

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22

Title Page

Title:

Thoracic Aorta Measurement Extraction from Computed Tomography Radiology Reports Using Instruction Tuned Large Language Models

Authors:

Ely Erez¹, Sedem Dankwa¹, McKenzie Tuttle¹, Afsheen Nasir¹, Prashanth Vallabhajosyula¹, Eric B. Schneider², Roland Assi¹, Chin Siang Ong^{2,3}

¹Division of Cardiac Surgery, Yale School of Medicine, New Haven, CT, USA

²Department of Surgery, Yale School of Medicine, New Haven, CT, USA

³Harvard T.H. Chan School of Public Health, Boston, MA, USA

Word Count: 4527

Corresponding Author:

Chin Siang Ong, MBBS, PhD
Assistant Professor of Surgery
Division of Surgical Outcomes
Surgery Center for Health Services and Outcomes Research
Department of Surgery, Yale School of Medicine
Phone: 203-432-4771
Email Address: chinsiang.ong@yale.edu

23 **Abstract:**

24 Chest computed tomography (CT) is essential for diagnosing and monitoring thoracic aortic
25 dilations and aneurysms, conditions that place patients at risk of complications such as aortic
26 dissection and rupture. However, aortic measurements in chest CT radiology reports are often
27 embedded in free-text formats, limiting their accessibility for clinical care, quality improvement
28 and research purposes. In this study, we developed a multi-method pipeline to extract structured
29 aortic measurements from radiology reports, and compared the performance of fine-tuned
30 BERT-based models with instruction-tuned Llama large language models (LLMs). Applying the
31 best-performing method to a real-world large chest CT radiology report database, we generated a
32 comprehensive aortic measurement dataset that facilitates big data aortic disease research.

33 **Introduction:**

34 Chest computed tomography (CT) is essential for diagnosing and monitoring thoracic aortic
35 dilations and aneurysms. These conditions, often asymptomatic, significantly increase the risk of
36 life-threatening complications such as aortic dissection and rupture¹. Frequently, thoracic aortic
37 dilations are detected incidentally on chest CTs performed for unrelated clinical indications.
38 Once diagnosed, chest CTs serve as the preferred modality for tracking aortic diameter changes
39 over time and guiding surgical decision-making¹. However, key diagnostic measures, including
40 aortic diameters, are predominantly documented in free-text CT radiology reports, rendering
41 them inaccessible in a structured format within most electronic health record systems.

42

43 Extracting these measures in a structured form could make a substantial clinical impact. These
44 structured data can be used to flag reports with incidental findings within the EHR software,
45 drawing the attention of physicians and facilitating timely referrals to surgeons. This can prevent
46 missed findings, improve the quality of care by monitoring referral rates, and adapt to evolving
47 evidence-based medicine (EBM) guidelines that adjust surgical referral thresholds. Additionally,
48 healthcare institutions can monitor surgical referral rates and ensure adherence to the most
49 current EBM, enhancing patient outcomes and healthcare quality. From the research perspective,
50 extracting these measures retrospectively in a large cohort of patients will enable researchers to
51 study rates of incidental aortic dilation detection and analyze how aortic diameters evolve over
52 time in affected patients. This structured data can also provide valuable epidemiological insights
53 into the prevalence and risk factors associated with thoracic aortic dilations.

54

55 The task of extracting aortic diameters falls within the broader domain of information extraction,
56 which is commonly addressed using natural language processing (NLP) techniques. Specifically,
57 it requires word- or token-level classification, akin to named entity recognition (NER), where
58 specific entities are identified and categorized within text. NER techniques have evolved rapidly
59 in recent years, transitioning from hand-crafted rule-based approaches to feature-driven
60 statistical models and, ultimately, to end-to-end deep learning models². As a key application of
61 NER, Biomedicine has been a central area of research³, with extensive work exploring its use in
62 medical imaging⁴. Early studies relied on hand-crafted rules to identify entity mentions, often
63 involving complex logic and requiring substantial domain expertise^{5,6}. Later work employed
64 more advanced classical machine learning techniques such as Hidden Markov Models⁷, Support
65 Vector Machines⁸, and Conditional Random Fields⁹. While these methods improved
66 performance, they still depended on manual feature extraction, a process heavily reliant on
67 domain knowledge, as a prerequisite to classification.

68
69 The advent of deep neural networks marked a paradigm shift in NER. Models such as Recurrent
70 Neural Networks (RNNs) and Bidirectional Long Short-Term Memory (BiLSTM) networks
71 achieved state-of-the-art performance while eliminating the need for manual feature engineering,
72 enabling end-to-end learning directly from raw data. These models, in turn, were rapidly
73 surpassed by pretrained language models like the Bidirectional Encoder Representations from
74 Transformers (BERT) model, which leverage encoder-based transformer architectures to set new
75 benchmarks in NER tasks. BERT and its domain specific variants, such as PubMedBERT¹⁰ and
76 BioBERT¹¹, have become the gold standard in NER. By fine-tuning on pre-labeled datasets, these
77 models achieve high performance in NER tasks with minimal reliance on additional domain

78 knowledge. These models have been widely adopted by the medical community for information
79 extraction³, including significant effort to improve information extraction from radiology
80 reports⁴. For instance, Khurshid et al.¹² developed Bio+Discharge Summary BERT, a fine-tuned
81 BERT model for NER, to extract vital signs such as height, weight, and blood pressure from
82 unstructured electronic health record (EHR) notes, reducing vital sign data missingness by 31%.
83 Similarly, Singh et al.¹³ addressed a challenge closely related to ours, using a fine-tuned BERT
84 model to extract 21 quantitative measures from cardiac magnetic resonance imaging. Their
85 approach achieved a macro-average F1 score of 0.957, highlighting the strong performance of
86 these models with minimal labeling effort.

87
88 In recent years generative LLMs such as OpenAI's GPT-3.5 and GPT-4¹⁴, have transformed the
89 field of natural language processing, breaking performance records across numerous
90 benchmarking tasks. These models utilize a Transformer architecture scaled to hundreds of
91 billions of parameters, enabling unprecedented contextual understanding and fluency in text
92 generation. The introduction of smaller, open-source LLMs such as Meta's Llama¹⁵, expanded
93 research opportunities by making these technologies accessible to the research community
94 enabling the development of tools for fine-tuning and inference on more modest hardware. The
95 ability of pretrained LLMs to perform zero-shot¹⁶ and few-shot¹⁷ learning has allowed them to
96 excel in many NLP tasks without extensive labeled data. Recognizing this potential, researchers
97 have sought to apply similar approaches to biomedical NER^{18,19}, hoping to reduce reliance on
98 labeled training data while maintaining competitive performance.

