

The attached DRAFT document (provided here for historical purposes), released on January 27, 2016, has been superseded by the following publication:

Publication Number: **NIST Special Publication (SP) 800-90B**

Title: **Recommendation for the Entropy Sources Used for
Random Bit Generation**

Publication Date: **January 2018**

- Final Publication: <https://doi.org/10.6028/NIST.SP.800-90B> (which links to <http://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.800-90B.pdf>).
- Related Information on CSRC:
Final: <https://csrc.nist.gov/publications/detail/sp/800-90B/final>
Draft (attached): <https://csrc.nist.gov/publications/detail/sp/800-90b/archive/2016-01-27>
- Information on other NIST Computer Security Division publications and programs can be found at: <https://csrc.nist.gov>

2 **Recommendation for the Entropy**
3 **Sources Used for Random Bit**
4 **Generation**

7 Meltem Sönmez Turan
8 Elaine Barker
9 John Kelsey
10 Kerry A. McKay
11 Mary L. Baish
12 Mike Boyle
13

14
15
16
17 This publication is available free of charge from:
18 <http://dx.doi.org/10.6028/NIST.SP.XXX>
19

20
21 **C O M P U T E R S E C U R I T Y**
22

24 (Second DRAFT) NIST Special Publication 800-90B

25 **Recommendation for the Entropy**
26 **Sources Used for Random Bit**
27 **Generation**

28 Meltem Sönmez Turan

29 Elaine Barker

30 John Kelsey

31 Kerry McKay

32 *Computer Security Division*

33 *Information Technology Laboratory*

34 Mary L. Baish

35 Mike Boyle

36 *National Security Agency*

37 *Fort Meade, MD*

38 This publication is available free of charge from:

39 <http://dx.doi.org/10.6028/NIST.SP.XXX>

40
41
42
43
44
45 January 2016



47 U.S. Department of Commerce

48 *Penny Pritzker, Secretary*

49 National Institute of Standards and Technology

50 *Willie May, Under Secretary of Commerce for Standards and Technology and Director*

51
52
53
54

55

Authority

56 This publication has been developed by NIST in accordance with its statutory responsibilities under the
57 Federal Information Security Modernization Act (FISMA) of 2014, 44 U.S.C. § 3541 *et seq.*, Public Law
58 (P.L.) 113-283. NIST is responsible for developing information security standards and guidelines, including
59 minimum requirements for federal information systems, but such standards and guidelines shall not apply
60 to national security systems without the express approval of appropriate federal officials exercising policy
61 authority over such systems. This guideline is consistent with the requirements of the Office of Management
62 and Budget (OMB) Circular A-130.

63 Nothing in this publication should be taken to contradict the standards and guidelines made mandatory and
64 binding on federal agencies by the Secretary of Commerce under statutory authority. Nor should these
65 guidelines be interpreted as altering or superseding the existing authorities of the Secretary of Commerce,
66 Director of the OMB, or any other federal official. This publication may be used by nongovernmental
67 organizations on a voluntary basis and is not subject to copyright in the United States. Attribution would,
68 however, be appreciated by NIST.

69 National Institute of Standards and Technology Special Publication 800-90B
70 Natl. Inst. Stand. Technol. Spec. Publ. 800-90B, 66 pages (January 2016)
71 CODEN: NSPUE2

72 This publication is available free of charge from:
73 <http://dx.doi.org/10.6028/NIST.SP.XXX>

74 Certain commercial entities, equipment, or materials may be identified in this document in order to describe an
75 experimental procedure or concept adequately. Such identification is not intended to imply recommendation or
76 endorsement by NIST, nor is it intended to imply that the entities, materials, or equipment are necessarily the best
77 available for the purpose.

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

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

87 **Public comment period: *January 25, 2016 through May 9, 2016***

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

89 National Institute of Standards and Technology
90 Attn: Computer Security Division, Information Technology Laboratory
91 100 Bureau Drive (Mail Stop 8930) Gaithersburg, MD 20899-8930
92 Email: rbg_comments@nist.gov

93

94

Reports on Computer Systems Technology

95 The Information Technology Laboratory (ITL) at the National Institute of Standards and
96 Technology (NIST) promotes the U.S. economy and public welfare by providing technical
97 leadership for the Nation’s measurement and standards infrastructure. ITL develops tests, test
98 methods, reference data, proof of concept implementations, and technical analyses to advance the
99 development and productive use of information technology. ITL’s responsibilities include the
100 development of management, administrative, technical, and physical standards and guidelines for
101 the cost-effective security and privacy of other than national security-related information in
102 Federal information systems. The Special Publication 800-series reports on ITL’s research,
103 guidelines, and outreach efforts in information system security, and its collaborative activities with
104 industry, government, and academic organizations.

105

Abstract

106 This Recommendation specifies the design principles and requirements for the entropy sources
107 used by Random Bit Generators, and the tests for the validation of entropy sources. These entropy
108 sources are intended to be combined with Deterministic Random Bit Generator mechanisms that
109 are specified in SP 800-90A to construct Random Bit Generators, as specified in SP 800-90C.

110

Keywords

111 Conditioning functions; Entropy source; health testing; IID testing; min-entropy; noise source;
112 predictors; random number generators

113

Acknowledgements

114 The authors of this Recommendation gratefully acknowledge and appreciate contributions by their
115 colleagues at NIST, Apostol Vassilev and Timothy A. Hall, and also Aaron H. Kaufman and Darryl
116 M. Buller of the National Security Agency, for assistance in the development of this
117 Recommendation. NIST also thanks the many contributions by the public and private sectors.

118

Conformance Testing

119 Conformance testing for implementations of this Recommendation will be conducted within the
120 framework of the Cryptographic Algorithm Validation Program (CAVP) and the Cryptographic
121 Module Validation Program (CMVP). The requirements of this Recommendation are indicated by
122 the word “**shall**.” Some of these requirements may be out-of-scope for CAVP or CMVP validation
123 testing, and thus are the responsibility of entities using, implementing, installing or configuring
124 applications that incorporate this Recommendation.

125

126

127

Note to Reviewers

128 To facilitate public review, we have compiled a number of open issues for which we would like
129 reviewer input. Please keep in mind that it is not necessary to respond to all questions listed below,
130 nor is review limited to these issues. Reviewers should also feel free to suggest other areas of
131 revision or enhancement to the document as they see fit.

132

133 - *Post-processing functions (Section 3.2.2)*: We provided a list of approved post-processing
134 functions. Is the selection of the functions appropriate?

135 - *Entropy assessment (Section 3.1.5)*: While estimating the entropy for entropy sources using
136 a conditioning component, the values of n and q are multiplied by the constant 0.85. Is the
137 selection of this constant reasonable?

138 - *Multiple noise sources*: The Recommendation only allows using multiple noise sources if
139 the noise sources are independent. Should the use of dependent noise sources also be
140 allowed, and how can we calculate an entropy assessment in this case?

141 - *Health Tests*: What actions should be taken when health tests raise an alarm? The minimum
142 allowed value of a type I error for health testing is selected as 2^{-50} . Is this selection
143 reasonable?

144

145 **Table of Contents**

146 **1 Introduction 1**

147 1.1 Scope..... 1

148 1.2 Symbols..... 1

149 1.3 Organization 2

150 **2 General Discussion..... 3**

151 2.1 Min-Entropy 3

152 2.2 The Entropy Source Model 3

153 2.2.1 Noise Source..... 4

154 2.2.2 Conditioning Component..... 5

155 2.2.3 Health Tests 5

156 2.3 Conceptual Interfaces 5

157 2.3.1 GetEntropy: An Interface to the Entropy Source 5

158 2.3.2 GetNoise: An Interface to the Noise Source..... 6

159 2.3.3 HealthTest: An Interface to the Entropy Source 6

160 **3 Entropy Source Validation 7**

161 3.1 Validation Process 7

162 3.1.1 Data Collection 9

163 3.1.2 Determining the track: IID track vs. non-IID track..... 9

164 3.1.3 Initial Entropy Estimate..... 10

165 3.1.4 Restart Tests 10

166 3.1.5 Entropy Estimation for Entropy Sources Using a Conditioning

167 Component..... 12

168 3.1.6 Using Multiple Noise Sources..... 14

169 3.2 Requirements for Validation Testing 15

170 3.2.1 Requirements on the Entropy Source..... 15

171 3.2.2 Requirements on the Noise Source..... 16

172 3.2.3 Requirements on the Conditioning Component..... 17

173 3.2.4 Requirements on Data Collection 17

174 **4 Health Tests..... 18**

175 4.1 Health Test Overview..... 18

176 4.2 Types of Health Tests 18

177	4.3	Requirements for Health Tests	19
178	4.4	Approved Continuous Health Tests	20
179	4.4.1	Repetition Count Test.....	21
180	4.4.2	Adaptive Proportion Test.....	22
181	4.5	Vendor-Defined Alternatives to the Continuous Health Tests	24
182	4.6	Alternative Health Test Criteria	24
183	5	Testing the IID Assumption.....	24
184	5.1	Permutation Testing.....	25
185	5.1.1	Excursion Test Statistic	26
186	5.1.2	Number of Directional Runs	27
187	5.1.3	Length of Directional Runs	27
188	5.1.4	Number of Increases and Decreases	28
189	5.1.5	Number of Runs Based on the Median.....	28
190	5.1.6	Length of Runs Based on Median	29
191	5.1.7	Average Collision Test Statistic	29
192	5.1.8	Maximum Collision Test Statistic.....	30
193	5.1.9	Periodicity Test Statistic	30
194	5.1.10	Covariance Test Statistic.....	30
195	5.1.11	Compression Test Statistics	31
196	5.2	Additional Chi-square Statistical Tests.....	31
197	5.2.1	Testing Independence for Non-Binary Data	31
198	5.2.2	Testing Goodness-of-fit for non-binary data	32
199	5.2.3	Testing Independence for Binary Data	33
200	5.2.4	Testing Goodness-of-fit for Binary Data	34
201	5.2.5	Length of the Longest Repeated Substring Test	35
202	6	Estimating Min-Entropy.....	35
203	6.1	IID Track: Entropy Estimation for IID Data	35
204	6.2	Non-IID Track: Entropy Estimation for Non-IID Data.....	36
205	6.3	Estimators.....	36
206	6.3.1	The Most Common Value Estimate	36
207	6.3.2	The Collision Estimate.....	37
208	6.3.3	The Markov Estimate.....	38
209	6.3.4	The Compression Estimate	41

210	6.3.5 t -Tuple Estimate	43
211	6.3.6 Longest Repeated Substring (LRS) Estimate	43
212	6.3.7 Multi Most Common in Window Prediction Estimate	44
213	6.3.8 The Lag Prediction Estimate	46
214	6.3.9 The MultiMMC Prediction Estimate	47
215	6.3.10 The LZ78Y Prediction Estimate	50
216	6.4 Reducing the Sample Space	52

217
218

List of Appendices

219	Appendix A— Acronyms	54
220	Appendix B— Glossary	55
221	Appendix C— References	60
222	Appendix D— Min-Entropy and Optimum Guessing Attack Cost.....	61
223	Appendix E— Post-processing Functions.....	63
224	Appendix F— The Narrowest Internal Width	64
225	Appendix G— CBC-MAC Specification	64
226	Appendix H— Different Strategies for Entropy Estimation	65
227	H.1 Entropic Statistics	65
228	H.1.1 Approximation of $F(1/z)$	66
229	H.2 Predictors.....	66

230
231

232

List of Figures

233 Figure 1 Entropy Source Model..... 4
234 Figure 2 Entropy Estimation Strategy 8
235 Figure 3 Entropy of the Conditioning Component 12
236 Figure 4 Generic Structure of Permutation Testing 25
237 Figure 5 Pseudo-code of the Fisher-Yates Shuffle..... 26

238

239

240

List of Tables

241 Table 1 The narrowest internal width and output lengths of the vetted conditioning
242 functions..... 13
243 Table 2 Adaptive proportion test on binary data for various entropy/sample levels with
244 $W=1024$ 23
245 Table 3 Adaptive proportion test on non-binary data for various entropy/sample levels
246 with $W=512$ 23

247

248 **1 Introduction**249 **1.1 Scope**

250 Cryptography and security applications make extensive use of random numbers and random bits.
251 However, the generation of random bits is problematic in many practical applications of
252 cryptography. The NIST Special Publication (SP) 800-90 series of Recommendations provides
253 guidance on the construction and validation of Random Bit Generators (RBGs) in the form of
254 Deterministic Random Bit Generators (DRBGs) or Non-deterministic Random Bit Generators
255 (NRBGs) that can be used for cryptographic applications. This Recommendation specifies how to
256 design and test entropy sources that can be used by these RBGs. SP 800-90A addresses the
257 construction of **approved** Deterministic Random Bit Generator (DRBG) mechanisms, while SP
258 800-90C addresses the construction of RBGs from the mechanisms in SP 800-90A and the entropy
259 sources in SP 800-90B. These Recommendations provide a basis for validation by NIST's
260 Cryptographic Algorithm Validation Program (CAVP) and Cryptographic Module Validation
261 Program (CMVP).

262
263 An entropy source that conforms to this Recommendation can be used by RBGs to produce a
264 sequence of random bits. While there has been extensive research on the subject of generating
265 pseudorandom bits using a DRBG and an unknown seed value, creating such an unknown seed
266 has not been as well documented. The only way for this seed value to provide real security is for
267 it to contain a sufficient amount of randomness, e.g., from a non-deterministic process referred to
268 as an entropy source. This Recommendation describes the properties that an entropy source must
269 have to make it suitable for use by cryptographic random bit generators, as well as the tests used
270 to validate the quality of the entropy source.

271
272 The development of entropy sources that construct unpredictable outputs is difficult, and providing
273 guidance for their design and validation testing is even more so. The testing approach defined in
274 this Recommendation assumes that the developer understands the behavior of the noise source
275 within the entropy source and has made a good-faith effort to produce a consistent source of
276 entropy. It is expected that, over time, improvements to the guidance and testing will be made,
277 based on experience in using and validating against this Recommendation.

278
279 This Recommendation was developed in concert with American National Standard (ANS) X9.82,
280 a multi-part standard on random number generation.

281 **1.2 Symbols**

282 The following symbols and functions are used in this Recommendation.

$A = \{x_1, x_2, \dots, x_k\}$ The set of all possible distinct sample outputs from a noise source, i.e. the
alphabet.

H The min-entropy of the samples from a (digitized) noise source or of the
output from an entropy source; the min-entropy assessment for a noise
source or entropy source.

H_I	Initial entropy estimate.
$\log_b(x)$	The logarithm of x with respect to base b .
$\max(a, b)$	A function that returns the maximum of the two values a and b .
k	The number of possible sample values, i.e., the size of the alphabet.
α	The probability of falsely rejecting the null hypothesis (type I error).
$ a $	A function that returns the absolute value of a .
p_i	The probability for an observation (or occurrence) of the sample value x_i in A .
p_{max}	The probability of observing the most common sample from a noise source.
$S=(s_1, \dots, s_L)$	A dataset that consists of an ordered collection of samples, where $s_i \in A$.
x_i	A possible output from the (digitized) noise source.
$[a,b]$	The interval of numbers between a and b , including a and b .
$\lceil x \rceil$	A function that returns the smallest integer greater than or equal to x ; also known as the <i>ceiling</i> function.
$\lfloor x \rfloor$	A function that returns the largest integer less than or equal to x ; also known as the <i>floor</i> function.
\parallel	Concatenation.
\oplus	Bit-wise exclusive-or operation.

283 1.3 Organization

284 Section 2 gives a general discussion on min-entropy, the entropy source model and the conceptual
 285 interfaces. Section 3 explains the validation process and lists the requirements on the entropy
 286 source, data collection, documentation, etc. Section 4 describes the health tests. Section 5 includes
 287 various statistical tests to check whether the entropy source outputs are IID (independent and
 288 identically distributed) or not. Section 6 provides several methods to estimate the entropy of the
 289 noise source. The appendices include a list of acronyms, a glossary, references, a discussion on
 290 min-entropy and the optimum guessing attack cost, descriptions of the post-processing functions,
 291 information about the narrowest internal width and the underlying information on different entropy
 292 estimation strategies used in this Recommendation.

293 **2 General Discussion**

294 The three main components of a cryptographic RBG are a source of random bits (an entropy
295 source), an algorithm for accumulating and providing random bits to the consuming applications,
296 and a way to combine the first two components appropriately for the cryptographic applications.
297 This Recommendation describes how to design and test entropy sources. SP 800-90A describes
298 deterministic algorithms that take an entropy input and use it to produce pseudorandom values. SP
299 800-90C provides the “glue” for putting the entropy source together with the algorithm to
300 implement an RBG.

301 Specifying an entropy source is a complicated matter. This is partly due to confusion in the
302 meaning of entropy, and partly due to the fact that, while other parts of an RBG design are strictly
303 algorithmic, entropy sources depend on physical processes that may vary from one instance of a
304 source to another. This section discusses, in detail, both the entropy source model and the meaning
305 of entropy.

306 **2.1 Min-Entropy**

307 The central mathematical concept underlying this Recommendation is *entropy*. Entropy is defined
308 relative to one’s knowledge of an experiment’s output prior to observation, and reflects the
309 uncertainty associated with predicting its value – the larger the amount of the entropy, the greater
310 the uncertainty in predicting the value of an observation. There are many possible types of entropy;
311 this Recommendation uses a very conservative measure known as *min-entropy*, which measures
312 the difficulty of guessing the most likely output of the entropy source.

313 In cryptography, the unpredictability of secret values (such as cryptographic keys) is essential. The
314 probability that a secret is guessed correctly in the first trial is related to the min-entropy of the
315 distribution that the secret was generated from. The min-entropy is closely related to the negative
316 logarithm of the maximum probability using the optimal guessing strategy [Cac97] (see Appendix
317 D for more information).

318 The min-entropy of an independent discrete random variable X that takes values from the set
319 $A = \{x_1, x_2, \dots, x_k\}$ with probability $\Pr(X=x_i) = p_i$ for $i = 1, \dots, k$ is defined as

$$320 \quad H = - \min_{1 \leq i \leq k} (-\log_2 p_i),$$

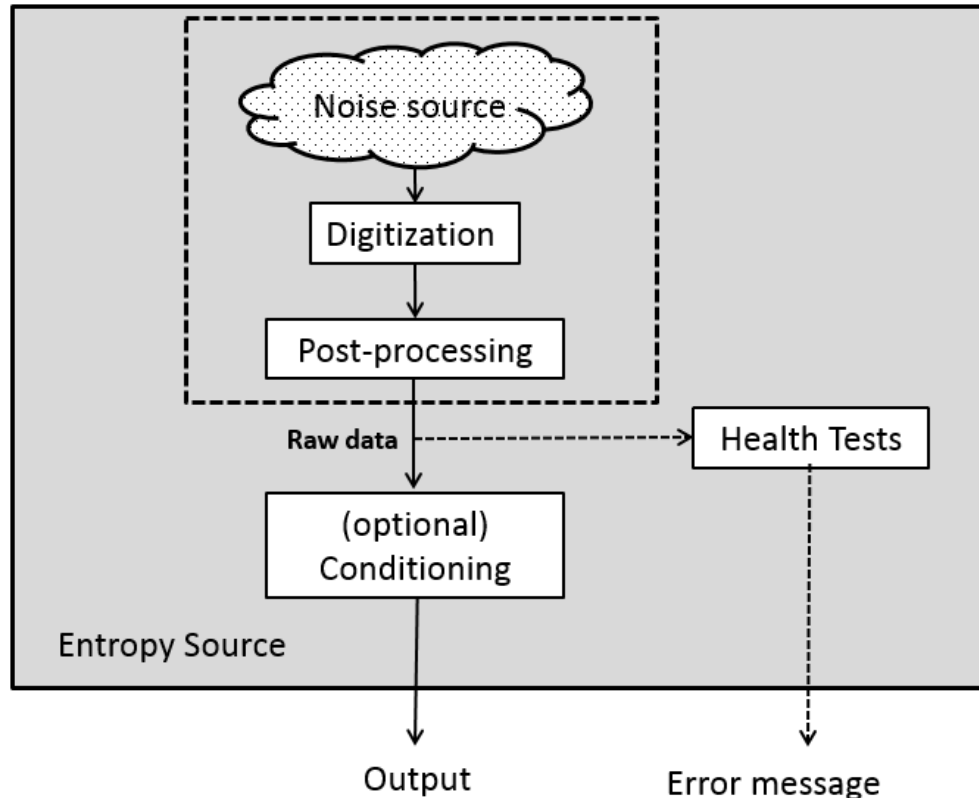
$$321 \quad = - \log_2 \max_{1 \leq i \leq k} p_i.$$

322 If X has min-entropy H , then the probability of observing any particular value for X is no greater
323 than 2^{-H} . The maximum possible value for the min-entropy of a random variable with k distinct
324 values is $\log_2 k$, which is attained when the random variable has a uniform probability distribution,
325 i.e., $p_1 = p_2 = \dots = p_k = 1/k$.

326 **2.2 The Entropy Source Model**

327 This section describes the entropy source model in detail. Figure 1 illustrates the model that this
328 Recommendation uses to describe an entropy source and its components: a noise source, an

329 optional conditioning component and a health testing component.



330

331

Figure 1 Entropy Source Model

332 2.2.1 Noise Source

333 The noise source is the root of security for the entropy source and for the RBG as a whole. This is
 334 the component that contains the non-deterministic, entropy-providing activity that is ultimately
 335 responsible for the uncertainty associated with the bitstrings output by the entropy source.

336 If the non-deterministic activity being sampled produces something other than binary data, the
 337 sampling process includes a *digitization* process that converts the output samples to bits. The noise
 338 source may also include some simple *post-processing operations* that can reduce the statistical
 339 biases of the samples and increase the entropy rate of the resulting output. The output of the
 340 digitized and optionally post-processed noise source is called the *raw data*.

341 This Recommendation assumes that the sample values obtained from a noise source consist of
 342 fixed-length bitstrings.

343 If the noise source fails to generate random outputs, no other component in the RBG can
 344 compensate for the lack of entropy; hence, no security guarantees can be made for the application
 345 relying on the RBG.

346 In situations where a single noise source does not provide sufficient entropy in a reasonable amount
 347 of time, outputs from multiple noise sources may be combined to obtain the necessary amount of

348 entropy. When multiple noise sources are used, the relationship between sources affects the
349 entropy of the outputs. If the noise sources are independent, their entropy assessments can be
350 added. Thermal noise and mouse movements can be given as examples of independent noise
351 sources (i.e., the output of the noise sources are independent). However, for some combinations of
352 noise sources, such as the ones based on dependent processes (e.g., packet arrival times in a
353 communication network and hard drive access times), the total entropy produced is harder to
354 estimate. This Recommendation only considers the use of independent noise sources.

