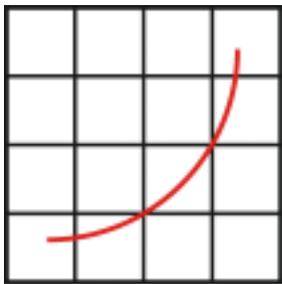


SPECpower Committee



spec[®]

Server Efficiency Rating Tool
(SERT)[™]

Design Document
3rd public draft

Standard Performance Evaluation Corporation

Table of Contents

1	Introduction	6
1.1	Summary	6
1.2	About SPEC	6
1.2.1	SPEC Membership	6
1.2.2	SPEC's General Development Guidelines	7
1.3	The EPA's ENERGY STAR for Computer Server Specification and SPEC	7
1.4	SERT's Differences from Conventional Benchmarks	7
1.5	Design Feedback Mechanism.....	8
1.6	Logistics	8
1.7	Trademark.....	8
2	Scope and Goals	9
2.1	Overview Summary.....	9
2.2	Sockets and Nodes.....	9
2.3	Scaling	9
2.4	Server Options and Expansion capabilities	10
2.5	IO Component.....	10
2.5.1	Storage IO.....	10
2.5.2	Network IO.....	10
2.6	Redundancy.....	11
2.7	Run Time.....	11
2.8	Platforms	11
2.8.1	Tested as Shipped.....	11
2.9	Implementation Languages.....	11
2.10	Load Levels.....	11
2.11	Worklets	12
2.12	Workload.....	12
2.13	Tentative Test Schedule	13
2.14	Schedule tradeoffs	13
3	SERT Architecture	14
3.1	System Overview	14
3.2	Execution of SERT.....	15
4	Worklet Execution Phases	16
5	Worklet Design Guidelines	19
5.1	Active Idle Worklet	19
5.2	CPU Worklet	19
5.3	Memory Worklet.....	19
5.4	Network IO Worklet.....	20
5.5	Storage IO Worklet	20
5.6	System Worklet.....	20
6	Power and Temperature Measurements.....	21
6.1	Environmental Conditions	21
6.2	Temperature Sensor Specifications.....	21
6.3	Power Analyzer Requirements	21
6.4	SPEC PTDaemon	22

6.5	Supported and Compliant Devices	22
6.6	Power Analyzer Setup	22
6.7	DC Line-Voltage.....	22
7	Metric/Score, Reporting, Logging	23
7.1	Metric/Score	23
7.2	Reporting and Output Files	23
7.2.1	Report 1: "Summary Report"	23
7.2.2	Report 2: "Power and Performance Data Sheet"	24
7.3	Validation / Verification	24
7.4	Logging	25
8	Future Enhancements / Stretch goals	26
8.1	Graphical User Interface (GUI)	26
8.2	Test Software	26
9	SERT and EPA ENERGY STAR for Server Version 2.0	27
9.1	Measurement	27
9.2	SERT Binaries and Recompilation.....	27
9.3	Manual Intervention	27
9.4	Fair Use of SERT information	27
9.4.1	Fair Use Rules	27
9.5	Accredited, Independent laboratory	27
9.6	Supply Voltage tolerance	28
10	Worklet Candidates	29
10.1	CPU Worklet: Compress.....	30
10.1.1	General Description	30
10.1.2	Sequence Execution Methods.....	30
10.1.3	Metric	30
10.1.4	Required Initialization	30
10.1.5	Configuration Parameters.....	30
10.1.6	Transaction Code	30
10.2	CPU Worklet: CryptoAES	31
10.2.1	General Description	31
10.2.2	Sequence Execution Methods.....	31
10.2.3	Metric	31
10.2.4	Required Initialization	31
10.2.5	Configuration Parameters.....	31
10.2.6	Transaction Code	31
10.3	CPU Worklet: FFT.....	32
10.3.1	General Description	32
10.3.2	Sequence Execution Methods.....	32
10.3.3	Metric	32
10.3.4	Required Initialization	32
10.3.5	Configuration Parameters.....	32
10.3.6	Transaction Code	32
10.4	CPU Workload: LU.....	33
10.4.1	General Description	33
10.4.2	Sequence Execution Methods.....	33

10.4.3	Metric	33
10.4.4	Required Initialization	33
10.4.5	Configuration Parameters.....	33
10.4.6	Transaction Code	33
10.5	CPU Workload: SOR.....	34
10.5.1	General Description	34
10.5.2	Sequence Execution Methods	34
10.5.3	Metric	34
10.5.4	Required Initialization	34
10.5.5	Configuration Parameters.....	34
10.5.6	Transaction Code	34
10.6	CPU Workload: XmlValidate	35
10.6.1	General Description	35
10.6.2	Sequence Execution Methods	35
10.6.3	Metric	35
10.6.4	Required Initialization	35
10.6.5	Configuration Parameters.....	35
10.6.6	Transaction Code	35
10.7	Memory Worklet: Flood.....	36
10.7.1	General Description	36
10.7.2	Sequence Execution Methods	36
10.7.3	Metric	36
10.7.4	Required Initialization	37
10.7.5	Configuration Parameters.....	37
10.7.6	Transaction Code	37
10.8	Memory Workload: XmlValidate.....	38
10.8.1	General Description	38
10.8.2	Sequence Execution Methods	38
10.8.3	Metric	38
10.8.4	Required Initialization	38
10.8.5	Configuration Parameters.....	38
10.8.6	Transaction Code	39
10.9	Storage IO Workload	40
10.9.1	General Description	40
10.9.2	Sequence Execution Methods	40
10.9.3	Metric	40
10.9.4	Required Initialization	40
10.9.5	Configuration Parameters.....	40
10.9.6	Transaction – Code 1 - RandomRead.....	41
10.9.7	Transaction – Code 1 - RandomWrite.....	41
10.9.8	Transaction – Code 2 – SequentialRead.....	41
10.9.9	Transaction – Code 2 – SequentialWrite.....	41
10.10	System Worklet: CSSJ.....	42
10.10.1	General Description	42
10.10.2	Sequence Execution Methods	42
10.10.3	Metric	42

10.10.4 Required Initialization	42
10.10.5 Configuration Parameters	42
10.10.6 New Order Transaction.....	42
10.10.7 Payment Transaction.....	43
10.10.8 Order Status Transaction.....	44
10.10.9 Delivery Transaction	44
10.10.10 Stock Level Transaction	44
10.10.11 Customer Report Transaction.....	44

1

1 Introduction

1.1 Summary

The EPA's ENERGY STAR development team is currently working on Version 2.0 of their Computer Server Specification¹. Version 2.0 aims to evolve the program by adding a means to measure the overall efficiency of the server while it is performing actual computing work via an Active Mode Efficiency Rating Tool.

The SPECpower committee is currently working on the design, implementation and delivery of the Server Efficiency Rating Tool (SERT)TM, a next generation tool set that will measure and evaluate the energy efficiency of computer servers. This public draft outlines the design of SERT for review by EPA stakeholders and their associates.

Please visit http://www.spec.org/sert/docs/SERT-Design_Doc.pdf for the latest updates.

1.2 About SPEC

The Standard Performance Evaluation Corporation (SPEC) was formed by the industry in 1988 to establish industry standards for measuring compute performance. SPEC has since become the largest and most influential benchmark consortium world-wide. Its mission is to ensure that the marketplace has a fair and useful set of metrics to analyze the newest generation of IT equipment.

The SPEC community has developed more than 30 industry-standard benchmarks for system performance evaluation in a variety of application areas and provided thousands of benchmark licenses to companies, resource centers, and educational institutions globally. Organizations using these benchmarks have published more than 20,000 peer-reviewed performance reports on SPEC's website (<http://www.spec.org/results.html>).

SPEC has a long history of designing, developing, and releasing industry-standard computer system performance benchmarks in a range of industry segments, plus peer-reviewing the results of benchmark runs. Performance benchmarking and the necessary work to develop and release new benchmarks can lead to disagreements among participants. Therefore, SPEC has developed an operating philosophy and range of normative behaviors that encourage cooperation and fairness amongst diverse and competitive organizations.

The increasing demand for energy-efficient IT Equipment has resulted in the need for power and performance benchmarks. In response, the SPEC community established SPECpower, an initiative to augment existing industry standard benchmarks with a power/energy measurement. Leading engineers and scientists in the fields of benchmark development and energy efficiency made a commitment to tackle this task. The development of the first industry-standard benchmark that measures the power and performance characteristics of server-class compute equipment started on January 26th 2006. In December of 2007, SPECpower_ssj2008 was released, which exercises the CPUs, caches, memory hierarchy and the scalability of shared memory processors on multiple load-levels. The benchmark runs on a wide variety of operating systems and hardware architectures. In version 1.10, which was released on April 15th 2009, SPEC augmented SPECpower_ssj2008 with multi-node support (e.g., blade-support).

1.2.1 SPEC Membership

SPEC membership is open to any interested company or entity. OSG members and associates are entitled to licensed copies of all released OSG benchmarks and unlimited publication of results on SPEC's public website. An initiation fee and annual fees are due for members. Nonprofit organizations and educational institutions have a reduced annual fee structure. Further details on membership information can be found on <http://www.spec.org/osg/joining.html> or requested at info@spec.org. Also a current list of SPEC members can be found here: <http://www.spec.org/spec/membership.html>.

¹ US Environmental Protection Agency – Energy Star Program Requirements for Computer Servers.
http://www.energystar.gov/index.cfm?c=revisions.computer_servers

51 1.2.2 SPEC's General Development Guidelines

52 SPEC's philosophy and standards of participation are the basis for the development of SERT. The tool
53 is being developed cooperatively by a committee representing diverse and competitive companies.
54 The following guides the committee in the development of a tool that will be useful and widely adopted
55 by the industry:

- 56 • Decisions are reached by consensus. Motions require a qualified majority to carry.
- 57 • Decisions are based on reality. Experimental results carry more weight than opinions. Data
58 and demonstration overrule assertion.
- 59 • Fair benchmarks allow competition among all industry participants in a transparent market.
- 60 • Tools and benchmarks should be architecture-neutral and portable.
- 61 • All who are willing to contribute may participate. Wide availability of results on the range of
62 available solutions allows the end user to determine the appropriate IT equipment.

63 Similar guidelines have resulted in the success and wide use of SPEC benchmarks in the performance
64 and power/performance industry and are essential to the success of SERT.

65

66 1.3 The EPA's ENERGY STAR for Computer Server Specification and SPEC

67 SPEC applauds the EPA for its goal to drive toward greater energy efficiency in IT Equipment, and
68 SPEC considers the EPA ENERGY STAR Program an industry partner in this effort. The
69 development of an Active Mode Efficiency Rating Tool is an essential component in the ongoing effort
70 to reduce worldwide energy consumption and paves the way for a successful ENERGY STAR for
71 Computer Servers program that has the potential to harmonize energy efficiency programs worldwide.

72 SPEC welcomes this opportunity to work with the EPA on SERT in support of the ENERGY STAR
73 Specification for Computer Server and is proudly looking forward to continuing our long-standing
74 association with the EPA ENERGY STAR development team.

75

76 1.4 SERT's Differences from Conventional Benchmarks

77 Performance benchmarks and energy efficiency benchmarks tend to focus on capabilities of computer
78 servers in specific business models or application areas. SERT is focused on providing a first order of
79 approximation² of energy efficiency across a broad range of application environments.

- 80 • The absolute score is less relevant for the end user, because it will not reflect specific
81 application capabilities.
- 82 • A rating tool that provides a pass-fail or a [Level 1/Level 2/Level 3] pass-fail rating is a better fit
83 for EPA's ENERGY STAR Environment for Computer Servers than a typical benchmark result
84 with multiple digits of precision in the metric.
- 85 • Marketing of the absolute scores will be disallowed, in order to encourage more participation
86 in the program

87 Benchmarks tend to focus on optimal conditions, including tuning options to customize the
88 configuration and software to the application of the benchmark business model. The need to achieve
89 competitive benchmark results often causes significant investment in the benchmark process. SERT is
90 designed to be more economical and easier to use, requiring minimal equipment and skills through:

- 91 • Highly automated processes and leveraging existing SPEC methods
- 92 • Focus on as-shipped default settings for the server
- 93 • Free from super-tuning

94 Where a benchmark represents a fixed reference point, ENERGY STAR programs are designed to
95 foster continuous improvement, with thresholds for success rising as the industry progresses. SERT
96 will be designed to match this paradigm, including:

- 97 • Quick adoption of new computing technologies
- 98 • Rapid turn-around for SERT version updates

² Andrew Fanara, Evan Haines, Arthur Howard

http://www.energystar.gov/ia/partners/prod_development/downloads/State_of_Energy_and_Performance_Benchmarking_for_Enterprise_Servers_Final.pdf

99 **1.5 Design Feedback Mechanism**

100 The SERT development team will evaluate input from a broad spectrum of industry experts during the
101 entire development process by utilizing its partnership with the EPA ENERGY STAR Program. The
102 team will collaborate on workload, metric and all other requirements of the EPA's Version 2.0
103 Framework.

104 Please provide your detailed feedback to the EPA via servers@energystar.gov. The EPA will collect,
105 sort, anonymize, and prioritize your feedback and pass it on to the SPEC development team.

106

107 **1.6 Logistics**

108 The licensee and price structure as well as the support and maintenance models that will be used for
109 SERT is work in progress.

110

111 **1.7 Trademark**

112 Product and service names mentioned herein may be the trademarks of their respective owners.

113

114

115 2 Scope and Goals

116 The current scope of Version 2.0 ENERGY STAR for Computer Servers includes servers with 1-4
117 processor sockets with a stated goal to expand to include blade technologies of similar scope. A
118 design goal of SERT is to accommodate these and larger technologies.

119 Among the issues involved with support of larger systems are the overall capacity of the system to
120 complete work, and the ability to design a workload that scales with the inclusion of additional
121 processors, memory, network interface cards, disk drives, etc. Different workload characteristics are
122 required to demonstrate effectiveness for each of these components. Providing a workload that fairly
123 represents their presence while not unfairly representing their absence is a challenge. These issues
124 are more prevalent with larger systems that have more expansion capabilities than smaller servers.

125 For these areas where it is concluded that the tool does not adequately represent the value of a
126 component compared to its power requirements, the tool will be designed to accommodate the
127 inclusion of “configuration power/performance modifiers”. A design goal is to automatically include this
128 additional information in the computation of the ENERGY STAR qualification results, including detailed
129 documentation that this was done.

