

# Clinical translation of machine learning algorithms for seizure detection in scalp electroencephalography: a systematic review

Nina Moutonnet<sup>1</sup>, Steven White<sup>4</sup>, Benjamin P Campbell<sup>5</sup>, Danilo Mandic<sup>6,7</sup>, and Gregory Scott<sup>2,7</sup>

<sup>1</sup>Department of Computing, Imperial College London

<sup>2</sup>Department of Brain Sciences, Imperial College London

<sup>4</sup>Department of Clinical Neurophysiology, National Hospital for Neurology & Neurosurgery, London, United Kingdom

<sup>5</sup>Department of Bioengineering, Imperial College London

<sup>6</sup>Department of Electrical and Electronic Engineering, Imperial College London

<sup>7</sup>UK Dementia Research Institute

## Abstract

Machine learning algorithms for seizure detection have shown great diagnostic potential, with recent reported accuracies reaching 100%. However, few published algorithms have fully addressed the requirements for successful clinical translation. For example, the properties of training data may critically limit the generalisability of algorithms, algorithms may be sensitive to variability across EEG acquisition hardware, and run-time processing costs may render them unfeasible for real-time clinical use cases. Here, we systematically review machine learning seizure detection algorithms with a focus on clinical translatability, assessed by criteria including generalisability, run-time costs, explainability, and clinically-relevant performance metrics. For non-specialists, we provide domain-specific knowledge necessary to contextualise model development and evaluation. Our critical evaluation of machine learning algorithms with respect to their potential real-world effectiveness can help accelerate clinical translation and identify gaps in the current seizure detection literature.

## Corresponding author

Dr Gregory Scott BEng MSc MBBS MRCP PhD Post-Doctoral, Post-CCT Research Fellow Honorary Consultant Neurologist UK DRI Care Research and Technology Centre 9th Floor, Sir Michael Uren Hub, Imperial College London, 86 Wood Lane, London. W12 0BZ. UK. Tel: +44 (0)7909 691484 Email: gregory.scott99@imperial.ac.uk

## Keywords




seizure detection; epilepsy; scalp electroencephalography; machine learning; deep learning; clinical translation

# 1 Introduction

A seizure is an abnormal synchronous excitation of one or more populations of neurons in the brain. Seizures are not a rare phenomenon, with a lifetime incidence of 2-5% [1]. Seizures may occur in a range of clinical contexts, in particular in patients with *epilepsy* - defined as a tendency to unprovoked seizures [2]. Other, provoking, causes of seizures include central nervous system (CNS) infections, metabolic abnormalities, traumatic brain injuries, and drug toxicity. Seizure symptoms vary widely, from abnormal sensations to convulsions and altered awareness. Seizures are usually self-limiting, in that they typically last less than two minutes. Occasionally, they may continue for more than five minutes and/or recur without full recovery, a state termed *status epilepticus*, a medical emergency [3].

The accurate and timely detection of seizures is an important healthcare challenge [4]. Seizure detection is straightforward when the seizure activity has a well-recognised clinical correlate like generalised convulsions or a subjective alteration in experience, termed an *aura*. However, seizures without clear symptoms may be missed or mistaken for other clinical phenomena. For example, ongoing seizure activity without a clinically obvious motor component, known as *non-convulsive status epilepticus* (NCSE), affects 8-20% of patients in intensive care units (ICUs) and can be fatal [5]. Despite being a medical emergency, it is critically underdiagnosed because it can go unnoticed by clinicians [6]. Patients with established epilepsy frequently under-report their seizures, as they may occur during sleep, or impair patients' awareness and memory of the events [4]. In contrast, patients with *non-epileptic attacks* (NEAs) - events appearing superficially similar to seizures but without the associated abnormal brain activity - are at risk of misdiagnosis and inappropriate treatment [7].

Consequently, there is a need for an automated and accurate method for seizure detection. Accelerometry, electromyography (EMG), and electroencephalography (EEG) have all been used for seizure detection [8, 4]. Of these methods, scalp EEG is advantageous, owing to its ability to detect subtle electrical abnormalities even in the absence of clinical correlates. For this reason it will be the focus of the rest of this paper.

Clinical context	Use cases	Specific challenges	Modality and duration
 Ward	Detection of seizures and status epilepticus, monitoring	Patient generalisable, easy to use, robust to noise, minimise false positive rate given abnormal brain activity	Scalp EEG recordings lasting from minutes to days
 Telemetry unit	Seizure detection, epilepsy diagnosis, pre-surgical planning	Robustness to seizure types, spatial localisation of onset and spread	Scalp or intracranial EEG recordings lasting days
 Community	Seizure detection and epilepsy diagnosis, patient safety alarm	Ambulatory patients, simple hardware, lack of clinical correlate	Scalp, intracranial, sub-cutaneous EEG recordings lasting days to months

**Figure 1. Potential use cases, challenges and recording modality for automated scalp EEG seizure detection [9].** Applications for seizure detection algorithms range widely, from (1) highlighting to clinicians sections of interest in long recordings to facilitate annotation offline; (2) real-time seizure detection using continuous EEG in ambulatory patients, telemetry units, or ICU; to (3) automatic recording of seizure diary entries.

EEG recorded from the scalp is used to non-invasively measure the electrical activity of the brain. Traditionally, scalp EEG traces are examined by a neurophysiologist who diagnoses seizure-related

disorders [10]. However, resource constraints as well as significant inter-rater variability currently limit both the availability and accuracy of clinical EEG [11, 12, 13]. Interpretation of EEG by a neurophysiologist is time consuming and costly. A health economics analysis concluded that continuous EEG monitoring accounted for an average of 5% of total hospital charges for patients monitored in intensive care [14]. In practice, EEG monitoring is rationed because there are not enough neurophysiologists to provide interpretation [15, 11, 16]. Automated EEG-based seizure detection algorithms could therefore improve the availability of clinical EEG, and therefore the care of patients with seizures.

Despite some successes, there are several challenges for automated seizure detection from scalp EEG. We have identified three main challenges: (i) data complexity, (ii) seizure definition, and (iii) data collection and labelling discrepancies. These challenges potentially impose an upper limit on the accuracy of seizure detection algorithms.

Firstly, scalp EEG data is multi-dimensional, non-stationary, imbalanced, and often corrupted by artefacts (e.g. muscle activity, eye movements). Combining this data complexity with the multiple manifestations and symptoms of seizures results in high variability of seizure EEG recordings. Without consideration of this variability, an algorithm designed for one application may fail if applied to a distinct context or patient population. Furthermore, seizures with a deep brain origin, such as the mesial temporal lobes, may be undetectable by standard interpretation of scalp EEG (so-called *scalp-negative* seizures [17]).

Secondly, there is no single benchmark for identifying a seizure. The nosology of epilepsy and seizures is an active area of debate in epileptology, with many contrasting systems of classification [18, 19, 20]. This is particularly relevant for the definition of *status epilepticus* and its subtypes [21, 20]. Another complication is the possibility of seizures being described as a non-binary phenomenon [22].

Thirdly, in manual clinical EEG interpretation and labelling practice there are significant discrepancies between what experts consider to be seizure activity. This results in labeling differences, in particular regarding the timings of seizure onset and offset [13, 23, 24]. There are a range of focal non-seizure-related EEG waveform abnormalities that may be observed clinically, where some may be mistaken for unambiguous seizure phenomena, and the clinical significance of others is still debated [25, 26]. These sources of variability is in addition to any variability arising from differences in EEG hardware and recording configurations.

The application of machine learning (ML) to seizure detection has shown great potential, with reported accuracies reaching 100% [27]. However, few published ML algorithms have fully addressed the three challenges outlined above, thereby reducing generalisability and clinical viability. Furthermore, patient-specific and cross-patient use cases are frequently conflated. Many algorithms are trained and evaluated on a single patient’s data (e.g. [28, 29, 30, 31]). Patient-specific approaches may achieve higher classification performance compared to cross-patient algorithms. However, the real-world use cases for patient specific seizure detection algorithms are far more limited than for algorithms that can generalise across patients (see Figure 1). For instance, patient-specific algorithms are well suited to long-term ambulatory seizure monitoring, but an algorithm used in the ICU should be capable of generalising across patients.

## 1.1 Structure of this review

In this article, we first provide key domain-specific knowledge about EEG recordings and seizures for non-specialists (Section 2). Then, we describe the key properties of publicly available datasets for seizure detection (Section 3). This is followed by an outline of the PRISMA guidelines used to select articles (Section 4). We then review the selected automated seizure detection algorithms, from data pre-processing to performance metrics (Section 5). Finally, we highlight potential research gaps and challenges (Section 6), alongside providing guidelines for the development of future EEG-based seizure detection algorithms (Section 7).

## 2 Domain-specific knowledge

### 2.1 EEG recording techniques

EEG records the electrical activity of the brain using electrodes. Whilst the scale and complexity of neuronal activity make microscopic brain dynamics difficult to measure, the currents generated by populations of neurons have an amplitude sufficiently high to record the potential between electrodes [32].

The electrodes can be placed on the scalp, or they can be implanted in the brain or below the skin. These two kinds of EEG recording arrangements are known as *extracranial* and *intracranial* techniques, respectively. Intracranial EEG can be recorded using methods such as *electrocorticography* (ECoG) or *stereotaxic EEG*. The former uses strips or grids of 50-100 electrodes implanted over an area of cortex via a *craniotomy* procedure to produce a skull window [33]. The latter uses 5-15 wired electrodes penetrating the brain via multiple small *burr holes* in the skull [34], potentially targeting deeper structures like the hippocampus.

Intracranial EEG is advantageous because electrodes can directly sample the brain regions of interest, whereas activity recorded extracranially consists of a superposition of various brain sources. Additionally, intracranial electrodes have a much higher signal-to-noise ratio than scalp electrodes and suffer from fewer artefacts (e.g. muscle activity) [35]. Consequently, algorithms trained on intracranial recording usually provide more accurate seizure detection than those trained on extracranial recordings [36]. However, fitting intracranial electrodes is costly and comes with considerable medical risk. This makes intracranial EEG recording relatively rare and generally confined to epilepsy surgery planning (see Figure 1). The locations of intracranial EEG electrodes are tailored to an individual patient’s case rather than defined by a standardised system, and since most surgical patients have seizures originating from temporal or frontal lobes, it is rare to find published intracranial EEG data from non-temporal and non-frontal sites. For all these reasons, we do not address ML algorithms for intracranial EEG in this review.

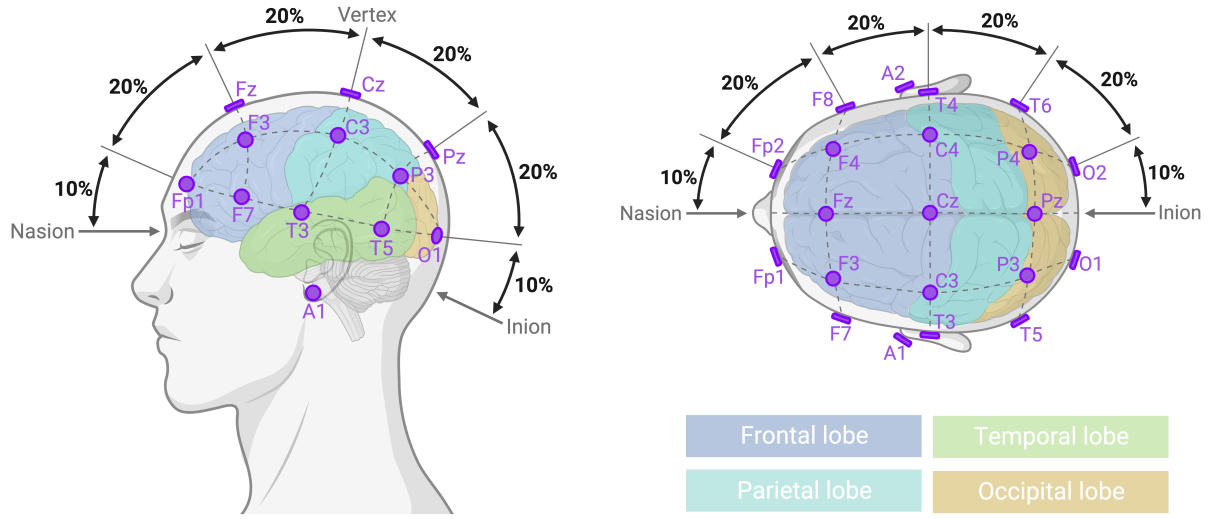
Another recording technique not discussed in this review is minimally invasive *sub-cutaneous EEG*. This has recently been developed to provide comparable signal quality to scalp EEG and allow high duration continuous recordings in ambulatory patients with seizures [37].

### 2.2 Scalp EEG

Scalp EEG is by far the most common clinical technique for recording brain activity in seizure-related disorders. In this review, we focus on scalp EEG recordings and ML algorithms developed for this modality.

Scalp EEG activity is primarily generated by cortical pyramidal neurons that are oriented perpendicularly to the surface of the brain. The amplitude of the EEG is about 10-100  $\mu\text{V}$  when measured on the scalp, or 1-2 mV when measured intracranially. Clinical EEG is usually sampled at a frequency of 256–512 Hz.

Scalp EEG typically uses  $\approx 20$  electrodes affixed across the skin in standardised locations typically using a cap of Ag-AgCl electrodes. The international system of electrode placement and naming is the *10-20 system*, which refers to the distances between adjacent electrodes being either 10% or 20% of the total distance across the skull, based on anatomical landmarks (see Figure 2).



**Figure 2. 10-20 electrode placement system with front-back (nasion to inion) 10% and 20% electrode distances.** The number of electrodes used for EEG recording varies. The spatial resolution of scalp EEG setups can range from 14 channels (low resolution) to 256 channels (high resolution). According to the standardised EEG electrode naming convention, each electrode has a name, composed of a character followed by a number. This combination describes the location of each electrode. The characters are Fp for frontal-polar, F for frontal, P for parietal, T for temporal, O for occipital, and C for the central region of the brain. Odd-numbered electrodes are on the left side of the brain, even-numbered electrodes on the right side, and Z electrodes along the midline of the scalp. [38] [39]

Voltage is an electric potential difference between two points as opposed to an absolute quantity. The electrodes can therefore be connected using different EEG *montages*. The term EEG montage refers to a specific arrangement of scalp electrodes as well as how the electrodes are connected to one another. In a *monopolar montage*, each electrode is connected to a common reference. This can be recorded by a reference electrode on the scalp or body, such as an auricular (ear) electrode. Alternatively, the reference can be obtained by averaging the recorded potential of all electrodes, called *average reference*. [40, 41]. In a *bipolar montage*, each electrode is connected to one or two neighbouring electrodes [42] (see Figure 3). Different montages are sensitive to different types of brain activity, thus producing different waveforms. It is possible to re-reference the data post recording with simple arithmetic steps which can be useful to bring new dynamics of interest to light.

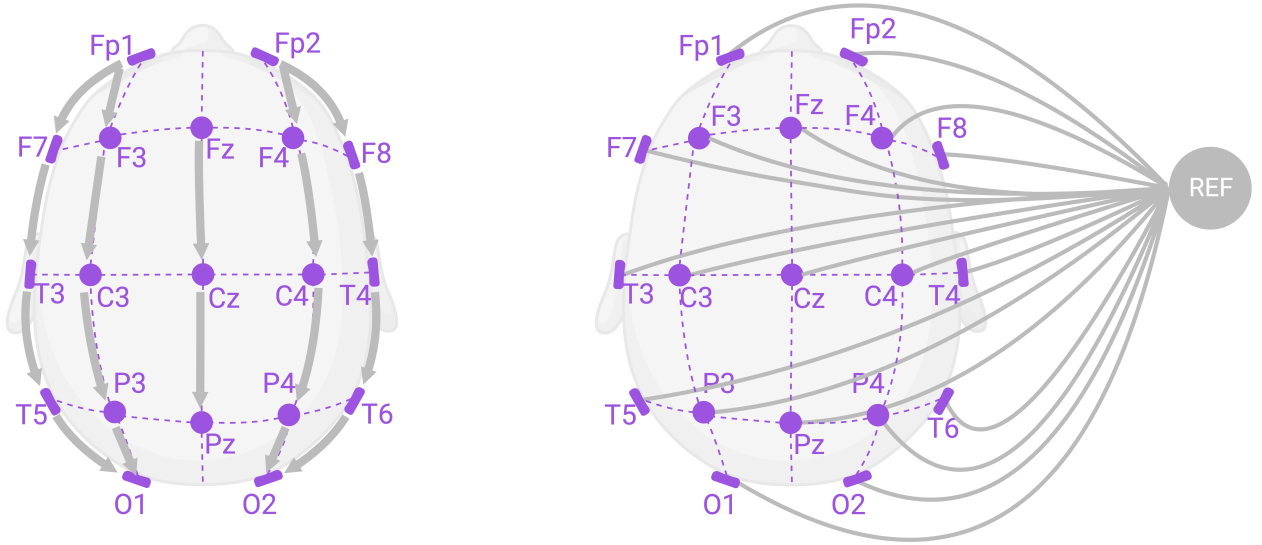
The duration of EEG recordings for the diagnosis and management of seizure-related disorders varies from minutes to days (see Figure 1). In practice, EEG recordings are often combined with video footage of the patient to provide information about their clinical status and behaviour, although video is not often available in public datasets. Video-EEG monitoring is typically used for diagnostic clarification of seizure-related disorders and for pre-surgical planning purposes [43].

