



NIST Interagency Report
NIST IR 8320E ipd

Hardware-Enabled Security

Confidential Computing of Data in Cloud Workloads

Initial Public Draft

Michael Bartock
Murugiah Souppaya
Timothy Knoll

This publication is available free of charge from:
<https://doi.org/10.6028/NIST.IR.8320E.ipd>

**NIST Interagency Report
NIST IR 8320E ipd**

Hardware-Enabled Security
Confidential Computing of Data in Cloud Workloads

Initial Public Draft

Michael Bartock
*Computer Security Division
Information Technology Laboratory*

Murugiah Souppaya
**Former NIST employee; all work for this publication was done while at NIST.*

Timothy Knoll
Intel Corporation

This publication is available free of charge from:
<https://doi.org/10.6028/NIST.IR.8320E.ipd>

May 2026



U.S. Department of Commerce
Howard Lutnick, Secretary

National Institute of Standards and Technology
Craig Burkhardt, Acting Under Secretary of Commerce for Standards and Technology and Acting NIST Director

Certain equipment, instruments, software, or materials, commercial or non-commercial, are identified in this paper in order to specify the experimental procedure adequately. Such identification does not imply recommendation or endorsement of any product or service by NIST, nor does it imply that the materials or equipment identified are necessarily the best available for the purpose.

There may be references in this publication to other publications currently under development by NIST in accordance with its assigned statutory responsibilities. The information in this publication, including concepts and methodologies, may be used by federal agencies even before the completion of such companion publications. Thus, until each publication is completed, current requirements, guidelines, and procedures, where they exist, remain operative. For planning and transition purposes, federal agencies may wish to closely follow the development of these new publications by NIST.

Organizations are encouraged to review all draft publications during public comment periods and provide feedback to NIST. Many NIST cybersecurity publications, other than the ones noted above, are available at <https://csrc.nist.gov/publications>.

NIST Technical Series Policies

[Copyright, Use, and Licensing Statements](#)

[NIST Technical Series Publication Identifier Syntax](#)

Publication History

Approved by the NIST Editorial Review Board on YYYY-MM-DD
Supersedes NIST Series XXX (Month Year) DOI

How to Cite this NIST Technical Series Publication:

Bartock M, Souppaya M, Knoll T (2026) Hardware-Enabled Security: Confidential Computing of Data in Cloud Workloads. (National Institute of Standards and Technology, Gaithersburg, MD), NIST Interagency Report (IR) NIST IR 8320E ipd. <https://doi.org/10.6028/NIST.IR.8320E.ipd>

NIST Author ORCID iDs

Michael Bartock: 0000-0003-0875-4555
Murugiah Souppaya: 0000-0002-8055-8527

Public Comment Period

May 29, 2026 – July 13, 2026

Submit Comments

hwsec@nist.gov

National Institute of Standards and Technology
Attn: Computer Security Division, Information Technology Laboratory
100 Bureau Drive (Mail Stop 8930) Gaithersburg, MD 20899-8930

Additional Information

Additional information about this publication is available at <https://csrc.nist.gov/pubs/ir/8320/e/ipd>, including related content, potential updates, and document history.

All comments are subject to release under the Freedom of Information Act (FOIA).

1 **Abstract**

2 Confidential computing addresses data security and privacy concerns for organizations that
3 move sensitive workloads to the cloud. It is a critical advancement that enables the encryption
4 of data both while it is being processed in memory and in active use. As cloud adoption
5 continues to grow, confidential computing will play a pivotal role in improving security and
6 privacy in cloud environments. This report describes an effective approach to protecting data
7 being acted upon by artificial intelligence workloads on cloud infrastructures so that the
8 datasets are protected from malware, data theft, and other security-related vulnerabilities.

9 **Keywords**

10 Confidential computing; cryptographic key; hardware-enabled security; machine identity;
11 machine identity management; trusted execution environment (TEE).

12 **Reports on Computer Systems Technology**

13 The Information Technology Laboratory (ITL) at the National Institute of Standards and
14 Technology (NIST) promotes the U.S. economy and public welfare by providing technical
15 leadership for the Nation’s measurement and standards infrastructure. ITL develops tests, test
16 methods, reference data, proof of concept implementations, and technical analyses to advance
17 the development and productive use of information technology. ITL’s responsibilities include
18 the development of management, administrative, technical, and physical standards and
19 guidelines for the cost-effective security and privacy of other than national security-related
20 information in federal information systems.

21

22

23

24 **Audience**

25 The primary audiences for this report are security professionals, such as security engineers and
26 architects; system administrators and other information technology (IT) professionals
27 responsible for securing physical or virtual platforms; and hardware, firmware, and software
28 developers who may be able to leverage hardware-enabled security techniques and
29 technologies to improve machine identity management and protection.

30

31 **Call for Patent Claims**

32 This public review includes a call for information on essential patent claims (claims whose use
33 would be required for compliance with the guidance or requirements in this Information
34 Technology Laboratory (ITL) draft publication). Such guidance and/or requirements may be
35 directly stated in this ITL Publication or by reference to another publication. This call also
36 includes disclosure, where known, of the existence of pending U.S. or foreign patent
37 applications relating to this ITL draft publication and of any relevant unexpired U.S. or foreign
38 patents.

39 ITL may require from the patent holder, or a party authorized to make assurances on its behalf,
40 in written or electronic form, either:

- 41 a) assurance in the form of a general disclaimer to the effect that such party does not hold
42 and does not currently intend holding any essential patent claim(s); or
- 43 b) assurance that a license to such essential patent claim(s) will be made available to
44 applicants desiring to utilize the license for the purpose of complying with the guidance
45 or requirements in this ITL draft publication either:
 - 46 i. under reasonable terms and conditions that are demonstrably free of any unfair
47 discrimination; or
 - 48 ii. without compensation and under reasonable terms and conditions that are
49 demonstrably free of any unfair discrimination.

50 Such assurance shall indicate that the patent holder (or third party authorized to make
51 assurances on its behalf) will include in any documents transferring ownership of patents
52 subject to the assurance, provisions sufficient to ensure that the commitments in the assurance
53 are binding on the transferee, and that the transferee will similarly include appropriate
54 provisions in the event of future transfers with the goal of binding each successor-in-interest.

55 The assurance shall also indicate that it is intended to be binding on successors-in-interest
56 regardless of whether such provisions are included in the relevant transfer documents.

