

Journal Pre-proof

Diagnostic Performance of Artificial Intelligence for Detection of Anterior Cruciate Ligament and Meniscus Tears: A Systematic Review

Kyle N. Kunze, M.D., David M. Rossi, B.S., Gregory M. White, M.D., Aditya V. Karhade, M.D., M.B.A., Jie Deng, Ph.D., Brady T. Williams, M.D., Jorge Chahla, M.D., Ph.D.

PII: S0749-8063(20)30744-1

DOI: <https://doi.org/10.1016/j.arthro.2020.09.012>

Reference: YJARS 57125

To appear in: *Arthroscopy: The Journal of Arthroscopic and Related Surgery*

Received Date: 5 April 2020

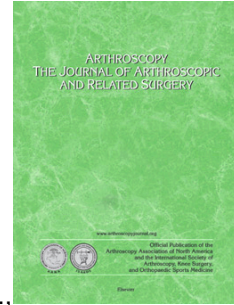
Revised Date: 2 September 2020

Accepted Date: 9 September 2020

Please cite this article as: Kunze KN, Rossi DM, White GM, Karhade AV, Deng J, Williams BT, Chahla J, Diagnostic Performance of Artificial Intelligence for Detection of Anterior Cruciate Ligament and Meniscus Tears: A Systematic Review, *Arthroscopy: The Journal of Arthroscopic and Related Surgery* (2020), doi: <https://doi.org/10.1016/j.arthro.2020.09.012>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2020 Published by Elsevier on behalf of the Arthroscopy Association of North America



Diagnostic Performance of Artificial Intelligence for Detection of Anterior Cruciate Ligament and Meniscus Tears: A Systematic Review

Kyle N. Kunze, M.D.¹

David M. Rossi, B.S.²

Gregory M. White, M.D.³

Aditya V. Karhade, M.D., M.B.A.⁴

Jie Deng, Ph.D.³

Brady T. Williams, M.D.²

Jorge Chahla, M.D., Ph.D.²

1. Department of Orthopaedic Surgery, Hospital for Special Surgery, New York, NY, USA
2. Department of Orthopaedic Surgery, Rush University Medical Center, Chicago, IL, USA
3. Department of Diagnostic Radiology and Nuclear Medicine, Rush University Medical Center, Chicago, IL, USA
4. Department of Orthopedic Surgery, Massachusetts General Hospital, Harvard Medical School, Boston, MA, USA

Correspondence:

Jorge Chahla, M.D., Ph.D.
Rush University Medical Center
1611 W. Harrison St. Suite 300
Chicago, IL, USA 60612
jachahla@msn.com

1 **Diagnostic Performance of Artificial Intelligence for**
2 **Detection of Anterior Cruciate Ligament and Meniscus**
3 **Tears: A Systematic Review**

Journal Pre-proof

4 ABSTRACT

5 **Purpose:** To (1) determine the diagnostic efficacy of artificial intelligence (AI) methods for
6 detecting anterior cruciate ligament (ACL) and meniscus tears and to (2) compare the efficacy
7 to human clinical experts.

8

9 **Methods:** PubMed, OVID/Medline, and Cochrane libraries were queried in November 2019 for
10 research articles pertaining to AI utilization for detection of ACL and meniscus tears.
11 Information regarding AI model, prediction accuracy/area under the curve (AUC), sample sizes
12 of testing/training sets, and imaging modalities were recorded.

13

14 **Results:** A total of 11 AI studies were identified: 5 investigated ACL tears, 5 investigated
15 meniscal tears, and 1 investigated both. The AUC of AI models for detecting ACL tears ranged
16 from 0.895-0.980, and the prediction accuracy ranged from 86.7%-100%. Of these studies,
17 three compared AI models to clinical experts. Two found no significant differences in diagnostic
18 capability, while one found that radiologists had a significantly higher sensitivity for detecting
19 ACL tears ($p=0.002$) and statistically similar specificity and accuracy. Of the 5 studies
20 investigating the meniscus, the AUC for AI models ranged from 0.847-0.910 and prediction
21 accuracy ranged from 75.0%-90.0%. Of these studies, 2 compared AI models to clinical experts.
22 One found no significant differences in diagnostic accuracy, while one found that the AI model
23 had a significantly lower specificity ($p=0.003$) and accuracy ($p=0.015$) than radiologists. Two
24 studies reported that the addition of AI models significantly increased the diagnostic
25 performance of clinicians compared to their efforts without these models.

26

27 **Conclusion:** AI prediction capabilities were excellent and may enhance the diagnosis of ACL and
28 meniscal pathology; however, AI did not outperform clinical experts.

29

30 **Clinical relevance:** AI models promise to improve diagnosing certain pathologies as well as or
31 better than human experts, are excellent for detecting ACL and meniscus tears, and may

32 enhance the diagnostic capabilities of human experts; however, when compared to these
33 experts, may not offer any significant advantage.

Journal Pre-proof

34 INTRODUCTION

35 The development and application of deep learning (DL) and machine learning (ML) algorithms
36 to generate prediction models from large datasets is an increasingly utilized statistical tool
37 which relies on pattern recognition and constrained feature selection. By “constraining” feature
38 selection, the algorithms limit the number of variables ultimately chosen and used in
39 subsequent analyses by selecting only those with the greatest predictive value from an initial
40 large pool of potential variables. A clinically relevant application of these artificial intelligence
41 (AI) methods pertains to its ability to diagnose injury and disease on medical imaging studies. AI
42 models learn to recognize disease patterns through repetition and learning, and it is thought
43 that such features may confer the ability to more quickly and accurately identify disease.

