scan a certain GPA range while exporting memory of a different GPA range.
- If TDH.MEM.SCAN.RANGE(DSCAN) encounters a busy SEPT entry due to a conflict with some other function, it skips
- that entry. The page will be processed by the next TDH.MEM.SCAN.RANGE(DSCAN) or TDH.MEM.SCAN.COMP(DCHECK). For details, see the [ABI Spec].
 - If TDH.EXPORT.MEM encounters a busy SEPT entry due to a conflict with some other function, it skips that entry. It writes the status into the GPA list; the host VMM may read that list and call TDH.EXPORT.MEM to export pages that have not been exported. For details, see the [ABI Spec].

20 8.5.1.5.2. Typical Final Export Round

The final export round ensures that the full up-to-date memory image (and non-memory state) is exported. It typically consists of the following steps:

- 1. Pause the TD using TDH.EXPORT.PAUSE. This requires no pages to be blocked and, if the TD is enabled for TDX Connects, no TDIs to be attached.
- 2. Export the TD's non-memory immutable state using TDH.EXPORT.STATE.TD and TDH.EXPORT.STATE.VP.
 - 3. Prepare a list remaining of pages to be exported in the final round by calling TDH.MEM.SCAN.COMP(DCHECK) one or more times until the whole GPA range has been scanned.
 - 3.1. TDH.MEM.SCAN.CONFIG should be called prior to TDH.MEM.SCAN.COMP to configure the GPA ranges for scanning.
 - 3.2. TDH.MEM.SCAN.COMP(DCHECK) can be called concurrently on multiple LPs.
 - 4. There is no need to do a TLB shootdown since the TD is paused and no TDI is bound to the TD.
 - 5. Start a new **migration epoch** by invoking TDH.EXPORT.TRACK. This is required for the TDX module on the destination to ensure proper ordering on import.
- 6. Call TDH.EXPORT.MEM with a list of pages. The TDX module determines the required export operation (MIGRATE,
- REMIGRATE or CANCEL) based on each page state.
 - 6.1. All the pages detected by TDH.MEM.SCAN.COMP(DCHECK) as requiring re-export must be exported.
 - 6.2. If post-copy is supported (for a TD where TDX Connect is not enabled), pages identified by TDH.MEM.SCAN.COMP(DCHECK) as never having been exported may be exported later during the out-of-order phase.
- 6.3. TDH.EXPORT.MEM may be called concurrently on multiple LPs, each exporting to a different migration stream.
- 7. End the in-order export phase by calling TDH.EXPORT.TRACK(DONE).

Concurrency and Order during Final Export

- TDH.MEM.SCAN.COMP(DCHECK) can run concurrently with TDH.EXPORT.MEM, as long as different pages are processed. E.g., it is possible to scan a certain GPA range while exporting memory of a different GPA range.
- Unlike TDH.MEM.SCAN.RANGE(DSCAN), if TDH.MEM.SCAN.COMP(DCHECK) encounters a busy SEPT entry due to a conflict with some other function, it returns to the host VMM indicating an interrupted operation. The host VMM is expected to resume TDH.MEM.SCAN.COMP(DCHECK). For details, see the [ABI Spec].

5

15

25

30

35

8.5.1.6. Interaction with Memory Management Operations

8.5.1.6.1. Memory Management Restrictions on Pages that Have Been Exported

Exported Page Blocking and Removal

A page that has been exported can still be blocked (TDH.MEM.RANGE.BLOCK) and/or removed (TDH.MEM.PAGE.REMOVE) by the host VMM, subject to other restrictions (e.g., no attached TDIs for TDs configured for 5 TDX Connect).

Synchronization with the destination is done by exporting a CANCEL operation for this page, so that the destination can remove it. A blocked or removed page is put in one of several special states (together called EXPORT_CANCEL_REQUIRED in the conceptual state diagram above). Pages in those states are detected by TDH.MEM.SCAN.RANGE(DSCAN) and TDH.MEM.SCAN.COMP(DCHECK) so that the host VMM can include them in the GPA list for TDH.EXPORT.MEM. Once

10 TDH.EXPORT.MEM has exported a CANCEL operation for such page, the page state is set to one of the NON EXPORTED states.

Exported Page Promotion or Demotion

A page that has been exported can't be promoted (TDH.MEM.PAGE.PROMOTE). Since an exported page mapping size is 4KB, it can't be demoted (TDH.MEM.PAGE.DEMOTE) either.

Even without migration, the host VMM should always be prepared for a Promote operation to fail. E.g., L1 may Note: set the attributes of some of the small pages to be merged differently than the other small pages, preventing promotion.

8.5.1.6.2. **Memory Management Restrictions During the Export Blackout Period**

20 The following memory management restriction apply during the PAUSED EXPORT state, which begins when the TD is paused by TDH.EXPORT.PAUSE and ends when a start token is generated by TDH.EXPORT.TRACK(DONE). These restrictions exist to facilitate export completeness tracking (for details, see 8.5.2.7 below).