57 Such statements should be addressed to: hwsec@nist.gov

58

59	Table of Contents	
60	1. Introduction	1
61	1.1. Purpose and Scope	1
62	1.2. Terminology	1
63	1.3. Document Structure.....	2
64	2. Challenges with Protecting Data in Use	3
65	2.1. Objective	3
66	2.2. Goals.....	4
67	3. Stage 0: Configure a TEE-Capable VM	5
68	4. Stage 1: Remote Attestation Service Configuration	7
69	5. Stage 2: Relying Party Service Configuration	8
70	6. Stage 3: Execute the Key-Release Workflow	9
71	References	11
72	Appendix A. Solution Architecture	12
73	Appendix B. Trusted Execution Environment Implementation	14
74	B.1. Creating a VM With Intel TDX on Microsoft Azure	14
75	B.2. Deploying Bastion	16
76	B.3. Configuring TDX Prerequisites	16
77	Appendix C. Key Management and Relying Party Implementation	17
78	C.1. Deploy the Key Broker Service	17
79	Appendix D. AI Model Encryption Implementation	19
80	D.1. Prepare the Sample AI Workload Image	19
81	D.2. Provision the Relying Party and Encrypt the AI Model.....	19
82	Appendix E. Remote Attestation Implementation	22
83	E.1. Execute the Key-Release Workflow	22
84	Appendix F. List of Symbols, Abbreviations, and Acronyms	25
85	List of Figures	
86	Fig. 1. Sample workload key release	9
87	Fig. 2. Cloud workload key release	12
88	Fig. 3. Creating a cloud VM	14
89	Fig. 4. Creating a cloud VM (continued)	15
90		

91 1. Introduction

92 1.1. Purpose and Scope

93 Confidential computing addresses one of the most pressing concerns for organizations that
94 move sensitive workloads to the cloud: data privacy. In traditional cloud environments, data is
95 typically encrypted while at rest and in transit but decrypted during processing, which exposes
96 potential vulnerabilities. Confidential computing fills this gap by enabling the encryption of data
97 even while it is being processed in memory, closing the gap that traditional at-rest and in-
98 transit encryption leaves unaddressed.

99 This innovation is important for industries that handle highly sensitive information, such as
100 finance, healthcare, and government. With confidential computing, organizations can
101 confidently move their sensitive workloads to the cloud, where the data being processed can be
102 made opaque to the cloud provider, thus reducing the risk of access by privileged infrastructure
103 software. This not only enhances trust in cloud environments but also opens up new
104 possibilities for collaboration and innovation, such as sharing data securely across different
105 organizations with greater protection against unauthorized access. As cloud adoption continues
106 to grow, confidential computing will play a pivotal role in improving security and privacy in
107 cloud environments.

108 This report describes an effective approach to protecting data being acted upon by artificial
109 intelligence (AI) workloads (i.e., data used for inference) on cloud infrastructures so that the
110 datasets are protected from malware, data theft, and other security-related vulnerabilities. The
111 reference implementations detailed in the appendices show how to perform inferences on
112 sample health data with an encrypted artificial intelligence model.

113 The prototype implementation presented in this report is only one possible way to solve the
114 security challenges. It is not intended to preclude the use of other products, services, or
115 techniques that can also solve the problem adequately, nor is it intended to preclude the use of
116 any cloud products or services that are not specifically mentioned in this report.

117 This report builds upon the terminology and concepts described in NIST Interagency Report (IR)
118 8320, *Hardware-Enabled Security: Enabling a Layered Approach to Platform Security for Cloud
119 and Edge Computing Use Cases* [1].

120 1.2. Terminology

121 For consistency with related NIST publications, this report uses the following definitions for
122 trust-related terms:

- 123 • **Trust:** “The confidence one element has in another that the second element will behave
124 as expected” [2]
- 125 • **Trusted:** An element that another element relies upon to fulfill critical requirements on
126 its behalf

- 127 • **Trustworthy:** Worthy of being trusted to fulfill whatever critical requirements may be
128 needed

129 1.3. Document Structure

130 This document is organized into the following sections and appendices:

- 131 • Section 2 defines the objective for the prototype implementation and the intermediate
132 goals to be met in order to achieve the objective.
- 133 • Sections 3 through 6 describe the four stages of prototype implementation.
- 134 • The References section lists the references cited throughout this document.
- 135 • Appendix A provides an overview of the solution architecture.
- 136 • Appendix B describes how to enable a trusted execution environment (TEE).
- 137 • Appendix C describes how to set up key management and relying party services.
- 138 • Appendix D describes a reference AI model encryption implementation.
- 139 • Appendix E describes remote attestation of the TEE implementation.

140 2. Challenges with Protecting Data in Use

141 In traditional computing environments, data stored in memory would be unencrypted while
142 applications performed actions on it through the system’s processor. However, recent
143 technological advances have allowed mechanisms to encrypt memory space and restrict which
144 processes on a system can decrypt that memory while in use. The ultimate goal is to be able to
145 use hardware-enabled mechanisms to protect data being processed by an AI model so that it is
146 accessible only to authorized code within the trust domain. This involves encrypting the AI
147 model’s memory while in use and making sure that only compute resources that can attest to
148 their hardware capabilities are given the decryption key to access the memory space.

- 149 • **Stage 0: Configure a TEE-capable virtual machine (VM)**

150 Many cloud service providers offer virtual machines that have TEE capabilities enabled.
151 Section 3 describes the benefits of utilizing a TEE and how it can be configured within a
152 CSP.

- 153 • **Stage 1: Remote attestation service configuration**

154 Workloads that use confidential computing capabilities require a way to be trusted that
155 they are implemented and enabled. Section 4 describes how remote attestation of
156 workload achieves this goal.

- 157 • **Stage 2: Relying party configuration**

158 Once a workload has attested that it has confidential computing capabilities enabled,
159 policies can be applied for what kind of information it can be allowed to access. Section
160 5 describes how a relying party can make these types of determinations based on
161 attestation reports.

- 162 • **Stage 3: Execute the key-release workflow**

163 Section 6 brings the previous stages together to depict how a workload goes from
164 instantiation to processing sensitive information by using protected keys in a
165 confidential computing environment.