355 2.2.2 Conditioning Component

356 The optional conditioning component is a deterministic function responsible for reducing bias
357 and/or increasing the entropy rate of the resulting output bits (if necessary to obtain a target value).
358 There are various methods for achieving this. The developer should consider the conditioning
359 component to be used and how variations in the behavior of the noise source may affect the entropy
360 rate of the output. In choosing an approach to implement, the developer may either choose to
361 implement a cryptographic algorithm listed in Section 3.1.5.1.1 or use an alternative algorithm as
362 a conditioning component. The use of either of these approaches is permitted by this
363 Recommendation.

364 2.2.3 Health Tests

365 Health tests are an integral part of the entropy source design that are intended to ensure that the
366 noise source and the entire entropy source continue to operate as expected. When testing the
367 entropy source, the end goal is to obtain assurance that failures of the entropy source are caught
368 quickly and with a high probability. Another aspect of health testing strategy is determining likely
369 failure modes for the entropy source and, in particular, for the noise source. Health tests are
370 expected to include tests that can detect these failure conditions.

371 The health tests can be separated into three categories: *start-up tests* (on all components),
372 *continuous tests* (primarily on the noise source), and *on-demand tests*.

373 2.3 Conceptual Interfaces

374 This section describes three conceptual interfaces that can be used to interact with the entropy
375 source: **GetEntropy**, **GetNoise** and **HealthTest**. However, it is anticipated that the actual
376 interfaces used may depend on the entropy source employed.

377 These interfaces can be used by a developer when constructing an RBG as specified in SP 800-
378 90C.

379 2.3.1 GetEntropy: An Interface to the Entropy Source

380 The **GetEntropy** interface can be considered to be a command interface into the outer entropy
381 source box in Figure 1. This interface is meant to indicate the types of requests for services that an
382 entropy source may support.

383 A **GetEntropy** call could return a bitstring containing the requested amount of entropy, along
384 with an indication of the status of the request. Optionally, an assessment of the entropy can be

385 provided.

386

GetEntropy

Input:

bits_of_entropy: the requested amount of entropy

Output:

entropy_bitstring: The string that provides the requested entropy.

status: A Boolean value that is TRUE if the request has been satisfied, and is FALSE otherwise.

387

2.3.2 GetNoise: An Interface to the Noise Source

389 The **GetNoise** interface can be considered to be a command interface into the noise source
390 component of an entropy source. This could be used to obtain raw, digitized and optionally post-
391 processed outputs from the noise source for use in validation testing or for external health tests.
392 While it is not required to be in this form, it is expected that an interface be available that allows
393 noise source data to be obtained without harm to the entropy source. This interface is meant to
394 provide test data to credit a noise source with an entropy estimate during validation or for health
395 testing. It is permitted that such an interface is available only in “test mode” and that it is disabled
396 when the source is operational.

397 This interface is not intended to constrain real-world implementations, but to provide a consistent
398 notation to describe data collection from noise sources.

399 A **GetNoise** call returns raw, digitized, samples from the noise source, along with an indication of
400 the status of the request.

GetNoise

Input:

number_of_samples_requested: An integer value that indicates the requested number of samples to be returned from the noise source.

Output:

noise_source_data: The sequence of samples from the noise source with length *number_of_samples_requested*.

status: A Boolean value that is TRUE if the request has been satisfied, and is FALSE otherwise.

401

2.3.3 HealthTest: An Interface to the Entropy Source

403 A **HealthTest** call is a request to the entropy source to conduct a test of its health. Note that it may
404 not be necessary to include a separate **HealthTest** interface if the execution of the tests can be
405 initiated in another manner that is acceptable to FIPS 140 [FIPS140] validation.

406

HealthTest

Input:

type_of_test_requested: A bitstring that indicates the type or suite of tests to be performed (this may vary from one entropy source to another).

Output:

status: A Boolean value that is TRUE if the entropy source passed the requested test, and is FALSE otherwise.

407

408 **3 Entropy Source Validation**

409 Entropy source validation is necessary in order to obtain assurance that all relevant requirements
410 of this Recommendation are met. This Recommendation provides requirements and guidance that
411 will allow an entropy source to be validated for an entropy assessment that will provide evidence
412 that the entropy source produces bitstrings providing entropy at a specified rate. Validation
413 consists of testing by an NVLAP-accredited laboratory against the requirements of SP 800-90B,
414 followed by a review of the results by NIST's CAVP and CMVP. Validation provides additional
415 assurance that adequate entropy is provided by the source and may be necessary to satisfy some
416 legal restrictions, policies, and/or directives of various organizations.

417 The validation of an entropy source presents many challenges. No other part of an RBG is so
418 dependent on the technological and environmental details of an implementation. At the same time,
419 the proper operation of the entropy source is essential to the security of an RBG. The developer
420 should make every effort to design an entropy source that can be shown to serve as a consistent
421 source of entropy, producing bitstrings that can provide entropy at a rate that meets (or exceeds) a
422 specified value. In order to design an entropy source that provides an adequate amount of entropy
423 per output bitstring, the developer must be able to accurately estimate the amount of entropy that
424 can be provided by sampling its (digitized) noise source. The developer must also understand the
425 behavior of the other components included in the entropy source, since the interactions between
426 the various components may affect any assessment of the entropy that can be provided by an
427 implementation of the design. For example, if it is known that the raw noise-source output is
428 biased, appropriate conditioning components can be included in the design to reduce that bias to a
429 tolerable level before any bits are output from the entropy source.

430 **3.1 Validation Process**

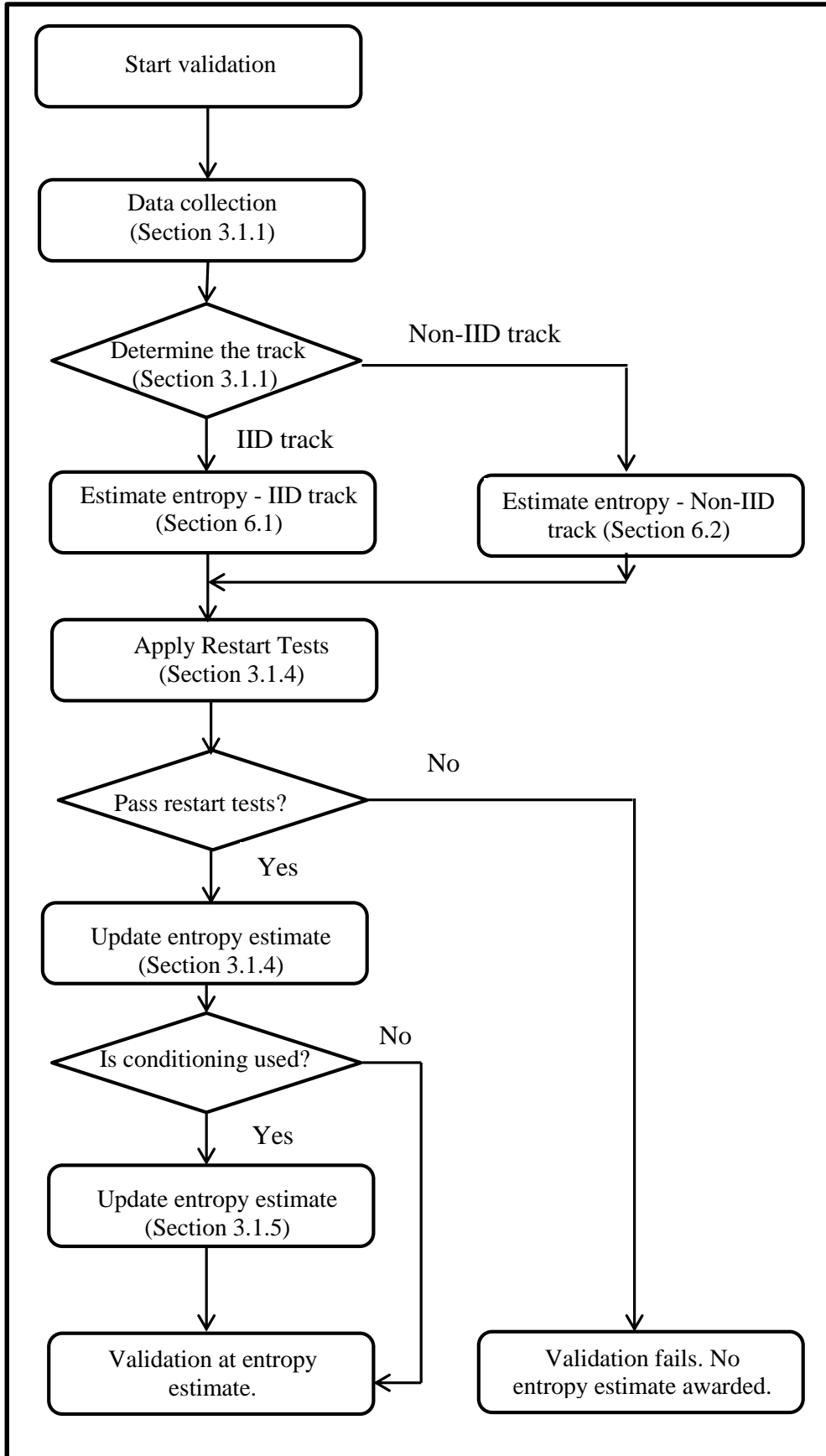
431 An entropy source may be submitted to an accredited lab for validation testing by the developer or
432 any entity with an interest in having an entropy source validated. After the entropy source is
433 submitted for validation, the labs will examine all documentation and theoretical justifications
434 submitted. The labs will evaluate these claims, and may ask for more evidence or clarification.

435 The general flow of entropy source validation testing is summarized in Figure 2. The following
436 sections describe the details of the validation testing process.

437

438

439



440

Figure 2 Entropy Estimation Strategy

441 **3.1.1 Data Collection**

442 The submitter provides the following inputs for entropy estimation, according to the requirements
443 presented in Section 3.2.4.

- 444 1. A *sequential* dataset of at least 1,000,000 consecutive sample values obtained directly from the
445 noise source (i.e., raw samples) **shall** be collected for validation¹. If the generation of 1,000,000
446 consecutive samples is not possible, the concatenation of several smaller sets of consecutive
447 samples (generated using the same device) is allowed. Smaller sets **shall** contain at least 1,000
448 samples. The concatenated dataset **shall** contain at least 1,000,000 samples. If multiple noise
449 sources are used, a dataset of at least 1,000,000 samples from each noise source **shall** be
450 collected.
- 451 2. If the entropy source includes a conditioning component that is not listed in Section 3.1.5.1.1,
452 a *conditioned sequential* dataset of at least 1,000,000 consecutive samples values obtained as
453 the output of the conditioning component **shall** be collected for validation. The output of the
454 conditioning component **shall** be treated as a binary string for testing purposes. Note that the
455 data collected from the noise source for validation may be used as input to the conditioning
456 component for the collection of conditioned output values.
- 457 3. For the restart tests (see Section 3.1.4), the entropy source must be restarted 1000 times; for
458 each restart, 1000 consecutive samples **shall** be collected directly from the noise source. This
459 data is stored in a 1000x1000 restart matrix M , where $M[i][j]$ represents the j th sample from
460 the i th restart.
- 461 4. If multiple noise sources are used, sequential and restart datasets from each noise source **shall**
462 be collected, as specified in item 1. If a conditioning component that is not listed in Section
463 3.1.5.1.1 is used, a single conditioned dataset **shall** be collected as an output of the entropy
464 source.

465 **3.1.2 Determining the track: IID track vs. non-IID track**

466 According to this Recommendation, entropy estimation is done using two different tracks: an IID-
467 track and a non-IID track. The *IID-track* (see Section 6.1) is used for entropy sources that generate
468 IID (independent and identically distributed) samples, whereas the *non-IID track* (see Section 6.2)
469 is used for noise sources that do not generate IID samples.

470 The track selection is done based on the following rules. The IID track **shall** be chosen only when
471 *all* of the following conditions are satisfied:

- 472 1. The submitter makes an IID claim on the noise source, based on his analysis of the design.
473 The submitter **shall** provide rationale for the IID claim.

¹ Providing additional data beyond what is required will result in more accurate entropy estimates. Lack of sufficient data may result in lower entropy estimates due to the necessity of mapping down the output values (see Section 6.4). It is recommended that, if possible, more data than is required be collected for validation. However, it is assumed in subsequent text that only the required data has been collected.

- 474 2. The sequential dataset described in item 1 of Section 3.1.1 is tested using the statistical
475 tests described in Section 5 to verify the IID assumption, and the IID assumption is verified
476 (i.e., there is no evidence that data is not IID).
- 477 3. The row and the column datasets described in item 3 of Section 3.1.1 are tested using the
478 statistical tests described in Section 5 to verify the IID assumption, and the IID assumption
479 is verified.
- 480 4. If a conditioning component that not listed in Section 3.1.5.1.1 is used, the conditioned
481 sequential dataset is tested using the statistical tests described in Section 5 to verify the IID
482 assumption, and the IID assumption is verified.

483 If any of these conditions are not met, the estimation process **shall** follow the non-IID track.

484 3.1.3 Initial Entropy Estimate

485 After determining the entropy estimation track, a min-entropy estimate per sample, denoted as
486 $H_{original}$, for the *sequential dataset* is calculated using the methods described in Section 6.1 (for the
487 IID track) or Section 6.2 (for the non-IID track). If the size of the sample space is greater than 256,
488 it **shall** be reduced to at most 256, using the method described in Section 6.4.

489 If the sequential dataset is not binary (i.e., the size of the sample space k is more than 2), an
490 additional entropy estimation (per bit), denoted $H_{bitstring}$, is determined (based on the entropy
491 estimation track, as specified in the previous paragraph), considering the sequential dataset as a
492 bitstring. The bits after the first 1,000,000 bits may be ignored. The entropy per sample is estimated
493 to be $n \times H_{bitstring}$ where n is the size of the fixed-length samples.

494 The submitter **shall** provide an entropy estimate for the noise source, which is based on the
495 submitter's analysis of the noise source (see Requirement 8 in Section 3.2.2). This estimate is
496 denoted as $H_{submitter}$.

497 The initial entropy estimate of the noise source is calculated as $H_I = \min(H_{original}, n \times H_{bitstring},$
498 $H_{submitter})$ for non-binary sources and as $H_I = \min(H_{original}, H_{submitter})$ for binary sources.

499 3.1.4 Restart Tests

500 The entropy estimate of a noise source, calculated from a single, long-output sequence, might
501 provide an overestimate if the noise source generates correlated sequences after restarts. Hence,
502 an attacker with access to multiple noise source output sequences after restarts may be able to
503 predict the next output sequence with much better success than the entropy estimate suggests. The
504 restart tests described in this section re-evaluate the entropy estimate for the noise source using
505 different outputs from many restarts of the source.

506 3.1.4.1 Constructing Restart Data

507 To construct restart data, the entropy source **shall** be restarted $r = 1000$ times; for each restart, $c =$
508 1000 consecutive samples **shall** be collected directly from the noise source. The output samples
509 are stored in an r by c matrix M , where $M[i][j]$ represents the j^{th} sample from the i^{th} restart.

510 Two datasets are constructed using the matrix M :

- 511 - The *row* dataset is constructed by concatenating the rows of the matrix M , i.e., the *row*
- 512 dataset is $M[1][1] \parallel \dots \parallel M[1][c] \parallel M[2][1] \parallel \dots \parallel M[2][c] \parallel \dots \parallel M[r][1] \parallel \dots \parallel M[r][c]$.
- 513 - The *column* dataset is constructed by concatenating the columns of the matrix M , i.e., the
- 514 *column* dataset is $M[1][1] \parallel \dots \parallel M[r][1] \parallel M[1][2] \parallel \dots \parallel M[r][2] \parallel \dots \parallel M[1][c] \parallel \dots \parallel M[r][c]$.

515 3.1.4.2 Validation Testing

516 The restart tests check the relations between noise source samples generated after restarting the
517 device, and compare the results to the initial entropy estimate, H_I (see Section 3.1.3).

518 First, the sanity check described in Section 3.1.4.3 is performed on the matrix M . If the test fails,
519 the validation fails and no entropy estimate is awarded.

520 If the noise source does not fail the sanity check, then the entropy estimation methods described
521 in Section 6.1 (for the IID track) or Section 6.2 (for the non-IID track) are performed on the *row*
522 and the *column* datasets, based on the track of the entropy source. Let H_r and H_c be the resulting
523 entropy estimates of the row and the column datasets, respectively. The entropy estimates from
524 the row and the column datasets are expected to be close to the initial entropy estimate H_I . If the
525 minimum of H_r and H_c is less than half of H_I , the validation fails, and no entropy estimate is
526 awarded. Otherwise, the entropy assessment of the noise source is taken as the minimum of the
527 row, the column and the initial estimates, i.e., $\min(H_r, H_c, H_I)$.

528 If the noise source does not fail the restart tests, and the entropy source does not include a
529 conditioning component, the entropy source will be validated at $\min(H_r, H_c, H_I)$. If the entropy
530 source includes a conditioning component, the entropy assessment of the entropy source is updated
531 as described in Section 3.1.5.

532 3.1.4.3 Sanity Check - Most Common Value in the Rows and Columns

533 This test checks the frequency of the most common value in the rows and the columns of the matrix
534 M . If this frequency is significantly greater than the expected value, given the initial entropy
535 estimate H_I calculated in Section 3.1.3, the restart test fails and no entropy estimate is awarded.

536 Given the 1000 by 1000 restart matrix M and the initial entropy estimate H_I , the test is performed
537 as follows:

- 538 1. Let α be $0.01/(k \times 2000)$, where k is the sample size.
- 539 2. For each row of the matrix, find the frequency of the most common sample value Fr_i , for $1 \leq i$
540 ≤ 1000 . Let F_R be the maximum of Fr_1, \dots, Fr_{1000} .
- 541 3. Repeat the same process for the columns of the matrix, i.e., find the frequency of the most
542 common sample value Fc_i for $1 \leq i \leq 1000$. Let F_C be the maximum of Fc_1, \dots, Fc_{1000} .
- 543 4. Let $F = \max(F_R, F_C)$.
- 544 5. Let $p = 2^{-H_I}$. Find the upper bound U of the $(1-\alpha)\%$ confidence interval for the frequency of
545 the most common value as $U = 1000p + Z_{(1-\alpha)}\sqrt{1000p(1-p)}$.

546 If F is greater than U , the test fails.

547 3.1.5 Entropy Estimation for Entropy Sources Using a Conditioning Component

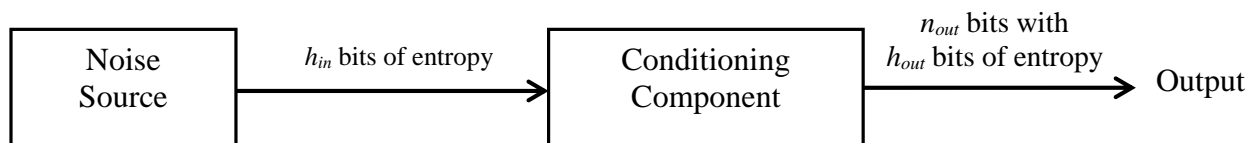
548 The optional conditioning component can be designed in various ways. Section 3.1.5.1.1 provides
549 a list of vetted cryptographic algorithms/functions for conditioning the noise source outputs.
550 Submitters are allowed to use other conditioning components; however, the entropy assessment
551 process differs from the case where a vetted conditioning component is used. If a conditioning
552 component from Section 3.1.5.1.1 is used, the entropy estimation is done as described in Section
553 3.1.5.1.2; if a non-listed algorithm is used, the entropy estimation is done as described in Section
554 3.1.5.2.

555 Let the amount of entropy in the input to the conditioning component be h_{in} bits. This input may
556 include multiple samples from one or more noise sources. For example, if the input includes w
557 samples from a noise source with h bits of entropy per sample, h_{in} is calculated as $w \times h$. If multiple
558 noise sources are used, h_{in} is calculated as the sum of amount of entropy from each noise source.

559 The submitter **shall** state the value of h_{in} , and the conditioning component **shall** produce output
560 only when at least h_{in} bits of entropy are available in its input.

561 Let the output size of the conditioning component be n_{out} (see Figure 3), and the narrowest internal
562 width within the conditioning component be q . Information on determining the narrowest internal
563 width is given in Appendix F. Denote the entropy of the output from the conditioning component
564 as h_{out} , i.e., h_{out} bits of entropy are contained within the n_{out} -bit output.

565 Since the conditioning component is deterministic, the entropy of the output is at most h_{in} .
566 However, the conditioning component may reduce the entropy of the output.



567

568 **Figure 3 Entropy of the Conditioning Component**

569 3.1.5.1 Using Vetted Conditioning Components

570 Both keyed and unkeyed algorithms have been vetted for conditioning. Section 3.1.5.1.1 provides
571 a list of vetted conditioning components. Section 3.1.5.1.2 discusses the method for determining
572 the entropy provided by a vetted conditioning component.

573 3.1.5.1.1 List of Vetted Conditioning Components

574 Three keyed algorithms have been vetted for a keyed conditioning component:

- 575 1. HMAC, as specified in FIPS 198, with any **approved** hash function specified in FIPS 180
576 or FIPS 202,
- 577 2. CMAC, as specified in SP 800-38B, with the AES block cipher (see FIPS 197), and

578 3. CBC-MAC, as specified in Appendix G, with the AES block cipher. This Recommendation
579 does not approve the use of CBC-MAC for purposes other than as a conditioning
580 component in an RBG.

581 The keys used by the keyed conditioning components **shall** be selected by the submitter in advance
582 (per implementation or per device). The submitter **shall** document how the selection is done, and
583 specify the key to test the correctness of the implementation.

584 Three unkeyed functions have been vetted for unkeyed conditioning component:

- 585 1. Any **approved** hash function specified in FIPS 180 or FIPS 202,
- 586 2. **Hash_df**, as specified in SP 800-90A, using any **approved** hash function specified in FIPS
587 180 or FIPS 202, and
- 588 3. **Block_Cipher_df**, as specified in SP800-90A using the AES block cipher (see FIPS 197).

589

590 The narrowest internal width and the output length for the vetted conditioning functions are
591 provided in the following table.

592 **Table 1 The narrowest internal width and output lengths of the vetted conditioning functions.**

Conditioning Function	Narrowest Internal Width (q)	Output Length (n_{out})
HMAC	hash-function output size	hash-function output size
CMAC	AES block size = 128	AES block size = 128
CBC-MAC	AES block size = 128	AES block size = 128
Hash Function	hash-function output size	hash-function output size
Hash_df	hash-function output size	hash-function output size
Block_Cipher_df	AES key size	AES key size

593

594 For HMAC, CMAC, CBC-MAC and the hash functions, the output length (n_{out}) specified in the
595 table is the “natural” output length of the function. For Hash_df and Block_cipher_df, the output
596 length indicated in the table **shall** be the value of *no_of_bit_to_return* used in the invocation of
597 Hash_df and Block_Cipher_df (see SP 800-90A).