No Blocked Pages

15

25

30

45

No page may be blocked while in the PAUSED EXPORT state. This condition is checked by TDH.EXPORT.PAUSE and TDH.MEM.SCAN.COMP(DCHECK). TDH.MEM.RANGE.BLOCK and TDH.MEM.RANGE.UNBLOCK are not allowed.

TDH.MEM.PAGE.DEMOTE and TDH.MEM.PAGE.RELOCATE are allowed without blocking in the PAUSED_EXPORT Note: state.

Other Restrictions

- Page addition and removal by TDH.MEM.PAGE.AUG and TDH.MEM.PAGE.REMOVE is not allowed while in the PAUSED EXPORT state. Note that TDH.MEM.PAGE.ADD is also not allowed since the TD build has been, by definition, finalized.
 - SEPT tree changes by TDH.MEM.PAGE.PROMOTE, TDH.MEM.SEPT.ADD, and TDH.MEM.SEPT.REMOVE are not allowed while in the PAUSED_EXPORT state.

8.5.1.7. **Exporting Pending Pages**

- The host VMM is not directly aware if a page is in a PENDING state or not: 35
 - The guest TD may accept a PENDING page by TDG.MEM.PAGE.ACCEPT at any time, converting it to a MAPPED page.
 - If supported, the guest TD may release a MAPPED page by TDG.MEM.PAGE.RELEASE, converting it to a PENDING page.

Thus, TDH.EXPORT.MEM may export a pending page. This is indicated by the GPA list entry, and no migration buffer is used since the page content is not exported. The page attributes (including the optional page attributes list entry) are 40 exported. On the destination platform, TDH.IMPORT.MEM creates the page in a PENDING state.

8.5.1.8. SEPT Cleanup after Export Abort

After an export session is aborted (by TDH.EXPORT.ABORT) the source TD is allowed to run. However, SEPT entries and, for partitioned TDs, L2 SEPT entries, that have been modified during the aborted export session keep their state. Such SEPT entries must be cleaned up by the host VMM before memory management operations that are not allowed on

Draft

migrated pages (such as TDH.MEM.PAGE.PROMOTE) are attempted on them, and/or before a new export session is attempted.

Cleanup is done either by TDH.EXPORT.RESTORE, which receives a list of up to 512 pages, or (if supported) by TDH.MEM.SCAN.RANGE(EXPORT_RESTORE), which scans a GPA range.

5 The TDX module holds a TDCS field called MIG_COUNT, which counts the number of exported pages that require cleanup. This field is readable by the host VMM, using TDH.MNG.RD. The counter is initialized to 0. To start a new migration session, its value must be 0.

8.5.2. Details of Non-Blocking Export

This section provides a detailed view of non-blocking export. Details may be of interest to host VMM programmers who require deeper understanding of TD Migration.

8.5.2.1. Details: L1 SEPT Leaf Entry Partial State Diagram for Mapped Page Export

To support non-blocking export of Mapped pages, the following SEPT states are added. Their main goal is to help keep the destination TD memory image in sync with the source TD memory image. The next section describes the states added for Pending pages.

15

Table 8.3: L1 SEPT Entry States Related to Mapped	Page Export
---	-------------

State Name	Description
EXPORTED	The page has been exported. This state indicates that any change in the page content or attributes would require a re-export.
EXPORTED_MODIFIED	The page was exported and later either identified as dirty based on the Dirty bit (which was atomically cleared by the same operation that checked its value) or some memory management operation changed the page attributes or state. This state indicates that the page must be re- exported, as a REMIGRATE operation.
EXPORTED_BLOCKED	The page was exported and later blocked by TDH.MEM.RANGE.BLOCK. This state indicates that the page must be re-exported, as a CANCEL operation.
EXPORTED_REMOVED	The page was exported and later removed by TDH.MEM.PAGE.REMOVE. This state indicates that the page should be re-exported, as a CANCEL operation.
EXPORTED_REMOVE_IN_PROGRESS	The page was exported and later partially removed by TDH.MEM.PAGE.REMOVE. This state indicates that the page should be re-exported, as a CANCEL operation.
	This state is only applicable for client platforms, which used ACT-based memory protection.

A partial SEPT entry state diagram for exporting Mapped pages using the Dirty bit method is shown below. The following sections describe the details. The color-coding convention used in this diagram is described in 8.2.



Figure 8.13: Partial L1 SEPT Leaf Entry State Diagram for MAPPED Page Non-Blocking Export

8.5.2.2. Details: L1 SEPT Leaf Entry Partial State Diagram for Pending Page Export

The state diagram for PENDING page export is very similar to the one above. Note that PENDING pages are considered non-present from the h/w perspective, so the CPU and DMA do not normally set the Dirty bit. An exception to this is pages converted to PENDING by TDG.MEM.PAGE.RELEASE where the h/w may hold a TLB translation.

State Name	Description
PENDING_EXPORTED	The page has been exported. This state indicates that any change in the page attributes would require a re-export.
PENDING_EXPORTED_MODIFIED	The page was exported and later identified as dirty based on the Dirty bit (which was atomically cleared by the same operation that checked its value). This state indicates that the page must be re-exported, as a REMIGRATE operation.