130

131 2.1 Overview Summary

132 The following table summarizes some of the design goals that SERT will and will not provide.

IS	IS NOT
Rating Tool for overall energy efficiency	A Benchmark nor a Capacity Planning Tool
Measuring tool for power, performance and inlet-temperature	Measuring tool for Airflow, Air pressure, outlet-temperature
General compute-environment measure	Specific application benchmark measure
Support of AC- powered servers	Support of DC-powered servers
Used in single OS instance per server environments	Intended to stress virtualization hypervisor technology ³
ENERGY STAR Rating Tool	Marketing Tool
Planned to be architecture and OS neutral	Planned to be implemented on architecture and/or OS environments where insufficient resource has been volunteered to accomplish development, testing, and support.

133

134 2.2 Sockets and Nodes

135 SERT 1.0.0.0 is designed to be scalable and will be tested up to a maximum of 8 sockets and a
136 maximum of 64 nodes (limited to a set of homogenous servers or blade servers). The server under
137 test (SUT) may be a single stand-alone server or a multi-node set of servers. A multi-node SUT will
138 consist of server nodes that cannot run independent of shared infrastructure such as a backplane,
139 power-supplies, fans or other elements. These shared infrastructure systems are commonly known as
140 “blade servers” or “multi-node server”. Only identical servers are allowed in a multi-node SUT
141 configuration.

142

143 2.3 Scaling

144 Since the server efficiency rating of a given server is the primary objective of SERT, one of the main
145 design goals for SERT is to be able to scale the performance on the system in proportion to the
146 system configuration. As more components (processors, memory, and disk storage) are added to the
147 server, the workloads should utilize the additional resources so that the resultant performance is
148 higher when compared to the performance on the same server with a lesser configuration. Similarly,
149 for a given server, when the components are upgraded with faster counterparts, the performance
150 should scale accordingly. This is a very important aspect of the tool since adding and upgrading
151 components typically increases the total power consumed by the server which will affect the overall
152 efficiency result of the server. Creating a tool that scales performance based on the number/speed of
153 CPUs is most readily achievable – for the other components, the complexity of implementing such a
154 tool increases substantially.

³ Virtualization can be an important tool for saving energy. In a first-order approximation tool, such as SERT, the impacts of virtualized environments can be determined by examining the results at higher load levels.

155 While SERT will be designed to scale performance with additional hardware resources of the SUT, if
156 there are performance bottlenecks in system components unrelated to the added hardware the SUT
157 itself may not be able to sustain higher performance. In such cases the addition of components to the
158 SUT will normally result in higher power consumption without a commensurate increase in
159 performance. It is also possible that the workload mix that is defined for smaller systems will not scale
160 well when examining larger systems.

161

162 **2.4 Server Options and Expansion capabilities**

163 A server may have many optional features that are designed to increase the breadth of applications.
164 These features not only absorb additional power, but also require more capacity in the power supplies
165 and cooling system. Some SERT workload components will be designed to demonstrate the
166 enhanced capabilities that these features provide. However, while the tool needs to credit these
167 capabilities for the expanded workloads that they will accommodate, it cannot penalize efficient
168 servers that are not designed with substantial expansion options. A balance must be struck between
169 providing enhanced ratings for enhanced configurations and avoiding easy qualification of servers by
170 simply adding features that may not be needed in all situations.

171 SERT's goal is to avoid unnecessarily penalizing servers that are designed for low expandability, while
172 crediting servers with greater expandability. For example a configuration with four I/O adapters in PCI
173 slots may execute the workload of the tool more effectively than a configuration with only one such
174 adapter. On the other hand it may only run the workload of the tool as effectively as a configuration
175 with two network adapters. Because the configuration with four adapters may run some real workloads
176 more effectively than configurations with only two adapters, the EPA may elect to allow for some form
177 of "configuration modifier" to provide credit for the power infrastructure needed to support the
178 additional PCI slots.

179 The tool will be designed and tested to ensure that, should "configuration power-performance modifier"
180 credits be included, the tool will accommodate them.

181

182 **2.5 IO Component**

183 Disk and Network IO components are strongly desired to provide a better-rounded picture of system
184 performance and power than a CPU-centric test. SPEC is in the early stages of evaluating IO
185 workloads for SERT, so this section provides many discussion points but not necessarily conclusions.

186 SPEC recognizes that some of the items in the next two sections may not be reasonable or practical to
187 test or measure in a meaningful way. In those cases we would suggest the use of "configuration
188 power-performance modifiers" to compensate for the extra power draw associated with extra
189 functionality. Other items under consideration include:

- 190 • Different types/quantities of IO for different server categories
- 191 • Self-calibrating performance measurements for the disk and network subsystem

192

193 **2.5.1 Storage IO**

194 Ideally the Storage IO component of SERT would give credit for:

- 195 • Higher performance storage subsystems
- 196 • Larger capacity storage subsystems
- 197 • Reliability and availability features (RAID, battery backed cache, etc)

198

199 **2.5.2 Network IO**

200 Ideally the network IO component of SERT would give credit for:

- 201 • Higher performance Network Interfaces
- 202 • Larger transfer speed Network Interfaces
- 203 • Reliability and availability features

204

205 2.6 Redundancy

206 Many servers have redundancy built in for power supplies and cooling fans. Some servers include
 207 different levels of redundancy for memory, disk, and even processors. A design goal is to include
 208 accommodation for redundant components, although no specific tests are planned for energy
 209 measurement under fault tolerant conditions when one of a redundant set of components is disabled.

210

211 2.7 Run Time

212 The right balance between high repeatability of the results, high sub-system coverage and low
 213 resource allocation is desirable. The run time will depend on the agreed set of worklets. The target
 214 run time is around 3 hours.

215

216 2.8 Platforms

217 SERT 1.0.0.0 will be implemented for and is planned to be tested on the following platform/OS/JVM
 218 combinations (64 bit only), pending resources. In some cases, SPEC recommend the use of more
 219 than one JVM, where more than one JVM is generally available and selecting one may unfairly
 220 penalize a specific processor architecture or operating system.

221

HW Platform	x86 AMD	x86 AMD	x86 AMD	x86 Intel	x86 Intel	x86 Intel	Itanium Intel	POWER IBM	POWER IBM	POWER IBM	SPARC Oracle	SPARC Fujitsu
OS	Windows Server 2008 R2	LINUX	Solaris	Windows Server 2008 R2	LINUX	Solaris	HP-UX 11i	AIX	IBM i	LINUX	Solaris	Solaris
JVM	IBM j9 Oracle HS	IBM j9 Oracle HS	Oracle	IBM j9 Oracle HS	IBM j9 Oracle HS	Oracle	HP HS	IBM - j9	IBM - j9	IBM - j9	Oracle	Oracle

222

223 Note: OS refers to versions (service pack and patch levels) that are current at the SERT release.

224

225 Platform/OS/JVM combinations currently not on the list have no resources allocated to them. If support
 226 for additional architectures or OSs is desired, then active participation from requesting entities is
 227 mandatory. The inclusion of a JVM is dependent on an agreement from the JVM provider for
 228 unrestricted use of their JVM for SERT. Companies dedicating additional resources to the SPECpower
 229 committee for development of SERT would relax the schedule constraints.

230

231 2.8.1 Tested as Shipped

232 To provide results that are representative of a customer environment, the goal is to test systems in an
 233 "as-shipped" state. No super tuning would be allowed, but rather a limited list of valid parameter
 234 changes for configuration and typical optimization be permitted. Other changes will cause the run to be
 235 marked as noncompliant. SERT will launch the JVM within the tool, to restrict additional tuning.

236

237 The list of allowable parameters will be included in a future version of this document and in the
 238 operational documentation of the tool. This list would be agreed with the EPA before SERT release,
 239 and would be clearly documented as part of the SERT Run Rules.

240

241 2.9 Implementation Languages

242 The main body of code is written in Java in order to lower the burden of cross-platform support.
 243 Regardless, the framework is designed to accommodate other language implementations as well.

244

245 2.10 Load Levels

246 Multiple load levels are a desired goal of SERT and the design will include support for multiple levels.
 247 The active idle load level as well as a 100% workload level (not max power) are already good
 248 candidates. Prototype testing will show which levels will be included and if any weighting will be
 249 necessary.

250

251 **2.11 Worklets**

252 Developing the workload in the traditional SPEC way based on real world applications would result in
 253 complex test environments and high run times, especially for the IO intensive workloads, e.g. many
 254 client systems would be required for network IO and large disk sub systems for storage IO. The
 255 resulting costs for running such tests could be prohibitive for a rating tool. Therefore the SERT
 256 workload will be a collection of synthetic worklets for a variety of different load scenarios.

257

258 **2.12 Workload**

259 The existing SPEC benchmarks are mainly based on tailored versions of real world applications
 260 representing a typical workload for one application area or a synthetic workload derived from the
 261 analysis of existing server implementations. These benchmarks are suitable to evaluate different sub-
 262 areas of the overall server performance or efficiency if power measurements are included. They are
 263 not designed to give a representative assessment of the overall server performance or efficiency.

264 The design goal for the SERT suite however is to include all major aspects of server architecture, thus
 265 avoiding any preference for specific architectural features which might make a server look good under
 266 one workload and show disadvantages with another workload.

267 The SERT workload will instead take advantage of different server capabilities by using various load
 268 patterns, which are intended to stress all major components of a server uniformly.

269 If some components cannot be stressed adequately by the respective load pattern this can be
 270 compensated by adjusting the threshold for these components, e.g. increasing the power allowance
 271 for additional components which are not used by the load pattern.

272 It is highly unlikely that a single workload can be designed which achieves the goals outlined above,
 273 especially given the time constraints of the schedule targeted for ENERGY STAR for Servers Version
 274 2.0 by the EPA. Therefore the SERT workload will consist of several different worklets each stressing
 275 specific capabilities of a server. This approach furthermore supports generating individual efficiency
 276 scores for the server components besides the overall system score.

277 Figure 1 describes the general structure of the SERT test suite and its components.

278

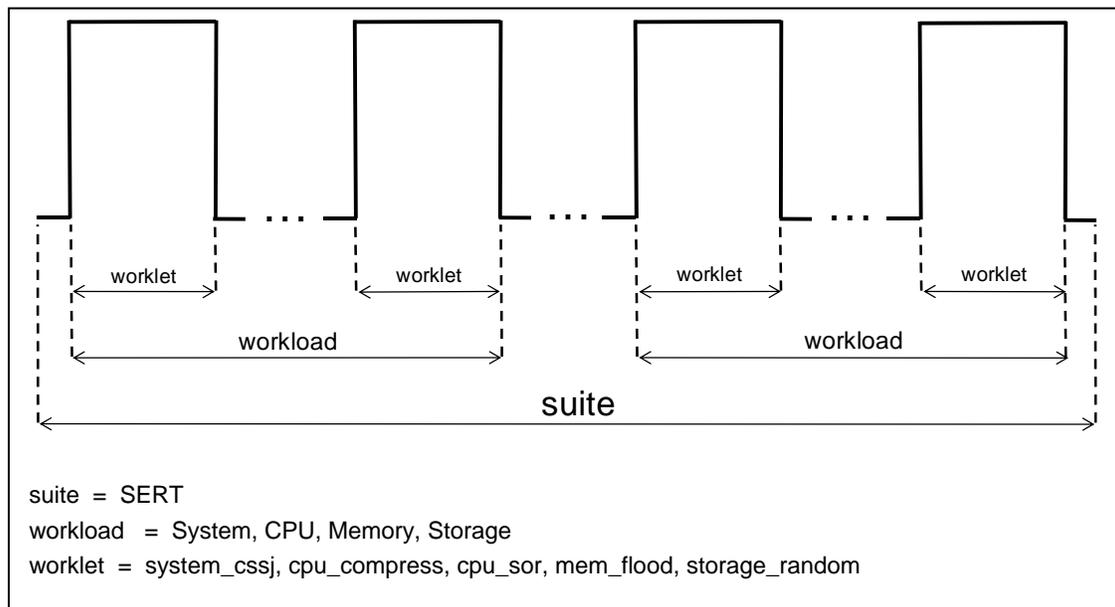


Figure 1: SERT Suite Components

281

282

283 **2.13 Tentative Test Schedule**

284 The alpha test phase is planned to start in March 2011 and the start of each phase requires
 285 successful completion of its predecessor. An estimated schedule can be created once we have
 286 decided on all design details.

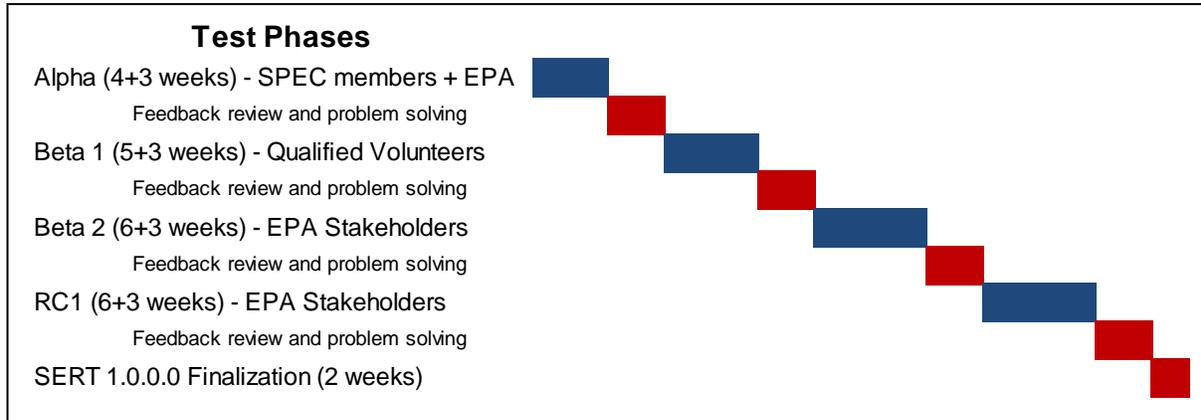


Figure 2: Tentative Test Schedule

287

288 **2.14 Schedule tradeoffs**

289 SPEC benchmarks are developed with the goal to generate results which are directly comparable for
 290 multiple hardware and software architectures to the extent this is possible. The same basic goal
 291 directs the design of SERT as specified in this document.

292 Even though SERT is designed with the goal of being architecture agnostic, code needs to be
 293 implemented for each of the workloads and the tool harness on all supported architectures.
 294 Furthermore this code must be tested intensively on all architectures in order to ensure a functionally
 295 equivalent set of binaries, which generate fair and comparable results. Simply using a portable
 296 programming language will not be sufficient to achieve these goals. Consequently significant
 297 complexity is added to the development process.

298 Given that SERT is designed as a first order approximation rating tool, comparability may be handled
 299 differently than with benchmarks (second order approximation tools) which are used for competitive
 300 marketing. Nevertheless it's essential to ensure a minimal level of comparability.