99

100 In practice, pretrained LLMs have struggled to match the performance of fine-tuned BERT-
101 based models on NER tasks^{20,21}. Instruction tuning, in contrast, has shown far greater
102 potential^{22,23}. Instruction tuning, or instruction fine-tuning, describes the process by which large
103 language models undergo supervised fine-tuning to better follow instructions and improve
104 performance on a given task²⁴. An example of this technique is a study by Keloth et al.²⁵, who
105 developed BioNER-Llama 2 by instruction-tuning Llama 2-7B on three publicly available
106 biomedical NER datasets focusing on diseases, chemicals and genes. BioNER-Llama 2 achieved
107 performance comparable to fine-tuned PubMedBERT, with F1 scores ranging from 0.949 to
108 0.956 on the test sets. Similarly, Bian et al.²⁶ created the VANER model by instruction-tuning
109 Llama 2-7B on eight biomedical NER datasets, reporting F1 scores between 0.77 and 0.94,
110 consistent with fine-tuned BERT-based models. Notably, both VANER and BioNER-Llama 2
111 demonstrated poor generalizability, with significantly reduced performance on previously unseen
112 datasets. Additionally, these studies primarily focus on benchmark tasks, which, while useful for
113 assessing baseline capabilities, may not reflect real-world complexities. The real-world
114 performance of LLMs in NER applications, particularly in the biomedical domain, remains
115 largely underexplored.

116
117 The primary objective of this study was to develop an automated machine learning pipeline for
118 extracting aortic measurements from chest CT radiology reports. To achieve this, we compared
119 the performance of fine-tuned BERT-based models with generative LLMs. A secondary
120 objective was to construct a comprehensive aortic measurement database by applying the
121 pipeline to a large cohort of chest CT radiology reports from our institution. The generated

122 dataset will enable future investigations into patterns of aortic dilation detection and the
123 progression of aortic disease in a hospital-based population.

124

125 **Methodology:**

126 This study was determined to be exempt from review by Yale University's Institutional Review
127 Board (IRB) under protocol number 2000037866 on May 3, 2024.

128

129 Dataset

130 We conducted an institutional search for chest CTs with corresponding radiology reports for
131 patients aged 18 and older, performed between January 2013 and December 2023. The search
132 encompassed 43 distinct CT protocols and yielded 363,423 radiology reports. Reports from CT
133 protocols with fewer than 2,000 instances and reports without narrative content were excluded,
134 resulting in a final dataset of 356,690 radiology reports across 16 CT protocols.

135

136 A subset of 1,506 radiology reports was selected for manual annotation using stratified random
137 sampling to ensure balanced representation across protocols. This subset was divided into 1,002
138 reports for training (sampled at a 1:356 ratio) and 504 reports for validation and testing (sampled
139 at a 1:712 ratio), ensuring that reports in the training set were distinct from the reports used for
140 validation and testing to prevent information leakage. The narratives were annotated using Label
141 Studio²⁷ by two medical students and a postdoctoral researcher with an MD. To ensure
142 consistency, all annotations were reviewed by the postdoctoral researcher. Thoracic aortic
143 diameters were labeled at eight anatomical sites: the annulus, sinus of Valsalva, sinotubular

144 junction, mid ascending, ascending proximal to the brachiocephalic, top of the arch, proximal
145 descending, and mid descending.

146

147 Due to the limited number of aortic diameter annotations identified in the 504 reports intended
148 for validation and testing (n=289 annotations), the reports designated for both validation and
149 testing were instead allocated exclusively for validation. To create the testing set, an additional
150 504 reports were sampled and annotated following the same protocol. Reports selected for
151 training and validation were excluded from the pool when selecting the test set to avoid data
152 leakage. This process resulted in 2,010 labeled reports, divided into training, validation, and
153 testing sets with a 50:25:25 split. Figure 1 provides a flowchart illustrating the dataset selection
154 process, while Table 1 presents the radiology reports and corresponding patient characteristics
155 for each set.

156

157 Data preprocessing

158 To account for BERT's token limit and ensure a fair comparison with Llama models, we opted to
159 perform fine-tuning and inference on individual sentences rather than complete narratives.

160 Because most sentences in the report narratives did not contain aortic measurements, splitting the
161 reports into sentences allowed us to exclude irrelevant content, thereby improving label balance
162 and significantly reducing fine-tuning and inference times. Sentence splitting was performed
163 using the Python package NLTK. The first sentence of each report, as well as any sentence
164 lacking a non-time or date numeric value, was excluded from analysis. Additionally, we retained
165 only sentences containing at least one aorta-related keyword.

166

167 For fine-tuning BERT-based models, each aortic measurement site was assigned a numeric label
168 between 1 and 8. Tokens within the span of a measurement were assigned the corresponding
169 numeric label, while all other tokens were labeled as 0. This numerical vector served as the label
170 for each sentence. For Llama models we used XML tags to delineate measurements in the target
171 output, such as <SOV> and </SOV> for the sinus of Valsalva. If no measurements were present,
172 the output remained identical to the input sentence. This annotation approach, similar to that
173 employed by Keltoh et al.²⁵, facilitates straightforward postprocessing. Figure 2 shows a sample
174 input and output for Llama models.

175

176 Baseline Model

177 We used Meta's Llama 3.1 instruction-tuned model with 8 billion (8B) parameters as a baseline
178 model²⁸. This 8B model was chosen for its strong benchmark performance relative to its size,
179 which represents the upper limit of our virtual machine's capacity for local fine-tuning. The
180 instruction-tuned version was chosen over the base (pre-trained) version because it underwent
181 several rounds of alignment, including supervised fine-tuning, rejection sampling, and direct
182 preference optimization, which improved its instruction-following capability, quality, and
183 safety²⁹. The baseline model's performance was further optimized through prompt engineering,
184 following the methodology described by Hu et al.¹⁹. The prompt included a task description,
185 labeling instructions based on the annotation guidelines, and additional instructions informed by
186 error analysis conducted on the training set. We evaluated the baseline model's zero-shot
187 performance as well as few-shot performance using three pre-selected annotated samples from
188 the training set. The selected prompt, as well as several annotated radiology report samples, are
189 publicly available on GitHub at <https://github.com/yalesurgeryresearch/RadTextExtractor>.