 Table 8.4: L1 SEPT Entry State Related to Pending Page Export

Draft

State Name	Description
PENDING_EXPORTED_BLOCKED	The page was exported and later blocked by TDH.MEM.RANGE.BLOCK. This state indicates that the page must be re-exported, as a CANCEL operation.

A partial SEPT entry state diagram for exporting Pending pages using the Dirty bit method is shown below. The following sections describe the details. The color-coding convention used in this diagram is described in 8.2.



Figure 8.14: Partial L1 SEPT Leaf Entry State Diagram for PENDING Page Non-Blocking Export

8.5.2.3. Details: TD Partitioning Considerations for Dirty Bit Operations

With TD partitioning, in addition to the main SEPT tree used by L1, each L2 VM has its own L2 SEPT tree. TD private memory writes by an L2 VM, by devices associated with that L2 VM or by TDX module flows that run in the context of that L2 VM set the Dirty bits in L2 SEPT entries.

10 Thus, when the logical Dirty bit is discussed, the actual operation may involve the L1 SEPT entry's Dirty bit and/or one or more of the L2 SEPT entries' Dirty bits. The following table describes the various operations done on the Dirty bit and their actual meaning for partitioned TDs.

Logical Operation	Actual Operation
Is the Dirty bit 0?	Logical Dirty is 0 if the Dirty bit is 0 in the L1 SEPT entry and all L2 SEPT entries.
Is the Dirty bit 1?	Logical Dirty is 1 if the Dirty bit is 1 in the L1 SEPT entry or any of the L2 SEPT entries.
Atomically clear the Dirty bit and return its previous value.	Atomically test and clear the Dirty bit in L1 and each L2 SEPT entries. Return true if the Dirty bit was 1 in any of the SEPT entries.
Clear the Dirty bit.	Clear the Dirty bit in L1 and each L2 SEPT entries.
Set the Access and Dirty bits as a result of reading or writing to the TD private page by a TDX module flow operating in the context of the TD, e.g., TDG.MR.REPORT.	Set the A and D bits as part of the soft SEPT walk done by the TDX module.
Set the Access and Dirty bits as a result of translating GPAs to HPAs and caching the HPA values to be used by the CPU (as VMCS fields) and by the TDX module, e.g., as part of $L1 \rightarrow L2$ entry.	Set the A and D bits as part of the soft SEPT walk done by the TDX module.

Table 8.5: Logical Dirty Bit vs. Actual L1 and L2 Dirty Bits

8.5.2.4. Details: Pending Pages Considerations

As described in 8.5.1.7 above, page state may change from PENDING to MAPPED (by TDG.MEM.PAGE.ACCEPT) and vice versa (by TDG.MEM.PAGE.RELEASE). Both operations are considered a change of page attributes and thus require re-export if the page has been exported.

8.5.2.5. Details: Blocked Pages Considerations

As described in 8.5.1.6.2 above, no blocked pages are allowed during the export blackout period (PAUSED_EXPORT phase). This applies to the following SEPT entry states: BLOCKED, EXPORTED_BLOCKED, PENDING_BLOCKED, PENDING_EXPORTED_BLOCKED.

There are two reasons for that:

- Blocked pages are not exported. The blocked state is always an interim state; the host VMM blocks a page when it
 intends to perform some other operation, such as page removal.
- Not allowing blocked pages during the PAUSED_EXPORT phase simplified the tracking of export completeness by TDH.MEM.SCAN.COMP(DCHECK) and TDH.EXPORT.MEM.

8.5.2.6. Details: Memory Management Considerations

8.5.2.6.1. Details: Host-Side Memory Management Considerations

The tables below summarize special consideration for the interaction of non-blocking export and TD private memory management operations done by the host VMM and guest TD.

20

5

10

15

 Table 8.6: Host-Side Memory Management Considerations of Non-Blocking Export

Operation	Details
TDH.MEM.PAGE.AUG	A page that has been exported and then removed is put in an EXPORTED_REMOVED state (see below). While in this state, the host VMM can call TDH.MEM.PAGE.AUG to add a new page at the same GPA. The page state is set to EXPORTED_PENDING_MODIFIED.

