

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21

Article type : Research Article

Running title: Assess 5 DBT systems parchment phantom

Assessment of task-based performance from five clinical DBT systems using an anthropomorphic breast phantom

Lynda C. Ikejimba^a, Jesse Salad^a, Christian G. Graff^a, Mitchell Goodsitt^b, Heang-Ping Chan^b, Hailing Huang^c, Wei Zhao^c, Bahaa Ghamraoui^a, Joseph Y. Lo^d, Stephen J. Glick^a

^aUS Food and Drug Administration, 10903 New Hampshire Ave, Silver Spring, MD USA 20993

^bMichigan Medicine, University of Michigan, 1500 East Medical Center Drive, Ann Arbor, MI USA 48109

^cStony Brook Medicine, Stony Brook University, 101 Nicolls Road, Stony Brook, NY USA 11794

^dMedical Physics Graduate Program, Duke University, 2424 Erwin Road, Durham, NC USA 27705

Corresponding author: Lynda Ikejimba

Mailing address: 10903 New Hampshire Ave, WO62-3022, Silver Spring, MD 20993

Email address: Lynda.ikejimba@fda.hhs.gov

ABSTRACT

Purpose: Digital breast tomosynthesis (DBT) is a limited-angle tomographic breast imaging modality that can be used for breast cancer screening in conjunction with full-field digital mammography (FFDM) or synthetic mammography (SM). Currently, there are five commercial DBT systems that have been approved by the U.S. FDA for breast cancer screening, all varying greatly in design and imaging protocol. Because the systems are different in technical specifications, there is a need for a quantitative approach for assessing them. In this study, the DBT systems are assessed using a novel methodology with an inkjet-printed anthropomorphic phantom and four alternative forced choice (4AFC) study scheme.

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1002/MP.14568](https://doi.org/10.1002/MP.14568)

30 **Method:** A breast phantom was fabricated using inkjet printing and parchment paper. The phantom contained 5
31 mm spiculated masses fabricated with potassium iodide (KI)-doped ink and microcalcifications made with
32 calcium hydroxyapatite. Images of the phantom were acquired on all five systems with DBT, FFDM, and SM
33 modalities where available using beam settings under automatic exposure control. A 4AFC study was conducted
34 to assess reader performance with each signal under each modality. Statistical analysis was performed on the data
35 to determine proportion correct (PC), standard deviations, and levels of significance.

36 **Results:** For masses, overall detection was highest with DBT. The difference in PC was statistically significant
37 between DBT and SM for most systems. A relationship was observed between increasing PC and greater gantry
38 span. For MCs, performance was highest with DBT and FFDM compared to SM. The difference between PC of
39 DBT and PC of SM was statistically significant for all manufacturers.

40 **Conclusions:** This methodology represents a novel approach for evaluating systems. This study is the first of its
41 kind to use an inkjet printed anthropomorphic phantom with realistic signals to assess performance of clinical
42 DBT imaging systems.

43
44 **Keywords:** digital breast tomosynthesis, iodine, mass, microcalcification, synthetic mammography

47 INTRODUCTION

48
49 The advent of screening mammography has been one of the driving forces that has resulted in a 39%
50 reduction in breast cancer mortality.¹ Nonetheless, sensitivity and specificity of mammography are
51 limited in certain population groups, most notably women with dense breast tissue and women on
52 hormone replacement therapy.² One of the challenges with conventional full field digital
53 mammography (FFDM) is the structural tissue overlap that can inherently obscure diagnostic features
54 in 2D imaging of the breast. Digital breast tomosynthesis (DBT) is a limited angle tomographic breast
55 imaging modality designed to reduce the superposition of breast tissue, commercially introduced in
56 2011 for combined usage with conventional 2D mammography. As of 2020, statistics from the
57 Mammography Quality Standards Act showed 69% of certified facilities in the US offered DBT.³
58 Recently, synthetic mammography (SM), a method designed to generate a mammography-like image
59 from the DBT stack of image slices, has been introduced with the goal of reducing radiation by
60 eliminating the standard 2D mammography acquisition.⁴

61 As of 2019, five commercial DBT systems have been approved by the U.S. Food and Drug
62 Administration (FDA). The GE Senographe Essential (SenoClaire) (GE Healthcare, Waukesha, WI)
63 has been approved for screening with DBT+SM, the GE Senographe Pristina (GE Healthcare,
64 Waukesha, WI) for DBT+SM, the Hologic Selenia Dimensions (Hologic, Bedford, MA) with
65 DBT+FFDM or DBT+SM, the Fuji ASPIRE Cristalle (Fujifilm, Stamford, CT) system for
66 DBT+FFDM or DBT+SM, and the Siemens MAMMOMAT Inspiration (Siemens, Erlangen,
67 Germany) for FFDM+DBT and DBT alone. Although these systems all perform tomosynthesis, they
68 have many design and operational differences including acquisition geometry, exposure techniques,
69 x-ray tube target and filter, detector type, use of varying reconstruction method, and different levels
70 of radiation dose to the breast. Although clinical trials have demonstrated that DBT can improve breast
71 cancer detection while reducing the false-positive recall rate,⁵⁻⁷ it is unclear how clinical performance
72 depends on the particular DBT system used, since to date no appropriately powered clinical studies
73 have been conducted to compare different commercial DBT systems. Unfortunately, this type of
74 comparison study is difficult to perform due to the high cost and complexity of conducting such a
75 clinical trial. To circumvent some of the limitations of clinical studies, phantom-based methodologies
76 are being developed for system evaluation.

77 Previous studies have described the development of breast phantoms and methodologies to assess
78 imaging systems, and in general the approaches can be broadly classified as virtual or physical. In
79 virtual clinical trials (VCTs), each component of the imaging chain is simulated, including the breast,
80 the imaging system, and the reader.⁸⁻¹⁰ Virtual clinical trials are becoming more common for assessing
81 new technology. Recently, a research group at the FDA conducted the VICTRE (Virtual Imaging
82 Clinical Trial for Regulatory Evaluation) trial demonstrating the use of VCTs in a regulatory application.⁸
83 While such approaches can be efficient and allow for a large number of subjects, they require accuracy
84 in modeling the imaging components. As a result, modeling clinical systems for use with VCTs can
85 prove challenging when processing software is proprietary. The other approach to assessing system
86 performance is with use of physical phantoms. For Quality Control (QC), standard phantoms include
87 the American College of Radiology (ACR) phantom¹¹ and CDMAM phantom,¹² approved for
88 accreditation purposes in the United States and Europe, respectively. While these phantoms are ideal
89 for quick or routine QC testing, they contain signals in a uniform background and thus may not be
90 sufficient for optimization studies, where system performance can change with anatomical

91 complexity.¹³ Physical structured phantoms exist for system optimization such as the Penn¹⁴ and
92 Duke¹³ phantoms, both based on anatomical properties and fabricated through additive manufacturing.
93 In addition, the phantom described by Cockmartin *et al.*¹⁵ consists of acrylic spheres of varying sizes
94 in a water bath. Masses and microcalcifications (MCs) can be added for task-based assessment of
95 mammography systems. Although these phantoms are very useful and unique in their approach, the
96 materials used can be somewhat limited in realism. There is a need for a realistic, anthropomorphic
97 physical breast phantom that can be used for QC, system optimization, and regulatory evaluation of
98 system effectiveness.

99 Our research group at the FDA has previously developed a methodology to objectively assess task
100 performance of breast x-ray imaging systems using a realistic anthropomorphic breast phantom.¹⁶ This
101 phantom uses a novel inkjet printing approach to fabricate a physical phantom based on a virtual breast
102 model. In addition, diagnostic features such as realistic MC clusters and extended masses can be
103 inserted into the phantom. Region of interest (ROI) and volume of interest (VOI) images containing
104 diagnostic features can then be extracted for use in performance assessment studies. A previous
105 proceedings paper¹⁷ from our group described very preliminary work. The present submission
106 represents a greatly expanded study with substantive changes. The present paper contains 5 more
107 figures, two additional tables, and a more rigorous statistical analysis. In addition, the present paper
108 includes a substantial additional reader study evaluating detection of microcalcification clusters. For
109 this, we have used a novel method for fabricating microcalcifications, as well as an approach for
110 generating a template modeling random microcalcification clusters that was inserted into the 3D paper
111 phantom. This study evaluating microcalcification detection was not included in the conference
112 proceedings.

113
114 The goal of this study was two-fold. First, we endeavored to explore whether this task-based phantom
115 assessment methodology was feasible for use on each currently available commercial DBT systems.
116 Each commercial system uses proprietary image processing software, and it is known that imaging of
117 some phantoms produces image artifacts that would not occur during imaging of patients.¹⁸ The
118 second goal was to use this task-based assessment methodology to compare performance achieved
119 with the five FDA-approved commercial DBT systems under DBT, FFDM, and synthetic
120 mammography (when available) imaging modes. Phantom images from each system were acquired

121 using the automatic exposure control settings for that system. Thus, all phantom acquisition settings,
122 and subsequent radiation dose levels, were dictated from the manufacturer settings. The results of this
123 comparison could provide insight into the different tradeoffs associated with varying operational and
124 design strategies used with different commercial DBT systems.

125 **METHODS**

126 | **2.1 Breast Phantom Fabrication**

127 The breast phantom used in the study was a custom-made, 3D parchment paper phantom fabricated
128 using an inkjet printing process described in detail previously.^{16, 19} A digital breast phantom was first
129 created through analytical modeling by making a shell for skin, then dividing the interior into
130 fibroglandular and adipose compartments, and finally adding ligaments and blood vessels.²⁰ The
131 digital phantom modeled a breast with 28% fibroglandular density, representing a dense breast with
132 an extensive parenchymal structure. The digital phantom was then compressed to 4 cm thickness
133 modeling the administration of the breast compression paddle. The model was sampled at 70 μm
134 isotropic voxel resolution to match the thickness of the paper onto which it would be printed. Fiducial
135 markers were inserted into every slice of the digital phantom to assist with proper registration of the
136 printed sheets. The fiducial markers were designed as rings placed on the medial and lateral sides of
137 the breast, separated by a fixed distance from each other and from the chest wall. These can be seen
138 in Figure 1.

