# Towards trustworthy seizure onset detection using workflow notes

Khaled Saab,[1] Siyi Tang,[1] Mohamed Taha,[2]
Christopher Lee-Messer,[3,*] Christopher Ré,[4,*] and Daniel L. Rubin[5,*]

[1]Department of Electrical Engineering, Stanford University, Stanford, CA, USA
[2]Department of Neurology, Stanford University, Stanford, CA, USA
[3]Department of Child Neurology, Stanford University, Stanford, CA, USA
[4]Department of Computer Science, Stanford University, Stanford, CA, USA
[5]Department of Biomedical Data Science, Radiology, and Medicine, Stanford University, Stanford, CA, USA

*Equal contribution.
Address correspondence to Khaled Saab (ksaab@stanford.edu)

## Abstract

A major barrier to deploying healthcare AI models is their trustworthiness. One form of trustworthiness is a model's robustness across different subgroups: while existing models may exhibit expert-level performance on aggregate metrics, they often rely on non-causal features, leading to errors in hidden subgroups. To take a step closer towards trustworthy seizure onset detection from EEG, we propose to leverage annotations that are produced by healthcare personnel in routine clinical workflows – which we refer to as workflow notes – that include multiple event descriptions beyond seizures. Using workflow notes, we first show that by scaling training data to an unprecedented level of 68,920 EEG hours, seizure onset detection performance significantly improves (+12.3 AUROC points) compared to relying on smaller training sets with expensive manual gold-standard labels. Second, we reveal that our binary seizure onset detection model underperforms on clinically relevant subgroups (e.g., up to a margin of 6.5 AUROC points between pediatrics and adults), while having significantly higher false positives on EEG clips showing non-epileptiform abnormalities compared to any EEG clip (+19 FPR points). To improve model robustness to hidden subgroups, we train a multilabel model that classifies 26 attributes other than seizures, such as spikes, slowing, and movement artifacts. We find that our multilabel model significantly improves overall seizure onset detection performance (+5.9 AUROC points) while greatly improving performance among subgroups (up to +8.3 AUROC points), and decreases false positives on non-epileptiform abnormalities by 8 FPR points. Finally, we propose a clinical utility metric based on false positives per 24 EEG hours and find that our multilabel model improves this clinical utility metric by a factor of 2× across different seizure onset detection recall and latency times. These results demonstrate the importance of leveraging additional cost-effective supervision to improve model robustness to classification errors in patient subgroups.

# 1 Introduction

The scalp electroencephalogram (EEG) is a non-invasive and valuable technique to measure the brain's electrical activity. Unlike other modalities that image the brain (e.g., fMRI, PET), EEG enables continuous analysis of rapid changes in the brain's electrical activity. In the intensive care unit (ICU), EEG is critical for the detection of seizures that may lack a behavioral correlate and worsen brain injury. Moreover, EEG is an essential tool to diagnose and care for epileptic patients of all ages[1].

While analyzing EEG data is a critical healthcare task, it poses several challenges. First, the continuous recording of hours of multi-channel EEG results in a vast amount of data that requires thorough interpretation, which is a highly time-consuming and costly task that demands deep neurologic-epileptologic understanding. Second, the gold-standard for EEG analysis is done by fellowship trained clinical neurophysiologists, who have not only been trained to identify seizure patterns, but also many common artifacts. For example, common artifacts on EEG signals may include muscle movement or environment noise, along with countless non-epileptiform abnormalities such as spikes and slowing. Finally, there is a shortage of EEG specialists, and as a result, low resource communities lack access to EEG interpretation[2]. Thus, there is a strong need to develop reliable tools that help clinicians analyze EEG data more efficiently.

Many studies have shown that deep learning (DL) techniques present great promise for automated seizure detection. There have been substantial efforts for curating large and publicly available EEG datasets, such as the Temple University Hospital Seizure Detection (TUSZ) corpus that includes thousands of EEGs from hundreds of patients[3,4]. The availability of large public datasets has enabled rapid progress in benchmarking and improved seizure detection models[5–10]. Recently, a DL model named SParCNet was trained on 6,097 EEGs from 2,711 patients, annotated independently by 20 fellowship-trained neurophysiologists, and was found to match or exceed most experts in classifying seizures[11].

Due to the high-stakes nature of healthcare, trustworthiness of DL models remains a major roadblock to clinical adoption[12,13]. Alarmingly, there has been a growing body of work revealing that healthcare models with "expert-level" performance often rely on non-generalizable features[14,15], resulting in unexpected drops in performance over hidden subgroups[16,17] or under data distribution shifts[18]. While many studies report impressive overall seizure detection performance[6,11], such studies lack the in-depth analysis needed to understand the clinically meaningful failure modes of existing models. For example, pediatric EEGs look drastically different from adult EEGs, different seizure types display unique EEG patterns, and there may be different types of abnormalities present in EEGs recorded from the ICU as compared to other clinical settings[1]; as a result, models may underperform on specific age groups, seizure subtypes, or ICU patients. Unfortunately, conducting an in-depth error analysis requires manual interpretation of both EEGs and model predictions over a diverse set of studies, making it a costly process. However, a clear understanding of a model's systematic errors is critical to provide trust in model predictions for clinical adoption.

In this work, we provide a strategy to scale training data, conduct a subgroup robustness analysis, and improve the trustworthiness of seizure onset detection models in a cost-effective manner. As opposed to relying on expensive gold-standard labels, which require a fellow-trained neurophysiologist to label EEGs outside existing clinical workflows, we propose to leverage seizure annotations that are produced by healthcare personnel within existing clinical workflows[5] – which we refer to as workflow notes. Since workflow notes are produced as part of routine clinical practice, we are able to train our DL models on an unprecedented scale of $68,920$ EEG hours. To conduct an in-depth error analysis we stratify the evaluation set of EEG recordings into clinically-relevant subgroups and analyze discrepancies in seizure onset detection performance in each subgroup. In particular, we use a combination of patient metadata (e.g., age), expert-provided subgroup labels (e.g., seizure types), along with numerous EEG
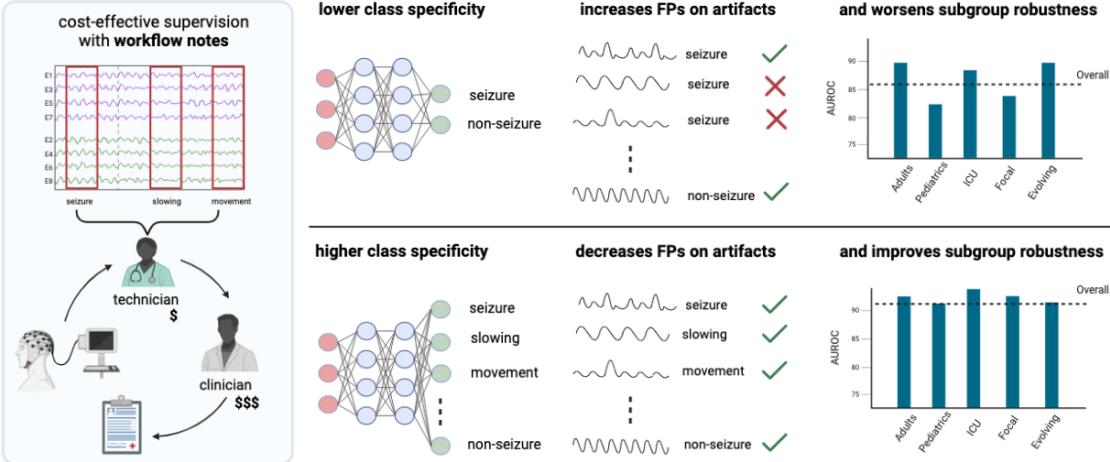
Figure 1: **Results overview:** We find that increasing class specificity by providing additional supervision decreases false positives on artifacts and improves subgroup robustness. Importantly, we supervise our models on large scale data (68,920 EEG hours) using readily available notes produced within clinical workflows (left panel).

attributes, such as spikes, slowing, movements, jerks, photoelectric stimulation, hyperventilation, and more (full list in Supplementary Table 4), that are readily available from workflow notes.