44
45 Machine learning describes statistical processes that exhibit the “learning” associated with
46 human intelligence and leverage this experiential learning to improve and refine programmed
47 algorithms to predict and outcome.¹ The algorithms are applied to a dataset of interest, and
48 self-train based on patterns in the dataset. Once trained, the algorithms can make specific
49 decisions when presented with data that it has not seen before. Each machine learning
50 algorithm makes decisions based on different sets of rules that are out-of-scope of the current
51 study, but allow them to come to decisions in different ways. Algorithms can be modified to
52 optimize their prediction capabilities, and ultimately the predictions made by the algorithm are
53 compared against the true outcome present in the data set to determine how accurate
54 predictions are.² This approach has become increasingly popular given the ability of these
55 algorithms to optimize prediction accuracy, whereas traditional statistics may sacrifice accuracy
56 at the cost of favoring interpretability. Machine learning has become of recent interest in
57 orthopaedics given these potential benefits, as evidenced by the recent increase in literature
58 which has applied this methodology.³⁻⁶

59
60 AI technology has been successfully applied in various clinical scenarios. Detection of diabetic
61 retinopathy through analyzing retinal fundus photographs^{7, 8} and skin cancer through
62 constructing deep neural networks based on imaging and disease labels⁹ have efficacy

63 comparable to, or better than, human experts. Within the field of orthopedic spine and
64 oncologic surgery specifically, AI algorithms are gaining popularity by aiding decision-making
65 and can be used in clinical settings.^{5,10,11} However, the performance of these models compared
66 to clinical experts in the field remains poorly understood. A recent systematic review by
67 Langerhuizen et al.¹² found that AI algorithms had excellent performance for fracture detection
68 in the orthopedic trauma literature and outperformed human examiners for detecting and
69 classifying hip and proximal humerus fractures. Although it appears that AI methods may
70 confer diagnostic benefits in other realms of orthopedic surgery, their performance and clinical
71 utility in sports medicine remains poorly defined.

72

73 Imaging-based detection of sports medicine injuries of the knee, specifically the use of
74 magnetic resonance imaging (MRI) for anterior cruciate ligament (ACL) and meniscus tears, is
75 the current gold standard for diagnosis. However, the diagnostic accuracy of MRI may be
76 decreased in several circumstances: (1) observer inexperience and bias, (2) small partial or
77 incomplete tears, (3) imaging artifacts, (4) incomplete MRI study, and (5) presence of
78 concomitant injuries. Application of AI methods may address these shortcomings by facilitating
79 clinical decision-making and improving patient management. As such, the purpose of the
80 current study was to (1) determine the diagnostic efficacy of AI methods for detecting ACL and
81 meniscus tears and to (2) compare the efficacy to human clinical experts. The authors
82 hypothesized that AI method performance would be excellent for detection of ACL and
83 meniscus tears and could outperform human examiners.

84 **METHODS**

85 *Identification and Selection of Articles*

86 A systematic search in accordance with the 2009 Preferred Reporting Items for Systematic
87 Review and Meta-Analysis (PRISMA) statement¹³ was conducted using PubMed, OVID/Medline,
88 and Cochrane libraries. The timeframe for the search was the conception of each online
89 database until November 8, 2019. The following Boolean search syntax was used to conduct the
90 search: (orthopedics OR orthopedic procedures OR ligament OR tear* OR (ligament* AND tear*
91 AND orthop*)) AND (artificial intelligence OR neural network* or deep learning OR machine
92 learning OR machine intelligence) AND (predict* OR predictive value of test OR score OR scores
93 OR scoring system OR scoring systems OR observ* OR observer variation OR detect* or
94 evaluat* OR analy* OR assess* OR measure*). The protocol for the current systematic review
95 was registered on PROSPERO prior to collection and analysis of the data (ID: *blinded for review*).

96
97 Articles populated from the above search met inclusion criteria if (1) the study methods and
98 analyses pertained to development or utilization of artificial intelligence or machine learning for
99 detecting or classifying the presence of an ACL or meniscus tear, and (2) was published in the
100 English language. Studies were excluded if (1) data was only published in the form of an
101 abstract, technique paper, cadaveric or animal study, or letter to the editor; or (2) pertained to
102 robotic-assisted surgery. Two observers (blinded for reviewer) independently screened the
103 abstracts and titles of potential articles. Full-text review was only performed during the study
104 selection process if necessary to determine if the articles satisfied inclusion and exclusion
105 criteria. Additionally, all references from the included studies were reviewed and reconciled to
106 verify that no relevant articles were missing from the systematic review. A total of 1,619
107 records were initially identified, and a total of 11 were ultimately included in the qualitative
108 synthesis (**Figure 1**).

109

110 Data Acquisition

111 All data were recorded into a custom spreadsheet using a modified information extraction
112 table.¹⁴ Categories for data collection for each full article included (1) article information; (2)

113 input features; (3) imaging plane; (4) size of training and testing samples; (5) ground truth label
114 assignments; (6) output classes; (7) AI models used; (8) use of pretrained Convolutional Neural
115 Networks (CNN); and (9) performance.

116

117 *Assessment of Heterogeneity and Methodological Quality*

118 A modified MINORS scoring criteria was used to assess quality as has been previously applied in
119 systematic reviews on AI in orthopedics¹² given that studies concerning AI methods are
120 classified as developmental rather than diagnostic. Quality appraisal focused on identification
121 of (1) a clear study aim, (2) description of inclusion and exclusion criteria for input features (all
122 eligible imaging examples included), (3) determination of ground truth (reference standards for
123 AI), (4) report of distribution of data set (training, validation, and testing phases), (5) described
124 how performance of AI model was assessed (area under the curve [AUC]/prediction), and (6)
125 clearly described AI model used. These criteria were applied to and quantified for each study.
126 For reference, the AUC is the quantitative output of a receiver operator curve (ROC) analysis of
127 discrimination. ROC and discrimination analysis is a common performance analysis in diagnostic
128 studies, which assesses the probability that the machine learning model will assign a greater
129 predicted probability to a randomly selected positive case (true positive case, i.e., a patient who
130 actually had an ACL or meniscus tear) relative to a randomly selected negative case (false
131 positive case, i.e., a patient who did not have an ACL or meniscus tear). Each study could score
132 a total of seven points, with a score of zero indicating poor methodological quality, and a score
133 of seven indicating the highest methodological quality. Two independent observers (blinded for
134 review) assessed all included studies. The inter-observer reliability was excellent at 0.97 (95%
135 Confidence interval, 0.93-0.99). Any discrepancies were resolved by consensus.