190

191 BERT-based Models

192 We fine-tuned six BERT-based models by combining three weight initialization strategies with
193 two tokenization schemes. The weight initialization strategies included: (1) the original BERT
194 model³⁰, (2) BERT-NER³¹, a variant fine-tuned on the English CoNLL-2003 Named Entity
195 Recognition dataset³², and (3) PubMedBERT¹⁰, pre-trained on PubMed abstracts and full-text
196 articles from PubMed Central. The two tokenization schemes tested were: (1) the standard BERT
197 tokenization and (2) a modified scheme in which numeric values were replaced with the unique
198 [NUM] pseudo-token³³.

199

200 Fine-tuning was performed using the `AutoModelForTokenClassification` class from the
201 HuggingFace transformers library³⁴, which adds a token classification head to the BERT
202 architecture for mapping hidden states to output labels. The fine-tuning process adhered to the
203 baseline scheme proposed by Mosbach et al.³⁵. Each model was fine-tuned on the selected
204 sentences from the training set over 20 epochs, using categorical cross-entropy loss.
205 Optimization was conducted with the AdamW³⁶ optimizer, employing a linear learning rate
206 scheduler and a 10% warm-up phase. Hyperparameter tuning was conducted using grid search,
207 testing learning rates of 1e-5, 2e-5, 5e-5 and 1e-4, along with batch sizes of 8, 16, and 32. A
208 random seed was fixed during training to ensure reproducibility. Loss on the validation set was
209 calculated after each epoch, and the epoch with the lowest validation loss was selected for each
210 fine-tuning iteration.

211

212 Llama Models

213 We compared the performance of three versions of Meta’s Llama models using instruction-
214 tuning: the Llama 2 chat-tuned model with 7 billion parameters, the Llama 3 instruction-tuned
215 model with 8 billion parameters, and the Llama 3.1 instruction-tuned model with 8 billion
216 parameters. To accommodate the instruction-tuning process on our virtual machine’s limited
217 GPU RAM, we employed 4-bit Quantized Low-Rank Adaptation (QLoRA)³⁷. QLoRA builds on
218 Low-Rank Adaptation (LoRA)³⁸, a technique that significantly reduces the number of trainable
219 parameters, by further quantizing model weights to 4-bit precision. This approach enables fine-
220 tuning of large language models on resource-constrained hardware while maintaining high
221 performance. Additionally, we utilized Unsloth³⁹, an open-source library that accelerates fine-
222 tuning through a custom backpropagation engine.

223
224 Instruction-tuning for each model used the same prompt as the zero-shot baseline Llama 3.1
225 model. Instruction-tuning was conducted on sentences from the training set over 5 epochs, using
226 cross-entropy loss and an AdamW³⁶ optimizer with a linear rate scheduler without a warm-up
227 phase. A batch size of 1 was used for all trials. During fine-tuning, we zeroed-out the loss on the
228 provided prompt, ensuring learning only on the model’s output. Validation loss was calculated
229 after each epoch, and the epoch with the lowest validation loss was selected for each instruction-
230 tuning trial. Hyperparameter tuning was performed using grid search to optimize the learning
231 rate, as well as LoRA’s rank, and alpha parameters. The learning rates tested were 2e-5, 5e-5,
232 and 1e-4. Rank values included 16, 32, and 64, with alpha values set to rank multiplied by 1 or 2
233 (e.g., for a rank of 16, the alpha values tested were 16 and 32). Sampling was disabled during
234 inference to ensure deterministic and reproducible results. When labeling sentences, the model
235 was provided with the prompt and a sentence as input and generated a labeled output.

236

237 Evaluation metrics

238 Model performance was evaluated on the validation set using exact match macro-averaged

239 precision, recall, and F1 scores across all aortic measurement sites, as well as site-specific

240 precision, recall, and F1 scores. To simplify terminology, macro-averaged precision, recall, and

241 F1 scores are referred to as "macro precision," "macro recall," and "macro F1," respectively.

242 Performance metrics were calculated across the entire validation set, including sentences not

243 selected for inference during preprocessing. The optimal model from the baseline models,

244 BERT-based models, and Llama models was selected based on the highest macro F1 score on the

245 validation set.

246

247 Ablation study

248 To evaluate the impact of training set size on model performance, we conducted an ablation

249 study in which the optimal model was fine-tuned using subsets of 10, 25, 50, and 100 randomly

250 selected sentences from the training set. To mitigate the effects of sentence selection, this process

251 was repeated five times, each time using a different random seed to generate distinct sentence

252 subsets. The median macro F1 score across the five trials was then compared to the performance

253 of the model trained on the full training set.

254

255 Following the ablation study, we evaluated the optimal model on the test set to assess its

256 generalizability and potential real-world performance. Finally, inference was conducted on the

257 entire radiology report cohort to create an aortic measurement database. This process was limited

258 to sentences selected according to the criteria outlined in the preprocessing phase. Measurement

259 extraction rates, aortic dilation rates, and aortic diameter median and interquartile range (IQR) at
260 each measurement site were analyzed as supplementary indicators of the model's real-world
261 performance.

262

263 Computational Resources and Framework

264 All analyses were conducted on a HIPAA-compliant virtual machine hosted by Yale's Spinup
265 service, utilizing Amazon Elastic Cloud Compute (EC2). The environment comprised an
266 Amazon AWS G4 instance with 4 vCPUs, 16 GB of RAM, and an NVIDIA T4 GPU with 16 GB
267 of GPU memory. GPU acceleration was facilitated using CUDA (v12.4). Fine-tuning and
268 inference were performed using Python (v3.9.13) with the following key libraries: PyTorch
269 (v2.3.0), Hugging Face Transformers (v4.43.3), and Unsloth (v2024.8).

270

271 The code for data preprocessing, model fine-tuning, and evaluation is available in a GitHub
272 repository at <https://github.com/yalesurgeryresearch/RadTextExtractor>. Due to the presence of
273 protected health information in the radiology reports, the datasets generated and analyzed in this
274 study, along with the fine-tuned models, are not publicly available. However, they can be
275 obtained from the corresponding author upon reasonable request, in accordance with institutional
276 policies and any applicable data access agreements.

277

278 **Results:**

279 Following preprocessing, the training dataset used for fine-tuning consisted of 214 out of 19,844
280 sentences (1.08%), of which 166 (77.6%) contained at least one annotation. The training set
281 included a total of 589 annotations, with a median of 52 annotations per measurement site

282 (range: 44 to 166). The validation set contained 103 out of 9,505 sentences (1.08%) selected for
283 inference, of which 71 (68.9%) included at least one annotation. All sentences with annotations
284 in the validation set were included in the inference subset, which had a total of 289 annotations,
285 with a median of 25.5 annotations per measurement site (range: 20 to 76). The test set comprised
286 91 out of 9,664 sentences (0.94%) selected for inference, with 68 (74.7%) containing at least one
287 annotation. Similar to the validation set, all sentences with annotations in the test set were
288 included in the inference subset, which contained 215 annotations, with a median of 18
289 annotations per measurement site (range: 13 to 69). A summary of the annotation characteristics
290 for the training, validation, and test datasets after preprocessing is provided in Table 2.