301 The resources available in the SPECpower committee are limited and a timely development of the tool
 302 for a single architecture will be challenging. Support for additional architectures will remove resources
 303 from the development of the basic test routines because they will be needed for porting the code.
 304 Furthermore additive testing effort is required not only for the new architectures but for the original
 305 implementation as well in order to ensure comparability. Therefore each extra architecture will add a
 306 currently undetermined amount of time to the schedule. The resource and schedule problems recur
 307 with the support of multiple operating systems. SERT will be initially implemented on selected
 308 Operating Systems (OS) per HW architecture.

309

310 **3 SERT Architecture**

311 **3.1 System Overview**

312 SERT shares design philosophies and elements from SPECpower_ssj2008 in its overall architecture.

313 SERT is composed out of multiple software components.

314

315 For the most basic SERT hardware measurement setup one of each of the following is required:

- 316 • system under test (SUT) – the actual system for which the measurements are being taken.
317 The controller and SUT are connected to each other via an Ethernet connection.
- 318 • controller (e.g. server, PC, laptop) – the system to which the power analyzer and temperature
319 sensor are connected.
- 320 • power analyzer – connected to the controller and used to measure the power consumption of
321 the SUT.
- 322 • temperature sensor – connected to the controller and used to measure the ambient
323 temperature where the SUT is located.

324

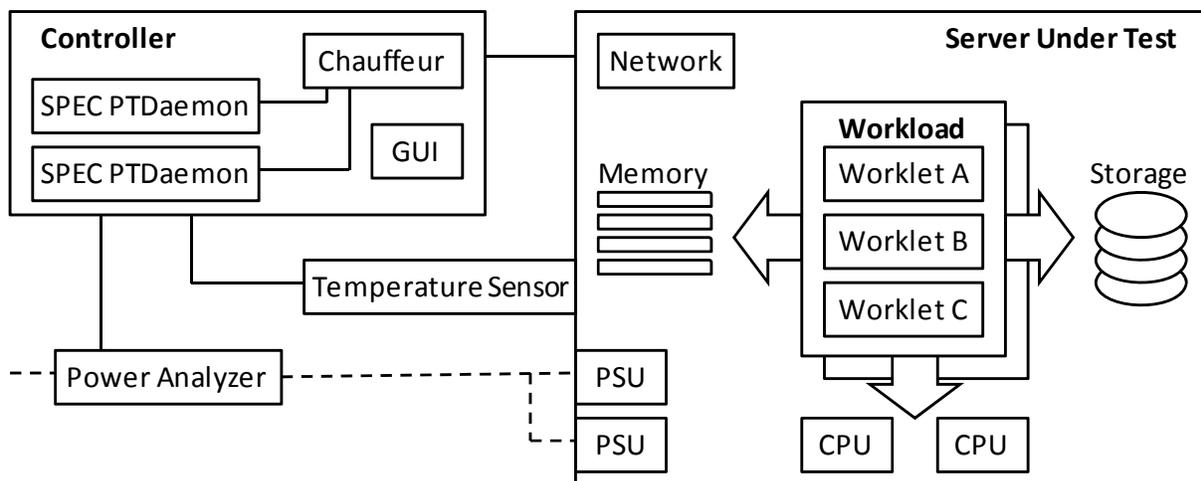
325 The SERT is composed of several elements including:

- 326 • the test harness (Chauffeur) – handles the logistical side of measuring and recording power
327 data along with controlling the software installed on the SUT and controller system itself.
- 328 • the director – instructs the SUT to execute the workload.
- 329 • the workload (a set of worklets) – exercises the SUT while the test harness collects the power
330 and temperature data.
- 331 • the SPEC PTDaemon – connects to the power analyzer and temperature sensor and gathers
332 their readings while the workload executes.
- 333 • the reporter – gathers the environmental, power and performance data after a run is complete
334 and compiles it into an easy to read format.
- 335 • Future versions of the kit will also include a GUI to ease setting up and executing the kit.

336

337 The basic system overview diagram shows these components in relationship to each other.

338



339

Figure 3: SERT Overview

340

341

342

3.2 Execution of SERT

These basic steps are needed in order to execute the SERT kit:

1. Setup the power analyzer and its associated SPEC PTDaemon script.
 - Configure the power analyzer to correctly measure the amperage and voltage of the SUT.
 - Edit the runpower.bat/sh script file to ensure that the proper power analyzer model is specified and the correct communication and network ports are used.
 - Ensure that the SPEC PTDaemon connects and communicates with the power analyzer.
2. Setup the temperature sensor and its associated SPEC PTDaemon script.
 - Edit the runtemp.bat/sh script file to ensure the proper temperature sensor model is specified and the correct communication and network ports are used.
 - Ensure that the SPEC PTDaemon connects and communicates with the temperature sensor.
3. Edit the Director script file.
 - Edit the director.bat/sh script file for the appropriate system configuration.
 - Ensure the proper Java path is specified.
 - Ensure the LOCAL_DIRECTOR variable contains the appropriate information.
 - Ensure the DIRECTOR_HOST variable contains the appropriate information.
4. Edit the SERT script file.
 - Edit the sert.bat/sh script file for the appropriate system configuration.
 - Ensure the proper Java path is specified.
 - Ensure the proper number of JVM's is specified.
 - Ensure the LOCAL_DIRECTOR variable contains the appropriate information.
 - Ensure the DIRECTOR_HOST variable contains the appropriate information.
5. Run the SPEC PTDaemon, Director and SERT scripts.
 - Execute the runpower.bat/sh, runtemp.bat/sh, director.bat/sh and sert.bat/sh scripts.

After the kit completes the run, there should be a results.xml file located in the \results\chauffeur-xxxx directory (where xxxx is the run iteration number).

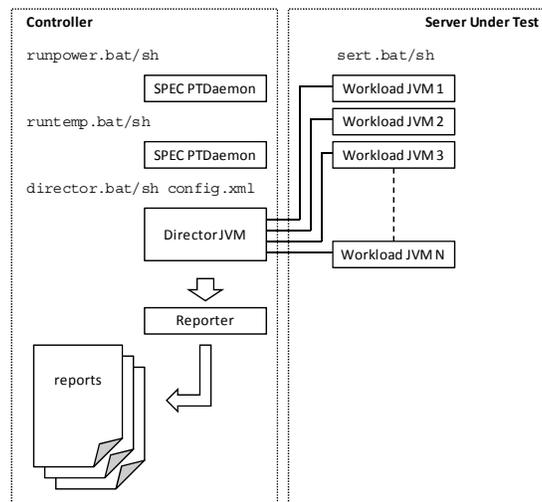


Figure 4: SERT Startup Procedure

372

373

374

404
405

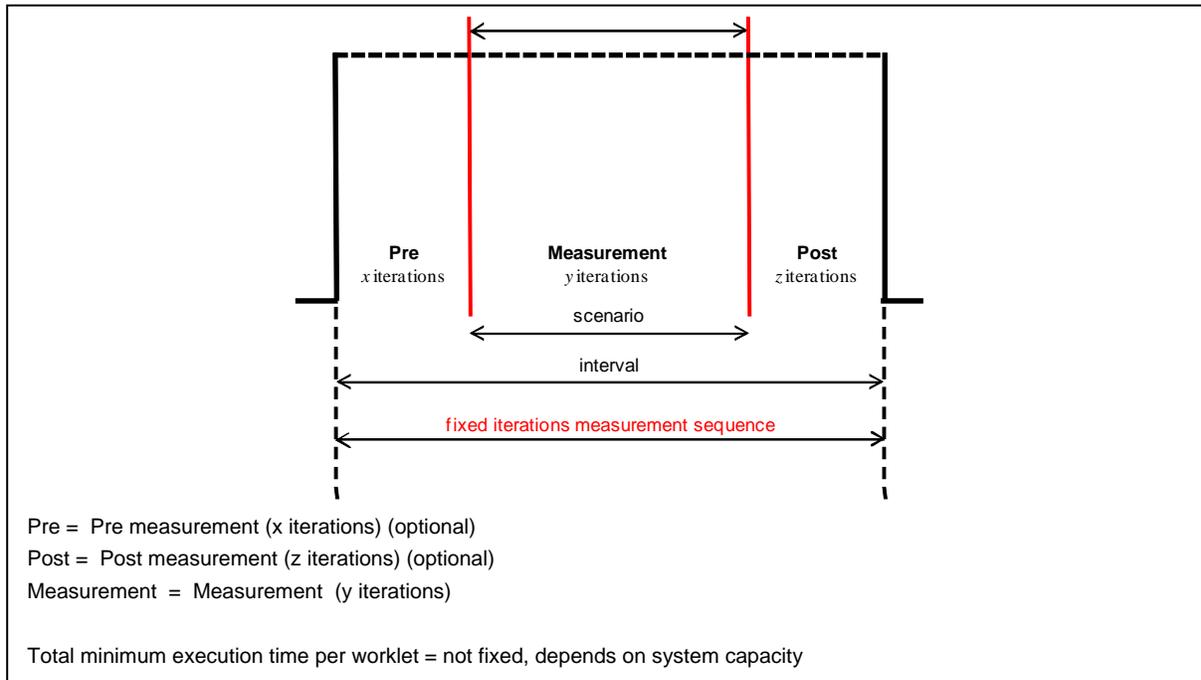


Figure 6: Fixed Iteration Execution

406
407
408
409
410
411
412
413

The fixed iteration execution scheme typically includes one sequence and one interval only. The duration of the interval is not predefined but determined by the capacity of the system, i.e. the time it takes to execute the fixed amount of work.

The number of intervals and scenarios can be defined in SERT configuration files for each worklet individually.

Currently only the Mem_Flood worklet uses this execution scheme

User	A User is a representation of a external agent that can initiate work (e.g. human being)
	Each User may maintain identifying information
	E.g. each User represents a Warehouse
	Each User may maintain state information
	Temporary information that persists from one transaction to another
	There may be multiple types of Users for a single Workload
Transaction	A transaction receives a User and transaction-specific Input as parameters
	It produces some Result
	Some transactions may be able to verify their results – this could be used for a small portion of transactions for auditing purposes
Scenario	A worklet is a set of transactions that can be executed by a particular type of User
	Workloads may contain multiple worklets
	Each worklet could represent a sequence of user interactions
	Think time may occur between transactions
Interval	Each interval in a sequence includes pre-measurement
	Within each interval, each User schedules the execution of worklets
	When a scenario's scheduled time arrives, it iterates through its transactions. Each transaction is submitted to a JVM-wide thread pool. The next transaction in the Worklet will be submitted after the current transaction completes

Sequence	Each phase consists of a sequence of intervals	
	The intervals in a sequence have something in common (though the "something" can vary based on workload or configuration)	
	no delay sequence	fixed time intervals running scenarios unrestricted
	graduated measurement sequence	fixed time intervals with controlled execution of scenarios
	fixed iterations measurement sequence	a predefined number of iterations per scenario is executed, the score is calculated from the execution time, the pre and post interval phases are optional and may be missing for some worklets
Worklet	A workload defines a set of Users and worklets	
	Execution of a Workload includes multiple phases: Warmup Calibration One or more measurement phases	
	Each of these phases are really a sequence of measurement intervals	
	Multiple measurement phases could be used for varying transaction mix, users, etc	

414
415
416
417

418 5 Worklet Design Guidelines

419 In order to achieve consistent results from all worklets and a broad coverage of technologies the
420 following guidelines should be observed:

- 421 • Each worklet must be adjustable to different performance levels, e.g. some predefined levels
422 between 100% (maximum load) and 0% (idle)
- 423 • Each worklet must calibrate to maximum performance level by itself, i.e. no definition of the
424 100% level by the test user
- 425 • Multiple programming languages may be used
- 426 • Precompiled binaries of the test programs should be used where possible.
- 427 • Each worklet should scale with the available hardware resources. More resources should
428 result in a higher performance score, e.g. more processor/memory/disk capacity or additional
429 processor/memory/disk modules yield a better result in the performance component of the
430 efficiency rating.
- 431 • Portable code that follows all SPEC rules for licensing, reuse and adaptation.
- 432 • Either architecture and OS agnostic or with “if-def” capability to accommodate different
433 architectures and/or OSs.
- 434 • The work accomplished by each worklet is clearly identifiable as “important” but is not required
435 to cover “all important” types of work.

436 In order to follow these guidelines the workloads will probably be based on batches of discrete work,
437 where each batch constitutes a transaction. The different load levels will be achieved by scheduling
438 the required number of transactions.

439

440 5.1 Active Idle Worklet

441 During active idle measurements, the SUT must be in a state in which it is capable of completing
442 workload transactions. The active idle worklet is treated in a manner consistent with all other worklets,
443 with the exception that no transactions occur during the active idle interval.

444

445 5.2 CPU Worklet

446 A combination of a wide variety of processor-intensive tasks, including string manipulation, task
447 management, Java “commercial” processing, C “commercial” processing, numeric processing, and
448 other tasks as identified and appropriate.

- 449 • Consistent processor characteristics per simulated “user” regardless of number of processors,
450 cores, enabled threads, etc.
- 451 • Bottleneck at 100% is the processor, not the storage or memory
- 452 • Able to schedule processor tasks or blocks of tasks in such a way that the load can be scaled
453 from 100% in graduated levels down to idle.
- 454 • The CPU worklets should measure a higher (better) performance score for:
 - 455 ○ higher #CPU, higher #core, higher #logical processors, higher frequency, larger
456 overall cache, lower latency, faster interconnect between CPU sockets

457

458 5.3 Memory Worklet

459 Combination of random and sequential reads and writes, small and large memory accesses.

- 460 • Consistent memory access characteristics per simulated “user” regardless of size and number
461 of memory DIMMs
- 462 • Bottleneck at 100% is the memory itself, not the processor or storage
- 463 • Able to schedule memory stress tasks or blocks of tasks in such a way that the load can be
464 scaled from 100% in graduated levels down to idle.
- 465 • The memory worklets should measure a higher (better) performance score based on memory
466 characteristics (e.g. higher bandwidth, lower latency, total memory size)

467

468

469

470 5.4 Network IO Worklet

471 Configuration power/performance modifier will be established in order to address Network IO.

- 472 • Avoid expensive and extensive external test system configurations
- 473 • Measurements show that there are no significant differences in power utilization between
- 474 100% and 0% network utilization for today's technology

475 5.5 Storage IO Worklet

476 Combination of random and sequential, reads and writes, small and large I/Os.

- 477 • Consistent I/O characteristics per simulated "user" regardless of system size and number of
- 478 disks or the installed memory
- 479 • Bottleneck at 100% is the storage subsystem, not the processor or memory
- 480 • Able to schedule I/O tasks or blocks of tasks in such a way that the load can be scaled from
- 481 100% in graduated levels down to idle.
- 482 • The storage worklets should measure a higher (better) performance score for a higher
- 483 bandwidth and lower latency

484 The measurements of power and performance of either optional add-in storage controller cards or

485 server blade enclosure storage are not in the scope of SERT.