139 Inkjet printing was used to realize the digital phantom. For the fibroglandular tissue, a custom ink
140 solution was formulated by combining a ratio of 2/3 dye ink to 1/3 iohexol with an iodine
141 concentration of 350 mg/mL. The final fibroglandular ink had an iodine concentration of 117 mg/mL.
142 For the fat tissue, parchment paper served as the background onto which the ink for fibroglandular
143 regions was printed. Printing was done on an Epson WF-3620 inkjet printer (Epson America, Long
144 Beach, CA) with refillable ink cartridges. Before printing, each channel was assessed to ensure it could
145 print a single color without “contamination” from other channels. To do this, a line pair pattern was
146 printed with a single color channel onto parchment paper which, being slightly hydrophobic, would
147 allow visualization of individual ink droplets. The samples were then examined under either a 5x
148 optical microscope or jeweler’s loupe. No droplets were observed from other color channels within
149 each line or along the line edges, where color mixing would be most evident. This was repeated for

150 every photographic setting available on the printer (“Photo glossy,” “High quality,” “Medium
151 quality,” etc.), totaling about 20 different settings. Printing proceeded only with a printer setting where
152 no mixing was observed that also provided the greatest print speed. For a 4 cm compressed breast, a
153 total of 571 sheets were printed. To house the sheets, a custom-made container was designed
154 consisting of a 6 mm sheet of acrylic as a base; two posts extending vertically from the base and
155 measuring roughly 6 mm in diameter and 76 mm in height, placed close to the chest wall; and an
156 additional 6 mm thick sheet of acrylic on top to provide minor compression. The diameter of each post
157 was set to match the diameter of the fiducial marker ring.

158 As previously mentioned, the fiducial markers were printed with each slice to ensure proper
159 registration of the sheets. A hole was punched in each sheet through the center of the ring, using a
160 custom hole punch designed for this purpose. The posts were then passed through the holes as each
161 sheet was stacked. The holes were visually inspected and tested for fit on each sheet before proceeding
162 to the next sheet.

163 | 2.2 Insertion of Masses

164 Masses were fabricated in a similar inkjet printing manner. A three-dimensional mass was first
165 digitally created using an implementation of lesion modeling software.²¹ These masses had an
166 approximate diameter of 5 mm with spiculations emanating outward from the center (see Figure 1(a)).
167 Select parameters used in the mass generation included *number of initial segments* = 1358, *maximum*
168 *number of neighborhood segments* = 0.98, and *mean radius decrease* of 0.89 as described in the work
169 by de Sisternes *et. al.*²¹ The mass was duplicated to create an insert consisting of 18 identical masses
170 arranged in a grid with 20 mm spacing center-to-center in the x- and y-directions, centered within the
171 same z-slice. The mass insert spanned 70 pages in the z-direction, as the thickness of each sheet was
172 70 μm and the masses were 5 mm in diameter. This design maximized the number of ROIs that could
173 fit within the breast while allowing sufficient space between each mass. The central slice of the mass
174 insert contained three markers for BBs, used for automating ROI and VOI extraction. This insert was
175 placed via pixel substitution into the central portion of the digital breast phantom, the region with the
176 greatest area. To obtain more samples, four sets of mass inserts were created by shifting the entire grid
177 of masses and BB markers in unison by 2 mm to 5 mm in the x- and y-directions, remaining centered

178 in the same slice. As a result, a total of 72 ROIs containing a mass were created with unique
179 background locations.

180 To print the masses, a new type of ink was synthesized. When higher concentrations of iohexol were
181 used in the ink, the print heads were prone to clogging. To reduce clogging, a different ink/iodine
182 solution was required to achieve a sufficiently high iodine concentration for the masses. A saltwater
183 solution was made by dissolving potassium iodide (KI) in water at a concentration of 300 mg/mL.
184 This was then mixed with ink at a ratio of 2/3 KI saltwater to 1/3 dye ink. The resulting ink for the
185 masses had a KI concentration of 200 mg/mL. The mixture was placed in a separate color cartridge in
186 the printer, different from the one containing the iohexol-based ink. Prior to printing, the
187 fibroglandular and mass tissues were recolored in the digital model to match the cartridge color they
188 would be printed with. This was done using GIMP (GIMP v 2.10.10, <http://gimp.org>), a freely
189 available image manipulation program. The program allowed selected pixel values to be printed from
190 a specific cartridge, enabling simultaneous printing of the two tissue types. Each slice containing the
191 mass was printed onto a sheet of paper, yielding a subsection stack with 70 sheets of paper roughly 5
192 mm thick in total. The four 70-sheet mass inserts were each printed, producing four physical stacks of
193 mass inserts. Figure 1 shows a side-by-side comparison of (a) the 3D lesion model, (b) the arrangement
194 of mass lesions positioned within the breast with the three markers for BB locations, (c) the center
195 slice of the breast phantom with inserted masses and the fiducial rings for sheet registration, and (d)
196 a printed sheet of the same slice with the BBs in place and posts through the rings.

197

198 | 2.3 Microcalcifications

199 A new MC insert was created using an improved method similar to one previously described.¹⁹ Using
200 MATLAB, a template was first designed to mark locations where clusters would be placed. The
201 template consisted of 5-row by 9-column grid of 5-mm-diameter circles, each with an MC cluster.
202 Within each cluster, the locations of the specks were randomly generated with a buffer around each
203 speck to ensure they are placed within the 5 mm circle and do not overlap. The specks were made by
204 combining calcium hydroxyapatite (HA) powder with the binding agent polyvinylpyrrolidone and
205 compressing the mixture into a tablet using a mechanical press. The tablets were then crushed and
206 separated by size using differential sieving. The resulting specks ranged in size between 150 μ m and

207 180 μm . Five specks were placed into each of the pre-designated locations in the 5 mm circle,
208 described above. The clusters were spaced 15 mm apart in x- and y-directions, forming a total of 45
209 clusters. Fiducial markers were included in the template to assist with extraction. The insert was sealed
210 with double-sided tape and more parchment paper. The excess paper and tape were trimmed to allow
211 the insert to fit within the breast boundary of the printed phantom. The process of fabricating and
212 inserting the MC template can be seen in Figure 2.

213

214 | 2.4 Image Acquisitions

215 Images of the phantom were acquired on five commercially available DBT systems: Hologic Selenia
216 Dimensions at Sibley Memorial Hospital in Washington, DC; GE Senographe Essential (SenoClaire)
217 at University of North Carolina-Chapel Hill (UNC) in Chapel Hill, NC; GE Senographe Pristina at
218 University of Michigan in Ann Arbor, MI; Siemens MAMMOMAT Inspiration at the State University
219 of New York (SUNY) in Stonybrook, NY; and Fujifilm Aspire Cristalle in Stamford, CT. Imaging
220 was performed with FFDM, DBT, and SM modalities, with the exception of the GE Senographe
221 Essential which lacked SM capability. The technical specifications of the systems are provided in
222 Table 1.

223 The imaging parameters were determined by the beam conditions used under automatic exposure
224 control (AEC) for each system. As a result, the x-ray tube settings ranged from 29 kVp to 34 kVp for
225 tube voltage, 40 mAs to 180 mAs for total mAs for all projections, and 1.2 mGy to 2.3 mGy for
226 average glandular dose (AGD) as reported from the system display. A summary of the acquisition
227 parameters is given in Table 2. Note that these are the dose levels reported by the vendor. Because
228 manufacturers may use very different assumptions for their calculation of displayed AGD, the system-
229 displayed doses might have limited comparability. An independent dose calculation was performed
230 using the IEC recommended²² procedure with the Dance method.²³⁻²⁵ Using the beam settings, the
231 AGD was calculated to a reference phantom of 40 mm PMMA with the equation

$$232 \quad D_T = K_E g c s T$$

233 where K_E is the entrance surface kerma; g , c , and s are factors dependent on the target, filter, and half
234 value layer; and T is variable for commercially available DBT systems. This calculation is provided
235 in the last column under Ref. AGD.

236

237 To determine the appropriate beam parameters, an image of the phantom was taken under AEC
238 conditions. The phantom was positioned with its chest wall extended off the edge of the detector cover,
239 and posts flush with the compression paddle, as illustrated in Figure 3. This configuration was
240 necessary to accommodate the height of the posts. Once the beam parameters were determined on
241 each system, imaging was performed with the phantom repositioned with the overhanging lip against
242 the detector cover and the compression paddle sitting atop the posts. The acquisition parameters were
243 manually set and used to image the phantom with and without signals.

244

245 To image the masses, a central subsection of the phantom was replaced by the 70-sheet stack with
246 printed masses. All four 70-sheet stacks containing masses were inserted and imaged one after the
247 other in this manner on all systems, with the exception of the GE Essential since the fourth stack was
248 not completed at the time of imaging. To image the MCs, the MC insert was placed between the central
249 two slices of the phantom. Multiple acquisitions were taken in order to sample different background
250 locations where the MC insert was shifted in a random manner between shots. For the signal-absent
251 data, a single shot was taken of the phantom sheets without the inserted masses or calcifications.
252 Additional scans were not necessary since the phantom background would remain the same.

253 ROIs and VOIs were extracted automatically from the images using a custom MATLAB program.
254 Background ROIs were randomly selected from within the breast volume. This was achieved by
255 extracting overlapping ROIs within the breast boundary in a raster fashion, producing on average
256 between 350 and 450 ROIs. From these, enough background ROIs were randomly selected to equal
257 3x the number of signal present ROIs for a case cohort, considered all the cases for a given vendor,
258 modality, and signal (e.g. Hologic FFDM masses). With this method, it is highly unlikely that an exact
259 background ROI would be selected that corresponded to any signal present locations. To prevent
260 learning the backgrounds, each signal absent ROI was randomly rotated by 90, 180, or 270 degrees.
261 Mass ROIs were also randomly rotated by 90, 180, or 270 degrees and had an additional search
262 component, whereby the center of the lesion could be located anywhere within a 10 mm x 10 mm area
263 within the middle of the ROI. For the MCs, ROIs were rejected if the cluster was outside of the breast
264 or too close to the breast boundary. The ROIs measured 15 mm x 15 mm for MCs and 20 mm x 20

265 mm for masses. DBT VOIs consisted of 9 reconstructed DBT slices of 1-mm thickness. For MCs,
266 between 92 and 106 ROIs were extracted per modality per system, and for masses between 44 and 62
267 ROIs.