To improve model robustness to non-epileptiform abnormalities and hidden subgroups, we utilize the workflow notes to increase class specificity. Specifically, as opposed to training a binary classification model (seizure or no seizure onset), we train a multilabel model to classify 25 classes in addition to seizure onset, such as spikes, slowing, and hyperventilation. Additionally, we study how our improvements in seizure onset detection robustness translate to clinical utility by tracking the false positives per 24 hours for different deployment settings.

## 2 Results

We first describe how we utilize workflow notes to scale supervision to 68,920 EEG hours (4,125,225 60-sec EEG clips) in a cost-effective manner, and show that training a model to detect seizure onset using workflow notes greatly improves performance compared with a model trained with a smaller set of gold-standard, expert-labeled EEG clips (Section 2.1). We further utilize the workflow notes to reveal that even with large-scale training, our binary seizure onset detection model underperforms on clinically-relevant subgroups of patients, and has higher false positive rates for non-seizure EEG clips with abnormal patterns (Section 2.2). To improve our model's performance across subgroups, we train a multilabel model to classify 25 attributes extracted from the workflow notes, in addition to seizure onset (Section 2.3, Figure 1). Finally, we propose a metric of clinical utility to assess the degree to which the multilabel model improves clinical utility over a range of settings (Section 2.4).

### 2.1 Scaling training data with workflow notes

**Seizure onset detection task.** Following previous studies[5,19], our task of interest is to classify the existence of a seizure onset in a 60 second EEG clip. Each EEG contains 19 electrodes that sample voltage readings at 200 Hz, therefore the input to the model is a 60-sec EEG clip $x \in \mathbb{R}^{12,000 \times 19}$ and the output is a binary label $y \in \{0, 1\}$ indicating the existence of a seizure onset in that clip. To evaluate and compare the performance of deep learning models on the task of seizure-onset detection,
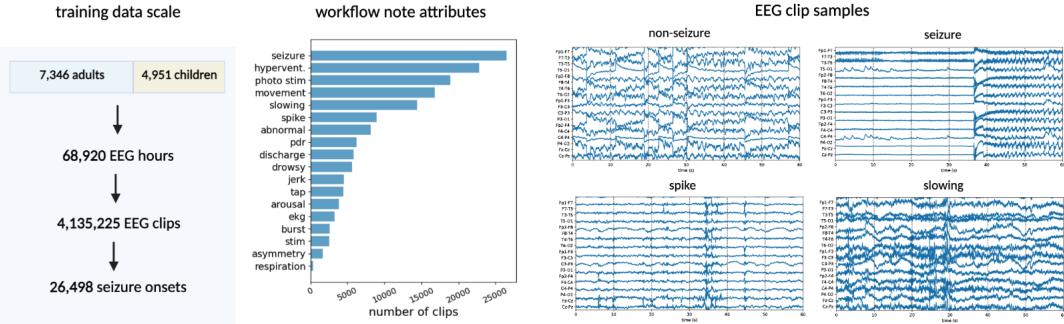
Figure 2: **Training dataset overview:** In the left panel, we provide statistics on the scale our training dataset of EEG recordings aggregated from adult and pediatric hospitals. In the middle panel, we plot the histogram of atribute labels extracted from workflow notes. In the right panel, we visualize four EEG clips, three of which are non-seizure EEG clips. The non-seizure EEG clips exhibit significant differences in temporal features, motivating the opportunity to use them to increase class specificity.

we curated a gold-standard evaluation set of 626 EEG hours (37,588 60-sec EEG clips) labeled by two fellowship-trained EEG readers (details in Section 4.1).

**Workflow notes.** Since acquiring gold-standard labels for all 68,920 hours of EEG (or 4,135,225 clips) would be extremely costly, we used a cost-effective technique that leverages workflow notes proposed by Saab et al.[5]. As visualized in Figure 1, a standard clinical workflow for EEG analysis starts with (1) EEG data collection, after which, (2) a mixed group of (mostly) technicians, fellows, and students are tasked with annotating any salient events that may be useful for the final stage, where (3) fellowship trained clinical neurophysiologists give a final diagnosis in a written report. Importantly, while the annotators in the second stage are less skilled and may therefore make mistakes, their annotations are well suited for ML supervision due to their fine-grained structure and temporal specificity. In particular, the annotations consist of repetitive standard descriptions of events such as seizures, spikes, and movements, and are produced with precise timestamps of when each event occured. Since the annotations only contain the start times of events, we only label the clip in which that event began, and we do not label subsequent clips unless another (or different) event occurs. The workflow notes are readily available for all EEGs in both our adult and pediatric hospitals, allowing us to scale our training data to unprecedented levels.

Each EEG recording is accompanied by a table of workflow notes, where each row is a logged note containing the text describing an event along with a timestamp representing the start of the event. We found 26 relevant attributes from our manual analysis, and wrote simple regular expressions to extract the unique attributes from the workflow notes (e.g., considering synonyms and case-insensitivity; details in Section 4.1). Figure 2 displays a histogram of the 18 most frequent attributes, where for example we have seizure onset annotations for 26,498 EEG clips, spike annotations for 8,942 EEG clips, and movement artifact annotations for 16,806 EEG clips.

**Impact of scaling labeled training data with workflow notes.** We hypothesize that even though workflow notes may contain errors and our regular expressions may extract noisy labels, leveraging workflow notes to scale the training data results in better performing models compared to training models using a much smaller subset of gold-standard labels. To test our hypothesis, we randomly split our gold-standard labeled dataset into train (50%), validation (10%), and test (40%) sets, stratified by patients (i.e., there are no overlapping patients among the three splits). We then trained two classification models, where the first model was trained using the gold-labeled train set (containing 16,058 EEG clips, of which 408 contained a seizure onset), and the second model was trained using the entire training

set that was not gold-labeled, resulting in 4,097,637 EEG clips, of which 25,254 contained seizure onset labels extracted from the workflow notes. Details on model architecture and training procedure can be found in Section 4. To evaluate seizure onset detection performance, we assessed the Area Under the Receiver Operating Characteristic curve (AUROC) on the held-out test set, and report the 95% confidence intervals.