The context of EEG recording not only affects the inclusion of video data but also the duration and technical properties of the EEG (see Figure 1). This includes variables such as the number of

channels, the likelihood of seizure occurrence, and the nature of artefacts. For example, interference of muscle activity may corrupt the EEG of freely-moving awake patients more than that of motionless and comatose patients. The majority of clinical EEGs are acquired in a hospital, either on in-patients (often on an intensive care unit), patients undergoing video-EEG monitoring (in a video telemetry unit), or out-patients (a *routine* EEG). Rarely, patients can be in the community and wear a portable *ambulatory* EEG [44].

### 2.3 Canonical EEG frequency bands

The analysis of EEG has traditionally focused on five canonical EEG frequency bands that have been identified. The frequency bands widely used are:  $\delta$  (0.5–4 Hz),  $\theta$  (4–8 Hz),  $\alpha$  (8–13 Hz),  $\beta$  (13–30 Hz), and  $\gamma$  (30–80 Hz). Correlation of these bands with different behaviours has been established. Briefly,  $\delta$  waves are observed in deep sleep;  $\theta$  waves are linked to states of relaxation;  $\alpha$  waves are prominent in eyes-closed restfulness;  $\beta$  waves are associated with mental effort; finally,  $\gamma$  waves are believed to play an important role in perception and conscious awareness [45]. Note that the boundary values of those five frequency bands vary across different studies. For instance, while most sources define  $\gamma$  as 30–80 Hz, some define it as 30–100 Hz or even (30–200 Hz).



**Figure 3. Bipolar and monopolar EEG montage.** (Left) Double banana bipolar montage, where each electrode is referenced by the one behind it. There is an outside temporal chain and an inside parasagittal chain on each side of the scalp, and a unique central chain in the middle. (Right) Monopolar montage, where all the electrodes are referenced by a single point. [39]

### 2.4 Seizure manifestations in scalp EEG

Domain knowledge of the different manifestations of seizures is important for anyone developing ML algorithms for seizure detection. Failure to incorporate this knowledge is likely to limit the robustness and generalisability of the developed algorithms. The following section summarises typical seizure characteristics for non-specialists.

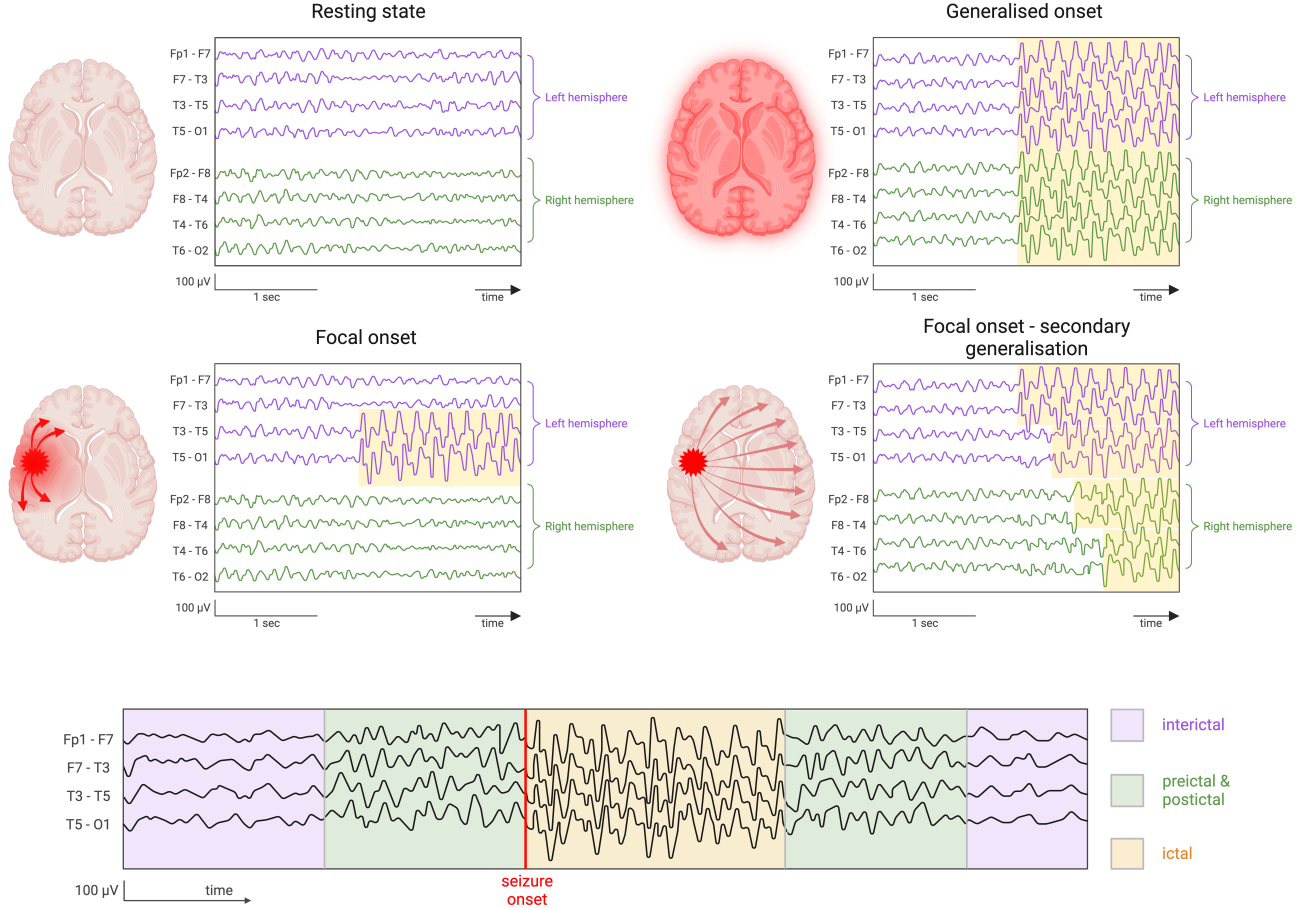


### 2.4.1 Seizure types, duration, and evolution

A seizure is characterised by abnormal and excessive electrical activity from one or more populations of cortical neurons. Seizures are broadly classified into two categories: *focal* (or partial) seizures and *generalised* seizures. Focal seizures involve specific cortical region(s) confined to one hemisphere, while generalised seizures involve widespread electrical activity involving both hemispheres (see Figure 4).

Seizures vary in duration. Most seizures last from a few seconds to two minutes, and are self limiting. The rate of seizures within an individual varies widely (e.g. from  $<1$  seizure/year to  $>5$ /day) [46]. Rarely, seizures continue for more than five minutes or occur recurrently without recovery of normal consciousness in between, termed *status epilepticus*.

A seizure can spatially evolve over time. The starting location is known as the *seizure onset zone (SOZ)*. A *focal-onset* seizure is one that begins in specific brain region(s) in one hemisphere. A *generalised-onset* seizure begins with widespread activity in both hemispheres. A focal-onset seizure with *secondary generalisation* spreads to both hemispheres (see Figure 4). The spread of a seizure can happen on the order of seconds or minutes. Typical patterns of the spatial spread is outlined in [47].



**Figure 4. Main seizure types and some EEG characteristics.** (Top) Normal brain activity, focal seizure and focal onset seizure with secondary generalisation alongside their EEG correlate. (Bottom) Nomenclature of seizure phases including demonstrative EEG segments of inter-ictal, pre-ictal, ictal and post-ictal activity. [48]

### 2.4.2 Seizure-related phases

Every seizure has four periods or phases, which have distinct clinical and electrographic properties (see Figure 4) [49]:

(1) *Inter-ictal* - the period absent from any seizure, when there is no seizure activity and the patient and EEG are at *baseline*. This baseline may contain abnormalities which depend on the underlying seizure disorder and clinical status of the patient. For instance, the EEG of individuals with epilepsy may contain *inter-ictal epileptiform discharges (IEDs)* that can occur several times per minute [50].

(2) *Pre-ictal* - the period immediately before a seizure, where the EEG may not be at *baseline*. This phase is particularly relevant for EEG-based seizure *prediction* algorithms (as opposed to *detection*) that aim to anticipate the subsequent ictal phase. In this period, a patient may experience forewarning symptoms of an impending seizure.

(3) *Ictal* - the period of seizure activity. The symptoms, and EEG characteristics, will depend on the brain region(s) affected by seizure activity. Any spreading of the seizure activity will occur during this phase, causing changes in symptoms and EEG characteristics [50].

(4) *Post-ictal* - the period immediately after a seizure, typically lasting seconds to minutes. The patient may be confused or drowsy, with EEG abnormalities distinct from those in the pre-ictal and ictal phases [51].

Despite seizures having four phases, the transition between phases is not always sudden or clearly defined. This poses challenges for the development of ML algorithms for seizure detection [52].

### 2.4.3 Epileptiform activity

*Epileptiform activity* refers to EEG patterns commonly associated with seizures and underlying seizure disorders, as illustrated in Figure 4.

The distribution of epileptiform activity across EEG channels provides information about the spatial localisation of underlying brain activity. In generalised-onset or secondary generalised seizures, epileptiform activity is present in channels covering both hemispheres. Typically, all channels will display epileptiform activity in the form of generalised spike-and-wave or polyspike-and-wave patterns (see Figure 4).

Whilst epileptiform activity is a hallmark of the ictal phase of a seizure, it is not pathognomic. Epileptiform activity can be seen during any phase of an EEG recording, including the inter-ictal phase. Seizure detection algorithm designers should therefore not assume that the presence/absence of epileptiform activity constitutes a perfect seizure/non-seizure distinction.

Similarly, the *background* or *baseline* EEG of a patient not having a seizure can sometimes display highly-abnormal looking EEG with non-epileptiform focal patterns. Those patterns may share some of the characteristics of epileptiform activity. Examples include periodic lateralised epileptiform activity (PLEDs), now commonly referred to as *lateralised periodic discharges* (LPDs); and generalised rhythmic delta activity (GRDA). However, despite the patient not displaying typical seizure symptoms, the clinical significance of these patterns is still debated [53]. Some argue that there exists somewhat arbitrarily-defined cut-offs for what constitutes true ictal activity [52]. For algorithm designers, the unfortunate - but important - takeaways are (1) that not all focal abnormalities constitute epileptiform activity, but some do, (2) not all epileptiform activity constitutes seizure activity.

### 2.4.4 Seizure subtypes

Beyond the focal/generalised distinction, seizures can be subdivided according to clinical, behavioral and electrographic manifestations. A number of different classification systems exist (e.g. [18, 19]).

In the publicly available Temple University Hospital (TUH) seizure corpus, the seizure subtype labels and their definitions, from most- to least- frequent, are: *focal seizure* (FNSZ), typical duration 10s-2 min, with or without altered awareness, and with or without secondary generalisation; *simple partial seizure* (SPSZ), a focal seizure without loss of awareness; *complex partial seizure* (CPSZ), a focal seizure with loss of awareness; *generalised tonic-clonic seizure* (TCSZ), a generalised seizure with convulsions; *absence seizure* (ABSZ), a generalised seizure with altered awareness but no

convulsions. Rarer subtypes include *myoclonic seizure* (MYSZ) and *tonic seizure* (TNSZ). The acronyms used above reproduce the naming convention used in the TUH seizure corpus.

Most datasets however, do not include labelling of seizure subtypes. This potentially limits the clinical validity of seizure detection algorithms as they might be applied to data where the distribution of seizure subtype differs from that of the training data. For example, absence seizures and tonic seizures are most frequently observed in children with genetic epilepsies, whereas the most common type of seizure in a hospitalised adult is a focal seizure with or without secondary generalisation. Ensuring that the algorithm is capable of detecting both type of seizure activity is critical.

#### 2.4.5 Non convulsive status epilepticus

An important use-case for clinical EEG is diagnosing *status epilepticus* (previously defined), particularly non-convulsive status epilepticus (NCSE). This can be challenging to diagnose because the seizure activity often lacks a clinical (observable) correlate. The lack of available EEG recording and interpretation therefore makes NCSE underdiagnosed.

Epileptiform activity for NCSE can be different to the epileptiform activity observed during ‘self-limiting’ seizures [21]. Furthermore, the EEG abnormality during a seizure, particularly during the protracted duration of status epilepticus (from minutes to hours), can change over time. Consequently, a classifier trained on data from isolated seizures may not accurately detect status epilepticus.

#### 2.4.6 Scalp EEG corrupting artefacts

Scalp EEG is prone to a range of recording artefacts. Whilst their impact can be reduced using EEG pre-processing techniques [54], an awareness of the characteristics of EEG artefacts and their relation to seizures is important for the development of seizure detection algorithms.

- Electrical and environmental interference: EEG recordings can be affected by external sources of electrical interference, such as nearby electronic devices or strong electromagnetic fields (e.g. 50Hz or 60Hz power line noise).
- Ocular artefacts: these result from eye movements and blinks, producing a large electrical potential around the eyes which spreads across the scalp and contaminates the EEG recordings, particularly in frontal electrodes. The amplitude of blinking artefacts is generally much larger (in the order of hundreds of microvolts) compared to that of background EEG activity (a few to tens of microvolts).
- Muscle artefacts: these result from electrical activity and motion associated with muscle contraction. This includes eye movement, jaw clenching, facial grimacing, and large movements of the body such as convulsions during a generalised tonic-clonic seizure. Muscle artefacts are more prevalent in awake patients and often appear as high-amplitude, high-frequency signals.
- Cardiac and respiratory artefacts: electrocardiographic (ECG) artefacts can be detected in some EEG channels. Despite considerable distance between the head and the heart, these artefacts appear as regular spikes or waves in the recording. Breathing related artefacts

can occur due to chest movement or changes in scalp impedance during respiration. These artefacts often appear as rhythmic fluctuations in the EEG.

- Other: sweating will decrease the impedance between the electrode and the skin. The conductive properties of the gel between the skin and the electrode will deteriorate over time, increasing the impedance between the electrode and the skin.

The impact of artefacts varies according to the clinical context (e.g. out-patient facilities versus intensive care units), the recording equipment and channels used, and the clinical status of the patient. Similarly, the artefacts can vary during the different phases of a seizure. In particular, motion artefacts, such as convulsions, can often predominate the measured brain activity in the EEG recording. It is therefore possible that a seizure detection algorithm may inadvertently be trained to detect correlated artefacts rather than the presence of genuine ictal activity.

### 3 Datasets

In this section, we highlight important properties of a range of available datasets for seizure detection applications. We focus exclusively on publicly accessible datasets containing scalp EEG recordings of ictal and non-ictal activity. The majority of published seizure detection algorithms have been trained on one or more of these datasets. However, some algorithms have been developed using proprietary scalp datasets or exclusively intracranial datasets (e.g. [55, 56, 57]). These have been omitted from this review as their clinical implementation is far more limited. For publicly available intracranial EEG datasets, see [58, 59, 60, 61, 62], reviewed in [63].

At the time of writing, there are six relevant datasets: the Children’s Hospital Boston Massachusetts Institute of Technology scalp EEG database [64] (CHB-MIT), the Neurology and sleep centre Hauz Khas database (NSC-HK) [65], the TUH EEG seizure corpus (TUSZ) [66], the Helsinki University Hospital EEG database (HUH) [67], the Siena scalp EEG database (SSE) [68], and the American University of Beirut medical center dataset [69]. Tables 1 and 2 summarise the dataset size and patient demographics; and the basic EEG data characteristics, respectively.

**Table 1. General information regarding public seizure scalp EEG datasets.** In cases where published summaries lacked specific information, we examined the datasets to retrieve the information. N/A denotes that we failed to retrieve the information, even after manually inspecting the dataset.

Dataset	Acronym	Publication date	Number of subjects	Population
<a href="#">TUH EEG Seizure Corpus 2.0.0</a>	TUSZ	2022	675	N/A
<a href="#">American university of Beirut Medical Center</a>	AUB-MC	2021	6	N/A
<a href="#">Siena Scalp EEG</a>	SSE	2020	14	20 - 71 Y
<a href="#">Helsinki University Hospital EEG</a>	HUH	2018	79	35 - 45 w
<a href="#">Neurology and Sleep Centre Hauz Khas</a>	NSC-HK	2016	10	N/A
<a href="#">CHB-MIT Scalp EEG Database</a>	CHB-MIT	2010	22	1.5 - 22 Y

**Table 2. Characteristics and recording parameters of the EEG segments of publicly available EEG datasets.** Note that the NSC-HK dataset has one EEG channel and no information on how it was selected or which placement it corresponds to. In cases where published summaries lacked specific information, we examined the datasets to retrieve the information.

<sup>1</sup> The recording lengths vary between segments.

<sup>2</sup> This is the minimum sampling frequency.

Dataset	Segment duration	Number of seizure events	Number of segments	Sampling frequency (Hz)	Continuous data	Number of channels
TUSZ	<sup>1</sup>	4029	7377	250 <sup>2</sup>	yes	23 - 31
AUB-MC	<sup>1</sup>	35	3895	500	no	19
SSE	<sup>1</sup>	47	41	512	yes	20 or 29
HUH	<sup>1</sup>	460	79	256	yes	19
NSC-HK	5.12s	50	150	200	no	1
CHB-MIT	1 - 4 h	198	664	256	yes	18 - 26

Important differences not presented in the tables include EEG electrode numbers, locations (see section 2.2), montage, EEG recording reference, and the initial pre-processing steps applied to the raw data. This heterogeneity poses significant challenges for making direct comparisons of algorithms trained on different datasets.

Here we expand further on specific characteristics that vary across datasets, of which algorithm designers should be aware.

*Seizure labelling and annotations:* A range of labelling systems have been used across the datasets to label segments of EEG. Following Section 2.4.2, non-seizure periods could correspond to pre-ictal, post-ictal or inter-ictal phases. However, only one dataset (HUH) distinguishes between pre-ictal, ictal and inter-ictal segments, with the rest using a binary seizure/non-seizure labelling system. For datasets with continuous data, it would be possible to extract pre-ictal and post-ictal activity by selecting segments surrounding the period labelled as seizure.