291

292 Model Performance Comparison

293 The performance of the baseline Llama 3.1 models, fine-tuned BERT-based models, and
294 instruction-tuned Llama models on the validation set is summarized in Table 3. The few-shot
295 Llama 3.1 was the best performing baseline model, besting the zero-shot model with a macro F1
296 score of 0.838 compared to 0.663. However, both baseline models were surpassed by all fine-
297 tuned BERT-based models and instruction-tuned Llama models.

298

299 The best performance among the BERT-based models was achieved by the fine-tuned
300 PubMedBERT with [NUM] tokenization, which attained a macro F1 score of 0.945. Numeric
301 tokenization using a [NUM] pseudo-token consistently outperformed the standard tokenization
302 technique across all three models, and the fine-tuned PubMedBERT was the best performing
303 model in both the standard and [NUM] tokenization.

304

305 The instruction-tuned Llama 3.1 delivered the highest overall performance, achieving a near-
306 perfect macro F1 score of 0.992, with a macro precision of 0.993 and macro recall of 0.992. The
307 instruction-tuned Llama 3 followed closely with a macro F1 score of 0.982. Notably, the
308 instruction-tuned Llama 2 achieved a macro F1 of 0.839, performance comparable to the few-
309 shot baseline model and significantly lower than the top-performing BERT-based models.

310
311 Figure 3 compares the distributions of F1 scores across measurement sites for the few-shot
312 Llama 3.1, PubMedBERT with [NUM] tokenization, and the instruction-tuned Llama 3.1. The
313 few-shot Llama 3.1 and PubMedBERT models showed significant variability in F1 scores across
314 different measurement sites. In contrast, the instruction-tuned Llama 3.1 demonstrated consistent
315 performance, achieving an F1 score of at least 0.971 across all sites.

316

317 Impact of Training Set Size

318 Figure 4 shows the results of the ablation study, with macro F1 scores of the instruction-tuned
319 Llama 3.1 model as a function of the number of training sentences. Performance improved
320 rapidly with the initial increase in training sentences but slowed and eventually plateaued as the
321 number grew. Median macro F1 scores were 0.800 with 10 training sentences, 0.880 with 20,
322 0.903 with 50, and 0.971 with 100, compared to 0.992 when using the full training set of 214
323 sentences. The results also highlight significant variability in performance due to the random
324 selection of sentence subsets, which decreased as the training set size increased.

325

326 Assessing Model Generalizability

327 The performance of the instruction-tuned Llama 3.1 model was evaluated on the test set to assess
328 its generalizability to unseen data sampled from the same source distribution. Table 4 compares
329 the model's performance across aortic measurement sites between the validation and test sets.
330 The macro F1 score on the test set was 0.970, slightly lower than the 0.992 achieved on the
331 validation set. Measurement site F1 scores ranged from 0.923 to 1.000 on the test set, compared
332 to 0.971 to 1.000 on the validation set.

333

334 Insights from Full Dataset Extraction

335 The complete chest CT radiology report dataset consisted of 356,690 reports from 140,645
336 unique patients. Following preprocessing, 74,483 sentences out of 6,960,729 (approximately
337 1.07%) were selected for extraction, consistent with the proportions observed in the labeled
338 datasets. After extraction using the instruction-tuned Llama 3.1, 49,387 radiology reports
339 (13.85%) contained at least one aortic measurement, showing higher rates of aortic measurement
340 reporting in males (18.51%) compared to females (9.50%). Table 5 summarizes the extraction
341 results by aortic measurement site. Measurement extraction rates across the aortic sites were
342 similar to those in the labeled datasets (Table 2).

343

344 The largest median diameters were observed at the mid ascending aorta and the sinus of
345 Valsalva, measuring 39 mm (IQR 36–42) and 36 mm (IQR 32–40), respectively. Median
346 diameters decreased distally along the aorta, measuring 31 mm (IQR 28–34) at the aortic arch,
347 30 mm (IQR 27–34) at the proximal descending aorta, and 29 mm (IQR 25–34) at the mid
348 descending aorta. Ascending aortic dilation of at least 40 mm was reported in 8.69% of patients

349 (12,228/140,645), with 2.27% (3,193/140,645) reported to have a dilation of at least 45 mm,
350 0.66% (925/140,645) at least 50 mm, and 0.28% (393/140,645) at least 55 mm.

351

352 **Discussion:**

353 In this study, we describe our experiences developing a machine learning pipeline for extracting
354 aortic measurements from chest CT radiology reports. Among the models evaluated, the
355 instruction-tuned Llama 3.1 outperformed both the BERT-based models and the pretrained
356 Llama 3.1 baseline, achieving macro F1 scores of 0.992 on the validation set and 0.970 on the
357 test set. PubMedBERT achieved the best performance among the BERT-based models,
358 suggesting that pre-training on medical literature, making it better suited for understanding and
359 processing medical texts, such as chest CT radiology reports. The effectiveness of [NUM]
360 tokenization is likely attributed to its consistent numerical tokenization compared as compared to
361 the standard BERT WordPiece tokenizer, which fragments numerical expressions requiring that
362 all fragments be correctly tagged³³. Among the Llama-based models, the instruction-tuned Llama
363 3.1 significantly outperformed the Llama 2 chat-tuned model and was marginally better than the
364 Llama 3 instruction-tuned model. Meta attributes Llama 3.1's superior performance to its
365 enhanced reasoning capabilities and improved context length²⁹. These improvements appear to
366 have carried over in instruction-tuning, which allowed it to better handle the complexities of the
367 dataset and achieve higher accuracy in extracting aortic measurements from chest CT radiology
368 reports.

369

370 When applied to our extensive radiology report database, the model successfully extracted aortic
371 measurements from 13.85% of reports, a rate consistent with both our labeled subset and prior

372 work assessing aortic measurement reporting in CT radiology reports⁴⁰. The extracted aortic
373 measurements and dilation rates similarly aligned with findings from previous studies^{41,42}. The
374 resulting aortic measurement database is one of the largest of its kind, encompassing
375 measurements from nearly 50,000 CT scans and over 28,000 patients, representing a valuable
376 resource for advancing the study of aortic disease.