598

599 3.1.5.1.2 Entropy Assessment using Vetted Conditioning Components

600 When using a conditioning component listed in Section 3.1.5.1.1 (given the assurance of correct
601 implementation by CAVP testing), the entropy of the output is estimated as

602
$$h_{out} = \begin{cases} \min(h_{in}, 0.85n_{out}, 0.85q), & \text{if } h_{in} < 2 \min(n_{out}, q) \\ \min(n_{out}, q), & \text{if } h_{in} \geq 2 \min(n_{out}, q) \end{cases}$$

603 When the input entropy is at least $2 \times \min(n_{out}, q)$, n_{out} full-entropy output bits are produced.
 604 Otherwise, the size of the output and the narrowest internal width are multiplied by the constant²
 605 0.85 for a conservative estimate.

606 If validation testing of the vetted algorithm indicates that it has not been implemented correctly,
 607 the conditioning component **shall** be treated as not vetted, and the procedure described in Section
 608 3.1.5.2 **shall** be followed.

609 The entropy source will be validated at the min-entropy per conditioned output, h_{out} , computed
 610 above.

611 Note that it is acceptable to truncate the outputs from a vetted conditioning component. If this is
 612 done, the entropy estimate is reduced to a proportion of the output (e.g., if there are six bits of
 613 entropy in an eight-bit output and the output is truncated to six bits, then the entropy is reduced to
 614 $3/4 \times 6 = 4.5$ bits).

615 3.1.5.2 Using Non-vetted Conditioning Components

616 For non-vetted conditioning components, the entropy in the output depends, in part, on the entropy
 617 of the input (h_{in}), the size of the output (n_{out}), and the size of the narrowest internal width (q). The
 618 size of the output and the narrowest internal width is multiplied by the constant 0.85 for a
 619 conservative estimate, as was done for the vetted conditioning functions listed in Section 3.1.5.1.1.
 620 However, an additional parameter is needed: the entropy of the *conditioned sequential dataset* (as
 621 described in item 2 of Section 3.1.1), which **shall** be computed using the methods described in
 622 Section 6.1 and Section 6.2 for IID and non-IID data, respectively. Let the obtained entropy
 623 estimate per bit be h' .

624 The output of the conditioning component (n_{out}) **shall** be treated as a binary string, for purposes of
 625 the entropy estimation.

626 The entropy of the conditioned output is estimated as

$$627 \quad h_{out} = \min(h_{in}, 0.85n_{out}, 0.85q, h' \times n_{out}).$$

628 The entropy source will be validated at the min-entropy per conditioned output, h_{out} , computed
 629 above.

630 Note that truncating subsequent to the use of a non-vetted conditioning component **shall not** be
 631 performed before providing output from the entropy source.

632 3.1.6 Using Multiple Noise Sources

633 If multiple independent noise sources are used, the sum of the entropies provided by each noise
 634 source is used as the entropy input to the conditioning component. For example, if the conditioning
 635 component inputs w_1 samples from Noise Source 1 with an entropy of h_1 bits per sample, and w_2

² The constant 0.85 used in the equation was selected after some empirical studies.

636 samples from Noise Source 2 with an entropy of h_2 bits per sample, then the h_{in} is calculated as
637 $w_1h_1+w_2h_2$.

638 3.2 Requirements for Validation Testing

639 In this section, high-level requirements (on both submitters and testers) are presented for validation
640 testing.

641 3.2.1 Requirements on the Entropy Source

642 The intent of these requirements is to assist the developer in designing/implementing an entropy
643 source that can provide outputs with a consistent amount of entropy and to produce the required
644 documentation for entropy source validation.

- 645 1. The entire design of the entropy source **shall** be documented, including the interaction of the
646 components specified in Section 2.2. The documentation **shall** justify why the entropy source
647 can be relied upon to produce bits with entropy.
- 648 2. Documentation **shall** describe the operation of the entropy source, including how the entropy
649 source works, and how to obtain data from within the entropy source for validation testing.
- 650 3. Documentation **shall** describe the range of operating conditions under which the entropy
651 source is claimed to operate correctly (e.g., temperature range, voltages, system activity, etc.).
652 Analysis of the entropy source's behavior at the edges of these conditions **shall** be documented,
653 along with likely failure modes.
- 654 4. The entropy source **shall** have a well-defined (conceptual) security boundary, which **should**
655 be the same as or be contained within a FIPS 140 cryptographic module boundary. This
656 security boundary **shall** be documented; the documentation **shall** include a description of the
657 content of the security boundary. Note that the security boundary may extend beyond the
658 entropy source itself (e.g., the entropy source may be contained within a larger boundary that
659 also contains a DRBG); also note that the security boundary may be logical, rather than
660 physical.
- 661 5. When a conditioning component is not used, the output from the entropy source is the output
662 of the noise source, and no additional interface is required. In this case, the noise-source output
663 is available during both validation testing and normal operation.
- 664 6. When a conditioning component is included in the entropy source, the output from the entropy
665 source is the output of the conditioning component, and an additional interface is required to
666 access the noise-source output. In this case, the noise-source output **shall** be accessible via the
667 interface during validation testing, but the interface may be disabled, otherwise. The designer
668 **shall** fully document the method used to get access to the raw noise source samples. If the
669 noise-source interface is not disabled during normal operation, any noise-source output using
670 this interface **shall not** be provided to the conditioning component for processing and eventual
671 output as normal entropy-source output.
- 672 7. The entropy source **may** restrict access to raw noise source samples to special circumstances
673 that are not available to users in the field, and the documentation **shall** explain why this

674 restriction is not expected to substantially alter the behavior of the entropy source as tested
675 during validation.

676 An optional, recommended feature of the entropy source is as follows:

- 677 8. The entropy source **may** contain multiple noise sources to improve resiliency with respect to
678 degradation or misbehavior. Only independent noise sources are allowed by this
679 Recommendation. When multiple noise sources are used, the requirements specified in Section
680 3.2.2 **shall** apply to each noise source.
- 681 9. If multiple noise sources are used, documentation **shall** specify whether all noise sources will
682 be available operationally; datasets obtained from noise sources that will not be available in
683 the field **shall not** be used for entropy assessment.

684 3.2.2 Requirements on the Noise Source

685 The entropy source will have no more entropy than that provided by the noise source, and as such,
686 the noise source requires special attention during validation testing. This is partly due to the
687 fundamental importance of the noise source (if it does not do its job, the entropy source will not
688 provide the expected amount of security), and partly because the probabilistic nature of its behavior
689 requires more complicated testing.

690 The requirements for the *noise source* are as follows:

- 691 1. The operation of the noise source **shall** be documented; this documentation **shall** include a
692 description of how the noise source works and rationale about why the noise source provides
693 acceptable entropy output, and **should** reference relevant, existing research and literature.
694 Documentation **shall** also include why it is believed that the entropy rate does not change
695 significantly during normal operation.
- 696 2. Documentation **shall** provide an explicit statement of the expected entropy rate and provide a
697 technical argument for why the noise source can support that entropy rate. This can be in broad
698 terms of where the unpredictability comes from and a rough description of the behavior of the
699 noise source (to show that it is reasonable to assume that the behavior is stable).
- 700 3. The noise source state **shall** be protected from adversarial knowledge or influence to the
701 greatest extent possible. The methods used for this **shall** be documented, including a
702 description of the (conceptual) security boundary's role in protecting the noise source from
703 adversarial observation or influence.
- 704 4. Although the noise source is not required to produce unbiased and independent outputs, it **shall**
705 exhibit random behavior; i.e., the output **shall not** be definable by any known algorithmic rule.
706 Documentation **shall** indicate whether the noise source produces IID data or non-IID data. This
707 claim will be used in determining the test path followed during validation. If the submitter
708 makes an IID claim, documentation **shall** include rationale for the claim.
- 709 5. The noise source **shall** generate fixed-length bitstrings. A description of the output space of
710 the noise source **shall** be provided. Documentation **shall** specify the fixed sample size (in bits)
711 and the list (or range) of all possible outputs from each noise source.

712 6. An ordered ranking of the bits in the n -bit samples **shall** be provided. A rank of ‘1’ **shall**
713 correspond to the bit assumed to be contributing the most entropy to the sample, and a rank of
714 n **shall** correspond to the bit contributing the least amount. If multiple bits contribute the same
715 amount of entropy, the ranks can be assigned arbitrarily among those bits. The algorithm
716 specified in Section 6.4 **shall** be used to assign ranks.

717 7. The noise source **may** include simple post-processing functions to improve the quality of its
718 outputs. When a post-processing function is used, the noise source **shall** use only one of the
719 **approved** post-processing functions: Von Neumann’s method, the linear filtering method, or
720 the length-of-runs method. The descriptions of these methods are given in Appendix E. If other
721 post-processing functions are approved in the future, they will be included in the
722 implementation guidance [IG140-2].

723 3.2.3 Requirements on the Conditioning Component

724 The requirements for the *conditioning component* are as follows:

- 725 1. If the entropy source uses a vetted conditioning component as listed in Section 3.1.5.1.1, the
726 implementation of that conditioning component **shall** be tested to obtain assurance of
727 correctness.
- 728 2. For entropy sources containing a conditioning component that is not listed in Section 3.1.5.1.1,
729 a description of the conditioning component **shall** be provided. Documentation **shall** state the
730 narrowest internal width and the size of the output blocks from the conditioning component.
- 731 3. Documentation **shall** include the minimum amount of entropy h_{in} in the input of the
732 conditioning component.

733 3.2.4 Requirements on Data Collection

734 The requirements on data collection are listed below:

- 735 1. The data collection for entropy estimation **shall** be performed in one of the three ways
736 described below:
 - 737 • By the submitter with a witness from the testing lab, or
 - 738 • By the testing lab itself, or
 - 739 • Prepared by the submitter in advance of testing, along with the following documentation:
740 a specification of the data generation process, and a signed document that attests that the
741 specification was followed.
- 742 2. Data collected from the noise source for validation testing **shall** be raw output values
743 (including digitization and optional post-processing).
- 744 3. The data collection process **shall not** require a detailed knowledge of the noise source or
745 intrusive actions that may alter the behavior of the noise source (e.g., drilling into the device).

- 746 4. Data **shall** be collected from the noise source and any conditioning component that is not listed
747 in Section 3.1.5.1.1 (if used) under normal operating conditions (i.e., when it is reasonable to
748 expect entropy in the outputs).
- 749 5. Data **shall** be collected from the entropy source under validation. Any relevant version of the
750 hardware or software updates **shall** be associated with the data.
- 751 6. Documentation on data collection **shall** be provided so that a lab or submitter can perform (or
752 replicate) the collection process at a later time, if necessary.

753 4 Health Tests

754 Health tests are an important component of the entropy source, as they aim to detect deviations
755 from the intended behavior of the noise source as quickly as possible and with a high probability.
756 Noise sources can be fragile, and hence, can be affected by the changes in operating conditions of
757 the device, such as temperature, humidity, or electric field, which might result in unexpected
758 behavior. Health tests take the entropy assessment as input, and characterize the expected behavior
759 of the noise source based on this value. Requirements on the health tests are listed in Section 4.3.

760 4.1 Health Test Overview

761 The health testing of a noise source is likely to be very technology-specific. Since, in the vast
762 majority of cases, the noise source will not produce unbiased, independent binary data, traditional
763 statistical procedures (e.g., randomness tests described in NIST SP 800-22) that test the hypothesis
764 of unbiased, independent bits will almost always fail, and thus are not useful for monitoring the
765 noise source. In general, tests on the noise source have to be tailored carefully, taking into account
766 the expected statistical behavior of the correctly operating noise source.

767 The health testing of noise sources will typically be designed to detect failures of the noise source,
768 based on the expected output during a failure, or to detect a deviation from the expected output
769 during the correct operation of the noise source. Health tests are expected to raise an alarm in three
770 cases:

- 771 1. When there is a significant decrease in the entropy of the outputs,
772 2. When noise source failures occur, or
773 3. When hardware fails, and implementations do not work correctly.

774 4.2 Types of Health Tests

775 Health tests are applied to the outputs of a noise source before any conditioning is done. (It is
776 permissible to also apply some health tests to conditioned outputs, but this is not required.)

777 *Start-up health tests* are performed after powering up or rebooting. They ensure that the entropy
778 source components are working as expected before they are used during normal operating
779 conditions, and nothing failed since the last time that the start-up tests were run. The samples
780 drawn from the noise source during the startup tests **shall not** be available for normal operations
781 until the tests are completed; after testing, these samples may simply be discarded.

782 *Continuous health tests* are run indefinitely while the entropy source is operating. Continuous tests
783 focus on the noise source behavior and aim to detect failures as the noise source runs. The purpose
784 of continuous tests is to allow the entropy source to detect many kinds of failures in its underlying
785 noise source. These tests are run continuously on all digitized samples obtained from the noise
786 source, and so tests must have a very low probability of raising a false alarm during the normal
787 operation of the noise source. In many systems, a reasonable false positive probability will make
788 it extremely unlikely that a properly functioning device will indicate a malfunction, even in a very
789 long service life. Note that continuous tests are resource-constrained – this limits their ability to
790 detect noise source problems, so that only gross failures are likely to be detected.

791 Note that the continuous health tests operate over a stream of values. These sample values may be
792 output as they are generated; there is no need to inhibit output from the noise source or entropy
793 source while running the test. It is important to understand that this may result in poor entropy
794 source outputs for a time, since the error is only signaled once significant evidence has been
795 accumulated, and these values may have already been output by the entropy source. As a result, it
796 is important that the false positive probability be set to an acceptable level. In the following
797 discussion, all calculations assume that a false positive probability of approximately once in 2^{40}
798 samples generated by the noise source is acceptable; however, the formulas given can be adapted
799 for different false positive probabilities selected by the submitter.

800 *On-demand health tests* can be called at any time. This Recommendation does not require
801 performing any particular on-demand testing during operation. However, it does require that the
802 entropy source be *capable* of performing on-demand health tests. Note that resetting, rebooting, or
803 powering up are acceptable methods for initiating an on-demand test if the procedure results in the
804 immediate execution of the start-up tests. Samples collected from the noise source during on-
805 demand health tests **shall not** be available for use until the tests are completed, and may simply be
806 discarded.

807 **4.3 Requirements for Health Tests**

808 Health tests on the noise source are a required component of an entropy source. The health tests
809 **shall** include both continuous and startup tests.

- 810 1. The submitter **shall** provide documentation that specifies all entropy source health tests and
811 their rationale. The documentation **shall** include a description of the health tests, the rate and
812 conditions under which each health test is performed (e.g., at start-up, continuously, or on-
813 demand), and rationale indicating why each test is believed to be appropriate for detecting one
814 or more failures in the entropy source.
- 815 2. The developer **shall** document any known or suspected noise source failure modes, and **shall**
816 include vendor-defined continuous tests to detect those failures.
- 817 3. Appropriate health tests tailored to the noise source **should** place special emphasis on the
818 detection of misbehavior near the boundary between the nominal operating environment and
819 abnormal conditions. This requires a thorough understanding of the operation of the noise
820 source.
- 821 4. The submitter **shall** provide source code for any tests implemented as an alternative or in
822 addition to those listed in this Recommendation.
- 823 5. Health tests **shall** be performed on the noise source samples before any conditioning is done.

- 824 6. Additional health tests **may** be performed on the outputs of the conditioning function. Any
825 such tests **shall** be fully documented.
- 826 7. In the case where a sufficiently persistent failure is detected, the entropy source **shall** notify
827 the consuming application (e.g., the RBG) of the error condition. The entropy source **may**
828 detect intermittent failures and react to them in other ways, e.g., by inhibiting output for a short
829 time, before notification of the error. The submitter **shall** describe the conditions for
830 intermittent and persistent failures.
- 831 8. The expected false positive probability of the health tests signaling a major failure to the
832 consuming application **shall** be documented.
- 833 9. The continuous tests **shall** include either:
- 834 a. The approved continuous health tests, described in Section 4.4, or
- 835 b. Some vendor-defined tests that meet the requirements to substitute for those approved
836 tests, as described in Section 4.5. If vendor-defined health tests are used in place of any
837 **approved** health tests, the tester **shall** verify that the implemented tests detect the failure
838 conditions detected by the **approved** continuous health tests, as described in Section 4.4.
839 The submitter can avoid the need to use the two approved continuous health tests by
840 providing convincing evidence that the failure being considered will be reliably detected
841 by the vendor-defined continuous tests. This evidence may be a proof or the results of
842 statistical simulations.
- 843 10. If any of the **approved** continuous health tests are used by the entropy source, the false positive
844 probability for these tests **shall** be set to at least 2^{-50} . The submitter **shall** specify and document
845 a false positive probability suitable for their application.
- 846 11. The continuous tests **may** include additional tests defined by the vendor.
- 847 12. The entropy source's startup tests **shall** run the continuous health tests over at least 4096
848 consecutive samples.
- 849 13. The samples subjected to startup testing **may** be released for operational use after the startup
850 tests have been passed.
- 851 14. The startup tests **may** include other tests defined by the vendor.
- 852 15. The entropy source **shall** support on-demand testing.
- 853 16. The entropy source **may** support on-demand testing by restarting the entropy source and
854 rerunning the startup tests, or by rerunning the startup tests without restarting the entropy
855 source. The documentation **shall** specify the approach used for on-demand testing.
- 856 17. The entropy source's on-demand testing **may** include other testing.

857 **4.4 Approved Continuous Health Tests**

858 This recommendation provides two **approved** health tests: the Repetition Count test, and the
859 Adaptive Proportion test. If these two health tests are included among the continuous health tests
860 of the entropy source, no other tests are required. However, the developer is allowed to include
861 additional health tests.

862 Both tests are designed to require minimal resources, and to be computed on-the-fly, while noise
863 source samples are being produced, possibly conditioned, and output. Neither test delays the
864 availability of the noise source samples.

865 Like all statistical tests, both of these tests have a false positive probability – the probability that a
866 correctly functioning noise source will fail the test on a given output. A reasonable choice for the
867 false positive probability in many applications is $\alpha = 2^{-40}$; this value will be used in all the
868 calculations in the rest of this section. The submitter of the entropy source must determine a
869 reasonable false positive probability, given the details of the entropy source and its consuming
870 application. In order to ensure that these tests have enough power to detect major failures, the
871 lowest allowed false positive probability for these **approved** tests is $\alpha = 2^{-50}$.

872 4.4.1 Repetition Count Test

873 The goal of the repetition count test is to quickly detect catastrophic failures that cause the noise
874 source to become "stuck" on a single output value for a long period of time. It can be seen as an
875 update of the "stuck test" which was previously required for random number generators within
876 FIPS-approved cryptographic modules.

877 Given the assessed min-entropy H of a noise source, the probability³ of that source generating n
878 identical samples consecutively is at most $2^{-H(n-1)}$. The test declares an error if a sample is repeated
879 more than the cutoff value C , which is determined by the acceptable false-positive probability α
880 and the entropy estimate H . The cutoff value of the repetition count test is calculated as:

$$881 \quad C = \left\lceil 1 + \frac{-\log_2 \alpha}{H} \right\rceil.$$

882 This value of C is the smallest integer satisfying the inequality $\alpha \geq 2^{-H(C-1)}$, which ensures that the
883 probability of obtaining a sequence of identical values from C consecutive noise source samples
884 is no greater than α . For example, for $\alpha = 2^{-40}$, an entropy source with $H = 2.0$ bits per sample
885 would have a repetition count test cutoff value of $\lceil 1 + 40/2.0 \rceil = 21$.

886 Given a dataset of noise source observations, and the cutoff value C , the repetition count test is
887 performed as follows:

- 888 1. Let A be the current sample value.
- 889 2. Initialize the counter B to 1.
- 890 3. If the next sample value is A , increment B by one.
891 If B is equal to C , return an error.
892 else:
893 Let A be the next sample value.

³ This probability can be obtained as follows. Let a random variable take possible values with probabilities p_i , for $i=1, \dots, k$, where $p_1 \geq p_2 \geq \dots \geq p_k$. Then, the probability of producing any C identical consecutive samples is $\sum p_i^C$. Since, $\sum p_i^C$ is less than or equal to $p_1 \cdot p_1^{C-1} + p_2 \cdot p_2^{C-1} + \dots + p_k \cdot p_k^{C-1} = (p_1 + \dots + p_k) p_1^{C-1} = p_1^{C-1} = 2^{-H(C-1)}$.

894 Initialize the counter B to 1.

895 Repeat Step 3.

896 Running the repetition count test requires enough memory to store:

897 A : the most recently observed sample value,

898 B : the number of consecutive times that the sample A has been observed, and

899 C : the cutoff value.

900 This test's cutoff value can be applied to any entropy estimate, H , including very small and very
901 large estimates. However, it is important to note that this test is not very powerful – it is able to
902 detect only catastrophic failures of a noise source. For example, a noise source evaluated at eight
903 bits of min-entropy per sample has a cutoff value of six repetitions to ensure a false-positive rate
904 of approximately once per one trillion samples generated. If that noise source somehow failed to
905 the point that each sample had a 1/16 probability of being the same as the previous sample, so that
906 it was providing only four bits of min-entropy per sample, it would still be expected to take about
907 sixteen million samples before the repetition count test would notice the problem.

908 4.4.2 Adaptive Proportion Test

909 The adaptive proportion test is designed to detect a large loss of entropy that might occur as a
910 result of some physical failure or environmental change affecting the noise source. The test
911 continuously measures the local frequency of occurrence of a sample value in a sequence of noise
912 source samples to determine if the sample occurs too frequently. Thus, the test is able to detect
913 when some value begins to occur much more frequently than expected, given the source's assessed
914 entropy per sample.

915 The test counts the number of times the current sample value is repeated within a window of size
916 W . If the sample is repeated more frequently than a cutoff value C , which is determined by the
917 false positive probability α and the assessed entropy/sample of the source, H , the test declares an
918 error. The window size W is selected based on the alphabet size, and **shall** be assigned to 1024 if
919 the noise source is binary (that is, it produces only two distinct values) and 512 if the noise source
920 is not binary (that is, it produces more than two distinct values).

921 Given a sequence of noise source observations, the cutoff value C and the window size W , the test
922 is performed as follows:

923 1. Let A be the current sample value.

924 2. Initialize the counter B to 1.

925 3. For $i = 1$ to $W-1$

926 If the next sample is equal to A , increment B by 1.

927 4. If $B > C$, return error.