Leveraging the workflow notes improved the model's performance, where the model trained on the smaller gold-labeled dataset achieved an AUROC of $73.3 \pm 3.2$, and the model trained on the much larger workflow-labeled dataset achieved an AUROC of $85.6 \pm 0.9$.
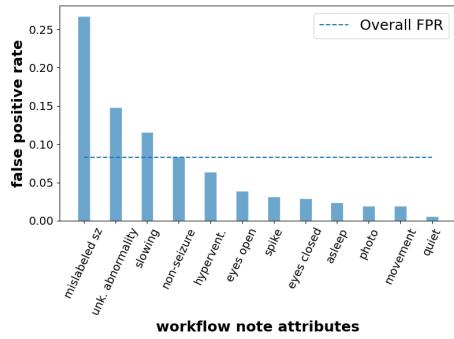
## 2.2 Revealing underperforming subgroups

To evaluate whether our models performed less well in certain patient subgroups, we performed a subgroup analysis where we evaluated the change in model performance across multiple clinically-relevant subgroups. We carried out the subgroup analysis by using a collection of patient metadata, gold-labeled seizure types, and the attributes from the workflow notes.

For patient subgroups, we recorded whether the patient was from the adult or pediatric hospital, and whether a patient's EEG recordings were collected in the ICU. For seizure subtypes, we analyzed performance differences among the focal spike-and-wave, evolving rhytmic slowing, and generalized spike-and-wave types (more details in Section 4.1).

| | Subgroup | AUROC |
|---|---|---|
| | Overall | $85.6 \pm 0.9$ |
| patient subgroups | Adults | $89.4 \pm 1.1$ |
| | Adults outside ICU | $89.4 \pm 1.7$ |
| | Adults from ICU | $88.5 \pm 1.2$ |
| | Pediatrics | $82.9 \pm 1.5$ |
| seizure subgroup | Focal spike-and-wave | $84.3 \pm 2.6$ |
| | Evolving rhythmic slowing | $89.8 \pm 3.3$ |
| | Generalized spike-and-wave | $85.5 \pm 4.0$ |

a



b

Figure 3: **Subgroup analysis:** (a): Model classification performance (AUROC with 95% confidence intervals) for both patient and seizure subgroups. (b): False positive rate among workflow attributes.

From our subgroup analysis on patient and seizure types in Figure 3a, we find that our model performed better for patients from the adult hospital with a 6.5 AUROC point difference compared to patients from the pediatric hospital. There were also differences in the performance of the model for various seizure types, with a 5.5 AUROC point difference between focal spike-and-wave and evolving rhythmic slowing seizures. From our subgroup analysis on workflow attributes in Figure 3b, our model had significantly higher false positive rates (FPR) with respect to seizure onset detection for EEG clips with non-seizure abnormalities (FPR of 0.27) compared to overall EEG clips (FPR of 0.08). Details on metrics can be found in Section 4.3.

## 2.3 Improving subgroup robustness with class specificity

We hypothesize that our model underperforms on clinically-relevant subgroups as a result of the target task being underspecified. Since we train our model to only classify whether an EEG clip contains a seizure onset or not, all abnormal patterns and artifacts are grouped together with normal brain activity patterns (in the non-seizure class). As a result, unlike the training protocols of expert EEG readers, our

model does not learn to differentiate among normal activity, abnormal seizure-like activity, and actual seizures, which we hypothesize causes the systematic errors displayed in Figure 3.

To combat task underspecification, we propose to train a multilabel model, where instead of outputing a binary class (seizure or non-seizure), the model identifies multiple attributes from an EEG clip, such as spikes, slowing, and movement. Importantly, since the workflow notes provide these attributes, we are able to train our multilabel model at no additional annotation cost, and training the model to recognize the additional attributes provides class specificity that we hypothesize can improve model performance. To test our hypothesis, we compared the overall and subgroup performances of a model supervised with binary seizure/non-seizure labels, which we will refer to as the binary model, to the same model trained on the same data but trained to classify the 26 attributes (including seizure onset) extracted from workflow notes, which we will refer to as the multilabel model. While the multilabel model outputs probabilities for all 26 attributes, we only consider the probability of seizure onset for evaluation, and calculate the AUROC with respect to the gold-labeled test set for each subgroup.

| | | binary model | multilabel model | p-value |
|---|---|---|---|---|
| | Overall | 85.6 ± 0.9 | **91.5 ± 0.9** | 1.9e-24 |
| patient subgroups | Adults | 89.4 ± 1.1 | **92.7 ± 1.1** | 7.9e-07 |
| | Adults outside ICU | 89.4 ± 1.7 | **91.7 ± 1.7** | 0.036 |
| | Adults from ICU | 88.5 ± 1.2 | **94.0 ± 1.2** | 1.1e-09 |
| | Pediatrics | 82.9 ± 1.5 | **91.2 ± 1.3** | 5.1e-20 |
| seizure subgroup | Focal spike-and-wave | 84.3 ± 2.6 | **92.0 ± 1.7** | 9.0e-06 |
| | Evolving rhythmic slowing | 89.8 ± 3.3 | 93.2 ± 3.0 | 0.10 |
| | Generalized spike-and-wave | 85.5 ± 4.0 | 90.0 ± 4.3 | 0.11 |

Table 1: **Improving subgroup robustness with class specificity:** Increasing class specificity improves overall model performance along with robustness to hidden subgroups. We stratified our evaluation set by patient and seizure subgroups, where the patient subgroups included patients from the adult hospital, pediatric hospital, or adults within or outside the ICU. We report the average AUROC along with 95% confidence intervals. Rows highlighted in blue indicate subgroups that the binary model underperformed on. We estimated the p-value using the DeLong test, which evaluates how statistically significant the improvements of the multilabel model are compared to the binary model.

As shown in Table 1, the multilabel model has significant improvements in both overall performance and subgroup performance (except for 2 of the seizure subgroups). The overall performance improved by 5.9 AUROC points, while the performance on patients from the pediatric hospital improved by 8.3 points, and 7.7 AUROC points for focal spike-and-wave seizure types. Importantly, the improvements in performance significantly minimized the gaps in performance among subgroups. Additionally, we compared the FPRs for each attribute (shown in Supplementary Figure 5) and found that the overall FPR decreased from 0.08 to 0.02. The top 3 attributes with the highest FPR, which correspond to abnormal attributes (mislabeled seizure, unknown abnormality, and slowing), all decreased significantly (e.g., FPR for EEG clips with unknown abnormal patterns decreased from 0.15 to 0.08). We further compared the 2D projected embeddings of the binary and multilabel models in Supplementary Figure 6, which shows that the embeddings of the multlabel model of abnormal EEG clips cluster more tightly than the embeddings of the binary model, reaffirming that the multilabel model can better differentiate EEG abnormalities.

We also investigated the impact of training a multilabel model on different subsets of the workflow attributes on subgroup robustness. We choose two additional subsets of classes: classifying seizures along with two abnormalities highly relevant to seizures (spikes and slowing), and classifying seizures along with only abnormal attributes (i.e., we remove the following attributes: drowsy, jerk, tap, respiration, eyes open/closed, asleep, ekg, arousal). As shown in Supplementary Table 3, we first

found that all multilabel models improved overall seizure detection performance over the binary model. Interestingly, training a multilabel model for detecting seizure onset along with only abnormal attributes performed similarly to the multilabel model trained on all attributes, indicating that increasing class specificity with the abnormal attributes is the most important.