136

137 *Statistical Analysis*

138 All data was qualitatively synthesized and reported in both narrative fashion in addition to table
139 format. Extracted data was presented as means and ranges when appropriate with associated
140 p-values given the degree of heterogeneity between studies. All studies considered a p-value

141 <0.05 to indicate statistical significance. All data extraction and analyses were performed in
142 Microsoft Excel (Microsoft Corporation, Washington, USA).

Journal Pre-proof

143 **RESULTS**

144 A total of 11 studies were identified in the search.¹⁵⁻²⁵ All 11 studies were included in the
145 qualitative data analysis. Of these 11 studies, five investigated the use of AI to detect ACL
146 tears,^{15, 17, 20, 21, 24} five investigated the use of AI to detect tears of the meniscus,^{18, 19, 22, 23, 25} and
147 one investigated both.¹⁶

148

149 *AI Model Performance: ACL Tear Detection*

150 All six studies that investigated the performance of AI on ACL tear detection utilized knee MRI
151 acquired in standard imaging planes. Four (66.6%) of the studies analyzed sagittal-plane images
152 only,^{15, 20, 21, 24} one study analyzed coronal images only,¹⁷ and one study analyzed sagittal,
153 coronal, and axial images.¹⁶ A total of four (66.6%) studies reported AUC data for complete ACL
154 tear detection. The AUC for these AI models ranged from 0.895-0.980 (**Table 1**). Additionally,
155 four studies report AI model prediction accuracy for specifically detecting complete ACL tear
156 (range 86.7%-100%).

157

158 Štadjuhar et al.²⁴ utilized two different feature extraction techniques: Histogram of Oriented
159 Gradient (HOG) and Generalized Search Tree (GIST). These feature extraction techniques were
160 subsequently paired with two commonly used machine learning models: Support Vector
161 Machine (SVM) and Random Forest (RF). They found that their best performing machine
162 learning model that combined HOG with linear-kernel SVM (HOG+lin-SVM) performed the best,
163 producing an AUC of 0.894 for differentiating between an injured ACL and healthy ACL and an
164 AUC of 0.943 for detecting completely ruptured ACL cases only.

165

166 Abdullah et al.¹⁵ described a diagnostic system consisting of image pre-processing, feature
167 extraction, and finally classification. The authors utilized k-Nearest Neighbor (k-NN) and Back
168 Propagation Artificial Neural Network (BP-ANN) classifiers to determine the best accuracy for
169 ACL tear classification. They found that BP-ANN produced a higher classification accuracy of
170 94.44%, compared to 87.33% for k-NN.

171 Chang et al.¹⁷ evaluated multiple customized CNN models with variations in the input fields-of-
172 view (i.e. full slice, cropped slice, dynamic patch-based sampling) and dimensionality (single
173 slice, three slices, five slices) for detection of complete ACL tears. They determined that the
174 model created to dynamically sample random cropped patches of images of the ACL
175 performed the best in terms of detecting ACL tears when compared to a similar model that
176 utilized the entire uncropped MRI slices. The model that utilized dynamic sampling had an
177 accuracy of 96.7% and AUC of 0.971.

178

179 Bien et al.¹⁶ used a fully automated deep learning CNN model with logistic regression for
180 predicting the presence or absence of ACL tears on MRI after image pre-processing
181 (intact=normal, mucoid degeneration, ganglion cyst, sprain; tear=low-grade partial tear with
182 <50% fibers torn, high-grade partial tear with >50% of fibers torn, or complete tear). They
183 reported that their best performing model produced an AUC of 0.965 (95% CI 0.938-0.993) for
184 ACL tear detection. The specificity, sensitivity, and accuracy of the model were also reported as
185 0.968 (95% CI 0.890-0.991), 0.759 (95% CI 0.635-0.850), and 0.867 (95% CI 0.794-0.916),
186 respectively.

187

188 *AI Model Performance Compared with Human Observers: : ACL Tear Detection*

189 Three (50.0%) studies compared the performance of an AI model for ACL tear detection with
190 human medical experts.^{16, 20, 21}

191

192 Liu et al.²⁰ trained multiple CNNs and applied them to a test set of 50 MRI images of full
193 thickness ACL tears and 50 MR images with intact ACLs. They found that their model with the
194 best overall diagnostic performance for detecting the presence or absence of a full thickness
195 ACL tear produced an AUC of 0.98 (95% CI: 0.93-1.00, p <0.001) However, there was no
196 statistically significant difference in diagnostic performance found between the AI model and
197 clinical radiologist performance (Radiologist 0.90 (95% CI: 0.95-1.00); Fellow 0.90 (95% CI: 0.95-
198 1.00); Resident1 0.93 (95%CI 0.88-0.98); Resident2 0.97 (95%CI 0.94-1.00); Resident3 0.98 (95%
199 CI 0.95-1.00).

200

201 Mazlan et al.²¹ tested an SVM algorithm on 60 samples from MR images of 100 non-injured
202 ACLs, 100 partially-torn ACLs, and 100 completely-torn ACLs that underwent pre-processing.
203 They reported that the SVM model had an accuracy of 100% for classifying ACL MRI samples as
204 normal, partial-tear, or complete-tear. The authors also sought to compare the diagnostic
205 capability of their AI model to that of two medical experts. No statistically significant
206 differences between the AI model and radiologists were found in terms of diagnostic
207 capabilities, as the SVM and both medical experts correctly identified all 10 samples with 100%
208 accuracy.