486

487 5.6 System Worklet

488 A combination of a wide variety of processor and memory-intensive tasks

- 489 • Bottleneck at 100% is the processor and memory
- 490 • Able to schedule processor tasks or blocks of tasks in such a way that the load can be scaled
- 491 from 100% in graduated levels down to idle.
- 492 • The system worklets should measure a higher (better) performance score for:
 - 493 ○ higher #CPU, higher #core, higher #logical processors, higher frequency, larger
 - 494 overall cache, lower latency, faster interconnect between CPU sockets
 - 495 ○ higher bandwidth, lower latency, total memory size

496

497

498 **6 Power and Temperature Measurements**

499 SERT provides the ability to automatically gather measurement data from accepted power analyzers
500 and temperature sensors and integrate that data into the SERT result. It will be required that the
501 analyzers and sensors must be supported by the measurement framework, and be compliant with the
502 specifications in this section.

503

504 **6.1 Environmental Conditions**

505 Power measurements need to be taken in an environment representative of the majority of usage
506 environments. The intent is to discourage extreme environments that may artificially impact power
507 consumption or performance of the server, before and during the SERT run.

508 The following environmental conditions need to be met:

- 509 • Ambient temperature lower limit: 20°C
- 510 • Ambient temperature upper limit: within documented operating specification of SUT
- 511 • Elevation: within documented operating specification of SUT
- 512 • Humidity: within documented operating specification of SUT
- 513 • Overtly directing air flow in the vicinity of the measured equipment in a way that would be
514 inconsistent with normal data center practices is not allowed.

515

516 **6.2 Temperature Sensor Specifications**

517 Temperature must be measured no more than 50mm in front of (upwind of) the main airflow inlet of the
518 SUT. To ensure comparability and repeatability of temperature measurements, SPEC requires the
519 following attributes for the temperature measurement device used during the benchmark:

- 520 • Logging - The sensor must have an interface that allows its measurements to be read by the
521 benchmark harness. The reading rate supported by the sensor must be at least 4 samples per
522 minute.
- 523 • Accuracy - Measurements must be reported by the sensor with an overall accuracy of +/- 0.5
524 degrees Celsius or better for the ranges measured during the benchmark run.

525

526 **6.3 Power Analyzer Requirements**

527 To ensure comparability and repeatability of power measurements, the following attributes for the
528 power measurement device are required for SERT. Please note that a power analyzer may meet
529 these requirements when used in some power ranges but not in others, due to the dynamic nature of
530 power analyzer Accuracy and Crest Factor. The usage of power analyzer's auto-ranging function is
531 not permitted.

- 532 • Measurements - the analyzer must report true RMS power (watts) and at least two of the following
533 measurement units: voltage, amperes and power factor
- 534 • Accuracy - Measurements must be reported by the analyzer with an overall uncertainty of 1% or
535 better for the ranges measured during the benchmark run. Overall uncertainty means the sum of
536 all specified analyzer uncertainties for the measurements made during the benchmark run.
- 537 • Calibration - the analyzer must be able to be calibrated by a standard traceable to NIST (U.S.A.)
538 (<http://nist.gov>) or a counterpart national metrology institute in other countries. The analyzer must
539 have been calibrated within the past year.
- 540 • Crest Factor - The analyzer must provide a current crest factor of a minimum value of 3. For
541 Analyzers which do not specify the crest factor, the analyzer must be capable of measuring an
542 amperage spike of at least 3 times the maximum amperage measured during any 1-second
543 sample of the benchmark run.
- 544 • Logging - The analyzer must have an interface that allows its measurements to be read by the
545 SPEC PTDaemon. The reading rate supported by the analyzer must be at least 1 set of
546 measurements per second, where set is defined as watts and at least 2 of the following readings:
547 volts, amps and power factor. The data averaging interval of the analyzer must be either 1
548 (preferred) or 2 times the reading interval. "Data averaging interval" is defined as the time period

549 over which all samples captured by the high-speed sampling electronics of the analyzer are
550 averaged to provide the measurement set.

551

552 Examples:

553 An analyzer with a vendor-specified accuracy of +/- 0.5% of reading +/- 4 digits, used in a test with a
554 maximum power value of 200W, would have "overall" accuracy of $((0.5\% * 200W) + 0.4W) = 1.4W/200W$
555 or 0.7% at 200W.

556 An analyzer with a wattage range 20-400W, with a vendor-specified accuracy of +/- 0.25% of range +/-
557 4 digits, used in a test with a maximum power value of 200W, would have "overall" accuracy of
558 $((0.25\% * 400W) + 0.4W) = 1.4W/200W$ or 0.7% at 200W.

559

560 6.4 SPEC PTDaemon

561 SPEC PTDaemon (also known as power/temperature daemon, PTD or ptd) is used by SERT to
562 offload the work of controlling a power analyzer or temperature sensor during measurement intervals
563 to a system other than the SUT. It hides the details of different power analyzer interface protocols and
564 behaviors from the SERT software, presenting a common TCP-IP-based interface that can be readily
565 integrated into different benchmark harnesses.

566 The SERT harness connects to PTDaemon by opening a TCP port and using the simple commands
567 detailed in the API section of this document. For larger configurations, multiple IP/port combinations
568 can be used to control multiple devices.

569 PTDaemon can connect to multiple analyzer and sensor types, via protocols and interfaces specific to
570 each device type. The device type is specified by a parameter passed locally on the command line on
571 initial invocation of the daemon.

572 The communication protocol between the SUT and PTDaemon does not change regardless of device
573 type. This allows SERT to be developed independently of the device types to be supported.

574

575 6.5 Supported and Compliant Devices

576 SERT will utilize SPEC's accepted measurement devices list and SPEC PTDaemon update process.
577 See Device List (http://www.spec.org/power_ssj2008/docs/device-list.html) for a list of currently
578 supported (by the SPEC PTDaemon) and compliant (in specifications) power analyzers and
579 temperature sensors.

580

581 6.6 Power Analyzer Setup

582 The power analyzer must be located between the AC Line Voltage Source and the SUT. No other
583 active components are allowed between the AC Line Voltage Source and the SUT.

584 Power analyzer configuration settings that are set by the SPEC PTDaemon must not be manually
585 overridden.

586

587 6.7 DC Line-Voltage

588 SPEC PTDaemon is neither supported nor tested with DC loads today and currently no resources are
589 devoted to including this support. We are in favor of including DC support if new resources from
590 companies whose focus is DC computing become available to the SPECpower committee to address
591 the development and support opportunity.

592 Additional, comparing servers powered by AC against servers powered by DC is not fair, since the
593 AC-DC conversion losses are not included in DC-powered server. Therefore we recommend creating
594 a separate category for DC-powered servers.

595

596

597

598

599 **7 Metric/Score, Reporting, Logging**

600 **7.1 Metric/Score**

601 While SERT is not intended to be a benchmark, nevertheless as a rating tool it must produce a metric
602 or score indicative of the efficiency of the server under test. That metric must combine both the
603 performance of the SUT as well as its power consumption in a way that allows comparison among all
604 systems subjected to it. The desired outcome of that comparison is a quantitative measure of the
605 relative power-performance efficiencies of the systems. The system which produces the higher metric
606 should have greater power-performance efficiency than the system which produces the lower metric.

607 Since different architectures perform differently on different workloads, SERT is composed of several
608 discrete worklets to ensure architecture neutrality. Each worklet will produce a measure representing
609 the performance achieved by the SUT, which then must be combined with the measures produced by
610 the other worklets to yield a metric indicative of the overall performance of the SUT on all worklets
611 used in the tool. SPEC recommends that the multiple performance measures produced in this manner
612 be combined into a single metric as the geometric mean of the individual measures.

613 The geometric mean of individual worklet performance may be used whether the individual worklets
614 are run sequentially or simultaneously. Depending on the worklets chosen and the magnitudes of
615 their individual measures, we intend indexing the measures to a set of reference scores before
616 combining them into the single metric as the geometric mean. These techniques have the advantages
617 of rendering the single metric unit-less, and of keeping the scale of the individual measures within
618 similar ranges, so that a worklet with large magnitude individual measure does not overwhelm the
619 result from a workload with a smaller measure.

620 Once determined, the overall performance must be combined with the measured power consumption
621 of the SUT in a way that demonstrates the power-performance efficiency of the system. This will be a
622 complex calculation automatically performed by SERT to take into account the power-performance
623 efficiency of the SUT at different utilization levels.

624 The metric that is produced by SERT is separate from the ENERGY STAR rating. The EPA will
625 determine criteria for ENERGY STAR acceptance of which the SERT scores may be only a part. It's
626 anticipated that the top 25% of tested units will achieve ENERGY STAR qualification. A "gold-level"
627 ENERGY STAR qualification may be available for units achieving in the top 5% of results. Additionally
628 the EU has proposed a system of graduated achievement in power-performance efficiency with levels
629 A through F, for which they will determine the overall criteria.

630 Server under test may be placed in different categories by the EPA. The EPA will decide how to apply
631 these categories and whether units in a particular category may be compared to units in another
632 category.

633

634 **7.2 Reporting and Output Files**

635 SERT will produce two reports and a set of log files. The reports will be created in XML format, in
636 order to reduce the effort for both EPA and the partner in displaying and or storing the desired
637 information. We will take steps in order to ensure authenticity (e.g. encryption) of the reports.

638

639 **7.2.1 Report 1: "Summary Report"**

640 This report will contain a placeholder for a "pass or fail" notice for the tested platform, to be provided
641 by the EPA. A test run is marked non-compliant if the test completes with technical errors. In such a
642 case, error messages and/or warnings will be automatically included in the report. The information in
643 this report is public and could be used for marketing purpose.

644

645 Items included in this report are:

- 646 • EPA Partner name and EPA Partner ID
- 647 • EPA ENERGY STAR Category of the tested platform
- 648 • Test Date and Location (plus "Tested by")
- 649 • Tested Platform Manufacturer and Model Number
- 650 • Placeholder for "Pass/Fail"

- 651 • Warnings or Error Notices if applicable
- 652 • System Configuration information (Redundant components to be marked appropriately):
 - 653 ○ form factor
 - 654 ○ number and type of processors
 - 655 ○ available processor sockets
 - 656 ○ memory size, type, # memory DIMMs, # DIMM Slots, Max Memory Capacity
 - 657 ○ available expansion slots
 - 658 ○ number of and make-model of power supply, output rating, min/max
 - 659 ○ Input power
 - 660 ○ OS supported / OS used for test
 - 661 ○ number of and make-model of storage controller
 - 662 ○ number of and make-model of mass storage devices
 - 663 ○ number of and make-model of network interface cards (NICs)
 - 664 ○ Management Controller or Service Processor Installed? [Yes/No]
 - 665 ○ Other Hardware Features / Accessories

666

667 **7.2.2 Report 2: “Power and Performance Data Sheet”.**

668 This report will contain all the information the EPA requires and that is deemed necessary by SPEC.
669 The Power and Performance Data Sheet will be public, but marketing use is prohibited by SERT Fair
670 Usage Rules. The information is intended to be delivered to the EPA in a form most expeditious for
671 EPA review.

672 This report will contain all the data from the “Summary Report” with the following additional detail
673 sections:

- 674 • Overall Result / Score
- 675 • All target load level results
- 676 • Hardware and Software Configuration
- 677 • Power Measurement Summary
- 678 • Environmental information

679

680 **7.3 Validation / Verification**

681 SERT software components will implement software checks wherever possible to increase information
682 accuracy, verify user input, monitor run-time data collection, and validate results with the intent of
683 improving accuracy and remedying user errors, preventing invalid reports to the EPA.

684 When conditions or results do not meet specific criteria, warnings will be displayed and error
685 messages will appear in the SERT reports.

686 These features will make it easy for the EPA Partner to generate compliant results and prevent
687 submission of erroneous reports to the EPA.

688 Examples of compliance checking are:

- 689 • Verify input properties (parameters) and run-time duration of load levels.
- 690 • Temperature out of range will be reported.
- 691 • Power and Temperature read errors must be under a chosen threshold.

692 All the SERT software components will perform validation checks within the domain of their functions,
693 e.g. warnings of connection problems, log measurement errors and out-of-range conditions, warning
694 the user of missing or incomplete information and check the validity of some entered data.

695 Other new validation methods will be considered as the SERT software design and implementation
696 progresses.

697

7.4 Logging

A set of log files will be produced for each test run.

- The information in the log files is intended to be “non-public”.
- These files will be identified by a run serial number such that multiple consecutive test runs produce multiple log file sets.
- Each log file will be a record of actions from the software during the various phases of the testing, including errors and warnings.
- The intent of the log files is for auditing and support purpose.
 - Problems or failures can be more easily resolved with this low level detail record. If any issues arise with regard to the accuracy or veracity of the partner reports, these log files (potentially encrypted) should be adequate to resolve most issues.
- Examples of log file content are:
 - Handshake validation messages among various components
 - Error or warning messages
 - State change messages/notifications.
 - ‘Transaction’ instantaneous/periodic summary information
 - ‘Transaction’ response times

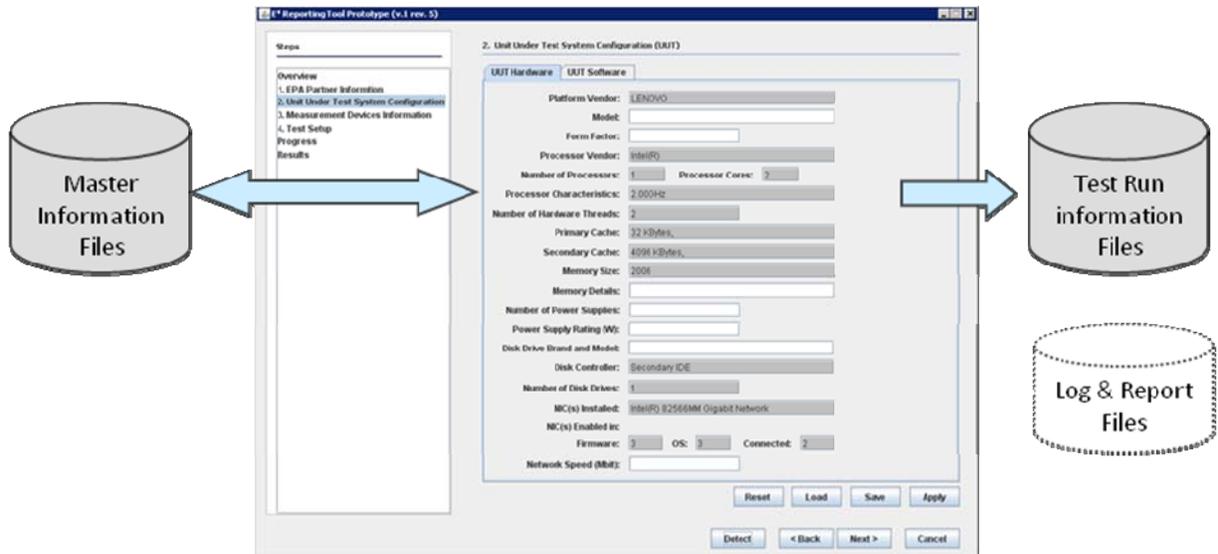
The EPA may require that any or all of the above outputs be delivered prior to ENERGY STAR qualification. Regardless, the partner must commit to archiving all output from any results submitted to the EPA.