Operation	Details
TDH.MEM.PAGE.DEMOTE	It is possible to demote a non-exported 2MB or 1GB page whose state is MAPPED, BLOCKED or PENDING or PENDING_BLOCKED.
	If the page state before the demote operation was MAPPED and a non-blocking export session is in the LIVE_EXPORT phase, the TDX module sets all the demoted pages' Dirty bits ¹ . Otherwise, it clears them.
	A page that has been exported is mapped at 4KB, so it can't be demoted.
TDH.MEM.PAGE.PROMOTE	It is possible to promote a 2MB or 1GB GPA range if all the small pages within that range are non-exported whose state is MAPPED or PENDING. The TDX module sets the promoted page's Dirty bit ² .
	It is not possible to promote a page that has been exported. Supporting that would complicate the SEPT state machine (e.g., handle the case where some 4KB pages have different states than other, add new states that indicate the need to export a CANCEL operation etc.). Note that the host VMM should already be able to handle situations where promotion is not possible, e.g., due to different L2 attributes set by L1.
TDH.MEM.PAGE.RELOCATE	A blocked page in the EXPORTED_BLOCKED or PENDING_EXPORTED_BLOCKED state can be relocated. This unblocks the page, and the page state is set to EXPORTED_MODIFIED or PENDING_EXPORTED_MODIFIED respectively.
TDH.MEM.PAGE.REMOVE	A page that has been exported can be removed. The page is put in an EXPORTED_REMOVED state to track the need for exporting a CANCEL operation for it, after which the page state is set to the regular FREE state.
TDH.MEM.SEPT.REMOVE	A page that has been exported and then removed is put in an EXPORTED_REMOVED state (see above). While in that state, the L1 SEPT entry mapping the page is not free. Thus, the L1 SEPT page where it resides cannot be removed.
TDH.MEM.RANGE.BLOCK	A page that has been exported can be blocked (and later removed or unblocked). Depending on whether the page is pending, the page is put in an EXPORTED_BLOCKED or PENDING_EXPORTED_BLOCKED state to track the need for exporting a CANCEL operation for it, after which the page state is set to the regular BLOCKED or PENDING_BLOCKED state respectively.
TDH.MEM.RANGE.UNBLOCK	A blocked page in the EXPORTED_BLOCKED or PENDING_EXPORTED_BLOCKED state can be unblocked. The page state is set to EXPORTED_MODIFIED or PENDING_EXPORTED_MODIFIED respectively.

8.5.2.6.2. Details: Guest-Side Memory Management Considerations

The guest TD is not directly aware of memory export; memory management operations are still available to it.

Table 8.7: Guest-Side Memory Management Considerations of Non-Blocking Export

Operation	Details
TDG.MEM.PAGE.ACCEPT	The guest TD may accept a PENDING_EXPORTED or a PENDING_EXPORTED_MODIFIED page. The page state becomes EXPORTED_MODIFIED so that TDH.MEM.SCAN.RANGE(DSCAN) or TDH.MEM.SCAN.COMP(DCHECK) will detect the need to re-export the page.

¹ With block-less demote, the h/w may still set the Dirty bit of the large page's SEPT entry while TDH.MEM.PAGE.DEMOTE is running and before it updated that SEPT entry to be a non-leaf entry. Thus, setting the Dirty bity of the small pages SEPT entries is necessary.

² This is done for simplification of implementation, since the page will have to be demoted in any case for migration.
Operation	Details			
TDG.MEM.PAGE.ATTR.WR	The guest TD may update the attributes of an EXPORTED or EXPORTED_MODIFIED. The page state becomes EXPORTED_MODIFIED so that TDH.MEM.SCAN.RANGE(DSCAN) or TDH.MEM.SCAN.COMP(DCHECK) will detect the need to re-export the page.			
	Similarly, the guest TD may update the attributes of a PENDING_EXPORTED, or PENDING_EXPORTED_MODIFIED page. The page state becomes PENDING_EXPORTED_MODIFIED.			
TDG.MEM.PAGE.RELEASE	If supported, the guest TD may release an EXPORTED or an EXPORTED_DIRTY page. The page state becomes PENDING_EXPORTED_MODIFIED so that TDH.MEM.SCAN.RANGE(DSCAN) or TDH.MEM.SCAN.COMP(DCHECK) will detect the need to re-export the page.			

8.5.2.7. Details: Export Completeness Tracking

8.5.2.7.1. Overview

After the TD is paused by TDH.EXPORT.PAUSE, the TDX module needs to enforce memory export completeness before a start token can be exported to enable the TD to run on the destination platform. This means that the exported memory image reflects the up-to-date state of the source TD memory image. Note that the migration protocol ensures that the imported memory image is in sync with the exported one.

The host VMM is expected to do the following:

- 1. Run a comprehensive scan of the TD's GPA space, by calling TDH.MEM.SCAN.COMP(DCHECK).
- 10 2. Export all memory pages reported by TDH.MEM.SCAN.COMP(DCHECK) as export candidates, using TDH.EXPORT.MEM.

The above operations can be done concurrently.

8.5.2.7.2. Memory Counters and Checking by TDH.EXPORT.TRACK(DONE)

The TDX module holds the following counters in TDCS. They are reset when a migration session starts (TDH.EXPORT.STATE.IMMUTABLE).

- **MEM_COUNT** Counts the number of TD private memory pages, in multiples of 4KB. MEM_COUNT is incremented when a new page is added to the TD, and decremented when a page is removed from the TD.
- **MIG_COUNT** Counts of the number of 4KB pages that have been exported at least once in the current migration session.
- 20 **DIRTY_COUNT** Counts the number of 4KB pages that have been exported but either their content or their attributes have changed since, so a re-export is required. Counting is in units of 4KB. DIRTY_COUNT is incremented by any function which transitions the page state out of EXPORTED or PENDING_EXPORTED; it is decremented by TDH.EXPORT.MEM.