209

210 Bien et al.¹⁶ compared their MRNet model's performance for detecting ACL tears to three
211 musculoskeletal (MSK) radiologists on a testing set of 120 knee MR images, with the majority
212 vote of 3 musculoskeletal (MSK) radiologists serving as the reference standard. They also
213 evaluated changes in the diagnostic performance of clinical experts when the AI model
214 predictions were provided to the radiologists during interpretation. Their model for detecting
215 complete ACL tear produced an AUC of 0.968 (95% CI 0.890-0.991) compared to radiologist
216 specificity of 0.933 (95% CI 0.906-0.953). However, results were not statistically significant (p-
217 value=0.441). Radiologists achieved significantly higher sensitivities for tear diagnosis than the
218 AI model (AUC 0.906 vs. 0.759, p-value=0.002). The AI model accuracy for ACL tear detection
219 was 0.867 (95% CI 0.794-0.916), which was lower than the MSK radiologist accuracy of 0.920
220 (95% CI 0.900-0.937), which was not statistically significant (p-value=0.075). When provided
221 with model assistance, there was a statistically significant increase (4.8%, p<0.001) in the
222 clinical experts' specificity in identifying ACL tears. They reported that because the testing set
223 consisted of 62 exams that were negative for ACL tear, the represented increase in specificity in
224 the optimal clinical setting would potentially translate to the avoidance of three unnecessary
225 surgeries for suspected ACL tears.

226

227 *AI Model Performance: Meniscus Tear Detection*

228 All six studies that investigated the performance of AI models on meniscus tear detection
229 utilized MRI. Five (83.3%) of the studies analyzed sagittal-plane images only,^{18, 19, 22, 23, 25} and
230 one study analyzed sagittal, coronal, and axial images.¹⁶ Four (66.67%) studies reported AUC
231 data for meniscus tear detection, ranging from 0.847-0.910 (**Table 2**).

232

233 Fu et al.¹⁹ compared the performance of two SVM models to detect meniscus tears. One model
234 was created to select relevant meniscus MR features, while the other model implemented the
235 SVM model without feature selection. The SVM model without feature selection produced an
236 AUC of 0.727 for meniscus tear detection, while their model with feature selection yielded an
237 AUC value of 0.912 for meniscus tear detection.

238

239 Couteaux et al.¹⁸ used an R-CNN model for tear detection (tear in any meniscus) and localization
240 (anterior or posterior). The anterior meniscus was classified as torn when at least one network
241 had detected a torn anterior meniscus and the posterior meniscus was classified as torn when
242 the strict majority of the networks had detected a torn posterior meniscus. When they applied
243 their model to a test set of 700 MRIs, the authors found that the model produced a weighted
244 AUC score of 0.906.

245

246 Roblot et al.²³ used a three-step AI model where an image was transferred into a R-CNN trained
247 to detect menisci as torn or normal, meniscus tear location, and whether the tear was
248 horizontal or vertical. The model was tested on a dataset of 700 MRI images to perform
249 detection of meniscus tear presence, position, and orientation. The AI model produced an AUC
250 of 0.94 for presence of a meniscal tear, 0.92 for detection of the position of the two meniscal
251 horns, and 0.83 for orientation of the tear. The overall combined AUC was 0.90.

252

253 Pedoia et al.²² created a deep-learning model that combined meniscus segmentation and a 3D
254 CNN for the detection and severity staging of meniscus lesions. The model was first built to
255 discriminate between the presence of a lesion versus no lesion (including no lesion and
256 intrasubstance abnormalities), and subsequently lesion severity (severe lesion=maceration of

257 the meniscus; mild-moderate lesion=non-displaced tears and displaced and complex tears
258 without deformity; no lesion= lesion absence and intrasubstance abnormalities). This model
259 produced a lesion versus no lesion AUC of 0.89 on the test dataset and accuracies of 80.74%,
260 78.02%, and 75.00% for determining severe lesion versus mild-moderate lesion versus no
261 lesion, respectively.

262

263 Bien et al.¹⁶ also investigated the ability of their CNN models to detect meniscus tears following
264 their investigations of diagnostic capabilities for ACL tears. For meniscus tear diagnosis, this
265 group reported an accuracy of 0.725 (95% CI 0.639-0.797) and an AUC of 0.847 (95% CI 0.780-
266 0.914).

267

268 Fazel-Zarandi et al.²⁵ used AI for MR image segmentation followed by the application of a
269 Perceptron Neural Network (PNN) for classification of meniscal tears. A testing dataset of 50
270 MRIs were fed into the PNN and resulted in meniscus tear versus no meniscus tear accuracy of
271 90%. Classification rate (precision %) was also reported for five different settings of meniscus
272 tear including: (1) medial anterior horn and posterior horn normal (88.82%), (2) lateral anterior
273 horn and posterior horn normal (92.13%), (3) medial anterior horn normal and posterior horn
274 torn (84.24%), (4) lateral anterior horn normal and posterior horn torn (91.96%) and (5) lateral
275 anterior horn torn and posterior horn normal (87.64%).

276

277 *AI Model Performance Compared with Humans: Meniscus*

278 Two (33.3%) studies compared the performance of using an AI model for meniscus tear
279 detection with human medical experts.^{16, 22}

280

281 Bien et al.¹⁶ compared the performance of their AI model with unassisted MSK radiologists for
282 detecting meniscus tear (intact=normal, degenerative changes without tear, postsurgical
283 changes without tear; tear=increased signal reaching the articular surface on at least two slices
284 or morphologic deformity). When compared to the MSK radiologists in the study, the AI model
285 had a statistically significant lower specificity (AUC 0.882, 95% CI 0.847-0.910 versus AUC 0.741,

286 95% CI 0.616-0.837; p-value=0.003) and accuracy (accuracy (0.849, 95% CI 0.823-0.871 versus
287 0.725, 95% CI 0.639-0.797, p=0.015). The sensitivity (0.820, 95% CI 0.781-0.853 versus 0.710,
288 95% CI 0.587-0.808; p=0.504) was also shown to be lower for the AI model compared to MSK
289 radiologists, although this was not statistically significant.