## 2.4  Measuring clinical utility

A major barrier for technicians and neurophysiologists who have access to commercial seizure detection models is the high number of false alarms[13,20], which results in alarm fatigue and in clinicians not utilizing model predictions. Therefore, a good metric to assess clinical utility is the average number of false positives after scanning 24 hours of EEG (FPs/24hr). In particular, we look at two parameters that directly impact the number of false positives:

- Recall (or sensitivity): Specifying the desired recall implicitly determines the threshold used to binarize the seizure probabilities. While having a higher desired recall is advantageous (since we miss fewer seizures), it is in direct tension with false positives, where number of false positives increase as we increase recall. In some settings, such as counting the precise number of occurring seizures, it may be critical to have a high recall. While in other settings, where the model is used as an assistant to prioritize which parts of the EEG to read first, having a high recall is not as critical. For these reasons, we look at the FPs/24hr for a recall of 0.5, 0.8, and 0.9.

- Delay tolerance ($\Delta_t$): we define the delay tolerance to be the maximum amount of time allowed between the actual seizure onset and the predicted seizure onset. In other words, if the time between actual and predicted seizure onset ($T$) is greater than $\Delta_t$, we count the predicted seizure as a false positive; however, if $T < \Delta_t$ then we count the predicted seizure as a true positive. The delay tolerance is an important parameter because not only does it impact how we determine the difference between a true or false positive, but it is also implicitly related to seizure detection latency — the speed in which our model flags seizures. Seizure detection latency may be critical in some settings, for example if we would like to precisely localize the seizure onset region for patients in the epilepsy monitoring unit, it is critical we accurately analyze the EEG near the true onset zone before spreading occurs. In other settings, such as counting number of seizures, seizure detection latency is not a critical parameter. For these reasons, we look at the FPs/24hr for a delay tolerance of 1 minute and 5 minutes.

In Figure 4, we compared the FPs/24hr for 6 different settings while varying recall and delay tolerance, and observed that the multilabel model improved our clinical utility metric by a factor of roughly 2× across all settings.

## 3  Discussion

In this work, we presented a strategy to improve the trustworthiness of seizure detection models by scaling training data and class specificity in a cost-effective manner. Unlike existing techniques that require fellowship-trained neurophysiologists to annotate thousands of EEGs[11], we instead leveraged annotations that provide class specificity and are generated in existing clinical workflows[5], allowing us to scale training data to an unprecedented level of 68,920 EEG hours at no additional annotation cost. In addition to bypassing expert labeling of the training set, workflow notes can also facilitate the ongoing training of healthcare models as additional data are accumulated over time, leading to significant cost savings in terms of upfront and maintenance expenses.
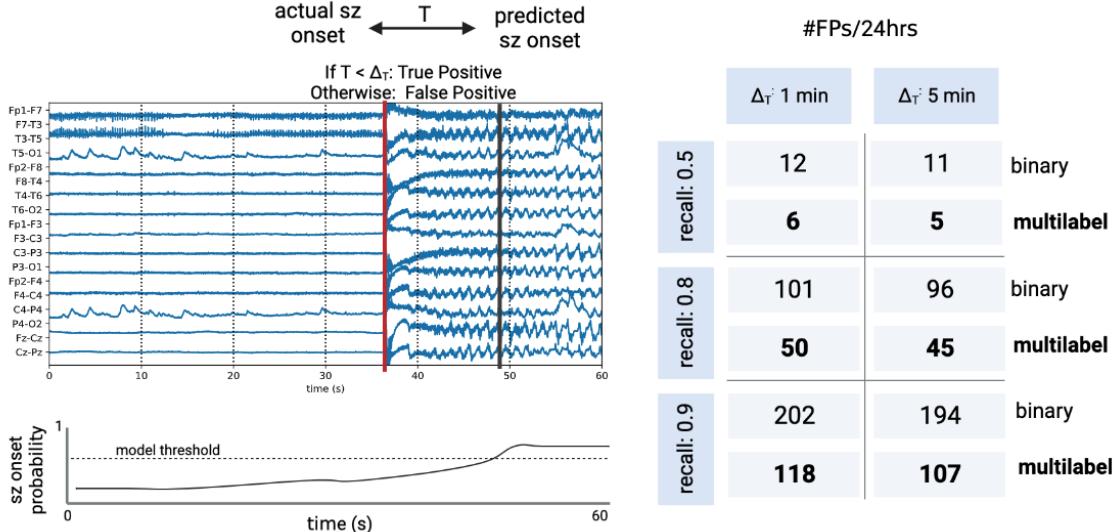
Figure 4: **Clinical utility metric:** On the left, is an EEG in which the red line indicates the actual seizure (sz) onset, and the black line indicates the predicted seizure onset by the model. The time elapsed between the actual and predicted onset is T, which is used to decide whether the predicted seizure onset is a true positive or false positive (depending on the delay tolerance for the clinical setting). The bottom left plot displays the model seizure onset probabilities across time, where the dashed line indicates the model threshold that is determined by the desired recall value. On the right, we compare the average number of false positives that occur every 24 hours of EEG in six different settings: a threshold such that we achieve a recall of 0.5, 0.8, or 0.9, with either of two values of $\Delta_T$, which is a proxy to seizure detection latency (the maximum time between the ground truth and predicted seizure onset we tolerate).

Aside from annotation costs, a major roadblock to successfully deploying healthcare AI is the limited understanding of their errors within hidden subgroups of patients, leading to a lack of trust [12,13,20]. Currently, the gold-standard technique to conduct an in-depth error analysis requires a clinician to manually interpret multiple data samples that the model classified incorrectly and find patterns that combine errors into clinically-relevant subgroups. Instead, we proposed to utilize patient metadata, gold-labeled seizure types, and multiple attributes describing EEG events to stratify the evaluation set and analyze differences in model performance. Apart from the gold-labeled seizure types, we are able to extract the attribute labels from the workflow notes, allowing us to greatly improve performance of our model with no additional costs. From our subgroup error analysis, we found that binary seizure classification models may have large performance gaps among patient age groups (-6.5 AUROC points on pediatrics compared to adults), seizure types (-5.5 AUROC points on focal spike-and-wave versus evolving rhythmic slowing), and has significantly higher false positives (+19 FPR points) for non-seizure EEG clips with abnormal brain activity compared to non-seizure clips. Identifying underperforming subgroups is a critical step in our goal towards trustworthy seizure classification models.

Our core hypothesis is that our binary classification model has high false positives on abnormal non-seizure EEG clips as a result of task underspecification. Since fellows are not only trained to differentiate seizure from non-seizure activity, but also to identify multiple artifacts and abnormalities to rule out seizure [21], we reason that a model should similarly be trained. To reduce high false positive rates and systematic errors, we leveraged attribute labels extracted from workflow notes and trained a multilabel model that learns to classify 26 EEG events such as seizures, spikes, slowing, and movement. We found that such a multilabel model significantly improves overall performance (+5.9 AUROC points), along with closing the performance gap among subgroups, and decreased the false positive rate on abnormal non-seizure clips by 8 FPR points, compared to the binary classification model. We believe this

general direction of increasing the specificity of the supervision task is a promising approach to improve model subgroup robustness. Other successful approaches within this direction include increasing spatial specificity for radiology[17] (e.g., segmentation) and training a chest X-ray model with a comprehensive class ontology[22].

In our investigation of seizure detection models, we also establish a metric of clinical utility. We report the average number of false positives per 24 hours of EEG for different recall and latency settings. We found that across different clinical settings, increasing class specificity reduces the FPs/24hr by a factor of $2\times$, suggesting that our improvements in subgroup robustness may translate to improvements in clinical utility.