928 5. Go to Step 1.

929 Running the test requires enough memory to store

930 A : the sample value currently being counted,

931 B : the number of times that A has been seen in the current window,

932 W : the window size,

933 i : the counter for the number of samples examined in the current window, and

934 C : the cutoff value at which the test fails.

935 The cutoff value C is chosen such that the probability of observing more than C identical samples
936 in a window size of W is at most α . Mathematically, C satisfies the following equation⁴.

937
$$\Pr(B > C) = \alpha,$$

938 where $p = 2^{-H}$. The following tables give cutoff values for various min-entropy estimates per
939 sample and window sizes with $\alpha = 2^{-40}$. For example, the cutoff value for binary sources with
940 $H=0.4$ is 867, and the probability of detecting a loss of 50% of the entropy using 1024 samples is
941 0.86, and the probability of detecting the same failure is almost 1 during the startup tests that use
942 at least 4096 samples. Note that the noise source failures whose probability of detection is listed
943 in the tables are of a very specific form – some value becomes much more common than it should
944 be, given the source’s entropy estimate, so that the maximum probability p_{max} is much higher, and
945 thus $h = -\log_2(p_{max})$ is much lower than claimed by the noise source’s entropy estimate.

946 **Table 2 Adaptive proportion test on binary data for various entropy/sample levels with $W=1024$**

H	Cutoff value	Probability of detecting noise source failure			
		50% entropy loss		33% entropy loss	
		in one window	in startup	in one window	in startup
0.2	960	0.25	0.69	0	0
0.4	867	0.86	≈ 1	0.06	0.23
0.6	779	0.81	≈ 1	0.29	0.74
0.8	697	0.76	≈ 1	0.50	0.94
1	624	0.71	0.99	0.56	0.96

947 **Table 3 Adaptive proportion test on non-binary data for various entropy/sample levels with $W=512$**

H	Cutoff value	Probability of detecting noise source failure			
		50% entropy loss		33% entropy loss	
		in one window	in startup	in one window	in startup
0.2	491	0.25	0.69	0	0.0
0.5	430	0.43	0.99	0	0.02
1	335	0.70	≈ 1	0.7	0.44
2	200	0.50	≈ 1	0.23	0.88
3	122	0.35	0.97	0.18	0.79
4	77	0.25	0.90	0.10	0.57
5	50	0.18	0.79	0.5	0.35

⁴ This probability can be computed using widely-available spreadsheet applications. In Microsoft Excel, Open Office Calc, and iWork Numbers, the calculation is done with the function =CRITBINOM(). For example, in Microsoft Excel, C would be computed as =CRITBINOM(W , power(2,(- H)), $1-\alpha$).

6	34	0.12	0.66	0.2	0.16
7	25	0.9	0.52	0.1	0.04
8	18	0.6	0.40	0	0.02

948

949 **4.5 Vendor-Defined Alternatives to the Continuous Health Tests**

950 Designer-defined tests are always permitted in addition to the two **approved** tests listed in Section
951 4.4. Under some circumstances, the vendor-defined tests may take the place of the two **approved**
952 tests. The goal of the two **approved** continuous health tests in Section 4.4, is to detect two
953 conditions:

- 954 a. Some value is consecutively repeated many more times than expected, given the assessed
955 entropy per sample of the source.
- 956 b. Some value becomes much more common in the sequence of noise source outputs than
957 expected, given the assessed entropy per sample of the source.

958 The designer of the entropy source is in an excellent position to design health tests specific to the
959 source and its known and suspected failure modes. Therefore, this Recommendation also permits
960 designer-defined alternative health tests to be used in place of the **approved** tests in Section 4.4,
961 so long as the combination of the designer-defined tests and the entropy source itself can guarantee
962 that these two conditions will not occur without being detected by the source with at least the same
963 probability.

964 **4.6 Alternative Health Test Criteria**

965 For concreteness, these are the criteria that are required for any alternative continuous health tests:

- 966 a. If a single value appears more than $\lceil 100/H \rceil$ consecutive times in a row in the sequence of
967 noise source samples, the test **shall** detect this with probability of at least 99%.
- 968 b. Let $P = 2^{-H}$. If the noise source's behavior changes so that the probability of observing a
969 specific sample value increases to at least $P^* = 2^{-H/2}$, then the test **shall** detect this with a
970 probability of at least 50% when examining 50,000 consecutive samples from this degraded
971 source.

972 The submitter can avoid the need to use the two **approved** continuous health tests by providing
973 convincing evidence that the failure being considered will be reliably detected by the vendor-
974 defined continuous tests. This evidence may be a proof or the results of statistical simulations.

975

976 **5 Testing the IID Assumption**

977 The samples from a noise source are considered to be *independent and identically distributed* (IID)
978 if each sample has the same probability distribution as every other sample, and all samples are
979 mutually independent. The IID assumption significantly simplifies the process of entropy
980 estimation. When the IID assumption does not hold, i.e., the samples are either not identically
981 distributed or are not independently distributed (or both), estimating entropy is more difficult and
982 requires different methods.

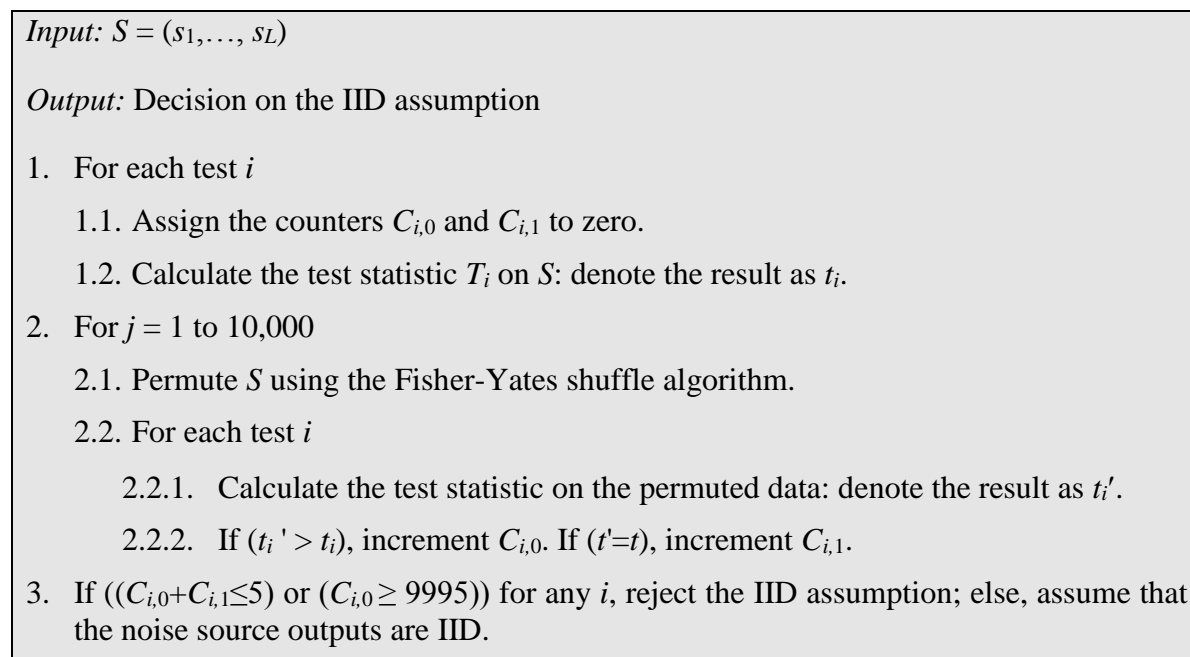
983 This section includes statistical tests that are designed to find evidence that the samples are not IID

984 and if no evidence is found that the samples are non-IID, then it is assumed that the samples are
 985 IID (see Section 3.1.1). These tests take the sequence $S = (s_1, \dots, s_L)$, where $s_i \in A = \{x_1, \dots, x_k\}$, as
 986 input, and test the hypothesis that the values in S are IID. If the hypothesis is rejected by any of
 987 the tests, the values in S are assumed to be non-IID.

988 Statistical tests based on permutation testing (also known as shuffling tests) are given in Section
 989 5.1. Five additional chi-square tests are presented in Section 5.2.

990 5.1 Permutation Testing

991 Permutation testing is a way to test a statistical hypothesis in which the actual value of the test
 992 statistic is compared to a reference distribution that is inferred from the input data, rather than a
 993 standard statistical distribution. The general approach of permutation testing is summarized in
 994 Figure 4. This is repeated for each of the test statistics described in Sections 5.1.1 – 5.1.11. The
 995 shuffle algorithm of step 2.1 is provided in Figure 5.



996 **Figure 4 Generic Structure of Permutation Testing**

997 If the samples are IID, permuting the dataset is not expected to change the value of the test statistics
 998 significantly. In particular, the original dataset and permuted datasets are expected to be drawn
 999 from the same distribution; therefore, their test statistics should be similar. Unusually high or low
 1000 test statistics are expected to occur infrequently. However, if the samples are not IID, then the
 1001 original and permuted test statistics may be significantly different. The counters $C_{i,0}$ and $C_{i,1}$ are
 1002 used to find the ranking of the original test statistics among permuted test statistics (i.e., where a
 1003 statistic for the original dataset fit within an ordered list of the permuted datasets). Extreme values
 1004 for the counters suggest that the data samples are not IID. If the sum of $C_{i,0}$ and $C_{i,1}$ is less than 5,
 1005 it means that the original test statistic has a very high rank; conversely, if $C_{i,0}$ is greater than 9995,
 1006 it means that the original test statistics has a very low rank. The cutoff values for $C_{i,0}$ and $C_{i,1}$ are

1007 calculated using a type I error⁵ of 0.001.

1008 The tests described in the following subsections are intended to check the validity of the IID
1009 assumption. Some of the tests (e.g., the compression test) are effective at detecting repeated
1010 patterns of particular values (for example, strings of sample values that occur more often than
1011 would be expected by chance if the samples were IID), whereas some of the other tests (e.g., the
1012 number of directional runs test and the runs based on the median test) focus on the association
1013 between the numeric values of the successive samples in order to find an indication of a trend or
1014 some other relation, such as high sample values that are usually followed by low sample values.

Input: $S = (s_1, \dots, s_L)$

Output: Shuffled $S = (s_1, \dots, s_L)$

1. $i = L$
2. While ($i > 1$)
 - a. Generate a random integer j that is uniformly distributed between 0 and i .
 - b. Swap s_j and s_i
 - c. $i = i - 1$

Figure 5 Pseudo-code of the Fisher-Yates Shuffle

1015
1016 For some of the tests, the number of distinct sample values, denoted k (the size of the set A),
1017 significantly affects the distribution of the test statistics, and thus the type I error. For such tests,
1018 one of the following conversions is applied to the input data, when the input is binary, i.e., $k = 2$.

1019 • *Conversion I* partitions the sequences into 8-bit non-overlapping blocks, and counts the
1020 number of ones in each block. For example, let the 20-bit input be
1021 (1,0,0,0,1,1,1,0,1,1,0,1,1,0,1,1,0,0,1,1). The first and the second 8-bit blocks include four
1022 and six ones, respectively. The last block, which is not complete, includes two ones. The
1023 output sequence is (4, 6, 2).

1024 • *Conversion II* partitions the sequences into 8-bit non-overlapping blocks, and calculates
1025 the integer value of each block. For example, let the input message be (1,0,0,0,1,1,1,0,
1026 1,1,0,1,1,0,1,1,0,0,1,1). The integer values of the first two blocks are 142, and 219. Zeroes
1027 are appended when the last block has less than 8 bits. Then, the last block becomes
1028 (0,0,1,1,0,0,0,0) with an integer value of 48. The output sequence is (142, 219, 48).

1029 Descriptions of the individual tests will provide guidance on when to use each of these conversions.

1030 5.1.1 Excursion Test Statistic

1031 The excursion test statistic measures how far the running sum of sample values deviates from its
1032 average value at each point in the dataset. Given $S = (s_1, \dots, s_L)$, the test statistic T is the largest

⁵ A type I error occurs when the null hypothesis is true, but is rejected by the test.

1033 deviation from the average and is calculated as follows:

1034 1. Calculate the average of the sample values, i.e., $\bar{X} = (s_1 + s_2 + \dots + s_L) / L$

1035 2. For $i = 1$ to L

1036 Calculate $d_i = | \sum_{j=1}^i s_j - i \times \bar{X} |$.

1037 3. $T = \max (d_1, \dots, d_L)$.

1038 *Example 1:* Let the input sequence be $S = (2, 15, 4, 10, 9)$. The average of the sample values is 8,
1039 and $d_1 = |2-8| = 6$; $d_2 = |(2+15) - (2 \times 8)| = 1$; $d_3 = |(2+15+4) - (3 \times 8)| = 3$; $d_4 = |(2+15+4+10) -$
1040 $(4 \times 8)| = 1$; and $d_5 = |(2+15+4+10+9) - (5 \times 8)| = 0$. Then, $T = \max(6, 1, 3, 1, 0) = 6$.

1041 *Handling Binary data:* The test can be applied to binary data, and no additional conversion steps
1042 are required.

1043 5.1.2 Number of Directional Runs

1044 This test statistic determines the number of runs constructed using the relations between
1045 consecutive samples. Given $S = (s_1, \dots, s_L)$, the test statistic T is calculated as follows:

1046 1. Construct the sequence $S' = (s'_1, \dots, s'_{L-1})$, where

1047
$$s'_i = \begin{cases} -1, & \text{if } s_i > s_{i+1} \\ +1, & \text{if } s_i \leq s_{i+1} \end{cases}$$

1048 for $i = 1, \dots, L-1$.

1049 2. The test statistic T is the number of runs in S' .

1050 *Example 2:* Let the input sequence be $S = (2, 2, 2, 5, 7, 7, 9, 3, 1, 4, 4)$; then $S' = (+1, +1, +1, +1,$
1051 $+1, +1, -1, -1, +1, +1)$. There are three runs: $(+1, +1, +1, +1, +1, +1)$, $(-1, -1)$ and $(+1, +1)$, so T
1052 $= 3$.

1053 *Handling Binary data:* To test binary input data, first apply Conversion I to the input sequence.

1054 5.1.3 Length of Directional Runs

1055 This test statistic determines the length of the longest run constructed using the relations between
1056 consecutive samples. Given $S = (s_1, \dots, s_L)$, the test statistic T is calculated as follows:

1057 1. Construct the sequence $S' = (s'_1, \dots, s'_{L-1})$, where

1058
$$s'_i = \begin{cases} -1, & \text{if } s_i > s_{i+1} \\ +1, & \text{if } s_i \leq s_{i+1} \end{cases}$$

1059 for $i = 1, \dots, L-1$.

1060 2. The test statistic T is the length of the longest run in S' .

1061 *Example 3:* Let the input sequence be $S = (2, 2, 2, 5, 7, 7, 9, 3, 1, 4, 4)$; then $S' = (+1, +1, +1, +1,$

1062 +1, +1, -1, -1, +1, +1). There are three runs: (+1, +1, +1, +1, +1, +1), (-1, -1) and (+1, +1), so T
1063 = 6.

1064 *Handling Binary data:* To test binary input data, first apply Conversion I to the input sequence.

1065 **5.1.4 Number of Increases and Decreases**

1066 This test statistic determines the maximum number of increases or decreases between consecutive
1067 sample values. Given $S = (s_1, \dots, s_L)$, the test statistic T is calculated as follows:

1068 1. Construct the sequence $S' = (s'_1, \dots, s'_{L-1})$, where

$$1069 \quad s'_i = \begin{cases} -1, & \text{if } s_i > s_{i+1} \\ +1, & \text{if } s_i \leq s_{i+1} \end{cases}$$

1070 for $i = 1, \dots, L-1$.

1071 2. Calculate the number of -1's and +1's in S' ; the test statistic T is the maximum of these
1072 numbers, i.e., $T = \max(\text{number of -1's}, \text{number of +1's})$.

1073 *Example 4:* Let the input sequence be $S = (2, 2, 2, 5, 7, 7, 9, 3, 1, 4, 4)$; then $S' = (+1, +1, +1, +1,$
1074 $+1, +1, -1, -1, +1, +1)$. There are eight +1's and two -1's in S' , so $T = \max(\text{number of +1s},$
1075 $\text{number of -1s}) = \max(8, 2) = 8$.

1076 *Handling Binary data:* To test binary input data, first apply the Conversion I to the input sequence.

1077 **5.1.5 Number of Runs Based on the Median**

1078 This test statistic determines the number of runs that are constructed with respect to the median of
1079 the input data. Given $S = (s_1, \dots, s_L)$, the test statistic T is calculated as follows:

1080 1. Find the median \tilde{X} of $S = (s_1, \dots, s_L)$.

1081 2. Construct the sequence $S' = (s'_1, \dots, s'_L)$ where

$$1082 \quad s'_i = \begin{cases} -1, & \text{if } s_i < \tilde{X} \\ +1, & \text{if } s_i \geq \tilde{X} \end{cases}$$

1083 for $i = 1, \dots, L$.

1084 3. The test statistic T is the number of runs in S' .

1085 *Example 5:* Let the input sequence be $S = (5, 15, 12, 1, 13, 9, 4)$. The median of the input sequence
1086 is 9. Then, $S' = (-1, +1, +1, -1, +1, +1, -1)$. The runs are (-1), (+1, +1), (-1), (+1, +1), and (-1).
1087 There are five runs, hence $T = 5$.

1088 *Handling Binary data:* When the input data is binary, the median of the input data is assumed to
1089 be 0.5. No additional conversion steps are required.

1090 5.1.6 Length of Runs Based on Median

1091 This test statistic determines the length of the longest run that is constructed with respect to the
1092 median of the input data and is calculated as follows:

1093 1. Find the median \tilde{X} of $S = (s_1, \dots, s_L)$.

1094 2. Construct a temporary sequence $S' = (s'_1, \dots, s'_L)$ from the input sequence $S = (s_1, \dots, s_L)$, as

$$1095 \quad s'_i = \begin{cases} -1, & \text{if } s_i < \tilde{X} \\ +1, & \text{if } s_i \geq \tilde{X} \end{cases}$$

1096 for $i = 1, \dots, L$.

1097 3. The test statistic T is the length of the longest run S' .

1098 *Example 6:* Let the input sequence be $S = (5, 15, 12, 1, 13, 9, 4)$. The median for this data subset
1099 is 9. Then, $S' = (-1, +1, +1, -1, +1, +1, -1)$. The runs are (-1) , $(+1, +1)$, (-1) , $(+1, +1)$, and (-1) .
1100 The longest run has a length of 2; hence, $T=2$.

1101 *Handling Binary data:* When the input data is binary, the median of the input data is assumed to
1102 be 0.5. No additional conversion steps are required.

1103 5.1.7 Average Collision Test Statistic

1104 The average collision test statistic counts the number of successive sample values until a duplicate
1105 is found. The average collision test statistic is calculated as follows:

1106 1. Let C be a list of the number of the samples observed to find two occurrences of the same
1107 value in the input sequence $S = (s_1, \dots, s_L)$. C is initially empty.

1108 2. Let $i = 1$.

1109 3. While $i < L$

1110 a. Find the smallest j such that (s_i, \dots, s_{i+j-1}) contains two identical values. If no such j
1111 exists, break out of the while loop.

1112 b. Add j to the list C .

1113 c. $i = i + j + 1$

1114 4. The test statistic T is the average of all values in the list C .

1115 *Example 7:* Let the input sequence be $S = (2, 1, 1, 2, 0, 1, 0, 1, 1, 2)$. The first collision occurs for
1116 $j = 3$, since the second and third values are the same. 3 is added to the list C . Then, the first three
1117 samples are discarded, and the next sequence to be examined is $(2, 0, 1, 0, 1, 1, 2)$. The collision
1118 occurs for $j = 4$. The third sequence to be examined is $(1, 1, 2)$, and the collision occurs for $j = 2$.
1119 There are no collisions in the final sequence (2) . Hence, $C = [3, 4, 2]$. The average of the values in
1120 C is $T = 3$.

1121 *Handling Binary data:* To test binary input data, first apply Conversion II to the input sequence.

1122 **5.1.8 Maximum Collision Test Statistic**

1123 The maximum collision test statistic counts the number of successive sample values until a
1124 duplicate is found. The maximum collision test statistic is calculated as follows:

- 1125 1. Let C be a list of the number of samples observed to find two occurrences of the same value
1126 in the input sequence $S = (s_1, \dots, s_L)$. C is initially empty.
- 1127 2. Let $i = 1$.
- 1128 3. While $i < L$
 - 1129 a. Find the smallest j such that (s_i, \dots, s_{i+j-1}) contains two identical values. If no such j
1130 exists, break out of the while loop.
 - 1131 b. Add j to the list C .
 - 1132 c. $i = i + j + 1$
- 1133 4. The test statistic T is the maximum value in the list C .

1134 *Example 8:* Let the input data be (2, 1, 1, 2, 0, 1, 0, 1, 1, 2). $C = [3, 4, 2]$ is computed as in Example
1135 7. $T = \max(3, 4, 2) = 4$.

1136 *Handling Binary data:* To test binary input data, first apply Conversion II to the input sequence.

1137 **5.1.9 Periodicity Test Statistic**

1138 The periodicity test aims to determine the number of periodic structures in the data. The test takes
1139 a lag parameter p as input, where $p < L$, and the test statistic T is calculated as follows:

- 1140 1. Initialize T to zero.
- 1141 2. For $i = 1$ to $L - p$
1142 If $(s_i = s_{i+p})$, increment T by one.

1143 *Example 9:* Let the input data be (2, 1, 2, 1, 0, 1, 0, 1, 1, 2), and let $p = 2$. Since $s_i = s_{i+p}$ for five
1144 values of i (1, 2, 4, 5 and 6), $T = 5$.

1145 *Handling Binary data:* To test binary input data, first apply Conversion I to the input sequence.

1146 The test is repeated for five different values of p : 1, 2, 8, 16, and 32.

1147 **5.1.10 Covariance Test Statistic**

1148 The covariance test measures the strength of the lagged correlation. The test takes a lag value $p <$
1149 L as input. The test statistic is calculated as follows:

- 1150 1. Initialize T to zero.
- 1151 2. For $i = 1$ to $L - p$
1152 $T = T + (s_i \times s_{i+p})$

1153 *Example 10:* Let the input data be (5, 2, 6, 10, 12, 3, 1), and let p be 2. T is calculated as $(5 \times 6) +$
1154 $(2 \times 10) + (6 \times 12) + (10 \times 3) + (12 \times 1) = 164$.

1155 *Handling Binary data:* To test binary input data, first apply Conversion I to the input sequence.