290

291 Pedoia et al.²² analyzed 1,478 MRI studies and utilized automatic segmentation of cartilage and
292 meniscus using 2D U-Net. Detection and severity staging of meniscus and cartilage lesion was
293 performed with a 3D CNN. Comparisons were made between their model and experts, where
294 they sought to determine the inter-rater variability between three MSK radiologists (expert 1:
295 >20 years of experience, expert 2: 10 years of experience, <1 year of experience) for
296 determining meniscus lesion severity on selected cases. They found an average agreement
297 between the three experts of 86.27% for no meniscus lesion, 66.48% for mild-moderate lesion,
298 and 74.66% for severe lesion, while the best AI model obtained accuracies of 80.74% for no
299 meniscus lesion, 78.02% for mild-moderate lesion, and 75.00% for severe lesion.

300

301 *Quality Assessment*

302 The average modified MINORS score among all studies was 4.9 ± 1.0 (**Table 3**), indicating
303 moderate to high methodological quality of the included studies on average. The most common
304 reasons for loss of quality points was failure to describe both the inclusion and exclusion criteria
305 of input features including patient and imaging selection (n=4, 36.4%) and failure to describe
306 ground truth assignment (n=4, 36.4%). In the absence of clearly defined inclusion/exclusion
307 criteria, selection bias cannot be excluded for the four studies. Failure to clearly describe the
308 ground truth assignment risks publishing data from poorly trained AI models that may be
309 inaccurate. The remaining limitations to quality was failure to describe the distribution of data,
310 which also potentiates selection bias and misinterpretation of conclusions.

311 DISCUSSION

312 The main finding of the current study was that the AUC and prediction accuracy of AI models
313 for detecting ACL tears ranged from 0.895-0.980 and 86.7%-100%, while the AUC and
314 prediction accuracy for detecting meniscus tears ranged from 0.847-0.910 and 75.0%-90.0%.
315 Additionally, in two studies that compared AI models to clinical experts, one found no
316 significant differences in diagnostic accuracy, while one found that the AI model had a
317 significantly lower specificity and accuracy than radiologists. Two studies reported that the
318 addition of AI models significantly increased the diagnostic performance of clinicians compared
319 to their efforts without these models. However, the heterogeneity of the studies and
320 methodology identified in this systematic review suggests several areas for improvement and
321 makes interpretation across studies challenging.

322

323 AI models are mathematical computing algorithms trained to integrate big data and
324 autonomously assign labels to unseen data. Through multiple statistical iterations and pattern
325 recognition, these models can apply learned features from training data sets and apply them to
326 test sets to detect or classify lesions on many imaging modalities. Discrimination is also
327 employed in conjunction with these AI models through generating a receiver operator curve
328 (ROC) and generating a c-statistic (area under the curve, AUC). An AUC of 1.0 indicates perfect
329 discrimination, while an AUC of 0.5 indicates discrimination similar to chance.²⁶ The current
330 study found that the AUC for detecting ACL tears was near perfect ranging from 0.90-0.98 and
331 that for detecting meniscus tears was excellent at 0.85-0.91.

332

333 Perhaps most importantly, the current study found that a combination of AI and human experts
334 outperformed human experts or AI alone for diagnosis of ACL and meniscal tears, similar to the
335 results of a prior systematic review of natural and artificial intelligence in neurosurgery.²⁷ AI
336 methods have been previously applied to achieve or exceed human-level performance for tasks
337 ranging from detection of distal radius fractures and hip fractures, to malignant pulmonary
338 nodules.²⁸⁻³² The clinical relevance of AI applications for ACL and meniscal tears may be divided
339 into the following categories: (1) human-level performance or better on routine tasks and (2)

340 human-level performance or better on difficult tasks. During algorithm development, cohort
341 curation, and study design, the studies included in this systematic review did not distinguish
342 between these ultimate goals. If the intended purpose of AI algorithms is to diagnose lesions
343 that are difficult for humans (partial tears, poor imaging), the AI algorithms should be
344 specifically trained for this purpose. On the other hand, if the purpose of AI algorithms is to
345 diagnose simple lesions, the output of these algorithms should be designed into clinically
346 relevant categories (“definitely normal”, “definitely abnormal”, “not sure”) to improve the
347 efficiency and workflow of diagnosticians.

348

349 Interestingly, none of the included studies compared the use of AI to the gold standard of
350 confirming ACL and meniscus lesions, which is arthroscopy. This is likely a function of the
351 designs of the included studies, in which those that made a comparison to human experts were
352 intended to do so. In the two studies that compared AI models to clinical experts, one found no
353 significant differences in diagnostic accuracy and one found that AI was inferior in diagnostic
354 accuracy. If it is assumed that clinical radiologist experts with only images at their disposal are
355 less accurate at diagnosing these lesions than the gold standard of arthroscopy, then by
356 transitive property, AI may be less accurate than arthroscopy as well. Future studies are
357 warranted to determine the accuracy of AI for diagnosis of these lesions in comparison to
358 arthroscopy as a ground truth label.

359

360 Another area for improvement is the requirement for adherence to peer-reviewed AI-specific
361 guidelines. Efforts are underway to update the TRIPOD guidelines and develop a standardized
362 system for AI applications in healthcare.³³ Several of the studies included in this analysis did not
363 report measures of model performance such as precision-recall curves and Brier score that are
364 key to interpreting diagnostic studies, particularly when outcomes are not balanced.^{34, 35} AI is
365 often criticized for the “black-box” nature of transformations required to take input data and
366 produce meaningful outcomes. This block-box limits our ability to understand the specific
367 imaging features an AI method utilized to produce its probability outputs. However, prior
368 studies have provided explainable algorithms where output not only predicted probabilities,

369 but also explanations in the form of augmented input images with heat maps highlighting the
370 regions of interest for the specific diagnosis task.²⁹ Success by Lindsey et al. in applying these
371 explainable AI techniques for distal radius fracture detection should also become the norm for
372 sports medicine.²⁹