Our proposed supervision strategies for improving trustworthiness of seizure detection models have limitations. First, while workflow notes offer a great alternative to manual expert labeling, the resulting labels come from personnel that are instructed to bias their reading to not miss abnormalities since final diagnosis is reviewed by an interpreting physician, which results in false positive labels and sub-optimal supervision. Additionally, our regular expressions to extract labels from the workflow notes may not correctly identify some of the labels, or they may produce errors or not apply to other institutions. Second, while we consider many clinically-relevant subgroups, our analysis can be more comprehensive by including many other important groups such as patient demographics, more seizure types, and finer-grained abnormal events. Third, we do not investigate other important robustness settings that include common distribution shifts, such as different EEG devices and patients from multiple hospitals. Other settings for improving trust may also include proper model calibration, calibration scores, and out-of-distribution detection. We believe it is critical to investigate robustness on a comprehensive list of settings before claiming a model to be trustworthy for deployment.

Future work is needed for improving the robustness of seizure detection models. Further scaling training data to include diverse patients can be done by combining our hospital datasets with existing publicly available datasets such as the TUSZ corpus[3,4]. In a similar spirit, we can utilize publicly available EEG-based models that classify seizures, sleep staging, and brain states[6,23], to either label relevant attributes or enable transfer learning. Another exciting direction is self-supervised and generative AI, where models do not rely on labeled training data to learn useful data representation. For example, recent work has shown that pretraining to forecast EEG signals boosts performance on rare seizure types[19]. We also envision models that generate text reports from EEG[24] may prove to have more robust representations due to learning finer-grained concepts.

In summary, our work provides evidence that scaling training data using labels from workflow notes and increasing class specificity are promising techniques to improve robustness of models to detect seizure onset. We believe that combating robustness challenges through in-depth error analyses, and assessing detection performance of models as well as clinical utility metrics, will be critical to continue improving upon the trustworthiness of AI tools for clinical deployment.

# 4  Methods

## 4.1  Dataset description

Our dataset consists of all EEGs recorded in both the Stanford Hospital and Lucile Packard Children's Hospital from 2006 to 2017. In total, our dataset contains 68,920 EEG hours from 12,297 patients. Our dataset is diverse, where patients span all ages, come from different hospital locations (ICU, epilepsy monitoring unit, and ambulatory), and have different seizure types and etiologies. More details on the statistics of our diverse patients can be found in Figure 2 and Supplementary Figure 2 in Saab et al.[5].

To prepare input data samples from long-form EEG recordings, we segment each recording into

non-overlapping 60 second clips (i.e., stride is 60 seconds). In total, our dataset contained 4,125,225 clips. To ensure consistent information across patients, we only considered the 19 electrodes from the standard 10-20 International EEG configuration, and exclude premature infants or patients with small heads that prevent the full deployment of the 19 electrodes. We further normalize each EEG clip across the temporal dimension using the global average and standard deviations for each channel. Such normalization of input samples is standard practice in deep learning and we find this improves training.

**Gold-standard annotations.** Two fellowship-trained EEG readers (M.T. and C.L.M.) interpreted a randomly selected subset of EEG recordings, annotating for seizure onset. This resulted in an evaluation set of 37,588 60-sec EEG clips (or 626 EEG hours), of which 1,244 clips contain seizures from 395 patients. Patients in the the evaluation set are excluded from the training set. C.L.M. labeled or supervised the labeling of each EEG clip according to the seizure type as defined by EEG ictal patterns; specifically, whether a seizure was a focal spike-and-wave, evolving rhythmic slowing, generalized spike-and-wave, paroxysmal fast acivity, polyspike-and-wave (myocolonic), or electrographically silent, for a subset of 358 patients from the gold-labeled EEGs. However, due to the low frequency of some seizure types, our evaluations only included focal spike-and-wave, evolving rhytmic slowing, and generalized spike-and-wave types (more details can be found in Supplementary Table 2).

**Extracting labels from workflow notes.** Each EEG recording is accompanied with a table of workflow notes with each row indicating an event description along with the event start time. However, the event descriptions are free-form text, and while the workflow annotators use repetitive and standard descriptions, there may be slight deviations. M.T. and C.L.M. analyzed the most common 1,000 event descriptions and by consensus determined a set of unique attributes that (1) are visibly detectable on EEGs, and (2) are typically used when searching for seizures (e.g., common artifacts that must be ruled out such as movement, or other abnormalities such as spikes and slowing). From this manual analysis, we found 26 class attributes of interest. To classify whether one of the 26 class attributes is mentioned in the event description, we produce simple regular expressions that take into account synonyms or acronyms the annotators may use. For example, an annotator may write "seizure", "sz", "spasm", or "absence"; another example is the description of an unknown abnormality, which may simply be indicated by "x", or "xx". We provide a full list of the regular expressions used in Supplementary Table 4.

## 4.2 Model architecture and training

There have been many deep learning model architectures proposed for seizure classification, such as convolutional models (CNNs)[5,25–27], recurrent neural networks (RNNs)[28–30], graphical neural networks (GNNs)[19,31,32], and more[6,9,33–35]. In our work, we study the impact of training data scale and the specificity of the supervision task on seizure classification performance, and not model architecture. However, due to inherent advantages of some architectures, such as simplicity and computational efficiency, we chose S4, a recently proposed convolutional-based model motivated by principles in signal processing[36].

**Deep state space sequence model (S4).** The global architecture of S4 follows a similar deep learning architecture as the transformer encoder, in which each layer is composed of multiple filters, where each filter is a sequence-to-sequence mapping (mixing across time), followed by a non-linear activation function, followed by a linear layer (mixing across filters), and finally a residual connection. The major deviation from the transformer encoder is the sequence-to-sequence filter, which as opposed to an attention mechanism, is a one-dimensional convolutional filter parametrized by linear state-space models (SSMs). An SSM is a fundamental model to represent signals and is ubiqutious across a range of signal

processing and control applications[37,38]. A discrete SSM, which maps observed inputs $u_k$ to hidden states $x_k$, before projecting back to observed outputs $y_k$, has the following recurrent form:

$$x_{k+1} = \boldsymbol{A}x_k + \boldsymbol{B}u_k \tag{1}$$

$$y_k = \boldsymbol{C}x_k + \boldsymbol{D}u_k \tag{2}$$

Where $\boldsymbol{A} \in \mathbb{R}^{d \times d}$, $\boldsymbol{B} \in \mathbb{R}^{d \times 1}$, $\boldsymbol{C} \in \mathbb{R}^{1 \times d}$, and $\boldsymbol{D} \in \mathbb{R}$ are learnable SSM parameters, and $d$ is the dimension of the hidden state $x$. Importantly, we can also compute the SSM as a 1-D convolution, which unlike recurrent models, enables parallelizable inference and training. To see how, if we assume the initial state $x_0 = 0$, and follow equations 1 and 2, we arrive at the following by induction:

$$y_k = \sum_{j=0}^{k-1} \boldsymbol{C}\boldsymbol{A}^{k-1-j}\boldsymbol{B}u_j \tag{3}$$

We can thus compute the output $y_k$ as a 1-D convolution with the following filter:

$$\boldsymbol{F} = (\boldsymbol{C}\boldsymbol{B}, \boldsymbol{C}\boldsymbol{A}\boldsymbol{B}, \boldsymbol{C}\boldsymbol{A}^2\boldsymbol{B}, \dots, \boldsymbol{C}\boldsymbol{A}^{\ell-1}\boldsymbol{B}) \tag{4}$$

$$y_k = (\boldsymbol{F} * \boldsymbol{u})_k \tag{5}$$

Following prior work on sequence model classification[36], we simply use the output squences from the last layer to project from the number of filters to the number of classes (e.g., 2 classes for the binary model and 26 classes for the multilabel model), and perform mean pooling over the temporal dimension before a softmax to output class logits.