1156 The test is repeated for five different values of p : 1, 2, 8, 16, and 32.

1157 **5.1.11 Compression Test Statistics**

1158 General-purpose compression algorithms are well adapted for removing redundancy in a character
1159 string, particularly involving commonly recurring subsequences of characters. The compression
1160 test statistic for the input data is the length of that data subset after the samples are encoded into a
1161 character string and processed by a general-purpose compression algorithm. The compression test
1162 statistic is computed as follows:

- 1163 1. Encode the input data as a character string containing a list of values separated by a single
1164 space, e.g., “ $S = (144, 21, 139, 0, 0, 15)$ ” becomes “144 21 139 0 0 15”.
- 1165 2. Compress the character string with the bzip2 compression algorithm provided in [BZ2].
- 1166 3. T is the length of the compressed string, in bytes.

1167 *Handling Binary data:* The test can be applied directly to binary data, with no conversion required.

1168 **5.2 Additional Chi-square Statistical Tests**

1169 This section includes additional chi-square statistical procedures to test independence and
1170 goodness-of-fit. The independence tests attempt to discover dependencies in the probabilities
1171 between successive samples in the (entire) sequence submitted for testing (see Section 5.2.1 for
1172 non-binary data and Section 5.2.3 for binary data); the goodness-of-fit tests attempt to discover a
1173 failure to follow the same distribution in ten data subsets produced from the (entire) input sequence
1174 submitted for testing (see Section 5.2.2 for non-binary data and Section 5.2.4 for binary data). The
1175 length of the longest repeated substring test is provided in Section 5.2.5.

1176 **5.2.1 Testing Independence for Non-Binary Data**

1177 Given the input $S = (s_1, \dots, s_L)$, where $s_i \in A = \{x_1, \dots, x_k\}$, the following steps are initially performed
1178 to determine the number of bins q needed for the chi-square tests.

- 1179 1. Find the proportion p_i of each x_i in S , i.e., $p_i = \frac{\text{number of } x_i \text{ in } S}{L}$. Calculate the expected number
1180 of occurrences of each possible pair (z_i, z_j) in S , as $e_{i,j} = p_i p_j (L - 1)$.
- 1181 2. Allocate the possible (z_i, z_j) pairs, starting from the smallest $e_{i,j}$, into bins such that the
1182 expected value of each bin is at least 5. The expected value of a bin is equal to the sum of the
1183 $e_{i,j}$ values of the pairs that are included in the bin. After allocating all pairs, if the expected
1184 value of the last bin is less than 5, merge the last two bins. Let q be the number of bins
1185 constructed after this procedure.

1186 After constructing the bins, the Chi-square test is executed as follows:

1187 1. For each pair (s_j, s_{j+1}) , $1 \leq j \leq L-1$, count the number of observed values for each bin, denoted
1188 as o_i , ($1 \leq i \leq q$).

1189 2. The test statistic is calculated as $T = \sum_{i=1}^q \frac{(o_i - E(\text{Bin}_i))^2}{E(\text{Bin}_i)}$. The test fails if T is greater than the
1190 critical value of the Chi-square test statistic with $q-1$ degrees of freedom when the type I error
1191 is chosen as 0.001.

1192 *Example 11:* Let S be (2, 2, 3, 1, 3, 2, 3, 2, 1, 3, 1, 1, 2, 3, 1, 1, 2, 2, 2, 3, 3, 2, 3, 2, 3, 1, 2, 2, 3, 3,
1193 2, 2, 2, 1, 3, 3, 3, 2, 3, 2, 1, 3, 2, 3, 1, 2, 2, 3, 1, 1, 3, 2, 3, 2, 3, 1, 2, 2, 3, 3, 2, 2, 2, 1, 3, 3, 3, 2, 3,
1194 2, 1, 2, 2, 3, 3, 3, 2, 3, 2, 1, 2, 2, 2, 1, 3, 3, 3, 2, 3, 2, 1, 3, 2, 3, 1, 2, 2, 3, 1, 1). The sample space
1195 consists of $k=3$ values $\{1, 2, 3\}$; and p_1, p_2 , and p_3 are 0.21, 0.41 and 0.38, respectively. With
1196 $L=100$, the sorted expected values are calculated as:

(z_i, z_j)	(1,1)	(1,3)	(3,1)	(1,2)	(2,1)	(3,3)	(2,3)	(3,2)	(2,2)
$e_{i,j}$	4.41	7.98	7.98	8.61	8.61	14.44	15.58	15.58	16.81

1197 The pairs can be allocated into $q = 8$ bins.

Bin	Pairs	$E(\text{Bin}_i)$
1	(1,1), (1,3)	12.39
2	(3,1)	7.98
3	(1,2)	8.61
4	(2,1)	8.61
5	(3,3)	14.44
6	(2,3)	15.58
7	(3,2)	15.58
8	(2,2)	16.81

1198
1199 The frequencies for the bins are calculated as 13, 9, 8, 8, 10, 19, 18 and 14 respectively, and the
1200 test statistics is calculated as 3.2084. The hypothesis is not rejected, since the test statistics is less
1201 than the critical value 24.322.

1202
1203 **5.2.2 Testing Goodness-of-fit for non-binary data**

1204 The test checks whether the distribution of samples are identical for different parts of the input.
1205 Given the input $S = (s_1, \dots, s_L)$, where $s_i \in A = \{x_1, \dots, x_k\}$, perform the following steps to calculate
1206 the number of bins q for the test.

- 1207 1. Let c_i be the number of occurrences of x_i in S , and let $e_i = c_i/10$, for $1 \leq i \leq k$. Note that c_i is
1208 divided by ten because S will be partitioned into ten data subsets.
- 1209 2. Let $List[i]$ be the sample value with the i^{th} smallest e_i (e.g., $List[1]$ has the smallest value for
1210 e_i ; $List[2]$ has the next smallest value, etc.)
- 1211 3. Starting from $List[1]$, allocate the sample values into bins. Assign consecutive $List[i]$ values to
1212 a bin until the sum of the e_i for those binned items is at least five, then begin assigning the
1213 following $List[i]$ values to the next bin. If the expected value of the last bin is less than five,
1214 merge the last two bins. Let q be the number of bins constructed after this procedure.
- 1215 4. Let E_i be the expected number of sample values in Bin i ; E_i is the sum of the e_i for the listed
1216 items in that bin. For example, if Bin 1 contains $(x_1, x_{10}$ and $x_{50})$, then $E_1 = e_1 + e_{10} + e_{50}$.

1217 *Example 12:* Let the number of distinct sample values k be 4; and let $c_1=43$, $c_2=55$, $c_3=52$ and
1218 $c_4=10$. After partitioning the input sequence into 10 parts, the expected value of each sample
1219 becomes $e_1=4.3$, $e_2=5.5$, $e_3=5.2$ and $e_4=1$. The sample list starting with the smallest expected value
1220 is formed as $List = [4, 1, 3, 2]$. The first bin contains sample 4 and 1, and the expected value of
1221 Bin 1 becomes 5.3 ($= e_4+e_1$). The second bin contains sample 2, and the last bin contains sample
1222 3. Since the expected value of the last bin is greater than five, no additional merging is necessary.

1223 The chi-square goodness-of-fit test is executed as follows:

- 1224 1. Partition S into ten non-overlapping sequences of length $\lfloor \frac{L}{10} \rfloor$, where $S_d =$
1225 $(s_{d\lfloor L/10 \rfloor + 1}, \dots, s_{(d+1)\lfloor L/10 \rfloor})$ for $d = 0, \dots, 9$. If L is not a multiple of 10, the remaining bits are
1226 not used.
- 1227 2. $T = 0$.
- 1228 3. For $d = 0$ to 9
- 1229 3.1. For $i = 1$ to q
- 1230 3.1.1. Let o_i be the number of sample values from Bin i in the data subset S_d .
- 1231 3.1.2. $T = T + \frac{(o_i - E_i)^2}{E_i}$

1232 The test fails if the test statistic T is greater than the critical value of chi-square with $9(q-1)$ degrees
1233 of freedom when the type I error is chosen as 0.001.

1234 5.2.3 Testing Independence for Binary Data

1235 This test checks the independence assumption for binary data. A chi-square test for independence
1236 between adjacent bits could be used, but its power is limited, due to the small output space (i.e.,
1237 the use of binary inputs). A more powerful check can be achieved by comparing the frequencies
1238 of m -bit tuples to their expected values that are calculated by multiplying the probabilities of each
1239 successive bit, i.e., assuming that the samples are independent. If nearby bits are not independent,
1240 then the expected probabilities of m -bit tuples derived from their bit probabilities will be biased
1241 for the whole dataset, and a chi-square test statistic will be much larger than expected.

1242 Given the input binary data $S = (s_1, \dots, s_L)$, the length of the tuples, m , is determined as follows:

- 1243 1. Let p_0 and p_1 be the proportion of zeroes and ones in S , respectively, i.e., $p_0 =$
 1244 $\frac{\text{number of zeroes in } S}{L}$, and $p_1 = \frac{\text{number of ones in } S}{L}$.
- 1245 2. Find the maximum integer m such that $(p_0)^m > 5/L$ and $(p_1)^m > 5/L$. If m is greater than 11,
 1246 assign m to 11. If m is 1, the test fails. For example, for $p_0 = 0.14$, $p_1 = 0.86$, and $L = 1000$,
 1247 the value of m is selected as 3.

1248 The test is applied if $m \geq 2$.

- 1249 1. Initialize T to 0.
- 1250 2. For each possible m -bit tuple (a_1, a_2, \dots, a_m)
- 1251 a. Let o be the number of times that the pattern (a_1, a_2, \dots, a_m) occurs in the input
 1252 sequence S . Note that the tuples are allowed to overlap. For example, the number
 1253 of times that $(1, 1, 1)$ occurs in $(1, 1, 1, 1)$ is 2.
- 1254 b. Let w be the number of ones in (a_1, a_2, \dots, a_m) .
- 1255 c. Let $e = p_1^w (p_0)^{m-w} (L - m + 1)$.
- 1256 d. $T = T + \frac{(o-e)^2}{e}$.

1257 The test fails if the test statistic T is greater than the critical value of chi-square with $2^m - 1$ degrees
 1258 of freedom, when the type I error is chosen as 0.001.

1259 5.2.4 Testing Goodness-of-fit for Binary Data

1260 This test checks the distribution of the number of ones in non-overlapping intervals of the input
 1261 data to determine whether the distribution of the ones remains the same throughout the sequence.
 1262 Given the input binary data $S = (s_1, \dots, s_L)$, the test description is as follows:

- 1263 1. Let p be the proportion of ones in S , i.e., $p = (\text{the number of ones in } S) / L$.
- 1264 2. Partition S into ten non-overlapping sub-sequences of length $\lfloor \frac{L}{10} \rfloor$, where $S_d =$
 1265 $(s_{d\lfloor L/10 \rfloor + 1}, \dots, s_{(d+1)\lfloor L/10 \rfloor})$ for $d = 0, \dots, 9$. If L is not a multiple of 10, the remaining bits
 1266 are discarded.
- 1267 3. Initialize T to 0.
- 1268 4. Let the expected number of ones in each sub-sequence S_d be $e = p \lfloor \frac{L}{10} \rfloor$.
- 1269 5. For $d = 0$ to 9
- 1270 a. Let o be the number of ones in S_d .
- 1271 b. $T = T + \frac{(o-e)^2}{e}$.

1272 T is a chi-square random variable with 9 degrees of freedom. The test fails if S is larger than the
 1273 critical value at 0.001, which is 27.88.

1274 **5.2.5 Length of the Longest Repeated Substring Test**

1275 This test checks the IID assumption using the length of the longest repeated substring. If this length
1276 is significantly longer than the expected value, then the test invalidates the IID assumption. The
1277 test can be applied to binary and non-binary inputs.

1278 Given the input $S = (s_1, \dots, s_L)$, where $s_i \in A = \{x_1, \dots, x_k\}$,

- 1279 1. Find the proportion p_i of each possible input value x_i in S , i.e., $p_i = \frac{\text{number of } x_i \text{ in } S}{L}$.
- 1280 2. Calculate the collision probability as $p_{col} = \sum_{i=1}^k p_i^2$.
- 1281 3. Find the length of the longest repeated substring W , i.e., find the largest W such that, for at
1282 least one $i \neq j$, $s_i = s_j$, $s_{i+1} = s_{j+1}, \dots, s_{i+W-1} = s_{j+W-1}$.
- 1283 4. The number of overlapping subsequences of length W in S is $L - W + 1$, and the number of pairs
1284 of overlapping subsequences is $\binom{L - W + 1}{2}$.
- 1285 5. Let E be a binomially distributed random variable with parameters $N = \binom{L - W + 1}{2}$ and a
1286 probability of success p_{col} . Calculate the probability that E is greater than or equal to 1, i.e.,
1287 $\Pr(E \geq 1) = 1 - \Pr(E = 0) = 1 - (1 - p_{col})^N$.

1288 The test fails if $\Pr(E \geq 1)$ is less than 0.001.

1289 **6 Estimating Min-Entropy**

1290 One of the essential requirements of an entropy source is the ability to reliably create random
1291 outputs. To ensure that sufficient entropy is input to an RBG construction in SP 800-90C, the
1292 amount of entropy produced per noise source sample must be determined. This section describes
1293 generic estimation methods that will be used to test the noise source and also the conditioning
1294 component, when non-vetted conditioning components are used.

1295 Each estimator takes a sequence $S = (s_1, \dots, s_L)$ as its input, where each s_i comes from an output
1296 space $A = \{x_1, \dots, x_k\}$ that is specified by the submitter. The estimators presented in this
1297 Recommendation follow a variety of strategies, which cover a range of assumptions about the data.
1298 For further information about the theory and origins of these estimators, see Appendix H. The
1299 estimators that are to be applied to a sequence depend on whether the data has been determined to
1300 be IID or non-IID. For IID data, the min-entropy estimation is determined as specified in Section
1301 6.1, whereas for non-IID data, the procedures in Section 6.2 are used.

1302 The estimators presented in this section work well when the entropy-per-sample is greater than
1303 0.1. For alphabet sizes greater than 256, some of the estimators are not very efficient. Therefore,
1304 for efficiency purposes, the method described in Section 6.4 can be used to reduce the sample space
1305 of the outputs.

1306 **6.1 IID Track: Entropy Estimation for IID Data**

1307 For sources with IID outputs, the min-entropy estimation is determined using the *most common*

1308 *value* estimate described in Section 6.3.1. It is important to note that the estimate provides an
1309 overestimation when the samples from the source are not IID.

1310 **6.2 Non-IID Track: Entropy Estimation for Non-IID Data**

1311 Many viable noise sources fail to produce IID outputs. Moreover, some sources may have
1312 dependencies that are beyond the ability of the tester to address. To derive any utility out of such
1313 sources, a diverse and conservative set of entropy tests are required. Testing sequences with
1314 dependent values may result in overestimates of entropy. However, a large, diverse battery of
1315 estimates minimizes the probability that such a source's entropy is greatly overestimated.

1316 For non-IID data, the following estimators are calculated on the outputs of the noise source, outputs
1317 of any conditioning component that is not listed in Section 3.1.5.1.1 and outputs of any vetted
1318 conditioning function that hasn't been validated as correctly implemented, and the minimum of all
1319 the estimates is taken as the entropy assessment of the entropy source for this Recommendation:

- 1320 • The Most Common Value Estimate (Section 6.3.1),
- 1321 • The Collision Estimate (Section 6.3.2),
- 1322 • The Markov Estimate (Section 6.3.3),
- 1323 • The Compression Estimate (Section 6.3.4),
- 1324 • The t -Tuple Estimate (Section 6.3.5),
- 1325 • The Longest Repeated Substring (LRS) Estimate (Section 6.3.6),
- 1326 • The Multi Most Common in Window Prediction Estimate (Section 6.3.7),
- 1327 • The Lag Prediction Estimate (Section 6.3.8),
- 1328 • The MultiMMC Prediction Estimate (Section 6.3.9), and
- 1329 • The LZ78Y Prediction Estimate (Section 6.3.10).

1330 **6.3 Estimators**

1331 **6.3.1 The Most Common Value Estimate**

1332 This method first finds the proportion \hat{p} of the most common value in the input dataset, and then
1333 constructs a confidence interval for this proportion. The upper bound of the confidence interval is
1334 used to estimate the min-entropy per sample of the source.

1335 Given the input $S = (s_1, \dots, s_L)$, where $s_i \in A = \{x_1, \dots, x_k\}$,

- 1336 1. Find the proportion of the most common value \hat{p} in the dataset, i.e.,

$$1337 \quad \hat{p} = \max_i \frac{\#\{x_i \text{ in } S\}}{L}.$$

- 1338 2. Calculate an upper bound on the probability of the most common value p_u as

1339
$$p_u = \min \left(1, \hat{p} + 2.576 \sqrt{\frac{\hat{p}(1-\hat{p})}{L}} \right).$$

1340 3. The estimated min-entropy is $-\log_2(p_u)$.

1341 *Example:* If the dataset is $S = (0, 1, 1, 2, 0, 1, 2, 2, 0, 1, 0, 1, 1, 0, 2, 2, 1, 0, 2, 1)$, with $L = 20$, the
1342 most common value is 1, with $\hat{p} = 0.4$. $p_u = 0.4 + 2.576\sqrt{0.012} = 0.6822$. The min-entropy
1343 estimate is $-\log_2(0.6822) = 0.5518$.

1344 6.3.2 The Collision Estimate

1345 The collision estimate, proposed by Hagerty and Draper [HD12], measures the mean number of
1346 samples to the first collision in a dataset, where a collision is any repeated value. The goal of the
1347 method is to estimate the probability of the most-likely output value, based on the collision times.
1348 The method will produce a low entropy estimate for noise sources that have considerable bias
1349 toward a particular output or value (i.e., the mean time until a collision is relatively short), while
1350 producing a higher entropy estimate for a longer mean time to collision.

1351 Given the input $S = (s_1, \dots, s_L)$, where $s_i \in A = \{x_1, \dots, x_k\}$,

- 1352 1. Set $v = 1$, $index = 1$.
- 1353 2. Beginning with s_{index} , step through the input until any observed value is repeated; i.e., find
1354 the smallest j such that $s_i = s_j$, for some i with $index \leq i < j$.
- 1355 3. Set $t_v = j - index + 1$, $v = v + 1$, and $index = j + 1$.
- 1356 4. Repeat steps 2-3 until the end of the dataset is reached.
- 1357 5. Set $v = v - 1$.
- 1358 6. If $v < 1000$, map down the noise source outputs (see Section 6.4), based on the ranking
1359 provided, and retest the data.
- 1360 7. Calculate the sample mean \bar{X} , and the sample standard deviation $\hat{\sigma}$, of t_i as

1361
$$\bar{X} = \frac{1}{v} \sum_{i=1}^v t_i, \quad \hat{\sigma} = \sqrt{\frac{1}{v} \sum_{i=1}^v (t_i - \bar{X})^2}.$$

- 1362 8. Compute the lower-bound of the confidence interval for the mean, based on a normal
1363 distribution with a confidence level of 99%,

1364
$$\bar{X}' = \bar{X} - 2.576 \frac{\hat{\sigma}}{\sqrt{v}}.$$

- 1365 9. Let k be the number of possible values in the output space. Using a binary search, solve
1366 for the parameter p , such that

1367
$$\bar{X}' = pq^{-2} \left(1 + \frac{1}{k} (p^{-1} - q^{-1}) \right) F(q) - pq^{-1} \frac{1}{k} (p^{-1} - q^{-1}).$$

1368 where

$$1369 \quad q = \frac{1-p}{k-1},$$

$$1370 \quad p \geq q,$$

$$1371 \quad F(1/z) = \Gamma(k+1, z)z^{-k-1}e^z,$$

1372 and $\Gamma(a, b)$ is the incomplete Gamma function⁶.

1373 10. If the binary search yields a solution, then the min-entropy estimation is the negative
1374 logarithm of the parameter, p :

$$1375 \quad \text{min-entropy} = -\log_2(p).$$

1376 If the search does not yield a solution, then the min-entropy estimation is:

$$1377 \quad \text{min-entropy} = \log_2(k).$$

1378 *Example:* Suppose that $S = (2, 2, 0, 1, 0, 2, 0, 1, 2, 1, 2, 0, 1, 2, 1, 0, 0, 1, 0, 0, 0)$. After step 5, $v=6$,
1379 and the sequence (t_1, \dots, t_v) is $(2, 3, 4, 4, 4, 3)$. For purposes of illustration, step 6 is skipped in this
1380 example. Then $\bar{X} = 3.3333$, $\hat{\sigma} = 0.7454$, and $\bar{X}' = 2.5495$. The solution to the equation is $p =$
1381 0.7063 , giving an estimated min-entropy of 0.5016 .

1382 6.3.3 The Markov Estimate

1383 In a first-order Markov process, the next sample value depends only on the latest observed sample
1384 value; in an n^{th} -order Markov process, the next sample value depends only on the previous n
1385 observed values. Therefore, a Markov model can be used as a template for testing sources with
1386 dependencies. The Markov estimate provides a min-entropy estimate by measuring the
1387 dependencies between consecutive values from the input dataset. The min-entropy estimate is
1388 based on the entropy present in any subsequence (i.e., chain) of outputs, instead of an estimate of
1389 the min-entropy per output.

1390 The key component in estimating the entropy of a Markov process is the ability to accurately
1391 estimate the transition matrix probabilities of the Markov process. The main difficulty in making
1392 these estimates is the large data requirement necessary to resolve the dependencies. In particular,
1393 low-probability transitions may not occur often in a “small” dataset; the more data provided, the
1394 easier it becomes to make accurate estimates of transition probabilities. This method, however,
1395 avoids large data requirements by overestimating the low-probability transitions; as a
1396 consequence, an underestimate of min-entropy is obtained with less data.

⁶ The equation presented here uses the incomplete gamma function, which has known efficient computation methods and can be found in many software libraries. An efficient approximation for the incomplete Gamma function is provided in Appendix H. For additional representations of the \bar{X}' calculation in step 9, see [HD12].

1397 The data requirement for this estimation method depends on the number of output samples k (i.e.,
1398 the alphabet size); the largest k accommodated by this test is 2^6 . An alphabet size greater than 2^6
1399 cannot be accommodated, since an unreasonable amount of data would be required to accurately
1400 estimate the matrix of transition probabilities – far more than is specified in Section 3.1.1⁷. For
1401 16-bit samples, for instance, a transition matrix of size $2^{16} \times 2^{16}$, containing 2^{32} sample entries,
1402 would have to be approximated, and the data requirement for this would be impractical.

1403 As an alternative for datasets with a number of samples greater than 64, the method in Section 6.4
1404 for mapping noise source outputs (based on a ranking of the bits in the output) **shall** be
1405 implemented. This will reduce the data requirement to a more feasible quantity.