268 | 2.5 Reader Study

269 A 4AFC reader study was conducted to evaluate the detection of masses and MCs with each modality.
270 Reading was conducted in a dark room designed for human reader studies. Ambient light was
271 minimized to model conditions in a clinical reading room. Reading was performed on a 30" 6MP
272 Coronis Fusion (6MP DL MDCC-6130, Barco NV, Kortrijk, Belgium) medical display calibrated to
273 DICOM grayscale standard display function. The display contained an active screen area of 26" × 16"
274 with 3,280 × 2,048 pixels. The display was operated in Diagnostic mode grayscale standard display
275 function (GSDF) with 300 cd/m² maximum luminance. The lighting in the room was kept low to
276 model that of a clinical reading room. Ambient light from the ceiling did not cause glare on the
277 monitor. The illuminance from the display was measured to be 2.71 lux. ROIs were displayed at a 1:1
278 magnification.

279 Scoring was performed by seven non-radiologist readers familiar with the type of images and done in
280 a dark room adapted for such studies. Reading was facilitated by the Foursquares software.²⁶ The
281 program consists of four windows, each displaying an ROI or VOI. Only one of the windows contains
282 a signal-present ROI or VOI, and the three others contained background images only. It was the
283 objective of the reader to select the correct window. A "cue" image was presented next to the
284 Foursquares program with a mass or one MC cluster in a uniform background; this provided the reader
285 with an example of the true signal. For each experimental condition, ROIs within the case cohort were
286 presented in a randomized order with respect to location within the breast phantom. As previously
287 mentioned, a case cohort consisted of all the ROIs for a given vendor, modality, and signal – for
288 example 62 ROIs for Hologic DBT masses. Readers were medical physicists that were experienced in
289 the given tasks with the 3D paper phantom. The task for mass detection was a signal location unknown,
290 and readers were instructed to perform some search. The ROIs for the masses were 20 mm x 20 mm,
291 and the center of the mass could be located anywhere within the central 10 mm x 10 mm area of the
292 ROI. The task for MC cluster detection represented a signal known exactly, and readers were
293 instructed to look within the center of the ROIs for the clusters. For DBT, observers were instructed
294 to scroll through all slices of the VOIs before finalizing their selection. When displaying the ROIs, the

295 Foursquares program calculated a default window width and level (W/L) based on the range of pixel
296 values within the image. All ROIs had 16-bit unsigned integers. Observers were allowed to adjust the
297 display W/L as needed. The study was performed over multiple sessions, and breaks were encouraged
298 after 25 minutes of reading to avoid observer fatigue. Prior to the study, reader training was conducted.
299 During this training, the reader was familiarized with the study objective and software interface. The
300 reader scored images from training data independent from the testing data, with supervision from the
301 investigators and feedback provided after each response. During the study, the reader scored training
302 images for each signal, modality, and vendor before scoring the corresponding testing images.
303 Feedback was given after every user response during both training and testing phases. A summary of
304 the number of ROIs scored for each modality is presented in Table 3.

306 Results were computed using the iMRMC²⁷ package in R Studio (Version 1.1.463). The proportion
307 correct (PC) was calculated as the ratio of correctly selected ROIs to the total number of ROIs scored.
308 In a 4-AFC, the PC for random guessing would be 0.25. The variance of the PC was calculated directly
309 from each trial, and accounts for all correlations across readers and cases. In summary, the variance
310 of PC was computed using u-statistics in the iMRMC package. To do this, the constituent parts of the
311 unbiased variance estimate were first calculated. From these, the statistical moments and their
312 associated coefficients may be derived. Finally, the variance is computed as the inner product of the
313 moments and coefficients. More details of this approach may be found in Gallas et al.²⁸ Using the
314 variance, the 95% confidence intervals were then calculated as a product of ± 1.96 and the standard
315 error. An estimate of detection relative to FFDM as a baseline was computed as ΔPC , defined as ΔPC
316 = $PC_i - PC_{FFDM}$, where PC_i is the PC of a given modality i (either DBT or SM). Since all ΔPC are
317 relative to FFDM, a value of $\Delta PC > 0$ indicates an improvement in signal detection over FFDM, while
318 a value of $\Delta PC < 0$ indicates a reduction.

319 To determine if differences in PC were statistically significant, p-values and significance level α were
320 required. The p-value was derived via t-table and calculation of the test statistic, computed for every
321 pairwise comparison: signal, vendor, modality. Then, the p-value can be compared with a Bonferonni-
322 corrected α to determine significance. Computation of the p-values via t-table required an estimate of
323 the number of degrees of freedom (df). The df was estimated as the number of readers, under the

324 assumption that the readers would contribute most to the variability in the results. In addition, having
325 fewer df yields a more conservative estimate of p-values, reducing the likelihood of Type I errors.
326 While a threshold value of $\alpha = 0.05$ is typically used to reject the null hypothesis, the Bonferonni
327 correction is needed in order to account for the increased likelihood of finding statistical significance
328 when there are multiple experiments. In MRMC study design, multiple experiments can arise from
329 comparing the PC across different modalities, vendors, or signals. Thus, having multiple comparisons
330 may require determination of a new threshold for significance α/m , where m is the number of
331 *independent* comparisons. In this study, there were three pairwise comparisons per signal for each
332 vendor with DBT, FFDM, and SM (DBT vs FFDM, DBT vs SM, and FFDM vs SM), resulting in a
333 threshold of $\alpha = 0.05/3 = 0.0166$. The GE Essential (SenoClaire) system did not have SM. Only one
334 pairwise comparison was made (DBT vs FFDM), so the threshold remained $\alpha = 0.05$.

335

336

RESULTS

337 Sample images are presented in Figure 4 for the masses; side-by-side comparisons are given for DBT,
338 FFDM and SM (unless unavailable) for each of the five systems. Arrows indicate the locations of the
339 masses. Some masses become difficult to detect when going from 3D to 2D, especially comparing
340 DBT to SM. For DBT, the masses appear most conspicuous for the systems with the largest gantry
341 spans, namely Siemens at 50°, GE Essential at 25°, and GE Pristina also at 25°.

342

343 The reader scores are presented in Table 4 for all systems, modalities, and signals. The values
344 displayed are the reader-averaged PC with the 95% confidence interval (CI_{95}) in brackets. The reader
345 averaged scores are presented for masses in Figure 5 and for MCs in Figure 6. The error bars represent
346 one standard deviation accounting for all sources of variability (reader and case). Results are given for
347 DBT with the red bars, for FFDM with green, and for SM with blue. The gantry span is indicated in
348 degrees above each subplot, and the average glandular dose is given below each bar. The pair-wise
349 comparisons with an asterisk denote statistical significance with the Bonferroni correction.

350

351 For masses, the highest performance was achieved overall with DBT. Furthermore, a relationship was
352 observed between overall PC and gantry span. The PC increased from 0.72 ± 0.05 for Hologic with

353 15° (and a similar score for Fuji) to 0.91 ± 0.03 for Siemens with 50°. The difference in performance
354 between DBT and FFDM was found to be statistically significant for all systems except Fuji and
355 Hologic—for GE Essential (SenoClaire) the difference between DBT and FFDM had a p-value of
356 0.05, right on the threshold of $\alpha = 0.05$ as the Bonferroni correction was not necessary. The difference
357 between DBT and SM was statistically significant for all systems except Fuji. The difference in
358 performance between FFDM and SM was not found to be statistically significant. For FFDM, the PC
359 for mass detection varied minimally across systems and different dose levels ranging from 0.61 ± 0.06
360 (Fuji) to 0.64 ± 0.04 (Siemens). Comparably, for SM the PC ranged from 0.52 ± 0.05 (Hologic) to
361 0.65 ± 0.05 (Fuji). Although the SM image is typically produced using the DBT dataset, no trend was
362 observed between the scores for SM and the gantry span of the system.

363

364 For MCs the highest PCs were observed with FFDM and DBT, both having similar scores. Overall
365 scores for MCs were observed to be higher than those of the masses; with DBT, the PC ranged from
366 0.84 ± 0.02 (Siemens) to 0.95 ± 0.02 (GE Pristina), while FFDM ranged from 0.78 ± 0.03 (Siemens)
367 to 0.94 ± 0.02 (Hologic). Performance with SM was lowest, with PC ranging from 0.39 ± 0.04
368 (Siemens) to 0.63 ± 0.05 (Fuji).

369

370 The ΔPC relative to FFDM is provided in Figure 7 for all vendors. Results are given for both masses
371 and MCs side-by-side, with DBT in red and SM in teal. For masses, DBT consistently yielded a
372 positive ΔPC greater than 0.10. This indicated that detection of masses was higher with DBT than
373 with FFDM regardless of system configuration, for the present task. For MCs, however, moderate
374 improvement was observed with DBT relative to FFDM, with all $\Delta PC \leq 0.06$. Conversely, SM yielded
375 negative ΔPC in all but one comparison, for both mass and MC detection. Moreover, the greatest
376 difference was observed for MCs, indicating that for this size of MCs worse performance will be
377 obtained with SM compared to FFDM.

378

379

DISCUSSION

380

381 The commercial systems investigated in this study varied greatly in design and how they operate.
382 Differences in x-ray spectra, detector type (direct versus indirect-conversion), detector pixel and
383 reconstructed voxel size, acquisition geometry, reconstruction method, image post-processing
384 methods, and step-and-shoot versus continuous gantry motion could have affected performance
385 depending on the task. In addition, the GE SenoClaire and Pristina systems utilize anti-scatter grids
386 while acquiring DBT projections, while the other systems do not. Using acquisition parameters
387 determined from the AEC software of each system, the estimated average glandular dose (AGD)
388 varied between systems, sometimes substantially. Owing to the many system parameters affecting
389 image quality, it is difficult to determine which factors affect performance the most. Therefore, it is
390 difficult to make conclusions on how specific design and acquisition parameters affect the results here
391 and the current study should not be considered a vendor comparison. Nevertheless, certain general
392 trends can be observed in this study. Furthermore, we believe that the resulting phantom images and
393 computed task performance demonstrate that this methodology can be utilized on all clinically
394 available DBT systems.