373

374 AI has significant implications for the future of diagnosis in orthopaedic sports medicine, but
375 clinicians must be informed and critical consumers of this rapidly evolving field.³⁶ The focus of
376 this review was on diagnostic applications of AI in orthopaedic sports medicine and a similar
377 analysis for prognostic applications of AI remains to be undertaken. Interestingly, AI did not
378 outperform human experts, which is potentially a result of the early applications of AI in this
379 field and the need to further refine and modify algorithms. It is possible that as these
380 algorithms continue to be trained with more data, that prediction accuracies and image
381 recognition improves and may eventually outperform these experts. As described above, the
382 drawbacks of current AI applications for diagnosis of ACL and meniscal tears should inspire
383 future studies to follow standardized guidelines to allow for reliable and reproducible research.
384 However, advantages of AI include the potential for the rapid and accurate identification and
385 diagnosis of pathology, such as ligamentous and meniscal tears, which may initially be missed
386 by the human eye. Eventually, AI and related technology may progress to the point where
387 fewer working personnel are required to perform these tasks (i.e., the development of
388 pathology-specific AI algorithms, where only one attending radiologist is required to double-
389 check the finding made by the algorithm, as opposed to the current use of teams of multiple
390 radiologists who are burdened with large numbers of images with multiple views to read). This,
391 in turn, may increase timeliness of reads and decrease healthcare costs. Ultimately, ensuring
392 clinical relevance at every step in algorithm conception, design, and development will lead to
393 true progress for AI in orthopaedic sports medicine. Applications of AI in orthopaedic surgery
394 are rapidly growing and periodic updates will be required to appropriately represent the state
395 of the literature in the years to come. At present, inadequate reference standards to train and
396 test AI is the biggest hurdle to overcome prior to integration into clinical workflows.

397

398 *Limitations*

399 There are a number of limitations that must be acknowledged to appropriately interpret the
400 results of this study. This was a systematic review that followed the PRISMA guidelines, but did
401 not include a more formal quantitative meta-analysis due to study heterogeneity. Despite the
402 comprehensive search, the total number of studies included in this analysis was relatively small.
403 Another limitation is that no studies included diagnostic arthroscopy as the gold standard
404 reference to diagnose ACL or meniscus lesions, which may limit the applicability of the findings
405 to clinical practice. Finally, studies were inherently heterogeneous given the AI models used,
406 inclusion/exclusion criteria, ground truth label assignments, and imaging protocols (**Tables 1**
407 **and 2**).

408 **CONCLUSION**

409 AI prediction capabilities were excellent and may enhance the diagnosis of ACL and meniscus
410 pathology; however, AI did not outperform clinical experts.

Journal Pre-proof

411 REFERENCES

412

- 413 1. Bini SA. Artificial Intelligence, Machine Learning, Deep Learning, and Cognitive
414 Computing: What Do These Terms Mean and How Will They Impact Health Care? *J*
415 *Arthroplasty*. 2018;33:2358-2361.
- 416 2. Helm JM, Swiergosz AM, Haeberle HS, et al. Machine Learning and Artificial Intelligence:
417 Definitions, Applications, and Future Directions. *Curr Rev Musculoskelet Med*.
418 2020;13:69-76.
- 419 3. Kunze KN, Karhade AV, Sadauskas AJ, Schwab JH, Levine BR. Development of Machine
420 Learning Algorithms to Predict Clinically Meaningful Improvement for the Patient-
421 Reported Health State After Total Hip Arthroplasty. *J Arthroplasty*. 2020.
- 422 4. Karhade AV, Ahmed AK, Pennington Z, et al. External validation of the SORG 90-day and
423 1-year machine learning algorithms for survival in spinal metastatic disease. *Spine J*.
424 2020;20:14-21.
- 425 5. Karhade AV, Schwab JH, Bedair HS. Development of Machine Learning Algorithms for
426 Prediction of Sustained Postoperative Opioid Prescriptions After Total Hip Arthroplasty.
427 *J Arthroplasty*. 2019;34:2272-2277 e2271.
- 428 6. Karhade AV, Shah AA, Bono CM, et al. Development of machine learning algorithms for
429 prediction of mortality in spinal epidural abscess. *Spine J*. 2019;19:1950-1959.
- 430 7. Ogunyemi O, Kermah D. Machine Learning Approaches for Detecting Diabetic
431 Retinopathy from Clinical and Public Health Records. *AMIA Annu Symp Proc*.
432 2015;2015:983-990.
- 433 8. Gulshan V, Peng L, Coram M, et al. Development and Validation of a Deep Learning
434 Algorithm for Detection of Diabetic Retinopathy in Retinal Fundus Photographs. *JAMA*.
435 2016;316:2402-2410.
- 436 9. Esteva A, Kuprel B, Novoa RA, et al. Dermatologist-level classification of skin cancer with
437 deep neural networks. *Nature*. 2017;542:115-118.
- 438 10. Thio Q, Karhade AV, Ogink PT, et al. Development and Internal Validation of Machine
439 Learning Algorithms for Preoperative Survival Prediction of Extremity Metastatic
440 Disease. *Clin Orthop Relat Res*. 2019.
- 441 11. Thio Q, Karhade AV, Ogink PT, et al. Can Machine-learning Techniques Be Used for 5-
442 year Survival Prediction of Patients With Chondrosarcoma? *Clin Orthop Relat Res*.
443 2018;476:2040-2048.
- 444 12. Langerhuizen DWG, Janssen SJ, Mallee WH, et al. What Are the Applications and
445 Limitations of Artificial Intelligence for Fracture Detection and Classification in
446 Orthopaedic Trauma Imaging? A Systematic Review. *Clin Orthop Relat Res*.
447 2019;477:2482-2491.
- 448 13. Moher D, Liberati A, Tetzlaff J, Altman DG, Group P. Preferred reporting items for
449 systematic reviews and meta-analyses: the PRISMA statement. *BMJ*. 2009;339:b2535.
- 450 14. Harris JD, Quatman CE, Manring MM, Siston RA, Flanigan DC. How to write a systematic
451 review. *Am J Sports Med*. 2014;42:2761-2768.
- 452 15. Abdullah AA A-ZN. Design of an Intelligent Diagnostic System for Detection of Knee
453 Injuries. *Applied Mechanics and Materials*. 2013;339:219-224.