1406 Samples are collected from the noise source, and specified as d -long chains of samples. From this
1407 data, probabilities are determined for both the initial state and transitions between any two states.
1408 Any values for which these probabilities cannot be determined empirically are overestimated to
1409 guarantee that the eventual min-entropy estimate is a lower bound. These probabilities are used to
1410 determine the highest probability of any particular d -long chain of samples. The corresponding
1411 maximum probability is used to determine the min-entropy present in all such chains generated by
1412 the noise source. This min-entropy value is particular to d -long chains and cannot be extrapolated
1413 linearly; i.e., chains of length wd will not necessarily have w times as much min-entropy present
1414 as a d -long chain. It may not be possible to know what a typical output length will be at the time
1415 of testing. Therefore, although not mathematically correct, in practice, calculating an entropy
1416 estimate per sample (extrapolated from that of the d -long chain) provides estimates that are close.

1417 The following algorithm uses output values as list indices. If the output set does not consist of
1418 consecutive values, then the values are adjusted before this algorithm is applied. This can be done
1419 without altering entropy estimates, as the data is categorical. For example, if the output set is $\{0,$
1420 $1, 4\}$, and the observed sequence is $(0, 0, 4, 1, 0, 4, 0, 1)$, 0 can stay the same, 1, can stay the same,
1421 but 4 must be changed to 2. The new set is $\{0, 1, 2\}$, and the new sequence is $(0, 0, 2, 1, 0, 2, 0,$
1422 $1)$.

1423 Given the input $S = (s_1, \dots, s_L)$, where $s_i \in A = \{x_1, \dots, x_k\}$,

- 1424 1. Define the confidence level to be $\alpha = \min(0.99^{k^2}, 0.99^d)$, where $d = 128$ is the assumed
1425 length of the chain.
- 1426 2. Estimate the initial probabilities for each output value. Let P be a list of length k . For i from
1427 1 to k :

$$1428 \quad P_i = \min \left\{ 1, \frac{o_i}{L} + \varepsilon \right\},$$

1429 where o_i denotes the number of times that value x_i has occurred in S , and ε is defined by:

⁷ This statement assumes that the output space is defined such that it contains all 2^6 (or more) possible outputs; if, however, the output space is defined to have 2^6 or less elements, regardless of the sample size, the test can accurately estimate the transition probabilities with the amount of data specified in Section 3.1.1.

$$\varepsilon = \sqrt{\frac{\log_2\left(\frac{1}{1-\alpha}\right)}{2L}}$$

1430

1431 3. Let $o_{s_L} = o_{s_L} - 1$. This step removes one from the count of the last symbol of the
1432 sequence, which is necessary to compute sample proportions in the next step.

1433

1434 4. Let T be a $k \times k$ matrix. Estimate the probabilities in the bounding matrix T ,
1435 overestimating where

1436

$$T_{i,j} = \begin{cases} 1 & \text{if } o_i = 0 \\ \min\left\{1, \frac{o_{i,j}}{o_i} + \varepsilon_i\right\} & \text{otherwise,} \end{cases}$$

1437

and $o_{i,j}$ is the number of transitions from state x_i to state x_j observed in the sample, and ε_i is
1438 defined to be

1438

$$\varepsilon_i = \sqrt{\frac{\log_2\left(\frac{1}{1-\alpha}\right)}{2o_i}}$$

1439

1440 5. Using the bounding matrix T , find the probability of the most likely sequence of outputs,
1441 \hat{p}_{max} , using a dynamic programming algorithm as follows:

1442

a. For j from 1 to $d - 1$:

1443

i. Let h be a list of length k .

1444

ii. Find the highest probability for a sequence of length $j+1$ ending in each
1445 sample value. For c from 1 to k :

1446

1. Let P' be a list of length k .

1447

2. For i from 1 to k :

1448

a. $P'_i = P_i \times T_{i,c}$

1449

3. $h_c = \max_{i=1..k}(P'_i)$

1450

iii. Store the highest probability for a sequence of length $j+1$ ending in each
1451 value in P . For all $i \in \{1, \dots, k\}$, set $P_i = h_i$.

1452

b. $\hat{p}_{max} = \max_{i=1..k}(P_i)$

1453

6. The min-entropy estimate is the negative logarithm of the probability of the most likely
1454 sequence of outputs, \hat{p}_{max} :

1455

$$\text{min-entropy} = -1/d \log_2(\hat{p}_{max}).$$

1456

Example: Suppose that $k = 3$, $L = 21$ and $S = (2, 2, 0, 1, 0, 2, 0, 1, 2, 1, 2, 0, 1, 2, 1, 0, 0, 1, 1, 0,$
1457 $0)$. In step 1, $\alpha = \min(0.99^9, 0.99^d) = 0.2762$. After step 2, $\varepsilon = 0.0877$, $P_1 = 0.4687$ $P_2 =$
1458 0.4211 , and $P_3 = 0.3734$. After step 4, the bounding matrix T has values:

1459

	0	1	2
0	0.4376	0.7233	0.2948
1	0.5805	0.2948	0.5805
2	0.6641	0.4974	0.3308

1460

1461 After step 5a, the loop iteration for $j=1$ completes, $P_1 = 0.2480$, $P_2 = 0.3390$, and $P_3 = 0.2444$.
 1462 This represents the most probable sequence of length two ending in $x_1=0$, $x_2=1$, and $x_3=2$,
 1463 respectively. After step 6, the highest probability of any chain of length 128 generated by this
 1464 bounding matrix is 1.7372×10^{-24} , yielding an estimated min-entropy of 0.6166.

1465 6.3.4 The Compression Estimate

1466 The compression estimate, proposed by Hagerty and Draper [HD12], computes the entropy rate of
 1467 a dataset, based on how much the dataset can be compressed. This estimator is based on the Maurer
 1468 Universal Statistic [Mau92]. The estimate is computed by generating a dictionary of values, and
 1469 then computing the average number of samples required to produce an output, based on the
 1470 dictionary. One advantage of using the Maurer statistic is that there is no assumption of
 1471 independence. When output with dependencies is tested with this statistic, the compression rate is
 1472 affected (and therefore the entropy), but an entropy estimate is still obtained. A calculation of the
 1473 Maurer statistic is efficient, as it requires only one pass through the dataset to provide an entropy
 1474 estimate.

1475 Given a dataset from the noise source, the samples are first partitioned into two disjoint groups.
 1476 The first group serves as the dictionary for the compression algorithm; the second group is used
 1477 as the test group. The compression values are calculated over the test group to determine the mean,
 1478 which is the Maurer statistic. Using the same method as the collision estimate, the probability
 1479 distribution that has the minimum possible entropy for the calculated Maurer statistic is
 1480 determined. For this distribution, the entropy per sample is calculated as the lower bound on the
 1481 entropy that is present.

1482 The following algorithm uses output values as list indices. If the output set does not consist of
 1483 consecutive values, then the values must be adjusted before this algorithm is applied. This can be
 1484 done without altering entropy estimates, as the data is categorical. For example, if the output set is
 1485 $\{0, 1, 4\}$, and the observed sequence is $(0, 0, 4, 1, 0, 4, 0, 1)$, 0 can stay the same, 1 can stay the
 1486 same, but 4 must be changed to 2. The new set is $\{0, 1, 2\}$, and the new sequence is $(0, 0, 2, 1, 0,$
 1487 $2, 0, 1)$.

1488 Given the input $S = (s_1, \dots, s_L)$, where $s_i \in A = \{x_1, \dots, x_k\}$,

- 1489 1. Partition the dataset into two disjoint groups. These two groups will form the dictionary
 1490 and the test data.
 - 1491 a. Create the dictionary from the first $d = 1000$ observations, (s_1, s_2, \dots, s_d) .
 - 1492 b. Use the remaining $v = L - d$ observations, (s_{d+1}, \dots, s_L) , for testing.

- 1493 2. Initialize the dictionary *dict* to an all zero array of size *k*. For *i* from 1 to *d*, let $dict[s_i] = i$.
 1494 *dict*[*s*_{*i*}] is the index of last occurrence of each *s*_{*i*} in the dictionary.
- 1495 3. Run the test data against the dictionary created in Step 2.
- 1496 a. Let D_i be a list of length *v*.
- 1497 b. For *i* from *d* + 1 to *L*:
- 1498 i. If *dict*[*s*_{*i*}] is non-zero, then $D_{i-d} = i - dict[s_i]$. Update the dictionary with the
 1499 index of the most recent observation, $dict[s_i] = i$.
- 1500 ii. If *dict*[*s*_{*i*}] is zero, add that value to the dictionary, i.e., $dict[s_i]=i$. Let $D_{i-d} =$
 1501 *i*.
- 1502 4. Let $b = \lceil \log_2(\max(x_1, \dots, x_k)) \rceil + 1$, the number of bits needed to represent the largest
 1503 symbol in the output alphabet. Calculate the sample mean, \bar{X} , and sample standard
 1504 deviation⁸, $\hat{\sigma}$, of $(\log_2(D_1), \dots, \log_2(D_v))$.

$$\bar{X} = \frac{\sum_{i=1}^v \log_2 D_i}{v},$$

$$c = 0.7 - \frac{0.8}{b} + \frac{\left(4 + \frac{32}{b}\right) v^{-3/b}}{15}$$

1508 and

$$\hat{\sigma} = c \sqrt{\frac{\sum_{i=1}^v (\log_2 D_i)^2}{v} - \bar{X}^2}.$$

- 1510 5. Compute the lower-bound of the confidence interval for the mean, based on a normal
 1511 distribution using

$$\bar{X}' = \bar{X} - \frac{2.576\hat{\sigma}}{\sqrt{v}}.$$

- 1513 6. Using a binary search, solve for the parameter *p*, such that the following equation is true:

$$\bar{X}' = G(p) + (n - 1)G(q),$$

1515 where

$$G(z) = \frac{1}{v} \sum_{t=d+1}^L \sum_{u=1}^t \log_2(u) F(z, t, u),$$

⁸ Note that a correction factor is applied to the standard deviation, as described in [Maurer]. This correction factor reduces the standard deviation to account for dependencies in the D_i values.

$$F(z, t, u) = \begin{cases} z^2(1-z)^{u-1} & \text{if } u < t \\ z(1-z)^{t-1} & \text{if } u = t \end{cases},$$

1518 and

$$q = \frac{1-p}{k-1}.$$

1520 7. If the binary search yields a solution, then the min-entropy is the negative logarithm of the
1521 parameter, p :

$$1522 \quad \text{min-entropy} = -\log_2(p).$$

1523 If the search does not yield a solution, then the min-entropy estimation is:

$$1524 \quad \text{min-entropy} = \log_2(k).$$

1525 *Example:* For illustrative purposes, suppose that $d = 10$ (instead of 1000), $k = 3$, $L = 21$ and $S = (2,$
1526 $2, 0, 1, 0, 2, 0, 1, 2, 1, 2, 0, 1, 2, 1, 0, 0, 1, 0, 0, 0)$. The dictionary sequence is $(2, 2, 0, 1, 0, 2, 0,$
1527 $1, 2, 1)$, and the testing sequence is $(2, 0, 1, 2, 1, 0, 0, 1, 0, 0, 0)$. $v = 11$. After the dictionary is
1528 initialized, $dict[0] = 7$, $dict[1] = 10$, and $dict[2] = 9$. In Step 4, b is calculated as 2. After processing
1529 the test sequence, $\bar{X} = 1.098$, $\hat{\sigma} = 0.2620$ and $\bar{X}' = 0.8944$. The value of p that solves the
1530 equation is 0.7003, and the min-entropy estimate is 0.5139.

1531 6.3.5 t -Tuple Estimate

1532 This method examines the frequency of t -tuples (pairs, triples, etc.) that appears in the input dataset
1533 and produces an estimate of the entropy per sample, based on the frequency of those t -tuples. The
1534 frequency of the t -tuple (x_1, x_2, \dots, x_t) in $S = (s_1, \dots, s_L)$ is the number of i 's such that $s_i = x_1, s_{i+1} =$
1535 $x_2, \dots, s_{i+t-1} = x_t$. It should be noted that the tuples can overlap.

1536 Given the input $S = (s_1, \dots, s_L)$, where $s_i \in A = \{x_1, \dots, x_k\}$,

- 1537 1. Find the largest t such that the number of occurrences of the most common t -tuple in S is
1538 at least 35.
- 1539 2. Let $Q[i]$ store the number of occurrences of the most common i -tuple in S for $i=1, \dots, t$. For
1540 example, in $S=(2, 2, 0, 1, 0, 2, 0, 1, 2, 1, 2, 0, 1, 2, 1, 0, 0, 1, 0, 0, 0)$, $Q[1] =$
1541 $\max(\#0\text{'s}, \#1\text{'s}, \#2\text{'s}) = \#0\text{'s} = 9$, and $Q[2] = 4$ is obtained by the number of 01's in S .
- 1542 3. For $i=1$ to t , an estimate for p_{max} is computed as
1543 a. Let $P[i]=Q[i] / (L-i+1)$, and compute an estimate on the maximum individual
1544 sample value probability as $P_{max}[i] = P[i]^{1/i}$.
- 1545 4. The entropy estimate is calculated as $-\log_2 \max (P_{max}[1], \dots, P_{max}[t])$.

1546 6.3.6 Longest Repeated Substring (LRS) Estimate

1547 This method estimates the collision entropy (sampling without replacement) of the source, based
1548 on the number of repeated substrings (tuples) within the input dataset. Although this method

1549 estimates collision entropy (an upper bound on min-entropy), this estimate handles tuple sizes that
1550 are too large for the t -tuple estimate, and is therefore a complementary estimate.

1551 Given the input $S = (s_1, \dots, s_L)$, where $s_i \in A = \{x_1, \dots, x_k\}$,

- 1552 1. Find the smallest u such that the number of occurrences of the most common u -tuple in S
1553 is less than 20.
- 1554 2. Find the largest v such that the number of occurrences of the most common v -tuple in S is
1555 at least 2 and the most common $(v+1)$ -tuple in S occurs once. In other words, v is the largest
1556 length that a tuple repeat occurs. If $v < u$, this estimate cannot be computed.
- 1557 3. For $W=u$ to v , compute the estimated W -tuple collision probability

$$1558 \quad P_W = \frac{\sum_i \binom{C_i}{2}}{\binom{L-W+1}{2}},$$

1559 where C_i is the number of occurrences of the i^{th} unique W -tuple.

- 1560 4. For each P_W , compute the estimated average collision probability per string symbol

$$1561 \quad P_{max,W} = P_W^{1/W}.$$

1562 The collision entropy estimate is calculated as $-\log_2 \max(P_{max,u}, \dots, P_{max,v})$.

1563 6.3.7 Multi Most Common in Window Prediction Estimate

1564 The Multi Most Common in Window (MultiMCW) predictor contains several subpredictors, each
1565 of which aims to guess the next output, based on the last w outputs. Each subpredictor predicts the
1566 value that occurs most often in that window of w previous outputs. The MultiMCW predictor keeps
1567 a scoreboard that records the number of times that each subpredictor was correct, and uses the
1568 subpredictor with the most correct predictions to predict the next value. In the event of a tie, the
1569 most common sample value that has appeared most recently is predicted. This predictor was
1570 designed for cases where the most common value changes over time, but still remains relatively
1571 stationary over reasonable lengths of the sequence.

1572 Given the input $S = (s_1, \dots, s_L)$, where $s_i \in A = \{x_1, \dots, x_k\}$,

- 1573 1. Let window sizes be $w_1=63$, $w_2=255$, $w_3=1023$, $w_4=4095$, and $N = L - w_1$. Let *correct* be
1574 an array of N Boolean values, each initialized to 0.
- 1575 2. Let *scoreboard* be a list of four counters, each initialized to 0. Let *frequent* be a list of four
1576 values, each initialized to *Null*. Let *winner* = 1.
- 1577 3. For $i = w_1 + 1$ to L :
1578 a. For $j = 1$ to 4,
1579 i. If $i > w_j$, let *frequent* _{j} be the most frequent value in $(s_{i-w_j}, s_{i-w_j+1}, \dots, s_{i-1})$. If
1580 there is a tie, then *frequent* _{j} is assigned to the most frequent value that has
1581 appeared most recently.
1582 ii. Else, let *frequent* _{j} = *Null*.

- 1583 b. Let $prediction = frequent_{winner}$.
- 1584 c. If ($prediction = s_i$), let $correct_{i-w_1} = 1$.
- 1585 d. Update the *scoreboard*. For $j = 1$ to 4,
- 1586 i. If ($frequent_j = s_i$)
- 1587 1. Let $scoreboard_j = scoreboard_j + 1$
- 1588 2. If $scoreboard_j \geq scoreboard_{winner}$, let $winner = j$
- 1589 4. Let C be the number of ones in *correct*.
- 1590 5. Calculate a 99% upper bound on the predictor's global performance $P_{global} = \frac{C}{N}$ as:

1591
$$P'_{global} = P_{global} + 2.576 \sqrt{\frac{P_{global}(1-P_{global})}{N-1}}$$

- 1592 6. Calculate the predictor's local performance, based on the longest run of correct predictions.
1593 Let r be one greater than the length of the longest run of ones in *correct*. Use a binary
1594 search to solve the following for P_{local} :

1595
$$0.99 = \frac{1 - P_{local}x}{(r + 1 - rx)q} \times \frac{1}{x^{N+1}},$$

1596 where $q = 1 - P_{local}$ and $x = x_{10}$, derived by iterating the recurrence relation

1597
$$x_j = 1 + qP_{local}^r x_{j-1}^{r+1}$$

1598 for j from 1 to 10, and $x_0 = 1$.

- 1599 7. The min-entropy is the negative logarithm of the greater performance metric

1600
$$min-entropy = -\log_2(\max(P'_{global}, P_{local})).$$

1601 *Example:* Suppose that $S = (1, 2, 1, 0, 2, 1, 1, 2, 2, 0, 0, 0)$, so that $L = 12$. For the purpose of the
1602 example, suppose that $w_1=3, w_2=5, w_3=7, w_4=9$ (instead of $w_1=63, w_2=255, w_3=1023, w_4=4095$).
1603 Then $N=9$. In step 3, the values are as follows:

i	<i>frequent</i>	<i>scoreboard</i> (step 3b)	<i>Winner</i> (step 3b)	<i>prediction</i>	s_i	$correct_{i-w_1}$	<i>scoreboard</i> (step 3d)
4	(1, --, --, --)	(0, 0, 0, 0)	1	1	0	0	(0, 0, 0, 0)
5	(0, --, --, --)	(0, 0, 0, 0)	1	0	2	0	(0, 0, 0, 0)
6	(2, 2, --, --)	(0, 0, 0, 0)	1	2	1	0	(0, 0, 0, 0)
7	(1, 1, --, --)	(0, 0, 0, 0)	1	1	1	1	(1, 1, 0, 0)
8	(1, 1, 1, --)	(1, 1, 0, 0)	2	1	2	0	(1, 1, 0, 0)
9	(1, 2, 2, --)	(1, 1, 0, 0)	2	2	2	1	(1, 2, 1, 0)
10	(2, 2, 2, 2)	(1, 2, 1, 0)	2	2	0	0	(1, 2, 1, 0)
11	(2, 2, 2, 2)	(1, 2, 1, 0)	2	2	0	0	(1, 2, 1, 0)
12	(0, 0, 2, 0)	(1, 2, 1, 0)	2	0	0	1	(2, 3, 1, 1)

1604

1605 After all of the predictions are made, $correct = (0, 0, 0, 1, 0, 1, 0, 0, 1)$. Then, $P_{global} = 0.3333$,
1606 $P'_{global} = 0.7627$, $P_{local} = 0.5$, and the resulting min-entropy estimate is 0.3909.

1607 6.3.8 The Lag Prediction Estimate

1608 The lag predictor contains several subpredictors, each of which predicts the next output, based on
1609 a specified lag. The lag predictor keeps a scoreboard that records the number of times that each
1610 subpredictor was correct, and uses the subpredictor with the most correct predictions to predict the
1611 next value.

1612 Given the input $S = (s_1, \dots, s_L)$, where $s_i \in A = \{x_1, \dots, x_k\}$,

1613 1. Let $D = 128$, and $N = L - 1$. Let lag be a list of D values, each initialized to *Null*. Let $correct$
1614 be a list of N Boolean values, each initialized to 0.

1615 1. Let $scoreboard$ be a list of D counters, each initialized to 0. Let $winner = 1$.

1616 2. For $i = 2$ to L :

1617 a. For $d = 1$ to D :

1618 i. If $(d < i)$, $lag_d = s_{i-d}$,

1619 ii. Else $lag_d = Null$.

1620 b. Let $prediction = lag_{winner}$.

1621 c. If $(prediction = s_i)$ let $correct_{i-1} = 1$.

1622 d. Update the $scoreboard$. For $d = 1$ to D :

1623 i. If $(lag_d = s_i)$

1624 1. Let $scoreboard_d = scoreboard_d + 1$.

1625 2. If $scoreboard_d \geq scoreboard_{winner}$, let $winner = d$.

1626 3. Let C be the number of ones in $correct$.

1627 4. Calculate a 99% upper bound on the predictor's global performance $P_{global} = \frac{C}{N}$ as:

$$1628 \quad P'_{global} = P_{global} + 2.576 \sqrt{\frac{P_{global}(1-P_{global})}{N-1}}.$$

1629 5. Calculate the predictor's local performance, based on the longest run of correct predictions.
1630 Let r be one greater than the length of the longest run of ones in $correct$. Use a binary
1631 search to solve the following for P_{local} :

$$1632 \quad 0.99 = \frac{1 - P_{local}x}{(r + 1 - rx)q} \times \frac{1}{x^{N+1}},$$

1633 where

$$1634 \quad q = 1 - P_{local}$$

1635 and $x = x_{10}$, derived by iterating the recurrence relation

$$1636 \quad x_j = 1 + qP_{local}^r x_{j-1}^{r+1}$$

1637 for j from 1 to 10, and $x_0=1$.

