

Multi-task weak supervision enables automated abnormality localization in whole-body FDG-PET/CT

Sabri Eyuboglu^{*1}, Geoffrey Angus^{*1}, Bhavik N. Patel², Anuj Pareek², Guido Davidzon², Jared Dunnmon,^{**1} Matthew P. Lungren ^{**2}

¹ Department of Computer Science, Stanford University, Stanford, CA 94305, USA

² Department of Radiology, Stanford University, Stanford, CA 94305, USA

* ** Equal contribution;

1 **Computational decision support systems could provide clinical value in whole-body FDG-**
2 **PET/CT workflows. However, limited availability of labeled data combined with large imag-**
3 **ing exams make it challenging to apply existing supervised machine learning systems. Lever-**
4 **aging recent advancements in natural language processing, we describe a weak supervision**
5 **framework that extracts imperfect, yet highly granular, regional abnormality labels from the**
6 **free text radiology reports. Our framework automatically labels each region in a custom on-**
7 **tology of anatomical regions, providing a structured profile of the pathologies in each imag-**
8 **ing exam. This framework is an attention-based, multi-task CNN architecture that can be**
9 **trained using these generated labels to localize abnormalities in whole-body scans. We exper-**
10 **imentally demonstrate that our multi-task representation is critical for strong performance**
11 **on rare abnormalities with limited training data and for accurately predicting mortality from**
12 **imaging data (AUROC 80%), which could assist palliative care teams .**

13 Introduction

14 The availability of large, labeled datasets has fueled recent progress in machine learning for med-
15 ical imaging. Breakthrough studies in chest radiograph diagnosis, skin and mammographic lesion
16 classification, and diabetic retinopathy detection have relied on hundreds of thousands of labeled
17 training examples [1–4]. However, for many diagnostic interpretation tasks, labeled datasets of this
18 size are not readily available either because (1) the task includes rare diagnoses for which training
19 examples are hard to find, and/or (2) the diagnoses are not recorded in a structured way within elec-
20 tronic medical records, requiring physicians to manually reinterpret exams or extract labels from
21 free-text reports. Even if datasets are manually annotated, common changes in the underlying data
22 distribution (e.g. scanner type, imaging protocol, post-processing techniques, patient population,
23 or clinical classification schema) could rapidly render the models they support obsolete. Further, it
24 is commonly the case that medical datasets are labeled using incomplete label ontologies, leading
25 to undesirable variation in performance on unlabeled subsets of the data (e.g. rare disease types),
26 a problem commonly referred to as hidden stratification [5].

27 These challenges are particularly apparent when working with whole-body fluorodeoxyglucose-
28 positron emission tomography/computed tomography (FDG-PET/CT), a medical imaging modal-
29 ity with a critical role in the staging and treatment response assessment of cancer. FDG-PET/CT
30 combines the high sensitivity of PET, which enables the detection of disease at picomolar concen-
31 trations, with the higher specificity of CT and thus is more accurate than either technology alone.
32 Indeed, FDG-PET/CT demonstrates equal sensitivity and better specificity even in the deep nodal
33 regions of the abdomen and the mediastinum [6, 7]. As a result of these improved capabilities and
34 the clinical value they provide, over the last five years there has been a 7% year-over-year increase
35 in the number of FDG-PET/CT examinations in the past decade within the United States [8].
36 Unfortunately, this increase in imaging demand has coincided with a reduction in both imaging
37 reimbursements and radiologist specialist availability [9].

38 While machine learning could provide clinical value within FDG-PET/CT workflows, this
39 potential remains largely untapped due to a lack of properly annotated data. While reading an FDG-
40 PET/CT scan, a radiologist will typically report the anatomical regions with abnormal metabolic
41 activity. This task, which we call *abnormality localization*, is clinically important in FDG-PET/CT
42 studies. However, manually labeling a dataset for abnormality localization is particularly painstaking-

43 ing because it requires either performing pixel-level annotations or resolving abnormalities into a
44 hierarchy of anatomical regions [10]. Existing approaches rely on clinical experts manually seg-
45 menting and annotating several thousand training examples [11]. In practice, radiologists report
46 abnormalities in many regions, even those in which abnormalities are exceedingly rare. Collecting
47 a sufficient number of training examples to train machine learning models that can detect abnor-
48 malities even in these low-prevalence regions can be extremely difficult.

49 To address these challenges, we present a machine learning framework for training abnor-
50 mality detection and localization models for large medical images without manual labeling or seg-
51 mentation. In particular, (1) we introduce a weak-supervision framework for rapidly and systemat-
52 ically generating abnormality localization labels from radiology reports describing FDG-PET/CT
53 findings, and (2) we combine multi-task learning with task-specific spatial attention mechanisms
54 to achieve high levels of automated detection and localization performance even on rare FDG-
55 PET/CT diagnoses for which we have little data. Our framework enables anatomically-resolved
56 abnormality detection in FDG-PET/CT scans, and models based on our framework could poten-
57 tially provide clinical value in four ways: (1) by automatically screening negatives and enabling
58 the prioritization of exams, (2) by enabling future diagnostic aids that could help radiologists lo-
59 calize abnormalities within the scan, flag rare abnormalities and quickly draft reports , and (3) by
60 yielding a learned representation of each scan that can enable more accurate mortality predictions,
61 which is important for palliative care planning.

62 We first curate a dataset comprised of 8,144 FDG-PET/CT scans, their accompanying re-
63 ports, and prospective normal-versus-abnormal clinical summary codes recorded by the radiologist
64 at the time of interpretation, similar to those reported by [12]. We leverage a custom ontology of
65 94 anatomical regions (e.g. left lung, liver, inguinal lymph nodes) to train an expressive language
66 model that extracts from each report the anatomical regions with abnormal FDG uptake. Using the
67 language model predictions as weak labels in the style of [13], we then train a 3D multi-task convo-
68 lutional neural network (CNN) to detect abnormal FDG uptake, leveraging the FDG-PET and CT
69 modalities jointly as input. The model localizes abnormalities by making binary abnormality pre-
70 dictions for a comprehensive set of 26 anatomical regions in which disease was present in at least
71 10% of the samples of our dataset. We evaluate our approach on a held-out dataset of 814 FDG-
72 PET/CT scans manually annotated by radiologists. While resource intensive, this fine-grained

73 manual labeling of the test set helps mitigate hidden stratification effects in our analysis [5].

74 We find that models trained with our multi-task weak supervision framework can detect
75 lymph node, lung and liver abnormalities – three of the pathologies with the highest prevalence
76 in our dataset – with median areas under the ROC curve of 87%, 85% and 92% respectively. In
77 a clinical screening setting, sorting a worklist by our models’ predictions would enable a team to
78 achieve sensitivities of 95%, 98%, and 96% respectively after only reading the first three-quarters
79 of exams.

80 We perform a comprehensive set of ablation studies to assess the utility of our multi-task
81 weak supervision framework. In one ablation, we find that using a language-model based labeler,
82 as opposed to a rule-based one like that of [14], improves the correctness of our weak labels. In
83 another, we show that the weak labels generated by our framework can be used to train models
84 that are statistically superior to fully supervised models trained on small, hand-labeled datasets.
85 Finally, we demonstrate that performance continues to increase as additional weakly-labeled train-
86 ing examples are added, a significant result given that the marginal cost of labeling more training
87 data is negligible in our framework while adding more hand-labeled data is costly.

88 We next evaluate the importance of training FDG-PET/CT scan classifiers in a massively
89 multi-task fashion (i.e. training them to detect abnormalities in many parts of the body simultane-
90 ously), an approach described by [15]. We find that multi-task models produce an effective learned
91 representation of a whole-body FDG-PET/CT scan that can facilitate training models for more
92 difficult abnormality localization tasks. For instance, we show that as the relative prevalence of ab-
93 normality in a given region decreases, the performance gains afforded by our multi-task model over
94 single-task models increase. Additionally, when compared to single-task models, our multi-task
95 model improves abnormality detection in regions with particularly low abnormality prevalence in
96 our dataset; we observe a 20-point improvement in the adrenal gland (AUROC 78%), and a 22-
97 point improvement in the pancreas (AUROC 78%), two regions with fewer than 3 abnormalities
98 per 100 exams. We also find that multi-task learning halves the computational cost of training
99 abnormality localization models.

100 We further explore whether multi-task abnormality localization pretraining (i.e. training a
101 model in a multi-task abnormality localization setting before fine-tuning on a different task) can
102 facilitate the development of models for other clinical tasks. Specifically, we find that our 3D-

103 convolutional neural network, when pretrained on multi-task abnormality localization, can be fine-
104 tuned to predict mortality within a 90-day period (AUROC 80%), a critical challenge for hospitals
105 looking to identify patients that could benefit from palliative care. While we do not argue that
106 our image-based mortality prediction models should be used without additional clinical covariates
107 drawn from EHR data, the predictive utility of these features suggests that they could be useful in
108 augmenting existing mortality prediction models [16, 17].

109 Finally, we evaluate the utility of the task-specific spatial attention mechanism that our mod-
110 els use to produce fine-grained representations for each anatomical region. We find that spatial
111 attention not only improves performance, but also produces attention distributions that align well
112 with human anatomy. These distributions are useful for model interpretability and serve as a sanity
113 check that the model is attending to the correct regions of the scan.

114 By bridging the feasibility gap associated with the lack of detailed hand-labeled training data
115 in volumetric imaging, our study lays a foundation for building clinically viable machine learning
116 models for FDG-PET/CT without burdensome labeled data requirements. While we focus on
117 FDG-PET/CT, we expect the methods and findings to have broad applicability across machine
118 learning on volumetric imaging modalities . Our multi-task weak supervision framework and
119 attention-based modeling approach for abnormality localization could enable applications in other
120 clinical settings where it is important to both diagnose and localize pathology in high-dimensional
121 images.

122 **Results**

123 **Combining weakly supervised multi-task learning with spatial attention mech-**
124 **anisms enables automated abnormality detection and localization in FDG-**
125 **PET/CT**

126 To circumvent the need for costly pixel-level annotations, we develop a weak supervision frame-
127 work for abnormality localization in large medical images. The locations of abnormalities in FDG-
128 PET/CT scans are buried in free-text radiology reports, so extracting labels for abnormality loca-
129 lization is nontrivial. Complex sentences like “there is a mildly FDG avid lingular pulmonary
130 nodule, which has slightly increased in size and metabolic activity,” indicate where abnormalities
131 appear in the scan. We propose using expressive pre-trained language models fine-tuned on a small
132 dataset of hand labeled sentences to extract abnormality localization labels from the reports. Our
133 labeling approach consists of three steps and is illustrated in Figure 1a-d: (1) we automatically tag
134 all mentions of anatomical regions (e.g. right lung, pancreas) in the report, (2) we use the language
135 model to predict whether or not each tagged mention in any given sentence is abnormal, and (3)
136 we reconcile these predictions into one large anatomical ontology where each region has a single
137 probability of abnormality. This approach allows us to extract anatomically resolved abnormality
138 labels with only a small amount of hand-labeled data.

139 Using these probabilistic labels, we train an attention-based, multi-task 3D convolutional
140 neural network (CNN) to detect abnormal metabolic activity in each of the 26 anatomical regions
141 with highest prevalence of abnormality in our dataset. Our model consists of a shared CNN encoder
142 (Figure 1e) that supports 26 binary classification task heads, one for each region. Each of these
143 task heads is equipped with a task-specific soft attention mechanism and a linear classifier (Figure
144 1f). The model is initially trained on all tasks jointly. Then, single-task instances of the model
145 are fine-tuned for each task, a technique commonly employed in Natural Language Processing
146 (NLP) [18]. A detailed description of our supervision and modeling procedures can be found in
147 the *Methods* section.

148 We evaluate our multi-task weak supervision framework through a series of experiments
149 meant to (1) evaluate its potential clinical utility and (2) quantify the contributions of our modeling
150 choices to the performance levels we observe.

We first evaluate the potential clinical utility of our models along two different axes: accurately localizing abnormalities and effectively screening abnormal exams. In Fig. 2a, we illustrate the capacity of weakly supervised multi-task FDG-PET/CT models to detect abnormal metabolic activity in the twenty-six anatomical regions with the highest prevalence of abnormality. These models are supervised with imperfect labels generated by our weak supervision framework — they use no manually annotated image data. In 22 of the 26 regions, our models achieves a mean AUROC > 75% over five random seeds. In 10 of the regions, including the lungs, liver, and thoracic lymph nodes, we achieve a mean AUROC > 85%.

In Figure 2b-d, we provide a deeper analysis of our model’s sensitivity (true positive rate) in three clinically important regions: the chest, the liver, and the inguinal lymph nodes. The purpose of this analysis is to explore how our weakly supervised abnormality localization framework could be used to develop triage tools for workflow prioritization in PET-CT screening. In order to evaluate potential model efficacy in this use case, we frame the test dataset as a worklist, which we sort using the predicted probability of abnormality output by our model. When we sort the predictions by decreasing predicted abnormality, if we only read the first three-quarters of exams, we can still achieve a sensitivity of 95%, 98%, and 96% in the chest, the liver, and the inguinal lymph nodes, respectively. In contrast, when reading exams in the order that they are received we expect to achieve a sensitivity of 75% if we interpret only the first three-quarters of the work-list. This analysis helps contextualize the potential utility of worklist triage in clinical interpretation. We can verify that the model is using information from the appropriate regions using 3D saliency maps, one of which is shown in Figure 2e.