395 For the mass detection study, task performance was higher with DBT than with FFDM or SM, and the
396 difference in PC was statistically significant for two systems. This finding concurs with other phantom
397 studies using structured backgrounds. For example, Cockmartin *et al.*¹⁵ showed clear improvement of
398 DBT over FFDM in detecting masses of various sizes. From a subjective visual impression comparing
399 mass lesions with DBT and FFDM, we concluded that mass conspicuity is generally improved with
400 DBT over all systems (see Figure 4), which likely relates to the improved diagnostic performance of
401 DBT reported in many lab based clinical studies.²⁹ Additionally, mass detection performance of the
402 DBT systems trended with increased gantry span, with lowest PC from Hologic and Fuji (both 15°),
403 then both GE systems (25°), and highest for Siemens (50°). A subjective visual impression clearly
404 shows that the wider angle DBT systems provide improved mass lesion conspicuity. This finding
405 concurs with other phantom studies that show a strong correlation between mass visualization and
406 increased gantry angle.³⁰⁻³⁵ In particular, previous reserachers³¹⁻³³ examined the relationship between
407 mass detectability and DBT gantry span from 15° or 16° up to 60°, along with a number of other
408 acquisition variables. Within the context of that work, the results in the present paper align with the
409 trends observed when matched with a similar gantry span and number of projections. Similar reader
410 scores for the mass detection task were observed between FFDM and SM, a finding that was in

411 agreement with other studies. For example, Mackenzie *et al.*³⁶ conducted a virtual clinical trial that
412 showed similar mass detection performance with SM and FFDM. Although there are many variations
413 of SM algorithms implemented by different vendors, clinical studies also seem to suggest similar
414 performance of SM and FFDM for the detection and diagnosis of mass lesions.^{4, 37, 38} For FFDM, the
415 performance of mass detection did not appear to be impacted by dose with PCs of 0.61, 0.64, 0.64,
416 0.62, 0.61 measured for reference AGD values of 0.84, 0.90, 0.93, 1.09, and 1.63 mGy respectively.
417 This finding was in agreement with previously published detection studies using hybrid clinical data,
418 i.e. normal patient data inserted with simulated lesions, by Svahn *et al.*³⁹ and Timberg *et al.*⁴⁰

419 For the MC detection task, performance was similar with DBT and FFDM, with both modalities
420 providing improved performance over SM. Unlike results with mass detection, no clear trend was
421 determined between reader performance and the system geometry. Chan *et al.*⁴¹ observed a trend of
422 decreasing MC detection sensitivity and conspicuity with increasing scan angle for acquisition with
423 uniform angular increments; however, the DBTs at all scan angles were acquired with a step-and-
424 shoot system, the same x-ray spectrum, dose, and detector in that study. The differences in the many
425 factors among the DBT systems for the current study may have reduced the dependence of MC
426 detection on the scan angle. Due to their size, detection of MCs are probably more limited by quantum
427 noise, whereas the detection of mass lesions are probably more limited by overlapping structure, thus
428 explaining the greater improvement of DBT over FFDM for the mass detection task. This observation
429 was discussed in detail by Burgess *et al.*⁴² showing that smaller microcalcification-like objects have
430 different contrast-detail characteristics than larger mass-like objects. Of course, it is difficult to be
431 certain of this trend because other factors such as detector type and pixel size might also contribute to
432 differences in performance. Unlike the mass detection task, performance of DBT and FFDM with
433 MCs was significantly higher even with the Bonferroni correction for most of the systems tested
434 compared to SM. This result concurs with the findings of Mackenzie *et al.*⁴³ who showed that detection
435 of subtle microcalcifications was significantly reduced with SM as compared to DBT and FFDM using
436 simulated lesions inserted into clinical data. Most clinical studies to date have shown that detection of
437 MCs is comparable between SM and FFDM, both alone and with DBT.^{37, 44, 45} However, it is important
438 to note that often these studies include a range of MC sizes present in clinical patient data. In the
439 present study, the sizes of the MCs were restricted to a range of 150 μm to 180 μm to interrogate
440 performance with a challenging task. Results showed that detection of MCs were inferior for SM, the

441 lowest PC scores for this size range. While some clinicians have indicated preference for SM when
442 viewing MCs, maybe due in part to over-enhancement for larger MCs, it is possible that the smallest,
443 more subtle MCs may not have been conspicuous on SM.

444

445 In Figure 5 and Figure 6, statistically significant differences are indicated with denoted asterisks. To
446 account for multiple comparisons, Bonferroni correction was used. Increasing the number of
447 comparisons m yields a lower threshold for significance α , representing a more conservative approach
448 to determine statistical significance. While this decreases the chance of Type I errors when rejecting
449 the null hypothesis, it also increases the chance of Type II errors. In this study, a high m value can be
450 justified if comparisons were made across many systems and modalities. However, a value of $m = 2$
451 or $m = 3$ may be appropriate since the comparisons were mainly made between 2 or 3 modalities from
452 a single manufacturer. In practice, the number of independent comparisons can be difficult to
453 determine, since the same object is imaged across the systems and the same readers are used for
454 assessing the images. It is important to use proper judgement when applying the correction.

455 The present study examined mass and MC detection viewed by each modality alone. In clinical
456 practice, the current standard of care for breast cancer screening in the US is to observe either a
457 conventional 2D FFDM study alone, a combination of a 2D plus 3D DBT study, or a 3D DBT study
458 alone in the case of Siemens. However, there is a growing trend towards screening with a single
459 modality to minimize radiation exposure to patients and to reduce reading time for exams. For this
460 reason, it is of interest to evaluate detectability of a single modality. The results reported herein
461 suggest that subtle MCs could be missed if reading SM images alone without also reading DBT
462 images.

463 In this study, systems were compared at the AEC dose levels for each machine, rather than at a fixed
464 dose across the systems. This was done for two reasons. Firstly, according to each vendor, the AEC
465 setting represents the optimal beam conditions for imaging of a specific breast. That is, the AEC
466 parameters are optimized to achieve a certain image quality on a given system. Secondly, operating at
467 an *a priori* fixed dose could result in an advantage or disadvantage for the system, depending on if the
468 AEC dose is respectively lower or higher than the selected dose. In this study, the reference AEC
469 doses ranged from 0.89 mGy to 1.86 mGy for DBT and from 0.84 mGy to 1.63 mGy for FFDM,

470 representing almost a two-fold increase from the lowest to the highest dose. If the dose was fixed to
471 an arbitrary value, it is possible that results would change in a way that was different for each system.
472 This could yield scores of reader performance that may not be reliable, particularly for dose-limited
473 tasks such as MC detection. For these reasons, the imaging data was collected under the standard AEC
474 conditions for each system.

475 The methodology presented can be a useful tool for routine QC, system optimization, or comparing
476 task-based performance between different imaging systems.⁴⁶ Currently, a widely accepted approach
477 for QC involves using the ACR mammography accreditation phantom or similar phantoms.^{47, 48} The
478 ACR phantom requires subjective reading of signals by human observers (i.e., the reader knows
479 beforehand that the signals are present, and they are asked to record whether the signal is visualized).
480 Less subjective, automated methods for reading the ACR phantom are available;⁴⁹⁻⁵¹ however, they
481 are not typically used. CNR measurements can also be used for QC, but CNR does not account for
482 pixel size or task. If model observers are developed, they may be utilized in the present methodology
483 for a fast, task-based quantitative approach to QC. For QC purposes the phantom can be
484 designed to incorporate various known signal types for detection (e.g. fibers, specks, and so on),
485 uniform regions for noise power spectra, additional BBs for point spread functions and azimuthal
486 spread functions, and other features for automatic analysis. This methodology can lend itself to system
487 optimization whereby optimal acquisition or reconstruction parameters may be determined based on
488 mass or MC detectability. This has the advantage that the results are task-specific. To improve with
489 the imaging workflow for optimization studies, the current design of the support plate can be modified
490 to push the posts towards the chest wall, removing the need for separate shots with the posts extending
491 from the detector edge. Lastly, this methodology has promise in regulatory applications for potentially
492 expediting the review process.

493 While this study was large in its scope, there were a few limitations. First, the phantom modeled a
494 single patient anatomy. A greater range of background anatomy was simulated by placing the signals
495 into different regions of breast parenchyma, but a future study could involve the use of multiple breast
496 phantoms to increase background variability even more. This would then yield a higher number of
497 ROIs for both MCs and masses. As previously mentioned, the phantom modeled a dense breast, so it
498 is possible that overall system performance may differ for a fatty breast more typical of the general

499 population. In addition, the phantom used here was not strictly based on patient data, although it is
500 visually similar. However, it is unclear if phantom realism would impact reader performance
501 differently on different modalities. Another limitation is that the signals represented only one size and
502 shape for the masses and one size range for the MCs, and the print density of the masses was adjusted
503 to make the task challenging. Future studies could investigate performance with both benign- and
504 malignant-appearing lesions, match KI attenuation with known tissue attenuation in the mass, and
505 fabricate MCs comprising different materials⁵² other than calcium hydroxyapatite.⁵³ Additionally, all
506 the MCs were contained within one slice. It is not clear if this would benefit a particular system
507 because of differences in axial resolution. Another potential limitation is that the ROIs were displayed
508 at a 1-to-1 magnification. Because the systems have different voxel sizes, this resulted in the ROIs
509 being displayed at different physical sizes. Consequently, it is possible that displayed ROIs could be
510 smaller than what a radiologist would use if she did not employ additional image magnification.
511 Finally, the readers were not radiologists. Using non-radiologist readers may not have affected the
512 scores since the detection tasks were relatively simple.⁵⁴ Nonetheless, human readers can suffer from
513 observer fatigue, and the presence of different skill levels can cause intra- and inter-observer
514 variability. More experienced readers were also observed to perform better than less experienced
515 readers for some tasks. To circumvent variability due to readers, model observers are being developed
516 for these types of tasks and can be used to minimize variance and reduce the reading time.⁵⁵ Still,
517 challenges associated with the wide scale use of phantoms of this type include: reproducibility of the
518 phantom printing process, similarity of multiple phantoms, long-term use of ink with high salt
519 concentration on desktop inkjet printers, and reproducibility across multiple printer brands.