There are also differences in the granularity of seizure classes, ranging from a binary seizure/non-seizure label (in the CHB-MIT and SSE datasets) to the eight seizure subtypes listed above, used in the TUH dataset [70]. Notably, none of the datasets have annotated segments explicitly labelling EEG from *status epilepticus*, even though the electrographic manifestations of *status epilepticus* may differ to that of isolated seizures.

Most often, labelling is applied to the entire recording, known as *term-based* labelling, and not at an individual channel level, known as *event-based* labelling. We will refer to this distinction as the *spatial labelling granularity*. The fine-grained event-based labelling is useful to chart the spatial and temporal propagation of seizure activity. It also has implications for the development of seizure detection algorithms [71].

In addition to spatial granularity, seizure labels can be applied at different temporal resolutions.

The TUH dataset provides precise demarcation of seizure onset/offset, other datasets label pre-defined intervals of EEG segments. For example, a 30 seconds segment containing a 10 second seizure would be labelled as ictal. Again, this has implications for algorithm design and performance evaluation.

*Pre-processing:* There are discrepancies between the pre-processing steps applied to different datasets. For example, the Haus Khas and AUB-MC datasets are band-pass filtered between 0.5-70 Hz and 1/1.6-70 Hz respectively, the AUB-MC dataset has been low-pass filtered at 50 Hz, and the TUH dataset contains raw recordings.

*Balancing and partitioning:* In most datasets, ictal periods constitute a small proportion of the overall data, compared to non-ictal periods [72] (Table 2). However, two datasets contain balanced data. The Haus Khas dataset contains 50 pre-ictal, ictal and inter-ictal segments respectively. Similarly, the AUB-MC dataset contains 3,895 1 second ictal and non-ictal segments respectively. The AUB-MC and TUSZ datasets have pre-partitioned data into independent train and testing subsets, thereby simplifying the comparison of algorithm performance.

A range of techniques have been used to compensate for imbalanced datasets. Firstly, Jiang et al. and Jiang, Xu, and Chen use random undersampling of the majority non-seizure class to provide an equal number of seizure and non-seizure segments. Secondly, a weighted loss function can be used to alleviate the imbalance in the two classes and minimise biased learning [75, 76, 77, 78]. Thirdly, additional non-ictal segments can be created using generative adversarial networks (GAN) Zhang et al.



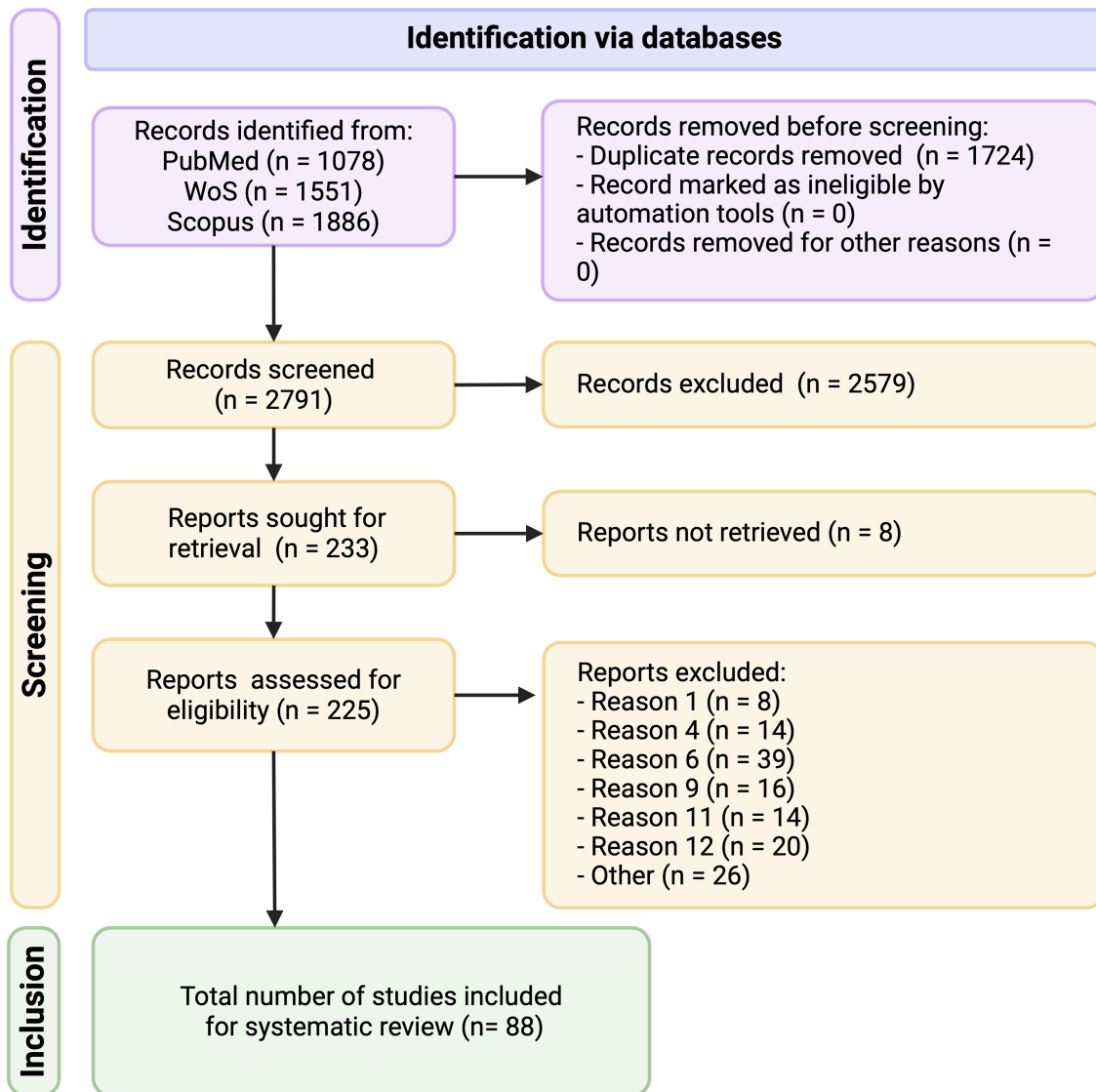
## 4 PRISMA: a systematic search strategy

We systematically selected 88 articles following the Preferred Reporting Items for Systematic Review and Meta-Analysis (PRISMA) guidelines [80]. The strategy is summarised in Figure 5. The three databases searched were PubMed, Web of Science, and Scopus. The initial search was done on February 22, 2023 and the search strings for each database is shown in Table 12 (Appendix). Due to the high number of relevant published material, the search was limited to content published from 2019 onward.

For each document identified, the title, abstract, keywords and author list were imported into Rayyan [81], software which facilitates the screening process. Out of 4515 initial entries, 1724 duplicates were removed, leaving 2791. Exclusion criteria included: (1) studies published before 2019, (2) studies that include other signals than EEG such as ECG or video data, (3) seizure prediction studies, (4) studies using only intracranial EEG data, (5) hardware implementation and testing studies, (6) patient specific algorithms, (7) clinical studies, (8) studies focusing on something else other than seizures, e.g., sleep staging, (9) studies exclusively using private data, (10) studies looking at a single type of epilepsy, (11) studies that do not include any ML, (12) studies that focus on data selected from specific patients without providing reasoning.

An initial screening to exclude articles was done using a keyword search and manual inspection of the results. Keywords screened for were: Bonn, prediction, clinical, sleep, video, neonatal, rats, rodent, survey, systematic review, review, mice, mouse, genetic, gene, Parkinson, MRI, stroke, IoT, antibodies, heart rate, patient specific, personalised, and COVID.

After this screening, 233 articles remained as potentially relevant and their full text manuscripts were retrieved. Out of these, 8 were not be retrieved as they were not free to access. Finally, a total of 88 records were included in the review.



**Figure 5. PRISMA flowchart.** Reasons for exclusion of reports are numbered following the list of exclusion criteria provided above. ‘Other’ includes retracted articles, studies that are not seizure detection and studies that did not meet reasonable quality standards. [82]

## 5 A review of automated seizure detection

Here we review the automated seizure detection ML literature. The stages of a traditional ML pipeline are summarised in Figure 6. It usually begins with the application of data pre-processing steps, and ends with model evaluation. In this section we comprehensively evaluate each possible stage of the pipeline. However, it is worth noting that not all algorithms necessarily adhere to every step of this pipeline.

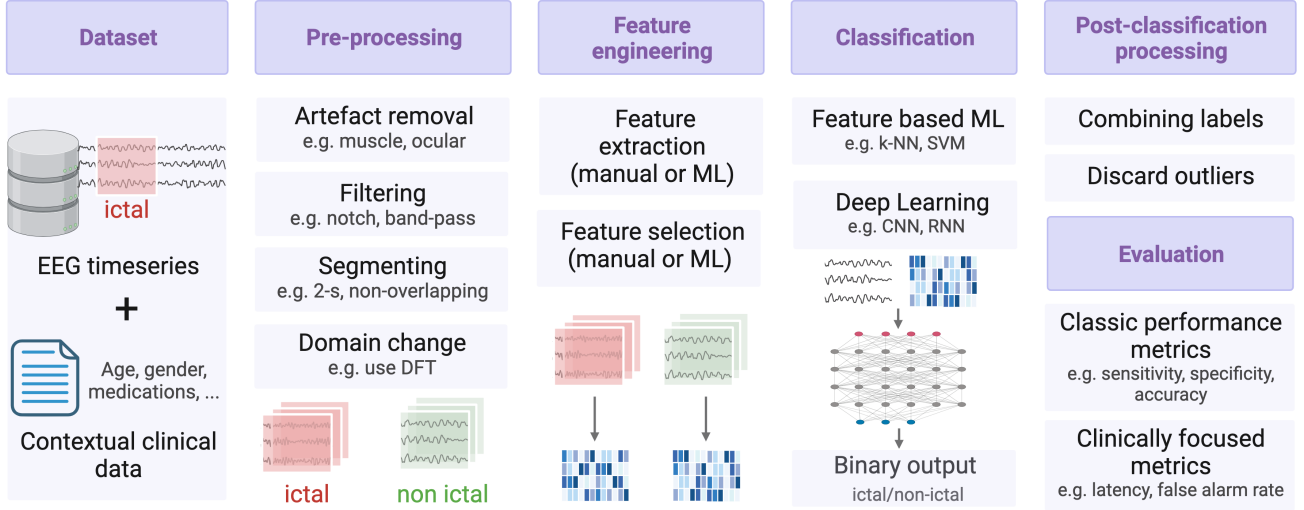


Figure 6. Standard pipeline for automated seizure detection using ML algorithms [9]

### 5.1 Data pre-processing

Pre-processing encompasses a number of possible steps (see Figure 6). Scalp EEG data may be corrupted by a wide variety of noise sources and artefacts. Removing contaminated data or filtering out artefacts can therefore enhance the signal-to-noise ratio and data consistency, potentially improving seizure detection performance. However, pre-processing is a double-edged sword: there is a risk of removing informative signal, such as epileptiform activity. Furthermore, the run time of some pre-processing steps can reduce an algorithm’s potential to operate in real-time.

#### 5.1.1 Noise reduction and artefact removal

A range of techniques have been applied to reduce the impact of noise, although none of the reviewed studies benchmark the effectiveness of their techniques on classification performance.

*Filtering* techniques are applied to attenuate selected frequencies in the EEG. A filter can be designed to remove frequencies above a threshold (low-pass), below a threshold (high-pass), frequencies that are not within a chosen interval (band-pass), or to remove those that are (notch).

High-pass filtering is commonly employed to remove the effect of slow ( $< 1$  Hz) presumed non-neural baseline signal drift. A low-pass filter is often used at the time of recording to prevent aliasing, which occurs when the sampling frequency is lower than twice the maximum frequency in the signal. Low-pass filtering is also often used offline to reject frequencies above 30-45 Hz, considered to mainly contain irrelevant noise and artefacts. However, this may result in the

exclusion of informative neurally-driven high-frequency signals [83, 84]. In practice, band-pass filtering is generally applied to attenuate both low- and high-frequency components (e.g., 0.3-60 Hz Humairani et al.). Notch filters are commonly used to remove power line interference (e.g. 60 Hz in North America and 50 Hz in Europe [86, 87, 88, 76, 89]).

There is no consensus on the selection of frequencies to attenuate, with studies generally excluding frequencies outside the range of  $\sim 0.5$ -40 Hz.

Beyond differences in bandwidths, a range of digital filtering techniques are employed. The Butterworth and the Savitzky–Golay filters are most commonly used. Less commonly encountered are adaptive filters, used for artefact removal. Their ability to continuously adapt to the input makes them well-suited for dealing with non-stationary data and hardware degradation. Notably, Poorani and Balasubramanie use an adaptive vector median filter to remove artefacts from the CHB-MIT data [90].

A range of *blind source separation* (BSS) techniques have been applied to EEG to recover the original source signals from their observation (mixed with non-neural sources). BSS techniques require the different source signals to have different statistical properties. Commonly encountered BSS techniques are principal component analysis (PCA), independent component analysis (ICA), non-negative matrix factorisation (NNMF), and common spatial pattern (CSP).

BSS is commonly applied to remove stereotypical artefacts. One example is blinking artefacts which has non-Gaussian statistical properties independent of other EEG source signals, making it appropriate for removal using ICA [91]. BSS techniques are typically applied after filtering. For example, Raghu et al. uses a notch filter and a 0.5-40 Hz band-pass filter before ICA-based artefact removal [92]. These techniques can, however, remove useful EEG information.

Another technique used for pre-processing, but also for feature extraction, is the *wavelet transform*. A wavelet, unlike the sinusoids used in the Fourier transform, is a decaying, wave-like oscillation. This makes it better suited to non-stationary signal analysis [93]. The wavelet transform evaluates how much of a wavelet is present in a signal for a specific scale and time shift, using convolution. It is sometimes used to separate the EEG into the canonical frequency bands (see definitions above) [94, 45].

Wavelet decomposition can help identify artefacts from each channel individually, as opposed to ICA/PCA approaches that require multi-channel data. Wavelet transforms can also be used as a set of new features. For example, Nemati and Meshgini use the discrete wavelet transform (DWT) to decompose the EEG signal and subsequently produce a correlation map which is used as the input to a classifier [95].

Breaking down signals into various components can be done using BSS techniques, Fourier transform, and Wavelet transforms. It can also be performed using empirical mode decomposition (EMD) and variational mode decomposition (VMD). These methods decompose a signal into various time domain signals called intrinsic mode functions (IMFs) [96]. Whilst EMD decomposes the signal adaptively in a data-driven way, VMD uses a non-recursive decomposition technique. Even though both methods can handle non-linear data, EMD is the only one capable of working on non-stationary signals [97]. For these reasons, EMDs are widely used in biological signal analysis [98, 99]. Raw IMFs or their derived features can be used as input to a ML algorithm. For example, Ru et al. employ VMD to partition the signal originating from electrodes FZ-CZ and CZ-PZ into five distinct IMFs. These IMFs’ phase synchrony index is then used as input for a detection model

[100]. In a related study, Davidson et al. compare three distinct source isolation techniques for artefact removal: empirical mode decomposition (EMD), empirical wavelet transform (EWT), and variational mode decomposition (VMD) [101]. The EMD method showed the best performance.

Numerous EEG pre-processing approaches exist, which can all be combined in multiple ways. There is no standard method for pre-processing EEG data. We recommend that researchers apply their detection algorithm on raw data, data that has undergone minimal filtering, and data that has been pre-processed with more complex techniques like ICA or wavelet decomposition. This way, performance improvements due to different pre-processing steps can be clearly measured in the context of seizure detection. For example, it is frequently assumed that DL techniques handle raw data better than feature-based ML, however, we found no study that has tested this assumption quantitatively for seizure detection.

### 5.1.2 Channel selection

Scalp EEG records the electrical activity of the brain typically using  $\sim 20$  electrodes. Selecting a subset of these electrodes for analysis and disregarding others can help an algorithm reach better performances and improve training time due to a reduction in the input dimension. However, there is also a risk of dismissing clinically relevant information, hereby decreasing the algorithm’s performance and generalisability. Here we outline the different approaches to channel selection.

Channel selection can be performed based on someone’s critical judgment and assessment of which electrodes are, in theory, most useful to keep. For instance, Wei and Mooney select 4 channels of the TUSZ data capable of recording different brain regions: F8-T4 (temporal lobe), T5-O1 (occipital lobe), FP1-F3 (frontal lobe) and F3-C3 (central cerebral hemisphere) [89]. On the other hand, Ru et al. only use channels FZ-CZ and CZ-PZ [100]. The examples above illustrate theory-based channel selection.

Alternatively, some studies rely on quantitative approaches to channel selection in order to minimise the loss of information and drop in performance. For instance, Nandini et al. use PCA to perform channel selection [102]. Although useful to reduce algorithmic complexity, this data-centric selection process is often not applicable in real-time, which limits its relevance to clinical applications. Additionally, seizure activity may propagate to different brain regions, which might be missed if the channel selection process is performed using a short segment of EEG.

It is also important to identify and state the motivation behind channel selection, whether it is to decrease the run time, to reduce redundancy, or to make the clinical implementation more realistic. The purpose of channel selection is the identification of a subset of channels that are patient independent and contain enough information to identify all types of seizures. We recommend the channel selection process to be repeated across different datasets to ensure that the identified channels are not dataset dependent.