166 2.1. Objective

167 A reference implementation that protects AI workloads with confidential computing introduces
168 a powerful layer of security by leveraging encryption keys for the AI model memory space for
169 the data it processes. In the context of AI, both the model and the data it uses are valuable
170 assets. The model — often developed through significant research and computational effort —
171 holds intellectual property and represents a competitive advantage for organizations. Similarly,
172 the data, which could contain sensitive or proprietary information, is crucial for accurate and
173 effective AI processing. Confidential computing can safeguard both of these elements by
174 isolating them in secure enclaves where they are encrypted, even when processed in the cloud
175 or while in active memory [3].

176 Encrypting the AI model provides strong protection for both the model itself and the data it
177 processes. By securing the model with an encryption key that is only accessible to the deploying
178 entity, organizations can restrict model execution to a trusted, controlled environment. This
179 significantly reduces the risk of unauthorized access or tampering, even from privileged
180 infrastructure software, when the platform is patched and its configuration is verified through
181 attestation. Because the model is encrypted, any data it processes is also protected from
182 indirect exposure. Attackers face a substantially higher barrier to extracting insights,
183 parameters, or hidden representations from the model to infer sensitive details about the input
184 data. Instead, decryption and execution occur within a secure enclave, where both the model
185 and the data can interact safely within hardware-enforced isolation.

186 This dual protection is especially important in regulated industries, where both intellectual
187 property (i.e., the model) and sensitive information (i.e., the data) must be secured. Encrypting
188 the model preserves proprietary algorithms and reinforces compliance with data protection
189 regulations by providing hardware-enforced controls that are designed to prevent information
190 handled by the model from being misused or exposed outside of authorized computations.

191 By combining model encryption with confidential execution, organizations can confidently
192 deploy AI in the cloud while safeguarding both their innovations and the sensitive data that
193 those models process.

194 **2.2. Goals**

195 The goal of this paper is to explore and detail the architecture that implements the protection
196 of AI workloads using confidential computing, with a particular focus on the use of encryption
197 keys for the AI model and implications for the data it processes. Examining this architecture will
198 illustrate how confidential computing can be effectively applied to safeguard both the
199 intellectual property embedded in AI models and the sensitive data they process, even when
200 operating in potentially untrusted cloud environments.

201 This discussion will delve into the technical specifics of how the architecture isolates and
202 encrypts AI models and data within secure enclaves designed so that only authorized
203 computations can access and process them. Additionally, this paper will discuss the rationale
204 behind using encryption keys and the benefits of this approach in terms of enhanced security
205 and compliance with data protection regulations. Through this exploration, readers will gain a
206 clear understanding of how this architecture is implemented and the practical implications for
207 securing AI workloads in the cloud.

208

209 **3. Stage 0: Configure a TEE-Capable VM**

210 TEE attestation is crucial because it provides verifiable proof that an application or workload is
211 running within a secure and isolated environment. In the context of cloud computing, where
212 workloads often run on shared infrastructure, the ability to trust that data and computations
213 are securely isolated is essential. TEE attestation allows cloud providers to demonstrate to
214 users that their workloads are running in a TEE whose hardware and software configuration
215 matches expected values. This assurance is key to building trust between cloud providers and
216 users, especially when handling sensitive or critical data.

217 Furthermore, TEE attestation is important for maintaining the integrity and confidentiality of
218 the data being processed. By using attestation, organizations can verify that the code running in
219 the TEE matches expected measurements and has not been modified in unexpected ways. This
220 is particularly important in scenarios where the integrity of the data and the computation is
221 paramount, such as in financial services, healthcare, and government applications. TEE
222 attestation provides a mechanism for verifying the trustworthiness of the execution
223 environment to securely and reliably process sensitive workloads in potentially untrusted cloud
224 environments.

225 Hardware-based confidential computing provides robust protection for AI data by leveraging
226 secure enclaves within the processor to create isolated environments in which sensitive
227 computations can occur without exposing data to the rest of the system under normal
228 operation. These secure enclaves are designed to protect code and data from access by the
229 operating system and hypervisor, with isolation enforced in hardware. When AI models and the
230 data they process are loaded into these enclaves, they are encrypted and only decrypted within
231 the enclave's boundaries, significantly raising the barrier to unauthorized access from external
232 threats, including privileged software and other tenants.

233 This protection is particularly crucial for AI workloads, where both the data and the models are
234 valuable assets. AI models are significant intellectual property that may use highly sensitive
235 information, such as personal data or proprietary business insights. By isolating these
236 components in a hardware-protected environment, confidential computing is designed to
237 protect both the AI model and the data from unauthorized access or tampering during
238 processing.

239 Additionally, hardware-based confidential computing enhances trust in AI systems, especially in
240 cloud environments where resources are shared among multiple users. Ensuring that data and
241 models remain encrypted and secure even when processed is particularly important in
242 regulated industries like finance and healthcare, where compliance with data protection
243 regulations is critical. Overall, hardware-based confidential computing provides a powerful
244 solution for protecting AI data by enabling secure, efficient, and compliant AI processing in the
245 cloud.

246 Many cloud service providers offer VMs that have TEE capabilities enabled. In order to use
247 them, they must be launched through the corresponding CSP's portal/console, for which an

248 account and payment method is typically required. Once a TEE-capable VM has been
249 instantiated, certain steps are required to connect a relying party and attesor:

- 250 • Dockerfile to build the Intel Key Broker Service (KBS) relying party container
- 251 • **KBS.env** environment file to configure the KBS, modified to suit a particular
252 environment

253 The steps in a provisioning workflow are:

- 254 • Create an RSA key pair.
- 255 • Push the public key to the relying party.
- 256 • The replying party creates a wrapped key for encryption.
- 257 • Encrypt the ML model with the private key.

258

259 **4. Stage 1: Remote Attestation Service Configuration**

260 Remote attestation is a critical component of ensuring trust in confidential VMs in cloud
261 computing environments. It allows a remote party (e.g., user, trusted workload orchestrator) to
262 verify that a VM is running in a secure and expected environment and in a hardware and
263 software configuration that matches expected values. This is especially important in multi-
264 tenant cloud scenarios where sensitive workloads may be deployed on shared infrastructure.
265 Remote attestation provides cryptographic evidence that a VM is running on hardware
266 equipped with TEE capabilities and that the VM's memory and state are isolated and protected
267 from both other tenants and the hypervisor on a patched and attested platform.