377
378 Initial enthusiasm for the potential of large language models (LLMs) in named entity recognition
379 (NER) has recently been tempered. The expectation that general-domain LLMs could achieve
380 domain-specific NER through in-context learning has been challenged by multiple studies, where
381 fine-tuned BERT-based models consistently outperform LLMs^{20,21,23}. Even with instruction-
382 tuning, LLMs have, at best, matched the performance of BERT-based models—a disappointing
383 outcome given their significantly larger parameter counts and the associated higher costs of fine-
384 tuning and inference^{23,25,26}. Researchers have suggested that the relatively poor NER
385 performance of LLMs may stem from the limitations of their decoder-only transformer
386 architecture and next-token prediction pretraining objective, compared to BERT’s encoder-only
387 architecture and masked language modeling pretraining objective⁴³. Our findings, however,
388 challenge this hypothesis. In our study, the instruction-tuned Llama 3.1 achieved near-perfect
389 performance on the NER task, surpassing the fine-tuned BERT-based models by what we
390 consider a substantial margin. In addition to its fantastic performance, Llama 3.1 offers several
391 additional advantages over BERT, including a much larger context length for analyzing longer
392 text segments and more human-interpretable outputs, which streamline error analysis.

393

394 Despite these excellent results, additional work is needed to discern whether instruction-tuned
395 generative LLMs have become the new gold standard for NER. Our findings are based on a
396 single dataset, and the observed differences might be attributed to suboptimal fine-tuning of the
397 BERT models, rather than the inherent superiority of the Llama 3 architecture. Further studies
398 replicating these results across additional NER datasets is essential to substantiate these claims.
399 Nonetheless, the ability of Llama 3 models to achieve this level of performance suggests that
400 instruction-tuned generative LLMs hold significant promise for NER and could play a valuable
401 role in clinical NER.

402
403 A significant advantage of our proposed methodology is its adaptability. The framework is
404 agnostic to both the entities being extracted and the domain, enabling straightforward adaptations
405 to various NER tasks. Instruction-tuning requires relatively few annotated samples, and open-
406 source annotation tools such as Label Studio facilitate efficient, collaborative annotation
407 processes. Frameworks like Hugging Face’s Transformers library offer well-developed pipelines
408 for instruction-tuning general-domain LLMs, making them easily adaptable to diverse tasks.
409 However, several barriers remain to the broader adoption of these techniques. LLMs still demand
410 substantial computational resources for training and inference. For clinical projects, the
411 additional requirement for HIPAA-compliant hardware introduces further costs and complexity.
412 While existing pipelines are robust, they often require advanced coding and machine learning
413 expertise, which may be beyond the scope of many clinical researchers. As the field of LLMs
414 continues to evolve, these barriers are likely to diminish. Companies such as Microsoft and
415 OpenAI are actively developing HIPAA-compliant implementations of their LLMs, and costs are
416 expected to decrease as competition increases and the technology matures. If these trends persist,

417 we anticipate that access to LLM instruction-tuning will become increasingly democratized,
418 empowering clinical researchers to leverage these powerful tools.

419
420 Our study has several limitations. Both the validation and test sets are relatively small, with few
421 annotations, making the results susceptible to variability as one or two errors can significantly
422 impact model performance. Additionally, selecting a subset of sentences for inference may have
423 led to the omission of relevant sentences when extracting measurements from the complete
424 radiology report dataset. The BERT-based models used in our analysis are known to be sensitive
425 to seed values⁴⁴, which may have influenced their performance. Another limitation is the
426 potential lack of generalizability to newly collected data. The validation and test sets share a
427 temporal distribution with the training set, and medical data is prone to domain drift over time⁴⁵.
428 This could limit the applicability of our findings to datasets collected in different time periods or
429 settings. We believe these limitations do not detract significantly from the overall value of our
430 findings. Replication of our study in different datasets and settings is needed to validate our
431 results and confirm the generalizability of our approach.

432

433 **Conclusion:**

434 In this study, we developed and evaluated a machine learning pipeline for extracting aortic
435 measurements from chest CT radiology reports. The instruction-tuned Llama model achieved the
436 best performance, surpassing state-of-the-art BERT-based models. Using this pipeline, we
437 created a large, comprehensive database of aortic measurements from radiology reports, offering
438 a valuable resource for aortic research. Our results highlight the potential of instruction-tuned
439 generative LLMs in the NER domain, with a generalizable workflow that requires few labeled

440 samples and modest computational resources. As the technology matures, this process is
441 expected to become even more streamlined, enabling broader adoption in clinical research.

442

443 **Acknowledgments**

444 R.A. discloses support for the research of this work from the Yale Department of Surgery and the
445 National Heart, Lung, and Blood Institute (NHLBI) of the National Institutes of Health (NIH)
446 [grant number R01HL168473].

447

448 **Competing Interests**

449 All authors declare no financial or non-financial competing interests.

450

451 **Author Contributions**

452 Conceptualization: EE, CSO, RA.

453 Data Curation: EE, SD, MT, AN.

454 Formal Analysis: EE.

455 Investigation: EE.

456 Methodology: EE, CSO.

457 Project Administration: EE, CSO.

458 Resources: EBS, CSO.

459 Supervision: RA, CSO.

460 Writing – Original Draft: EE, CSO.

461 Writing – Review & Editing: CSO, RA, PV, EBS.

462 All authors reviewed the manuscript.

463

464 **Data Availability**

465 Study data are available upon reasonable request from the corresponding author, in accordance
466 with institutional policies and any applicable data sharing or data use agreements.

467

468 **Code Availability**

469 The code used in this study is open source and freely available under the MIT license. It can be
470 accessed on GitHub at <https://github.com/yalesurgeryresearch/RadTextExtractor/>. This
471 repository includes detailed documentation and examples to facilitate reproducibility and
472 adaptation for related research.

473

474 **References:**

- 475 1. Isselbacher, E. M. *et al.* 2022 ACC/AHA Guideline for the Diagnosis and Management of
476 Aortic Disease: A Report of the American Heart Association/American College of
477 Cardiology Joint Committee on Clinical Practice Guidelines. *Circulation* **146**, e334–e482
478 (2022).
- 479 2. Munnangi, M. A Brief History of Named Entity Recognition. Preprint at
480 <https://doi.org/10.48550/arXiv.2411.05057> (2024).
- 481 3. Nunes, M., Bone, J., Ferreira, J. C. & Elvas, L. B. Health Care Language Models and Their
482 Fine-Tuning for Information Extraction: Scoping Review. *JMIR Med. Inform.* **12**, e60164
483 (2024).
- 484 4. Hu, M. *et al.* Advancing medical imaging with language models: featuring a spotlight on
485 ChatGPT. *Phys. Med. Biol.* **69**, 10TR01 (2024).
- 486 5. Tsuruoka, Y. & Tsujii, J. Boosting precision and recall of dictionary-based protein name
487 recognition. in *Proceedings of the ACL 2003 workshop on Natural language processing in*