TDH.EXPORT.TRACK(DONE) checks the counter values as follows:

- TDH.MEM.SCAN.COMP(DCHECK) was called and completed a comprehensive scan of the TD's GPA space.
 - DIRTY_COUNT is 0 i.e., all memory pages that required a re-export were indeed re-exported.
 - If the TDX module does not support post-copy for non-blocking export, or if the TD is enabled for TDX Connect (TD_PARAMS.CONFIG_FLAGS.TDX_CONNECT), MIG_COUNT is equal to MEM_COUNT – i.e., all memory pages have been exported.
- 30 Note: If the TDX module supports post-copy for non-blocking export and the TD is not enabled for TDX Connect, the host VMM is allowed not to export some of the memory image at this point. Typically, unexported pages are exported on-demand during the post-copy phase. Failure to export (but not re-export) memory is considered a denial of service, which is not prevented by TDX.

The host VMM can read the above counters, if required, using TDH.MNG.RD.

Draft

Tracking of export completeness is based on the following properties, which are true during the PAUSED_EXPORT state, once the TD is paused by TDH.EXPORT.PAUSE and before a start token is generated by TDH.EXPORT.TRACK(DONE):

- TD private memory is immutable.
 - The TD is paused. 0
 - No TDIs are bound. 0
 - Host-side memory management operations are restricted, as described in 8.5.1.6.2. Specifically, no page may 0 be blocked. Note that BLOCKED and PENDING BLOCKED are intermediate states, as preparations for some memory management operation. Not allowing blocked pages simplifies tracking for correctness.
- TDH.MEM.SCAN.COMP(DCHECK) scans each TD private page exactly once.
- TDH.EXPORT.MEM can only export a page once. 10

Backward Compatibility 8.5.2.7.3.

The migrated TD may had been created by an older TDX module that didn't support non-blocking export, and later the TDX module was updated using TD-preserving update. In this case, TDH.MEM.SCAN may be less efficient; some metadata maintained by a TDX module that supports non-blocking export during the TD lifecycle for optimizing the scan may be missing.

15

5

8.6. Memory Import

8.6.1. **Host VMM Perspective**

This section describes memory import from the host VMM perspective. This is a simplified view; a more detailed view is discussed later.

20 8.6.1.1. In-Order Import Phase

8.6.1.1.1. **Overview of In-order Import**

During the in-order import phase, a page may be imported multiple times. In addition, a page import may be cancelled. Ordering is maintained by the MBMD's MB COUNTER and the requirement that a page can only be imported once per migration epoch.

8.6.1.1.2. **Memory Management During In-Order Import** 25

During the in-order import phase, no blocking and no TLB tracking is required, since the destination TD is not running yet.

Addition and removal of Secure EPT pages are allowed during the in-order phase – they are required as part of building the TD on the destination platform. To import a page that has L2 pages mappings, the host VMM on the destination platform must have built the L2 SEPT for the applicable L2 VMs, using TDH.MEM.SEPT.ADD, down to the proper page mapping level.

30

Page management operations are prohibited during the in-order import phase. These include:

- Page addition by TDH.MEM.PAGE.ADD and TDH.MEM.PAGE.AUG
- Page removal by TDH.MEM.PAGE.REMOVE
- Page promotion and demotion by TDH.MEM.PAGE.PROMOTE and TDH.MEM.PAGE.DEMOTE
- Page relocation by TDH.MEM.PAGE.RELOCATE
 - GPA range blocking and unblocking by TDH.MEM.RANGE.BLOCK and TDH.MEM.RANGE.UNBLOCK •

8.6.1.2. **Out-of-Order import Phase**

8.6.1.2.1. **Overview of Out-of-Order Import**

40

35

When the out-of-order import phase begins, any pages that have been imported are designed to be up to date. A page may only be imported if its GPA mapping does not exist yet (SEPT entry's state is FREE). An import attempt of a page that has been imported before during the out-of-order phase is dropped but not considered an error; this is normally the result of the same page being migrated on a high-priority queue. Source memory state is immutable, so ordering is not required.

8.6.1.2.2. Memory Management During Out-of-Order Import

During the out-of-order import phase, TLB tracking is required in the LIVE_IMPORT OP_STATE, since the TD may be running on the destination platform.

Secure EPT Management

5 Addition and removal of Secure EPT pages are allowed during the out-of-order phase – they are required as part of building the TD on the destination platform.

Page Addition

TDH.MEM.PAGE.ADD is prohibited (since the TD had been finalized before the migration session began). TDH.MEM.PAGE.AUG is allowed after TDH.IMPORT.COMMIT, when the TD is allowed to run on the destination platform while the import session continues.

10

If a page was not imported but was added locally (TDH.MEM.PAGE.AUG), this is equivalent to the VMM removing a page without coordinating with the TD, then adding a new page. The TD should not accept (TDG.MEM.PAGE.ACCEPT) such a page since from its point of view this is a page that already exists in its GPA space. The secure EPT entry state for the locally added page is PENDING, and if a page is imported to the same GPA, import will fail.

15 **Promotion and Demotion**

Page promotion and demotion are allowed during the out-of-order phase.