Weak supervision reduces labeling costs for automated abnormality detection and localization

Weak supervision underpins our modeling framework, as it allows us to overcome a lack of labeled data and train multi-task abnormality localization models. Here, we evaluate our weak supervision approach and compare to alternatives. First, to assess the capacity of our labeling framework to accurately extract abnormality locations from radiology reports, we compare the labels generated by our framework to the hand-labels manually extracted from the reports by radiologists. Figure 3a compares the performance of our language model labeling framework to that of a regular expres-

180 sion baseline. The regular expression baseline is a labeling procedure that, for each exam, scans the
181 associated radiologist report for mentions of predefined anatomical regions. If there are no words
182 suggesting a negative finding (e.g. ‘physiologic’, ‘without’, ‘unremarkable’) in a sentence men-
183 tioning a region, the procedure assigns an abnormal label to that region and all of its parents in our
184 regional ontology. This formulation is similar to other, domain-specific tools [19, 20]. We observe
185 a mean AUROC of 90.0% when evaluating weak labels emitted by our labeling framework against
186 hand labels extracted from the same reports by radiologists (see Methods); further, we find a 28%
187 improvement in mean F1 score (52.2% to 66.6%) with respect to a regular expression baseline.
188 With a threshold of 0.5, our model achieves a mean sensitivity of 87.0% and a mean specificity of
189 82.2%, outperforming the regular expression baseline, which achieved a sensitivity of 83.3% and
190 a specificity of 66.9%. Note that we compare our labeler to the regular expression baseline using
191 F1-score (not AUROC) because the baseline outputs binary predictions rather than probabilistic
192 scores. The algorithm for the regular expression baseline can be found in Supplementary Note 4.

193 To characterize the trade-off between weak and full supervision, we train the same multi-
194 task abnormality localization model described above on a manually labeled training dataset of 400
195 exams. This labeling procedure, which requires annotating twenty-six anatomical regions for the
196 presence of hypermetabolic abnormality, took 16 board-certified radiologist hours. In comparison,
197 our weak supervision approach required only 12 non-expert hours to label a dataset of 6,530 ex-
198 ams. Despite the reduced labeling effort, our model trained on weak labels outperforms the model
199 trained on hand labels in every anatomical region (Figure 3b and Supplementary Table 3). On
200 average, weak supervision enables a 22 point increase in AUROC over the fully-supervised model
201 ($p = 0.0000$, paired permutation test). In the liver, for example, the weakly supervised model
202 detects abnormalities with an AUROC of 92.6%(91.8%, 92.4%), while the fully supervised model
203 detects these same abnormalities with an AUROC of 57.4% (52.4%, 61.2%). In nine anatomical
204 regions, including the neck, iliac lymph nodes, and abdomen, the fully-supervised model is no
205 better than random. Indeed, the increased dataset size enabled by weak supervision leads to con-
206 siderable performance gains over fully supervised approaches that use the same multi-task schema,
207 and equivalent or greater amounts of labeling resources.

208 We also evaluate a hybrid approach that combines weak supervision and manual labeling.
209 This model is trained on a mixed dataset consisting of 6,530 weakly labeled exams and 400 hand

210 labeled exams. On average, this hybrid approach enables an improvement of 1 AUROC point
211 over weak supervision alone – a statistically insignificant difference ($p = 0.2719$, two-sided paired
212 permutation test). While we do see performance gains in some regions (e.g. lungs), the boost in
213 performance is usually small.

214 Second, we compare our weakly supervised whole-body abnormality detection model with
215 its fully supervised counterpart. For each exam in our dataset, we have a summary code recorded by
216 the interpreting radiologist at the time of the exam that flags the presence of abnormality anywhere
217 in the FDG-PET/CT scan. We use these summary codes to train a fully-supervised, single-task
218 model for whole-body abnormality detection. We also generate analogous weak-labels using our
219 framework and train a weakly-supervised model.

220 In Fig. 3c, we explore how the relative performance of our weakly and fully supervised mod-
221 els vary as we increase the number of training examples. We find that our weakly supervised model
222 performs on par with the fully-supervised model across a wide range of training set sizes (100,
223 1,000, 2,000, 4,000, and 6,530). When trained on all 6,530 exams, our weakly supervised model
224 achieves a mean AUROC of 80.3%, which is statistically equivalent to the fully-supervised model
225 (78.3%, $p = 0.1933$, two-sided paired permutation test). These results suggest that our weak
226 supervision approach can dramatically reduce labeling cost while supporting high-performance
227 abnormality detection and localization models.

228 **Multi-task learning improves automated localization performance, reduces
229 computational cost, and facilitates mortality prediction**

230 Next, we show that multi-task learning enables strong performance on tasks for which we have very
231 few positive examples and significantly reduces the computational resources required to train ab-
232 normality localization models. We also show that multi-task abnormality localization pre-training
233 improves patient mortality prediction performance.

234 Our labeling framework generates labels for 94 anatomical regions, many of which are so
235 refined that we lack enough positive examples in our dataset to train performant single-task mod-
236 els. In Figure 4a, we show that multi-task learning mitigates some of the issues associated with
237 this large class imbalance and enables us to train performant scan models, even on some of the
238 regions for which we have the fewest positive examples. To demonstrate this, we chose four clin-

ically relevant regions beyond the main 26: celiac lymph nodes (68 positive training examples, 84th most out of 94 regions), kidney (182 positives, 64th most), adrenal gland (196 positives, 60th most), and pancreas (201 positives, 58th most), and compare model performance across various types of supervision in Figure 4b. On all four, our multi-task models substantially outperform single-task models pre-trained on the Kinetics activity detection dataset [21]. In terms of mean AUROC, we see a 56% improvement (50.2% to 78.0%, $p = 0.0003$ paired permutation test) in detecting abnormalities in the celiac lymph nodes, a 41% improvement (63.5% to 89.9%, $p = 0.0022$ paired permutation test) in the kidneys, a 34% improvement (58.0% to 77.9%, $p = 0.0004$ paired permutation test) in the adrenal glands, and a 40% improvement (55.9% to 78.2%, $p = 0.0000$ paired permutation test) in the pancreas. Additionally, we demonstrate that it is the multi-task formulation, not just in-domain pre-training, that enables these performance gains: compared to models pre-trained using weak labels on the task of binary, whole-body abnormality detection, our multi-task model achieves a 15.8% mean AUROC improvement (67.3% to 78.0%, $p = 0.0139$ paired permutation test) in detecting abnormalities in the celiac lymph nodes, a 25.4% improvement (71.7% to 89.9%, $p = 0.0180$ paired permutation test) in the kidneys, a 28.9% improvement (58.2% to 77.9%, $p = 0.0000$ paired permutation test) in the adrenal glands, and a 26.2% improvement (62.0% to 78.2%, $p = 0.0021$ paired permutation test) in the pancreas. We also compare to pre-training on summary codes and see similar improvements (Supplementary Table 7).

In Figure 4c, we demonstrate how weakly supervised multi-task pre-training reduces the computational resources required to fine-tune abnormality localization models until convergence. In particular, we compare the number of epochs of fine-tuning required for convergence after pre-training on multi-task abnormality localization versus pre-training on Kinetics. Our multi-task approach leads to reduced training times across all 26 anatomical regions. 68% of models pre-trained with our approach achieved their best performance within the first 4 epochs. On the other hand, only 7% of models pre-trained on Kinetics converged as quickly, while 60% required 8 or more epochs to train to convergence.

We also assess the capacity of our multi-task representation to generalize to other clinically relevant tasks: specifically, predicting mortality from FDG-PET/CT imaging data alone. We frame mortality prediction as a binary classification task, where the model predicts whether the patient’s date of death is within x days of the study given only the FDG-PET/CT scan as input. In our

269 experiments, we set x to be 45, 90, 180, and 365 days. Our mortality prediction model is first
270 pre-trained on multi-task abnormality detection in 26 anatomical regions . Then, the model is fine-
271 tuned to predict mortality within x days using the subset of FDG-PET/CT exams in our data for
272 which we have a recorded date of death. We fine-tune a separate model for each time point. Our
273 multi-task model predicts mortality within 90 days with an AUROC of 79.9% (76.2%, 83.8%). At
274 this threshold, multi-task pre-training enables a 15.3% improvement in mean AUROC over pre-
275 training on Kinetics (69.3% to 79.9%, $p = 0.0271$ paired permutation test). Similarly, compared
276 to pretraining on FDG-PET/CT summary codes in a single-task setting (see *Methods*), multi-task
277 pre-training enables a 15.3% improvement in predicting mortality within 90 days (69.3% to 79.9%,
278 $p = 0.0226$ paired permutation test). Figure 4d compares our model to the two baselines at each
279 of the four thresholds.

280 To understand why multi-task abnormality localization pretraining improves mortality pre-
281 diction performance, we perform a survival analysis that explores the relationship between ab-
282 normality localization predictions and longevity. We fit a Cox proportional hazard model on
283 our mortality data using PET/CT abnormality localization predictions as covariates (see *Meth-*
284 *ods* for details, Supplementary Table 18 for hazard ratios). Via the likelihood ratio test, we confirm
285 that the fit is significantly better than a null model with just a baseline hazard and no covariates
286 ($p = 1.1 \times 10^{-12}$, likelihood ratio test). The covariates corresponding to hypermetabolic abnor-
287 malities in the liver, spine, skeleton, and the hilar, inguinal, and thoracic lymph nodes exhibit statisti-
288 cally significant hazard ratios (Wald test). In Supplementary Figures 1c-d, we show Kaplan-Meier
289 curves stratified by the predictions for the liver and hilar lymph nodes.

290 The locations of hypermetabolic abnormalities appear to be predictive of mortality, but is
291 this information still useful when we have access to other more readily available covariates like
292 age or disease type? To explore this question, we fit a Cox proportional hazards model on our mor-
293 tality data using abnormality localization predictions, age, indication, and exam summary codes
294 as covariates (see *Methods* for details, Supplementary Table 16 for hazard ratios). We also fit a
295 nested Cox model that includes all covariates except the abnormality localization predictions (see
296 Supplementary Table 17 for hazard ratios). Via the likelihood ratio test, we show that the fit with
297 abnormality predictions is significantly better than the fit without ($p = 3.7 \times 10^{-11}$) [22]. We also
298 compare the out-of-sample predictive power of these Cox models by fitting on the validation set and

299 computing the concordance index on the test set. A Cox model fit just on age, summary code and
300 indication achieves a concordance index of 0.609, whereas a model that incorporates abnormality
301 location prediction achieves a concordance index of 0.720. Indeed, the locations of hypermetabolic
302 abnormalities seem to provide useful signal not present in other, more simply-attained covariates
303 such as age or indication.

304 To evaluate the importance of each covariate, we compute the difference between the covari-
305 ate's Wald χ^2 statistic and the covariate's degrees of freedom as described in [23] (Supplementary
306 Figure 1a). The most important covariate in the model, according to this metric, is the abnormality
307 prediction in the liver. Other important covariates include the patient's age, abnormality prediction
308 in the hilar lymph node, and the scan's indication.

309 **Spatial attention mechanisms improve automated localization performance 310 and facilitate model interpretation**

311 In order to train a model with the flexibility to make effective use of the granular labels produced by
312 our framework, we incorporate an attention mechanism into each of the task heads, which learns
313 to attend to the specific anatomical region of that head's task. We compare a multi-task FDG-
314 PET/CT model that uses a spatial, soft attention mechanism in each task head (see *Methods*) with
315 a near-identical multi-task FDG-PET/CT model that uses sum reduction in place of soft attention.
316 On 22 of the 26 tasks, the soft attention model outperforms the sum reduction model. In Figure
317 5a, we show overall performance (far left), as well as region specific performance. Models trained
318 with the attention mechanism see modest empirical performance improvement over a naive sum
319 reduction, with a 3 point increase in mean AUROC taken over all tasks and all seeds ($p = 0.0002$
320 paired permutation test). Additionally, our attention mechanism enables potentially useful new
321 ways of interpreting model predictions. As shown in Figure 5b-c, we can project the model's
322 attention distributions onto the original scan to highlight the voxels that most informed the model's
323 predictions. These can be used in conjunction with saliency maps, like those shown in Figure 5d-e,
324 to quickly locate abnormalities in the original scan, and confirm that each region-specific task head
325 is attending to the correct part of the image.

326 Refer to Supplementary Tables 2-9 to see a full numerical breakdown of the results described
327 in this section.

328 Discussion

329 FDG-PET/CT abnormality localization models could ultimately lower the burden on nuclear medicine
330 specialists and improve the quality of scan interpretation. However, the lack of properly annotated
331 data combined with the sheer size of each FDG-PET/CT exam makes it challenging to apply ex-
332 isting supervised machine learning systems to whole-body scans. There are on average 60 million
333 voxels in a whole-body FDG-PET/CT examination, only a tiny fraction of which represent an ab-
334 normality. When combined with relatively high signal-to-noise ratios in PET and the wide variety
335 of possible disease presentations, the size and complexity of each exam make it extremely chal-
336 lenging to train models that can localize abnormalities in whole-body scans. Instead of operating
337 on the full scan, existing approaches reduce the problem by only classifying a few slices at a time.
338 While this makes for an easier machine learning task, it demands that nuclear medicine specialists
339 painstakingly segment abnormalities on a large number of FDG-PET/CT slices [11, 24]. If the
340 scanner type, patient population or clinical classification schema changes, training data will need
341 to be relabeled.

342 From a technical perspective, we demonstrate in this work a set of repeatable techniques to
343 rapidly build weakly supervised, anatomically resolved classifiers for high-dimensional volumetric
344 imaging. The combination of multi-task learning, weak supervision, and attention mechanisms
345 that we propose enables us to identify and localize abnormalities in whole-body FDG-PET/CT
346 scans with almost no hand labeled data. This capability is critical for building machine learning
347 systems that are amenable to routine updates in the course of clinical practice. We show that
348 our weak supervision framework produces both whole-body abnormality detection models that are
349 statistically similar to their fully-supervised counterparts and abnormality localization models that
350 substantially outperform those trained using a small hand-labeled dataset. Moreover, we find that
351 multi-task learning enables significant performance gains for abnormality localization, particularly
352 on rare pathologies for which we have little training data. Finally, we find that task-specific spatial
353 attention mechanisms improve performance and enable new ways of interpreting model predic-
354 tions. We further justify the design of our framework through a series of ablation studies showing
355 that our language model labeler improves accuracy in radiologist report parsing with respect to a
356 rule-based system, that multi-task learning increases AUROC over comparable models and reduces
357 train time, and that our attention mechanisms outperform simple sum reduction.