### 5.1.3 Data Segmentation

Segmenting data is important for the development of online algorithms that are capable of making real-time decisions. There are two parameters involved in the segmentation process: the segment length and amount of overlap between consecutive segments. We identify three important clinical implications for choosing these parameters.

Firstly, given the non-stationarity of seizure activity, the segment length used influences the patterns the model learns. Most algorithms are fed EEG segments shorter than 10 seconds as input, as this is generally short enough to consider the segment as quasi-stationary.

Secondly, the amount of overlap between the segments can be leveraged to (1) reduce the imbalance in the training data, and (2) reduce detection latency. Introducing overlap between the extracted segments increases the number of samples. This can help for training the model, but can also introduce redundancy and overfitting. A small number of studies have used non-overlapping EEG segments lasting at least 60 seconds [103, 104, 105]. The longer the time between the analysis of two consecutive EEG segments, the larger the seizure detection alarm delay can be. The length of the EEG segments and whether they are overlapping is shown in Table 4 to 11.

Thirdly, a segment length which is too small can remove the historical context required for an accurate classification. However, this can be incorporated by some ML algorithms with memory such as RNNs, or, by using some post-classification analysis (Section 5.5). For instance, the detection of status epilepticus depends on having 5 minutes of ictal activity, which a single 5 second segment does not provide.

#### 5.1.4 Validation strategy

Model training and evaluation requires the data to be divided into either a train-set and a test-set or into a train-set, validation-set and test-set. Each of these sets serves a different purpose in the model building pipeline. The train-set is the data used during model training. The validation-set is used after each training epoch to evaluate the performance of the model. It further enables tuning of the model’s hyperparameters without overfitting to the train data. The test-set is the final data set and is used to evaluate the performance of the model on unseen data after training and tuning.

There are a number of ways to construct either the train-test or the train-validation-test datasets. Common examples are (1) random split, (2) k-fold cross-validation and (3) leave-one-out validation. All of these approaches can be used for train-test or train-validation-test splits. Tables 4-6 and 7-11 summarise the validation strategy of the reviewed literature.

A random split is the simplest data splitting strategy. In this regime, the data is randomly divided into 2 or 3 different sets, often with 10% for testing, and the remaining 90 % solely for training, or for training and validation.

k-fold cross-validation (CV) partitions the dataset in k evenly-sized pieces, each of which is called a fold. k is often chosen as 5 or 10. On each iteration, one of the k pieces is taken as the test set and the remaining data, that is, the other k-1 pieces, are taken as the training set, or the training and validation sets. A model is then trained and tested using these datasets and the performance recorded. The performances of the k models are then averaged, which gives a more reliable estimate of the model’s performance on unseen data.

Finally, the leave-one-out (LOO) strategy, also called leave-one-patient-out in the context of seizure detection, is the same in principle as k-fold CV, but where one takes the data of a single patient as the test set, and the remaining data as the training or training and validation sets. In particular, it is a special case of k-fold CV. We refer to this as LOO cross validation (LOO-CV). As an example, the CHB-MIT dataset contains 22 subjects, and so using LOO-CV would train 22 models, each being trained using the data of 21 patients and tested on the data of the patient that was ‘left

out’. Each model is tested on a different patient. Patient independent performance of the model is obtained by averaging the performance of the 22 models. This is a rigorous validation strategy in the field of seizure detection as it prevents overfitting and is patient independent.

## 5.2 Feature engineering

A feature of a time-series is any calculated measure, for example, the power in set frequency band, entropy values or Hjorth parameters. This section starts by describing the different EEG domains that features can be extracted from, followed by how this set of features can be reduced to an adequate subset for classification purposes.

### 5.2.1 Data domains

The domains, or representations, which EEG time-series data can be mapped to are: time, frequency, time-frequency, and network. In this subsection, we briefly review the four different domains for EEG data representation.

The *time domain* is simply the amplitude as a function of time, and is the domain of raw EEG recordings.

The *frequency domain* refers to the representation of signals in terms of their frequency content. Switching from the time to frequency domain can be achieved using the Fourier transform or Laplace transform (for the complex frequency domain).

The *time-frequency domain* represents how the frequency components of a signal change over time. This can be calculated by segmenting the data into consecutive segments whose frequency components are extracted one after the other. The output is 2-dimensional and known as a spectrogram. Various techniques exist to obtain a time-frequency representation of a signal, but common ones include the short-time Fourier transform (STFT) and continuous wavelet transform (CWT). Spectrograms are often used as input to CNNs and are widely used in seizure detection.

Finally, EEG data can be represented in a *network domain*, capturing *functional connectivity*, or measures of how brain regions interact with one another over time. Typically, a network is constructed as a weighted connected graph, where each EEG channel is represented by a node and edges between nodes are assigned numerical quantities corresponding to the connectivity between the signals from pairs of EEG channels. For instance, Li et al. used the cross-correlation between channels as edge values [106], but a range of connectivity-based measures can be used, including phase synchrony and the phase lag index. Once the EEG data are represented as a connected graph, a range of graph theoretical metrics, such as the degree and clustering coefficient, can be used as features.

### 5.2.2 Features

Table 3 shows commonly extracted features from EEG for seizure detection. Importantly, EEG features are either univariate, bivariate, or multivariate, that is, involving one, two, or more than two channels in the analysis, respectively [107].

Although features are typically applied to EEG data in a certain domain, as shown in the table, they can be applied to other domains, provided that the dimension of the input is compatible with the arithmetic steps of the feature calculation. For example, Hjorth parameters are defined for



the time domain. However, it is possible to calculate them for any series data, such as a power spectrum which belongs to the frequency domain. Of course, the interpretation of features will differ according to the domain of the data to which they are applied.

**Table 3. Example of popular time-series features used to discriminate between ictal and non-ictal EEG segments.** Even though features can be extracted from data belonging to any domain, providing it has the correct dimensions, they are often used in one domain only. The features in the table below are grouped by the domain within which they are most commonly applied.

Domain	Features
Time domain	Absolute energy, energy ratio
	Mean and variance over the autocorrelation for different lags, partial autocorrelation
	Autoregressive coefficients
	Hjorth activity, mobility, and complexity, line length
	Mean, median, mode, maximum, minimum, mean, standard deviation, variance, quantile range, kurtosis
	Permutation entropy, Shannon entropy, spectral entropy, approximate entropy, Sample entropy, Fuzzi entropy, Renyi entropy
Frequency domain	Fourier transform aggregate and coefficients, band power
Time-frequency domain	Continuous wavelet coefficients and peaks
Network domain	Weighted degree, clustering coefficient, harmonic centrality, modularity, closeness centrality, eigenvector centrality

### 5.2.3 Feature selection

It is common to encounter studies where a large number of features from different domains are extracted. Only using a subset of those features as input to a ML algorithm can reduce computational complexity, redundancy, and bias.

Broadly, there are four strategies for feature selection. Firstly, *dimensionality reduction* is commonly achieved through principal component analysis (PCA). PCA can be applied prior to training the algorithm to reduce the dimensionality of either the EEG input signal directly or the input feature space, as done by Jiang et al. [73]. Other algorithms reducing the dimensionality of the feature space include minimal redundancy maximal relevance (mRMR). mRMR emphasizes selecting a feature subset that optimally balances relevance to the target variable while minimizing redundancy among the selected features.

Secondly, after an algorithm has been trained, feature selection can be informed by an evaluation of *feature importance*, that is, the relative weights of each feature in a classification process. For example, Wu, Zhou, and Li extracted a total of 798 features from the raw EEG signal and several IMFs obtained after complete ensemble empirical mode decomposition (CEEMD) [108]. Using an XGBoost algorithm allowed them to quantify which of the 798 features were allocated the highest discriminative power, or rank. Among the 20 features with the highest rank obtained when using the CHB-MIT dataset, 14 belong to the time domain (70%), 4 belong to the frequency domain



(20%), 1 belongs to the time-frequency domain (5%) and 1 is an entropy-based feature (5%). A similar analysis was performed on the Bonn dataset and revealed that although some features were similarly high ranking across both datasets, others were not. For instance, the autoregressive coefficients of the raw signal and IMF10 trace both played an important role in the classification of the Bonn dataset. However, they were not among the 20 most important features when analysing the CHB-MIT dataset. This analysis shows that features from multiple domains can contribute to differentiating ictal and non-ictal EEG segments. It also demonstrates the heterogeneity of publicly available datasets, and the challenges of developing a patient and hardware independent algorithm. Another method for quantifying feature importance is Shapley’s additive explanations (SHAP) [109]. Another feature selection method based on feature importance is sequential forward feature selection (SFFS). Features are iteratively added to the model and the corresponding improvement of performance metrics like accuracy or precision is evaluated at each step.

Thirdly, features can be chosen according to the required *explainability* of the algorithm. The correspondence of some features to ictal activity is intuitive: for example, it is easy to understand that, generally, ictal activity is associated with an increase in the power of the EEG signal. However, many features discriminate in less intuitive ways, e.g. auto-regressive coefficients. Using features where the meaning is difficult to understand, or using too many features, can make the results difficult to interpret or risk overfitting the data. This could generate scepticism amongst end-users and pose a barrier to medical deployment.

Fourthly, choosing which features to use can be based upon their *computational complexity* and, hence, their execution time. For example, the complexity of approximate entropy and sample entropy are both  $\mathcal{O}(n^2)$ , where  $n$  is the length of the time series. As such, they are less suitable for real-time applications using EEG segments with a large number of data points [110]. On the other hand, the fast Fourier transform has a computational complexity of  $\mathcal{O}(n \log(n))$ . This makes features such as the power present in the canonical bands much quicker to calculate than entropy measures for large  $n$ .

### 5.3 Feature-based ML algorithms for classification

Once the data pre-processing, feature extraction, and feature selection steps are completed, the data is ready for classification. It is worth noting that deep learning (DL) does not necessarily require extracted features as input. This section reviews three common feature-based ML algorithms for automated seizure detection. We focus on DL models in the following section. Tables 4, 5, 6 summarise the feature-based ML models encountered in the literature.

Firstly, *k-nearest neighbours* (*k-NN*) is one of the simplest supervised machine learning algorithms used in seizure detection. For classification tasks, the algorithm assigns the label to the test data that is most common amongst its  $k$ -nearest neighbours. As this algorithm necessitates data points associated with labels as input, it does not support time series data. Therefore, features of the EEG segments are used as input for the classification of ictal and non-ictal activity. For instance, Jumaah, Shihab, and Farhan performs a discrete cosine transform on the EEG segments. The energies of the sub-bands of the power spectrum are the input to a  $k$ -NN algorithm which achieved an accuracy of 93.64% using 5 of the 23 available channels of the CHB-MIT dataset [111]. Jiang, Xu, and Chen extract features using a dual-tree discrete wavelet transform and use a feature selection method known as transfer component analysis [74]. The binary  $k$ -NN had an overall accuracy of 74.03% on the CHB-MIT dataset. A downside of the  $k$ -NN algorithm is that

it is a non-parametric method, meaning that it does not fit a model during training. Instead, it ‘memorises’ the training data. Identifying the  $k$  nearest neighbours of a new data point requires computing its distance to every training data points available. This computational inefficiency makes  $k$ -NN too slow for online applications using large amounts of data [112, 113].

A second popular algorithm which can handle larger datasets is the *support vector machine (SVM)*. It is a parametric method that finds the optimal hyperplane separating the data into various groupings. Humairani et al. compares the use of Shannon entropy and Renyi entropy as input to a SVM [85]. They achieve a 92.96% accuracy using Renyi’s entropy on the CHB-MIT dataset. Similarly, Raghu et al. use a DWT based sigmoid entropy as input to a SVM to perform binary classification on the CHB-MIT data and achieve a 94.21% accuracy [92]. Fathima, Rahna, and Gaffoor use more features as input to an SVM, they first compute the DWT decomposition at level 5 using Daubechies 4 wavelet. Note that the chosen wavelet should be the one whose profile best matches the characteristic shape we are looking to detect in the EEG. The Daubechies wavelets are well known for their resemblance to biological signals, and are often used in DWT decomposition of EEG. The level, or scale, of decomposition refers to how ‘stretched’ the wavelet is. The higher the level, the wider the wavelet. Choosing the scale of decomposition is also arbitrary and depends on the nature of the signal. Back to Fathima, Rahna, and Gaffoor, who calculate the mean, RMS, inter-quartile range, and entropy of the decomposed signal [114]. The reported accuracy is 98.6%, however, there is no justification behind the restriction of the feature extraction to the level 5 decomposition of the EEG signal. Some papers use algorithms for feature extraction and selection prior to using them as input to the SVM for classification. This is the case of Shariat et al., who compare two feature selection techniques; sequential forward feature selection (SFFS) and minimal redundancy maximal relevance (mRMR) [115]. They achieve a maximal accuracy of 99.87% using SFFS followed by a SVM on all 23 patients of the CHB-MIT.

The third common feature-based ML algorithm for seizure detection is the *gradient boosting machine (GBM)*. This is an ensemble learning method that consists of multiple weak prediction models, typically decision trees, that are combined to form a stronger prediction model. An initial weak sub-learner is trained and, subsequently, a second sub-learner is constructed to fit the residuals of the first one, et cetera. Nandini et al. calculate seven time domain features (min/max, median, mean, skewness, standard deviation and kurtosis) of EEG signals decomposed using wavelet transforms. For classification, they use an extreme gradient boosting (XGBoost) algorithm, and achieve a classification accuracy of 94.46% [102].

In addition to the three main classification algorithms, other less common methods include random forest classifiers (RF) and linear discriminant analysis (LDA) [116, 117, 118, 119, 120, 100]. For instance, Xiong et al. use a RF classifier on the CHB-MIT and Siena datasets using network features such as weighted degree and clustering coefficient. Their reported accuracy is 84.83% [116]. Rohira et al. compares two inputs to a RF classifier, (i) spectral-based measure of functional connectivity, and (ii) the power of 6 frequency bands [117]. Using 8 channels they achieve an accuracy of 90.87% and 95.73% when using coherence coefficients and the power of different frequency bands respectively.

**Table 4. Feature-based ML methods from systematic review of literature for seizure detection in scalp EEG data**

Classifier	Feature(s)	Dataset(s)	Performance	Validation	Segment Length	Year	Reference
k-NN	Energy of signal after DCT	CHB-MIT (21 patients, 5 electrodes)	acc: 93.64% sen: 94.77% spe: 92.21% F-score: 93.12% FPR: 0.07 FNR: 0.05 Error: 0.06	10-fold CV	1s, no overlap	2019	[111]
fuzzy k-NN	GNMF decomposed SSTFT maps	CHB-MIT, Bonn	acc: 98.99%, sen: 99.27 %, spe: 98.53%	10-fold CV	1s, no overlap	2023	[121]
Neutrosophic logic k-NN	$\theta, \beta, \delta, \alpha$ power bands of four wavelet bands and $\alpha$ to $\delta$ power band ratio	Bonn, CHB-MIT (18 channels, no patient 12), AIIMS	acc: 89.06%, sen: 85%, spe: 89.06%	LOO	1s, no overlap	2020	[122]
k-NN (feature selection), RF (classification)	Weighted degree, clustering coefficient	CHB-MIT, Siena scalp	CHB-MIT: F1: 86.69% AUC: 84.33% acc: 84.83% pre:85.60% sen :87.81% spe: 81.01%	5-fold CV	4s	2023	[116]
TCA (feature mapping to latent sub-space) and k-NN (classification)	Separable features from dual-tree discrete wavelet parameters	CHB-MIT	acc: 74.03% F1: 0.7473 AUC: 0.8204	LOO	3s, 2.5s overlap	2019	[74]
SVM vs k-NN	Hurst exponent, Logarithmic HFD	CHB-MIT	SVM acc: 99.81%, rec: 100%, TNR: 0.99 k-NN acc: 93.21%, rec: 92.56%, TNR: 0.92	10-fold CV	300s, 240s overlap	2022	[103]
SVM, LDA, k-NN	Power, mean, std extracted after using DTCWT	CHB-MIT	acc: 100%	10-fold CV	8s, overlap	2020	[119]
SVM, ELM (SVM is best)	Weighted FPE complexity-based feature (W-FPE-F)	CHB-MIT (12 patients), Bonn	acc: 98.9883% spe: 89.3300% sen: 94.1650%	10-fold CV	4s, 3s overlap	2019	[123]
SVM	Shannon Entropy Renyi's Entropy after DWT	CHB-MIT	acc: 92.96% (using Renyi)	N/A	10 mins	2022	[85]
SVM	Std, mean absolute deviation, RMS, min, interquartile range, skewness, entropy and max were extracted over wavelet coefficients	CHB-MIT	spe: 100% sen: 97.2% acc: 98.6%	N/A	5s	2020	[114]
SVM	DWT based sigmoid entropy (in time and frequency domain)	CHB-MIT, Bonn, RMCH	sen: 94.21%	LOO	1s	2019	[92]
SVM	Successive decomposition index (SDI)	RMCH, CHB-MIT, TUSZ	CHB-MIT: sen: 97.28% FA: 0.57/h median latency: 1.7s TUSZ: sen: 95.80% FA: 0.49/h median latency: 1.5s	LOO	1s, no overlap	2019	[124]
SVM	Kurtosis, skewness, line length, quartile values, correlation coefficient matrix of the frequency energy between any two channels. PCA reduction of the dimensionality.	CHB-MIT, Siena	acc: 96.67%, spe: 95.62%, sen: 97.72%	Bootstrap	1s, 0.5s overlap	2023	[73]

**Table 5. continued - Feature-based ML methods from systematic review of literature for seizure detection in scalp EEG data**