Page Removal

Page removal (TDH.MEM.PAGE.REMOVE) is allowed during the out-of-order import phase. However, the page's SEPT entry is not marked as FREE when the page is removed. Instead, the SEPT entry state is set to REMOVED. The REMOVED state is equivalent to the FREE state, except for the following limitations that apply after TDH IMPORT COMMIT, when

20

state is equivalent to the FREE state, except for the following limitations that apply after TDH.IMPORT.COMMIT, when the TD is allowed to run on the destination platform while the import session continues, until TDH.IMPORT.END:

- Page import is not allowed to this GPA.
- Removal of the parent SEPT page is not allowed.

Page Relocation

25 Page relocation is supported at any stage without any changes.

8.6.1.3. In-Place Import

In-place import repurposes the physical pages holding the imported data as private memory pages that hold the decrypted data. This saves the host VMM on the destination platform the need to allocate memory for the imported data, at the cost of a small fixed-sized intermediate buffer that needs to be held by Intel TDX Module, and some other complications. In-place import may be selected for each page imported for the first time, or following a previous CANCEL, but not for re-import of a new version of a previously imported page.

30

8.6.2. Updated: Details of Memory Import

8.6.2.1. Updated: Details: In-Order Import Phase

8.6.2.1.1. Updated: Details: Overview of In-Order Import

35 Once out-of-order import phase begins, any pages that have been imported are designed to be up to date.

The state diagram below shows the L1 SEPT entry state during in-order import. The color-coding convention used in this diagram is described in 8.2.





Figure 8.15: Updated: Partial L1 SEPT Entry State Diagram for Page In-Order Import Phase

An SEPT page can only be removed if all its entries are FREE; specifically, it can't be removed if any entry state is REMOVED.

5 8.6.2.1.2. Details: Enforcing a Single Import Operation per Migration Epoch

When a page is imported during the in-order phase, the current migration epoch is recorded. Page re-import and import cancel operations check the recorded migration epoch. For the import to succeed, it should be older than the current migration epoch.

When a page import is cancelled during the in-order phase, the physical page is removed but its SEPT entry is put into a
REMOVED state, and the current migration epoch is recorded. For partitioned TDs, any L2 SEPT entries that map the page become L2_FREE. A page first-time import operation checks the recorded migration epoch. For the import to succeed, it should be older than the current migration epoch.

8.6.2.1.3. Updated: Importing L2 Page Mappings During In-Order Import

For partitioned TDs, L2 page mappings are imported as part of a page import, and follow these rules:

- When a page is imported or re-imported, its L2 page mappings, if any, are created or updated in the L2 SEPT leaf entry.
 - When a page is re-imported, and a previously imported page mapping does not exist in the new import, the L2 SEPT leaf entry is set to L2_FREE.
 - When a page import is cancelled, its L2 SEPT leaf entries, if any, are set to L2_FREE.
- The state diagram below shows the L2 SEPT entry state during in-order import. The color-coding convention used in this diagram is described in 8.2.





Figure 8.16: Updated: Partial L2 SEPT Leaf Entry State Diagram for L2 Page Mapping Import During the In-Order Phase

8.6.2.2. Details: Out-of-Order import Phase

5 8.6.2.2.1. Details: Overview of Out-of-Order Import

The state diagram below shows the L1 SEPT entry state during out-of-order import. The color-coding convention used in this diagram is described in 8.2.



Figure 8.17: Partial L1 SEPT Entry State Diagram for Page Out-of-Order Import Phase

An SEPT page can only be removed if all its entries are FREE; specifically, it can't be removed if any entry state is REMOVED.

5 8.6.2.2.2. Details: Importing L2 Page Mappings During Out-of-Order Import

During out-of-order import, L2 page mappings may be created as part of a page import. L2 SEPT leaf entries don't need the REMOVED state; thus, the L2 state transitions are very simple. The state diagram below shows the L2 SEPT entry state during out-of-order import. The color-coding convention used in this diagram is described in 8.2.

Draft



Figure 8.18: Partial L2 SEPT Leaf Entry State Diagram for Page Out-of-Order Import Phase

8.7. Secure EPT Concurrency Considerations

8.7.1. Overview

5 To support high performance migration, memory migration interface functions are allowed to run concurrently on multiple LPs. However, no concurrent operation is allowed on any single Secure EPT entry. Interface functions that work on specific Secure EPT entries acquire an exclusive lock to that entry.

8.7.2. New: GPA List Processing Implications

On the source side, export functions that process a GPA list (TDH.EXPORT.BLOCKW, TDH.EXPORT.MEM, 10 TDH.EXPORT.RESTORE) that encounter a busy SEPT entry skip the processing of that GPA. In this case, they set the STATUS field of the GPA list entry to SEPT_ENTRY_BUSY_HOST_PRIORITY and the OPERATION field to NOP. The host VMM can scan the GPA list to detected skipped entries and retry the operation later.

On the destination side, if TDH.IMPORT.MEM encounters a busy SEPT entry, its behavior depends on the import phase. During the in-order phase, any SEPT entry error causes the import session to fail. In the out-of-order phase, TDH.IMPORT.MEM skips the busy entry, similar to TDH.EXPORT.MEM.