358 From a clinical perspective, the shortage of radiologist expertise, coupled with a stark rise
359 in utilization of FDG-PET/CT imaging, suggest that automation for triage and reporting tasks
360 may become an important part of radiologist workflows in the near future. Automated pathology
361 localization could decrease time required for interpretation and highly specific models could op-
362 erationalize preliminary interpretations, draft reports, and improve workflow for radiologists and
363 the clinical referring services. Localization, from a modeling standpoint, is key in that it mitigates
364 hidden stratification, or performance variability in subgroups of data due to imprecise labels [5].
365 In our work, we pursue abnormality localization in part to provide explicit supervision over many
366 regions of the body, in contrast to our single-task models trained on summary codes describing the
367 whole body. Our approach ultimately leads to better performance on both whole-body abnormal-
368 ity detection and other downstream tasks. One of these tasks, mortality prediction, is of immense
369 importance in palliative care. The National Palliative Care registry estimates that less than half of
370 all patients that need palliative care actually receive it [16]. Statistical models that integrate the
371 abnormality localization predictions we propose in this work alongside traditional covariates could
372 eventually play a quantitative role in efficiently identifying patients who may benefit most from
373 palliative care , and who might have otherwise been missed by current care models. In the context
374 of our work, both identifying abnormalities in a triage setting and better flagging patients who are
375 likely to die within clinically important time frames (e.g. 90 days) may represent new capabilities
376 that could improve patient care.

377 Importantly, our approach drastically reduces the costs associated with training FDG-PET/CT
378 abnormality localization models. In order to train the text classification model that underpins our
379 labeling framework, only 1,279 sentences (< 0.3% of the sentences in our full dataset) were la-
380 beled. This took fewer than a dozen non-expert person-hours. In contrast, manually labeling the
381 800 exams in our test set took 32 radiologist-hours – manually labeling our full training set for
382 abnormality localization would require over 250 radiologist-hours, and this would only worsen for
383 larger datasets. Our weak supervision framework empowers clinical machine learning practition-
384 ers to rapidly update their labels to best fit clinical workflows and requirements [19, 25]. Training
385 multi-task models also front-loads the computational cost of training FDG-PET/CT models with-
386 out compromising on model performance. Fine-tuning our multi-task models takes a fraction of
387 the time that would be required to train an equivalent model from scratch. This further lowers the

388 barriers to consistent model refinement and retraining. Our framework provides reductions in both
389 labeling costs and training times, bringing us closer to machine learning systems that can handle
390 the variability and dynamism inherent in clinical practice. Further, our spatial attention mecha-
391 nism could serve as another tool for model interpretation and auditing. Analysis and visualization
392 of attention distributions, such as those shown in Figure 5b-c, can be paired with analysis of two-
393 dimensional and three-dimensional saliency maps (shown in Figure 5d-g and Fig 2e, respectively)
394 in order to assess model behavior.

395 While our results are promising in several respects, it is important to emphasize that deliv-
396 ering production-ready models in their final clinical form is beyond the scope of this study, and
397 that there are several additional steps that must be taken before deploying such models in clinical
398 practice.

399 First, there are several approaches through which one could improve performance. As we
400 show in Figure 3c, increasing the size of the training dataset is likely to provide additional per-
401 formance improvements (other studies in clinical weak supervision have also documented this
402 trend [26]). While our 6,530-sample dataset is large in FDG-PET/CT terms, it is still relatively
403 small compared to most other medical imaging datasets in machine learning, it is drawn from a
404 single center, and scans were performed on devices from a single vendor using our standard insti-
405 tutional technique for image acquisition; each of these realities limits the generalizability of the
406 models we present here. Since the marginal cost of adding training examples using our framework
407 is negligible, a deployed model would ideally be trained on a more diverse dataset an order of mag-
408 nitude larger than ours by drawing on unlabeled scans from multiple academic centers or hospitals.
409 This is particularly important given that performance on some regions (e.g. spine, skeleton and
410 pelvic skeleton) is low: increasing the number of weakly labeled training examples would be sim-
411 ple and effective way to boost performance on these regions. One could also consider deploying
412 an ensemble model comprised of multiple models, each trained from a different weight initializa-
413 tion, in order to account for variation inherent in the stochastic optimization procedure [27]. For
414 example, we ensemble the five seeds of our multi-task model to show performance gains of up to
415 3 AUROC points on each task, a result that can be found in Supplementary Table 19.

416 Second, although training abnormality localization models with our framework does not re-
417 quire manually annotating large amounts of training data, like with all deployed clinical machine

418 learning frameworks, routine validation on a gold-standard, prospective dataset is still critical for
419 detecting degradation in model performance caused by shifting patient populations or changing
420 imaging protocols. Weak labels would not suffice for these evaluations, which means that deploy-
421 ing models in practice will require some manual labeling of test data. This critical step will likely
422 require significant effort: hand-labeling our test set of 800 exams took 32 radiologist-hours. That
423 being said, if an evaluation indicates that model performance has degraded, our framework makes
424 it easier to integrate new training data to address the issue. Note that these evaluations would be
425 required for any automated interpretation system that would be deployed in clinical practice, and
426 thus the requirement is not specific to the models we propose

427 Third, in this study, we show how our scan model can be fine-tuned to predict mortality from
428 FDG-PET/CT scans. We chose to train mortality prediction models on imaging data alone not
429 because doing so would necessarily make sense in clinical practice, but because it demonstrates
430 the utility of this data source in mortality prediction. In practice, palliative care decisions should be
431 based on a number of different factors, but our work shows for the first time that automated analysis
432 of PET/CT data could be useful in this setting. Patients were not prospectively followed for clinical
433 outcomes in this retrospective study design. Dates of death were acquired retrospectively and
434 patients without a date of death were censored in our analyses. Future studies should evaluate
435 mortality prediction models that incorporate learned representations of FDG-PET/CT data in a
436 controlled clinical trial population with specific follow-up protocols designed for the purposes of
437 this research question.

438 In conclusion, we have developed a framework for rapidly training machine learning models
439 to detect and localize abnormalities in whole-body FDG-PET/CT scans with sparse hand labeled
440 data. We provide experiments justifying each of modeling choices behind our framework – from
441 the decision to weakly supervise models using a large, unlabeled dataset to our use of multi-task
442 learning to the inclusion of spatial attention mechanisms. We also show how we can make accurate
443 mortality predictions from PET-CT data alone and explore how the locations of abnormalities relate
444 to a patient’s survival function. Importantly, the techniques we’ve outlined can be adapted to train
445 models for other imaging modalities where labeled data is scarce. We hope that our work can
446 serve as a stepping stone for future research in automated analysis of FDG-PET/CT and other
447 volumetric imaging modalities, both through its technical contributions and through its potential

⁴⁴⁸ to create direct clinical impact.

449 **Methods**

450 **Dataset.** We use a dataset of 8,144 FDG-PET/CT exams from 4,691 patients (see Figure 1a).
451 The exams were administered at Stanford hospital between 2003 and 2010. Each exam includes
452 PET and CT axially-oriented image sequences that span from the upper thigh to base of skull. Each
453 FDG-PET/CT exam also includes a summary code and a free-text report written by the interpreting
454 radiologist at the time of the examination. The summary code describes overall patient status and
455 takes on a value of 1, 2, 4, or 9, where 1 indicates no evidence of abnormality, 2 indicates that there
456 are abnormalities but they were known prior to the exam, 4 indicates that there is at least one new
457 abnormality not previously known, and 9 indicates the presence of an abnormality that requires
458 emergent attention. Details about the dataset and the data preprocessing procedure can be found
459 in Supplementary Notes 1 and 2, respectively. Details regarding exam and patient metadata can be
460 found in Supplementary Tables 12-15.

461 FDG-PET/CT is a combination of two different imaging modalities. The FDG-PET scan,
462 which captures metabolic activity within the body, is composed of a sequence of 128×128 pixel
463 images. The CT scan, which captures anatomical structure, consists of 512×512 pixel images.
464 Figure 1a shows an example FDG-PET/CT scan, visualized such that the image frames are stacked
465 vertically into a three-dimensional volume. The length of the image sequences used in this study
466 range from 69 to 307 ($\mu = 212.39$, $\sigma = 23.77$) images.

467 An FDG-PET/CT report is an unstructured text document describing the clinical context of
468 the patient’s exam and the findings of the interpreting radiologist. The report typically consists of 4
469 sections: (1) the “clinical history” section, which specifies the indication of the study as well as any
470 prior disease and/or treatment, (2) the “procedure” section, which describes the techniques used in
471 administering the exam, (3) the “findings” section, which details clinically significant observations
472 made in each region of the body, and (3) the “impression” section, which serves as a summary of
473 the most significant observations made in the report [28].

474 This study was approved by the Stanford Institutional Review Board. The Stanford Research
475 Repository (STARR, formerly STRIDE) is Stanford Medicine’s approved resource for working
476 with clinical data for research purposes [29]. FDG-PET/CT imaging data and radiology reports for
477 patients in the cohort were acquired with STARR and deidentified. Additionally, dates of death for
478 a subset of the cohort were retrieved using STARR, which integrates medical records at Stanford

479 with the Social Security Death Index (SSDI). Of the 4,691 patients in our dataset, 867 (18%) had a
480 date of death recorded in STARR. Patients were not prospectively followed for clinical outcomes
481 in this retrospective study design. Patients without a date of death recorded in STARR were right
482 censored at the time of their last PET/CT exam in our data.

483 The exams were split into training, validation, and test sets. The validation and test sets were
484 sampled at random with uniform probability to capture the class distribution expected in a clinical
485 setting. The exams were split by patient, ensuring that there is no patient overlap between the
486 training, validation, and test sets. The validation set was used to evaluate prototypes during model
487 development, to tune hyperparameters, and for early stopping during model training. All reported
488 metrics were computed on the test set, unless otherwise specified. For experiments comparing
489 supervision strategies (see Figure 3b and Supplementary Table 3), we use a hand labeled test set of
490 423 exams from 235 patients, a validation set of 800 exams from 469 patients, and a training set
491 of 6921 exams from 3987 patients, of which 391 exams from 235 patients are hand labeled. For
492 all other experiments, we use a validation set of 800 exams from 469 patients, a hand-labeled test
493 set of 814 exams from 470 patients, and a training set of 6,530 exams from 3,752 patients.

494 **Weak supervision.** Recent work has established the effectiveness of weak supervision: a machine
495 learning paradigm wherein supervised machine learning models are trained with imperfect, yet
496 cheaply generated training labels. With weak supervision we can reduce or even eliminate the
497 need for costly hand labeled data [13, 30, 31].

498 In this study, we work from within the weak supervision paradigm and implement a labeling
499 framework that ingests radiology reports — rich, yet unstructured bodies of text describing
500 the findings of the interpreting radiologist — and outputs probabilistic labels for each anatomical
501 region in the FDG-PET/CT scan (e.g. lungs, liver). Our labeling framework leverages (1) a custom
502 ontology of anatomical regions relevant to FDG-PET/CT (See Figure 1d), (2) programmatic
503 functions that tag anatomical regions in the reports (See Figure 1b), and (3) a text classification
504 model, which we call the *report model*, that determines whether the tagged regions are described
505 as metabolically abnormal in the report (See Figure 1c). We then use the labels generated from
506 the reports to train a large convolutional neural network, which we call the *scan model*, to predict
507 abnormalities in the full FDG-PET/CT scans. Although we depend on reports to generate training
508 labels, at test time the trained scan model can detect abnormalities in FDG-PET/CT scans without

509 an accompanying report.

510 **Regional ontology.** We construct an ontology of 94 anatomical regions relevant to FDG-PET/CT.
511 Anatomical regions include coarse, high-level regions like “chest”, “abdomen”, and “pelvis” as
512 well as fine-grained regions like “left inguinal lymph node” and “upper lobe of the right lung”.
513 Our ontology is a directed-acyclic graph where nodes represent regions and edges connect regions
514 to sub-regions. For example, edges lead from “lungs” to “left lung” and from “thoracic lymph
515 node” to “hilar lymph node”. To determine which anatomical regions to include in the ontology,
516 we perform a systematic analysis of region mentions in FDG-PET/CT reports. Specifically, we
517 compute k -gram counts ($k = 1, 2, 3$) across all reports in our training dataset. Then, if a k -gram
518 refers to an anatomical region and appears at least 35 times in our dataset, we add the anatomical
519 region to the ontology. The edges connecting regions to sub-regions were added in consultation
520 with nuclear medicine and radiology specialists (G.D. and M.L.). For a visualization of the full
521 ontology, see Supplementary Fig. 3. Note, when training a multi-task scan model (see **Scan model**
522 **training.** below) we use only the 26 regions for which there is at least one positive for every nine
523 negative examples.

524 **Tagging functions.** Each anatomical region in the ontology is accompanied by a set of tagging
525 functions that search a report for mentions of the region (see Figure 1b). A tagging function
526 could be a simple regular expression query, or a complex set of rules that capture more elaborate
527 descriptions of the region. Given a report and a region from the ontology, tagging functions allow
528 us to extract a list of sentences that mention that region. In our experiments, we run the tagging
529 functions on the findings and impression sections of each report.

530 Note that the mention of an anatomical region does not necessarily imply that there is an
531 abnormality in that region. In fact, the majority of mentions in our dataset occur in the context of
532 a neutral finding. Existing approaches depend on words of negation and uncertainty (e.g. “not”,
533 “no”, “unlikely”) to classify mentions as neutral or negative. However, compared to other modal-
534 ities, the language in whole-body FDG-PET/CT reports is quite nuanced, so these approaches
535 produce a large number of false positives (a baseline labeler that uses this approach achieves a
536 sensitivity of 83.3%, but a specificity of only 66.9%). For example, it requires a nuanced under-
537 standing of FDG-PET/CT language to know that the sentence “intense physiologic uptake in the
538 cerebral cortex” is describing a neutral finding.