268 The process of remote attestation typically begins with the generation of a hardware-backed
269 attestation report by the TEE. This report includes measurements of the VM's initial state (e.g.,
270 firmware, bootloader, operating system, application code) and a unique identifier for the
271 platform. These measurements are securely signed by the TEE's attestation key, which is rooted
272 in hardware and vouched for by the processor manufacturer. The VM or a trusted software
273 component inside it transmits this report to a remote verifier, which may be a tenant or an
274 external trust service.

275 The remote verifier then checks the attestation report's authenticity and integrity by validating
276 the cryptographic signature and comparing the measurements against known good or expected
277 values. These measurements can be forwarded to a relying party to evaluate whether the
278 requesting VM or service should be given access to secrets or keys.

279

280 **5. Stage 2: Relying Party Service Configuration**

281 Relying parties play a critical role in enabling TEEs for confidential computing in cloud
282 environments. TEEs are isolated regions of memory within a CPU that provide a secure enclave
283 for code and data to protect them from access by privileged system software, including the
284 hypervisor and OS. Relying parties (e.g., application owners, external verifiers) depend on these
285 assurances and interact with TEEs to establish trust before processing sensitive data. As
286 organizations increasingly migrate sensitive workloads to the cloud, concerns about data
287 confidentiality and integrity — especially from potentially untrusted infrastructure providers —
288 have driven the adoption of TEEs.

289 Relying parties often enable the use of TEEs through remote attestation, which is a
290 cryptographic protocol that allows the TEE to prove to a relying party that the code running
291 inside the secure enclave matches expected measurements. The relying party verifies this
292 attestation (e.g., using endorsements from hardware vendors) to ensure that the execution
293 environment meets predefined security policies. Only after successful attestation does the
294 relying party provision sensitive data or secrets into the enclave, so that the data is not exposed
295 outside of the secure boundary under normal operation.

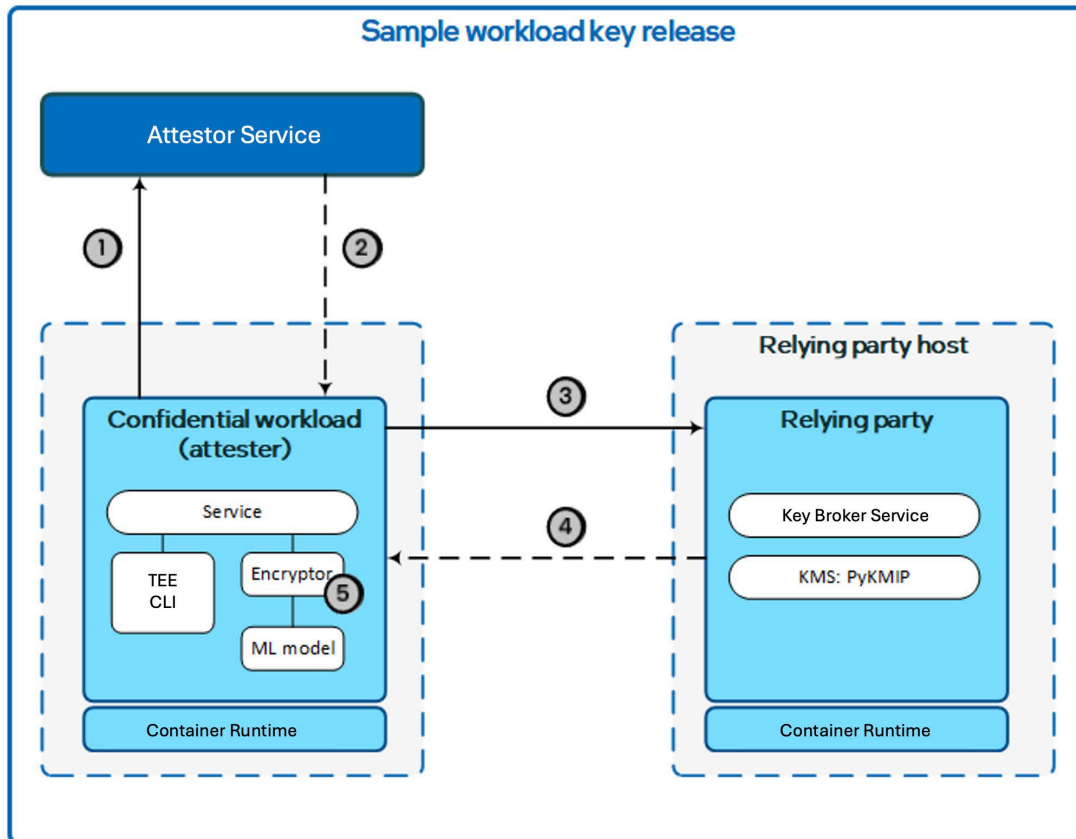
296 Relying parties also influence the life cycle of confidential workloads by defining policies for
297 enclave code versions, revocation statuses, and the integrity of data pipelines. For instance,
298 they may require specific measurements (i.e., cryptographic hashes) of the application code
299 before trusting the enclave. This level of control enables fine-grained governance and
300 accountability, particularly in regulated industries like healthcare and finance. Moreover, by
301 integrating with cloud-native tools and confidential computing platforms, relying parties can
302 automate these trust evaluations and maintain confidentiality without significantly impacting
303 performance or scalability.

304 Relying parties are foundational to the trust model of confidential computing in the cloud. Their
305 role in verifying and enforcing the security properties of TEEs provides meaningful assurance
306 that data confidentiality can be maintained during processing, even in untrusted or multi-
307 tenant cloud environments. As TEEs and attestation frameworks mature, the collaboration
308 between cloud providers and relying parties will become even more critical to securing the next
309 generation of cloud-native applications.

310

311 6. Stage 3: Execute the Key-Release Workflow

312 In a confidential computing model, the key-release workflow is central to enabling secure data
313 inference, particularly in scenarios where sensitive data must be processed within a TEE. This
314 workflow (see Fig. 1) is designed to release encryption keys that are used to decrypt model
315 inputs, parameters, or outputs only to verified and trusted execution environments, supporting
316 confidentiality and integrity.



317

318

Fig. 1. Sample workload key release

319 The process typically begins with remote attestation, where the TEE presents cryptographic
320 evidence (e.g., attestation report, quote) that describes its identity, the code it is running (e.g.,
321 a hash, measurement), and its security state. This evidence is sent to a key management or
322 attestation service, which compares the attestation data against predefined trust policies set by
323 the data owner or a relying party. These policies may require the enclave to run a specific,
324 approved version of the machine learning model; use up-to-date TEE firmware; and be
325 deployed on verified hardware.