1638 6. The min-entropy is the negative logarithm of the greater performance metric

$$1639 \quad \text{min-entropy} = -\log_2(\max(P'_{global}, P_{local})).$$

1640 *Example:* Suppose that $S = (2, 1, 3, 2, 1, 3, 1, 3, 1, 2)$, so that $L = 10$ and $N = 9$. For the purpose of
1641 the example, suppose that $D = 3$ (instead of 128). The following table shows the values in step 3.

i	lag	$scoreboard$ (step 3b)	$Winner$ (step 3b)	$prediction$	s_i	$correct_{i-1}$	$scoreboard$ (step 3d)
2	(2, --, --)	(0, 0, 0)	1	2	1	0	(0, 0, 0)
3	(1, 2, --)	(0, 0, 0)	1	1	3	0	(0, 0, 0)
4	(3, 1, 2)	(0, 0, 0)	1	3	2	0	(0, 0, 1)
5	(2, 3, 1)	(0, 0, 1)	3	1	1	1	(0, 0, 2)
6	(1, 2, 3)	(0, 0, 2)	3	3	3	1	(0, 0, 3)
7	(3, 1, 2)	(0, 0, 3)	3	2	1	0	(0, 1, 3)
8	(1, 3, 1)	(0, 1, 3)	3	1	3	0	(0, 2, 3)
9	(3, 1, 3)	(0, 2, 3)	3	3	1	0	(0, 3, 3)
10	(1, 3, 1)	(0, 3, 3)	2	3	2	0	(0, 3, 3)

1642

1643 After all of the predictions are made, $correct = (0, 0, 0, 1, 1, 0, 0, 0, 0)$. Then, $P_{global} = 0.2222$,
1644 $P'_{global} = 0.6008$, $P_{local} = 0.6667$, and the resulting min-entropy estimate is 0.5850.

1645 6.3.9 The MultiMMC Prediction Estimate

1646 The MultiMMC predictor is composed of multiple Markov Model with Counting (MMC)
1647 subpredictors. Each MMC predictor records the observed frequencies for transitions from one
1648 output to a subsequent output (rather than the probability of a transition, as in a typical Markov
1649 model), and makes a prediction, based on the most frequently observed transition from the current
1650 output. MultiMMC contains D MMC subpredictors running in parallel, one for each depth from 1
1651 to D . For example, the MMC with depth 1 creates a first-order model, while the MMC with depth
1652 D creates a D^{th} -order model. MultiMMC keeps a scoreboard that records the number of times that
1653 each MMC subpredictor was correct, and uses the subpredictor with the most correct predictions
1654 to predict the next value.

1655 Given the input $S = (s_1, \dots, s_L)$, where $s_i \in A = \{x_1, \dots, x_k\}$,

- 1656 1. Let $D = 16$, and $N = L - 2$. Let $subpredict$ be a list of D values, each initialized to *Null*. Let
1657 $correct$ be an array of N Boolean values, each initialized to 0.
- 1658 2. For $d = 1$ to D , let M_d be a list of counters, where $M_d[x, y]$ denotes the number of observed
1659 transitions from output x to output y for the d^{th} -order MMC.

- 1660 3. Let *scoreboard* be a list of D counters, each initialized to 0. Let *winner* = 1.
- 1661 4. For $i=3$ to L :
- 1662 a. For $d = 1$ to D :
- 1663 i. If $d < i-1$, increment $MMC_d[(s_{i-d-1}, \dots, s_{i-2}), s_{i-1}]$ by 1.
- 1664 b. For $d = 1$ to D :
- 1665 i. Find the y value that corresponds to the highest $M_d[(s_{i-d}, \dots, s_{i-1}), y]$ value, and
- 1666 denote that y as y_{max} . Let $subpredict_d = y_{max}$. If all possible values of M_d
- 1667 $[(s_{i-d}, \dots, s_{i-1}), y]$ are 0, then let $subpredict_d = Null$.
- 1668 c. Let $prediction = subpredict_{winner}$.
- 1669 d. If $(prediction = s_i)$, let $correct_{i-2} = 1$.
- 1670 e. Update the scoreboard. For $d = 1$ to D :
- 1671 i. If $(subpredict_d = s_i)$
- 1672 1. Let $scoreboard_d = scoreboard_d + 1$.
- 1673 2. If $scoreboard_d \geq scoreboard_{winner}$, let $winner = d$.
- 1674 5. Let C be the number of ones in *correct*.
- 1675 6. Calculate a 99% upper bound on the predictor's global performance $P_{global} = \frac{C}{N}$ as:

1676
$$P'_{global} = P_{global} + 2.576 \sqrt{\frac{P_{global}(1-P_{global})}{N-1}}.$$

- 1677 7. Calculate the predictor's local performance, based on the longest run of correct predictions.
- 1678 Let r be one greater than the length of the longest run of ones in *correct*. Use a binary
- 1679 search to solve the following for P_{local} :

1680
$$0.99 = \frac{1 - P_{local}x}{(r + 1 - rx)q} \times \frac{1}{x^{N+1}},$$

1681 where

1682
$$q = 1 - P_{local}$$

1683 and $x = x_{10}$, derived by iterating the recurrence relation

1684
$$x_j = 1 + qP_{local}^r x_{j-1}^{r+1}$$

1685 for j from 1 to 10, and $x_0=1$.

- 1686 8. The min-entropy is the negative logarithm of the greater performance metric

1687
$$min-entropy = -\log_2(\max(P'_{global}, P_{local})).$$

1688 *Example:* Suppose that $S = (2, 1, 3, 2, 1, 3, 1, 3, 1)$, so that $L = 9$ and $N = 7$. For the purpose of

1689 example, further suppose that $D=3$ (instead of 16). After each iteration of step 4 is completed, the
1690 values are:

i	<i>subpredict</i>	<i>scoreboard</i> (step 4c)	<i>Winner</i> (step 4c)	<i>prediction</i>	s_i	<i>correct_{i-2}</i>	<i>scoreboard</i> (step 4e)
3	(Null, Null, Null)	(0, 0, 0)	1	Null	3	0	(0, 0, 0)
4	(Null, Null, Null)	(0, 0, 0)	1	Null	2	0	(0, 0, 0)
5	(1, Null, Null)	(0, 0, 0)	1	1	1	1	(1, 0, 0)
6	(3, 3, Null)	(1, 0, 0)	1	3	3	1	(2, 1, 0)
7	(2, 2, 2)	(2, 1, 0)	1	2	1	0	(2, 1, 0)
8	(3, Null, Null)	(2, 1, 0)	1	3	3	1	(3, 1, 0)
9	(2, 2, Null)	(3, 1, 0)	1	2	1	0	(3, 1, 0)

1691
1692 Let $\{x \rightarrow y : c\}$ denote a nonzero count c for the transition from x to y . Models M_1 , M_2 , and M_3 are
1693 shown below after step 4a (the model update step) for each value of i .

i	M_1	M_2	M_3
3	{2→1:1}	--	--
4	{1→3:1}, {2→1:1}	{(2, 1)→3:1}	--
5	{1→3:1}, {2→1:1}, {3→2:1}	{(1, 3)→2:1}, {(2, 1)→3:1}	{(2, 1, 3)→2:1}
6	{1→3:1}, {2→1:2}, {3→2:1}	{(1, 3)→2:1}, {(2, 1)→3:1}, {(3, 2)→1:1}	{(1, 3, 2)→1:1}, {(2, 1, 3)→2:1}
7	{1→3:2}, {2→1:2}, {3→2:1}	{(1, 3)→2:1}, {(2, 1)→3:2}, {(3, 2)→1:1}	{(1, 3, 2)→1:1}, {(2, 1, 3)→2:1}, {(3, 2, 1)→3:1}
8	{1→3:2}, {2→1:2}, {3→1:1}, {3→2:1}	{(1, 3)→1:1}, {(1, 3)→2:1}, {(2, 1)→3:2}, {(3, 2)→1:1}	{(1, 3, 2)→1:1}, {(2, 1, 3)→1:1}, {(2, 1, 3)→2:1}, {(3, 2, 1)→3:1}
9	{1→3:3}, {2→1:2}, {3→1:1}, {3→2:1}	{(1, 3)→1:1}, {(1, 3)→2:1}, {(2, 1)→3:2}, {(3, 1)→3:1}, {(3, 2)→1:1}	{(1, 3, 1)→3:1}, {(1, 3, 2)→1:1}, {(2, 1, 3)→1:1}, {(2, 1, 3)→2:1}, {(3, 2, 1)→3:1}

1694
1695 After the predictions are all made, $correct = (0, 0, 1, 1, 0, 1, 0)$. Then, $P_{global} = 0.4286$, $P'_{global} =$
1696 0.9490 , $P_{local} = 0.6667$, and the resulting min-entropy estimate is 0.0755.

1697 **6.3.10 The LZ78Y Prediction Estimate**

1698 The LZ78Y predictor is loosely based on LZ78 encoding with the Bernstein's Yabba scheme
1699 [Sal07] for adding strings to the dictionary. The predictor keeps a dictionary of strings that have
1700 been added to the dictionary so far, and continues adding new strings to the dictionary until the
1701 dictionary has reached its maximum capacity. Each time that a sample is processed, every
1702 substring in the most recent B samples updates the dictionary or is added to the dictionary.

1703 Given the input $S = (s_1, \dots, s_L)$, where $s_i \in A = \{x_1, \dots, x_k\}$,

1704 1. Let $B = 16$, and $N = L - B - 1$. Let *correct* be an array of N Boolean values, each initialized
1705 to 0. Let $maxDictionarySize = 65536$.

1706 2. Let D be an empty dictionary. Let $dictionarySize = 0$.

1707 3. For $i=B+2$ to L :

1708 a. For $j=B$ down to 1:

1709 i. If $(s_{i-j-1}, \dots, s_{i-2})$ is not in D , and $dictionarySize < maxDictionarySize$:

1710 1. Let $D[s_{i-j-1}, \dots, s_{i-2}]$ be added to the dictionary.

1711 2. Let $D[s_{i-j-1}, \dots, s_{i-2}][s_{i-1}] = 0$.

1712 3. $dictionarySize = dictionarySize + 1$

1713 ii. If $(s_{i-j-1}, \dots, s_{i-2})$ is in D ,

1714 1. Let $D[s_{i-j-1}, \dots, s_{i-2}][s_{i-1}] = D[s_{i-j-1}, \dots, s_{i-2}][s_{i-1}] + 1$.

1715 b. Use the dictionary to predict the next value, s_i . Let $prediction = Null$, and let
1716 $maxcount = 0$. For $j = B$ down to 1:

1717 i. Let $prev = (s_{i-j}, \dots, s_{i-1})$.

1718 ii. If $prev$ is in the dictionary, find the $y \in \{x_1, \dots, x_k\}$ that has the highest
1719 $D[prev][y]$ value.

1720 iii. If $D[prev][y] > maxcount$:

1721 1. $prediction = y$.

1722 2. $maxcount = D[prev][y]$.

1723 c. If $(prediction = s_i)$, let $correct_{i-B-1} = 1$.

1724 4. Let C be the number of ones in *correct*. Calculate a 99% upper bound on the predictor's
1725 global performance $P_{global} = \frac{C}{N}$ as:

1726
$$P'_{global} = P_{global} + 2.576 \sqrt{\frac{P_{global}(1-P_{global})}{N-1}}$$

1727 5. Calculate the predictor's local performance, based on the longest run of correct predictions.
1728 Let r be one greater than the length of the longest run of ones in *correct*. Use a binary
1729 search to solve the following for P_{local} :

1730

$$0.99 = \frac{1 - P_{local}x}{(r + 1 - rx)q} \times \frac{1}{x^{N+1}},$$

1731

where $q = 1 - P_{local}$ and $x = x_{10}$, derived by iterating the recurrence relation

1732

$$x_j = 1 + qP_{local}^r x_{j-1}^{r+1}$$

1733

for j from 1 to 10, and $x_0=1$.

1734

6. The min-entropy is the negative logarithm of the greater performance metric

1735

$$min-entropy = -\log_2 \left(\max(P'_{global}, P_{local}) \right).$$

1736

Example: Suppose that $S = (2, 1, 3, 2, 1, 3, 1, 3, 1, 2, 1, 3, 2)$, and $L=13$. For the purpose of example, suppose that $B=4$ (instead of 16), then $N=8$.

1737

i	Add to D	$prev$	Max $D[prev]$ entry	$prediction$	s_i	$correct_{i-B-1}$
6	$D[2, 1, 3, 2][1]$	(1, 3, 2, 1)	Null	Null	3	0
	$D[1, 3, 2][1]$	(3, 2, 1)	Null			
	$D[3, 2][1]$	(2, 1)	Null			
	$D[2][1]$	(1)	Null			
7	$D[1, 3, 2, 1][3]$	(3, 2, 1, 3)	Null	Null	1	0
	$D[3, 2, 1][3]$	(2, 1, 3)	Null			
	$D[2, 1][3]$	(1, 3)	Null			
	$D[1][3]$	(3)	Null			
8	$D[3, 2, 1, 3][1]$	(2, 1, 3, 1)	Null	3	3	1
	$D[2, 1, 3][1]$	(1, 3, 1)	Null			
	$D[1, 3][1]$	(3, 1)	Null			
	$D[3][1]$	(1)	3			
9	$D[2, 1, 3, 1][3]$	(1, 3, 1, 3)	Null	1	1	1
	$D[1, 3, 1][3]$	(3, 1, 3)	Null			
	$D[3, 1][3]$	(1, 3)	1			
	$D[1][3]$	(3)	1			
10	$D[1, 3, 1, 3][1]$	(3, 1, 3, 1)	Null	3	2	0
	$D[3, 1, 3][1]$	(1, 3, 1)	3			
	$D[1, 3][1]$	(3, 1)	3			
	$D[3][1]$	(1)	3			
11	$D[3, 1, 3, 1][2]$	(1, 3, 1, 2)	Null	1	1	1
	$D[1, 3, 1][2]$	(3, 1, 2)	Null			
	$D[3, 1][2]$	(1, 2)	Null			
	$D[1][2]$	(2)	1			
12	$D[1, 3, 1, 2][1]$	(3, 1, 2, 1)	Null	3	3	1
	$D[3, 1, 2][1]$	(1, 2, 1)	Null			
	$D[1, 2][1]$	(2, 1)	3			
	$D[2][1]$	(1)	3			

13	$D[3, 1, 2, 1][3]$	(1, 2, 1, 3)	<i>Null</i>	1	2	0
	$D[1, 2, 1][3]$	(2, 1, 3)	1			
	$D[2, 1][3]$	(1, 3)	1			
	$D[1][3]$	(3)	1			

1738

1739 After the predictions are all made, $correct = (0, 0, 1, 1, 0, 1, 1, 0)$. Then, $P_{global} = 0.5$, $P'_{global} =$
 1740 0.9868 , $P_{local} = 0.6667$, and the resulting min-entropy estimate is 0.0191.

1741 6.4 Reducing the Sample Space

1742 It is often the case that the data requirements for a test on noise source samples depends on the
 1743 number of possible different bitstrings from the noise source (i.e., the size of the alphabet A). For
 1744 example, consider two different noise sources. The first source outputs four-bit samples, and thus
 1745 has a possible total of $2^4 = 16$ different outputs, and the second source outputs 32-bit samples, for
 1746 a possible total of 2^{32} different outputs.

1747 In many cases, the variability in the output that contributes to the entropy in a sample may be
 1748 concentrated among some portion of the bits in the sample. For example, consider a noise source
 1749 that outputs 32-bit high-precision clock samples that represent the time it takes to perform some
 1750 system process. Suppose that the bits in a sample are ordered in the conventional way, so that the
 1751 lower-order bits of the sample correspond to the higher resolution measurements of the clock. It is
 1752 easy to imagine that in this case, the low-order bits would contain most of the variability. In fact,
 1753 it would seem likely that some of the high-order bits may be constantly 0. For this example, it
 1754 would be reasonable to truncate the 32-bit sample to a four-bit string by taking the lower four bits,
 1755 and perform the tests on the four-bit strings. Of course, it must be noted that in this case, only a
 1756 maximum of four bits of min-entropy per sample could be credited to the noise source.

1757 The description below provides a method for mapping the n -bit samples, collected as specified in
 1758 Section 3.1.1, to m -bit samples, where $n \geq m$. The resulting strings can be used as input to tests
 1759 that may have infeasible data requirements if the mapping were not performed. Note that after the
 1760 mapping is performed, the maximum amount of entropy possible per n -bit sample is m bits.

1761 Given a noise source that produces n -bit samples, where n exceeds the bit-length that can be
 1762 handled by the test, the submitter **shall** provide the tester with an ordered ranking of the bits in the
 1763 n -bit samples (see Section 3.2.2). The rank of '1' corresponds to the bit assumed to be contributing
 1764 the most entropy to the sample, and the rank of n corresponds to the bit contributing the least
 1765 amount. If multiple bits contribute the same amount of entropy, the ranks can be assigned
 1766 arbitrarily among those bits. The following algorithm, or its equivalent, is used to assign ranks.

1767 **Input:** A noise source and corresponding statistical model with samples of the form $X = a_1a_2\dots a_n$,
 1768 where each a_i is a bit.

1769 **Output:** An ordered ranking of the bits a_1 through a_n , based on the amount of entropy that each
 1770 bit is assumed to contribute to the noise source outputs.

1771 1. Set $M = \{a_1, a_2, \dots, a_n\}$.

- 1772 2. For $i = 1$ to n :
- 1773 a. Choose an output bit a from M such that no other bit in S is assumed to contribute
- 1774 more entropy to the noise source samples than a .
- 1775 b. Set the rank of a to i .
- 1776 c. Remove a from M .
- 1777 Given the ranking, n -bit samples are mapped to m -bit samples by simply taking the m -bits of
- 1778 greatest rank in order (i.e., bit 1 of the m -bit string is the bit from an n -bit sample with rank 1, bit
- 1779 2 of the m -bit string is the bit from an n -bit sample with rank 2, ... and bit m of the m -bit string is
- 1780 the bit from an n -bit sample with rank m).
- 1781 Note that for the estimators in Section 6, a reference to a sample in the dataset will be interpreted
- 1782 as a reference to the m -bit subsets of the sample when the test necessitates processing the dataset
- 1783 as specified in this section.
- 1784

1785 **Appendix A—Acronyms**

1786 Selected acronyms and abbreviations used in this paper are defined below.

AES	Advanced Encryption Standard
ANS	American National Standard
CAVP	Cryptographic Algorithm Validation Program
CMVP	Cryptographic Module Validation Program
DRBG	Deterministic Random Bit Generator
FIPS	Federal Information Processing Standard
HMAC	Hash-based Message Authentication Code
IID	Independent and Identically Distributed
LRS	Longest Repeated Substring
NIST	National Institute of Standards and Technology
NRBG	Non-deterministic Random Bit Generator
NVLAP	National Voluntary Laboratory Accreditation Program
RBG	Random Bit Generator
SP	NIST Special Publication

1787

1788

Appendix B—Glossary

<i>Alphabet</i>	A finite set of two or more symbols.
<i>Alphabet size</i>	See sample size.
<i>Algorithm</i>	A clearly specified mathematical process for computation; a set of rules that, if followed, will give a prescribed result.
<i>Approved</i>	FIPS-approved or NIST-Recommended.
<i>Array</i>	A fixed-length data structure that stores a collection of elements, where each element is identified by its integer index.
<i>Assessment (of entropy)</i>	An evaluation of the amount of entropy provided by a (digitized) noise source and/or the entropy source that employs it.
<i>Biased</i>	A value that is chosen from a sample space is said to be biased if one value is more likely to be chosen than another value. (Contrast with unbiased.)
<i>Binary data (from a noise source)</i>	Digitized and possibly post-processed output from a noise source that consists of a single bit; that is, each sampled output value is represented as either 0 or 1.
<i>Bitstring</i>	A bitstring is an ordered sequence of 0's and 1's. The leftmost bit is the most significant bit.
<i>Collision</i>	An instance of duplicate sample values occurring in a dataset.
<i>Conditioning (of noise source output)</i>	A method of processing the raw data to reduce bias and/or ensure that the entropy rate of the conditioned output is no less than some specified amount.
<i>Confidence interval</i>	An interval, [<i>low</i> , <i>high</i>], where the true value of a parameter <i>p</i> falls within that interval with a stated probability. E.g., a 95% confidence interval about an estimate for <i>p</i> yields values for <i>low</i> and <i>high</i> such that $low \leq p \leq high$ with probability 0.95.
<i>Continuous test</i>	A type of health test performed within an entropy source on the output of its noise source in order to gain some level of assurance that the noise source is working correctly, prior to producing each output from the entropy source.
<i>Consuming application (for an RBG)</i>	An application that uses the output from an approved random bit generator.
<i>Dataset</i>	A sequence of sample values. (See Sample.)

<i>Deterministic Random Bit Generator (DRBG)</i>	An RBG that includes a DRBG mechanism and (at least initially) has access to a source of entropy input. The DRBG produces a sequence of bits from a secret initial value called a seed, along with other possible inputs. A DRBG is often called a Pseudorandom Number (or Bit) Generator.
<i>Developer</i>	The party that develops the entire entropy source or the noise source.
<i>Dictionary</i>	A dynamic-length data structure that stores a collection of elements or values, where a unique label identifies each element. The label can be any data type.
<i>Digitization</i>	The process of generating bits from the noise source.
<i>DRBG mechanism</i>	The portion of an RBG that includes the functions necessary to instantiate and uninstantiate the RBG, generate pseudorandom bits, (optionally) reseed the RBG and test the health of the DRBG mechanism. Approved DRBG mechanisms are specified in SP 800-90A.
<i>Entropy</i>	A measure of the disorder, randomness or variability in a closed system. Min-entropy is the measure used in this Recommendation.
<i>Entropy rate</i>	The rate at which a digitized noise source (or entropy source) provides entropy; it is computed as the assessed amount of entropy provided by a bitstring output from the source, divided by the total number of bits in the bitstring (yielding assessed bits of entropy per output bit). This will be a value between zero (no entropy) and one.
<i>Entropy source</i>	The combination of a noise source, health tests, and an optional conditioning component that produce random bitstrings to be used by an RBG.
<i>Estimate</i>	The estimated value of a parameter, as computed using an estimator.
<i>Estimator</i>	A technique for estimating the value of a parameter.
<i>False positive</i>	An erroneous acceptance of the hypothesis that a statistically significant event has been observed. This is also referred to as a type 1 error. When “health-testing” the components of a device, it often refers to a declaration that a component has malfunctioned – based on some statistical test(s) – despite the fact that the component was actually working correctly.