539 **Report model.** Rather than rely on hard-coded rules to classify mentions, we leverage an ex-
540 pressive language model that can capture the complexity of FDG-PET/CT reports. Our language
541 model, which we call in this section the *report model*, is trained to predict whether an anatom-
542 ical region is mentioned in the context of an abnormal finding (see Figure 1c). With a trained
543 report model, we can assign a probability of abnormality to each mention returned by the tagging
544 functions.

545 Recent work in natural language processing has shown that pre-training large-scale language
546 models on lengthy corpora of unlabeled text can enable strong performance on down-stream tasks
547 with relatively little labeled training data [32–35]. One such language model known as BERT
548 (bidirectional encoder representations from transformers) learns word representations by condi-
549 tioning on context to both the left and the right of the word [32]. Our report model is based on the
550 BERT language model.

551 The model accepts as input one or more sentences of natural language. As a preprocessing
552 step, we split up the sentences into a list of *wordpieces* [36]. A wordpiece is a sequence of a few
553 characters that make up part or all of a word. The model treats each wordpiece as an indivisible
554 unit. Wordpieces allow us to operate on rare, out-of-vocabulary words. For example, the sentence:

555 “FDG uptake in subcarinal and contralateral mediastinal lymph nodes.”

556 might be represented with the following wordpieces:

557 [FDG uptake in sub -carinal and contra -lateral media -stinal lymph nodes .]

558 Notice that rare, complex words like “subcarinal” and “contralateral” are split into wordpieces
559 while common words like “lymph” and “uptake” are kept intact. Wordpieces are particularly useful
560 for FDG-PET/CT reports where prefixes like “sub-” and “hyper-” often play important semantic
561 roles.

562 BERT is typically used with a vocabulary of 30,000 wordpiece tokens optimized for generic
563 text (BookCorpus and English Wikipedia) [37]. Because FDG-PET/CT reports are filled with
564 highly specialized vocabulary, using BERT’s out-of-the-box wordpiece tokens will mean splitting
565 up many important, domain-specific words. To account for the specialized FDG-PET/CT text,
566 we use a greedy algorithm to find the set of 3,000 wordpieces that minimize the number of tokens
567 required to reconstruct the reports in our training dataset [36,38]. Of these 3,000, 1,675 are already

568 in the BERT vocab. We add the remaining 1,343 wordpiece tokens to the original vocabulary
 569 by replacing the ”unusedX” and non-ASCII tokens provided in the BERT implementation. This
 570 allows us to leverage the initial BERT weights while also introducing domain specific tokens. Note,
 571 it is important to exclude any testing data when generating word pieces to ensure that evaluations
 572 on the test set provide a fair estimate of generalization error. We use Google’s SentencePiece
 573 library to generate these wordpiece tokens.

574 We can formalize our report model as a function \mathcal{G} (parameterized by $\theta_{\mathcal{G}}$) that maps a se-
 575 quence of wordpiece tokens $(x_1, x_2 \dots x_n)$ to an equal-length sequence of hidden representations
 576 $(\mathbf{z}_1, \mathbf{z}_2 \dots \mathbf{z}_n)$,

$$\mathcal{G}(x_1, x_2 \dots x_n; \theta_{\mathcal{G}}) = (\mathbf{z}_1, \mathbf{z}_2 \dots \mathbf{z}_n). \quad (1)$$

577 The report model has a Transformer architecture, which uses self-attention to draw relations be-
 578 tween tokens in the input. Because we use an implementation identical to the original, we refer the
 579 reader to the original Transformer manuscript for details on the architecture [39].

580 To predict whether or not a token x_i occurs in the context of an abnormal finding, we pass
 581 its hidden representation \mathbf{z}_i through a single fully-connected layer with a sigmoid activation. This
 582 gives us a probability

$$P(x_i \text{ is abnormal} | (x_1, x_2 \dots x_n)) = \sigma(\mathbf{w}^T \mathbf{z}_i + b) \quad (2)$$

583 where \mathbf{w} is a weight vector and b is a bias term.

584 **Report model training.** To train our report model, we need a dataset of sentences that mention
 585 anatomical regions and are labeled as negative, neutral or positive. Formally, we can train the
 586 model with labeled examples $((x_1, x_2 \dots x_n), (y_1, y_2 \dots y_n))$ where $y_i \in \{0, 1, -\}$ takes on a value of
 587 “0” if token x_i forms part of a mention that is neutral or negative, “1” if token x_i forms part of a
 588 mention that is positive, and “−” if x_i is not part of any mention. For example, the sentence

589 $\mathbf{x} = [\text{Abnormal FDG uptake in the left lung}]$

590 would be labeled

591 $\mathbf{y} = [- \ - \ - \ - \ 1 \ 1]$

592 because only the words “left” and “lung” are part of the anatomical mention.

During training, we use a loss function that ignores wordpiece tokens that are not part of any anatomical mention (i.e. labeled with “–”). Formally, our optimization objective is

$$\theta_{\mathcal{G}}^* = \operatorname{argmin}_{\theta_{\mathcal{G}}} \sum_{k=1}^M \frac{1}{n} \sum_{i=1}^n \mathbf{1}[y_i^{(k)} \neq -] \mathcal{L}(y_i^{(k)}, \hat{y}_i^{(k)}) \quad (3)$$

593 where \mathcal{L} is cross-entropy loss and $\hat{y}_i = \sigma(\mathbf{w}^T \mathbf{z}_i^{(k)} + b)$, m is the number of examples in the training
 594 set, and a superscript (k) is a reference to the k^{th} training example.

595 To efficiently generate a mini-dataset of labeled mentions, we implement a lightweight, la-
 596 beling GUI. Using it, two non-expert annotators (G.A. and S.E.) were able to label the mentions in
 597 a sample of 1,279 sentences (< 0.3% of the sentences in our full dataset) over the course of twelve
 598 person-hours. The GUI integrates with Jupyter Notebooks and could be used to easily label data
 599 for a different task.

600 Prior to training, we pre-train our report model with masked language modeling (MLM)
 601 and next-sentence prediction (NSP), the two pre-training tasks proposed in the original BERT
 602 manuscript [32]. In MLM, we randomly mask wordpiece tokens in the input and task the model
 603 with recovering the masked token using just the context around the mask. To do so, we pass each
 604 hidden representations \mathbf{z}_i through a MLM task head. In NSP, we feed the model two sentences
 605 and task it with predicting whether or not the first sentence preceded the second in the corpus. We
 606 take the weights from $\text{BERT}_{\text{base}}$ pre-trained on BookCorpus and English Wikipedia [40]. Then, we
 607 perform domain-specific pre-training on our training dataset of FDG-PET/CT reports.

608 We train all of our models with an Adam optimizer and an initial learning rate of $\alpha = 0.0001$
 609 [41]. We anneal the learning rate by half every 20 epochs. We do not employ any regularization.
 610 We train the report model with a batch size of 16. In each epoch, we sample 100 mentions with
 611 replacement from the training set. We train the model for 85 epochs. We perform early stopping
 612 with AUROC on a validation set of 61 sentences. Visualizations of the BERT attention heads can
 613 be found in Supplementary Fig. 2.

614 **Generating Labels.** Using a trained report model, we can generate labels for our training and
 615 validation datasets. For each report, we run the tagging functions from the region ontology. This
 616 gives us a list of sentences that mention anatomical regions of interest. We then tokenize each
 617 of those sentences into a sequence of wordpieces (x_1, x_2, \dots, x_n) and feed them through our report

model. This yields a sequence of predictions $(\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n)$ where \hat{y}_i represents the probability that token x_i occurs in the context of an abnormal finding (i.e. $P(x_i \text{ is abnormal} | (x_1, x_2, \dots, x_n))$), see Eq. 2). Because mentions can span multiple tokens, we reduce multiple predictions to a single probability for the whole mention by taking the mean probability across the tokens in the mention.

Note that some anatomical regions may be mentioned more than once in a report and others not at all. To reconcile the predictions made by the report model into a single probability for each anatomical region we leverage the regional ontology (see Figure 1d). Specifically, by propagating probabilities up the regional ontology we collect for each region t a list of probabilities $(p_1, p_2, \dots, p_{n_t})$ for all mentions of t and its children. We then compute the probability, assuming independence, that at least one of those mentions occurred in the context of an abnormal finding,

$$P(t \text{ is abnormal}) = 1 - \prod_{i=1}^{n_t} (1 - p_i) \quad (4)$$

The label propagation process is illustrated in Figure 1d.

After label propagation we are left with a single probability $\bar{y}_t = P(t \text{ is abnormal})$ for each anatomical region t in the ontology. Below, we show how we can use these predictions as probabilistic labels and train a model to detect abnormalities in FDG-PET/CT scans. The class balance for each of the anatomical regions in the training, validation, and test sets can be found in Supplementary Table 1.

Multi-task learning. A considerable body of research has focused on using multi-task learning to reduce generalization error in computer vision and natural language machine learning models. Multi-task learning has proven particularly useful in settings where labeled training examples are scarce [42, 43]. In this work, we leverage the noisy, probabilistic labels generated by the report model to train a *scan model* that maps a whole-body FDG-PET/CT scan to a probability of abnormality in one or more regions of our ontology. We use a simple multi-task architecture comprised of a shared encoder module \mathcal{F} (parameterized by $\theta_{\mathcal{F}}$), and T region-specific decoders $\{\mathcal{D}_1, \dots, \mathcal{D}_T\}$ (parameterized by $\{\theta_{\mathcal{D}_1}, \dots, \theta_{\mathcal{D}_T}\}$) [44]. For each region t , the model outputs the probability that there is a metabolic abnormality in that region. We can formalize the prediction

for some input scan $\mathbf{X} \in \mathbb{R}^{2 \times 224 \times 224 \times l}$ and region t as

$$P(t \text{ is abnormal} | \mathbf{X}) = \mathcal{D}_t(\mathcal{F}(\mathbf{X}; \theta_{\mathcal{F}}); \theta_{\mathcal{D}_t}). \quad (5)$$

To train the model, we perform the following optimization

$$\theta^* = \operatorname{argmin}_{\theta} \sum_{k=1}^M \frac{1}{T} \sum_{t=1}^T \mathcal{L}(\bar{y}_t^{(k)}, \hat{y}_t^{(k)}) \quad (6)$$

where M is the number of samples in our dataset, \mathcal{L} is cross-entropy loss, \bar{y}_t is the probability of abnormality output by the report model (See Eq. 4), and $\hat{y}_t^{(k)}$ is the probability of abnormality output by the scan model (See Eq. 9).

Scan model. For our shared encoder module \mathcal{F} , we use an Inflated Inception V1 3D CNN (I3D) pre-trained on the Kinetics dataset with optical flow [45]. We remove the final classification layer so that the encoder outputs a 3-dimensional encoding of the input scan. The encoding consists of $d = 1024$ channels, each of shape $7 \times 7 \times \lceil \frac{l}{6} \rceil$, where l is the number of slices in the original exam. A schematic illustration of the encoder is provided in Figure 1e. Formally, the encoder module $\mathcal{F}(\mathbf{X}; \theta_{\mathcal{F}})$ outputs a tensor $\mathbf{A} \in \mathbb{R}^{d \times 7 \times 7 \times \lceil \frac{l}{6} \rceil}$. The encoding \mathbf{A} can be viewed as a volume where each voxel is a vector $\mathbf{a}_{i,j,k} \in \mathbb{R}^d$. We visualize this encoding on the right hand side of Figure 1f.

Each region-specific decoder \mathcal{D}_t is composed of a soft-attention mechanism and a single linear classification layer. Intuitively, the attention mechanism allows each task head to “focus” on specific regions of the scan (see Figure 1f). To perform soft-attention, we compute the dot product between each voxel $\mathbf{a}_{i,j,k} \in \mathbb{R}^d$ in the encoding and a learned weight vector $\mathbf{w} \in \mathbb{R}^d$

$$s_{i,j,k} = \mathbf{w}^T \mathbf{a}_{i,j,k} \quad (7)$$

yielding a score $s_{i,j,k}$ for each voxel. We apply softmax across the scores $\alpha = \text{Softmax}(s_{i,j,k})$, and use them to compute a linear combination of all the voxels in the scan encoding

$$\mathbf{a} = \sum_{i,j,k} \alpha_{i,j,k} \mathbf{a}_{i,j,k}. \quad (8)$$

Intuitively, the larger the $\alpha_{i,j,k}$, the more attention is paid to the voxel at coordinates (i, j, k) . The

linear combination $\mathbf{a} \in \mathbb{R}^d$ is then fed to a final linear classification layer with a sigmoid activation. Altogether, each region-specific decoder outputs a single probability of abnormality

$$P(t \text{ is abnormal} | \mathbf{A}) = \mathcal{D}_t(\mathbf{A}; \theta_{\mathcal{D}_t}). \quad (9)$$

650 **Scan model training.** We train the scan model with an Adam optimizer and an initial learning rate
651 of $\alpha = 0.0001$ [41]. We anneal the learning rate by half every 16 epochs. We do not employ any
652 regularization aside from data augmentation. Due to GPU memory constraints, we train the scan
653 model with batch size of only 2. In each epoch, we sample 2,000 exams with replacement from
654 the training set.