- 488 *biomedicine - Volume 13* 41–48 (Association for Computational Linguistics, USA, 2003).
489 doi:10.3115/1118958.1118964.
- 490 6. Hanisch, D., Fundel, K., Mevissen, H.-T., Zimmer, R. & Fluck, J. ProMiner: rule-based
491 protein and gene entity recognition. *BMC Bioinformatics* **6**, S14 (2005).
- 492 7. Collier, N., Nobata, C. & Tsujii, J. Extracting the names of genes and gene products with a
493 hidden Markov model. in *Proceedings of the 18th conference on Computational linguistics -*
494 *Volume 1* 201–207 (Association for Computational Linguistics, USA, 2000).
495 doi:10.3115/990820.990850.
- 496 8. Kazama, J., Makino, T., Ohta, Y. & Tsujii, J. Tuning support vector machines for biomedical
497 named entity recognition. in *Proceedings of the ACL-02 workshop on Natural language*
498 *processing in the biomedical domain - Volume 3* 1–8 (Association for Computational
499 Linguistics, USA, 2002). doi:10.3115/1118149.1118150.
- 500 9. Settles, B. Biomedical named entity recognition using conditional random fields and rich
501 feature sets. in *Proceedings of the International Joint Workshop on Natural Language*
502 *Processing in Biomedicine and its Applications* 104–107 (Association for Computational
503 Linguistics, USA, 2004).
- 504 10. Gu, Y. *et al.* Domain-Specific Language Model Pretraining for Biomedical Natural
505 Language Processing. *ACM Trans. Comput. Healthc.* **3**, 1–23 (2022).
- 506 11. Lee, J. *et al.* BioBERT: a pre-trained biomedical language representation model for
507 biomedical text mining. *Bioinformatics* **36**, 1234–1240 (2020).
- 508 12. Khurshid, S. *et al.* Cohort design and natural language processing to reduce bias in electronic
509 health records research. *Npj Digit. Med.* **5**, 1–14 (2022).

- 510 13. Singh, P. *et al.* One Clinician Is All You Need—Cardiac Magnetic Resonance Imaging
511 Measurement Extraction: Deep Learning Algorithm Development. *JMIR Med. Inform.* **10**,
512 e38178 (2022).
- 513 14. OpenAI *et al.* GPT-4 Technical Report. Preprint at
514 <https://doi.org/10.48550/arXiv.2303.08774> (2024).
- 515 15. Touvron, H. *et al.* LLaMA: Open and Efficient Foundation Language Models. Preprint at
516 <https://doi.org/10.48550/arXiv.2302.13971> (2023).
- 517 16. Kojima, T., Gu, S. (Shane), Reid, M., Matsuo, Y. & Iwasawa, Y. Large Language Models
518 are Zero-Shot Reasoners. *Adv. Neural Inf. Process. Syst.* **35**, 22199–22213 (2022).
- 519 17. Brown, T. B. *et al.* Language Models are Few-Shot Learners. Preprint at
520 <https://doi.org/10.48550/arXiv.2005.14165> (2020).
- 521 18. Agrawal, M., Hegselmann, S., Lang, H., Kim, Y. & Sontag, D. Large Language Models are
522 Few-Shot Clinical Information Extractors. Preprint at
523 <https://doi.org/10.48550/arXiv.2205.12689> (2022).
- 524 19. Hu, Y. *et al.* Improving large language models for clinical named entity recognition via
525 prompt engineering. *J. Am. Med. Inform. Assoc.* ocad259 (2024) doi:10.1093/jamia/ocad259.
- 526 20. Wang, S. *et al.* GPT-NER: Named Entity Recognition via Large Language Models. Preprint
527 at <https://doi.org/10.48550/arXiv.2304.10428> (2023).
- 528 21. Xie, T. *et al.* Empirical Study of Zero-Shot NER with ChatGPT. Preprint at
529 <https://doi.org/10.48550/arXiv.2310.10035> (2023).
- 530 22. Chen, Q. *et al.* A systematic evaluation of large language models for biomedical natural
531 language processing: benchmarks, baselines, and recommendations. Preprint at
532 <https://doi.org/10.48550/arXiv.2305.16326> (2024).

- 533 23. Xu, D. *et al.* Large Language Models for Generative Information Extraction: A Survey.
534 Preprint at <https://doi.org/10.48550/arXiv.2312.17617> (2024).
- 535 24. Wei, J. *et al.* Finetuned Language Models Are Zero-Shot Learners. Preprint at
536 <https://doi.org/10.48550/arXiv.2109.01652> (2022).
- 537 25. Keloth, V. K. *et al.* Advancing entity recognition in biomedicine via instruction tuning of
538 large language models. *Bioinformatics* **40**, btae163 (2024).
- 539 26. Biana, J., Zhai, W., Huang, X., Zheng, J. & Zhu, S. VANER: Leveraging Large Language
540 Model for Versatile and Adaptive Biomedical Named Entity Recognition. Preprint at
541 <https://doi.org/10.48550/arXiv.2404.17835> (2024).
- 542 27. Tkachenko, M., Malyuk, M., Holmanyuk, A. & Liubimov, N. Label Studio: Data labeling
543 software. (2020).
- 544 28. Dubey, A. *et al.* The Llama 3 Herd of Models. Preprint at
545 <https://doi.org/10.48550/arXiv.2407.21783> (2024).
- 546 29. Introducing Llama 3.1: Our most capable models to date. *Meta AI*
547 <https://ai.meta.com/blog/meta-llama-3-1/>.
- 548 30. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of Deep
549 Bidirectional Transformers for Language Understanding. Preprint at
550 <https://doi.org/10.48550/arXiv.1810.04805> (2019).
- 551 31. Lim, David. dslim/bert-base-NER · Hugging Face. <https://huggingface.co/dslim/bert-base->
552 [ner](https://huggingface.co/dslim/bert-base-ner) (2024).
- 553 32. Sang, E. F. T. K. & Meulder, F. D. Introduction to the CoNLL-2003 Shared Task: Language-
554 Independent Named Entity Recognition. Preprint at
555 <https://doi.org/10.48550/arXiv.cs/0306050> (2003).