521 CONCLUSION

522
523
524 In this work, we demonstrated the use of an anthropomorphic breast phantom to objectively assess
525 task-based performance of different commercial breast imaging systems. The phantom was imaged
526 on five commercially available DBT systems across four states, and scans were collected with masses
527 and MCs inserted. A 4AFC observer study was conducted to assess performance with FFDM, DBT,

528 and SM. For masses, overall detection was highest using DBT, with an improvement observed with
529 increased gantry span. For MCs, performance was highest with DBT and FFDM and worse with SM.
530 This study is the first of its kind to use a physical inkjet-printed anthropomorphic phantom to assess
531 clinical performance of all commercially available breast imaging systems.

532 ACKNOWLEDGEMENTS

533 The authors would like to acknowledge the help of Dr. Guo Zhang with the Foursquare software
534 and the help of Andrea Kim with lesion model visualization.

535 FUNDING

536 This work was supported by a Critical Path grant from the Center for Devices and Radiological
537 Health, with a fellowship administered by the Oak Ridge Institute for Science and Education through
538 an interagency agreement between the U.S. Department of Energy and the U.S. Food and Drug
539 Administration. The mention of commercial products, their sources, or their use in connection with
540 material reported herein is not to be construed as either an actual or implied endorsement of such
541 products by the Department of Health and Human Services.

542 DISCLOSURES

543 H.-P. C. and M.G. have research collaboration with GE through an institutional grant not related to
544 the current study. All other authors have no conflicts of interest to disclose.

545 REFERENCES

- 546 1. DeSantis CE, Ma J, Goding Sauer A, Newman LA, Jemal A. Breast cancer statistics, 2017, racial disparity in
547 mortality by state. *CA: A Cancer Journal for Clinicians*. 2017;67(6):439-448. doi:10.3322/caac.21412
- 548 2. Carney PA, Miglioretti DL, Yankaskas BC, et al. Individual and Combined Effects of Age, Breast Density, and
549 Hormone Replacement Therapy Use on the Accuracy of Screening Mammography. *Annals of Internal Medicine*.
550 2003;138(3):168-175. doi:10.7326/0003-4819-138-3-200302040-00008
- 551 3. Mammography Quality Standards Act (MQSA) National Statistics. 2020. [https://www.fda.gov/radiation-](https://www.fda.gov/radiation-emitting-products/mqsa-insights/mqsa-national-statistics)
552 [emitting-products/mqsa-insights/mqsa-national-statistics](https://www.fda.gov/radiation-emitting-products/mqsa-insights/mqsa-national-statistics)
- 553 4. Garayoa J, Chevalier M, Castillo M, et al. Diagnostic value of the stand-alone synthetic image in digital breast
554 tomosynthesis examinations. *European radiology*. 2018;28(2):565-572.

- 555 5. Haas BM, Kalra V, Geisel J, Raghu M, Durand M, Philpotts LE. Comparison of tomosynthesis plus digital
556 mammography and digital mammography alone for breast cancer screening. *Radiology*. 2013;269(3):694-700.
- 557 6. McDonald ES, Oustimov A, Weinstein SP, Synnestvedt MB, Schnall M, Conant EF. Effectiveness of digital
558 breast tomosynthesis compared with digital mammography: outcomes analysis from 3 years of breast cancer screening.
559 *JAMA oncology*. 2016;2(6):737-743.
- 560 7. Skaane P, Bandos AI, Gullien R, et al. Comparison of digital mammography alone and digital mammography
561 plus tomosynthesis in a population-based screening program. *Radiology*. 2013;267(1):47-56.
- 562 8. Badano A, Graff CG, Badal A, et al. Evaluation of digital breast tomosynthesis as replacement of full-field digital
563 mammography using an in silico imaging trial. *JAMA network open*. 2018;1(7):e185474-e185474.
- 564 9. Bakic PR, Barufaldi B, Higginbotham D, et al. Virtual clinical trial of lesion detection in digital mammography
565 and digital breast tomosynthesis. In: *Proc. SPIE 10573, Med. Imaging. Phys. Med. Imaging*; 2018:1057306.
- 566 10. Elangovan P, Warren LM, Mackenzie A, et al. Development and validation of a modelling framework for
567 simulating 2D-mammography and breast tomosynthesis images. *Physics in Medicine & Biology*. 2014;59(15):4275.
- 568 11. American College of Radiology. *Mammography Quality Control Manual, Medical physicist's section*. 1999:225-
569 330.
- 570 12. Bijkerk K, Lindeijer J, Thijssen M. The CDMAM-Phantom: a contrast-detail phantom specifically for
571 mammography. *Radiology*. 1993;185:395.
- 572 13. Ikejimba LC, Glick SJ, Choudhury KR, Samei E, Lo JY. Assessing task performance in FFDM, DBT, and
573 synthetic mammography using uniform and anthropomorphic physical phantoms. *Medical Physics*. 2016;43(10):5593-
574 5602.
- 575 14. Carton AK, Bakic P, Ullberg C, Derand H, Maidment AD. Development of a physical 3D anthropomorphic breast
576 phantom. *Medical Physics*. 2011;38(2):891-896.
- 577 15. Cockmartin L, Marshall NW, Zhang G, et al. Design and application of a structured phantom for detection
578 performance comparison between breast tomosynthesis and digital mammography. *Physics in Medicine & Biology*.
579 2017;62(3):758.
- 580 16. Ikejimba LC, Graff CG, Rosenthal S, et al. A novel physical anthropomorphic breast phantom for 2D and 3D
581 x-ray imaging. *Medical Physics*. 2017;44(2):407-416.
- 582 17. Ikejimba LC, Salad J, Graff CG, et al. Assessment of task-based performance from five clinical DBT systems
583 using an anthropomorphic breast phantom. In: *Proc. SPIE Vol. 11513, 15th Int Workshop on Breast Imaging SPIE*;
584 2020:1151305.
- 585 18. Personal communication with Andy Smith of Hologic Inc.
- 586 19. Ikejimba LC, Salad J, Graff CG, et al. A four-alternative forced choice (4AFC) methodology for evaluating
587 microcalcification detection in clinical full-field digital mammography (FFDM) and digital breast tomosynthesis (DBT)
588 systems using an inkjet-printed anthropomorphic phantom. *Medical Physics*. 2019;46(9):3883-3892.
589 doi:10.1002/mp.13629

- 590 20. Graff CG. A new open-source multi-modality digital breast phantom. In: *Proc. SPIE 9783, Med. Imaging. Phys.*
591 *Med. Imaging*; 2016:978309.
- 592 21. de Sisternes L, Brankov JG, Zysk AM, Schmidt RA, Nishikawa RM, Wernick MN. A computational model to
593 generate simulated three-dimensional breast masses. *Medical physics*. 2015;42(2):1098-1118.
- 594 22. IEC. 61223-3-6 EVALUATION AND ROUTINE TESTING IN MEDICAL IMAGING DEPARTMENTS -
595 Parts 3-6: Acceptance and constancy tests of mammographic X-ray equipment used in a mammographic tomosynthesis
596 mode of operation.
- 597 23. Dance D, Skinner C, Young K, Beckett J, Kotre C. Additional factors for the estimation of mean glandular breast
598 dose using the UK mammography dosimetry protocol. *Physics in medicine & biology*. 2000;45(11):3225.
- 599 24. Dance D, Young K, Van Engen R. Further factors for the estimation of mean glandular dose using the United
600 Kingdom, European and IAEA breast dosimetry protocols. *Physics in Medicine & Biology*. 2009;54(14):4361.
- 601 25. Dance D, Young K, Van Engen R. Estimation of mean glandular dose for breast tomosynthesis: factors for use
602 with the UK, European and IAEA breast dosimetry protocols. *Physics in Medicine & Biology*. 2010;56(2):453.
- 603 26. Zhang G, Cockmartin L, Bosmans H. A four-alternative forced choice (4AFC) software for observer performance
604 evaluation in radiology. *Proc SPIE*. 2016;9787
- 605 27. Obuchowski NA, Gallas BD, Hillis SL. Multi-reader ROC studies with split-plot designs: a comparison of
606 statistical methods. *Academic Radiology*. 2012;19(12):1508-1517.
- 607 28. Gallas BD, Bandos A, Samuelson FW, Wagner RF. A framework for random-effects ROC analysis: biases with
608 the bootstrap and other variance estimators. *Communications in Statistics—Theory and Methods*. 2009;38(15):2586-2603.
- 609 29. Vedantham S, Karellas A, Vijayaraghavan GR, Kopans DB. Digital breast tomosynthesis: state of the art.
610 *Radiology*. 2015;277(3):663-684.
- 611 30. Goodsitt MM, Chan H-P, Schmitz A, et al. Digital breast tomosynthesis: studies of the effects of acquisition
612 geometry on contrast-to-noise ratio and observer preference of low-contrast objects in breast phantom images. *Physics in*
613 *Medicine & Biology*. 2014;59(19):5883.
- 614 31. Gang GJ, Lee J, Stayman JW, et al. Analysis of Fourier-domain task-based detectability index in tomosynthesis
615 and cone-beam CT in relation to human observer performance. *Medical physics*. 2011;38(4):1754-1768.
- 616 32. Sechopoulos I, Ghetti C. Optimization of the acquisition geometry in digital tomosynthesis of the breast. *Medical*
617 *physics*. 2009;36(4):1199-1207.
- 618 33. Reiser I, Nishikawa R. Task-based assessment of breast tomosynthesis: Effect of acquisition parameters and
619 quantum noise a. *Medical physics*. 2010;37(4):1591-1600.
- 620 34. Samei E, Thompson J, Richard S, Bowsher J. A case for wide-angle breast tomosynthesis. *Academic radiology*.
621 2015;22(7):860-869.
- 622 35. Scaduto DA, Huang H, Liu C, et al. Impact of angular range of digital breast tomosynthesis on mass detection in
623 dense breasts. In: *Proc. Vol. 10718, 14th Int Workshop on Breast Imaging (IWBI 2018)*. SPIE; 2018:107181V.