655 Although our labeling framework generates labels for 94 anatomical regions, in practice we
656 find it challenging to train a multi-task model with 94 different tasks. Instead, we train a multi-
657 task scan model on the 26 regions for which there is a prevalence of at least one positive example
658 for every nine negative examples (i.e. fraction of positive examples $\geq 10\%$). We estimate the
659 prevalence using the entire dataset of 8,144 exams and the weak labels generated by our labeling
660 framework. Note the remaining 68 regions in our ontology are all sub-regions of one of these main
661 26. The multi-task model is trained for 15 epochs. We then perform single-task fine-tuning for
662 the 26 multi-task regions as well as four “rare” regions with class balance less than 10%. Through
663 this fine-tuning step, we are able to (a) improve performance in the main 26 regions and (b) train
664 performant models on regions beyond the 26 most prevalent (see Results). We perform fine-tuning
665 for 5 epochs. After each epoch, we evaluate on the validation set using weak labels. Note, that
666 fine-tuning typically converges after 1 – 2 epochs, so training for 5 epochs is not strictly necessary
667 in practice (see Figure 4d). After fine-tuning, we use the model weights from the epoch with the
668 highest validation AUROC. Although our labeling framework generates labels for 94 anatomical
669 regions, in practice we find it challenging to train a multi-task model with 94 different tasks.
670 Instead, we train a multi-task scan model on the 26 regions for which there is a prevalence of at
671 least one positive example for every nine negative examples (i.e. fraction of positive examples
672 $\geq 10\%$). We estimate the prevalence using the entire dataset of 8,144 exams and the weak labels
673 generated by our labeling framework. Note the remaining 68 regions in our ontology are all sub-
674 regions of one of these main 26. The multi-task model is trained for 15 epochs. We then perform
675 single-task fine-tuning for the 26 multi-task regions as well as four “rare” regions with class balance

676 less than 3%. Through this fine-tuning step, we are able to (a) improve performance in the main
677 26 regions and (b) train performant models on regions beyond the 26 most prevalent (see Results).
678 We perform fine-tuning for 5 epochs. After each epoch, we evaluate on the validation set using
679 weak labels. Note, that fine-tuning typically converges after 1 – 2 epochs, so training for 5 epochs
680 is not strictly necessary in practice (see Figure 4d). After fine-tuning, we use the model weights
681 from the epoch with the highest validation AUROC.

682 During training, we apply two basic data-augmentation transforms to each scan. We ran-
683 domly crop the image sequence to a 200×200 pixel region, then resize to 224×224 pixels,
684 downsampling the 512×512 CT images and upsampling the 128×128 PET images using bilin-
685 ear interpolation. We additionally jitter the brightness of the image sequence by adjusting bright-
686 ness throughout the sequence by a factor $\gamma \sim \text{Uniform}(0.0, 0.25)$. Further details regarding the
687 training setup can be found in Supplementary Note 3, and ablations that demonstrate the effect of
688 upsampling the FDG-PET images and excluding the CT modality during training can be found in
689 Supplementary Tables 10 and 11, respectively.

690 **Model evaluation.** Our test set of 800 exams was hand-labeled by four board-certified radiolo-
691 gists (G.D., B.P., A.P., M.L.) with experience ranging from 4 to 15 years. The test set was split
692 among the radiologists such that each radiologist labeled 200 exams. The radiologists labeled
693 each exam based only on the contents of its associated report (i.e. the actual FDG-PET/CT exam
694 was not reinterpreted). For each of 30 anatomical regions (twenty-six main regions plus four rare
695 regions we examine in detail in this work), the radiologists assigned a binary (i.e. $\{0, 1\}$) abnor-
696 mality label. A label of 1 was assigned to a given region when there was an explicit mention of
697 abnormal FDG uptake in that region in the *findings* or *impression* of the report. For example, if
698 the impression contained the sentence “There is intense radiotracer uptake in approximately 13
699 $\times 12$ mm left cervical level-ii lymph node (slice -220, SUV 6.5), the label would be 1. A label
700 of 0 was assigned when the impression or findings contained either an explicit mention of no ab-
701 normal FDG uptake, as in the sentence “no pathologically enlarged or hypermetabolic cervical or
702 supraclavicular lymphadenopathy on the current study, or no mention at all of FDG activity in the
703 region. Exams with ambiguous or uncertain wording were flagged and reviewed by a sub-specialist
704 PET/CT radiologist with 15 years of experience (G.D.). The labeling radiologists were trained un-
705 der the supervision of this sub-specialist PET/CT radiologist and were given a document providing

706 labeling guidelines.

707 Unless otherwise specified, all reported metrics in results were computed using data from
708 the test set and the hand labels described above. No hyperparameter tuning or model development
709 was performed with the test set.

710 Every model in our analysis was trained with five different random initializations (i.e. ran-
711 dom seeds). We evaluate each of the five trained models on the test set and report the mean of the
712 five resulting scores (e.g. AUROC). We account for uncertainty in this estimate of the mean score
713 with 95% confidence intervals computed via bootstrapping over the sample of five scores. When
714 comparing our models to baselines (e.g. weakly-supervised vs. fully-supervised), we test the null
715 hypothesis that the sample of scores of our model comes from a distribution with the same or lower
716 mean than the sample of scores of the other. This is done via a one-sided paired permutation test
717 with $n = 10,000$ iterations. Note when testing whether two models are statistically equivalent, we
718 use a two-sided paired permutation test.

719 We evaluate the performance of our label generation framework using positive predictive
720 value (precision), sensitivity (recall), F1-score, and area under the ROC curve (AUROC). We use
721 F1 score to make direct performance comparisons against a regular expression baseline. We eval-
722 uate the performance of our FDG-PET/CT model using positive predictive value (precision), sen-
723 sitivity (recall), F1-score, and area under the ROC curve (AUROC).

Model interpretation. We present three distinct ways to interpret the predictions of the FDG-
PET/CT model—the three-dimensional saliency map shown in Figure 2e, the two-dimensional
saliency map shown in Figure 5d-g, and the three-dimensional attention map shown in Figure 5b-c.
The three-dimensional saliency maps produce a high-level visualization of model behavior. They
are generated using Guided Backpropagation [46]. The gradients from each of the task heads
are often on very different numerical scales. In order to visualize the task gradients jointly, we
preprocess the task gradients. Let $\mathbf{X}'_t \in \mathbb{R}^{2 \times 224 \times 224 \times l}$ be the gradient of scan \mathbf{X} with respect to
the prediction of task head t . We first take the absolute value of \mathbf{X}'_t and then the maximum across
input channels yielding $\hat{\mathbf{X}}'_t \in \mathbb{R}^{224 \times 224 \times l}$, which we can think of as representing scalar saliency
scores for each voxel in the scan. We then normalize these scores by subtracting the minimum
value and dividing by the maximum value, yield $\tilde{\mathbf{X}}'_t \in [0, 1]^{224 \times 224 \times l}$. We manually set clipping
thresholds β_1, \dots, β_T such that for the scan gradient $\tilde{\mathbf{X}}'_t$ of each task head t , we create a saliency

map δ_t :

$$\delta_t = \begin{cases} \tilde{\mathbf{X}}'_t & \text{if } \tilde{\mathbf{X}}'_t \geq \beta_t \\ 0, & \text{otherwise} \end{cases} \quad (10)$$

724 The two-dimensional saliency maps produce a slice-wise interpretation of model behavior for a
725 single task head t . We use the aforementioned normalization scheme to compute $\tilde{\mathbf{X}}'_t$ and set
726 $\beta_t = \min \tilde{\mathbf{X}}'_t$ (i.e. we do not clip perform any clipping).

727 We can additionally visualize the attention scores computed per attention head for a single
728 task. The visualization simply maps the scalar values computed in α to their respective voxels in
729 a. An example of an attention visualization can be seen in Figure 5b-c. For an additional t-SNE
730 visualization of the model activations for each task head, see Supplementary Fig. 4.

731 **Training with Summary Codes.** Each exam in our dataset has a summary code that was assigned
732 by the interpreting radiologist at the time of the study. The summary code is a single digit number,
733 either 1,2,4, or 9, that indicates the degree of abnormality in the exam. A summary code of 1
734 indicates that there are no abnormalities anywhere in the exam and a summary code greater than
735 1 indicates that there is at least one abnormality in the scan. We derive binary abnormality labels
736 from these summary codes and train a single task model to detect abnormalities in the full scan.
737 This whole-body summary code model serves as a fully-supervised baseline against which we can
738 compare our weakly-supervised multi-task model.

739 We train our summary code model with an Adam optimizer and an initial learning rate of
740 $\alpha = 0.0001$ [41]. We anneal the learning rate by half every 16 epochs. We do not employ any
741 regularization aside from data augmentation. Due to GPU memory constraints, we train the scan
742 model with batch size of only 2. In each epoch, we sample 2,000 exams with replacement from
743 the training set.

744 **Mortality prediction and survival analysis.** To evaluate the capacity of our scan model to
745 predict mortality from PET/CT imaging data alone, we frame mortality prediction as a binary
746 classification task: predict whether the patient’s date of death falls within x days of the PET/CT
747 scan date. We fine-tune our abnormality localization models on this mortality prediction task by
748 optimizing the binary cross-entropy loss over a training dataset of the 1,195 exams in the training

749 set for which we have a recorded date of death. We report AUROC evaluated on the 161 exams
750 in the test set for which we have a recorded date of death. We fine-tune and evaluate a separate,
751 single-task model at four different thresholds: $x = 45$, $x = 90$, $x = 180$ and $x = 365$. We compare
752 pretraining on (1) multi-task abnormality localization, (2) single-task abnormality detection using
753 summary codes, and (3) Kinetics [21].

754 We also perform a more traditional survival analysis to explore whether the abnormality
755 localization predictions are predictive of mortality. We fit Cox proportional hazard models on our
756 mortality data using abnormality localization predictions, age, indication, and exam summary code
757 as covariates. Unlike the binary mortality prediction experiments described above, patients without
758 a date of death were included in the survival analysis and right censored at the time of their last
759 PET/CT exam in our data.

760 To get the abnormality localization predictions we apply our weakly-supervised, multi-task
761 abnormality localization model to the scans in the validation and test set. We do this over five ran-
762 dom seeds and take the median of the five scores. This gives us a single probability of abnormality
763 $\bar{y}_t \in [0, 1]$ for each anatomical region t . The probabilities for the 26 anatomical regions with the
764 highest prevalence of abnormality (the same 26 regions we train our multi-task scan model on) are
765 used as covariates in our Cox models.

766 The other covariates (age, indication and summary code) are retrieved from metadata in the
767 DICOM file of the FDG-PET/CT exam. We are unable to retrieve the age of one patient in our test
768 set. This patient is excluded from the analysis. Of the 42 different indications in our dataset, we
769 consider the 13 with more than two occurrences among patients with a recorded date of death in the
770 test set. These 13 indications are: breast cancer restaging, cervical cancer, colorectal cancer, head-
771 neck cancer, lung cancer, lymphoma, ovaries, diagnostic breast cancer, diagnostic lung cancer,
772 diagnostic head-neck cancer. Each patient's indication is one-hot encoded with 13 binary variables.
773 The summary codes, which are recorded by the interpreting radiologist at the time of the study,
774 describe overall patient status and take on a value of 1, 2, 4, or 9, where 1 indicates no evidence
775 of abnormality and 2, 4, and 9 indicate the presence of at least one abnormality. We binarize the
776 summary code $s \in \{1, 2, 4, 9\}$ as $\mathbf{1}[s > 1]$ before providing it as a covariate to the Cox model.

777 When fitting our Cox proportional hazard models, the baseline hazard $h_0(t)$ is modeled non-
778 parametrically using Breslow's method. We do not use a penalizer when fitting. We use the

779 Python implementation provided by lifelines v0.25.1 and the R implementation provided by rms
780 v6.0. [47, 48].

781 In our analysis, we fit Cox models on different subsets of the covariates described above.
782 We fit a multivariable model using all the covariates; single variable models for each covariate
783 separately; a multivariable model with just abnormality localization predictions as covariates; and
784 a multivariable model with just age, indication and summary code as covariates. For each model,
785 we report log hazard ratios with Wald 95% confidence intervals and p values (Supplementary
786 Tables 16-18).

787 When interpreting these hazard ratios, it is important to remember that there are correla-
788 tions among the covariates. This may lead to confounding that could obscure the true effect of the
789 variable on the response. Rather than rely on hazard ratios to evaluate the importance of each pre-
790 dictor in the Cox model, we compute the difference between the Wald χ^2 value and the predictor's
791 degrees of freedom as described in [23, 49]. The greater the difference, the more important the pre-
792 dictor is to the model. In Supplementary Figure 1a, we plot these differences for the multivariable
793 Cox model that uses all covariates. We use the rms R package to compute these differences [48].

794 Two Cox models are considered *nested* when the covariates in one model are a subset of the
795 covariates in the other. If two Cox models are nested, we can compare the fit of one model to
796 the fit of the other using the likelihood ratio test. In our analysis, we fit a model with a full co-
797 variate set Θ consisting of abnormality location predictions, ages, indications, and exam summary
798 codes. We also fit a model with a subset $\hat{\Theta}$ of those covariates, excluding the abnormality location
799 predictions. We compute the test statistic $2(\ell(\Theta) - \ell(\hat{\Theta}))$, where $\ell(\Theta)$ is the log-likelihood of
800 the data when using covariate set Θ . Under the null hypothesis, the test statistic is approximately
801 χ_k^2 distributed with degrees of freedom, k , equal to the difference in cardinality between the two
802 covariate sets [22]. This allows us to compute a p -value. Furthermore, to estimate out-of-sample
803 predictive power, we fit a Cox model on the validation set and compute the concordance index on
804 the test set. Since the Cox models use abnormality location predictions as covariates, we do not fit
805 the models on exams used to train the abnormality localization model.