8.8. Moved & Updated: Security Analysis: Achieving Memory Migration Security Objectives

8.8.1. Updated: General

The key security design goal for TD Private memory migration is to ensure integrity and freshness of the TD private memory at the destination TD after migration – this helps ensure that a malicious VMM cannot execute the TD after migration with any stale or modified data.

20

25

15

Integrity of memory includes the memory contents as well as the guest physical to host physical mapping and attributes that control TD access to private memory.

Using PAMT and Secure EPT, Intel TDX Module enforces the following properties for TD private GPA accesses:

Unique TD Association: A physical page used as a TD private page, Secure EPT page or a control structure can only be assigned to single guest TD.

Unique GPA Mapping: A TD private page or a Secure EPT page can be mapped at most by single guest TD GPA.

Section 2: TD Migration Architecture Specification

These security properties are maintained for a TD during migration with some additional functionality afforded to allow for live migration.

- Private TD pages and Secure EPT entries (for partitioned TDs, this includes L2 Secure EPT entries) are initialized in a single operation (via TDH.IMPORT.MEM) for pages migrated using TDH.EXPORT.MEM. Like the pre-conditions for the non-migration TDH.MEM.PAGE.ADD, the parent Secure EPT entry must be free (unmapped).
- When write-blocking based export is used on the source platform, a private page may be mapped as non-writable (a.k.a. blocked for writing) to allow for the page contents to be exported. For partitioned TDs, this any L2 mapping of a page is also mapped as non-writable. Following from previous security requirements, this mapping update also requires TLB tracking to help ensure that no active writable cached GPA address translations exist to the to-be-migrated GPA range.
- When non-blocking export is used on the source platform, a private page is detected as requiring export based on an atomic test-and-clear of the SEPT entry's Dirty bit set by the h/w. For partitioned TDs, this includes applicable L2 SEPT entries. Following from previous security requirements, this scan also requires TLB tracking prior to export to help ensure that no active writable cached GPA address translations prevent further updates of the SEPT entry's Dirty bit.
- 15

5

10

- For 1GB and 2MB pages, secure EPT mapping demotion (to a 4KB page size) is required as a pre-condition to exporting contents of a page for migration.
- The Migration key used for exporting and important TD memory and CPU state is distinct from keys used for other operations such as Paging.

20 8.8.2. Updated: Preventing Usage of Stale Memory Copies due to Mis-Ordering

Multiple versions of a page that are exported in a certain order are imported in the same order. This is achieved using the following mechanisms:

- Migration streams help ensure ordering within each stream.
- Migration epochs, limiting migration of a specific page to no more than one time per epoch, help ensure total ordering per page.
- 8.8.3. New: Enforcing Export of the Entire Memory Image

TDX Connect

When the TD is enabled for TDX Connect, while TDIs are bound and have DMA access to the TD memory, DMA must not fail. Therefore, TDX Connect requires that all the TD mapped pages remain pinned as long as TDI can access them. The entire TD private memory image must be migrated intact. This is enforced by the TDX module as described below.

30

45

25

Only non-blocking export is supported when TDX Connect is enabled for a TD. A counter of unexported pages (TDCS.UNEXPORTED_COUNT) is maintained by the TDX module at the source platform based on a comprehensive scan of the TD's GPA space by TDH.MEM.SCAN.COMP(DCHECK) and on TDH.EXPORT.MEM. If the value of that counter is not 0, then TDH.EXPORT.TRACK fails. See 8.5.2.7 for details.

35 Non TDX Connect

If the TD is not configured for TDX Connect, the TDX module does not enforce exporting all memory. If a page is not migrated, an EPT violation happens when the guest TD attempts to access it on the destination platform. This is normal behavior for the post-copy method of migration, and at the worst case is considered a denial of service.

8.8.4. New: Non-Blocking Export: Detecting Memory State Change

- 40 Memory state change is detected by scanning the SEPT tree and detecting pages where the Dirty bit is set.
 - 1. TDH.MEM.SCAN.RANGE(DSCAN) does the following with the SEPT entry exclusively locked (but Accessed/Dirty bits can be set by h/w):
 - 1.1. Detect that an EXPORTED page is dirty (either L1 or any L2 SEPT entries' Dirty bit is 1).
 - 1.2. Clears the Dirty bit in each of L1 and L2 SEPT entries.
 - 1.3. Update the page state to EXPORTED_MODIFIED.
 - 1.4. Records the TD epoch in the page's PAMT.
 - 2. TDH.EXPORT.MEM does the following with the SEPT entry exclusively locked (but Accessed/Dirty bits can be set by h/w):
 - 2.1. Find that the page state is MODIFIED.
- 50 2.2. Check the TLB tracking conditions.

2.3. Export the page.

5

20

2.4. Update the page state to EXPORTED.

The above sequence helps ensure the following:

- Setting the EXPORTED_MODIFIED state in step 1.3 means that this page will be detected later by DSCAN or DCHECK and will be re-exported as a precondition to start token generation (see below).
- The order of TLB tracking check (in step 2.2) before the page is exported (in step 2.3) means that any write access by either the TD s/w, by the TDX module working on behalf of the TD s/w or by a TDI will result in an SEPT walk and setting of the Dirty bit. The page will be detected later by DSCAN or DCHECK and will be re-exported again as a precondition to start token generation (see below).