- 624 36. Mackenzie A, Kaur S, Elangovan P, Dance D, Young K. Comparison of synthetic 2D images with planar and
625 tomosynthesis imaging of the breast using a virtual clinical trial. In: *Proc. SPIE 10577, Med. Imaging. Img Perception,*
626 *Obsv Perf, and Tech Assessment*; 2018:105770H.
- 627 37. Choi JS, Han B-K, Ko EY, et al. Comparison between two-dimensional synthetic mammography reconstructed
628 from digital breast tomosynthesis and full-field digital mammography for the detection of T1 breast cancer. *European*
629 *radiology*. 2016;26(8):2538-2546.
- 630 38. Zuley ML, Guo B, Catullo VJ, et al. Comparison of two-dimensional synthesized mammograms versus original
631 digital mammograms alone and in combination with tomosynthesis images. *Radiology*. 2014;271(3):664-671.
- 632 39. Svahn T, Hemdal B, Ruschin M, et al. Dose reduction and its influence on diagnostic accuracy and radiation risk
633 in digital mammography: an observer performance study using an anthropomorphic breast phantom. *The British journal*
634 *of radiology*. 2007;80(955):557-562.
- 635 40. Timberg P, Ruschin M, Båth M, et al. Potential for lower absorbed dose in digital mammography: A JAFROC
636 experiment using clinical hybrid images with simulated dose reduction. *Proc SPIE*. 2006;6146:614614.
- 637 41. Chan H-P, Goodsitt MM, Helvie MA, et al. Digital breast tomosynthesis: observer performance of clustered
638 microcalcification detection on breast phantom images acquired with an experimental system using variable scan angles,
639 angular increments, and number of projection views. *Radiology*. 2014;273(3):675-685.
- 640 42. Burgess AE, Jacobson FL, Judy PF. Human observer detection experiments with mammograms and power-law
641 noise. *Medical physics*. 2001;28(4):419-437.
- 642 43. Mackenzie A, Thomson E, Elangovan P, et al. An observer study to assess the detection of calcification clusters
643 using 2D mammography, digital breast tomosynthesis, and synthetic 2D imaging. In: *Proc. SPIE 01952, Med. Imaging.*
644 *Img Perception, Obsv Perf, and Tech Assessment*; 2019:109520U.
- 645 44. Choi JS, Han B-K, Ko EY, Kim GR, Ko ES, Park KW. Comparison of synthetic and digital mammography with
646 digital breast tomosynthesis or alone for the detection and classification of microcalcifications. *European radiology*.
647 2019;29(1):319-329.
- 648 45. Lai Y-C, Ray KM, Lee AY, et al. Microcalcifications detected at screening mammography: synthetic
649 mammography and digital breast tomosynthesis versus digital mammography. *Radiology*. 2018;289(3):630-638.
- 650 46. Makeev A, Ikejimba LC, Salad J, Glick SJ. Objective assessment of task performance: a comparison of two
651 FFDM detectors using an anthropomorphic breast phantom. *Journal of Medical Imaging*. 2019;6(4):043503.
- 652 47. VOXMAM phantom, Leeds Test Objects. <https://www.leedstestobjects.com/index.php/phantom/voxmam-phantom/>
- 653 48. TOR MAM, Leeds Test Objects. <https://www.leedstestobjects.com/index.php/phantom/tor-mam/>
- 654 49. Chakraborty DP, Eckert MP. Quantitative versus subjective evaluation of mammography accreditation phantom
655 images. *Medical Physics*. 1995;22(2):133-143. doi:10.1118/1.597463
- 656 50. Gagne RM, Gallas BD, Myers KJ. Toward objective and quantitative evaluation of imaging systems using images
657 of phantoms. *Medical Physics*. 2006;33(1):83-95. doi:10.1118/1.2140117
- 658

- 659 51. Gennaro G, Ferro F, Contento G, Fornasin F, Di Maggio C. Automated analysis of phantom images for the
660 evaluation of long-term reproducibility in digital mammography. *Physics in Medicine & Biology*. 2007;52(5):1387.
- 661 52. Warren L, Mackenzie A, Dance D, Young K. Comparison of the x-ray attenuation properties of breast
662 calcifications, aluminium, hydroxyapatite and calcium oxalate. *Physics in Medicine & Biology*. 2013;58(7):N103.
- 663 53. Makeev A, Ghamraoui B, Badal A, Graff CG, Glick SJ. Classification of breast calcifications in dual-energy
664 FFDM using a convolutional neural network: simulation study. In: *Proc. SPIE 11312, Med. Imaging*. Phys. Med. Imaging;
665 2020:113120M.
- 666 54. Elangovan P, Mackenzie A, Dance DR, Young KC, Wells K. Using non-specialist observers in 4AFC human
667 observer studies. In: *Proc. SPIE 10132, Med. Imaging*. Phys. Med. Imaging; 2017:1013256.
- 668 55. Petrov D, Marshall N, Young K, Bosmans H. Model and human observer reproducibility for detecting
669 microcalcifications in digital breast tomosynthesis images. In: *Proc. SPIE 10577, Med. Imaging*. Img Perception, Obsv
670 Perf, and Tech Assessment; 2018:105770B.

671
672 Figure Captions:

673 Figure 1. Insertion and printing of masses. (a) The 3D lesion model is shown with spiculations. (b) The mass was duplicated
674 and arranged into a grid with BB markers, with a 2D central slice shown. (c) The grid was inserted into the virtual breast
675 with ring fiducial markers and (d) printed with the corresponding type of ink. Portions of figure reprinted with permission
676 from Ikejimba et al. "Assessment of task-based performance from five clinical DBT systems using an anthropomorphic
677 breast phantom," *15th International Workshop on Breast Imaging (IWBI2020)*. Vol. 11513 (2020)

678
679 Figure 2. Fabrication and insertion of MC template. Clusters were made by manually placing MC specks within a 5 mm-
680 diameter circle, shown (a) from above and (b) as a close-up with visible specks. The completed template with BBs was
681 placed between the central sheets of the printed phantom, shown (c) from above and (d) from the side.

682
683 Figure 3. Phantom positioning during AEC imaging. The phantom is positioned with the posts outside the field of view to
684 allow for AEC estimation with accurate phantom height. The placement can be seen from the (a) top and (b) side views.

685
686 Figure 4. Example of masses within phantom. Regions containing masses are presented for each vendor. The same location
687 was selected in each image with arrows indicating the locations of signals.

688
689 Figure 5. Reader averaged PC for masses. PC was highest with DBT across all systems, while FFDM and SM had similar,
690 lowers PC scores. Asterisks indicate a statistically significant difference. Portions of figure reprinted with permission from
691 Ikejimba et al. "Assessment of task-based performance from five clinical DBT systems using an anthropomorphic breast
692 phantom," *15th International Workshop on Breast Imaging (IWBI2020)*. Vol. 11513 (2020)

693

694 Figure 6. Reader averaged PC for MCs. PC was highest with DBT and FFDM across all systems. Asterisks indicate a
695 statistically significant difference.

696

697 Figure 7. Δ PC for all systems. The Δ PC was computed relative to FFDM for DBT and SM, with results given for masses
698 and MCs.

699

700 Figure Legends:

701 Figure 5. Red – “DBT”. Green – “FFDM”. Blue – “SM”

702 Figure 6. Red – “DBT”. Green – “FFDM”. Blue – “SM”

Author Manuscript

Table 1. Summary of technical specifications.

	Hologic Selenia Dimensions	GE Senographe Essential (Sen oClaire)	GE Senographe Pristina	Siemens MAMMOMAT Inspiration	Fuji Aspire Cristalle
Version	AWS: 1.8.3.63, Cview: 2.0.1.1	Application ADS_56.21.3, RECON_01.10.4	Recon 02.8.7, SM: 2.3.0	VB60E\ VX20A SL21P21 syngo VH22B SL19P26	FDR- 3000AWS Mainsoft V9.0
Detector conversion	Direct	Indirect	Indirect	Direct	Direct
Anti-scatter grid in DBT	No	Yes	Yes	No	No
Detector version	CM862326	PLC0096_05	PXA0045_02	L03-00010	
Field of view (mm)*	217 x 267	239 x 306	239 x 285	238 x 299	236 x 296
Detector element size (μm)	70 ⁱ	100	100	85	50 ⁱⁱ
In-plane pixel size (μm)	FFDM: 65 DBT, SM: 105 ⁱⁱⁱ	FFDM, DBT: 100	FFDM, DBT:100, SM: 100	FFDM, DBT: 85, SM: 89	FFDM: 50, DBT, SM: 100
X-ray tube target	W	Mo or Rh	Mo or Rh	W	W
X-ray tube filtration	Al or Rh	Mo or Rh	Mo or Ag	Rh	Al or Rh
X-ray tube motion	Continuous	Step and shoot	Step and shoot	Continuous	Continuous
Angular range (deg)	15	25	25	50 ^{iv}	15
Number of projections	15	9	9	25	15
Source-to-imager distance (mm)	700	660	660	650	650
Reconstruction method	FBP	Iterative	Iterative	FBP	FBP

*Field of view in the FFDM acquisition.

ⁱDBT uses 2x2 pixel binning.

ⁱⁱPixels are hexagonal.

ⁱⁱⁱIn-plane resolution changes with slice number. This is the pixel size in the plane of focus.

^{iv}While gantry span is 50 degrees, acquisitions take place over 46 degrees.

Table 2. Summary of acquisition parameters.

Vendor	Modality	Gantry Span	Target/ Filter	Tube Voltage (kVp)	Current-time (mAs)	x-y Voxel Size (μm)	AGD (mGy)	Ref. AGD (mGy)
Hologic Selenia Dimensions	DBT*/SM	15°	W/Al	32	65	105	2.1	1.66
	FFDM		WRh	30	180	65	2.0	1.63
GE Senographe Essential (SenoClaire)	DBT	25°	Rh/Rh	29	71	100	1.4	0.89
	FFDM		Rh/Rh	29	71	100	1.4	0.90
GE Senographe Pristina	DBT/SM	25°	Rh/Ag	34	40	100	1.5	1.08
	FFDM		Rh/Ag	34	40	100	1.5	1.09
Siemens MAMMOMAT Inspiration	DBT/SM	50°	W/Rh	30	200	85	2.3	1.86
	FFDM		W/Rh	30	100	85	1.2	0.93
Fuji Aspire Cristalle	DBT*/SM	15°	W/Al	33	52	100	1.9	1.47
	FFDM		W/Rh	30	89	50	1.2	0.84

*Detector uses 2x2 binning in DBT mode

Table 3. Summary of ROIs used in 4AFC study.