326 If the attestation is successful and complies with policy, the key management service proceeds
327 with key release. This usually involves encrypting the decryption key in a way that only the TEE
328 can access, often leveraging hardware-protected key-wrapping or sealing techniques. The key is
329 then transmitted to the enclave, where it can be safely unwrapped and used to decrypt the
330 encrypted model or data inputs. The key is not exposed outside of the TEE in plaintext form

331 under normal operation, which protects it from potential compromise in multi-tenant cloud
332 environments when deployed on a patched, attested platform. Once the key is released and the
333 data is decrypted inside the enclave, the model can securely perform inference.

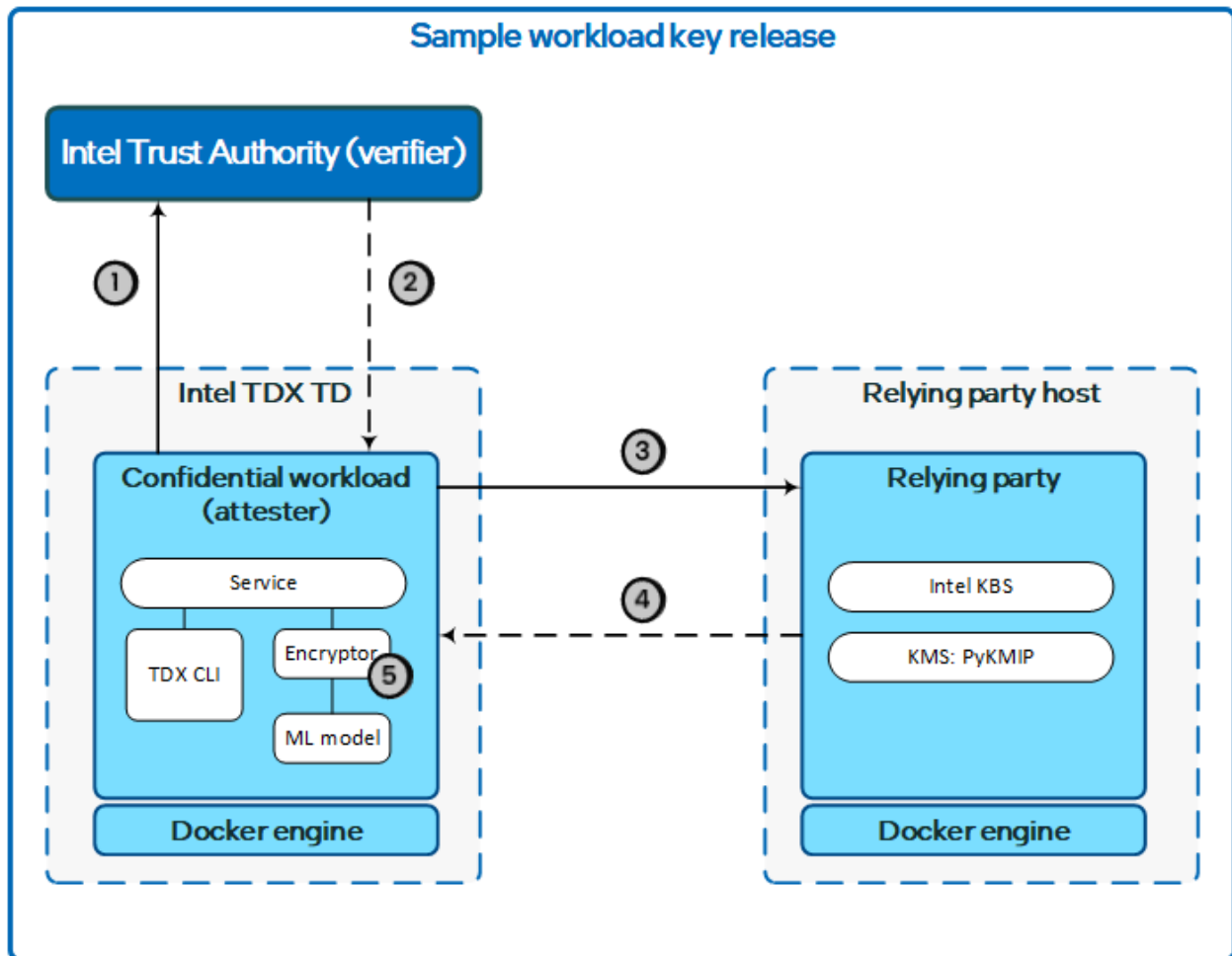
334 The output may also be encrypted before leaving the TEE, depending on the desired level of
335 protection. This end-to-end trust chain — from attestation through key release to data
336 processing — allows sensitive workloads (e.g., healthcare predictions, financial risk analysis,
337 personal data classification) to be run in the cloud while limiting the exposure of data to the
338 platform provider or other tenants. The key-release workflow is therefore a cornerstone of
339 secure, scalable confidential computing for AI inference.

340

341 **References**

- 342 [1] Bartock M, Souppaya M, Savino R, Knoll T, Shetty U, Cherfaoui M, Yeluri R, Malhotra A,
343 Banks D, Jordan M, Pendarakis D, Rao JR, Romness P, Scarfone KA (2022) Hardware-
344 Enabled Security: Enabling a Layered Approach to Platform Security for Cloud and Edge
345 Computing Use Cases. (National Institute of Standards and Technology, Gaithersburg,
346 MD), NIST Interagency or Internal Report (IR) NIST IR 8320.
347 <https://doi.org/10.6028/NIST.IR.8320>
- 348 [2] Polydys ML, Wisseman S (2009) Software Assurance in Acquisition: Mitigating Risks to
349 the Enterprise. Available at <https://apps.dtic.mil/sti/tr/pdf/ADA495389.pdf>
- 350 [3] Confidential Computing Consortium (2022) A Technical Analysis of Confidential
351 Computing. Available at [https://confidentialcomputing.io/wp-](https://confidentialcomputing.io/wp-content/uploads/sites/10/2023/03/CCC-A-Technical-Analysis-of-Confidential-Computing-v1.3_unlocked.pdf)
352 [content/uploads/sites/10/2023/03/CCC-A-Technical-Analysis-of-Confidential-](https://confidentialcomputing.io/wp-content/uploads/sites/10/2023/03/CCC-A-Technical-Analysis-of-Confidential-Computing-v1.3_unlocked.pdf)
353 [Computing-v1.3_unlocked.pdf](https://confidentialcomputing.io/wp-content/uploads/sites/10/2023/03/CCC-A-Technical-Analysis-of-Confidential-Computing-v1.3_unlocked.pdf)

354 **Appendix A. Solution Architecture**



355

356

Fig. 2. Cloud workload key release

357

358

359

1. The attesting workload obtains a quote from the Intel TDX TD by using the [Intel Trust Authority CLI for Intel TDX](#) ("TDX CLI"), which is also used to request an attestation [token](#) from the Intel Trust Authority.