- 556 33. Loukas, L. *et al.* FiNER: Financial Numeric Entity Recognition for XBRL Tagging. Preprint
557 at <https://doi.org/10.48550/arXiv.2203.06482> (2022).
- 558 34. Wolf, T. *et al.* HuggingFace’s Transformers: State-of-the-art Natural Language Processing.
559 Preprint at <https://doi.org/10.48550/arXiv.1910.03771> (2020).
- 560 35. Mosbach, M., Andriushchenko, M. & Klakow, D. On the Stability of Fine-tuning BERT:
561 Misconceptions, Explanations, and Strong Baselines. Preprint at
562 <https://doi.org/10.48550/arXiv.2006.04884> (2021).
- 563 36. Loshchilov, I. & Hutter, F. Decoupled Weight Decay Regularization. Preprint at
564 <https://doi.org/10.48550/arXiv.1711.05101> (2019).
- 565 37. Dettmers, T., Pagnoni, A., Holtzman, A. & Zettlemoyer, L. QLoRA: Efficient Finetuning of
566 Quantized LLMs. *Adv. Neural Inf. Process. Syst.* **36**, 10088–10115 (2023).
- 567 38. Hu, E. J. *et al.* LoRA: Low-Rank Adaptation of Large Language Models. Preprint at
568 <https://doi.org/10.48550/arXiv.2106.09685> (2021).
- 569 39. Daniel Han and Michael Han. unslothai/unsloth. Unsloth AI (2024).
- 570 40. Zamirpour, S. *et al.* Sex differences in ascending aortic size reporting and growth on chest
571 computed tomography and magnetic resonance imaging. *Clin. Imaging* **105**, 110021 (2024).
- 572 41. Benedetti, N. & Hope, M. D. Prevalence and Significance of Incidentally Noted Dilation of
573 the Ascending Aorta on Routine Chest Computed Tomography in Older Patients. *J. Comput.*
574 *Assist. Tomogr.* **39**, 109 (2015).
- 575 42. Mori, M. *et al.* Prevalence of Incidentally Identified Thoracic Aortic Dilations: Insights for
576 Screening Criteria. *Can. J. Cardiol.* **35**, 892–898 (2019).
- 577 43. Lu, Q. *et al.* Large Language Models Struggle in Token-Level Clinical Named Entity
578 Recognition. Preprint at <https://doi.org/10.48550/arXiv.2407.00731> (2024).

- 579 44. Dodge, J. *et al.* Fine-Tuning Pretrained Language Models: Weight Initializations, Data
580 Orders, and Early Stopping. Preprint at <https://doi.org/10.48550/arXiv.2002.06305> (2020).
- 581 45. Guo, L. L. *et al.* Systematic Review of Approaches to Preserve Machine Learning
582 Performance in the Presence of Temporal Dataset Shift in Clinical Medicine. *Appl. Clin.*
583 *Inform.* **12**, 808–815 (2021).
- 584

585 **Tables:**

586 **Table 1. Train, validation and test dataset radiology report characteristics**

	Train	Validation	Test
Radiology reports (N)	1002	504	504
CT type (n [%])			
CT chest without IV contrast	288 [28.7]	144 [28.6]	144 [28.6]
CT chest with IV contrast	196 [19.6]	98 [19.4]	98 [19.4]
CTA chest (PE) with IV contrast	172 [17.2]	86 [17.1]	86 [17.1]
CT chest, abdomen, pelvis with IV contrast	157 [15.7]	79 [15.7]	79 [15.7]
CT ED chest, abdomen, pelvis with IV contrast	27 [2.7]	14 [2.8]	14 [2.8]
CTA chest, abdomen, pelvis with and/or without IV contrast	24 [2.4]	12 [2.4]	12 [2.4]
CT chest, abdomen, pelvis without IV contrast	22 [2.2]	11 [2.2]	11 [2.2]
CTA chest with and/or without IV contrast	18 [1.8]	9 [1.8]	9 [1.8]
CT chest without IV contrast, high resolution	17 [1.7]	9 [1.8]	9 [1.8]
CTA chest, abdomen with and/or without IV contrast	16 [1.6]	8 [1.6]	8 [1.6]
CTA coronary	16 [1.6]	8 [1.6]	8 [1.6]
CT initial lung cancer screening	12 [1.2]	6 [1.2]	6 [1.2]
CT cardiac scoring without IV contrast	12 [1.2]	6 [1.2]	6 [1.2]
CT subsequent lung cancer screening	11 [1.1]	6 [1.2]	6 [1.2]
CT thoracic spine without IV contrast	7 [0.7]	4 [0.8]	4 [0.8]
CTA chest vascular with and/or without IV contrast/gated	7 [0.7]	4 [0.8]	4 [0.8]
Age (median [IQR])	66 [55-75]	65 [54-75]	66 [54-76]
Females (n [%])	490 [48.9]	272 [54.0]	259 [51.4]
Race			
White	778 [77.6]	384 [76.2]	382 [75.8]
Black or African American	117 [11.7]	67 [13.3]	78 [15.5]
Asian	11 [1.1]	4 [0.8]	8 [1.6]
American Indian or Native American	1 [0.1]	1 [0.2]	2 [0.4]
Native Hawaiian or Other Pacific Islander	2 [0.2]	3 [0.6]	1 [0.2]
Other	67 [6.7]	31 [6.2]	28 [5.6]
Missing	20 [2.0]	12 [2.4]	5 [1.0]

587

588 **Table 2. Train, validation and test dataset annotation characteristics following**
 589 **preprocessing**

	Train	Validation	Test
Sentences in analysis (N)	214	103	91
Sentences with at least one annotation (n [%])	166 [77.6]	71 [68.9]	68 [74.7]
Sentences with annotations by measurement site (n [%])			
Annulus	25 [11.7]	13 [12.6]	10 [11.0]
Sinus of Valsalva	43 [20.1]	23 [22.3]	17 [18.7]
Sinotubular junction	27 [12.6]	13 [12.6]	9 [9.9]
Mid ascending	137 [64.0]	60 [58.3]	59 [64.8]
Ascending proximal to brachiocephalic	27 [12.6]	13 [12.6]	9 [9.9]
Top of Arch	32 [15.0]	17 [16.5]	12 [13.2]
Proximal descending	34 [15.9]	12 [11.7]	12 [13.2]
Mid Descending	40 [18.7]	20 [19.4]	14 [15.4]
Total annotations (N)	589	289	215
Annotation counts by measurement site (n [%])			
Annulus	49 [8.3]	26 [9.0]	20 [9.3]
Sinus of Valsalva	120 [20.4]	64 [22.1]	47 [21.9]
Sinotubular junction	44 [7.5]	22 [7.6]	13 [6.0]
Mid ascending	166 [28.2]	76 [26.3]	69 [32.1]
Ascending proximal to brachiocephalic	44 [7.5]	22 [7.6]	13 [6.0]
Top of Arch	49 [8.3]	29 [10.0]	17 [7.9]
Proximal descending	55 [9.3]	20 [6.9]	17 [7.9]
Mid Descending	62 [10.5]	30 [10.4]	19 [8.8]