718
719
720
721

722 8 Future Enhancements / Stretch goals

723 8.1 Graphical User Interface (GUI)

724 One of the stretch goals is the incorporation of a graphical user interface (GUI) to facilitate
 725 configuration and setup of test runs, allow real-time monitoring of test runs and to review the results.
 726 The SERT GUI will lead the user through the steps of detecting or entering the hardware and software
 727 configuration, setting up a trial run or a valid test, displaying results reports and other functions
 728 common to the testing environment.



729

730

731 The SERT GUI will include several features to enable SERT testing with minimal training and enhance
 732 the accuracy of results:

- 733 • Easy Navigation with Tabbed Screens
- 734 • How to Use (in-line usage guidance and help)
- 735 • Configuration Discovery (Detect function) will automatically populate most fields about SUT
 736 and Controller hardware and software.
- 737 • The GUI will display, allow entry of and store required information about the test environment
 - 738 ○ For use in reports: e.g. Company Info, Platform Config, Run-Time parameters, etc.
 - 739 ○ Master and Test Run information files can be stored, enabling reuse, saving time with
 740 multiple platforms.
- 741 • Test Setup, Execution and Progress Display
 - 742 ○ Start measurements; Choose type of run (trial or final)
 - 743 ○ Display progress, warnings and errors.
- 744 • Display results and enable printing and capture of reports
- 745 • Provisions for redundant components and power and performance modifier.

746

747 8.2 Test Software

748 A “stretch goal” of SERT is to enable a “Live CD” approach to tool installation, for some environments
 749 – such that the entire tool suite along with the underlying operating system could all be run from a
 750 single bootable CD or DVD with no other operating system installed on the SUT. This should provide
 751 increased ease of installation and improve the adoption rate of the tool.

752 Possible issues with this approach include the lack of specific hardware drivers for newer devices, the
 753 potential lack of vendor specific power management, licensing and availability issues for some
 754 operating systems. Alternatives include allowing additional drivers to be installed during setup, or
 755 providing separate test installers with binaries for use with a vendor’s own as-shipped OS installation.

756

757 **9 SERT and EPA ENERGY STAR for Server Version 2.0**

758 In order to ensure that SERT is utilized in the intended matter, we recommend the inclusion of the
759 following items in the ENERGY STAR for Server Specification.

760

761 **9.1 Measurement**

762 The provided SERT test kit must be used to run and produce measured SERT results. The SERT
763 metric is a function of the SERT workload (see section 0). SERT results are not comparable to power
764 and performance metrics from any other application.

765

766 **9.2 SERT Binaries and Recompile**

767 Valid runs must use the provided binary files and these files must not be updated or modified in any
768 way.

769

770 **9.3 Manual Intervention**

771 No manual intervention or optimization for the SUT or its internal and external environment is allowed
772 during the test measurement, after initial setup is completed.

773

774 **9.4 Fair Use of SERT information**

775 A clear goal of the ENERGY STAR program is to have the broadest possible participation among
776 vendors. Experience in the computer industry's performance benchmark community demonstrates that
777 when performance details become available for marketing purposes, only vendors with superior (at the
778 time of publication) products are incented to publish results. To encourage broader participation
779 across the industry, a set of strong rules must be in place that will restrict marketing use of any of the
780 detailed information generated by the tool. No data besides the actual ENERGY STAR qualification
781 should be utilized in EPA Partners' marketing collateral. These rules will be stipulated in both the
782 license for the tool and the EPA Partner agreement.

783 Note that, while these rules are not strictly a part of the tool "design", the existence of these rules are
784 necessary to allow the flexibility of the design and the delivery of detailed consumer information that is
785 desired.

786

787 **9.4.1 Fair Use Rules**

- 788 • The only information provided by the tool that can be used for marketing collateral is the ENERGY
789 STAR qualification of a server configuration or server family
- 790 • The only information provided by the tool that can be used for public comparison is the ENERGY
791 STAR qualification of a server configuration or server family All other publicly available information
792 from the tool is made available to help to verify that the tests were run correctly and to allow
793 consumers to better understand how well the configurations tested match their specific needs.
- 794 • If the tool is used for research to generate information outside of the ENERGY STAR program, the
795 information may not be compared to the ENERGY STAR program results and competitive
796 comparisons may not be made using the data generated.
- 797 • The EPA ENERGY STAR Qualification is governed by EPA rules.

798

799 **9.5 Accredited, Independent laboratory**

800 The requirement to use accredited, independent laboratories may place a large burden on EPA
801 ENERGY STAR partners, especially smaller companies. We recommend the use of an independent
802 laboratory as an option, but not implementing this as a requirement.

803

9.6 Supply Voltage tolerance

805 In order to use a voltage within a 1% difference, an extra voltage source is needed. This will
806 unnecessarily increase the cost for the partner, especially smaller companies. We recommend the
807 tolerance be set to $\pm 5\%$.

808

809

810

811 **10 Worklet Candidates**

812

813 The following table shows the current Worklet candidates and their anticipated use in different SERT
 814 test phases. Worklet candidates included in early releases may change in subsequent releases. Early
 815 release test results may influence the inclusion of some worklets in future releases.

816

Workload	Worklet candidate	Alpha	Beta 1	Beta 2	RC1
CPU	CPU_Compress	Included	TBD	TBD	TBD
CPU	CPU_CryptoAES	Included	TBD	TBD	TBD
CPU	CPU_SOR	Included	TBD	TBD	TBD
CPU	CPU_FFT	Included	TBD	TBD	TBD
CPU	CPU_LU	Included	TBD	TBD	TBD
CPU	CPU_XMLvalidate	Included	TBD	TBD	TBD
Memory	Mem_Flood	Included	TBD	TBD	TBD
Memory	Mem_XMLvalidate1	Included	TBD	TBD	TBD
Memory	Mem_XMLvalidate2	Included	TBD	TBD	TBD
Storage	Storage_Random	-	TBD	TBD	TBD
Storage	Storage_Sequential	-	TBD	TBD	TBD
Storage	Storage_Mixed	Included	TBD	TBD	TBD
System	System_CSSJ	Included	TBD	TBD	TBD

817

818

819 **10.1 CPU Worklet: Compress**

820

821 **10.1.1 General Description**

822 The Compress workload implements a transaction that compresses and decompresses data using a
 823 modified Lempel-Ziv method (LZW). Essentially, it finds common substrings and replaces them with a
 824 variable size code. This is both deterministic and done on the fly. Thus, the decompression procedure
 825 needs no input table, but tracks the way the table was built. The algorithm is based on "A Technique
 826 for High Performance Data Compression", Terry A. Welch, IEEE Computer Vol 17, No 6 (June 1984),
 827 pp 8-19.

828

829 **10.1.2 Sequence Execution Methods**

830 Graduated Measurement Sequence

831

832 **10.1.3 Metric**

833 Transactions Per Second

834

835 **10.1.4 Required Initialization**

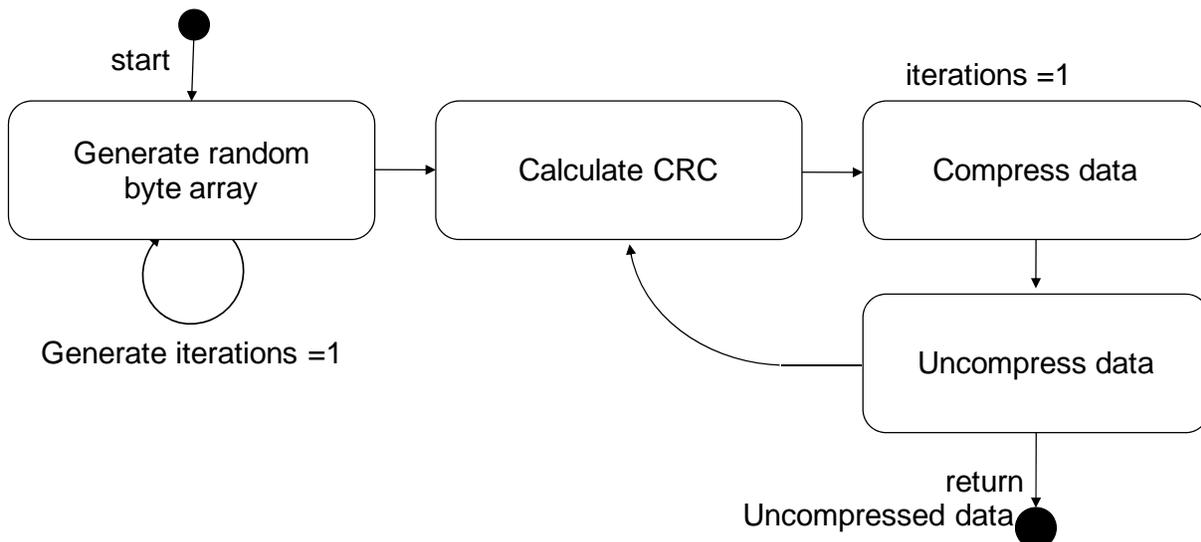
836 A constant size byte array is generated on the fly before for each transaction execution. The contents
 837 of the byte array are randomly generated.

838

839 **10.1.5 Configuration Parameters**

size	Size of the input byte array for each transaction execution.
enable-idc	Enables/disables memory scaling using input data caching (IDC). Must be set to false.
iterations	Number of executions per transaction.
debug-level	Value governs the volume of debug messages printed during execution.
input-generate-iterations	Number of random byte array assignment iterations.

840

841 **10.1.6 Transaction Code**

842

843

844

845

846

847

848 **10.2 CPU Worklet: CryptoAES**

849

850 **10.2.1 General Description**

851 The CryptoAES workload implements a transaction that encrypts and decrypts data using the AES (or
 852 DES) block cipher algorithms. Which algorithm is a configurable parameter, but the current candidate
 853 version uses AES with CBC and no PKCS5 padding. Encryption and decryption are done using the
 854 Java Cryptographic Extension (JCE) framework, and the Cipher class in particular.

855

856 **10.2.2 Sequence Execution Methods**

857 Graduated Measurement Sequence

858

859 **10.2.3 Metric**

860 Transactions Per Second

861

862 **10.2.4 Required Initialization**

863 A constant size byte array is generated on the fly before for each transaction execution. The contents
 864 of the byte array are randomly generated.

865

866 **10.2.5 Configuration Parameters**

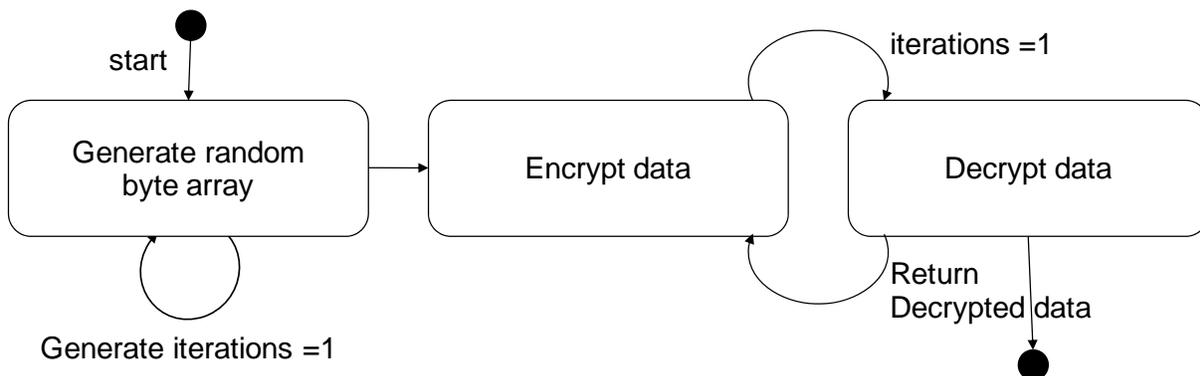
867

size	Size of the input byte array for each transaction execution.
key-generator	Key generator algorithm. (AES or DESede)
key-size	Key size. (128 for AES, 168 for DES)
algorithm	Encryption algorithm. (E.g., AES/CBC/NoPadding, AES/CBC/PKCS5Padding, DESede/CBC/NoPadding, DES/CBC/PKCS5Padding)
level	Number of times to perform the encryption.
enable-idc	Enables/disables memory scaling using input data caching (IDC). Must be set to false.
iterations	Number of executions per transaction.
debug-level	Value governs the volume of debug messages printed during execution.
input-generate-iterations	Number of random byte array assignment iterations.