There are many advantages of using deep SSMs for long sequence modeling described in recent work[36,39,40]. We highlight the following advantages for EEG modeling: since our EEG clips are of length 12,000, RNNs are slow to train, while CNNs fail to capture long-range dependencies due to limited filter lengths; on the other hand, SSMs are computationally efficient to train (due to their convolutional view), but are also able to capture long-range dependencies with structured initialization of the $\boldsymbol{A}$ matrix. Moreover, we do not need to learn graph structures among the EEG electrodes, which adds an additional layer of complexity in recent state-of-the-art EEG classification models[6]. Nevertheless, to validate that S4 is a well suited model architecture for seizure classification, we compared its performance to other architectures on the public TUSZ benchmark in Supplementary Table 5, and found that S4 is competitive with state-of-the-art models while being more computationally efficient.

**Training details.** We trained all models with the cross-entropy loss using the Adam optimizer in Pytorch[41], with randomly initialized weights. The learning rate was initially set at 0.004 and followed a cosine scheduler[42]. We used a weight decay of 0.1 and a dropout probability of 0.1. Since the training set is very large ($\sim$ 4 million samples) and highly unbalanced with just 0.6% of clips having seizure onset, we used a weighted random sampler with a 25-to-1 bias for positively labeled clips. For more frequent checkpointing, we randomly sampled a maximum of 150,000 clips for each epoch (with replacement), and trained for 200 epochs, while checkpointing on the validation set AUROC. The S4 model architectures had a parameter count of 366k for the binary classification model, and 379k for the multilabel model (due to larger output dimension). The model architecture contained 128 filters per layer for 4 layers with a hidden state dimension $d$ of 64, and the gaussian error linear unit for the non-linear activations. We performed a grid search for the initial learning rate, weight decay, and dropout values using our validation set. We used default values for the other hyperparameters, including model architecture.

## 4.3 Performance metrics

The two main classification metrics used to evaluate seizure classification performance are the the Area Under the Receiver Operating Characteristic curve (AUROC) and the false-positive rate (FPR). We chose the classification threshold such that the class balance of the model predictions matches the ground truth class balance. The ROC curve displays the tradeoff between the True Positive Rate (TPR) and FPR for different classification thresholds. Therefore, the AUROC summarizes the ROC curve in a single scalar value regardless of the specific classification threshold chosen. The FPR and TPR are defined as follows:

$$FPR = \frac{FP}{FP + TN} \tag{6}$$

$$TPR = \frac{TP}{TP + FN} \tag{7}$$

where true-positives (TP) are correct seizure classifications, true-negates (TN) are correct non-seizure classifications, false-positives (FP) are incorrect seizure classifications, and false-negatives (FN) are incorrect non-seizure classifications. To calculare 95% confidence intervals and p-values when comparing the AUROC of two models, we used the DeLong test[43].

# Data Availability

The Stanford clinical datasets used in this study are subject to restrictions regarding the availability of Protected Health Information. They were accessed with approval from the Institutional Review Board solely for the purpose of this specific study and are not accessible to the public.

# Code Availability

The code used to generate the main results in this manuscript can be found in the following github repository: `https://github.com/khaledsaab/eeg_robustness`.

# References

1. Donald L Schomer and Fernando Lopes Da Silva. *Niedermeyer's electroencephalography: basic principles, clinical applications, and related fields.* Lippincott Williams & Wilkins, 2012.

2. Jan Brogger, Tom Eichele, Eivind Aanestad, Henning Olberg, Ina Hjelland, and Harald Aurlien. Visual eeg reviewing times with score eeg. *Clinical Neurophysiology Practice*, 3:59–64, 2018.

3. Iyad Obeid and Joseph Picone. The temple university hospital eeg data corpus. *Frontiers in neuroscience*, 10:196, 2016.

4. Vinit Shah, Eva Von Weltin, Silvia Lopez, James Riley McHugh, Lillian Veloso, Meysam Golmohammadi, Iyad Obeid, and Joseph Picone. The temple university hospital seizure detection corpus. *Frontiers in neuroinformatics*, 12:83, 2018.

5. Khaled Saab, Jared Dunnmon, Christopher Ré, Daniel Rubin, and Christopher Lee-Messer. Weak supervision as an efficient approach for automated seizure detection in electroencephalography. *NPJ digital medicine*, 3(1):59, 2020.

6. Siyi Tang, Jared A Dunnmon, Liangqiong Qu, Khaled K Saab, Christopher Lee-Messer, and Daniel L Rubin. Spatiotemporal modeling of multivariate signals with graph neural networks and structured state space models. *arXiv preprint arXiv:2211.11176*, 2022.

7. Yang Li, Yu Liu, Wei-Gang Cui, Yu-Zhu Guo, Hui Huang, and Zhong-Yi Hu. Epileptic seizure detection in eeg signals using a unified temporal-spectral squeeze-and-excitation network. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 28(4):782–794, 2020.

8. Punnawish Thuwajit, Phurin Rangpong, Phattarapong Sawangjai, Phairot Autthasan, Rattanaphon Chaisaen, Nannapas Banluesombatkul, Puttaranun Boonchit, Nattasate Tatsaringkansakul, Thapanun Sudhawiyangkul, and Theerawit Wilaiprasitporn. Eegwavenet: Multiscale cnn-based spatiotemporal feature extraction for eeg seizure detection. *IEEE Transactions on Industrial Informatics*, 18 (8):5547–5557, 2021.

9. David Ahmedt-Aristizabal, Tharindu Fernando, Simon Denman, Lars Petersson, Matthew J Aburn, and Clinton Fookes. Neural memory networks for seizure type classification. In *2020 42nd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, pages 569–575. IEEE, 2020.

10. Meysam Golmohammadi, Amir Hossein Harati Nejad Torbati, Silvia Lopez de Diego, Iyad Obeid, and Joseph Picone. Automatic analysis of eegs using big data and hybrid deep learning architectures. *Frontiers in human neuroscience*, 13:76, 2019.

11. Jin Jing, Wendong Ge, Shenda Hong, Marta Bento Fernandes, Zhen Lin, Chaoqi Yang, Sungtae An, Aaron F Struck, Aline Herlopian, Ioannis Karakis, et al. Development of expert-level classification of seizures and rhythmic and periodic patterns during eeg interpretation. *Neurology*, 2023.

12. Jenna Wiens, Suchi Saria, Mark Sendak, Marzyeh Ghassemi, Vincent X Liu, Finale Doshi-Velez, Kenneth Jung, Katherine Heller, David Kale, Mohammed Saeed, et al. Do no harm: a roadmap for responsible machine learning for health care. *Nature medicine*, 25(9):1337–1340, 2019.

13. EEM Reus, GH Visser, MPJ Sommers-Spijkerman, JG van Dijk, and FME Cox. Automated spike and seizure detection: are we ready for implementation? *Seizure*, 2023.