Classifier	Feature(s)	Dataset(s)	Performance	Validation	Segment Length	Year	Reference
SVM (classification), Sequential Forward Feature Selection (SFFS) vs Minimal Redundancy Maximal Relevance (mRMR) (feature selection)	Covariance matrices of channels modified using Riemannian geometry	CHB-MIT (22 channels)	CHB-MIT using SFFS and 10-fold CV: acc: 99.87% sen: 99.91% spe: 99.82%	10-fold CV	2s, no overlap	2021	[115]
LS-SVM	Mean, std, variance, Shannon entropy, approximate entropy, spectral centroid, spectral speed, spectral flatness, spectral slope, spectral entropy, Hurst exponent, Katz fractal exponent	CHB-MIT	acc: 98.37%, sen: 91.11%, pre: 91.67%, spe: 91.46%, ADR: 91.28%, G-mean: 91.28%, AUC: 0.992	10-fold CV	N/A	2022	[125]
SVM, Quadratic Discriminant Analysis (QDA) and LDA	Maximum value, minimum value, mean value, variance, range, skewness, kurtosis, estimation of cross-correlation and various connectivity measures from graph theory	CHB-MIT (14 channels)	accuracies: SVM: 98.09%, QDA: 81.49%, LDA: 80.90%, SVM AUC: 99.7%, SVM sen: 98.1% SVM spe: 98.1%	10-fold CV	5s, no overlap	2022	[126]
Layered directed acyclic graph support vector machine (LDAG-SVM)	Entropy, largest Lyapunov exponent, correlation dimension	CHB-MIT, Bonn	acc: 95% sen: 99% spe: 96% run time: 98ms	50-50 train-test	N/A	2019	[127]
RF	Mean value and peak-to-peak value of wavelet energy obtained after performing PDWC	CHB-MIT, NICU, Pune_pat, Bonn	TPR: 99.42% PPV: 99.71% TNR: 99.71% NPV: 99.71% acc: 99.67% F1: 99.54%	80-20 train-test	4s	2021	[120]
RF	Improved sample entropy, phase synchronization index	CHB-MIT (2 electrodes)	acc: 91.78% sen: 91.27% spe: 93.61%	10-fold CV	2s	2022	[100]
RF	Weighted degree, clustering coefficient, harmonic centrality, modularity, closeness centrality and eigenvector centrality are extracted in 3 networks constructed by PCC, MI and permutation disalignment index	CHB-MIT, Siena	CHB-MIT: acc: 97.26%, sen: 96.89%, spe: 97.55%, F1: 97.11%, Siena: acc: 98.88%, sen: 98.36%, spe: 99.13%, F1: 98.75%	70-30 train-test, 5-fold CV	4s, 3s overlap	2022	[128]
RF	Power of 6 PSD brain wave bands, vs coherence coefficient	TUEP (8 channels)	Coherence coefficients: acc: 90.87%, PSD, acc: 95.73%	70-30 train-test	10s, no overlap	2023	[117]
RF	Time and frequency features	CHB-MIT, private dataset	acc: 99.36%, spe: 82.98%, sen: 99.41%, FPR: 0.57 times/h	Leave-5-patient-out	5s, 4s overlap	2019	[129]
RF	Standard deviation, the IQR and energy of each sub-band	CHB-MIT	acc: 94.04%, sen: 89.5%, pre: 98.4%	10-fold CV	10s, no overlap	2022	[130]
RF	Hjorth parameter, time correlation coefficient matrix, eigenvalues of correlation coefficient matrix, sub-band signal energy, frequency correlation coefficient matrix, fuzzy entropy	CHB-MIT	acc: 98.03%, spe: 99.04%, sen: 97.02%	Leave-5-patient-out	4s, 2s overlap	2023	[131]

**Table 6. continued - Feature-based ML methods from systematic review of literature for seizure detection in scalp EEG data**

Classifier	Feature(s)	Dataset(s)	Performance	Validation	Segment Length	Year	References
LDA (classification), bagging (feature selection)	Spectral edge frequencies, spectral edge powers, IQR, MAD, PCC	CHB-MIT (18 channels) AIIMS (private)	CHB-MIT: acc: 84.83% FDR: 1.2/hour mean latency: 1.43s	N/A	1s, no overlap	2021	[118]
LDA	Univariate features: kurtosis, mean absolute deviation, interquartile range, and semivariance are calculated after the DWT. Bivariate feature: measure of correlogram	CHB-MIT (14 patients)	sen: 100% FP/Hour: 0.59 spe: 99.8% acc: 99.6%	3-fold CV	1s, no overlap	2020	[132]
XGBoost	WAF-based hybrid extracted features, SSA and time-domain features	CHB-MIT (18 channels, 10 patients)	acc: 94.46% sen :88.61% spe: 88.61% precision: 99.81% MCC: 89.54% kappa: 89.03%	5-fold CV	6s, no overlap	2022	[102]
XGBoost	Mean, std, signal envelope, kurtosis, skewness, complexity, mobility, TKEO, fractal dimension, band power, sum of relative beta and gamma	TUSZ (4 channels)	sen: 20% FA/24h: 15.59	N/A	1s, 0.5 overlap	2020	[89]
Naive Bayes	Relative amplitude, spectral entropy, logarithmic band power, tonal power ratio, 1D local binary pattern, PSD, spectrogram	CHB-MIT, TUEP	TUEP: acc: >90%, sen: >85%, spe: >85%, CHBMIT: acc: 90%, sen: >92%, spe: >92%	90-10 train-test	N/A	2022	[133]
Naive Bayes	10 geometric features extracted in each frequency band $\theta, \beta, \delta, \alpha$	CHB-MIT	acc: 94.54%	10-fold CV	20s, 15s overlap	2022	[134]
Genetic algorithm - Binary Grey Wolf Optimisation	Std, Shannon entropy, kurtosis, Hjorth parameters, skewness, energy and nonlinear energy, Higuchi fractal dimension, Katz fractal dimension, spectral entropy	TUH	acc: 85%	N/A	1.8s, no overlap	2021	[135]
Hidden Markov Model	DMD power, sum of 2D PSD, variance, and KFD features	CHB-MIT, AIIMS	average CHB-MIT: acc: 99.60% MCC: 0.97 Kappa: 0.97 FPR: 0.12% NPV: 99.69% PPV: 98.73% Sen: 96.64% Spe: 99.88%	N/A	5s, no overlap	2020	[136]
NN	AM bandwidth, FM bandwidth, frequency, kurtosis, Hjorth complexity, Hjorth mobility, skewness, spectral centroid, spectral entropy, spectral peak, Spectral power for 8 IMFs	Bonn, NSC-HK	acc: 98.1%, sen: 98.21 %, spe: 97.65%	70-30 train-test	N/A	2022	[137]
Multi-layer perceptron	Riemannian tangent space map features	TUSZ (18 channels)	acc: 98.94%, Kappa: 0.916	N/A	6s, 3s overlap	2021	[138]

## 5.4 DL algorithms for classification

Whilst feature-based ML requires features as input, deep learning (DL) can uncover patterns and features from different types of data. The input of DL algorithms can be raw or filtered EEG data, any domain representation, or a set of features extracted from EEG signals. DL architectures commonly used for automated seizure detection include artificial neural networks (ANN), convolutional neural networks (CNN) and graph machine learning (GML). Different architectures will classify an EEG segment based on different properties of the signal. In this section we will review these four DL methods. It is worth noting that to leverage the strengths and compensate for the weaknesses of some ML architectures, some studies combine different DL architectures (in parallel or in series). Tables 7, 8, 9, 10 and 11 summarise the DL models encountered in the literature.

ANN are composed of a combination of node layers: an input layer, one or more hidden layers, and an output layer. Sallam et al. use the averaged frequency-spectrum values of the  $\gamma, \beta, \alpha, \theta$  and  $\delta$  frequency bands as input to a single layer ANN comprised of 10 hidden neurons. This simple architecture yields an overall accuracy of 93.5% when trained and tested using the CHB-MIT dataset.

CNNs are the most widely used DL technique in the field of seizure detection: 43 of the 54 deep learning papers included in this study used convolutional layers. The input of a CNN can be 1-dimensional, for example, a single channel EEG trace; 2-dimensional, a tensor where each row represents a channel and each column is a sample; or 3-dimensional, where each segment is represented as a 2D tensor, and samples are concatenated in the third dimension. For instance, Kaziha and Bonny and Huang, Chen, and Cao construct a 2D input tensor from 100 and 23 second raw EEG segments from 24 and 21 channels respectively [104, 140]. In contrast, Kumar, Janghel, and Sahu and Qiu, Wang, and Jiao use 1-dimensional EEG segments as input to a CNN [141, 110]. In Qiu, Wang, and Jiao’s work, the EEG segments of each channel are convoluted separately [110]. The input to CNNs is not restricted to the time-domain. Commonly, frequency domain representations are used, or the concatenation of frequency domain features across channels. Sharan and Berkovsky compute the power of 132 different frequency bins for each channel and store it as a 2-dimensional input tensor [142]. Similarly, Li and Chen concatenate the 1-dimensional FFT of each channel to form a 2-dimensional input tensor [143]. These approaches result in an accuracy of 97.25% and 98.47% respectively. Varlı and Yilmaz combine two different DL algorithms, each with a different input. They input spectrograms and scalograms (absolute value of the CWT of a signal, plotted as a function of time and frequency) into a CNN, and raw EEG traces into a RNN [144]. Combining both DL approaches potentially leverages their respective strengths. Surprisingly, the reported accuracy of 96.23% is lower than that obtained by other studies using a single CNN.

Unlike ANNs and CNNs, graph machine learning incorporates the spatial information of the electrode placement, which can provide insight on seizure dynamics. Examples of graph machine learning are graph attention networks (GAT), and graph convolution networks (GCN). Zhao et al. use a GAT with the CHB-MIT dataset, and obtain an accuracy, sensitivity and specificity of 98.89%, 97.10% and 99.63% respectively [145].

Some methods combine ML and DL. For example, Dalal, Paunwala, and Chapatwala use statistical features like kurtosis, variance and skew as the input to multiple CNNs [146]. The output of the CNNs are concatenated and used as input to a LSTM (a type of RNN) for binary classification. They achieve an accuracy of 94.6%, recall of 97.15%, and precision of 95.78%.

There are instances where increasing complexity of the respective DL algorithm does not lead to an increase in accuracy. For example, Yan et al. trained four different CNNs [147]. They observed that the weights of the CNN with the highest complexity (4 convolution blocks/10 layers) did not converge during training. Consequently, the CNN fails to perform the classification task. The authors suggest that after using four convolution blocks the extracted spectrographic features are not seizure specific anymore [147].

**Table 7. DL methods from systematic review of literature for seizure detection in scalp EEG data**

Classifier	Dataset(s)	Performance	Validation	Segment Length	Year	References
CNN	CHB-MIT	sen: 97.25%, spe: 97.25%, acc: 97.25%	10-fold CV	3s	2020	[142]
CNN	CHB-MIT	acc: 96.74% spe: 100% sen: 82.35%	5-fold CV	100s	2020	[104]
CNN	CHB-MIT	acc: 87.4% sen:88.10% spe:87.10% F1:87.40% pre: 86.98%	10-fold CV	8s	2021	[141]
CNN	CHB-MIT, Bonn	acc: 96.69% sen: 96.19% spe: 97.08%	k-fold CV	2s	2023	[148]
CNN	CHB-MIT, Bonn	acc: 98.80% sen: 98% spe: 98%	10-fold CV	N/A	2021	[149]
CNN + MIDS, CNN + data augmentation	CHB-MIT	CNN+MIDS: sen: 74.08% spe: 92.46% CNN+Data aug: sen: 72.11% spe: 95.89%	LOO	5s	2019	[150]
CNN aided factor graph	CHB-MIT	AUC-ROC: 90.23% AUC-PR: 76.77% F1: 90.42%	6 fold, leave 4 patients out	4s	2021	[87]
CNN aided factor graph	CHB-MIT	AUC-ROC: 83.8% AUC-PR: 50.38% F1: 93.42%	6 fold, leave 4 patients out	4s and 32s	2022	[86]
Attention-based CNN-BiRNN	CHB-MIT	No missing channel: spe: 93.94% sen: 92.88% 2 missing channels: spe: 90% sen: 95%	10-fold CV	23s	2019	[140]
Medium weight deep CNN	CHB-MIT	acc: 96%	10-fold CV	300ms, 20ms overlap	2022	[95]
Hybrid Probabilistic Graphical Model CNN (PGM-CNN)	CHB-MIT, Johns Hopkins Hospital (JHH)	TPR: 0.61 FPR: 0.0067 AUC:0.8 F1: 0.67 pre: 0.83	5-fold CV	1s	2019	[88]
GCN	CHB-MIT	acc: 98.35%	10-fold CV	60s	2020	[105]
3D-CNN	CHB-MIT, TUH	CHB-MIT: acc: 94.36% rec: 95.57% TUH: acc: 92.26% rec: 93.86%	N/A	2s	2021	[151]



Table 8. continued - DL methods from systematic review of literature for seizure detection in scalp EEG data

Classifier	Dataset(s)	Performance	Validation	Segment Length	Year	References
Deep CNN	CHB-MIT, Bonn	acc: 91.82% sen: 91.93% F1: 95.73% FRP: 0.005/hour	three-way holdout	5s no overlap	2021	[152]
CNN	CHB-MIT	acc: 93.4%	6-fold CV	1s	2022	[153]
2D-PCANet (feature extraction) SVM (classification)	CHB-MIT Bonn	acc: 98.47% sen: 98.28% spe: 98.50%	10-fold CV	1s	2021	[143]
GBDT, attention-based CNN-BiRNN, FC layer for classification	CHB-MIT	acc: 97.56% sen: 90.97% spe: 91.93%	train-val-test 70-15-15	20s	2021	[154]
CNN, LSTM	CHB-MIT (22 patients, 8 channels)	acc: 94.6% rec: 97.15% pre: 95.78%	10-fold CV	N/A	2022	[146]
1D CNN	CHB-MIT (21 channels)	acc: 97.09% sen: 96.49% spe: 97.09%	10-fold CV	2s, 1s overlap	2022	[110]
ResNet-based	TUSZ (20 channels)	acc: 69% segment level, acc: 61.67%	3-fold CV	1s, 0.75s overlap	2022	[70]
Attention based CNN	TUSZ (14 subjects)	sen: 97.4% spe: 88.1% acc: 80.5%	LOO, and 14-fold CV	1s, 0.5s overlap	2022	[155]
U-net(feature extraction), LSTM (classification)	TUSZ (16 channels)	sen: 12.37 Fas/24hr: 1.44 TAES score: 2.46	10-fold CV	20s	2020	[77]
CNN-SVM	CHB-MIT	acc: 98.31%	train-test-val: 70/15/15	N/A	2022	[156]
CNNs, FC layer	CHB-MIT (remove patient 12, 21 channels), TUSZ (28 patients)	CHB-MIT: acc: 96.17% sen: 56.83% spe: 96.97% F1: 38.26% TUSZ: acc: 67.68% sen: 59.21% spe: 75.30% F1: 47.55%	5-fold CV	4s, 1s overlap	2021	[78]
CNN	CHB-MIT (8 channels, 16 patients)	acc: 97.57% sen: 98.90% FPR:2.13% delay:10.46s	LOO	5s, 1s overlap	2023	[157]
CNN	NYP-WC, CHB-MIT	-	5-fold CV	120s, 119 overlap	2019	[147]
CNNs with an attention mechanism	TUH	acc: 86%, F1: 81%	LOO	3s, no overlap	2023	[158]

**Table 9. continued - DL methods from systematic review of literature for seizure detection in scalp EEG data**

Classifier	Dataset(s)	Performance	Validation	Segment Length	Year	References
CNN vs Xception	CHB-MIT	CNN: acc: 98.47%, pre: 99.79%, rec: 98.93%, F1: 98.51% Xception: acc: 95.52%, pre: 99.93%, rec: 98.63%, F1: 97.05%	CV	N/A	2022	[159]
ResNest18	TUSZ	sen: 42.05%, FAR/day : 5.78	CV	250 samples	2021	[75]
Multi-fuse reduced deep CNN (MF-RDCNN)	Bonn, CHB-MIT, Neurology Sleep Centre Delhi	CHB-MIT: acc: 99.29%, sen: 99.29%, spe: 99.86%, FPR: 0.71%	40-40-20 train-test-val	N/A	2022	[160]
Multilayer deep convolutional neural network (MDCNN)	CHB-MIT (18 subjects, all of which have 23 channels)	acc: 71.60%	LOO	1s, 0.5 overlap	2021	[161]
CNN using adversarial network methods	CHB-MIT (18 channels)	acc: 91.71%, sen: 91.09%, spe: 94.73%, FPR: 0.58/hr, latency: 4.45s	LOO	4s, no overlap	2022	[162]
CNN and RNN	CHB-MIT, Bonn, Bern-Barcelona	acc: 96.23%	8-fold CV	N/A	2023	[144]
AttVGGNet-RC	CHB-MIT (23 channels, remove patient 12)	sen: $93.84 \pm 0.63\%$ , spe: $95.84 \pm 0.74\%$ , acc: $95.12 \pm 0.20\%$	10-fold CV	1s, no overlap	2020	[163]
CNN (feature extraction), and ANN, LR, RF, SVM, GB, k-NN, SGD, Ensembles (classification)	CHB-MIT, Bonn	ANN: 94.4%, LR: 91.7%, RF: 92.4%, SVM: 95.7%, GB: 94.6%, k-NN: 96.8%, SGD: 87%, Ensembles: 97%	10-fold CV	5s, no overlap	2022	[164]
ANN	CHB-MIT	N/A	N/A	100s, no overlap	2019	[139]
Asymmetrical Back Propagation Neural Network (ABPN)	CHB-MIT	sen: 96.32%, spe: 95.12%, acc: 98.36%	N/A	N/A	2021	[90]
BERT (LLM)	TUSZ	acc: $\sim 77\%$	N/A	1s, no overlap	2022	[101]
CNViT (Convolutional Vision Transformer) that first uses multi-layer convolution to extract features, and then adopts transformer blocks	CHB-MIT	sen: 96.71%, spe: 97.23%, acc: 97.15%, AUC: 99.54%	N/A	2s, no overlap	2022	[165]
Graph isomorphism network (GIN)	CHB-MIT	acc: 96.2%, sen: 95.4%, spe: 97.0%	10-fold CV	20s, no overlap	2022	[166]