Modality	Pixel Size (μm)	Signal	Number of signal present ROIs	Number of signal absent ROIs	ROI Size (pixels)
DBT/ SM	85 - 105	MC	92 - 106	276 - 318	143 x 143 177 x 177
		Mass	44 - 62	132 - 186	190 x 190 235 x 235
FFDM	50 - 100	MC	99 - 102	297 - 306	150 x 150 231 x 231
		Mass	44 - 62	132 - 186	200 x 200 308 x 308

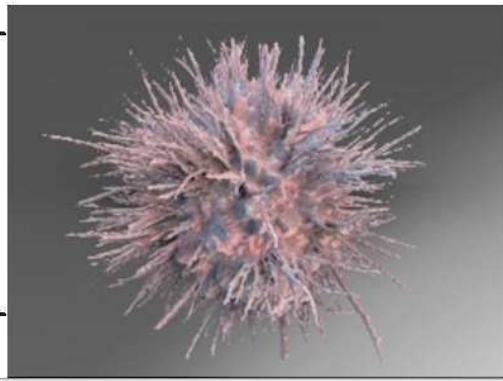
Author Manuscript

Table 4. PC scores for all systems with 95% confidence interval (CI₉₅) in brackets.

		Masses		Microcalcifications	
		PC	CI ₉₅	PC	CI ₉₅
Hologic	DBT	0.72	[0.62,0.81]	0.93	[0.89,0.97]
	FFDM	0.61	[0.51,0.70]	0.94	[0.91,0.97]
	SM	0.52	[0.42,0.61]	0.61	[0.54,0.68]
Fuji	DBT	0.73	[0.64,0.82]	0.87	[0.80,0.93]
	FFDM	0.61	[0.50,0.72]	0.84	[0.77,0.91]
	SM	0.65	[0.55,0.75]	0.63	[0.54,0.72]
GE Essential	DBT	0.80	[0.71,0.89]	0.84	[0.79,0.90]
	FFDM	0.64	[0.54,0.75]	0.79	[0.73,0.85]
GE Pristina	DBT	0.84	[0.76,0.91]	0.95	[0.92,0.98]
	FFDM	0.62	[0.54,0.70]	0.92	[0.88,0.96]
	SM	0.60	[0.50,0.70]	0.53	[0.44,0.62]
Siemens	DBT	0.91	[0.84,0.97]	0.84	[0.79,0.88]
	FFDM	0.64	[0.56,0.71]	0.78	[0.72,0.83]

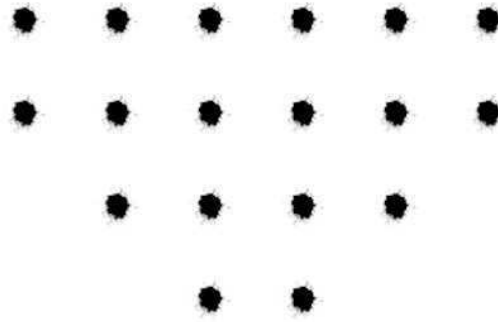
(a)

5 mm



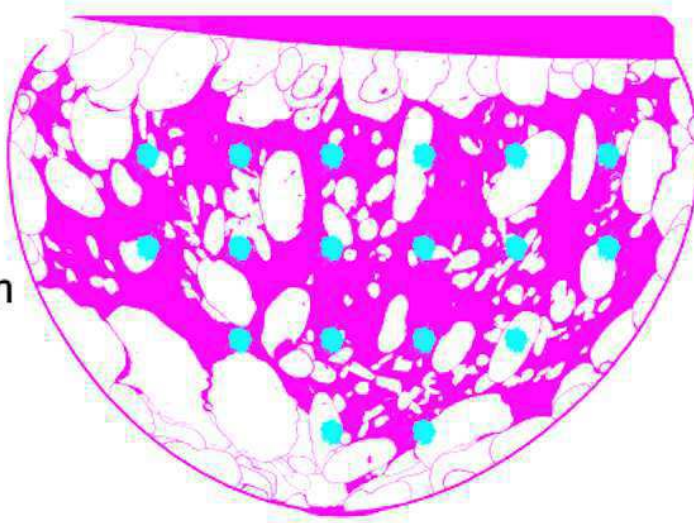
(b)

5 mm



(c)

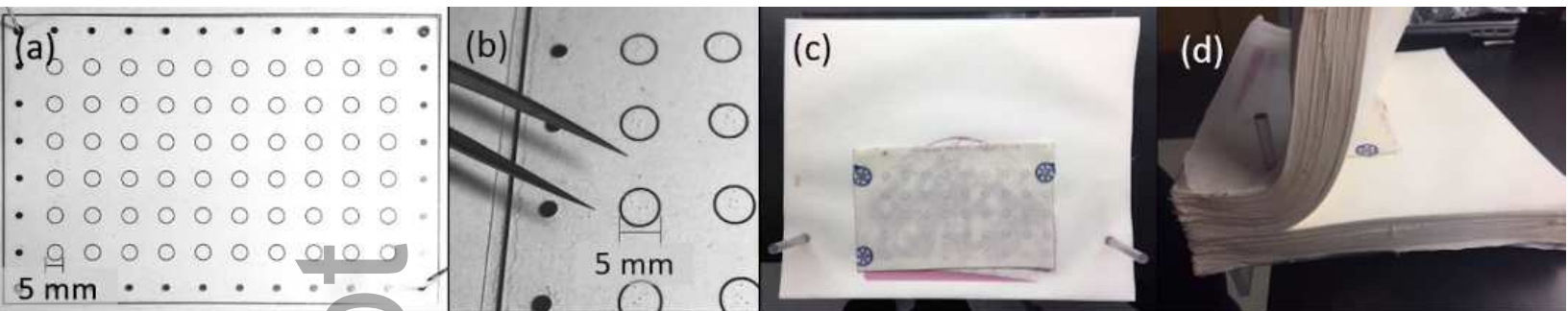
5 mm



(d)