360

361

362

363

364

- a. The workload service creates an RSA key pair. The private key is retained by the service, and the public key is base64-encoded and then copied to the user-data parameter in the token command. The user data is output in the **attester_held_data** claim in the attestation token. The public key is used by the KBS to encrypt (i.e., encapsulate) the key request response body.

365

2. The Intel Trust Authority evaluates the quote and returns an attestation token.

366

367

3. The workload forwards the attestation token to Intel KBS in a Representational State Transfer (REST) API request for the decryption key needed to decrypt the ML model.

368

369

370

4. The KBS evaluates the request and applies its key-release policy. An Intel KBS key-release policy compares values in the attestation token to reference values stored in the key-release policy. If the key-release policy evaluation is successful (i.e., all compared

371 values match), the KBS contacts the key management system (KMS) to obtain the
372 decryption key. The KMS will also apply its own key-release policy, which is independent
373 of the KBS policy. If the KMS policy allows, the KMS will return the requested key or
374 secret to the KBS. The KBS uses the public key passed in the attestation token to wrap
375 the KMS-supplied decryption key before sending it back to the workload.

376 a. Intel KBS creates a symmetric AES key to encrypt the KMS-supplied key in the
377 response body of a key request. This key is called the symmetric wrapping key
378 (SWK). For enhanced security, Intel KBS can also use asymmetric key
379 encapsulation with attestation tokens provided by Intel Trust Authority.
380 Asymmetric key encapsulation uses a public key provided by the attester to
381 encrypt the response body, including the SWK.

382 5. The workload service decrypts the Intel KBS response by using the RSA private key from
383 step 1.1. The service extracts the SWK from the response and uses it to decrypt the
384 KMS-supplied key, which can then be used to decrypt the ML model. The model is now
385 ready for use.

386 Appendix B. Trusted Execution Environment Implementation

387 B.1. Creating a VM With Intel TDX on Microsoft Azure

388 To create an Azure confidential VM with Intel TDX, create a VM with the following attributes:

- 389 • Security type: Trusted launch virtual machine
- 390 • Image: Ubuntu Server 22.04 LTS
- 391 • Size: DC2esv5

392 The following steps create an Azure VM with these attributes:

- 393 1. Sign in to Azure.
- 394 2. Select “Create a resource.”
- 395 3. Select “virtual machine.”

The screenshot shows the 'Instance details' section of the Azure portal. The form is filled with the following values:

- Virtual machine name ***: tdx-sample-vm
- Region ***: (US) East US 2
- Availability options**: Availability zone
- Availability zone ***: Zone 2
- Security type**: Confidential virtual machines
- Image ***: Ubuntu Server 22.04 LTS (Confidential VM) - x64 Gen2
- VM architecture**: x64 (selected)
- Run with Azure Spot discount**:
- Size ***: Standard_DC2es_v5 - 2 vcpus, 8 GiB memory (\$49.06/month)

Additional information visible in the form includes a note about selecting multiple zones and a warning that Arm64 is not supported with the selected image.

396

397

Fig. 3. Creating a cloud VM

398 4. Complete the following fields:

- 399 • Virtual machine name: Give your virtual machine a name.

- 400 • Region: Select **(US) West US 2**.
- 401 • Availability options: Select **Availability zone**.
- 402 • Availability zone: Select **zone 2**.¹
- 403 • Security type: Select **Trusted launch virtual machine**.
- 404 • Image: Select **Ubuntu Server 22.04 LTS - x64 Gen2**.
- 405 • Size: Select **DC2esv5**.²
- 406 • Public inbound: **None**.

Administrator account

Authentication type ⓘ SSH public key Password

i Azure now automatically generates an SSH key pair for you and allows you to store it for future use. It is a fast, simple, and secure way to connect to your virtual machine.

Username * ⓘ azureuser ✓

SSH public key source Generate new key pair ✓

Key pair name * tdx-sample-vm_key ✓

Inbound port rules

Select which virtual machine network ports are accessible from the public internet. You can specify more limited or granular network access on the Networking tab.

Public inbound ports * ⓘ None Allow selected ports

Select inbound ports Select one or more ports ✓

i All traffic from the internet will be blocked by default. You will be able to change inbound port rules in the VM > Networking page.

407

408

Fig. 4. Creating a cloud VM (continued)

409

- On the Networking page: For Public IP, select **None**.

410

5. Review the options and then create the VM. Deployment typically takes a few minutes.

¹ The availability of specific confidential virtual machine images and sizes in specific regions and availability zones is dynamic and may change. Check the Azure [Products by Region](#) page to find the regions and availability zones with available Confidential VM support.

² If not immediately listed, you may need to select “see all sizes” to choose **DC1s_v3** from the full list.

411 **B.2. Deploying Bastion**

- 412 1. Select the VM resource.
- 413 2. Select **Connect** and then select **Bastion**.
- 414 3. Select **Deploy Bastion**. Bastion deployment takes a few minutes.
- 415 4. Let Azure create a new SSH key pair. Download the private key to connect to the VM.
- 416 5. Connect using “SSH private Key from Local File.”
- 417 6. Enter the username.³ ()
- 418 7. Set the “local file” to the downloaded key <vm_name_key.pem>.

419 **B.3. Configuring TDX Prerequisites**

420 Verify that the /dev/tpmrm0 device exists.

```
421 ll /dev/tpmrm0  
422 crw-rw---- 1 tss tss 253, 65536 Jan 2 03:11 /dev/tpmrm0
```

423 This is the device used by Microsoft Azure to expose Intel TDX functionality. Other providers
424 and bare-metal hardware may use a different device name, usually /dev/tdx_guest.