<i>Health testing</i>	Testing within an implementation immediately prior to or during normal operation to determine that the implementation continues to perform as implemented and as validated.
<i>Independent</i>	Two random variables X and Y are independent if they do not convey information about each other. Receiving information about X does not change the assessment of the probability distribution of Y (and vice versa).
<i>Independent and Identically Distributed (IID)</i>	A sequence of random variables for which each element of the sequence has the same probability distribution as the other values, and all values are mutually independent.
<i>List</i>	A dynamic-length data structure that stores a sequence of values, where each value is identified by its integer index.
<i>Markov model</i>	A model for a probability distribution where the probability that the i^{th} element of a sequence has a given value depends only on the values of the previous n elements of the sequence. The model is called an n^{th} order Markov model.
<i>Min-entropy</i>	The min-entropy (in bits) of a random variable X is the largest value m having the property that each observation of X provides at least m bits of information (i.e., the min-entropy of X is the greatest lower bound for the information content of potential observations of X). The min-entropy of a random variable is a lower bound on its entropy. The precise formulation for min-entropy is $(\log_2 \max p_i)$ for a discrete distribution having probabilities p_1, \dots, p_k . Min-entropy is often used as a worst-case measure of the unpredictability of a random variable.
<i>Narrowest internal width</i>	The maximum amount of information from the input that can affect the output. For example, if $f(x) = \text{SHA-1}(x) \parallel 01$, and x consists of a string of 1000 binary bits, then the narrowest internal width of $f(x)$ is 160 bits (the SHA-1 output length), and the output width of $f(x)$ is 162 bits (the 160 bits from the SHA-1 operation, concatenated by 01).
<i>Noise source</i>	The component of an entropy source that contains the non-deterministic, entropy-producing activity. (e.g., thermal noise or hard drive seek times)
<i>Non-deterministic Random Bit Generator (NRBG)</i>	An RBG that has access to an entropy source and (when working properly) produces outputs that have full entropy (see SP 800-90C). Also called a <i>true random bit</i> (or <i>number</i>) <i>generator</i> . (Contrast with a DRBG)

<i>On-demand test</i>	A type of health test that is available to be run whenever a user or a relying component requests it.
<i>Output space</i>	The set of all possible distinct bitstrings that may be obtained as samples from a digitized noise source.
<i>P-value</i>	The probability that the chosen test statistic will assume values that are equal to or more extreme than the observed test statistic value, assuming that the null hypothesis is true.
<i>Predictor</i>	A function that predicts the next value in a sequence, based on previously observed values in the sequence.
<i>Probability distribution</i>	A function that assigns a probability to each measurable subset of the possible outcomes of a random variable.
<i>Probability model</i>	A mathematical representation of a random phenomenon.
<i>Pseudorandom</i>	A deterministic process (or data produced by such a process) whose output values are effectively indistinguishable from those of a random process, as long as the internal states and internal actions of the process are unknown. For cryptographic purposes, “effectively indistinguishable” means “not within the computational limits established by the intended security strength.”
<i>Random</i>	A non-deterministic process (or data produced by such a process) whose output values follow some probability distribution. The term is sometimes (mis)used to imply that the probability distribution is uniform, but no uniformity assumption is made in this Recommendation.
<i>Random Bit Generator (RBG)</i>	A device or algorithm that outputs a random sequence that is effectively indistinguishable from statistically independent and unbiased bits. An RBG is classified as either a DRBG or an NRBG.
<i>Raw data</i>	Digitized and possibly post-processed output of the noise source.
<i>Run (of output sequences)</i>	A sequence of identical values.
<i>Sample</i>	An observation of the raw data. Common examples of output values obtained by sampling are single bits, single bytes, etc. (The term “sample” is often extended to denote a sequence of such observations; this Recommendation will refrain from that practice.)

<i>Sample size</i>	The number of possible distinct values that a sample can have. May also be called <i>alphabet size</i> .
<i>Security boundary</i>	A conceptual boundary that is used to assess the amount of entropy provided by the values output from an entropy source. The entropy assessment is performed under the assumption that any observer (including any adversary) is outside of that boundary.
<i>Seed</i>	A bitstring that is used as input to (initialize) an algorithm. In this Recommendation, the algorithm using a seed is a DRBG. The entropy provided by the seed must be sufficient to support the intended security strength of the DRBG.
<i>Sequence</i>	An ordered list of quantities.
<i>Shall</i>	The term used to indicate a requirement that needs to be fulfilled to claim conformance to this Recommendation. Note that shall may be coupled with not to become shall not .
<i>Should</i>	The term used to indicate an important recommendation. Ignoring the recommendation could result in undesirable results. Note that should may be coupled with not to become should not .
<i>Start-up testing</i>	A suite of health tests that are performed every time the entropy source is initialized or powered up. These tests are carried out on the noise source before any output is released from the entropy source.
<i>Submitter</i>	The party that submits the entire entropy source and output from its components for validation. The submitter can be any entity that can provide validation information as required by this Recommendation (e.g., developer, designer, vendor or any organization).
<i>Testing laboratory</i>	An accredited cryptographic security testing laboratory
<i>Unbiased</i>	A value that is chosen from a sample space is said to be unbiased if all potential values have the same probability of being chosen. (Contrast with biased.)

1789
1790

1791

Appendix C—References

- [BZ2] BZIP2 Compression Algorithm. <http://www.bzip.org/>.
- [Cac97] C. Cachin, *Entropy Measures and Unconditional Security in Cryptography*, PhD Thesis, Reprint as vol.1 of ETH Series in Information Security and Cryptography, ISBN 3-89649-185-7, Hartung-Gorre Verlag, Konstanz, ETH Zurich, 1997.
- [Fel50] W. Feller, *An Introduction to Probability Theory and its Applications*, volume one, chapter 13, John Wiley and Sons, Inc., 1950.
- [FIPS140] Federal Information Processing Standard 140-2, *Security Requirements for Cryptographic Modules*, May 25, 2001.
- [FIPS180] Federal Information Processing Standard 180-4, *Secure Hash Standard (SHS)*, August 2015.
- [FIPS197] Federal Information Processing Standard 197, *Specification for the Advanced Encryption Standard (AES)*, November 2001.
- [FIPS198] Federal Information Processing Standard 198-1, *The Keyed-Hash Message Authentication Code (HMAC)*, July 2008.
- [FIPS202] Federal Information Processing Standard 202, *SHA-3 Standard: Permutation-Based Hash and Extendable-Output Functions*, August 2015.
- [HD12] P. Hagerty and T. Draper, *Entropy Bounds and Statistical Tests*, NIST Random Bit Generation Workshop, December 2012, http://csrc.nist.gov/groups/ST/rbg_workshop_2012/hagerty_entropy_paper.pdf. (Presentation available at http://csrc.nist.gov/groups/ST/rbg_workshop_2012/hagerty.pdf.)
- [IG140-2] National Institute of Standards and Technology Special Publication, Communications Security Establishment Canada, *Implementation Guidance for FIPS PUB 140-2 and the Cryptographic Module Validation Program*, September 15, 2015
- [Kel15] J. Kelsey, Kerry A. McKay, M. Sonmez Turan, *Predictive Models for Min-Entropy Estimation*, Proceedings of the Workshop on Cryptographic Hardware and Embedded Systems 2015 (CHES 2015), France
- [Mau92] U. Maurer, *A Universal Statistical Test for Random Bit Generators*, Journal of Cryptology, Vol. 5, No. 2, 1992, pp. 89-105.

- [Sal07] D. Salomon, *Data Compression: The Complete Reference*, Chapter 3, Springer, 2007
- [Shan51] C.E Shannon, *Prediction and Entropy of Printed English*, Bell System Technical Journal, volume 30, pp. 50-64, January 1951, <https://archive.org/details/bstj30-1-50>.
- [SP800-38B] National Institute of Standards and Technology Special Publication (SP) 800-38B *Recommendations for Block Cipher Modes of Operation: The CMAC Mode for Authentication*, May 2005
- [SP800-57] National Institute of Standards and Technology Special Publication (SP) 800-57 *Recommendation for Key Management – Part I: General (Revision 3)*, January 2016.
- [SP800-90A] National Institute of Standards and Technology Special Publication (SP) 800-90A, *Recommendation for Random Number Generation Using Deterministic Random Bit Generators*, June 2015.
- [SP800-90C] National Institute of Standards and Technology Special Publication (SP) 800-90C, *Recommendations for Random Bit Generator (RBG) Constructions*, Draft.
- [SP800-107] National Institute of Standards and Technology Special Publication (SP) 800-107, Revision 1, *Recommendations for Applications using Approved Hash Functions*, August 2012.

1792

Appendix D— Min-Entropy and Optimum Guessing Attack Cost

1793 Suppose that an adversary wants to determine at least one of several secret values, where each
 1794 secret value is independently chosen from a set of M possibilities, with probability distribution P
 1795 $= \{p_1, p_2, \dots, p_M\}$. Assume that these probabilities are sorted so that $p_1 \geq p_2 \geq \dots \geq p_M$. Consider a
 1796 guessing strategy aimed at successfully guessing as many secret values as possible. The adversary's
 1797 goal would be to minimize the expected number of guesses per successful recovery. Such a strategy
 1798 would consist of guessing a maximum of k possibilities for a given secret value, moving on to a
 1799 new secret value when either a guess is correct or k incorrect guesses for the current value have
 1800 been made. In general, the optimum value of k can be anywhere in the range $1 \leq k \leq M$, depending
 1801 on the probability distribution P . Note that when $k = M$, the M^{th} guess is considered a valid (though
 1802 trivial) guess. Regardless of the value of k chosen, it is clear that the k guesses selected for a given
 1803 secret value should be the k most likely possible values, in decreasing order of probability.

1804 The expected work per success can be computed for this attack as follows. For $1 \leq j \leq k - 1$, the
 1805 attacker will make exactly j guesses if the secret value is the j^{th} most likely value, an event having
 1806 probability p_j . The attacker will make exactly k guesses if the secret value is not any of the $k - 1$
 1807 most likely values, an event having probability $1 - \sum_{j=1}^{k-1} p_j$. Thus, the expected number of guesses
 1808 for the attack is given by the following:

$$1809 \quad p_1 + 2p_2 + \cdots + (k-1)p_{k-1} + k \left(1 - \sum_{j=1}^{k-1} p_j \right).$$

1810 Since this attack will be successful if and only if the secret value is one of the k most likely
 1811 possibilities, which is the case with probability $\sum_{j=1}^k p_j$, the expected number of times the attack
 1812 must be performed until the first success is the reciprocal of this probability. Multiplying this
 1813 reciprocal by the expected number of guesses per attack gives the following as the expected work
 1814 per success:

$$1815 \quad W_k(P) = \frac{p_1 + 2p_2 + \cdots + (k-1)p_{k-1} + k \left(1 - \sum_{j=1}^{k-1} p_j \right)}{\sum_{j=1}^k p_j}.$$

1816 It is not critical to determine the value k^* that minimizes $W_k(P)$, since the min-entropy of P leads
 1817 to an accurate approximation (and sometimes the exact value) of $W_{k^*}(P)$. Stated more precisely,
 1818 $W_1(P) = \frac{1}{p_1}$ is an upper bound of $W_{k^*}(P)$, and it can be shown that $W_k(P) \geq \frac{1}{2p_1} + \frac{1}{2}$ for all k
 1819 such that $1 \leq k \leq M$. Since the min-entropy of P is $-\log_2(p_1)$, these two bounds imply that the
 1820 error between the min-entropy of P and $\log_2(W_{k^*}(P))$ can be bounded as follows:

$$1821 \quad 0 \leq -\log_2 p_1 - \log_2(W_{k^*}(P)) \leq 1 - \log_2(p_1 + 1).$$

1822 Notice that since $\frac{1}{M} \leq p_1 \leq 1$, the upper bound on the error approaches 0 as $p_1 \rightarrow 1$, and
 1823 alternatively, this bound approaches 1 as $M \rightarrow \infty$ and $p_1 \rightarrow \frac{1}{M}$. In other words, the min-entropy of
 1824 P either corresponds to the exact expected work, measured in bits, needed to perform the optimum
 1825 guessing attack or over-estimates this work by at most one bit.

1826 In order to prove the claim that $W_k(P) \geq \frac{1}{2p_1} + \frac{1}{2}$, for $1 \leq k \leq M$, rewrite the expected work per
 1827 success as

$$1828 \quad W_k(P) = \frac{1 + (1 - p_1) + (1 - p_1 - p_2) + \cdots + (1 - p_1 - p_2 - \cdots - p_{k-1})}{p_1 + p_2 + \cdots + p_k}.$$

1829 Consider an alternative probability distribution on a set of M possibilities $P' =$
 1830 $\{p_1, p_1, \dots, p_1, r, 0, \dots, 0\}$, where p_1 occurs $t = \left\lfloor \frac{1}{p_1} \right\rfloor$ times and $r = 1 - tp_1$. It is straightforward to
 1831 see that $W_k(P) \geq W_k(P')$, since each term in the numerator of $W_k(P)$ is at least as large as the
 1832 corresponding term in $W_k(P')$, and the denominator of $W_k(P')$ is at least as large as the
 1833 denominator of $W_k(P)$.

1834 Now to show that $W_k(P') \geq \frac{1}{2p_1} + \frac{1}{2}$. Based on the above formula for $W_k(P)$, for $1 \leq k \leq t + 1$,
 1835 the numerator of $W_k(P')$ can be written as

$$1836 \quad \sum_{i=0}^{k-1} (1 - ip_1) = k - \frac{k(k-1)}{2} p_1 = kp_1 \left(\frac{1}{p_1} - \frac{k-1}{2} \right).$$

1837 Consider the following two cases where $1 \leq k \leq t$ and $k = t + 1$. These are the only cases to check,
1838 since if $M > t + 1$, then $W_k(P') = W_{t+1}(P')$ for $k > t + 1$, because the remaining probabilities are
1839 all zero. Furthermore, $r = 0$ if and only if $\frac{1}{p_1}$ is an integer, and when this happens, only the first
1840 case needs to be addressed since $W_{t+1}(P') = W_t(P')$.

1841 For $1 \leq k \leq t$, the denominator of $W_k(P') = kp_1$. Then,

$$1842 \quad W_k(P') = \frac{kp_1 \left(\frac{1}{p_1} - \frac{k-1}{2} \right)}{kp_1} = \frac{1}{p_1} - \frac{k-1}{2},$$

$$1843 \quad \geq \frac{1}{p_1} - \frac{1}{2} \left(\left\lfloor \frac{1}{p_1} \right\rfloor - 1 \right),$$

$$1844 \quad \geq \frac{1}{p_1} - \frac{1}{2} \left(\frac{1}{p_1} - 1 \right),$$

$$1845 \quad \geq \frac{1}{2p_1} + \frac{1}{2}.$$

1846 For $k = t + 1$, the denominator of $W_k(P')$ is $tp_1 + r = 1$. Let $x = \frac{1}{p_1} - \left\lfloor \frac{1}{p_1} \right\rfloor$, so $0 \leq x < 1$. This implies

1847

$$1848 \quad W_k(P') = kp_1 \left(\frac{1}{p_1} - \frac{k-1}{2} \right) = \left(\left\lfloor \frac{1}{p_1} \right\rfloor + 1 \right) p_1 \left(\frac{1}{p_1} - \frac{1}{2} \left\lfloor \frac{1}{p_1} \right\rfloor \right),$$

$$1849 \quad = \left(\frac{1}{p_1} - x + 1 \right) \left(\frac{1}{2} + \frac{p_1 x}{2} \right),$$

$$1850 \quad = \frac{1}{2p_1} + \frac{1}{2} + \frac{p_1 x(1-x)}{2},$$

$$1851 \quad \geq \frac{1}{2p_1} + \frac{1}{2}.$$

1852 Therefore, it has been shown that $W_k(P) \geq W_k(P') \geq \frac{1}{2p_1} + \frac{1}{2}$ for $1 \leq k \leq M$. Note that this lower
1853 bound is sharp, since $W_k(P)$ achieves this value when P is a uniform distribution.

1854 **Appendix E—Post-processing Functions**

1855 This section provides the details of the allowed post-processing functions for a noise source.

1856

1857 *Von Neumann's method:* This method produces independent unbiased random bits for a source
1858 that generates independent biased output bits. This method divides the sequence into pairs and
1859 applies the following mapping:

<i>input</i>	<i>output</i>
00	discard
01	1
10	0
11	discard

1860 For a source that produces independent biased random bits (s_1, s_2, \dots) , with $\Pr(s_i = 0) = p$, and $p \neq$
 1861 $\frac{1}{2}$, the method extracts approximately $np(1 - p)$ unbiased bits from n biased bits. Independent of
 1862 the value of p , the method throws away a pair of bits at least half of the time. It should be noted
 1863 that the bias in the correlated sources might increase after applying the technique.

1864 *Linear filtering method:* This method divides the sequence into non-overlapping blocks of w bits
 1865 and applies a linear function to each block. Mathematically, the output of the j^{th} block is calculated
 1866 as $f(s_{jw+1}, \dots, s_{(j+1)w}) = c_1 s_{jw+1} + \dots + c_w s_{(j+1)w}$, where the c_i values are predetermined binary constants.
 1867 A typical value of w may be between 16 and 64; this Recommendation does not put a restriction
 1868 on the selection of the block size w .

1869 *Length of runs method:* This method outputs the length of the runs in (s_1, s_2, \dots) , where the s_i 's are
 1870 bits.

1871

1872 Appendix F— The Narrowest Internal Width

1873 The narrowest internal width of a conditioning component is the maximum amount of information
 1874 from the input that can affect the output. It can also be considered as the logarithm of an upper
 1875 bound on the number of distinct outputs, based on the size of the internal state.

1876 *Example:* Let $F(X)$ be a function defined as follows:

- 1877 1. Let h_1 be the output of SHA256(X) truncated to 64 bits.
- 1878 2. Return SHA256($h_1 || h_1$) truncated to 128 bits.

1879 This function takes an arbitrarily-long input X and will yield 128-bit output value, but its internal
 1880 width is only 64 bits, because the value of the output only depends on the value of 64-bit h_1 .

1881 Appendix G—CBC-MAC Specification

1882 A conditioning component may be based on the use of CBC-MAC using a 128-bit **approved**
 1883 block-cipher algorithm. This CBC-MAC construction **shall not** be used for any other purpose than
 1884 as the algorithm for a conditioning component, as specified in Section 3.1.5.1.1. The following
 1885 notation is used for the construction.

1886 Let $\mathbf{E}(\text{Key}, \text{input_string})$ represent the **approved** encryption algorithm, with a *Key* and an
 1887 *input_string* as input parameters. The length of the *input_string* **shall** be an integer multiple of the
 1888 output length n of the block-cipher algorithm and **shall** always be the same length (i.e., variable
 1889 length strings **shall not** be used as input).

1890 Let n be the length (in bits) of the output block of the **approved** block cipher algorithm, and let w

1891 be the number of n -bit blocks in the *input_string*.

1892 Let *output_string* be the n -bit output of CBC-MAC.

1893 **CBC-MAC:**

1894 **Input:** bitstring *Key*, *input_string*.

1895 **Output:** bitstring *output_string*.

1896 **Process:**

1897 1. Let s_0, s_1, \dots, s_{w-1} be the sequence of blocks formed by dividing *input string* into n -bit
1898 blocks; i.e., each s_i consists of n bits.

1899 2. $V = 0$.

1900 3. For $i = 0$ to $w-1$

1901 $V = \mathbf{E}(\text{Key}, V \oplus s_i)$.

1902 4. Output V as the CBC-MAC output.

1903

1904 **Appendix H—Different Strategies for Entropy Estimation**

1905 Each of the estimation methods presented in Section 6 follows one of two approaches to estimating
1906 min-entropy. The first approach is based on entropic statistics, first described for IID data in
1907 [HD12], and later applied to non-IID data [HD12]. The most common value test estimates entropy
1908 by bounding the probability of the most-common output. In the IID case, the collision and
1909 compression estimators in Section 6.3 provide a lower bound on min-entropy by fitting the
1910 distribution to a near-uniform distribution, where one probability is highest, and the rest are all
1911 equal. Empirically, these estimators appear to be conservative for independent, but not necessarily
1912 identically distributed samples, as well. The final estimator proposed in [HD12] and specified in
1913 Section 6.3.3 constructs a first-order Markov model and estimates entropy from the most-probable
1914 sequence.

1915 **H.1 Entropic Statistics**

1916 The entropic statistics presented in [HD12], each designed to compute a different statistic on the
1917 samples, provide information about the structure of the data: collision, collection, compression,
1918 and Markov. While the estimators (except for the Markov) were originally designed for application
1919 to independent outputs, the tests have performed well when applied to data with dependencies.
1920 Given empirical evidence and the confidence level of the tests, their application to non-IID data
1921 will produce valid, although conservative, entropy estimates.

1922 The estimators assume that a probability distribution describes the output of a random noise source,
1923 but that the probability distribution is unknown. The goal of each estimator is to reveal information
1924 about the unknown distribution, based on a statistical measurement.

1925 The collision and compression estimators in Section 6 each solve an equation for an unknown
1926 parameter, where the equation is different for each estimator. These equations come from the target
1927 statistic's expected value using a near-uniform distribution, which provides a lower bound for min-
1928 entropy. A near-uniform distribution is an instance of a one-parameter family of probability
1929 distributions parameterized by p , P_p :

1930

$$P_p(i) = \begin{cases} p, & \text{if } i = 0 \\ \frac{1-p}{k-1}, & \text{otherwise} \end{cases}$$

1931 where k is the number of states in the output space, and $p \leq \frac{1-p}{k-1}$. In other words, one output state
1932 has the maximum probability, and the remaining output states are equally likely. For more
1933 information, see [HD12].

1934 H.1.1 Approximation of $F(1/z)$

1935 The function $F(1/z)$, used by the collision estimate (Section 6.3.2), can be approximated by the
1936 following continued fraction⁹:

1937

$$z + \frac{1}{z + \frac{-n}{1 + \frac{1}{z + \frac{1-n}{1 + \frac{2}{z + \frac{2-n}{1 + \frac{3}{z + \frac{3}{1 + \dots}}}}}}}}}$$

1938 H.2 Predictors

1939 Shannon first published the relationship between the entropy and predictability of a sequence in
1940 1951 [Shan51]. Predictors construct models from previous observations, which are used to predict
1941 the next value in a sequence. The prediction-based estimation methods in this Recommendation
1942 work in a similar way, but attempt to find bounds on the min-entropy of integer sequences
1943 generated by an unknown process (rather than N -gram entropy of English text, as in [Shan51]).

1944 The predictor approach uses two metrics to produce an estimate. The first metric is based on the
1945 global performance of the predictor, called *accuracy* in machine-learning literature. Essentially, a
1946 predictor captures the proportion of guesses that were correct. This approximates how well one
1947 can expect a predictor to guess the next output from a noise source, based on the results over a
1948 long sequence of guesses. The second metric is based on the greatest number of correct predictions
1949 in a row, which is called the local entropy estimate. This metric is useful for detecting cases where
1950 a noise source falls into a highly predictable state for some time, but the predictor may not perform
1951 well on long sequences. The calculations for the local entropy estimate come from the probability
1952 theory of runs and recurrent events [Fel50]. For more information about min-entropy estimation
1953 using predictors, see [Kel15].

⁹ Derived from Equation 8.9.2 at <http://dlmf.nist.gov/8.9>.