868

869 **10.2.6 Transaction Code**

870



871

872

873

874 **10.3 CPU Worklet: FFT**

875

876 **10.3.1 General Description**

877 The Fast Fourier Transform (FFT) workload implements a transaction that performs a one-dimensional
 878 forward transform of complex numbers. Its floating point computations exercise complex arithmetic,
 879 shuffling, non-constant memory references and trigonometric functions. The first section performs the
 880 bit-reversal portion (no flops) and the second performs the actual $N\log(N)$ computational steps.
 881 (Adapted from the NIST-developed Scimark benchmark.)
 882

883 **10.3.2 Sequence Execution Methods**

884 Graduated Measurement Sequence

885

886 **10.3.3 Metric**

887 Transactions Per Second

888

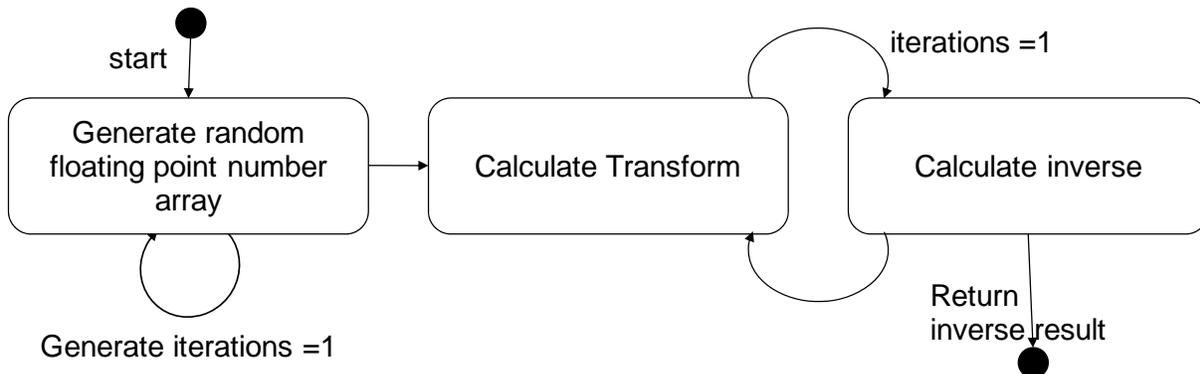
889 **10.3.4 Required Initialization**

890 A constant size floating point number array is generated on the fly before for each transaction
 891 execution. The contents of the array are randomly generated.
 892

893 **10.3.5 Configuration Parameters**

array-length	Size of the input floating point number array for each transaction execution.
enable-idc	Enables/disables memory scaling using input data caching (IDC). Must be set to false.
iterations	Number of executions per transaction.
debug-level	Value governs the volume of debug messages printed during execution.
input-generate-iterations	Number of random array assignment iterations.

894

895 **10.3.6 Transaction Code**

896

897

898

899

900 **10.4 CPU Workload: LU**

901
902 **10.4.1 General Description**

903 The LU workload implements a transaction that computes the LU factorization of a dense matrix using
904 partial pivoting. It exercises linear algebra kernels (BLAS) and dense matrix operations. The algorithm
905 is the right-looking version of LU with rank-1 updates. (Adapted from the NIST-developed Scimark
906 benchmark.)

907
908 **10.4.2 Sequence Execution Methods**

909 Graduated Measurement Sequence

910
911 **10.4.3 Metric**

912 Transactions Per Second

913
914 **10.4.4 Required Initialization**

915 A constant size matrix of floating point numbers is generated on the fly before for each transaction
916 execution. The contents of the matrix are randomly generated.

917
918 **10.4.5 Configuration Parameters**

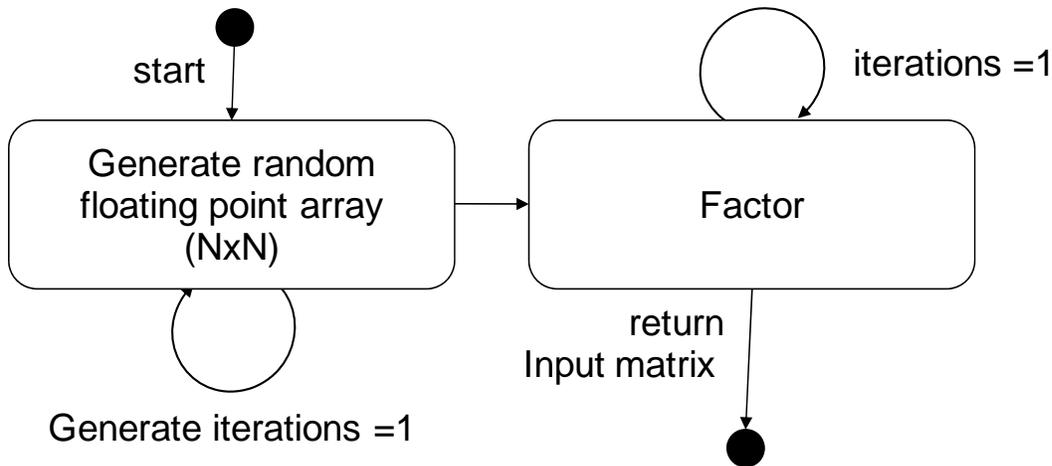
919

matrix-dimen	Dimension of the input floating point matrix for each transaction execution. (NxN)
enable-idc	Enables/disables memory scaling using input data caching (IDC). Must be set to false.
iterations	Number of executions per transaction.
debug-level	Value governs the volume of debug messages printed during execution.
input-generate-iterations	Number of random matrix assignment iterations.

920

921

922 **10.4.6 Transaction Code**



923

924

925

926

927

928 **10.5 CPU Workload: SOR**

929

930 **10.5.1 General Description**

931 The Jacobi Successive Over-relaxation (SOR) workload implements a transaction that exercises
 932 typical access patterns in finite difference applications, for example, solving Laplace's equation in 2D
 933 with Dirichlet boundary conditions. The algorithm excercises basic "grid averaging" memory patterns,
 934 where each $A(i,j)$ is assigned an average weighting of its four nearest neighbors. Some hand-
 935 optimizing is done by aliasing the rows of $G[][]$ to streamline the array accesses in the update
 936 expression. (Adapted from the NIST-developed Scimark benchmark.)

937

938 **10.5.2 Sequence Execution Methods**

939 Graduated Measurement Sequence

940

941 **10.5.3 Metric**

942 Transactions Per Second

943

944 **10.5.4 Required Initialization**

945 A constant size grid of floating point numbers is generated on the fly before for each transaction
 946 execution. The contents of the grid are randomly generated.

947

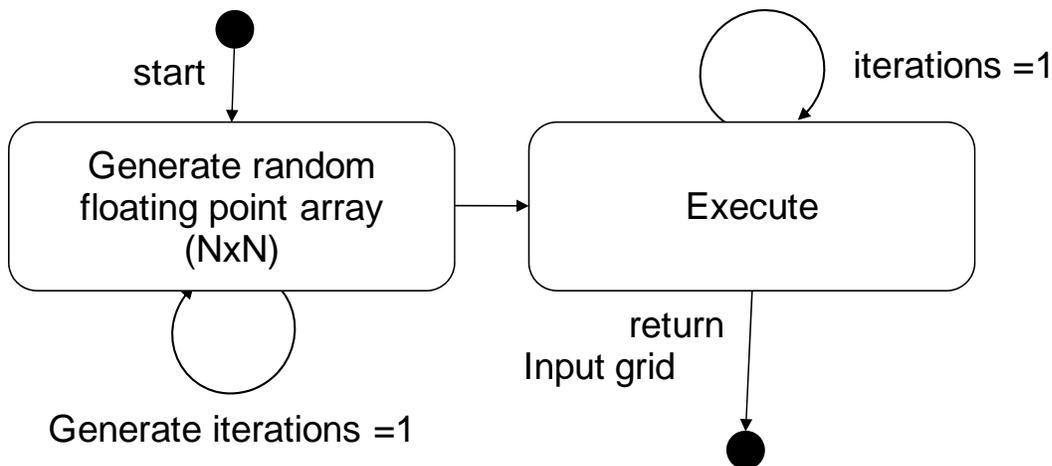
948 **10.5.5 Configuration Parameters**

949

grid-dimen	Dimension of the input floating point grid for each transaction execution. (NxN)
enable-idc	Enables/disables memory scaling using input data caching (IDC). Must be set to false.
iterations	Number of executions per transaction.
debug-level	Value governs the volume of debug messages printed during execution.
input-generate-iterations	Number of random grid assignment iterations.

950

951 **10.5.6 Transaction Code**



952

953

954

955

956

957

958 **10.6 CPU Workload: XmlValidate**

959

960 **10.6.1 General Description**

961 The XML validate workload implements a transaction that exercises Java's XML validation package
 962 javax.xml.validation. Using both SAX and DOM APIs, an XML file (.xml) is validated against an XML
 963 schemata file (.xsd). To randomize input data, an algorithm is applied that swaps the position of
 964 commented regions within the XML input data.

965

966 **10.6.2 Sequence Execution Methods**

967 Graduated Measurement Sequence

968

969 **10.6.3 Metric**

970 Transactions Per Second

971

972 **10.6.4 Required Initialization**

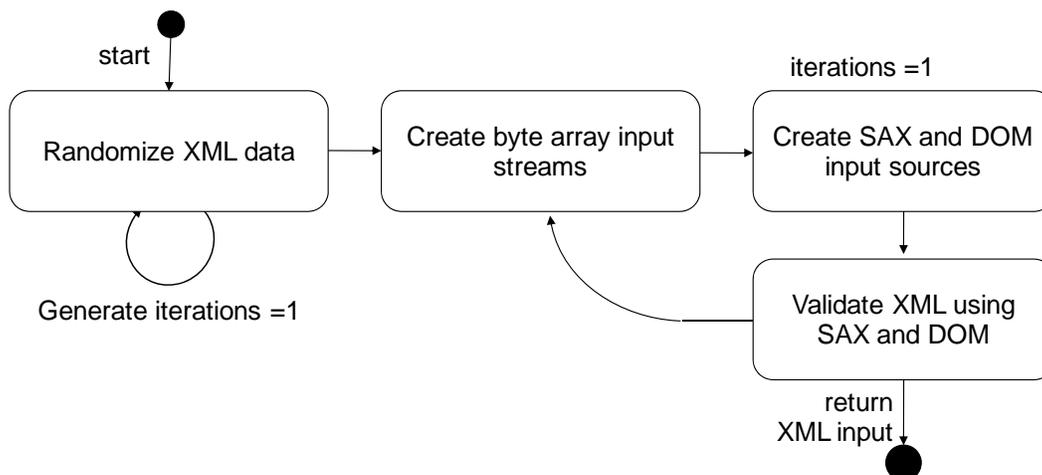
973 A initialization time, both XML and XML schemata files are read in from disk and saved in a buffer for
 974 future use. (There will be no further disk IO once this is completed.) A randomization algorithm is
 975 applied to the original XML data on the fly before for each transaction execution to create variations in
 976 parsing without modifying file size or complexity.

977

978 **10.6.5 Configuration Parameters**

xml-schema-dir	Specifies the directory of the XML schema file.
xml-schema-file	Specifies the name of the XML schema file.
xml-dir	Specifies the directory of the XML file.
xml-file	Specifies the name of the XML file.
enable-idc	Enables/disables memory scaling using input data caching (IDC). Must be set to false.
iterations	Number of executions per transaction.
debug-level	Value governs the volume of debug messages printed during execution.
input-generate-iterations	Number of XML file randomization iterations.

979

980 **10.6.6 Transaction Code**

981

982

983

984

985 date: 03/09/2011

page 35 of 45

Copyright © 2006-2011 SPEC

986 10.7 Memory Worklet: Flood

987

988 10.7.1 General Description

989 The Flood workload is based upon STREAM, a popular benchmark that measures memory bandwidth
990 across four common and important array operations. For the *long* (64-bit) integer arrays used in
991 Flood, the following amounts of memory are involved per assignment:

- 992 1. **COPY:** $a(i) = b(i)$
993 -- 8 bytes read + 8 bytes write per assignment = 16 bytes / assignment
- 994 2. **SCALE:** $a(i) = k * b(i)$
995 -- 8 bytes read + 8 bytes write per assignment = 16 bytes / assignment
- 996 3. **ADD:** $a(i) = b(i) + c(i)$
997 -- 16 bytes read + 8 bytes write per assignment = 24 bytes / assignment
- 998 4. **TRIAD:** $a(i) = b(i) + k * c(i)$
999 -- 16 bytes read + 8 bytes write per assignment = 24 bytes / assignment

1000 The Flood score is based upon the aggregate system memory bandwidth calculated from the average
1001 of these four tests multiplied by the amount of physical memory installed in the SUT. While Flood is
1002 based upon STREAM, it uses no STREAM code and is implemented wholly in Java.

1003 Flood enhances STREAM in a variety of important ways.

- 1004 1. Flood rewards systems with large memory configurations by scaling results based upon
1005 physical memory size.
- 1006 2. Flood is designed to fully exploit the memory bandwidth capabilities of modern multi-core
1007 servers. Flood is multithreaded and threads are scheduled to operate concurrently during
1008 bandwidth measurements ensuring maximum throughput and minimizing result variability.
- 1009 3. Flood requires little to no user configuration, yet automatically expands the data set under test
1010 to fully utilize available memory.

1011 Measuring aggregate system memory bandwidth on large servers with many cores and multiple
1012 memory controllers is challenging. In particular, run-to-run variability is often unmanageable with
1013 existing memory bandwidth benchmarks. Flood minimizes run-to-run variation by taking three
1014 memory bandwidth tests back-to-back and discarding the first and last tests. This ensures that all
1015 threads are running under fully concurrent conditions during the middle measurement which is
1016 used in Flood scoring calculations.

1017

1018 Flood scores scales linearly with a SUT's aggregate memory bandwidth as well as with the SUT's
1019 physical memory configuration. CPU, storage and network performance have little to no impact on
1020 Flood scores.

1021

1022 Since the Flood workload always deploys a fixed number of iterations and the amount of memory
1023 under test will automatically adjust to fully utilize installed DRAM, run time will vary depending upon
1024 system configuration. On a 2.2GHz, 24-core SUT with 24 threads and 48GB of physical memory,
1025 Flood takes about 20 minutes to complete. Run time varies proportionally with the amount of physical
1026 memory installed in the SUT. Run time is also impacted by the overall thread count.

1027

1028 10.7.2 Sequence Execution Methods

1029 *FixedIterationsDirectorSequence* – Flood is executed for a given set of iterations specified within
1030 *config.xml*.

1031

1032 10.7.3 Metric

1033 Score = aggregate system memory bandwidth (GB/s) * physical memory size (GB)

1034

1035 **10.7.4 Required Initialization**

1036 Flood calculates the amount of memory available to the thread and creates three 64-bit (*long*) integer
 1037 arrays, a[], b[] and c[], to completely utilize all available space. These arrays are initialized with
 1038 random data. To ensure full load concurrency during bandwidth measurements, a complete set of
 1039 pre-measurement tests is launched prior to an identical measurement period followed by identical
 1040 post-measurement tests. Only the test results for the measurement period are utilized for Flood score
 1041 generation.