806 **Related work.** Weak supervision is a broad term used to describe techniques for training machine
807 learning models without hand labeled data [50]. Distant supervision is one such technique that
808 leverages noisy labels in order to train models for a closely related task. This technique is com-

monly used in NLP due to frequent correlation between easily identifiable tokens and high-level semantic meaning [43, 51]. In medical imaging, distant supervision in the form of text-mining has been explored for the classification of pathology in CT and FDG-PET/CT, some incorporating a label ontology similar to our own [52–54]. Ratner et al. build upon the distant supervision paradigm and propose an unsupervised framework that uses generative modeling techniques to denoise labels derived from programmatic, coarse-grain labeling functions [55]. Their framework has been used to achieve state of the art performance on numerous NLP benchmarks, and has additionally proven useful in the medical imaging domain. Fries et al. effectively classify aortic valve malformations in unlabeled cardiac MRI sequences [30]. Dunnmon et al. identify abnormalities in 2-D chest radiographs (CXR), knee extremity radiographs, 3-D head CT scans (HCT), and electroencephalography (EEG) signals [56].

Our proposed weak-supervision framework is related to model distillation, a technique where one model’s predictions are used as probabilistic labels to train a smaller, less memory-intensive classifier [57]. The primary difference between the two methods lies in their higher-level objectives: our framework enables cross-modal learning (i.e. the transfer knowledge between different input modalities) whereas model distillation enables knowledge compression (i.e. learning the same task, but with fewer parameters). Using trained models to label unlabeled data, which is then used to train “student” models has been immensely successful in computer vision tasks such as object detection and human key-point detection [58].

Multi-task learning is a training technique which has been shown to enable learning more generalizable features in some settings, particularly in settings where the number of samples is small [44]. It is common to see multi-task learning employed in medical imaging in order to improve model performance on any one task by training a several related tasks simultaneously [59]. However, such approaches are often limited by the extensive labeling costs often associated with multi-task learning. Some approaches employ semi-supervised learning in order to address this bottleneck, for example through self-teaching segmentation masks [52]. However, our model requires no segmentation data and relies only on non-expert annotations for model training. Our model also incorporates an order of magnitude more tasks than most multi-task learning models, putting it in the realm of massively multi-task learning, a paradigm explored by [15].

Data availability. We will not be releasing our dataset of FDG-PET/CT exams and accompanying

839 reports due to concern for patient privacy. However, the raw metrics underlying Figures 2a, 3b, 3c,
840 4a, 4b, 4c, 4d and 5a are provided as Source Data files.

841 **Code availability.** A Python v3.6 package that includes an implementation of our weak supervi-
842 sion framework, preprocessing code, and experiments is available at: <https://github.com/seyuboglu/weakly-supervised-petct>.
843

844 Acknowledgements

845 Research reported in this publication was supported by the National Library Of Medicine of the
846 National Institutes of Health under Award Number R01LM012966. The content is solely the re-
847 sponsibility of the authors and does not necessarily represent the official views of the National
848 Institutes of Health.

849 This research used data or services provided by STARR, “Stanford Medicine Research Data
850 Repository,” a clinical data warehouse containing live Epic Clarity warehouse data from Stanford
851 Health Care (SHC), the Stanford Childrens Hospital (SCH), the University Healthcare Alliance
852 (UHA) and Packard Children’s Health Alliance (PCHA) clinics and other auxiliary data from
853 Hospital applications such as radiology PACS. STARR is made possible by Stanford School of
854 Medicine Research Office The study has received a grant from General Electric (GE Healthcare,
855 Waukesha, WI). The authors are solely responsible for the design and conduct of the study, all
856 study analyses, the drafting and editing of the manuscript, and its final contents.

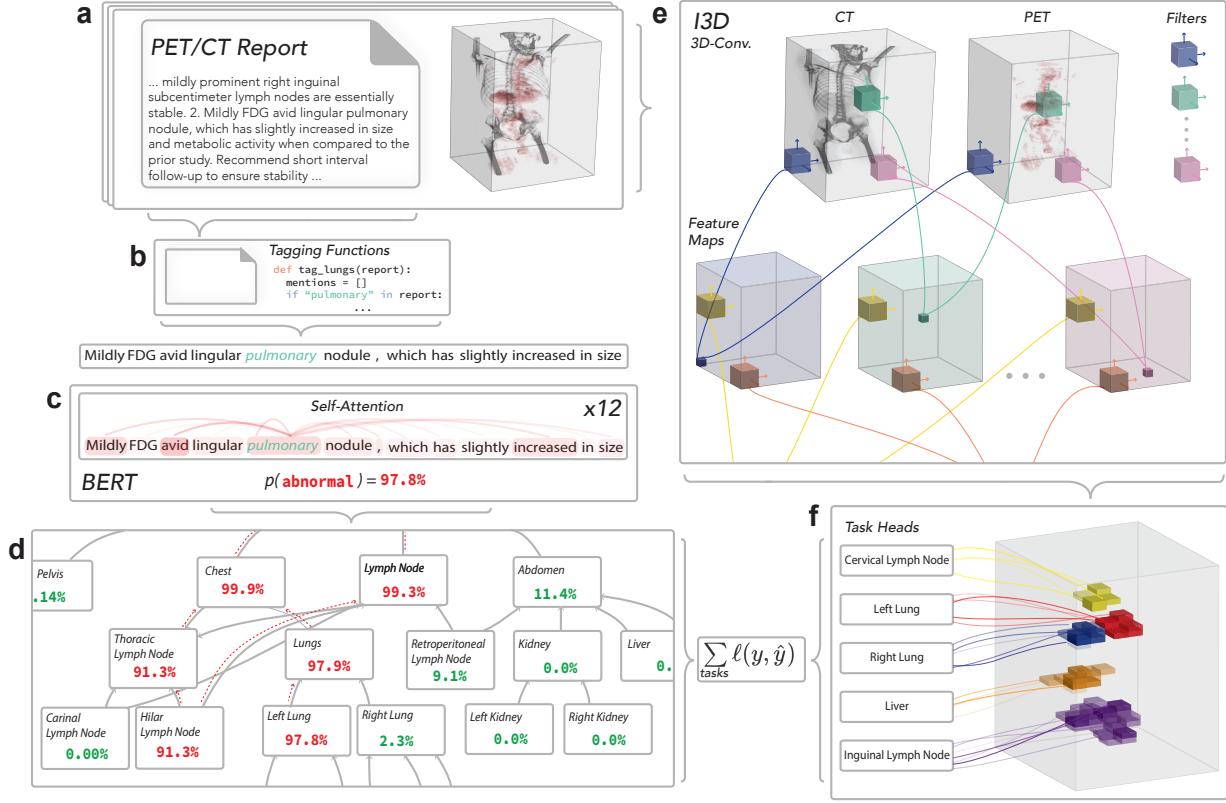


Figure 1: Weak supervision, multi-task learning, and spatial attention are combined to build convolutional neural networks (CNNs) for FDG-PET/CT analysis without hand-labeled training data. (a) Our dataset of 8,144 FDG-PET/CT examinations. Each exam consists of (1) a 3D scan consisting of approximately 250 FDG-PET and CT slices and (2) a natural language, unstructured report written by the interpreting radiologist at the time of the study. Critically, there are no structured, ground truth labels for metabolic abnormalities in each anatomical region. (b) Tagging functions powered by regular expressions extract sentences that mention anatomical regions. (c) A language model predicts whether there is a metabolic abnormality in the tagged anatomical region. (d) Conflicting language model predictions are reconciled into one probability of abnormality for each region. This is done by propagating probabilities up the directed acyclic graph of the ontology. (e) The entire whole-body PET/CT scan is encoded using a 3D-convolutional neural network. This encoding process produces a three-dimensional representation of the original scan. (f) Each task head uses an attention module to extract from the encoded scan the voxels most relevant to the anatomical region it is charged with. We apply a final linear classification layer to weighted-sum of those voxels to make a final binary prediction for each region. To train the scan model, we minimize the cross-entropy loss between the report model predictions (left) and the scan model predictions (right).

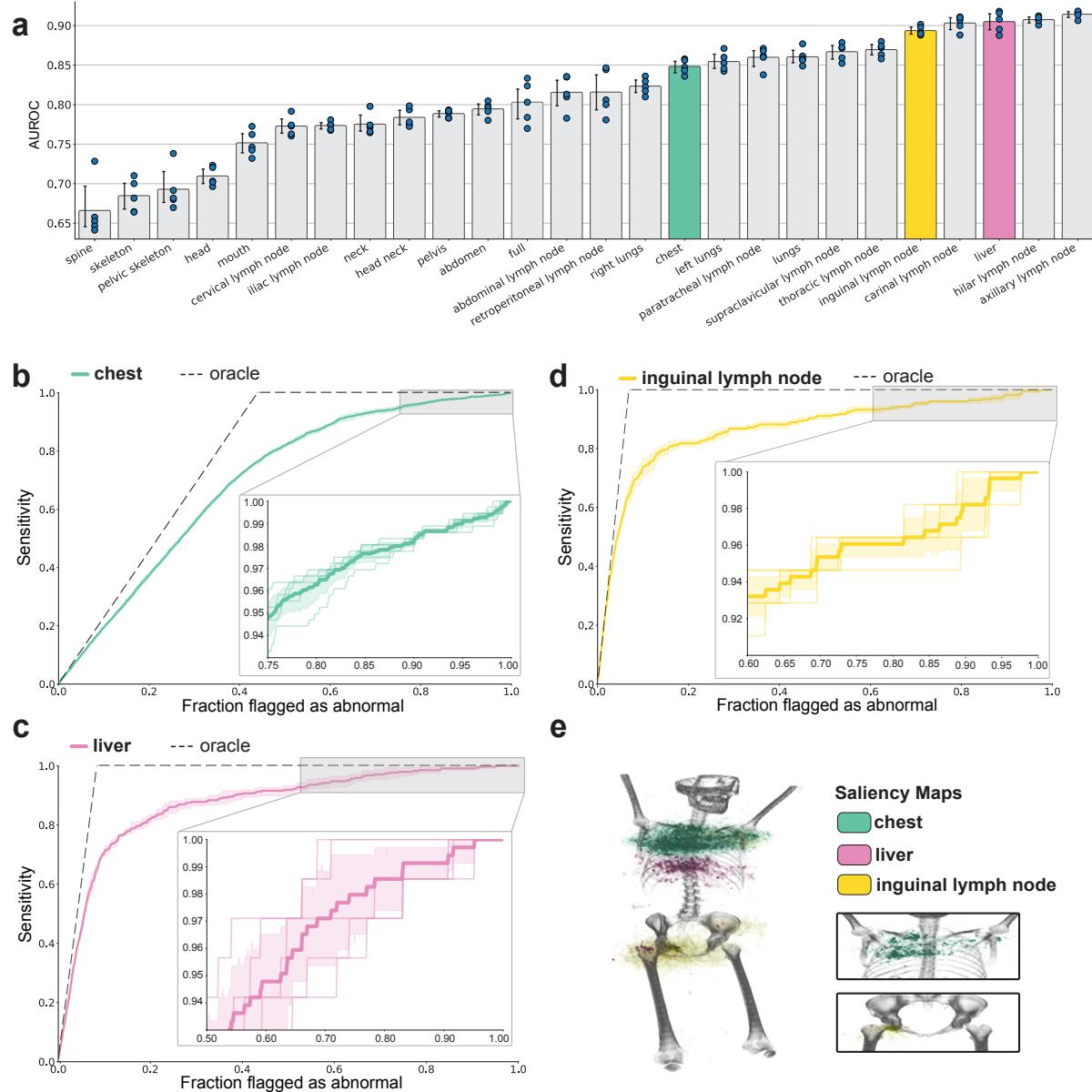


Figure 2: Combining weakly supervised multi-task learning with spatial attention mechanisms enables automated abnormality detection and localization in FDG-PET/CT. (a) Our model’s predictive performance across anatomical regions. Each bar indicates the model’s AUROC for detecting abnormal metabolic activity in a particular anatomical region. Confidence intervals (95%) were determined using bootstrapping with $n = 1,000$ samples from five random initializations. The individual AUROC result for each initialization is shown as a blue dot. Anatomical regions are sorted in order of increasing mean AUROC. (b-d) Sensitivity curve for abnormality detection in the (b) chest, (c) liver, and (d) inguinal lymph nodes. Each point on the curve indicates the sensitivity of our model at a prediction threshold where $x\%$ of the exams in our test set are flagged as abnormal. The sensitivity curve of a perfect detector (i.e. a model that ranks all abnormal examples above all normal examples) is shown in grey. The shaded region represents confidence intervals (95%) computed using bootstrapping with $n = 1,000$ samples from five random initializations. The sensitivity curves for each of those initializations are shown in light (b) green, (c) yellow, and (d) pink. Sensitivity curves illustrate the utility of the model in a potential screening application. With our model, clinicians could ignore around 15% of exams while maintaining 99% sensitivity in liver abnormality screening. (e) 3D saliency map for abnormality detection in the chest (green), liver (pink), and inguinal lymph nodes (yellow). Colored volumes indicate regions where small perturbations to the input scan most effect the model’s prediction for chest, liver, and inguinal lymph nodes.