14. Alex J DeGrave, Joseph D Janizek, and Su-In Lee. Ai for radiographic covid-19 detection selects shortcuts over signal. *Nature Machine Intelligence*, 3(7):610–619, 2021.

15. Marcus A Badgeley, John R Zech, Luke Oakden-Rayner, Benjamin S Glicksberg, Manway Liu, William Gale, Michael V McConnell, Bethany Percha, Thomas M Snyder, and Joel T Dudley. Deep learning predicts hip fracture using confounding patient and healthcare variables. *NPJ digital medicine*, 2(1):31, 2019.

16. Luke Oakden-Rayner, Jared Dunnmon, Gustavo Carneiro, and Christopher Ré. Hidden stratification causes clinically meaningful failures in machine learning for medical imaging. In *Proceedings of the ACM conference on health, inference, and learning*, pages 151–159, 2020.

17. Khaled Saab, Sarah Hooper, Mayee Chen, Michael Zhang, Daniel Rubin, and Christopher Ré. Reducing reliance on spurious features in medical image classification with spatial specificity. *Machine learning for healthcare*, 2022.

18. John R Zech, Marcus A Badgeley, Manway Liu, Anthony B Costa, Joseph J Titano, and Eric Karl Oermann. Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: a cross-sectional study. *PLoS medicine*, 15(11):e1002683, 2018.

19. Siyi Tang, Jared A Dunnmon, Khaled Saab, Xuan Zhang, Qianying Huang, Florian Dubost, Daniel L Rubin, and Christopher Lee-Messer. Self-supervised graph neural networks for improved electroencephalographic seizure analysis. *International Conference on Learning Representations*, 2022.

20. Andreea M Pavel, Janet M Rennie, Linda S de Vries, Mats Blennow, Adrienne Foran, Divyen K Shah, Ronit M Pressler, Olga Kapellou, Eugene M Dempsey, Sean R Mathieson, et al. A machine-learning algorithm for neonatal seizure recognition: a multicentre, randomised, controlled trial. *The Lancet Child & Adolescent Health*, 4(10):740–749, 2020.

21. William O Tatum IV. *Handbook of EEG interpretation*. Springer Publishing Company, 2021.

22. Jarrel CY Seah, Cyril HM Tang, Quinlan D Buchlak, Xavier G Holt, Jeffrey B Wardman, Anuar Aimoldin, Nazanin Esmaili, Hassan Ahmad, Hung Pham, John F Lambert, et al. Effect of a comprehensive deep-learning model on the accuracy of chest x-ray interpretation by radiologists: a retrospective, multireader multicase study. *The Lancet Digital Health*, 3(8):e496–e506, 2021.

23. Pouya Bashivan, Irina Rish, Mohammed Yeasin, and Noel Codella. Learning representations from eeg with deep recurrent-convolutional neural networks. *arXiv preprint arXiv:1511.06448*, 2015.

24. Siddharth Biswal, Cao Xiao, M Brandon Westover, and Jimeng Sun. Eegtotext: learning to write medical reports from eeg recordings. In *Machine Learning for Healthcare Conference*, pages 513–531. PMLR, 2019.

25. Alison O'Shea, Gordon Lightbody, Geraldine Boylan, and Andriy Temko. Neonatal seizure detection from raw multi-channel eeg using a fully convolutional architecture. *Neural Networks*, 123:12–25, 2020.

26. Shivarudhrappa Raghu, Natarajan Sriraam, Yasin Temel, Shyam Vasudeva Rao, and Pieter L Kubben. Eeg based multi-class seizure type classification using convolutional neural network and transfer learning. *Neural Networks*, 124:202–212, 2020.

27. Tomas Iešmantas and Robertas Alzbutas. Convolutional neural network for detection and classification of seizures in clinical data. *Medical & Biological Engineering & Computing*, 58:1919–1932, 2020.

28. Lasitha Vidyaratne, Alexander Glandon, Mahbubul Alam, and Khan M Iftekharuddin. Deep recurrent neural network for seizure detection. In *2016 International Joint Conference on Neural Networks (IJCNN)*, pages 1202–1207. IEEE, 2016.

29. Meysam Golmohammadi, Saeedeh Ziyabari, Vinit Shah, Eva Von Weltin, Christopher Campbell, Iyad Obeid, and Joseph Picone. Gated recurrent networks for seizure detection. In *2017 IEEE Signal Processing in Medicine and Biology Symposium (SPMB)*, pages 1–5. IEEE, 2017.

30. Ibrahim Aliyu, Yong Beom Lim, and Chang Gyoon Lim. Epilepsy detection in eeg signal using recurrent neural network. In *Proceedings of the 2019 3rd International Conference on Intelligent Systems, Metaheuristics & Swarm Intelligence*, pages 50–53, 2019.

31. Yogatheesan Varatharajah, Min Jin Chong, Krishnakant Saboo, Brent Berry, Benjamin Brinkmann, Gregory Worrell, and Ravishankar Iyer. Eeg-graph: a factor-graph-based model for capturing spatial, temporal, and observational relationships in electroencephalograms. *Advances in neural information processing systems*, 30, 2017.

32. Khuong Vo, Manoj Vishwanath, Ramesh Srinivasan, Nikil Dutt, and Hung Cao. Composing graphical models with generative adversarial networks for eeg signal modeling. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1231–1235. IEEE, 2022.

33. Khansa Rasheed, Adnan Qayyum, Junaid Qadir, Shobi Sivathamboo, Patrick Kwan, Levin Kuhlmann, Terence O'Brien, and Adeel Razi. Machine learning for predicting epileptic seizures using eeg signals: A review. *IEEE Reviews in Biomedical Engineering*, 14:139–155, 2020.

34. Mohammad Khubeb Siddiqui, Ruben Morales-Menendez, Xiaodi Huang, and Nasir Hussain. A review of epileptic seizure detection using machine learning classifiers. *Brain informatics*, 7(1):1–18, 2020.

35. Umar Asif, Subhrajit Roy, Jianbin Tang, and Stefan Harrer. Seizurenet: Multi-spectral deep feature learning for seizure type classification. In *Machine Learning in Clinical Neuroimaging and Radiogenomics in Neuro-oncology: Third International Workshop, MLCN 2020, and Second International Workshop, RNO-AI 2020, Held in Conjunction with MICCAI 2020, Lima, Peru, October 4–8, 2020, Proceedings 3*, pages 77–87. Springer, 2020.

36. Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021.

37. Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. 1960.

38. James D Hamilton. State-space models. *Handbook of econometrics*, 4:3039–3080, 1994.

39. Michael Zhang, Khaled K Saab, Michael Poli, Tri Dao, Karan Goel, and Christopher Ré. Effectively modeling time series with simple discrete state spaces. *arXiv preprint arXiv:2303.09489*, 2023.

40. Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher Ré. Combining recurrent, convolutional, and continuous-time models with linear state space layers. *Advances in neural information processing systems*, 34:572–585, 2021.

41. Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

42. Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.

43. Elizabeth R DeLong, David M DeLong, and Daniel L Clarke-Pearson. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics*, pages 837–845, 1988.