590

591 **Table 3. Comparison of model performance on validation set.**

Model	Macro-averaged evaluation metric		
	Precision	Recall	F1
<i>Baseline</i>			
Zero-shot Llama 3.1	0.898	0.535	0.663
Few-shot Llama 3.1	0.792	0.903	0.838
<i>Fine-tuned BERT</i>			
Fine-tuned BERT	0.832	0.880	0.851
Fine-tuned PubMedBERT	0.901	0.940	0.919
Fine-tuned BERT-NER	0.868	0.942	0.902
Fine-tuned BERT + [Num]	0.855	0.895	0.870
Fine-tuned PubMedBERT + [Num]	0.940	0.956	0.945
Fine-tuned BERT-NER + [Num]	0.923	0.944	0.931
<i>Instruction-tuned Llama</i>			
Instruction-tuned Llama 2	0.826	0.853	0.839
Instruction-tuned Llama 3	0.973	0.994	0.982
Instruction-tuned Llama 3.1	0.993	0.992	0.992

592 Numbers in bold represent best performance per evaluation metric.

593 **Table 4. Fine-tuned Llama 3.1 performance on validation and test sets by aortic**
 594 **measurement site.**

Measurement site	Validation set			Test set		
	Precision	Recall	F1	Precision	Recall	F1
Annulus	1.000	1.000	1.000	1.000	1.000	1.000
Sinus of Valsalva	1.000	1.000	1.000	1.000	1.000	1.000
Sinotubular junction	1.000	1.000	1.000	0.857	1.000	0.923
Mid ascending	1.000	0.973	0.986	1.000	0.985	0.992
Ascending proximal to brachiocephalic	1.000	1.000	1.000	0.929	1.000	0.963
Top of arch	1.000	0.962	0.980	1.000	1.000	1.000
Proximal descending	0.944	1.000	0.971	0.938	0.938	0.938
Mid descending	1.000	1.000	1.000	0.900	1.000	0.947
Macro- averaged	0.993	0.993	0.992	0.953	0.990	0.970

595 Numbers in bold represent best performance per evaluation metric.

596

597 **Table 5. Complete dataset inference results by aortic measurement site.**

Measurement Site	Reports with measurements (n)	Percent of sentences analyzed with measurements (%)	Max Diameter (mm, median [IQR])
Annulus	7,072	9.49%	28 [25.7-30.3]
Sinus of Valsalva	12,620	16.94%	36 [32-40]
Sinotubular junction	7,621	10.23%	34 [30-37]
Mid ascending	44,472	59.71%	39 [36-42]
Ascending proximal to brachiocephalic	8,056	10.82%	35 [32-39]
Top of arch	9,865	13.24%	31 [28-34]
Proximal descending	8,886	11.93%	30 [27-34]
Mid descending	13,939	18.71%	29 [25-34]

598

599 **Figure Legends:**

600 Figure 1: Flowchart illustrating the selection and annotation process of radiology reports, divided
601 into training, validation, and testing sets.

602

603 Figure 2: Sample input and output. (A) Zero-shot prompt and (B) Few-shot prompt for Llama
604 models, illustrating the use of XML tags to delineate aortic measurement sites.

605

606 Figure 3: F1 Scores by Model. F1 scores for few-shot Llama 3.1, finetuned PubMedBERT +
607 [NUM], and instruction-tuned Llama 3.1.

608

609 Figure 4: Ablation Study. Macro F1 scores by number of training samples: 10, 25, 50, 100, 214.

Chest CT reports 2013 to 2023 (n=363,423)

Excluded reports (n=6,733)

- CT protocols with <2,000 reports (n=6,697)
- CT reports without narratives (n=36)

Chest CT reports included in cohort (n=356,690)

- Total sentences: 6,960,729
- Sentences for inference: 74,483

Training set (n=1,002)

- Total sentences: 19,844
- Sentences in analysis: 214

Validation set (n=504)

- Total sentences: 9,505
- Sentences in analysis: 103

Test set (n=504)

- Total sentences: 9,664
- Sentences in analysis: 91

A) Zero-shot prompt

Input:

Prompt: You will be provided with CT scan narratives detailing CT scan findings. Your task is to tag diameter measurements of the thoracic aorta in the provided CT narrative using the provided tags. Format each measurement as follows: <tag>measurement</tag>. The tags for each aortic anatomic location are listed below...

User: The sinus of Valsalva is dilated, measuring 4.2 cm.

Output:

Assistant: The sinus of Valsalva is dilated, measuring <SOV>4.2</SOV> cm.

B) Few-shot prompt

Input:

Prompt: You will be provided with CT scan narratives detailing CT scan findings. Your task is to tag diameter measurements of the thoracic aorta in the provided CT narrative using the provided tags. Format each measurement as follows: <tag>measurement</tag>. The tags for each aortic anatomic location are listed below...

User: Stably ectatic ascending aorta measuring up to 3.8 cm.

Assistant: Stably ectatic ascending aorta measuring up to <ASC>3.8</ASC> cm.

User: There is aneurysmal dilatation of the descending thoracic aorta measuring up to 3.5 cm.

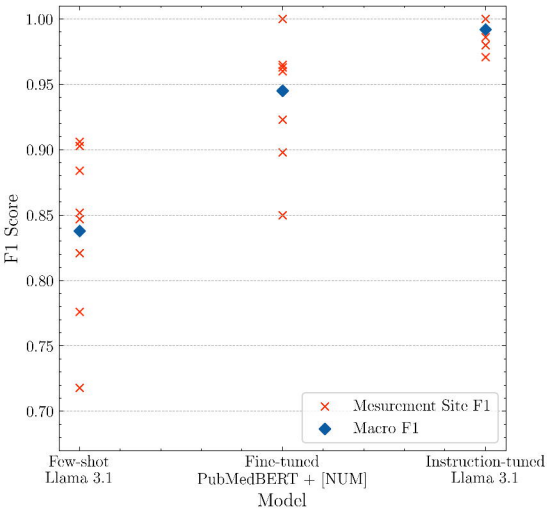
Assistant: There is aneurysmal dilatation of the descending thoracic aorta measuring up to <DSC>3.5</DSC> cm.

User: The sinus of Valsalva is dilated, measuring 4.2 cm.

Output:

Assistant: The sinus of Valsalva is dilated, measuring <SOV>4.2</SOV> cm.

F1 Scores by Model



Macro F1 Scores by Number of Training Samples