- 454 **16.** Bien N, Rajpurkar P, Ball RL, et al. Deep-learning-assisted diagnosis for knee magnetic
455 resonance imaging: Development and retrospective validation of MRNet. *PLoS Med.*
456 2018;15:e1002699.
- 457 **17.** Chang PD, Wong TT, Rasiej MJ. Deep Learning for Detection of Complete Anterior
458 Cruciate Ligament Tear. *J Digit Imaging.* 2019;32:980-986.
- 459 **18.** Couteaux V, Si-Mohamed S, Nempont O, et al. Automatic knee meniscus tear detection
460 and orientation classification with Mask-RCNN. *Diagn Interv Imaging.* 2019;100:235-
461 242.
- 462 **19.** Fu J LC, Wang C, Ou Y. Computer-aided diagnosis for knee meniscus tears in magnetic
463 resonance imaging. *Journal of Industrial and Production Engineering.* 2013;30:67-77.
- 464 **20.** Liu F GB, Zhou Z, Samsonov A, Rosas H, Lian K, Sharma R, Kanarek A, Kim J, Guermazi A,
465 Kijowski R. Fully Automated Diagnosis of Anterior Cruciate Ligament Tears on Knee MR
466 Images by Using Deep Learning. *Radiology: Artificial Intelligence.* 2019;1:1-10.
- 467 **21.** Mazlan SS AM, Kadir Bakti ZA. Anterior Cruciate Ligament (ACL) Injury Classification
468 System Using Support Vector Machine (SVM). *Proc. Int. Engin and Tech.* . 2017:1-5.
- 469 **22.** Pedoia V, Norman B, Mehany SN, Bucknor MD, Link TM, Majumdar S. 3D convolutional
470 neural networks for detection and severity staging of meniscus and PFJ cartilage
471 morphological degenerative changes in osteoarthritis and anterior cruciate ligament
472 subjects. *J Magn Reson Imaging.* 2019;49:400-410.
- 473 **23.** Roblot V, Giret Y, Bou Antoun M, et al. Artificial intelligence to diagnose meniscus tears
474 on MRI. *Diagn Interv Imaging.* 2019;100:243-249.
- 475 **24.** Stajduhar I, Mamula M, Miletic D, Unal G. Semi-automated detection of anterior
476 cruciate ligament injury from MRI. *Comput Methods Programs Biomed.* 2017;140:151-
477 164.
- 478 **25.** Zarandi MH, Khadangi A, Karimi F, Turksen IB. A Computer-Aided Type-II Fuzzy Image
479 Processing for Diagnosis of Meniscus Tear. *J Digit Imaging.* 2016;29:677-695.
- 480 **26.** Cook NR. Use and misuse of the receiver operating characteristic curve in risk
481 prediction. *Circulation.* 2007;115:928-935.
- 482 **27.** Senders JT, Arnaout O, Karhade AV, et al. Natural and Artificial Intelligence in
483 Neurosurgery: A Systematic Review. *Neurosurgery.* 2018;83:181-192.
- 484 **28.** Nam JG, Park S, Hwang EJ, et al. Development and Validation of Deep Learning-based
485 Automatic Detection Algorithm for Malignant Pulmonary Nodules on Chest Radiographs.
486 *Radiology.* 2019;290:218-228.
- 487 **29.** Lindsey R, Daluiski A, Chopra S, et al. Deep neural network improves fracture detection
488 by clinicians. *Proceedings of the National Academy of Sciences of the United States of*
489 *America.* 2018;115:11591-11596.
- 490 **30.** Topol EJ. High-performance medicine: the convergence of human and artificial
491 intelligence. *Nature medicine.* 2019;25:44-56.
- 492 **31.** Gale W, Oakden-Rayner L, Carneiro G, Bradley AP, Palmer LJ. Detecting hip fractures
493 with radiologist-level performance using deep neural networks. *arXiv preprint*
494 *arXiv:1711.06504.* 2017.
- 495 **32.** Lee H, Yune S, Mansouri M, et al. An explainable deep-learning algorithm for the
496 detection of acute intracranial haemorrhage from small datasets. *Nature biomedical*
497 *engineering.* 2019;3:173-182.

- 498 **33.** Collins GS, Moons KGM. Reporting of artificial intelligence prediction models. *Lancet*
499 *(London, England)*. 2019;393:1577-1579.
- 500 **34.** Ozenne B, Subtil F, Maucort-Boulch D. The precision–recall curve overcame the
501 optimism of the receiver operating characteristic curve in rare diseases. *Journal of*
502 *clinical epidemiology*. 2015;68:855-859.
- 503 **35.** Brier GW. Verification of forecasts expressed in terms of probability. *Monthly weather*
504 *review*. 1950;78:1-3.
- 505 **36.** Ramkumar PN, Kunze KN, Haeberle HS, et al. Clinical and Research Medical Applications
506 of Artificial Intelligence: Fundamentals for the Orthopaedic Surgeon. *Arthroscopy*. 2020.

507
508 **Figure legends:**

509 **Figure 1:** Preferred Reporting Items for Systematic Review and Meta-Analysis (PRISMA)
510 flowchart for included studies.