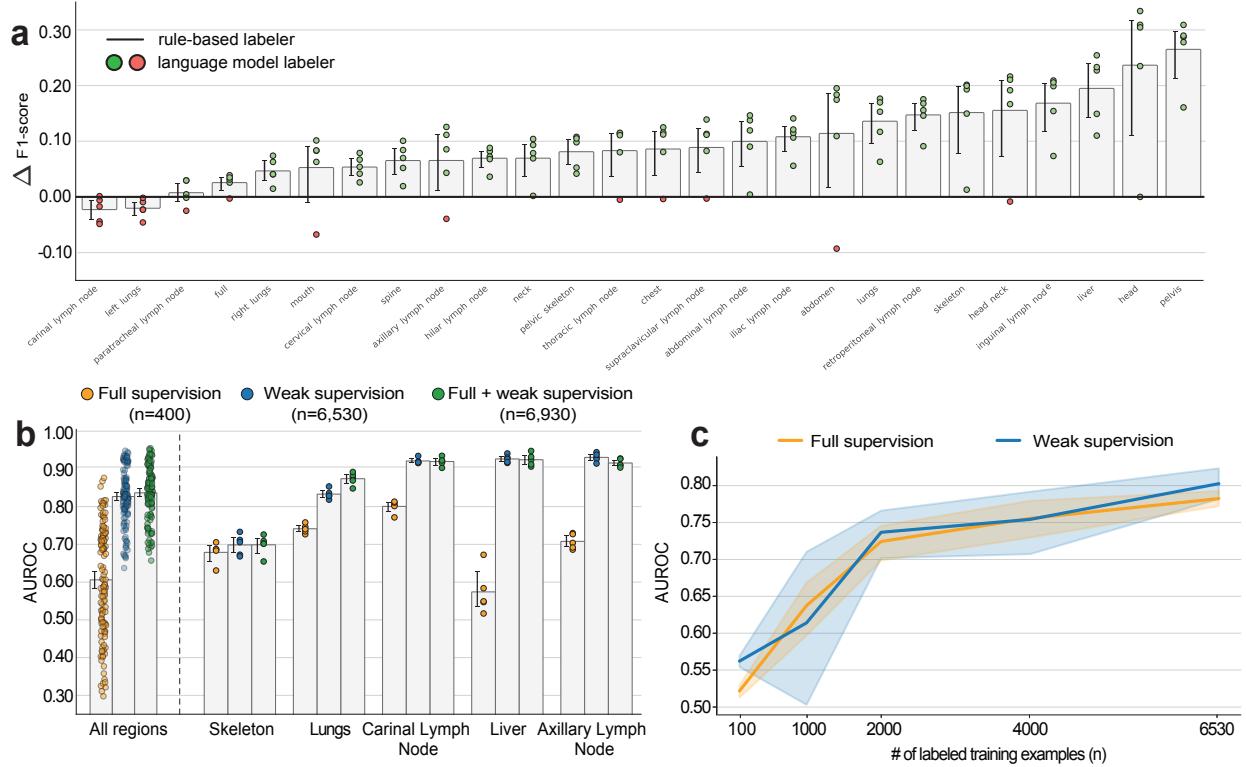


Figure 3: Weak supervision reduces labeling costs for automated abnormality detection and localization. (a) Our labeling framework outperforms a regular expression baseline on 24 of the 26 regions specified in our test set. The five dots for each region represent the difference in F1-score between the baseline and our labeling framework for five different random seeds. For each region, we also show the average difference across seeds and 95% confidence intervals computed via bootstrapping. (b) 26-region abnormality localization models trained on large datasets using weak supervision outperform models fully-supervised on a small dataset. We compare three models: (yellow, left) a fully supervised model trained using $n = 400$ manually annotated training examples, (blue, middle) a weakly supervised model trained using $n = 6,530$ automatically annotated training examples, and (green, right) a hybrid model trained using a combined dataset of $n = 6,930$ training examples, 400 of which are manually annotated. To the left, we show the distribution of AUROC across all 26 regions and five random seeds. Each point represents the AUROC on one region and one random seed. To the right, we compare the supervision approaches on five representative regions: skeleton, lungs, carinal lymph node, liver, and axillary lymph node. Bars represent mean AUROC across all seeds and 95% confidence intervals computed via bootstrapping. The numbers for all regions are provided in Supplementary Table 3. (c) Binary abnormality detection AUROC vs. the number of training data points. The performance of our weakly supervised model is statistically equivalent to that of a traditionally supervised baseline model (trained using summary codes) ($p = 0.1933$ paired permutation test). The shape of the curve suggests that more data, which can be rapidly annotated using our framework, would lead to a continued increase in performance. Line shows the mean AUROC across five different random seeds and shaded regions represent 95% confidence intervals computed via bootstrapping.

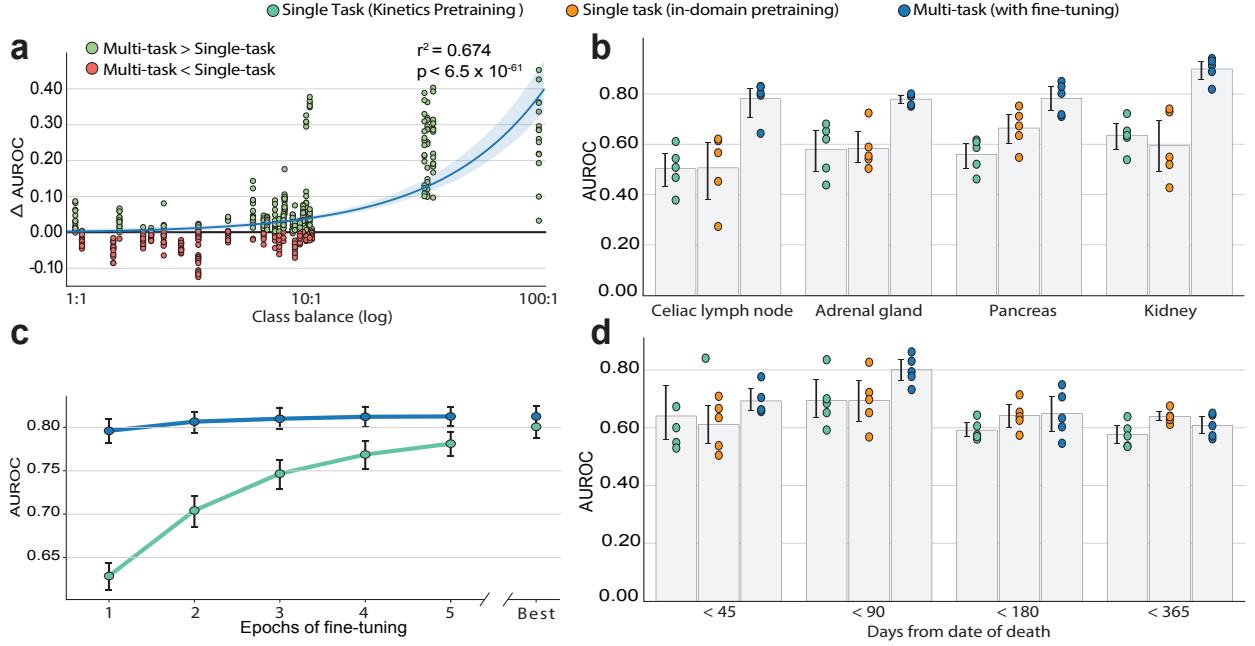


Figure 4: Multi-task learning improves automated localization performance, reduces computational cost, and facilitates mortality prediction. (a) Gains from multi-task learning are the greatest in regions with severe class imbalance. For each region, the difference in abnormality localization AUROC between a multi-task model and a single task model is plotted against the class balance for the region. If a region exhibits class balance 10:1, then for every exam with an abnormality in the region, there will be ten exams where region is normal. The single-task and multi-task models were both trained on five different random seeds and the pairwise differences are plotted. We also fit a linear model to the two variables and report Pearson correlation coefficient $r^2 = 0.674$ and two-sided p-value $p < 6.5 \times 10^{-61}$ computed using the Wald test. (b) Multi-task learning improves performance in regions in which hypermetabolic abnormality is rare. We fine-tune our 26-region multi-task abnormality localization model to detect abnormalities in four new regions beyond the main 26: celiac lymph nodes, the pancreas, the adrenal gland, and the kidneys. Each of these regions exhibit severe class imbalance with more than 30 positive examples for every negative example. We compare this multi-task approach (right, blue) to two other single task approaches: (teal, left) a single task model pre-trained on Kinetics and (orange, middle) a single task model pre-trained on full-body abnormality detection (in-domain pretraining). Shown for each region and modeling approach is AUROC across five different seeds. 95% confidence intervals about the mean are computed via bootstrapping. (c) Multi-task model reduces training complexity. We fine-tune our multi-task FDG-PET/CT model in a single-task setting for each of the twenty-six core anatomical regions. We do the same with a model trained on Kinetics [45], an out-of-domain dataset. The plot shows how mean AUROC across all anatomical regions improves with more epochs of fine-tuning. To the right, we show the best mean AUROC after training to convergence. Confidence intervals (95%) were determined using bootstrapping across five random initializations. (d) Multi-task learning improves performance on the task of predicting mortality from FDG-PET/CT scan alone. We frame mortality prediction as a binary classification task, where the model predicts whether the patients date of death is within x days of the study given only the scan as input. For $x = \{45, 90, 180, 365\}$ we show AUROC across five different random seeds. We compare three approaches: (teal, left) a single task model pre-trained on Kinetics, (orange, middle) a single task model pre-trained on full-body abnormality detection (in-domain pretraining), and (right, blue) multi-task abnormality localization model fine-tuned to predict patient's date of death.

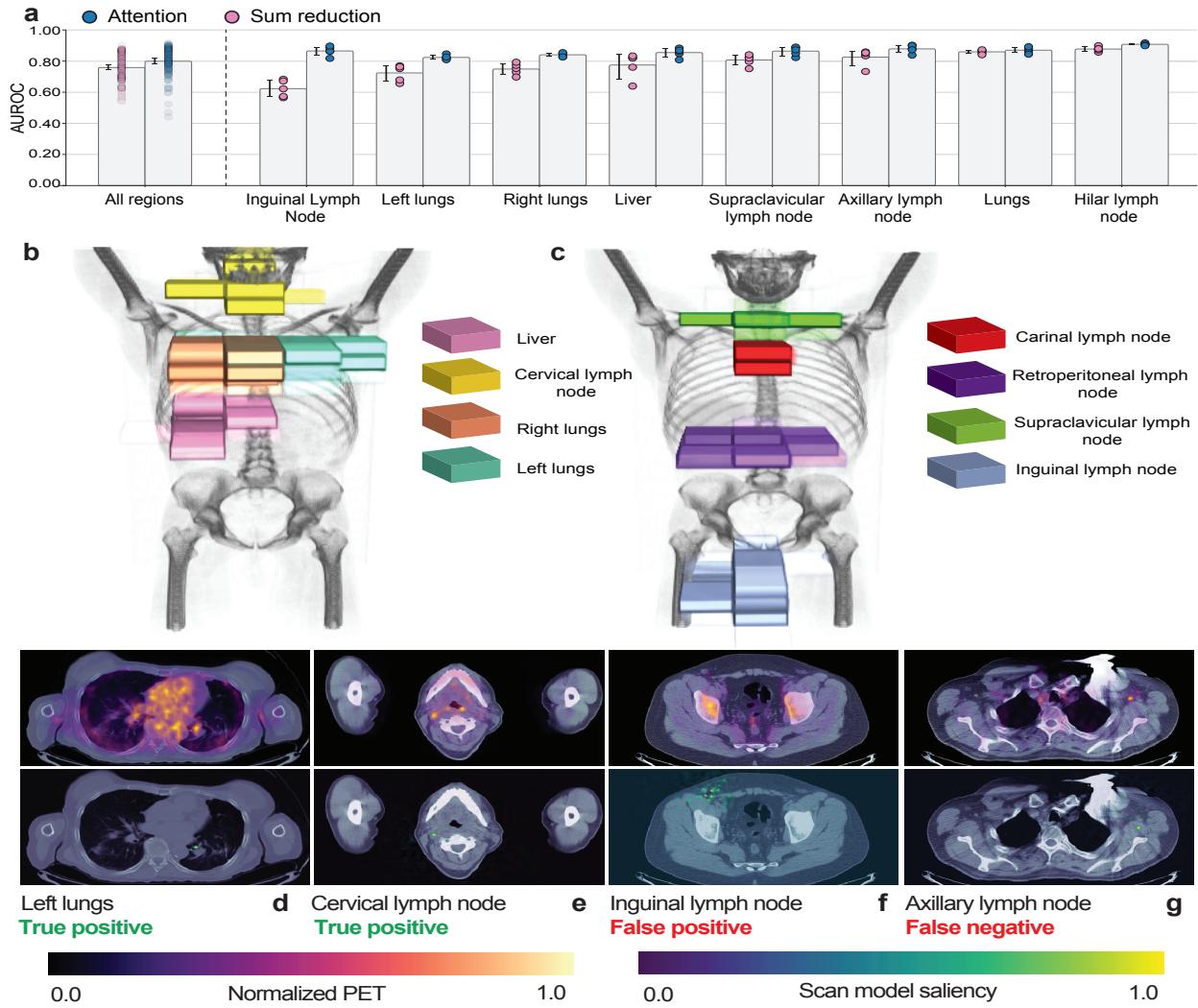


Figure 5: Spatial attention mechanisms improve automated localization performance and facilitate model interpretation. (a) Models trained with the attention mechanism see a modest improvement over a naive sum reduction, with a 3 point increase in mean AUROC taken over all tasks and all seeds ($p = 0.0002$ paired permutation test). We show overall performance (far left), as well as performance on a subset of regions. (b-c) The soft-attention mechanism, supervised only by weak abnormality-localization labels, learns to attend to the appropriate part of the scan. The scan encoder transforms a $2 \times 224 \times 224 \times 200$ PET-CT scan into a $1,024 \times 7 \times 7 \times 33$ encoding. The soft-attention mechanism reduces this encoding to a single 1,024-dimensional vector by computing a weighted sum of the voxels in the encoding (see Methods). The opacity of each voxel is proportional to the weight assigned to the voxel by the soft-attention mechanism. (d-g) Two-dimensional saliency maps facilitate the interpretation of correct (d-e) and incorrect (f-g) abnormality localization predictions. Shown in the first row are FDG-PET uptake values overlaid over CT. FDG-PET values are normalized by subtracting the minimum value in the slice and dividing by the maximum. Additionally, in order to decrease image scatter, we clip PET values at the value of the 60th percentile pixel in the slice. Shown in the second row are saliency maps overlaid over CT. Saliency maps are produced by computing the gradient of the scan model prediction with respect to the input scan and normalizing such that the saliency for each pixel is in the range [0, 1] (See Methods). The color bars for the normalized PET and saliency values are shown in the legend. (d) The model accurately detects a hypermetabolic pulmonary nodule along the medial aspect of the left lung shown in the PET (top). The model's saliency is highly concentrated in the pixels containing the nodule (bottom). (e) A right level II cervical hypermetabolic lymph node (top) is detected by the scan model. The saliency map is focused precisely on the abnormality (bottom). (f) The model detects abnormal mild FDG uptake in a subcentimeter right inguinal lymph node, around which there is a significant concentration of saliency (bottom panel). However, the clinical report does not mention this finding because, in the context of the patient's cancer, the lymph node is deemed non-specific and likely physiologic or inflammatory. Because the model and the clinical report are not congruent, this is a classified as a false positive, however the model is flagging FDG uptake that at least warrants review. (g) The model fails to detect an abnormality in a hypermetabolic left axillary lymph node of a patient with metastatic breast cancer and a pacemaker in the left chest wall (left panel). While the saliency map is concentrated on the abnormal node, the threshold for an abnormal prediction is not reached perhaps because of the the overlying beam hardening artifact from the pacemaker.