1042

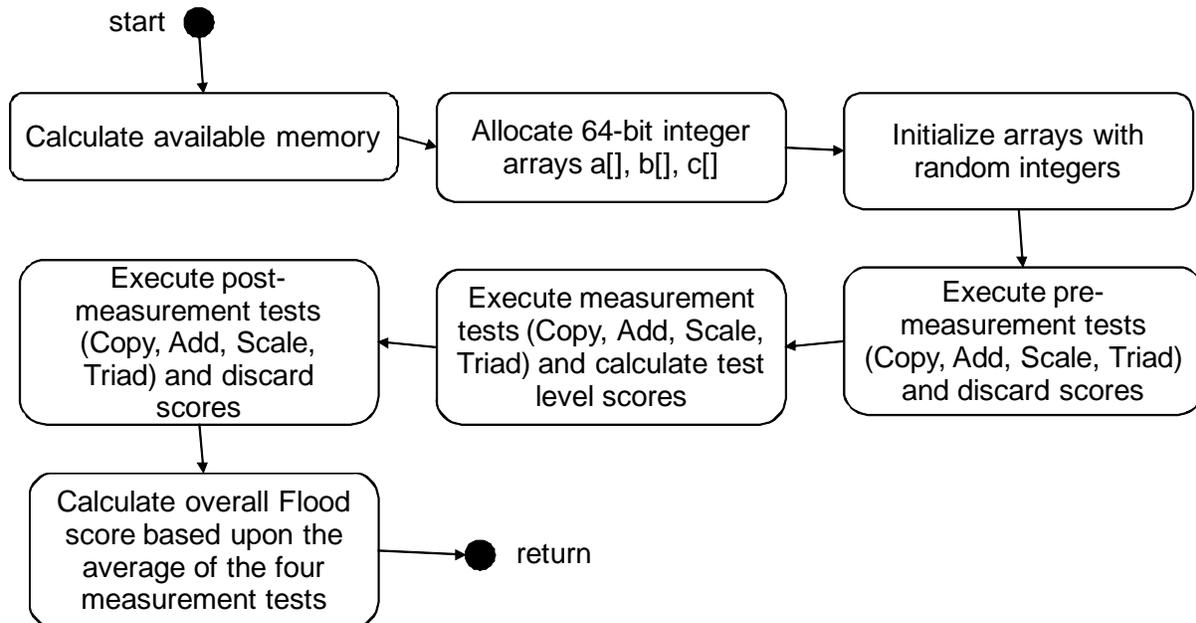
1043 **10.7.5 Configuration Parameters**

memory-under-test	The default value of “-1 MB” turns on automatic configuration of the data set size. However, the user can override this behavior and explicitly define the amount of memory to test per JVM. Valid values are (san quotation marks): “200 MB”, “1.1 GB”, “10000000 B”.
iterations	Flood internally iterates the number of memory bandwidth tests based upon the value of the iterations parameter. The default is 100.
debug-level	Detailed diagnostic information can be enable through the <i>debug</i> parameter. Valid values are 0 = no additional debug information (default), 1 = debug information turned on, 2 = detailed debug information.
return-bandwidth	The raw, aggregate system memory bandwidth calculated by Flood can be obtained by setting the parameter return-bandwidth to “true” in which case Flood will return measured memory bandwidth instead of a score. The default value is “false”.

1044

1045 **10.7.6 Transaction Code**

1046



1047

1048

1049

1050 **10.8 Memory Workload: XmlValidate**

1051

1052 **10.8.1 General Description**

1053 The XML validate workload implements a transaction that exercises Java's XML validation
 1054 package javax.xml.validation. Using both SAX and DOM APIs, an XML file (.xml) is validated
 1055 against an XML schemata file (.xsd). To randomize input data, an algorithm is applied that
 1056 swaps the position of commented regions within the XML input data.

1057 Memory scaling in XmlValidate is done through a scheme known as input data caching
 1058 (IDC). In IDC, the universe of possible input data (here, randomized XML file data) is pre-
 1059 computed and then cached within memory before the start of the workload. During workload
 1060 execution, the input data for a particular transaction instance is then chosen randomly and
 1061 retrieved from this cache rather than computed on the fly.

1062

1063 **10.8.2 Sequence Execution Methods**

1064 Graduated Measurement Sequence

1065

1066 **10.8.3 Metric**

1067 Transactions Per Second * Cache size * Cache size scaling factor

1068

1069 **10.8.4 Required Initialization**

1070 A initialization time, both XML and XML schemata files are read in from disk and saved in a buffer for
 1071 future use. (There will be no further disk IO once this is completed.) IDC initialization follows during
 1072 which all possible input data sets are pre-computed and cached in memory. For each input data set, a
 1073 randomization algorithm is applied to the original XML data to create variations in parsing without
 1074 modifying file size or complexity.

1075

1076 **10.8.5 Configuration Parameters**

1077

1078 XmlValidate parameters:

1079

xml-schema-dir	Specifies the directory of the XML schema file.
xml-schema-file	Specifies the name of the XML schema file.
xml-dir	Specifies the directory of the XML file.
xml-file	Specifies the name of the XML file.
enable-idc	Enables/disables memory scaling using input data caching (IDC). Must be set to false.
iterations	Number of executions per transaction.
debug-level	Value governs the volume of debug messages printed during execution.
input-generate-iterations	Number of XML file randomization iterations.

1080

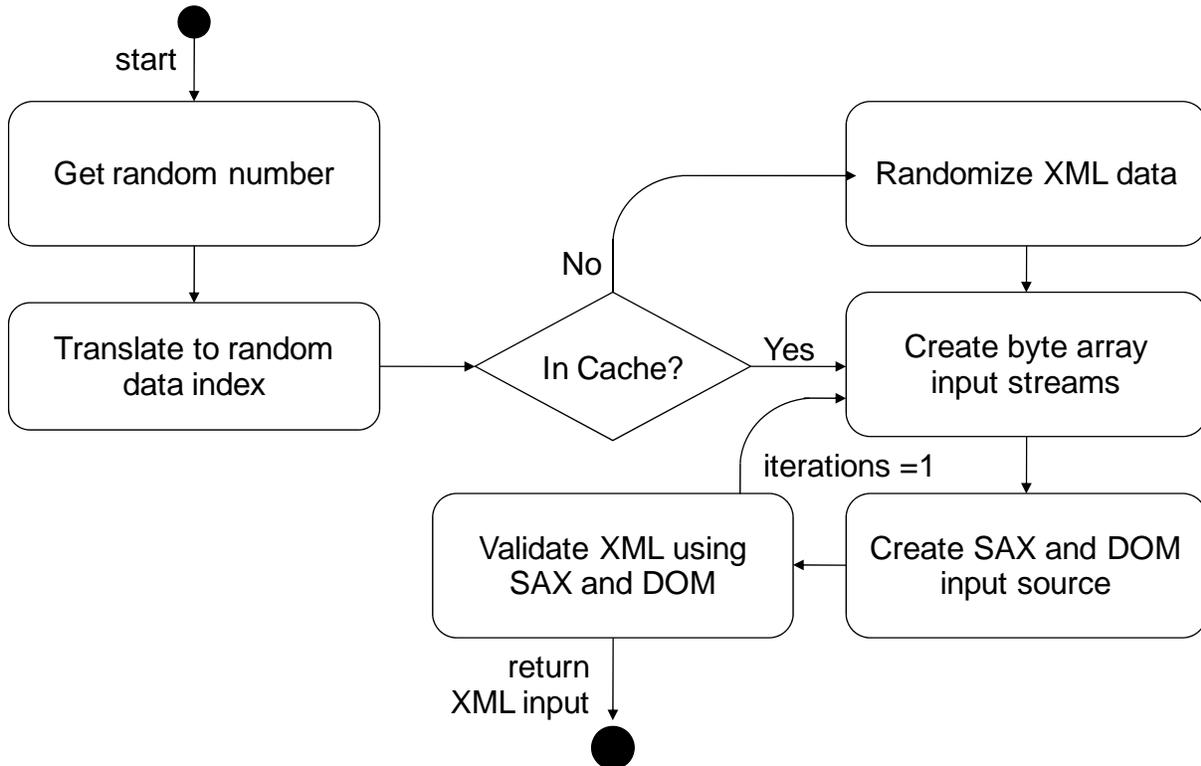
1081

1082 Additional IDC configuration parameters:
 1083

store-type	Specifies the algorithm to use in generating data when a cache miss occurs.
locality-distribution	Specifies the probability distribution to use when randomly choosing input data indices.
data-store-size	Specifies the size of the universe of possible input data.
data-cache-size	Specifies the size of the input data cache.
data-cache-report-interval	Governs the frequency of output messages on cache hit/miss ratio.
custom-score-policy	Specifies the algorithm to use in computing custom score reflecting cache size configuration.
data-cache-size-scale-factor	Specifies the scaling factor to use in the DataCacheSizeMultiplierGB custom scoring algorithm.
data-cache-to-heap-ratio	Ratio of cache size to JVM heap size used in automatic cache sizing.

1084
 1085 **10.8.6 Transaction Code**

1086
 1087
 1088



1089
 1090

1091 **10.9 Storage IO Workload**

1092

1093 **10.9.1 General Description**

1094 The Storage-Workload has four different transactions, two random and two sequential transaction-
 1095 pairs. Each pair has a write and a read transaction.

1096

1097 **10.9.2 Sequence Execution Methods**

1098 [Graduated Measurement Sequence] or [Fixed Iteration Measurement Sequence]

1099

1100 **10.9.3 Metric**

1101 Score name and definition of what the score value represent

1102

1103 **10.9.4 Required Initialization**

1104 A set of files is created before execution of the transaction

1105

1106 **10.9.5 Configuration Parameters**

1107

file-size-bytes	size of a file.
file-per-user	number of files opened by each user.
file-path	location of the files - In this example the path is "D:\data\", please note that the files always reside in a subfolder called "data".
max-count	amount of blocks that are accessed by the sequential transaction in one file before the next file is addressed.

1108

1109 Example:

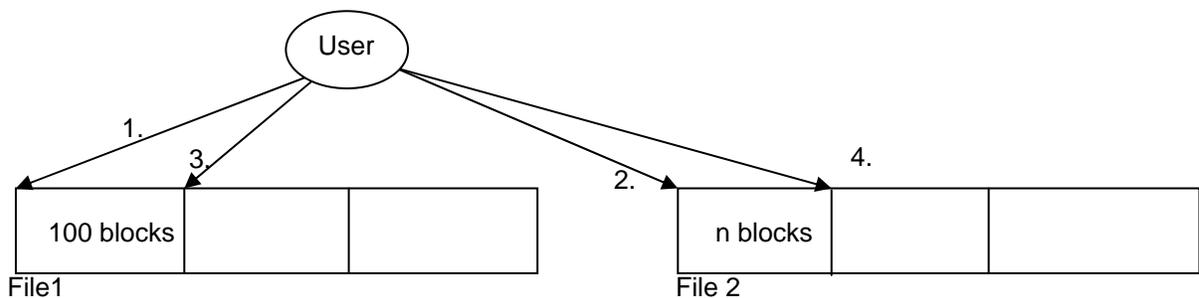
```
1110 <file-size-bytes>1000000</file-size-bytes>
```

```
1111 <file-path>D:\</file-path>
```

```
1112 <file-per-user>2</file-per-user>
```

```
1113 <max-count>100</max-count>
```

1114



1115

[Figure 7: File Example (2 files per user and max-count of 100)]

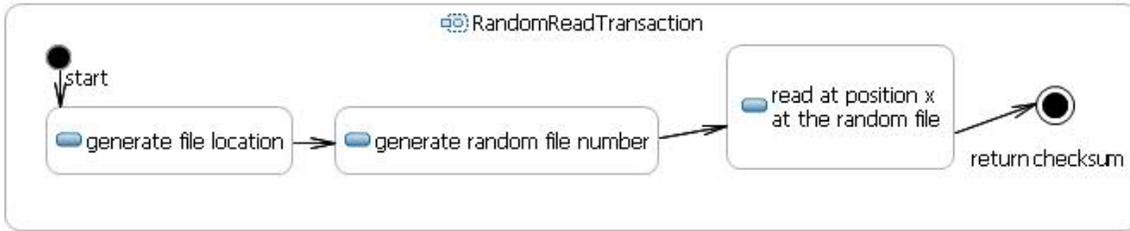
1116

1118

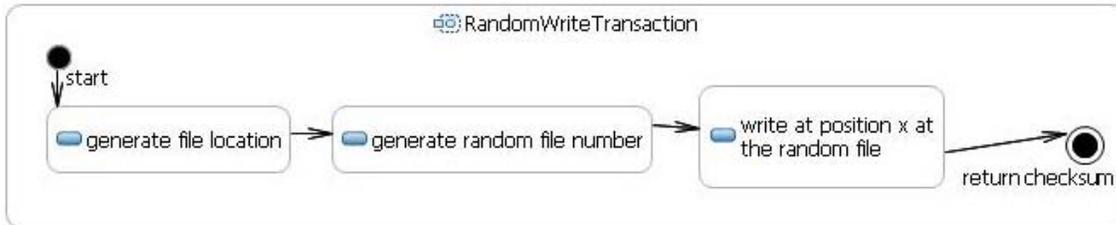
1119

1120

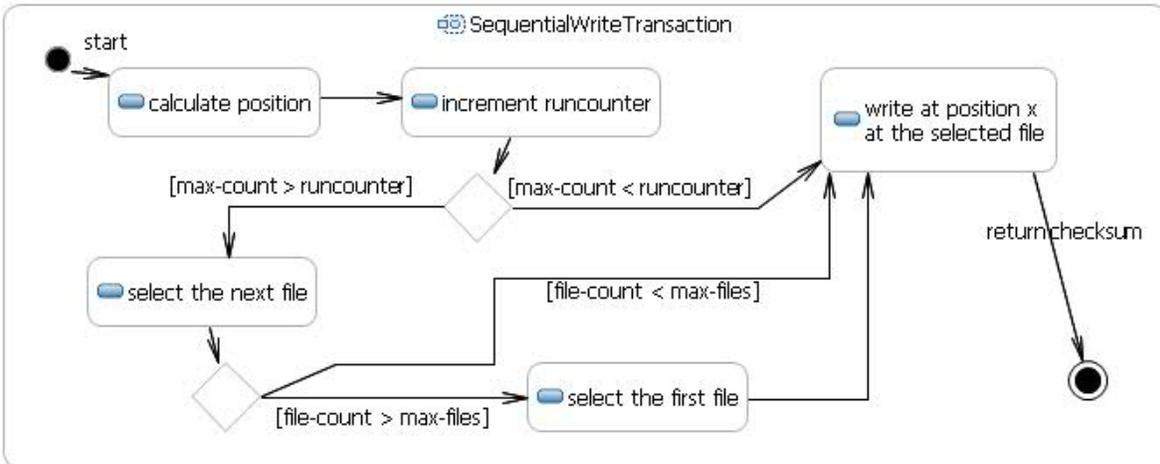
1121 **10.9.6 Transaction – Code 1 - RandomRead**