Table 10. continued - DL methods from systematic review of literature for seizure detection in scalp EEG data

Classifier	Dataset(s)	Performance	Validation	Segment Length	Year	References
Graph-generative neural network (GGN)	TUH	acc: 91%	70-30 train-test	5s, no overlap	2022	[106]
GAT and BiLSTM	CHB-MIT, TUH	CHB-MIT: acc: 98.52%, spe: 94.34%, sen: 97.75%, TUH: acc: 98.02%, spe: 99.06%, sen: 97.7%	5-fold CV	1s, 0.5s overlap	2022	[167]
AE (feature extraction), RF (classification)	Siena	F1 ictal: 91%, F1 non-ictal: 90.1%	Leave-2-out	6s, 1s overlap	2021	[168]
Deep convolutional Autoencoder Bi-LSTM	CHB-MIT	sen: 99.7%, acc: 99.8%, spe: 99.9%, precision: 99.9%, F1: 99.6%	10-fold CV	4s, no overlap	2023	[169]
Deep Stacked AE	CHB-MIT, TUEP	TUEP: acc: 91.5%, sen: 85.2%, spe: 86.0%, CHBMIT: acc: 91.4%, sen: 85.5%, spe: 85.3%	90-10 train-test	N/A	2022	[170]
AE	CHB-MIT (13 patients, at least 4 seizure/patient)	Sen: 0.86%, FPR/h: 0.08%	LOO	1s, no overlap	2022	[171]
CNN (feature extraction), LSTM (classification)	TUSZ	acc: 82%, pre: 71.69%, sen: 85%	LOO	N/A	2020	[76]
Scalp Swarm Algorithm (SSA) (feature selection), LSTM (classification)	TUSZ	sen: 98.99%, FDR: 98.43%, spe: 99.01%, acc: 99.2%, F1: 97.54%	80-20 train-test	1s, no overlap	2022	[172]
RNN	CHB-MIT, TUEP	TUEP: acc: 84.7%, sen: 89.2%, spe: 82.2%, CHBMIT: acc: 85.3%, sen: 93.0%, spe: 79.7%	90-10 train-test	N/A	2020	[173]
LSTM	CHB-MIT, Siena	Siena: acc: 92.59%, sen: 94.83%, spe: 96.82%, CHBMIT: acc: 89.88%, sen: 96.71%, spe: 89.88%	10-fold CV	N/A	2019	[174]

**Table 11. continued - DL methods from systematic review of literature for seizure detection in scalp EEG data**

Classifier	Dataset(s)	Performance	Validation	Segment Length	Year	References
ConvLSTM	TUEP	acc: 92.17%, sen: 93.27%, spe: 90.96%, pre: 91.23%, F1: 0.93	5-fold CV and LOO	3s, no overlap	2022	[175]
Convolution Attention Layer, BiRNN classification	CHB MIT (patient 1-11, 14, 20-24)	acc: 97.62%, sen: 96.69%, spe: 98.41%, F1: 97.38%	N/A	1s, no overlap	2022	[176]
AE (feature extraction), RF (classification)	Siena	acc: 97.22%	LOO	6s, 1s overlap	2022	[177]
2D-DCAE (feature extraction), Bi-LSTM (classification)	CHB-MIT (16 patients)	acc: $98.79 \pm 0.53\%$ , sen: $98.72 \pm 0.77\%$ , spe: $98.86 \pm 0.53\%$ , pre: $98.86 \pm 0.53\%$ , F1: $98.79 \pm 0.53\%$	10-fold CV	4s, no overlap	2021	[27]

## 5.5 Post-classification processing

We use the term post-classification processing to describe steps applied after a ML model has made predictions on the input data. These steps include refining, improving, or interpreting the model’s output before presenting it to the end-user. An important role for post-classification processing is to amalgamate classifier outputs. This amalgamation can be across channels, in the case of algorithms trained on single-channel inputs applied to multi-channel data, and across time. This can be advantageous for increasing accuracy. However, it may increase latency when multiple segments are required to produce a final output.

It is generally accepted that in a clinical setting, high sensitivity and low false detection rate should take precedence over low latency. Indeed, devices producing false alarms can lead to a phenomenon known as ‘alarm fatigue’, and are unlikely to be deployed in a clinical setting [178]. Post-classification steps help in reducing the rate of false alarms, and are therefore important for the clinical translation of seizure detection algorithm. However, a minority of the encountered studies used any post-processing.

As mentioned above, post-processing steps for seizure detection algorithms can be divided into two categories: spatial and temporal. This is linked to the distinction between segment based and event based performance.

1. Spatial concatenation: combining the predictions of different channels if the algorithm outputs a prediction per channel. For instance, only sending an alarm when 3 or more channels are classified as being in a seizure state within 5 seconds of each other.
2. Temporal concatenation: combining the predictions of consecutive EEG segments. This is applied by Khalkhali et al., who set  $BDmin$  and  $SDmin$  as a threshold for the minimum background and seizure duration, respectively [75]. All background events, in between seizure events, with a duration less than  $BDmin$  are converted to seizure events. All seizure events, in between background events, with a duration less than  $SDmin$  are converted to background

events. Using this method, the maximum seizure detection delay is defined as:  $BD_{min} + SD_{min}$ . Wei and Mooney, who arrived in fifth position at the 2020 Neureka<sup>TM</sup> challenge (a seizure prediction challenge using the v.1.5.2 TUSZ dataset), grouped seizure labels if they were less than 2 seconds apart and classified an event as a seizure if it was longer than 15 seconds [89].

Chatzichristos et al. won the 2020 Neureka<sup>TM</sup> competition [77]. They merged consecutive events and discarded short events following three rules. Firstly, seizure events less than 30 seconds apart are merged together. Secondly, merged seizure events, for which the mean probability of being a seizure is less than 82%, are rejected. Finally, seizure events of duration less than 15 seconds are rejected.

Most post-processing steps do not incorporate domain specific knowledge about seizures. For example, in patients on anti-seizure medication, the median length of a generalized tonic-clonic seizure is 79.5 seconds, whereas that of a focal tonic seizure is 15 seconds [46]. In the case of Chatzichristos et al., half of the tonic-clonic seizures would be missed.

Embedding medical and contextual information within the post-processing design could improve performance. For example, in the case of a seizure detection device used in the ICU, post-processing should lower the rate of false positive and increase the sensitivity to status epilepticus (see section 2.4.5). However, when keeping a seizure diary, logging shorter seizures may be useful. Indeed, seizure types like generalised tonic seizures have a median duration of 8 seconds, are important to record and monitor [46]. Overall, post-processing steps require tuning to work best in different clinical settings where different forms of seizures are expected, tolerated and monitored.

## 5.6 Metrics for performance evaluation

Metrics such as accuracy, sensitivity and specificity are good bench marking tools, and allow objective comparisons of algorithms. However, to provide an assessment of the clinical usability of an algorithm, it is important to use a combination of metrics that provide information about the expected behaviour of the algorithm in a clinical setting such as the rate of false positives. The majority of studies do not include such clinically relevant metrics.

For an exhaustive description of objective evaluation metrics in the field of automatic classification of EEG segments, see [179]. This section describes performance metrics and their clinical relevance with regards to seizure detection. The two-by-two confusion matrix in Figure 7 illustrates the binary classification results of seizure detection algorithms, from which several metrics are derived. Table 13 (Appendix) shows how accuracy, sensitivity, specificity, precision and F1 relate to the values of the confusion matrix.

### 5.6.1 Traditional evaluation metrics

*Accuracy* is the most common metric for evaluating classification algorithms. It is defined as the proportion of correctly classified instances out of the total number of instances in the dataset. While accuracy is the most commonly used metric for evaluating the performance of classification models, it can be misleading. For example, if the test data is imbalanced and contains 1% of ictal segments and 99% of non-ictal segments, a model classifying all the instances as being non-ictal will have an accuracy of 99%.

		Actual value	
		Ictal (+)	Non-ictal (-)
Predicted value	Ictal (+)	TP	FP
	Non-ictal (-)	FN	TN

**Figure 7. Confusion matrix of a binary seizure detection algorithm.** Ictal segments are labeled as being positive and non-ictal segments as negative. TP = true positive, FP = false positive, FN = false negative, TN = true negative. [9]

As discussed previously, EEG datasets used to test and train seizure detection algorithms are often imbalanced. Most studies balance the training data. However, this is not always the case for test data, where some authors try to mimic real likelihoods of ictal and non-ictal activity, i.e. reflect the proportion of seizure/non-seizure activity, thereby using imbalanced data. In this case, metrics other than accuracy should be reported.

Even with balanced test data, the accuracy metric weights all classification errors (false positives, false negatives) equally, which does not necessarily reflect the real-world implications of different errors. In many applications of seizure detection algorithms, incorrectly classifying an ictal segment as non-ictal is worse than incorrectly classifying a non-ictal segment as an ictal segment. As such, it is necessary to use additional metrics, to assess algorithm performance in a way that corresponds better to their clinical usability.

*Sensitivity*, also known as *true positive rate* or *recall*, is the proportion (or percentage) of positive instances that are correctly classified. A model with high sensitivity can effectively identify positive instances (ictal segments), i.e., true positives (TP). However, sensitivity is unrelated to how many negative instances are correctly labelled (TN). For example, if an algorithm assigns a positive label to all segments, it will have a sensitivity of 1. Sensitivity is independent to data imbalance as it only depends on the classification of the positive class, which is often the minority class.

*Specificity*, also known as *selectivity* or *true negative rate*, is complementary to sensitivity. Specificity is the proportion or percentage of the negative class (non-ictal) correctly predicted to be negative (TN).

*Precision* is the proportion of positive predictions that are true positives (TP). Unlike sensitivity and specificity, precision is affected by class imbalance as it considers the number of negative samples incorrectly labeled as positive.

The *F1 score* combines precision and recall. A high F1 score indicates that the model is good at identifying positive cases whilst avoiding false positives, while a low F1 score indicates that the model is not performing well in one or both of these areas. The formula for F1 score is given in Table 13 (Appendix).

The *receiver operating characteristic curve (ROC curve)* plots the true positive rate (sensitivity) over the false positive rate (1 - specificity) for different classification thresholds, depicting the trade-off between correctly classified positive samples and incorrectly classified negative samples. Some studies integrate the ROC curve to give the *area under the receiver operating characteristic curve (AUC-ROC)*, which is a metric between 0 and 1. A value of 0.5 for a binary classifier indicates a performance no better than random guessing, and a value of 1 indicates perfect classification. However, some argue that AUC-ROC is misleading when dealing with imbalanced data [180].

A more appropriate metric than the AUC-ROC when dealing with imbalanced datasets is the *area under the precision recall (AUC-PR) curve*. The AUC-PR plots the precision over recall for different classification thresholds. This area is more representative of the algorithm’s performance than AUC-ROC because precision is more sensitive to the presence of false positives than the FPR. In two papers from Salafian et al., the AUC-ROC is 13.79 [87] and 33.42 [86] percent higher than the AUC-PR, highlighting the importance of the metric choice on the apparent performance. Interestingly, both studies use the CHB-MIT dataset, but the ratio of ictal to non-ictal segments vary. In [87], for every second of seizure data, there are 6 seconds of non-seizure data whereas in [86] for every second of seizure data, there are 20 seconds of non-seizure data. The study using the most imbalanced dataset ([86]) is also the one displaying the highest gap between the AUC-ROC and AUC-PR.

### 5.6.2 Clinically-focused evaluation metrics

The metrics mentioned above are not sufficient to assess the performance of an algorithm developed for clinical application. There is a growing need for a set of standardised scoring metrics that reflects the needs of end-users.

*Latency* is often used to evaluate the performance of seizure detection algorithms. Latency is the duration between the labelled start time of the seizure and the triggering of the seizure onset alarm. Latency is relevant when using online, or ‘live’, seizure detection. Reducing latency enables quicker medical intervention, but often results in a higher rate of false positives – instances where the algorithm incorrectly signals a seizure onset. For instance, Khalkhali et al.’s model has a detection latency of 300ms but a 42.05% sensitivity and a rate of 5.78 false alarms per day [75]. Some studies point to a very low detection latency in support of their algorithm’s efficacy. However, the associated disadvantages for accuracy may reduce an algorithm’s clinical utility. Additionally, any “arms race” towards faster latency seizure detection should be tempered by the fact that it can be difficult to objectively define the onset of a seizure, with different experts labelling the start of the seizure at different times. Therefore, small reductions in latency, particularly if they compromise accuracy, do not necessarily represent a clinical advancement. Nevertheless, latency can be important for evaluating the clinical relevance of an algorithm and we encourage authors to provide this measure.

Another important metric for successful clinical deployment is the rate of false alarms (FAs) [181]. A high rate of false alarms can desensitise medical staff, a phenomenon known as “alarm fatigue” which leads clinicians to ignore critical alarms, potentially endangering patients [182, 183]. Lower rates of false positives are needed for clinical deployment. The false positive rate, often expressed as per 24 hours, is one of the most important metrics and should always be reported.

*Time-aligned event scoring (TAES)* is a metric that was developed for the Neureka<sup>TM</sup> epilepsy challenge. The metric combines sensitivity, the rate of false positives, and the proportion of channels used (out of 19 EEG channels). TAES is calculated as follows:

$$\text{TAES} = \text{Sens} - \alpha \times \text{FAs}_{24\text{hr}} - \beta \times N/19$$

Where, Sens is the sensitivity,  $\text{FAs}_{24\text{hr}}$  is the number of false positives every 24 hours, then  $\alpha$  and  $\beta$  are constants defined by the organising committee (set to 2.5 and 7.5 respectively). Incorporating an extra EEG channel may improve the algorithm’s sensitivity and reduce the rate of false alarms, but will be penalised for increased complexity. The higher the TAES, the better.



## 6 Outstanding research questions

In this section, we highlight the key research avenues for clinical translation of automated seizure detection algorithms identified after reviewing the selected articles.

### 6.1 Generalisability

The main challenge in the field of seizure detection remains the generalisability of algorithms. A robust algorithm should be able to generalise across individuals; patient demographics; EEG montages, electrode number and placement; datasets; and seizure types.

Whilst some research approaches imply that making an algorithm patient specific is necessary to achieve good performance, we do not find this is reflected in the recent literature. For instance, Qiu, Wang, and Jiao developed LightSeizureNet (LSN), a deep learning model for seizure detection [110]. They compared the patient dependent and patient independent version of LSN. The patient specific model achieves 99.77% accuracy, 97.11% sensitivity, and 99.78% specificity, with 113,800 parameters and 3.7 million multiply-accumulate operations (MACs, a computational complexity metric). In comparison, the patient independent model achieves 97.09% accuracy, 96.49% sensitivity, and 97.09% specificity, with 198,300 parameters and 6.2 million MACs.

A number of algorithms have been designed to remove inter-patient information. Zhang et al. proposed an adversarial network to decompose scalp EEG data into patient-related and seizure-related latent spaces [155], thereby reducing inter-patient contributions, allowing a patient independent analysis of seizure activity. Similarly, Nasiri and Clifford rely on adversarial network methods to learn patient invariant representations and reduce patient specific variations [162]. Jiang, Xu, and Chen use transfer component analysis (TCA) to reduce the impact of individuals on EEG characteristics. TCA maps the original high-dimensional feature space to a lower dimensional subspace, where the data of individuals all follow the same distribution. The low-dimensional features are then input to a k-NN algorithm for classification [74]. Thuwajit et al. employ transfer learning to fine-tune a patient independent model to better fit a specific individual's data [78]. By using approximately 1 hour of an individual's data, the models accuracy increases from 88.41% ( $\pm 15.23$ ) to 96.69% ( $\pm 3.59$ ) on that individual.

In summary, patient independent algorithms seem to reach performances similar to that of patient specific algorithms. It is also clear that the implementation of patient dependent algorithms is far more limited. Indeed, patient dependent algorithms can be useful for analysing long-term EEG recordings obtained using, for example, subcutaneous or intracranial EEG recordings. However, in a setting such as the ICU, algorithms should be patient independent as there is no time for training the model using the patient's data. Future research should focus on patient independent solutions, using a validation scheme mimicking a clinical environment such as the leave-one-out (LOO) validation strategy (see section 5.1.4).

Although inter-patient generalisability is addressed in the literature, only one of the studies addresses generalisability across datasets. Typically, algorithms are trained and tested on a single dataset. This can lead to overfitting, for example, to the EEG hardware and recording technique used in a particular dataset. This possibility was confirmed by a study conducted by Yan et al. [147]. Yan et al. trained a CNN using spectrograms obtained from the CHB-MIT dataset and tested it on (1) a held out subset of the CHB-MIT dataset and, (2) private data, which they collected at the New York Presbyterian – Weill Cornell Medical Center (NYP-WC). For seizures

visible on a spectrogram, they achieved sensitivity and specificity of  $>90\%$  using the CHB-MIT test set, and  $>90\%$  sensitivity but only 75–80% specificity using the NYP-WC test set. The reduction in performance from changing the test dataset suggests that the model was overfitting to some of the CHB-MIT’s recording properties. We recommend training and testing algorithms on multiple datasets.