857 **References**

- 858
859 1. Rajpurkar, P. *et al.* CheXNet: Radiologist-Level Pneumonia Detection on Chest X-Rays with
860 Deep Learning. *Tech. Rep.*
- 861 2. Esteva, A. *et al.* Dermatologist-level classification of skin cancer with deep neural networks.
862 *Nature* **542**, 115–118 (2017).
- 863 3. Gulshan, V. *et al.* Development and Validation of a Deep Learning Algorithm for Detection
864 of Diabetic Retinopathy in Retinal Fundus Photographs. *JAMA* **316**, 2402–2402 (2016).
- 865 4. Kooi, T. *et al.* Large scale deep learning for computer aided detection of mammographic
866 lesions. *Medical Image Analysis* **35**, 303–312 (2017).
- 867 5. Oakden-Rayner, L., Dunnmon, J., Carneiro, G. & R, C. Hidden stratification causes clinically
868 meaningful failures in machine learning for medical imaging (2019). [1909.12475](#).
- 869 6. Hutchings, M. *et al.* Position emission tomography with or without computed tomography in
870 the primary staging of Hodgkin’s lymphoma. *Haematologica* **91**, 482–489 (2006).
- 871 7. El-Galaly, T. C., Gormsen, L. C. & Hutchings, M. PET/CT for Staging; Past, Present, and
872 Future. *Seminars in Nuclear Medicine* **48**, 4–16 (2018).
- 873 8. Young, L. PET/CT drives PET scan volume to new heights (2019).
- 874 9. National Cancer Policy Forum, Board on Health Care Services, Institute of Medicine & Na-
875 tional Academies of Sciences, Engineering, and Medicine. *Appropriate Use of Advanced*
876 *Technologies for Radiation Therapy and Surgery in Oncology: Workshop Summary* (National
877 Academies Press, Washington, D.C., 2016).
- 878 10. Saab, K. *et al.* Doubly weak supervision of deep learning models for head ct. In *Inter-
879 national Conference on Medical Image Computing and Computer-Assisted Intervention*, 811–
880 819 (Springer, 2019).
- 881 11. Sible, L. *et al.* 18F-FDG PET/CT Uptake Classification in Lymphoma and Lung Cancer by
882 Using Deep Convolutional Neural Networks. *Radiology* **191114** (2019).
- 883 12. Dunnmon, J. A. *et al.* Assessment of Convolutional Neural Networks for Automated Classifi-
884 cation of Chest Radiographs. *Radiology* **290**, 537–544 (2019).
- 885 13. Ratner, A. J., De Sa, C. M., Wu, S., Selsam, D. & Ré, C. Data Programming: Creating Large
886 Training Sets, Quickly. In Lee, D. D., Sugiyama, M., Luxburg, U. V., Guyon, I. & Garnett, R.
887 (eds.) *Advances in Neural Information Processing Systems* 29, 3567–3575 (Curran Associates,
888 Inc., 2016).
- 889 14. Irvin, J. *et al.* CheXpert: A Large Chest Radiograph Dataset with Uncertainty Labels and
890 Expert Comparison. *arXiv:1901.07031 [cs, eess]* (2019). [1901.07031](#).
- 891 15. Ratner, A. J., Hancock, B. & Ré, C. The role of massively multi-task and weak supervision in
892 software 2.0. In *CIDR* (2019).
- 893 16. Avati, A. *et al.* Improving palliative care with deep learning. *BMC Medical Informatics and*
894 *Decision Making* **18** (2018).

- 895 17. Banerjee, I. *et al.* Probabilistic Prognostic Estimates of Survival in Metastatic Cancer Patients
896 (PPES-Met) Utilizing Free-Text Clinical Narratives. *Scientific Reports* **8** (2018).
- 897 18. Wang, A. *et al.* Superglue: A stickier benchmark for general-purpose language understanding
898 systems. *CoRR* **abs/1905.00537** (2019). [1905.00537](#).
- 899 19. Irvin, J. *et al.* Chexpert: A large chest radiograph dataset with uncertainty labels and expert
900 comparison. *arXiv preprint arXiv:1901.07031* (2019).
- 901 20. Peng, Y. *et al.* Negbio: a high-performance tool for negation and uncertainty detection in
902 radiology reports. *CoRR* **abs/1712.05898** (2017). [1712.05898](#).
- 903 21. Kay, W. *et al.* The Kinetics Human Action Video Dataset. *arXiv:1705.06950 [cs]* (2017).
904 ArXiv: 1705.06950.
- 905 22. Harrell, F. E. Overview of Maximum Likelihood Estimation. In Harrell, J., Frank E. (ed.)
906 *Regression Modeling Strategies: With Applications to Linear Models, Logistic and Ordinal*
907 *Regression, and Survival Analysis*, Springer Series in Statistics, 181–217 (Springer Interna-
908 tional Publishing, Cham, 2015).
- 909 23. Harrell, F. E. Describing, Resampling, Validating, and Simplifying the Model. In Harrell,
910 J., Frank E. (ed.) *Regression Modeling Strategies: With Applications to Linear Models, Lo-*
911 *gistic and Ordinal Regression, and Survival Analysis*, Springer Series in Statistics, 103–126
912 (Springer International Publishing, Cham, 2015).
- 913 24. Huang, B. *et al.* Fully Automated Delineation of Gross Tumor Volume for
914 Head and Neck Cancer on PET-CT Using Deep Learning: A Dual-Center Study.
915 <https://www.hindawi.com/journals/cmmi/2018/8923028/> (2018). Library Catalog:
916 www.hindawi.com Pages: e8923028.
- 917 25. Saab, K. *et al.* Doubly Weak Supervision of Deep Learning Models for Head CT. In Shen, D.
918 *et al.* (eds.) *Medical Image Computing and Computer Assisted Intervention – MICCAI 2019*,
919 vol. 11766, 811–819 (Springer International Publishing, Cham, 2019). Series Title: Lecture
920 Notes in Computer Science.
- 921 26. Dunnmon, J. *et al.* Cross-Modal Data Programming Enables Rapid Medical Machine Learn-
922 ing. *arXiv:1903.11101 [cs, eess, stat]* (2019). [1903.11101](#).
- 923 27. Goodfellow, I., Bengio, Y. & Courville, A. *Deep learning* (MIT press, 2016).
- 924 28. Niederkohr, R. D. *et al.* Reporting Guidance for Oncologic 18F-FDG PET/CT Imaging. *Jour-*
925 *nal of Nuclear Medicine* **54**, 756–761 (2013).
- 926 29. Lowe, H. J., Ferris, T. A., Hernandez, P. M. & Weber, S. C. STRIDE – An Integrated
927 Standards-Based Translational Research Informatics Platform. *AMIA Annual Symposium Pro-*
928 *ceedings* **2009**, 391–395 (2009).
- 929 30. Fries, J. A. *et al.* Weakly supervised classification of aortic valve malformations using unla-
930 beled cardiac MRI sequences. *Nature Communications* **10**, 1–10 (2019).
- 931 31. Reed, S. *et al.* Training Deep Neural Networks on Noisy Labels with Bootstrapping.
932 *arXiv:1412.6596 [cs]* (2014). [1412.6596](#).

- 933 32. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of Deep Bidirectional
934 Transformers for Language Understanding. *arXiv:1810.04805 [cs]* (2018). [1810.04805](#).
- 935 33. Phang, J., Févry, T. & Bowman, S. R. Sentence Encoders on STILTs: Supplementary Training
936 on Intermediate Labeled-data Tasks. *arXiv:1811.01088 [cs]* (2018). [1811.01088](#).
- 937 34. Radford, A., Narasimhan, K., Salimans, T. & Sutskever, I. Improving Language Understanding
938 by Generative Pre-Training 12.
- 939 35. Peters, M. E. *et al.* Deep contextualized word representations. *arXiv:1802.05365 [cs]* (2018).
940 [1802.05365](#).
- 941 36. Wu, Y. *et al.* Google’s Neural Machine Translation System: Bridging the Gap between Human
942 and Machine Translation. *arXiv:1609.08144 [cs]* (2016). [1609.08144](#).
- 943 37. Zhu, Y. *et al.* Aligning Books and Movies: Towards Story-Like Visual Explanations by Watch-
944 ing Movies and Reading Books. In *2015 IEEE International Conference on Computer Vision*
945 (*ICCV*), 19–27 (IEEE, Santiago, Chile, 2015).
- 946 38. Sennrich, R., Haddow, B. & Birch, A. Neural Machine Translation of Rare Words with Sub-
947 word Units. In *Proceedings of the 54th Annual Meeting of the Association for Computational*
948 *Linguistics (Volume 1: Long Papers)*, 1715–1725 (Association for Computational Linguistics,
949 Berlin, Germany, 2016).
- 950 39. Vaswani, A. *et al.* Attention Is All You Need. Tech. Rep.
- 951 40. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. Bert: Pre-training of deep bidirectional
952 transformers for language understanding. In *Proceedings of the 2019 Conference of the North*
953 *American Chapter of the Association for Computational Linguistics: Human Language Tech-*
954 *nologies, Volume 1 (Long and Short Papers)*, 4171–4186 (2019).
- 955 41. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. *arXiv:1412.6980 [cs]*
956 (2014). [1412.6980](#).
- 957 42. Luo, Y., Tao, D., Geng, B., Xu, C. & Maybank, S. J. Manifold Regularized Multitask Learning
958 for Semi-Supervised Multilabel Image Classification. *IEEE Transactions on Image Processing*
959 **22**, 523–536 (2013).
- 960 43. Rei, M. Semi-supervised Multitask Learning for Sequence Labeling. *arXiv:1704.07156 [cs]*
961 (2017). ArXiv: 1704.07156.
- 962 44. Caruana, R. Multitask Learning 35.
- 963 45. Carreira, J., Zisserman, A., Com, Z. & Deepmind, t. Quo Vadis, Action Recognition? A New
964 Model and the Kinetics Dataset. Tech. Rep.
- 965 46. Springenberg, J. T., Dosovitskiy, A., Brox, T. & Riedmiller, M. Striving for Simplicity: The
966 All Convolutional Net. *arXiv:1412.6806 [cs]* (2014). [1412.6806](#).
- 967 47. Davidson-Pilon, C. *et al.* Camdavidsonpilon/lifelines: v0.25.1 (2020).
- 968 48. Harrell, F. E. Rms: Regression Modeling Strategies (2020).

- 969 49. Asher, A. L. *et al.* An analysis from the Quality Outcomes Database, Part 2. Predictive model
970 for return to work after elective surgery for lumbar degenerative disease. *Journal of Neuro-*
971 *surgery: Spine* **27**, 370–381 (2017).
- 972 50. Saab, K. *et al.* Doubly weak supervision of deep learning models for head ct. In Shen, D.
973 *et al.* (eds.) *Medical Image Computing and Computer Assisted Intervention – MICCAI 2019*,
974 811–819 (Springer International Publishing, Cham, 2019).
- 975 51. Go, A., Bhayani, R. & Huang, L. Twitter Sentiment Classification using Distant Supervision 6.
- 976 52. Khosravan, N. & Bagci, U. Semi-supervised multi-task learning for lung cancer diagnosis. In
977 *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology*
978 *Society (EMBC)*, 710–713 (IEEE, 2018).
- 979 53. Singh, S. *et al.* Deep-learning-based classification of fdg-pet data for alzheimer's disease
980 categories. In *13th International Conference on Medical Information Processing and Analysis*,
981 vol. 10572, 105720J (International Society for Optics and Photonics, 2017).
- 982 54. Yan, K. *et al.* Holistic and comprehensive annotation of clinically significant findings on
983 diverse ct images: learning from radiology reports and label ontology. In *Proceedings of the*
984 *IEEE Conference on Computer Vision and Pattern Recognition*, 8523–8532 (2019).
- 985 55. Ratner, A. *et al.* Snorkel: Rapid Training Data Creation with Weak Supervision. *Proceedings*
986 *of the VLDB Endowment* **11**, 269–282 (2017). ArXiv: 1711.10160.
- 987 56. Dunnmon, J. *et al.* Cross-modal data programming enables rapid medical machine learning.
988 *CoRR abs/1903.11101* (2019). [1903.11101](#).
- 989 57. Hinton, G., Vinyals, O. & Dean, J. Distilling the Knowledge in a Neural Network.
990 *arXiv:1503.02531 [cs, stat]* (2015). ArXiv: 1503.02531.
- 991 58. Radosavovic, I., Dollar, P., Girshick, R., Gkioxari, G. & He, K. Data Distillation: Towards
992 Omni-Supervised Learning. In *2018 IEEE/CVF Conference on Computer Vision and Pattern*
993 *Recognition*, 4119–4128 (IEEE, Salt Lake City, UT, USA, 2018).
- 994 59. Moeskops, P. *et al.* Deep learning for multi-task medical image segmentation in multiple
995 modalities. In *International Conference on Medical Image Computing and Computer-Assisted*
996 *Intervention*, 478–486 (Springer, 2016).