# Acknowledgements

# Supplement

| seizure type | seizure count |
|---|---|
| focal spike-and-wave | 181 |
| evolving rhythmic slowing | 81 |
| generalized spike-and-wave | 53 |
| polyspike-and-wave (myoclonic) | 22 |
| paroxysmal fast activity | 17 |
| fast spiking | 7 |
| sz without clear electrographic change | 2 |
| electrographically silent | 1 |

Table 2: **Seizure types:** Seizure counts for different seizure types as defined by EEG ictal patterns in a subset of our gold-labeled evaluation set.

| Subgroups | sz only | sz / spikes / slowing | sz / all abnormal attributes | sz / all attributes |
|---|---|---|---|---|
| Overall | 85.6 ± 0.9 | **88.4 ± 1.0** | **91.4 ± 0.9** | **91.5 ± 0.9** |
| Adults | 89.4 ± 1.1 | 89.6 ± 1.5 | **92.8 ± 1.0** | **92.7 ± 1.1** |
| Adults w/o ICU | 89.4 ± 1.7 | 88.7 ± 2.5 | 91.1 ± 1.7 | **91.7 ± 1.7** |
| Adults w/ ICU | 88.5 ± 1.2 | **91.2 ± 1.5** | **94.7 ± 1.1** | **94.0 ± 1.2** |
| Pediatrics | 82.9 ± 1.5 | **87.8 ± 1.5** | **90.7 ± 1.4** | **91.2 ± 1.3** |
| Focal | 84.3 ± 2.6 | **89.4 ± 2.4** | **90.9 ± 2.0** | **92.0 ± 1.7** |
| Evolving slow | 89.8 ± 3.3 | 91.4 ± 4.1 | **94.8 ± 2.3** | 93.2 ± 3.0 |
| Generalized | 85.5 ± 4.0 | 73.7 ± 7.5 | 89.9 ± 3.5 | 90.0 ± 4.3 |

Table 3: **Subgroup robustness:** Increasing task specificity improves overall model performance along with robustness to hidden subgroups. We stratify our evaluation set by patient and seizure subgroups, where the patient subgroups include patients from the adult hospital, children hospital, or adults within or outside the ICU. We report the average AUROC for the multilabel seizure detection model along with 95% confidence intervals. Bolded numbers indicate statistically signficant lifts over the binary classification model. "sz" stands for seizure.
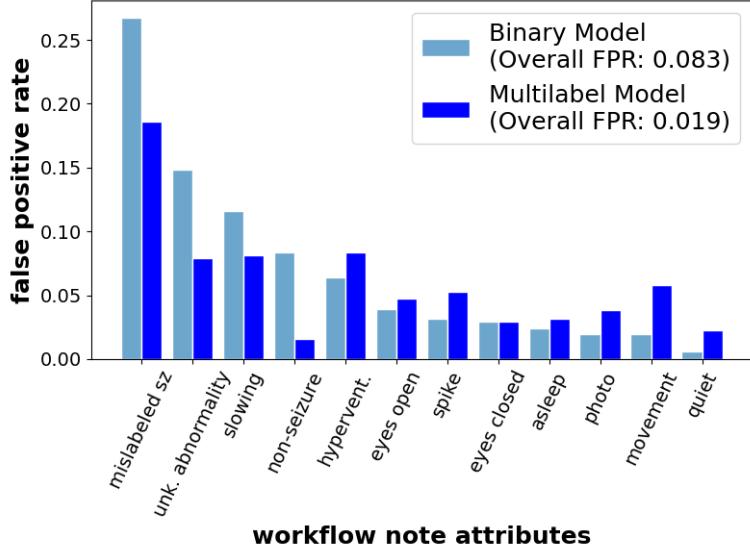
Figure 5: **FPR attribute analysis:** False positive rates with respect to seizure detection across subsets of our evaluation set stratified by the workflow attributes for the binary and multilabel model. The seizure detection threshold was chosen such that the class balance of the model predictions matched the ground truth class balance.
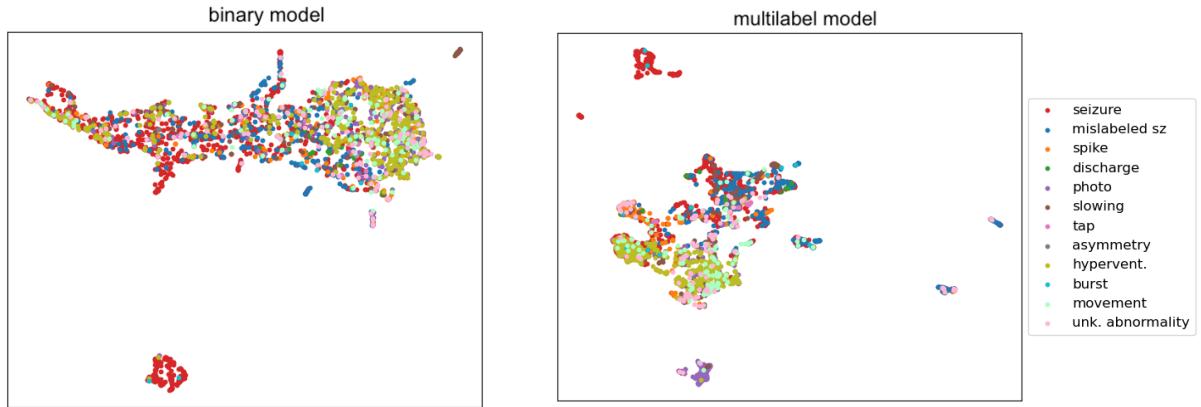


Figure 6: **Model embedding analysis:** UMap-projected embeddings show that the multilabel model embeddings cluster the abnormal attributes (mislabeled sz, spike, slowing, unknown abnormality) more tightly compared to the binary model embeddings, reaffirming that the multilabel model has learned to more effectively differentiate seizures from other EEG abnormalities.

| attribute | regular expression |
|---|---|
| seizure | "seizure \| sz \| absence \| spasm" |
| spike | "spike" |
| slowing | "slow" |
| photoelectric stimulation | "photo" |
| stimulation | "stim" |
| posterior dominant rhythm | "pdr" |
| unknown abnormality | "^x*$" |
| movement artifact | "movement \| mvt" |
| EKG artifact | "ekg" |
| discharge | "discharge \| discharges" |
| tapping artifact | "tap" |
| hyperventilation | "hv" |
| jerking | "jerk" |
| drowsy | "drowsy" |
| asymmetry | "asymmetry" |
| arousal | "arousal" |
| respiration | "respiration" |
| asleep | "asleep \| sleep" |
| awake | "awake" |
| burst | "burst" |
| quiet | "quiet" |
| suspicion in left hemisphere | "^L*$" |
| suspicion in right hemisphere | "^R*$" |
| eyes closed | "eyes closed" |
| eyes opened | "eyes opened" |

Table 4: **EEG attributes and regular expressions to extract EEG attributes from workflow notes:** For each regular expression, we turn off the case sensitivity flag.

| Model | LSTM | CNN-LSTM | Dense-CNN | DCRNN | Graphs4mer | S4 |
|---|---|---|---|---|---|---|
| TUH | $71.5 \pm 1.6$ | $68.2 \pm 0.3$ | $79.6 \pm 1.4$ | $80.4 \pm 1.5$ | $90.6 \pm 1.2$ | $87.7 \pm 1.1$ |

Table 5: Architecture comparisons on TUH v1.5.2[4] test set (AUROC). Our chosen architecture (S4) is competitive with SoTA seizure detection methods. Performance of baseline models (first five columns) are taken from Tang et al., 2022[6].