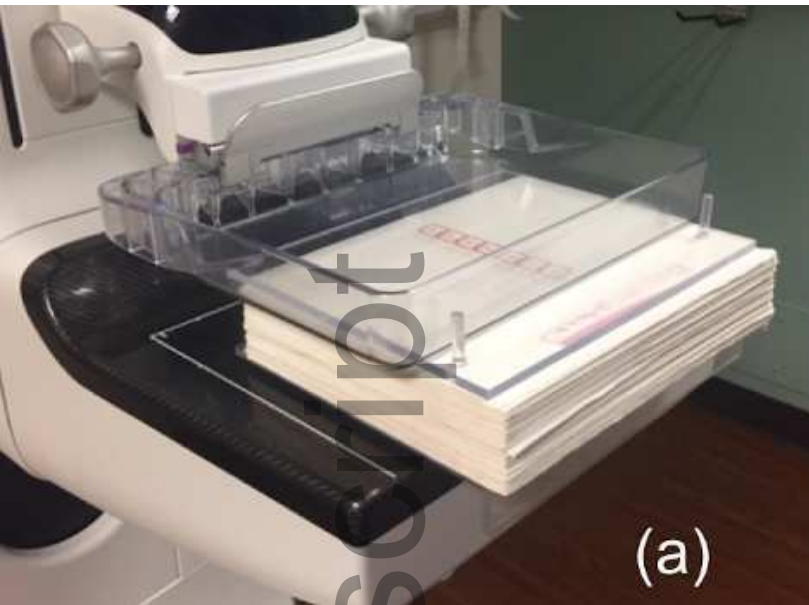
5 mm





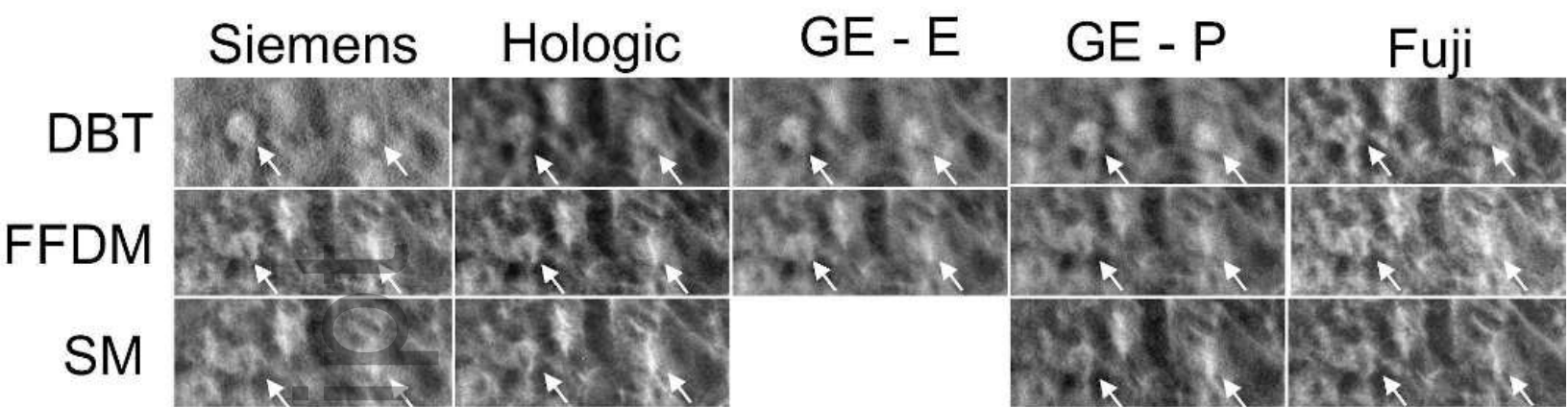
mp_14568_f2.eps

Author Manuscript



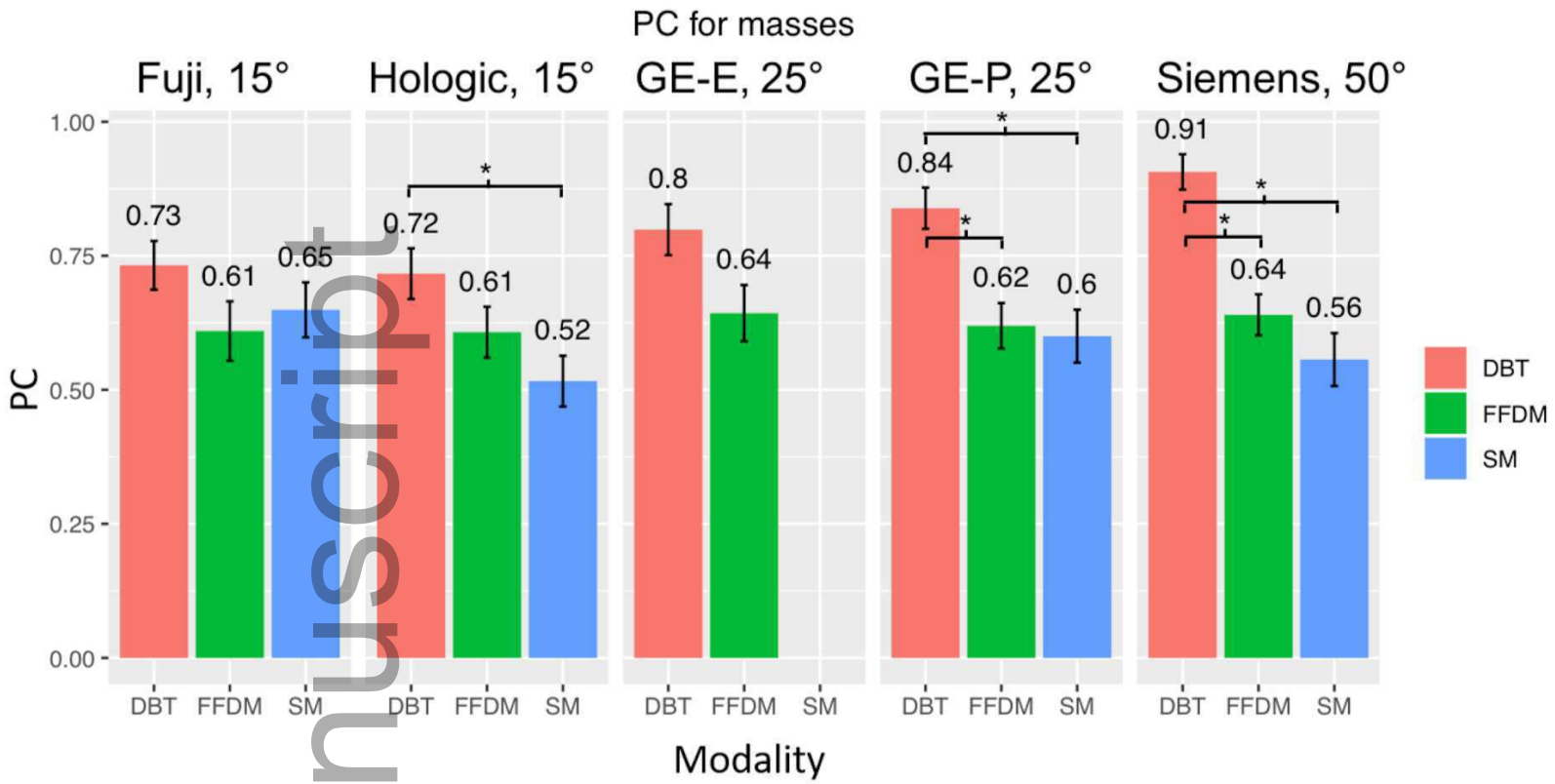
mp_14568_f3.eps

Author Manuscript



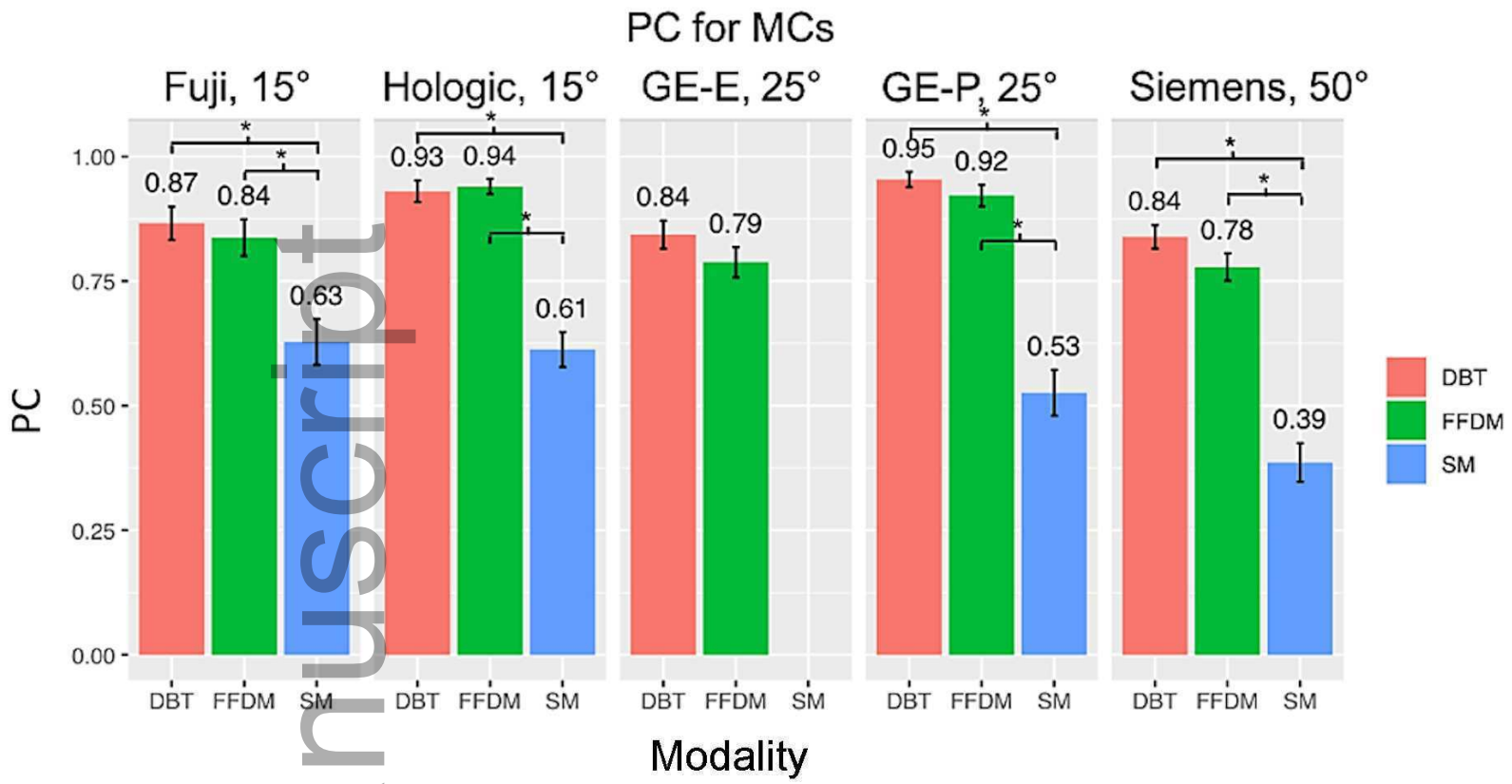
mp_14568_f4.eps

Author Manuscript

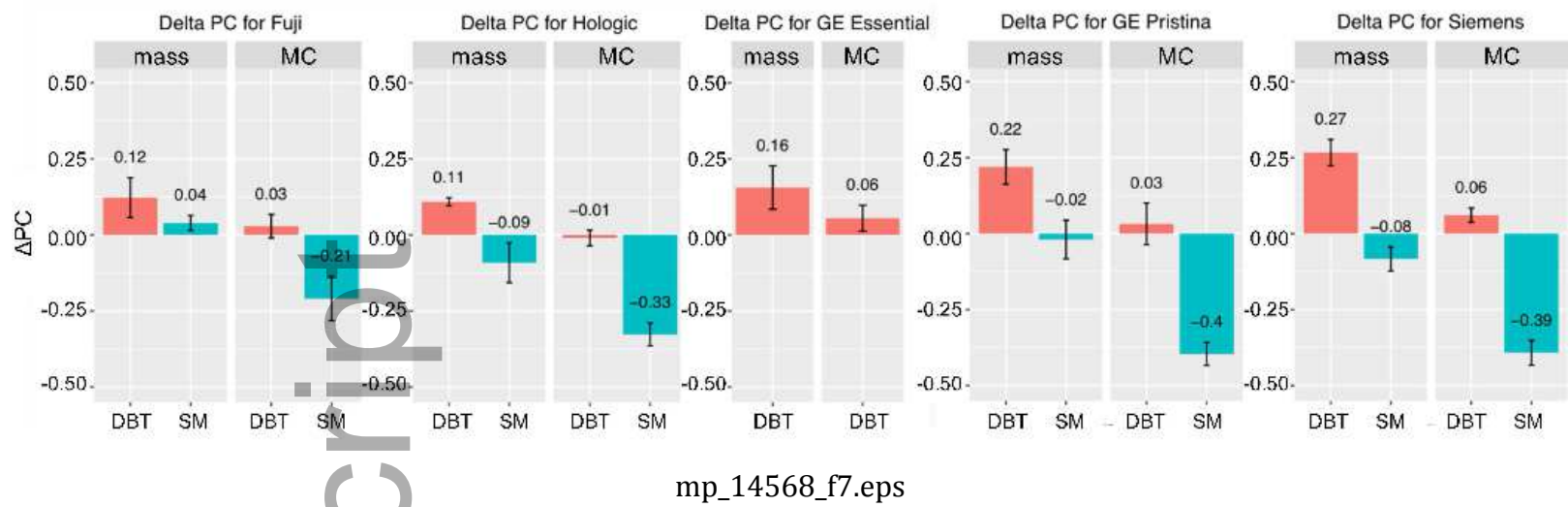


mp_14568_f5.eps

Author Manuscript



mp_14568_f6.eps



Author Manuscript