511 **Tables**512 **Table 1.** Artificial Intelligence and Methodology for Anterior Cruciate Ligament Studies

Study	Input Features	Imaging plane	Dataset size	Anatomic structure	Ground/truth label assignment	Output classes	AI models used	Pretrained CNN	Size training set	Size validation set/validation method	Size test set	Performance (accuracy/AUC)
Štadjuhar et al. 2017	MR	Sagittal	917	ACL	Two Radiologists	2	HOG+linSVM	NA	NA	10- fold cross-validation	NA	NA/0.894 (linear-kernel SVM+HOG: injury-detection problem)
							HOG+RF					NA/0.943 (linear-kernel SVM +HOG: complete rupture)
							GIST+rbfSVM					NA/0.884 (injury-detection)
							GIST+RF					NA/0.937 (complete rupture)
												NA/0.889 (injury-detection)
	NA/0.913 (complete rupture)											
												NA/0.880 (injury-detection)
												NA/0.895 (complete rupture)
Mazlan et al. 2017	MR	Sagittal	300	ACL	Two Medical Experts with >7 years of experience	3	SVM	NA	210 (70%)	30 (10%)	60 (20%)	100%/NA
Chang et al. 2019	MR	Coronal	260	ACL	Visual inspection by a board-certified subspecialist	2	CNN	ResNet-Derived, U-net	160	40/5-fold cross-validation	60	0.967/0.971

MSK radiologist												
Bien et al. 2018	MR	Sagittal, Coronal, Axial	1,370	ACL	Three MSK radiologists' majority vote, average 12 years in practice	3	CNN	AlexNet, MRNet	1,130	120	120	Model ACL tear: 0.867 [95%CI 0.794, 0.916]/0.965 [95% CI 0.938, 0.993] Model abnormality detection: NA/0.937 [95%CI 0.895,0.980]
Abdullah et al. 2013	MR	Sagittal	90	ACL	NA	3	BP ANN, k-NN	NA	72	NA/5-fold and 6-fold	18	BP ANN: 0.9444/NA k-NN: 0.878333/NA
Liu et al. 2019	MR	Sagittal	350	ACL	Orthopedic Surgeon + Fellowship trained MSK radiologist with >15 years of clinical experience	2	CNN	LeNet-5, YOLO, DenseNet, VGG16, AlexNet	200 (57%)	50 (14%)	100 (29%)	NA/0.98 (DenseNet 95% CI (0.93-1.0) p<0.001

513 MR, magnetic resonance; NA, not available; ACL, anterior cruciate ligament; CNN, convolutional neural network; ANN, artificial
 514 neural network; BP, back-propagation; k-NN, K-nearest neighbors; RF, random forest; HOG, histogram of oriented gradients; GIST,
 515 generalized search tree; rbf, radial basis function; MSK, musculoskeletal; YOLO, you only look once; VGG, visual graphics group (type
 516 of neural network architecture); AUC, area under the curve; AI, artificial intelligence.

517

518

519 **Table 2.** Artificial Intelligence and Methodology for Meniscus Studies

Study	Input Features	Imaging plane	Dataset size	Anatomic location	Ground truth label assignment	Output classes	AI models	Pretrained CNN	Size training set	Size validation set/validation method	Size test set	Performance (accuracy/AUC)
Pedoia et al. 2018	MR	Sagittal	1,478	Meniscus	Expert 1: >20 years experience, Expert 2: >10 years experience, Expert 3: <1 year training as a radiologist	2	2D U-Net, 3D CNN	NA	960 (65%)	295 (20%)	221 (15%)	Binary mode (lesion vs non-lesion): NA/0.89 (no lesion, mild-moderate lesion, severe lesion): 80.74%/NA, 78.02%/NA, 75.00%/NA
	MR	Sagittal	1,478	Meniscus	Expert 1: >20 years experience, Expert 2: >10 years experience, Expert 3: <1 year training as a radiologist	3	2D U-Net, 3D CNN, RF	NA	960 (65%)	295 (20%)	221 (15%)	80.74%/NA, 78.02%/NA, 75.00%/NA
Fazel-Zarandi et al. 2016	MR	Sagittal	248	Meniscus	NA	2	IT2FCM, IT2PCM, PNN	NA	198 (80%)	NA	50 (20%)	0 and 1 mode: 90%/NA Binary mode: 78%/NA
Couteaux et al. 2019	MR	Sagittal	1,128	Meniscus	NA	6	R-CNN, ConvNet	ResNet-101, ConvNet, R-CNN	246	54	700	NA/0.906
Fu et al. 2013	MR	Sagittal	166	Meniscus	NA	2	SVM	NA	NA	5-FCA	NA	SVM model: NA/0.727 SFFS+SVM: NA/0.912
Roblot et al. 2019	MR	Sagittal	2,246	Meniscus	CSV file	6	Fast RCNN, Faster RCNN	NA	1,123	NA	700	NA/0.90
Bien et al. 2018	MR	Sagittal, Coronal, Axial	1,370	Meniscus	Three MSK radiologists' majority vote, average 12 years in practice on an internal validation set of 120 exams	3	CNN	AlexNet, MRNet	1,130	120	120	Model Meniscal tear: 0.725 [95%CI 0.639, 0.797]/0.847 (95% CI 0.780-0.914);

520 MR, magnetic resonance; NA, not available; 2D, two-dimensional; 3D, three-dimensional; CNN, convolutional neural network; R-
521 CNN, regions with CNN; SVM, support vector machine; IT2FCM, Interval type-2 fuzzy c-means; PNN, probabilistic neural network;

522 SFSS, sequential floating forward selection; AI, artificial intelligence; AUC, area under the curve; MSK, musculoskeletal; CSV, comma-
523 separated values.
524

Journal Pre-proof

525 **Table 3.** Quality appraisal of included studies.

Study	Quality Appraisal Score
Abdullah et al. 2013	4 – Failure to describe both inclusion/exclusion criteria of input features and ground truth assignment
Fu et al. 2013	3 – Failure to describe inclusion/exclusion criteria of input features, ground truth assignment, and distribution of data
Fazel-Zarandi et al. 2016	4 – Failure to describe ground truth assignment and distribution of data
Mazlan et al. 2017	6
Štadjuhar et al. 2017	5 – Failure to describe distribution of data
Bien et al. 2018	5 – Failure to describe inclusion/exclusion criteria of input features
Pedoia et al. 2018	6
Chang et al. 2019	6
Couteaux et al. 2019	4 – Failure to describe inclusion/exclusion criteria of input features and ground truth assignment
Liu et al. 2019	6
Roblot et al. 2019	5 – Failure to describe inclusion/exclusion criteria of input features

526