425 Verify that the Intel TDX feature is active.

```
426 sudo dmesg | grep -i tdx  
427 [ 0.853076] Memory Encryption Features active: Intel TDX
```

428 In this example, a TEE is achieved using Intel TDX, as provided by Microsoft Azure’s Confidential
429 VMs.

430

³ The default username is “azureuser.”

431 **Appendix C. Key Management and Relying Party Implementation**

432 **C.1. Deploy the Key Broker Service**

- 433 • Prerequisites
- 434 • Makeself
- 435 • Go v1.21 or newer
- 436 • Create an attestation API key via the Intel Trust Authority web portal (requires
- 437 subscription).

438 Download the kbs.env file, and update the configuration.

```
439 wget https://raw.githubusercontent.com/intel/trustauthority-
```

441 When the relying party container is run, a startup script configures Intel KBS to use the
442 ADMIN_USERNAME and ADMIN_PASSWORD provided in the kbs.env file. An attestation API
443 key will also need to be provided for TRUSTAUTHORITY_API_KEY, which can be obtained from
444 the Intel Trust Authority portal. An HTTPS proxy may also need to be configured, depending on
445 the network configuration. In the following sample kbs.env file, replace the parameters marked
446 with angle brackets (e.g., <admin-username>).

```
447 LOG_LEVEL=INFO
448 KEY_MANAGER=kmip
449 ADMIN_USERNAME=<admin-username>
450 ADMIN_PASSWORD=<admin-password>
451 KMIP_VERSION=2.0
452 KMIP_SERVER_IP=127.0.0.1
453 KMIP_SERVER_PORT=5696
454 KMIP_CLIENT_KEY_PATH=/etc/pykmip/client_key.pem
455 KMIP_CLIENT_CERT_PATH=/etc/pykmip/client_certificate.pem
456 KMIP_ROOT_CERT_PATH=/etc/pykmip/ca_certificate.pem
457 PYKMIP_LOG_LEVEL=info
458 HTTPS_PROXY=
459 TRUSTAUTHORITY_BASE_URL=https://portal.trustauthority.intel.com
460 TRUSTAUTHORITY_API_URL=https://api.trustauthority.intel.com
461 TRUSTAUTHORITY_API_KEY=<API Key>
```

462

463 The kbs.env file contains secrets (e.g., Intel KBS administrator password, KMS password or
464 access token) and an Intel Trust Authority attestation API key. For a secure production system,
465 delete the kbs.env file after initial configuration to avoid compromising the configuration
466 secrets.

467 Run the following command to start the Intel KBS as a Docker container:

```
468     sudo docker run --name relying-party -d --env-file kbs.env --  
469     restart=always -p 9443:9443 relying-party:latest
```

470

471 **Appendix D. AI Model Encryption Implementation**

472 **D.1. Prepare the Sample AI Workload Image**

473 This section describes how to build and deploy the demonstration workload, which is an
474 application that runs a machine learning (ML) model in an Intel TDX trust domain. Since the ML
475 model is assumed to contain confidential data and parameters, it is encrypted when first
476 deployed. The demo workload (i.e., attester) must obtain an Intel Trust Authority attestation
477 token and include the token with a request to Intel KBS to retrieve the key needed to decrypt
478 the ML model. This workflow implements a secure key-release use case in passport attestation
479 mode.

480 Download the Dockerfile for building the Intel TDX demo workload.

```
481 wget https://raw.githubusercontent.com/intel/trustauthority-  
482 samples/main/deployment/sample-workload/Dockerfile
```

483 Run the following command to build the Docker image for the TDX demo workload:

```
484 sudo docker build --no-cache -t trustauthority-demo .
```

485 **D.2. Provision the Relying Party and Encrypt the AI Model**

486 Download the sample workload environment file.

```
487 wget https://raw.githubusercontent.com/intel/trustauthority-  
488 samples/main/deployment/sample-workload/workload.env
```

489 The workload.env file contains the settings needed for the sample workload to communicate
490 with Intel KBS and Intel Trust Authority. This file contains secrets. Embedding secrets in an
491 unencrypted file is not a recommended best practice for secure systems. It is only done here to
492 make the demo easy to configure.

493 The KBS_ADMIN and KBS_PASSWORD must match the ADMIN_USERNAME and
494 ADMIN_PASSWORD used in KBS.env.

495 The KBS_URL is of the form `https://<IP_address>:9443/kbs/v1`, where IP_address is the IP
496 address of the host machine on which the relying party (KBS) container is running. If the KBS is
497 running on the same TDVM as the workload, then provide 172.17.0.1 for the IP_address. Port
498 9443 is the default listening port for Intel KBS. If the port number for the Intel KBS container is
499 changed, a corresponding change will need to be made to KBS_URL.

500 Setting SKIP_TLS_VERIFICATION to true will skip the TLS server certificate verification. For this
501 demo, set it to true since a self-signed TLS certificates is being used for the KBS.

502 TRUSTAUTHORITY_API_KEY is the same as used in kbs.env.

```
503 KBS_ADMIN=<KBS-admin-username>  
504 KBS_PASSWORD=<KBS-admin-password>
```

```
505     KBS_URL= https://172.17.0.1:9443/kbs/v1
506     SKIP_TLS_VERIFICATION=true
507     PLATFORM=azure
508     # HTTPS_PROXY=<uncomment and enter proxy information if
509     necessary>
510     TRUSTAUTHORITY_API_URL=https://api.trustauthority.intel.com
511     TRUSTAUTHORITY_API_KEY=<API Key>
```

512 Run the following command to start the demo workload as a Docker container:

```
513     sudo docker run --name ita-demo -d --restart=always --env-file
514     workload.env --device=/dev/tpmrm0 -v /tmp:/tmp -p 12780:12780 --
515     group-add $(getent group tss | cut -d: -f3) trustauthority-demo:latest
```

516 The device option is used to set the Intel TDX device that the workload will use. The example
517 provided uses `/dev/tpmrm0`, which is the device used by Microsoft Azure Confidential VMs. If a
518 different Cloud provider or other hardware is used, set the device accordingly. Most non-Azure
519 providers use `/dev/tdx_guest`.

```
520     sudo docker run --name ita-demo -d --env-file workload.env --device=
521     /dev/tpmrm0 -p 12780:12780 --group-add $(getent group <user-group>
522     | cut -d: -f3) trustauthority-demo:latest
```

523 Where `<user-group>` is the group owning the `tdx_guest` or `tpmrm0` device:

```
524     crw-rw---- 1 <user> <user-group> 10, 123 /dev/tpmrm0
```

525 On a successful run, the container will execute the [encrypt-model.sh](#) script, which performs a
526 number of provisioning steps to finish configuring the demo environment:

- 527 • Authenticate to the KBS.
- 528 • Create a new Key Transfer Policy, which dictates when the KBS will release the
529 decryption key. In this demo, the policy requires Intel TDX and an updated attester (i.e.,
530 the demo workload running on the Azure Confidential VM) Trusted Compute base (TCB).
531 Because this is a “background check” model attestation flow, the relying party (i.e., the
532 KBS) is the enforcer of policy rather than the Intel Trust Authority (i.e., the verifier).
- 533 • Register the decryption key with the KBS and associate the wrapped key with the Key
534 Transfer Policy.
- 535 • Encrypt the `diabetes-linreg.model` file.