In addition to generalisability across patients and datasets, algorithms should generalise to all seizure types. Despite the qualitative differences between seizure types clinically (e.g. focal seizures versus generalised seizures), and their differences in terms of standard EEG interpretation, the possibility of variability of accuracy across all seizure types is often overlooked in the seizure detection literature. Although not the main focus of this review, a few studies evaluated classification performance in relation to specific seizures types. For example, using a type of wavelet decomposition, fractal analysis, and a SVM for classification, Tang, Zhao, and Wu achieve an accuracy  $>95\%$  for all seven seizure subtypes in the TUSZ dataset [184]. A clinically very important form of seizure is that of status epilepticus, a medical emergency (see section 2.4.5). None of the encountered public seizure datasets contain this form of seizure. Researchers should seek to include this type of seizures in their research. Given the scarcity of recordings of certain seizure types, using data augmentation techniques to increase the the relative number of EEG recordings associated with rare types of seizures may be one approach to improve the performance of algorithms.

## 6.2 Variability in the ground truth labels

There is subjectivity in the seizure labelling done by clinicians, which means that the ground truth labels used for training are subjective. In a 2009 study, Ronner et al. asked nine clinicians with different levels of experience to evaluate 90 epoch (10s each) of 30 EEG recordings of 23 different patients admitted to the ICU [13]. For each epoch, clinicians had to decide whether there was an electrographic seizure or not. The results show a limited inter-observer agreement. The labelling of the more experienced clinicians obtained a Kappa value of 0.5 (moderate agreement) and that of the less experienced clinicians a Kappa value of 0.29 (fair agreement). Quantifying the labelling variance would require access to large amounts of labelled data from a wide range of clinicians, but could help develop more robust algorithms.

In April 2023, Jing et al. asked 30 clinicians to scored varying numbers of ten-second EEG segments as seizures (SZ), generalized periodic discharges (GPDs), lateralized periodic discharges (LPDs), lateralized rhythmic delta activity (LRDA), generalized rhythmic delta activity (GRDA), or other [185]. In total, clinicians scored 50,697 EEG segments. Results show that the average percent agreement with group consensus (i.e. majority voting) was 65%. They suggest that the observed variations in labelling are due to variations in decision thresholds, rather than level of experience. In May 2023, they released what seemed to be a slightly modified version of this dataset as part of the harmful brain activity classification contest [185, 186]. The dataset contains 71,982 10 seconds long scalp EEG segments that have been independently annotated by 20 fellowship-trained neurophysiologists. This is the first dataset to enable quantification of labelling heterogeneity, which could be very useful to develop probabilistic classification algorithms. As of April 2024, to the best of our knowledge, it was retracted and has not yet been re-released as the authors are double-checking data quality.

### 6.3 Robustness to non-ictal activity

An emerging challenge for automated seizure detection research is the *ictal-interictal continuum (IIC)*, which introduces a new dimension of complexity for algorithm development. The IIC covers a range of EEG signatures typically associated with critically ill patients. IIC patterns can be sharp, rhythmic or periodic and risk being incorrectly classified incorrectly as ictal. For example, in critically ill patients one may encounter abnormal and often stereotyped bursts of striking epileptiform activity. If these bursts last less than 10 seconds, they are known as brief (potentially) ictal rhythmic discharges (BIRDs) and are not classified as seizures [187]. Critically ill patients may often suffer from both frequent seizures (>10 seconds) with BIRDs (<10 seconds) in between. BIRDs are an intermediate phenomenon between short runs of standard interictal epileptiform discharges and definitive electrographic seizures. They are regarded as a marker of cortical irritability or hyperexcitability, and are on the IIC. It is important the seizure detection algorithms are robust to the presence of IIC patterns, to ensure that IIC segments are not misclassified as ictal activity [188].

Finally, testing algorithms using scalp EEG recordings of healthy patients would be useful to assess their specificity. A number of publicly available datasets could be used for this purpose. For example, the Leipzig study for mind-body-emotion interactions (LEMON) dataset is composed of 227 healthy participants across a range of ages, making it a potentially suitable algorithmic control test set [189].

### 6.4 Virtual brain

All the ML algorithms encountered in this review use a data-based approach for identifying ictal/non-ictal instances. In contrast, there is increasing interest in a “bottom-up” modelling approach to seizures, based on biophysical simulations of neural activity coupled to an EEG forward model. The Virtual Brain (TVB) is a publicly available toolkit for such mesoscopic simulations [190]. In the future, these models may provide valuable insights into the relationship between neural activity and observable EEG patterns. Reproducing observed EEG signals/features using these models could shed light on the underlying dynamics of seizure genesis and propagation. The application of such models to individual patients is the basis of the so-called “digital twin” approach, where the TVB has seen validation in the context of epilepsy surgery planning [191, 192]. Additionally, a bottom-up modeling approach could supplement limited empirical data by generating synthetic EEG.

### 6.5 Wearable recording devices

Traditional clinical EEG is costly and requires experts to set up and interpret. Hence, EEG in most healthcare settings is not frequently used, even in ICUs, where seizure occurrence is high [193]. Using several continuous EEG monitoring studies published between 1994 and 2011, Westover et al. established the occurrence of seizures in the ICU to be between 8-34% [194]. As part of their study, Claassen et al. found that of 570 patients who had continuous EEG monitoring, 19% (110) had seizures, 92% of which were non-convulsive (102). Out of those 110 patients, 89% (98) were in ICU at the time of monitoring [195].

Technological advancements in wearable EEG devices is promising for addressing these issues and revolutionising EEG monitoring. One such portable device, the point-of-care EEG (POC-EEG) by

Ceribell, consists of a headband with ten electrodes connected to a small battery powered recorder equipped with a screen for real-time EEG streaming (see <https://ceribell.com>). In a single centre cohort study, Rajshekar et al. found that in 72% of patients monitored, POC-EEG was thought to have expedited diagnostic testing and/or treatment [196]. Another alternative is in-ear or behind-the-ear EEG recording devices. The simple setup of ear-EEG devices may enable more widespread use of EEG in hospital and community settings. For this reason, algorithm design and evaluation for such devices will become increasingly important. The ICASSP 2023 Seizure Detection Challenge is currently the only seizure detection competition using behind-the-ear EEG data [197].

## 7 Conclusion and algorithm development guidelines

The worldwide economic burden associated with epilepsy in 2019 was estimated at \$119.27 billion [198]. Whilst ML based algorithms have demonstrated significant potential for managing seizure disorders, there remain technological, clinical, and regulatory challenges for their successful translation.

In this review, we have provided a review of the key considerations for creating a robust, accurate, and clinically relevant EEG-based seizure detection system. Based on this analysis, we have developed the following recommendations for creating more clinically informed seizure detection algorithms.

### 7.1 Identify the intended clinical use(s) of the algorithm

Applications for seizure detection algorithms range widely, from (1) highlighting to clinicians sections of interest in long recordings to facilitate annotation offline; (2) real-time seizure detection using continuous EEG in ambulatory patients, telemetry units, or ICU; to (3) automatic recording of seizure diary entries. Identifying the specific intended use case(s) of the algorithm should be seen as a prerequisite for designing the algorithm.

### 7.2 Leverage multi-domain data representations

To optimally distinguish between ictal and non-ictal EEG segments, researchers should leverage a range of EEG data representations, i.e. in the time domain as well as spectral and network domains. Yan et al. showed that there is no guarantee that all seizures will be spectrographically visible (i.e. using spectral domain-based features), and subtle seizures, although visible on the EEG, may not be distinguishable from the spectrogram background. As such, we recommend using multiple domain representations as it is likely to improve performance.

A number of high throughput feature extraction libraries support multi-domain data representations. For example, the MATLAB package *hctsa* extracts 7,749 different features from time-series data [199]. This approach can help identify discriminating features for a given classification task, and provides an exhaustive overview of a signal’s properties. When using feature-based ML techniques, we recommend using similar tools during the feature selection process, or providing the rationale behind the choice of features.

### 7.3 Use relevant performance metrics

The performance metrics encountered in the literature are varied and reflect different aspects of algorithmic performance. We recommend providing all metrics in Table 13, latency, and the rate of false positives.

### 7.4 Explore post-classification processing

The input rate of EEG segments into an algorithm is typically less than 5 seconds. In post-processing, the use of multiple successive predictions increases the algorithm’s ability to correctly classify ictal/non-ictal activity. We suggest evaluating the effects of post-classification processing, and using multiple consecutive EEG segments during classification. We also encourage authors to

provide the rationale behind their choice of post-processing steps, as they may be influenced by the use-case of an algorithm.

## 7.5 Use diverse training data

We recommend training and testing algorithms on a range of datasets. This approach will enable greater understanding of which features and/or architectures are most independent of patients, seizure types and hardware recording devices. Finally, the results could be validated using additional seizure-free control datasets like the LEMON dataset (healthy participants) [189]. For generalised medical applications, using multiple datasets reduces sensitivity to recording equipment, cohort community, and provides a larger dataset.

## 7.6 Use a leave-one-out validation strategy

We recommend using a leave-one-patient-out cross-validation strategy as it simulates real-world scenarios effectively. Excluding the data of one patient from the training set and using cross-validation evaluates the ability of the model to generalise across the entire patient cohort. Using a random train/test split may introduce sampling bias and overfitting.

## 7.7 Respect the assumption of independence

A common pitfall in ML is the violation of the assumption of independence. The assumption of independence states that the data used for model training and evaluation is independent and identically distributed. This is a necessary condition to achieve unbiased estimates of performance metrics and generalisation error of the model. To avoid this common pitfall, oversampling, feature selection, and data augmentation should always be performed after the data was split into training and evaluation sets [200].

## 7.8 Model architecture and pre-processing steps

Regarding the choice of model architecture, segment length, channels to use and other pre-processing steps; it is challenging to identify an optimal strategy. In fact, there is likely no ‘one size fits all’ strategy, as the use case of the algorithm influences design choices drastically.

However, we found that segment lengths of 3-5 seconds (with a sampling rate  $>170$  Hz) were good candidates for seizure detection as they are quasi-stationary and carry sufficient information. We recommend using a notch filter to remove the power line interference and a band-pass filter between 0.3-60 Hz. For channel selection, we recommend using all of the available channels for seizure detection. Depending on the type of seizure and its location, any of the 10-20 electrodes can pick up seizure activity. Given the absence of a subset of channels that is independent of patient and seizure types, and that also contains all necessary information to replicate detection results achieved with the complete set of channels, we recommend against discarding any channel data.

Both feature-based ML and DL approaches can achieve high seizure detection performances. We found that the average detection accuracy using feature-based ML was 94.3% (95% confidence interval:  $94.3 \pm 2.3$ ), compared to 91.8% (95% confidence interval:  $91.8 \pm 2.3$ ) using deep learning architectures.

Finally, we highly recommend tailoring the algorithm to the intended clinical use. Whilst there is likely no ‘one size fits all’ solution, we recommend providing justification behind the choice of features, pre-processing steps, model architecture, validation strategy, evaluation metrics and dataset(s).

## 8 Appendix

**Table 12. Search strategy for each repository.** Each search string corresponds to equivalent inclusion criteria. They are formatted to fit the database specific search tool’s syntax.

Databases	Search string	Number of papers
Scopus	( TITLE-ABS-KEY ( eeg AND detection AND epilepsy OR seizure ) AND NOT TITLE-ABS-KEY ( animal OR neonatal OR infant OR intracranial OR iieg ) ) AND PUBYEAR > 2018 AND PUBYEAR < 2024 AND ( LIMIT-TO ( DOCTYPE , "cr" ) OR LIMIT-TO ( DOCTYPE , "cp" ) OR LIMIT-TO ( DOCTYPE , "ar" ) ) AND ( LIMIT-TO ( LANGUAGE , "English" ) )	1886
PubMed	((eeg AND detection) AND (epilepsy OR seizure) NOT (animal OR neonatal OR infant OR intracranial OR iEEG) AND (("2019"[Date - Publication] : "2023-02-22") AND (english[Language]))	1078
WoS	(EEG detection) AND (epilepsy OR seizure) NOT (animal OR neonat* OR infant OR intracranial OR iEEG) (Topic) and Article OR Book Chapter OR Book OR Proceedings Paper (Document Type) and English (Language)	1551
<b>TOTAL</b>		<b>4515</b>

**Table 13.** Table of common metrics used to evaluate the performance of ML algorithms. Note: sensitivity is also known as *True Positive Rate* and specificity as *True Negative Rate*

Name	Formula	Description
Accuracy	$\frac{TP + TN}{TP + TN + FN + FP}$	Proportion of correctly classified instances or correct predictions out of the total number of instances in the dataset
Sensitivity, Recall	$\frac{TP}{TP + FN}$	Evaluate how well a classification algorithms identifies positive instances
Specificity	$\frac{TN}{TN + FP}$	Evaluate how well a classification algorithms identifies negative instances
Precision	$\frac{TP}{TP + FP}$	Proportion of positive predictions that are actually true positives
F1	$\frac{2 * Precision * Recall}{Precision + Recall}$	It combines both precision and recall to provide a balanced view of a model’s performance



## References

- [1] Aidan Neligan, JW Sander, et al. “The incidence and prevalence of epilepsy”. In: *London: UCL Institute of Neurology* (2009).
- [2] Robert S Fisher et al. “ILAE official report: a practical clinical definition of epilepsy”. In: *Epilepsia* 55.4 (2014), pp. 475–482.
- [3] Jessica J Falco-Walter and Thomas Bleck. “Treatment of established status epilepticus”. In: *Journal of Clinical Medicine* 5.5 (2016), p. 49.
- [4] Christian E Elger and Christian Hoppe. “Diagnostic challenges in epilepsy: seizure under-reporting and seizure detection”. In: *The Lancet Neurology* 17.3 (2018), pp. 279–288.
- [5] Ikuko Laccheo et al. “Non-convulsive status epilepticus and non-convulsive seizures in neurological ICU patients”. In: *Neurocritical care* 22 (2015), pp. 202–211.
- [6] Peter W Kaplan. “Assessing the outcomes in patients with nonconvulsive status epilepticus: nonconvulsive status epilepticus is underdiagnosed, potentially overtreated, and confounded by comorbidity”. In: *Journal of Clinical Neurophysiology* 16.4 (1999), pp. 341–352.
- [7] Johannes Jungilligens, Rosa Michaelis, and Stoyan Popkirov. “Misdiagnosis of prolonged psychogenic non-epileptic seizures as status epilepticus: epidemiology and associated risks”. In: *Journal of Neurology, Neurosurgery & Psychiatry* 92.12 (2021), pp. 1341–1345.
- [8] Vaidehi Naganur et al. “Automated seizure detection with noninvasive wearable devices: A systematic review and meta-analysis”. In: *Epilepsia* 63.8 (2022), pp. 1930–1941.
- [9] *Created with BioRender*. 2023. URL: <https://www.biorender.com>.
- [10] Shelagh JM Smith. “EEG in the diagnosis, classification, and management of patients with epilepsy”. In: *Journal of Neurology, Neurosurgery & Psychiatry* 76.suppl 2 (2005), pp. ii2–ii7.
- [11] Mitesh Patel, Manny Bagary, and Dougall McCorry. “The management of convulsive refractory status epilepticus in adults in the UK: no consistency in practice and little access to continuous EEG monitoring”. In: *Seizure* 24 (2015), pp. 33–37.
- [12] Paula A Gerber et al. “Interobserver agreement in the interpretation of EEG patterns in critically ill adults”. In: *Journal of Clinical Neurophysiology* 25.5 (2008), pp. 241–249.
- [13] HE Ronner et al. “Inter-observer variability of the EEG diagnosis of seizures in comatose patients”. In: *Seizure* 18.4 (2009), pp. 257–263.
- [14] John P Ney et al. “Continuous and routine EEG in intensive care: utilization and outcomes, United States 2005–2009”. In: *Neurology* 81.23 (2013), pp. 2002–2008.
- [15] Nicholas S Abend, Alexis A Topjian, and Sankey Williams. “How much does it cost to identify a critically ill child experiencing electrographic seizures?” In: *Journal of clinical neurophysiology: official publication of the American Electroencephalographic Society* 32.3 (2015), p. 257.
- [16] *Health Committee, Written evidence from the Joint Epilepsy Council (ETWP 33)*. <https://publications.parliament.uk/pa/cm201213/cmselect/cmhealth/6/6vw28.htm>. Accessed: 2023-08-15.
- [17] Alice D Lam et al. “Widespread changes in network activity allow non-invasive detection of mesial temporal lobe seizures”. In: *Brain* 139.10 (2016), pp. 2679–2693.
- [18] Ingrid E Scheffer et al. “ILAE classification of the epilepsies: Position paper of the ILAE Commission for Classification and Terminology”. In: *Epilepsia* 58.4 (2017), pp. 512–521.
- [19] H Lüders et al. “Semiological seizure classification.” In: *Epilepsia* 39.9 (1998), pp. 1006–1013.