10 8.8.5. Updated: Preventing Usage of Stale Memory Copies due to Failure to Re-export

Running the destination TD with a stale copy of a memory page, because the source VMM failed to re-export a newer copy of a page, is prevented as described below.

Assume for example that the source VMM exported an older version (v1) of page but never re-exported a newer version (v2) of that page. In this case, generating a start token by TDH.EXPORT.TRACK is prevented.

- With write-blocking based export, a counter of dirty pages (TDCS.DIRTY_COUNT) is maintained by the TDX module at the source platform, based on tracking the TDH.EXPORT.MEM and TDX.EXPORT.UNBLOCKW operations. If the value of that counter is not 0, then TDH.EXPORT.TRACK fails. See 8.4.2.38.4.2.3 for details.
 - With non-blocking export, the same counter is maintained by the TDX module based on a comprehensive scan of the TD's GPA space by TDH.MEM.SCAN.COMP(DCHECK) and on TDH.EXPORT.MEM. If the value of that counter is not 0, then TDH.EXPORT.TRACK(DONE) fails. See 8.5.2.7 for details.

8.8.6. Updated: Preventing Usage of Missing or Stale Memory Copies due to Failure to Import

Running the destination TD with a missing or a stale copy of a memory page, because the destination VMM failed to import a copy of a page, is prevented as described below.

The in-order phase commitment protocol is designed to ensure that the export will fail, and the destination TD will not run. As discussed above, TDH.EXPORT.TRACK(DONE) generates a start token that is dependent of the exact export sequence; it checks that no unexported newer versions of previously exported pages remain. If the TD is enabled for TDX Connect, TDH.EXPORT.TRACK(DONE) checks that no unexported pages remain.

The start token is verified by TDH.IMPORT.TRACK; the out-of-order migration phase may start, and the destination TD may run only if the start token verifies correctly. For migration session control details see Ch. 6.

30 8.8.7. Preventing Usage of Stale Memory GPA Mapping and Attributes

The destination TD is preventing from running with a copy of a memory page with stale GPA mapping, access permissions and other attributes (for partitioned TD, this included L2 mapping) as follows:

- GPA mappings and attributes are migrated together with their respective pages.
- Any change to GPA mappings or attributes is considered a change to the page and requires re-migration.

35 8.8.8. Out-Of-Order Phase and Its Usage for Post Copy

In the in-order import phase the host VMM can import pages for GPA addresses that are free, and it may also reload newer versions of pages to previously imported and present GPA addresses. In the out-of-order import phase, import is only allowed to non-present GPA addresses. At this stage, all memory state on the source platform is designed to be immutable, and the latest version of all pages exported so far will be imported. Thus, the order of out-of-order import is not relevant a present that memory content available in order phase sent that memory exponent of a state of the imported during the interval.

40 not relevant – except that memory content exported during the in-order phase can't be imported during the out-of-order phase. This allows using a separate migration stream for high-priority, low-latency updates, e.g., to implement **post-copy** by allowing the TD to run and migrate memory pages on demand at a high priority, based on EPT violation.

8.9. Updated: Memory Migration Interface Functions Summary

45

This section provides a short summary of the memory migration interface functions. A detailed specification is provided in [TDX Module ABI Spec].

Name	Export Mode	Description	Preconditions
TDH.EXPORT.BLOCKW	Write Blocking	Block a list of 4KB pages for writing, as a preparation for export.	Write-blocking based export session is in the in-order phase and the TD has not been paused yet
TDH.EXPORT.MEM	Any	Export, re-export or sends an export cancellation request for a list of 4KB pages.	Export session is in progress
TDH.EXPORT.RESTORE	Any	Restore a list of 4KB pages after an export session abort.	Export session is not in progress
TDH.EXPORT.UNBLOCKW	Write Blocking	Unblock a single 4KB page that has been blocked for writing.	The TD is configured for write- blocking based export, and either an export session is in progress but committed export phase has not begun, or the TD is allowed to run
TDH.IMPORT.MEM	Any	Import, re-import or cancel a previous import for a list of 4KB pages.	Import session is in progress
TDH.MEM.SCAN.COMP(DCHECK)	Non- Blocking	Do a comprehensive scan of the TD private memory and detect pages that must be exported for the export session to complete.	Non-blocking export session is in the EXPORT_PAUSED state.
TDH.MEM.SCAN.CONFIG	Non- Blocking	Configure comprehensive scan parameters	
TDH.MEM.SCAN.RANGE(DSCAN)	Non- Blocking	Scan a TD private memory GPA range and detect pages that need to be exported.	
TDH.MEM.SCAN.RANGE (EXPORT_RESTORE)	Any	Scan the TD private memory and restore SEPT after an export session abort.	Export session is not in progress
TDH.MEM.SCAN.RESET	Non- Blocking	Reset comprehensive scan state	

Table 8.8: Memory Migration Interface Functions