536 The key transfer policy is:

```
537     {
538     "attestation_type": "TDX",
539     "tdx": {
```

```
540     "attributes": {  
541         "enforce_tcb_upto_date": false  
542     }  
543 }  
544 }
```

545 These steps prepare the workload and KBS for the attestation demo. The AI model is now
546 encrypted, and the decryption key is held in wrapped form by the KBS. The KBS has a policy that
547 restricts the release of the wrapped decryption key to requests that contain an attestation that
548 Intel TDX is in use by the attester and that the attester's TCB is current.

549 The container then generates an answer file at /tmp/execute_workload_flow.env. This answer
550 file contains the URLs needed to execute the demo flow as well as the KBS key ID generated
551 during provisioning.

```
552     WORKLOAD_URL=https://127.0.0.1:12780  
553     KBS_URL=https://172.17.0.1:9443/kbs/v1  
554     KBS_KEY_ID=<key ID>
```

555 In a real-world production environment, a systems administrator would perform these
556 provisioning steps before uploading the workload with the encrypted AI model (e.g., to their
557 image registry, Cloud provider). For the purposes of this demo, the steps are scripted for
558 simplicity.

559

560 **Appendix E. Remote Attestation Implementation**

561 **E.1. Execute the Key-Release Workflow**

562 Download the execute_workload_flow.sh script to execute the secure key-release workflow:

```
563 wget https://raw.githubusercontent.com/intel/trustauthority-  
564 samples/main/deployment/sample-  
565 workload/execute\_workload\_flow.sh
```

566 Run the following command to execute the workflow script:

```
567 bash execute_workload_flow.sh
```

568 The script will:

- 569 • Request an attestation token for the workload from the Intel Trust Authority using
570 evidence collected from Intel TDX
- 571 • Request the model decryption key from the KBS using the attestation token
- 572 • Use the decryption key to decrypt the AI model
- 573 • Execute the AI model against sample patient data and output the results
- 574 • Clear the decrypted AI model from memory

```
575 ✓ curl is installed  
576 ✓ jq is installed  
577 =====  
578 =====Using /tmp/execute_workload_flow.env  
579 environment file  
580 Using Workload Host url : https://127.0.0.1:12780  
581 Using KBS host url : https://172.17.0.1:9443/kbs/v1  
582 Using KBS Key id : 2fa1511e-6606-42d9-adfd-854bb6b37f5a  
583 =====  
584 =====Get Attestation token from Intel Trust Authority:  
585 Attesting TDX workload with Intel Trust Authority by submitting  
586 TDquote.  
587 Url | https://127.0.0.1:12780/taa/v1/token  
588 Method | GET  
589 Expected Code | 200  
590 Received Code | 200  
591 Response | {"attestation_token": "<Base64-encoded attestation  
592 token>"}
```

593 Decoded token
594 {<Decoded attestation token>
595 -----
596 Get decryption key from KBS: Requesting decryption key from KBS by
597 sending attestation token obtained in previous step.
598
599 Url | <https://127.0.0.1:12780/taa/v1/key>
600 Method | POST
601 Expected Code | 200
602 Request | {
603 "attestation_token": "<Base64-encoded attestation token>",
604 "key_transfer_url": "[https://172.17.0.1:9443/kbs/v1/keys/2fa1511e-
605 6606-42d9-adfd-854bb6b37f5a/transfer](https://172.17.0.1:9443/kbs/v1/keys/2fa1511e-6606-42d9-adfd-854bb6b37f5a/transfer)"
606 }
607 Received Code | 200
608 Response | {"wrapped_key": "<wrapped decryption key>"}
609 -----
610 Decrypt model: Decrypting Model using decryption key received from
611 KBS.
612
613 Url | <https://127.0.0.1:12780/taa/v1/decrypt>
614 Method | POST
615 Expected Code | 204
616 Request | {"wrapped_key": "<wrapped key>", "wrapped_swk": "<key
617 encryption key>"}
618 Received Code | 204
619 -----
620 Execute model: Executing model with sample test data.
621
622 Url | <https://127.0.0.1:12780/taa/v1/execute>
623 Method | POST
624 Expected Code | 200

```
625     Request | {
626         "pregnancies": "1",
627         "blood-glucose": "128",
628         "blood-pressure": "88",
629         "skin-thickness": "39",
630         "insulin": "110",
631         "bmi": "36.5",
632         "dbf": "1.057",
633         "age": "37"
634     }
635     Received Code | 200
636     Response | {"high-risk":1}
637     -----
638     Reset Model: Clearing decrypted model from memory.
639
640     Url | https://127.0.0.1:12780/taa/v1/reset
641     Method | POST
642     Expected Code | 204
643     Received Code | 204
644     -----
645
646
```

647 **Appendix F. List of Symbols, Abbreviations, and Acronyms**

648 **API**
649 Application Programming Interface

650 **CPU**
651 Central Processing Unit

652 **CSP**
653 Cloud Service Provider

654 **ID**
655 Identifier

656 **KBS**
657 Key Broker Service

658 **KMS**
659 Key Management System

660 **REST**
661 Representational State Transfer

662 **SSH**
663 Secure Shell

664 **SWK**
665 Symmetric Wrapping Key

666 **TEE**
667 Trusted Execution Environment

668 **TLS**
669 Transport Layer Security

670 **TPM**
671 Trusted Platform Module

672 **VM**
673 Virtual Machine

674 **VPN**
675 Virtual Private Network

676