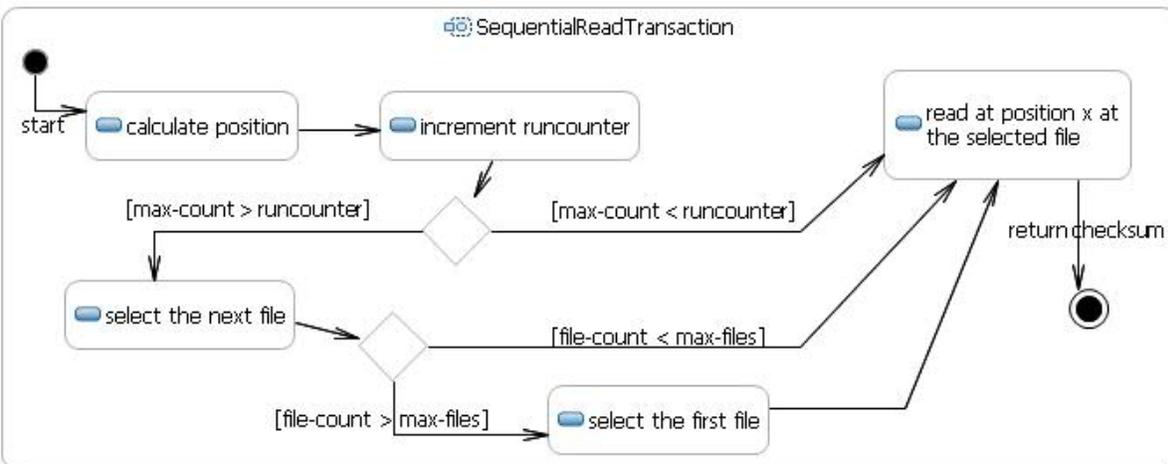
1122 **10.9.7 Transaction – Code 1 - RandomWrite**
 1123



1124 **10.9.8 Transaction – Code 2 – SequentialRead**
 1125



1126 **10.9.9 Transaction – Code 2 – SequentialWrite**
 1127



1128
 1129

1130 **10.10 System Worklet: CSSJ**

1131

1132 **10.10.1 General Description**

1133 CSSJ is an Online Transaction Processing (OLTP) workload, and represents a Server Side Java
1134 application. It is based on the SSJ workload in SPECpower_ssj2008, which was based on
1135 SPECjbb2005, which was inspired by the TPC-C specification; however, there are several differences
1136 between all of these workloads, and CSSJ results are not comparable to any of these other
1137 benchmarks.

1138

1139 The System Worklet exercises the CPU(s), caches, and memory of the UUT. The peak throughput
1140 level is determined by maximum number of transaction of the above type the system can perform per
1141 second. Once the peak value of the transactions is determined on a given system, the worklet is run
1142 from peak (100%) down to the system idle in a graduated manner.

1143 The performance of the System Worklet depends on the combination of the processor type, number of
1144 processors, their operating speed, and the latency and bandwidth of the memory subsystem of the
1145 system.

1146

1147 CSSJ includes 6 transactions, with the approximate frequency shown below:

- 1148 • New Order (30.3%) – a new order is inserted into the system
- 1149 • Payment (30.3%) – record a customer payment
- 1150 • Order Status (3.0%) – request the status of an existing order
- 1151 • Delivery (3.0%) – process orders for delivery
- 1152 • Stock Level (3.0%) – find recently ordered items with low stock levels
- 1153 • Customer Report (30.3%) – create a report of recent activity for a customer

1154

1155 **10.10.2 Sequence Execution Methods**

1156 Graduated Measurement Sequence

1157

1158 **10.10.3 Metric**

1159 Transactions per second

1160

1161 **10.10.4 Required Initialization**

1162 Each user represents a warehouse. During initialization, each warehouse is populated with a base set
1163 of data, including customers, initial orders, and order history.

1164

1165 **10.10.5 Configuration Parameters**

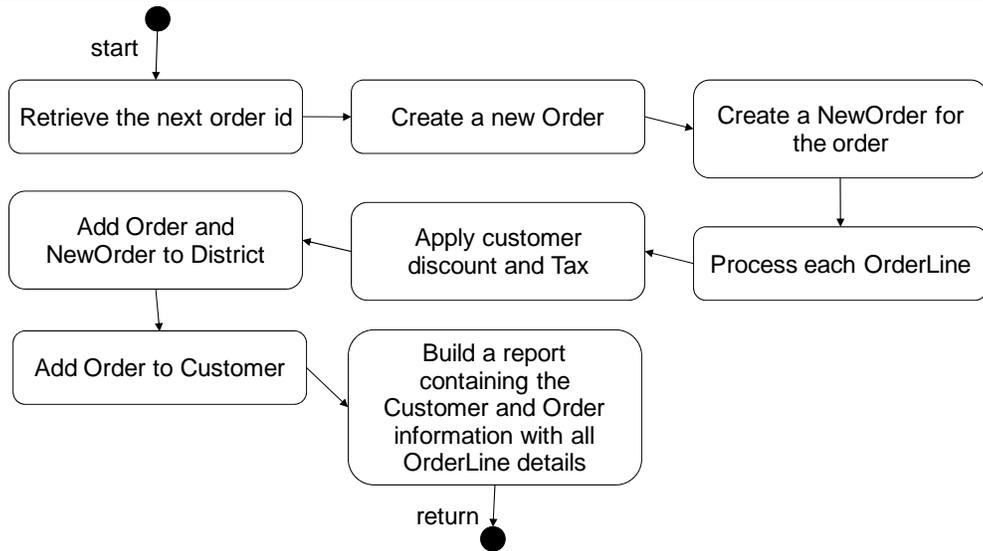
1166 The CSSJ workload does not have any supported configuration parameters.

1167

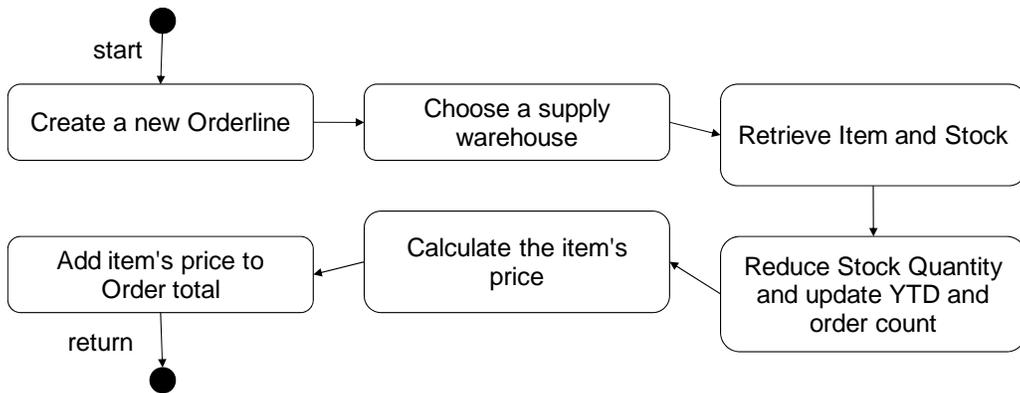
1168 **10.10.6 New Order Transaction**

1169 The input for a New Order Transaction consists of a random district and customer id in the user's
1170 warehouse, and a random number of orderlines between 5 and 15.

1171



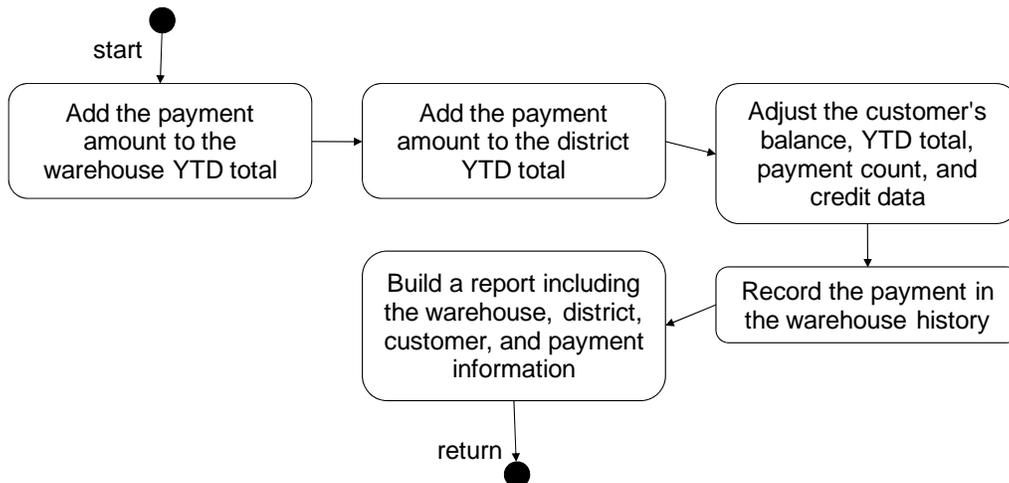
1172
1173
1174



1175
1176

1177 **10.10.7 Payment Transaction**

1178 The input for a Payment Transaction consists of a random district from the user's warehouse, a
1179 random customer id or last name (from either the user's warehouse or a remote warehouse) and a
1180 random payment amount.

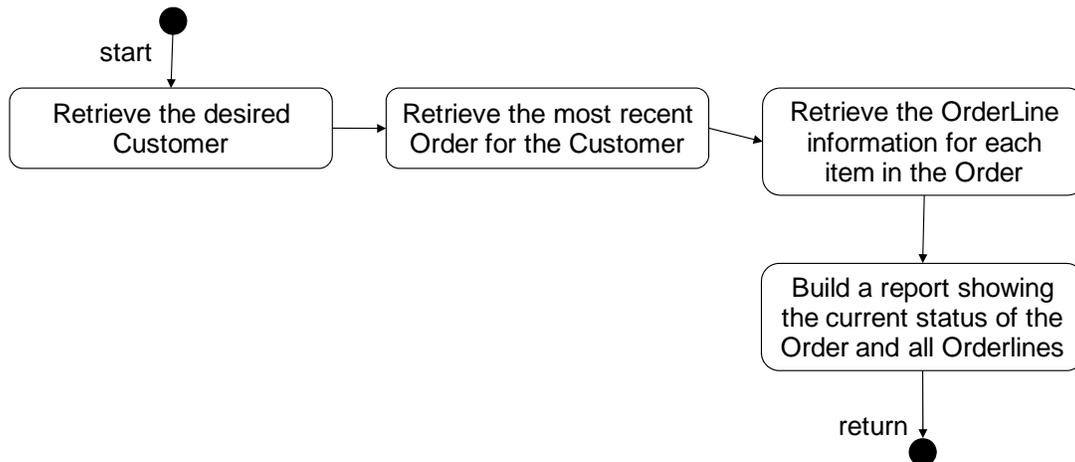


1181
1182

1183 **10.10.8 Order Status Transaction**

1184 The input for an Order Status Transaction consists of a random district and either a customer id or last
 1185 name from the user's warehouse.

1186



1187

1188

1189 **10.10.9 Delivery Transaction**

1190 The input for a Delivery transaction is a random carrier id.

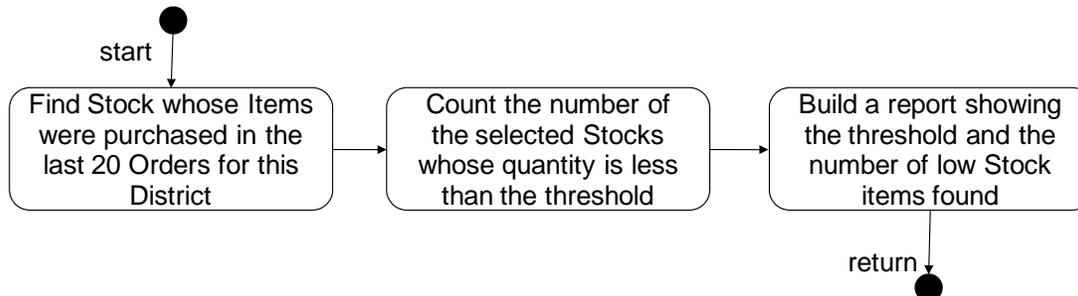
1191

1192 [The activity diagram is work in progress]

1193

1194 **10.10.10 Stock Level Transaction**

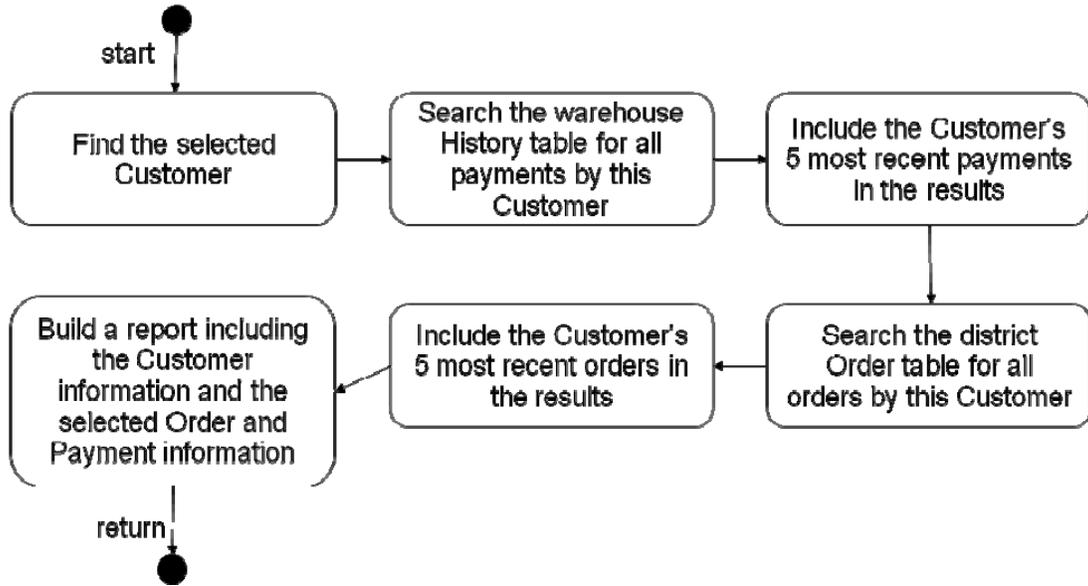
1195 The input for a Stock level transaction is a random district from the user's warehouse and a random
 1196 "low level" threshold between 10 and 20.



1197

1198 **10.10.11 Customer Report Transaction**

1199 The input for a Customer Report transaction consists of a random district from the user's warehouse
 1200 and a random customer id or last name (from either the user's warehouse or a remote warehouse).



1201
1202
1203
1204