- [20] Hans Lüders et al. “Critique of the 2017 epileptic seizure and epilepsy classifications”. In: *Epilepsia* 60.6 (2019), pp. 1032–1039.
- [21] Eugen Trinka and Markus Leitinger. “Which EEG patterns in coma are nonconvulsive status epilepticus?” In: *Epilepsy & behavior* 49 (2015), pp. 203–222.
- [22] Giridhar P Kalamangalam and Bernhard Pohlmann-Eden. “Ictal–interictal continuum”. In: *Journal of Clinical Neurophysiology* 35.4 (2018), pp. 274–278.
- [23] Arthur C Grant et al. “EEG interpretation reliability and interpreter confidence: a large single-center study”. In: *Epilepsy & Behavior* 32 (2014), pp. 102–107.
- [24] Mark L Scheuer et al. “Seizure detection: interreader agreement and detection algorithm assessments using a large dataset”. In: *Journal of Clinical Neurophysiology* 38.5 (2021), p. 439.
- [25] Brandon Foreman et al. “Generalized periodic discharges in the critically ill: a case-control study of 200 patients”. In: *Neurology* 79.19 (2012), pp. 1951–1960.
- [26] Philippe Gelisse et al. “Lateralized Periodic Discharges: Which patterns are interictal, ictal, or peri-ictal?” In: *Clinical Neurophysiology* 132.7 (2021), pp. 1593–1603.
- [27] Ahmed Abdelhameed and Magdy Bayoumi. “A deep learning approach for automatic seizure detection in children with epilepsy”. In: *Frontiers in Computational Neuroscience* 15 (2021), p. 650050.
- [28] Syed Muhammad Usman et al. “Detection of preictal state in epileptic seizures using ensemble classifier”. In: *Epilepsy Research* 178 (2021), p. 106818.
- [29] Afef Saidi, Slim Ben Othman, and Slim Ben Saoud. “A novel epileptic seizure detection system using scalp EEG signals based on hybrid CNN-SVM classifier”. In: *2021 IEEE Symposium on Industrial Electronics & Applications (ISIEA)*. IEEE. 2021, pp. 1–6.
- [30] Ozlem Karabiber Cura and Aydin Akan. “Epileptic EEG Classification Using Synchrosqueezing Transform and Machine Learning”. In: *2020 Medical Technologies Congress (TIPTE-KNO)*. IEEE. 2020, pp. 1–4.
- [31] Priya Das and Sarita Nanda. “A novel multivariate approach for the detection of epileptic seizure using BCS-WELM”. In: *International Journal of Information Technology* 15.1 (2023), pp. 149–159.
- [32] Michal Teplan et al. “Fundamentals of EEG measurement”. In: *Measurement science review* 2.2 (2002), pp. 1–11.
- [33] Ronald P Lesser, Nathan E Crone, and WRS Webber. “Subdural electrodes”. In: *Clinical neurophysiology* 121.9 (2010), pp. 1376–1392.
- [34] Josef Parvizi and Sabine Kastner. “Human intracranial EEG: promises and limitations”. In: *Nature neuroscience* 21.4 (2018), p. 474.
- [35] Simon Shorvon et al. *Oxford textbook of epilepsy and epileptic seizures*. OUP Oxford, 2012.
- [36] Marc J Casale et al. “The sensitivity of scalp EEG at detecting seizures—a simultaneous scalp and stereo EEG study”. In: *Journal of Clinical Neurophysiology* 39.1 (2022), pp. 78–84.
- [37] Pedro F Viana et al. “230 days of ultra long-term subcutaneous EEG: seizure cycle analysis and comparison to patient diary”. In: *Annals of clinical and translational neurology* 8.1 (2021), pp. 288–293.
- [38] Revati Shriram, M Sundhararajan, and Nivedita Daimiwal. “EEG based cognitive workload assessment for maximum efficiency”. In: *Int. Organ. Sci. Res. IOSR* 7 (2013), pp. 34–38.
- [39] Adapted from “*EEG Electrode Positioning (10/20 System)*”, by BioRender.com. 2023. URL: <https://app.biorender.com/biorender-templates>.

- [40] Shiang Hu et al. “How do reference montage and electrodes setup affect the measured scalp EEG potentials?” In: *Journal of neural engineering* 15.2 (2018), p. 026013.
- [41] Ernst Niedermeyer and FH Lopes da Silva. *Electroencephalography: basic principles, clinical applications, and related fields*. Lippincott Williams & Wilkins, 2005.
- [42] Jeffrey W Britton et al. “Electroencephalography (EEG): An introductory text and atlas of normal and abnormal findings in adults, children, and infants”. In: (2016).
- [43] Giulia Attard Navarro and Khalid Hamandi. “Lessons from the video-EEG telemetry unit”. In: *Practical Neurology* 22.4 (2022), pp. 301–310.
- [44] Andrew Lawley et al. “The role of outpatient ambulatory electroencephalography in the diagnosis and management of adults with epilepsy or nonepileptic attack disorder: a systematic literature review”. In: *Epilepsy & Behavior* 53 (2015), pp. 26–30.
- [45] Nikesh Bajaj. “Wavelets for EEG Analysis”. In: *Wavelet Theory* (2020), pp. 1–16.
- [46] Pirgit Meritam Larsen et al. “Duration of epileptic seizure types: A data-driven approach”. In: *Epilepsia* 64.2 (2023), pp. 469–478.
- [47] Viktor K Jirsa et al. “On the nature of seizure dynamics”. In: *Brain* 137.8 (2014), pp. 2210–2230.
- [48] *Adapted from "EEG Recording of Seizure (Generalized)" and "Seizure Classification in Epilepsy", by BioRender.com. 2023. URL: <https://app.biorender.com/biorender-templates>.*
- [49] Robert S Fisher, Helen E Scharfman, and Marco DeCurtis. “How can we identify ictal and interictal abnormal activity?” In: *Issues in clinical epileptology: A view from the bench* (2014), pp. 3–23.
- [50] Mustafa Aykut Kural et al. “Criteria for defining interictal epileptiform discharges in EEG: A clinical validation study”. In: *Neurology* 94.20 (2020), e2139–e2147.
- [51] Julia CM Pottkämper et al. “The postictal state—What do we know?” In: *Epilepsia* 61.6 (2020), pp. 1045–1061.
- [52] James X Tao, Xiaoxiao Qin, and Qun Wang. “Ictal-interictal continuum: a review of recent advancements”. In: *Acta Epileptologica* 2 (2020), pp. 1–10.
- [53] John R Hughes. “Periodic lateralized epileptiform discharges: Do they represent an ictal pattern requiring treatment?” In: *Epilepsy & Behavior* 18.3 (2010), pp. 162–165.
- [54] Xiao Jiang, Gui-Bin Bian, and Zean Tian. “Removal of artifacts from EEG signals: a review”. In: *Sensors* 19.5 (2019), p. 987.
- [55] Rongrong Fu et al. “Automatic detection of epileptic seizures in EEG using sparse CSP and fisher linear discrimination analysis algorithm”. In: *Journal of medical systems* 44 (2020), pp. 1–13.
- [56] Chengang Lyu et al. “Automatic epilepsy detection based on generalized convolutional prototype learning”. In: *Measurement* 184 (2021), p. 109954.
- [57] Chuancheng Song et al. “A Feature Tensor-Based Epileptic Detection Model Based on Improved Edge Removal Approach for Directed Brain Networks”. In: *Frontiers in Neuroscience* 14 (2020), p. 557095.
- [58] Ralph G Andrzejak et al. “Indications of nonlinear deterministic and finite-dimensional structures in time series of brain electrical activity: Dependence on recording region and brain state”. In: *Physical Review E* 64.6 (2001), p. 061907.
- [59] PHILIPPA KAROLY et al. “Melbourne NeuroVista Seizure Prediction Trial”. In: (Aug. 2018). DOI: [10.26188/5b6a999fa2316](https://doi.org/10.26188/5b6a999fa2316). URL: [https://melbourne.figshare.com/articles/dataset/Seizure\\_Data/6939809](https://melbourne.figshare.com/articles/dataset/Seizure_Data/6939809).

- [60] UPenn and Mayo Clinic’s Seizure Detection Challenge. 2014. URL: <https://www.kaggle.com/competitions/seizure-detection/data> (visited on 05/18/2023).
- [61] American Epilepsy Society Seizure Prediction Challenge. 2014. URL: <https://www.kaggle.com/competitions/seizure-prediction/data> (visited on 05/18/2023).
- [62] Melbourne University AES/MathWorks/NIH Seizure Prediction. 2016. URL: <https://www.kaggle.com/competitions/melbourne-university-seizure-prediction/data> (visited on 05/18/2023).
- [63] Sheng Wong et al. “EEG datasets for seizure detection and prediction—A review”. In: *Epilepsia Open* (2023).
- [64] Ali Hossam Shoeb. “Application of machine learning to epileptic seizure onset detection and treatment”. PhD thesis. Massachusetts Institute of Technology, 2009.
- [65] Piyush Swami et al. “EEG epilepsy datasets”. In: DOI: <https://doi.org/10.13140/RG.2.14280.32006> (2016).
- [66] Iyad Obeid and Joseph Picone. “The temple university hospital EEG data corpus”. In: *Frontiers in neuroscience* 10 (2016), p. 196.
- [67] Nathan J Stevenson et al. “A dataset of neonatal EEG recordings with seizure annotations”. In: *Scientific data* 6.1 (2019), pp. 1–8.
- [68] Paolo Detti, Giampaolo Vatti, and Garazi Zabalo Manrique de Lara. “Eeg synchronization analysis for seizure prediction: A study on data of noninvasive recordings”. In: *Processes* 8.7 (2020), p. 846.
- [69] Wassim Nasreddine. *Epileptic EEG Dataset*. Mendeley Data, V1, doi: 10.17632/5pc2j46cbc.1. 2021.
- [70] Hrishikesh Tiwary, Piyus Rout, and Arnav Bhavsar. “Deep Learning with Spatial and Channel Attention for Epileptic Seizure Type Classification Using Frequency Characterization”. In: *2022 IEEE 32nd International Workshop on Machine Learning for Signal Processing (MLSP)*. IEEE. 2022, pp. 1–6.
- [71] Vinit Shah et al. “The temple university hospital seizure detection corpus”. In: *Frontiers in neuroinformatics* 12 (2018), p. 83.
- [72] M Mostafizur Rahman and Darryl N Davis. “Addressing the class imbalance problem in medical datasets”. In: *International Journal of Machine Learning and Computing* 3.2 (2013), p. 224.
- [73] Lurong Jiang et al. “Seizure detection algorithm based on improved functional brain network structure feature extraction”. In: *Biomedical Signal Processing and Control* 79 (2023), p. 104053.
- [74] Xinyu Jiang, Ke Xu, and Wei Chen. “Transfer component analysis to reduce individual difference of EEG characteristics for automated seizure detection”. In: *2019 IEEE Biomedical Circuits and Systems Conference (BioCAS)*. IEEE. 2019, pp. 1–4.
- [75] Vahid Khalkhali et al. “Low latency real-time seizure detection using transfer deep learning”. In: *2021 IEEE Signal Processing in Medicine and Biology Symposium (SPMB)*. IEEE. 2021, pp. 1–7.
- [76] A Einizade et al. “A deep learning-based method for automatic detection of epileptic seizure in a dataset with both generalized and focal seizure types”. In: *2020 IEEE Signal Processing in Medicine and Biology Symposium (SPMB)*. IEEE. 2020, pp. 1–6.
- [77] Christos Chatzichristos et al. “Epileptic seizure detection in EEG via fusion of multi-view attention-gated U-net deep neural networks”. In: *2020 IEEE Signal Processing in Medicine and Biology Symposium (SPMB)*. IEEE. 2020, pp. 1–7.

- [78] Pun nawish Thuwajit et al. “EEGWaveNet: Multiscale CNN-based spatiotemporal feature extraction for EEG seizure detection”. In: *IEEE Transactions on Industrial Informatics* 18.8 (2021), pp. 5547–5557.
- [79] Yuan Zhang et al. “Epilepsy seizure prediction on EEG using common spatial pattern and convolutional neural network”. In: *IEEE Journal of Biomedical and Health Informatics* 24.2 (2019), pp. 465–474.
- [80] Matthew J Page et al. “The PRISMA 2020 statement: an updated guideline for reporting systematic reviews”. In: *International journal of surgery* 88 (2021), p. 105906.
- [81] Mourad Ouzzani et al. “Rayyan—a web and mobile app for systematic reviews”. In: *Systematic reviews* 5 (2016), pp. 1–10.
- [82] Adapted from “PRISMA flow diagram for systematic reviews”, by BioRender.com. 2023. URL: <https://app.biorender.com/biorender-templates>.
- [83] Jean Gotman. “High frequency oscillations: the new EEG frontier?” In: *Epilepsia* 51.Suppl 1 (2010), p. 63.
- [84] Alain de Cheveigné and Israel Nelken. “Filters: when, why, and how (not) to use them”. In: *Neuron* 102.2 (2019), pp. 280–293.
- [85] Annisa Humairani et al. “Wavelet-based Entropy Analysis on EEG Signal for Detecting Seizures”. In: *2022 10th International Conference on Information and Communication Technology (ICoICT)*. IEEE. 2022, pp. 93–98.
- [86] Bahareh Salafian et al. “CNN-aided factor graphs with estimated mutual information features for seizure detection”. In: *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE. 2022, pp. 8677–8681.
- [87] Bahareh Salafian et al. “Efficient epileptic seizure detection using CNN-aided factor graphs”. In: *2021 43rd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*. IEEE. 2021, pp. 424–429.
- [88] Jeff Craley, Emily Johnson, and Archana Venkataraman. “Integrating convolutional neural networks and probabilistic graphical modeling for epileptic seizure detection in multichannel EEG”. In: *Information Processing in Medical Imaging: 26th International Conference, IPMI 2019, Hong Kong, China, June 2–7, 2019, Proceedings 26*. Springer. 2019, pp. 291–303.
- [89] Lan Wei and Catherine Mooney. “Epileptic seizure detection in clinical EEGs using an XGBoost-based method”. In: *2020 IEEE Signal Processing in Medicine and Biology Symposium (SPMB)*. IEEE. 2020, pp. 1–6.
- [90] S Poorani and P Balasubramanie. “Seizure detection based on eeg signals using asymmetrical back propagation neural network method”. In: *Circuits, Systems, and Signal Processing* 40.9 (2021), pp. 4614–4632.
- [91] Tzyy-Ping Jung et al. “Removing electroencephalographic artifacts by blind source separation”. In: *Psychophysiology* 37.2 (2000), pp. 163–178.
- [92] Shivarudhrappa Raghu et al. “Performance evaluation of DWT based sigmoid entropy in time and frequency domains for automated detection of epileptic seizures using SVM classifier”. In: *Computers in biology and medicine* 110 (2019), pp. 127–143.
- [93] PS Kumar et al. “Removal of Ocular Artifacts in the EEG through Wavelet Transform without using an EOG Reference Channel (Translation Journals style)”. In: *J. Open Problems Compt. Math* 1.3 (2008).
- [94] R Shantha Selva Kumari and J Prabin Jose. “Seizure detection in EEG using time frequency analysis and SVM”. In: *2011 international conference on emerging trends in electrical and computer technology*. IEEE. 2011, pp. 626–630.



- [95] Nazanin Nemati and Saeed Meshgini. "A medium-weight deep convolutional neural network-based approach for onset epileptic seizures classification in EEG signals". In: *Brain and Behavior* 12.11 (2022), e2763.
- [96] Md Khademul Islam Molla et al. "Artifact suppression from EEG signals using data adaptive time domain filtering". In: *Neurocomputing* 97 (2012), pp. 297–308.
- [97] Haiyang Ju, Xinhua Wang, and Yizhen Zhao. "Variational specific mode extraction: a novel method for defect signal detection of ferromagnetic pipeline". In: *Algorithms* 13.4 (2020), p. 105.
- [98] Varun Bajaj and Ram Bilas Pachori. "Classification of seizure and nonseizure EEG signals using empirical mode decomposition". In: *IEEE Transactions on Information Technology in Biomedicine* 16.6 (2011), pp. 1135–1142.
- [99] Pushpendra Singh et al. "The Fourier decomposition method for nonlinear and non-stationary time series analysis". In: *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* 473.2199 (2017), p. 20160871.
- [100] Yandong Ru et al. "Epilepsy Detection Based on Variational Mode Decomposition and Improved Sample Entropy". In: *Computational Intelligence and Neuroscience* 2022 (2022).
- [101] S Davidson et al. "Seizure Classification Using BERT NLP and a Comparison of Source Isolation Techniques with Two Different Time-Frequency Analysis". In: *2022 IEEE Signal Processing in Medicine and Biology Symposium (SPMB)*. IEEE. 2022, pp. 1–7.
- [102] Durgesh Nandini et al. "Efficient Patient Independent Seizure Detection System using WAF based Hybrid Feature Extraction Method and XGBoost classifier". In: *2022 IEEE Delhi Section Conference (DELCON)*. IEEE. 2022, pp. 1–5.
- [103] Geetanjali Sharma and Amit M Joshi. "A Fractal based Machine Learning Method for Automatic Detection of Epileptic Seizures using EEG". In: *2022 International Conference on Engineering and Emerging Technologies (ICEET)*. IEEE. 2022, pp. 1–4.
- [104] Omar Kaziha and Talal Bonny. "A convolutional neural network for seizure detection". In: *2020 Advances in Science and Engineering Technology International Conferences (ASET)*. IEEE. 2020, pp. 1–5.
- [105] Xin Chen et al. "Epilepsy classification for mining deeper relationships between EEG channels based on GCN". In: *2020 international conference on computer vision, image and deep learning (CVIDL)*. IEEE. 2020, pp. 701–706.
- [106] Zhengdao Li et al. "Graph-generative neural network for EEG-based epileptic seizure detection via discovery of dynamic brain functional connectivity". In: *Scientific Reports* 12.1 (2022), p. 18998.
- [107] Ralph G Andrzejak et al. "Seizure forecasting: Where do we stand?" In: *Epilepsia* (2023).
- [108] Jiang Wu, Tengfei Zhou, and Taiyong Li. "Detecting epileptic seizures in EEG signals with complementary ensemble empirical mode decomposition and extreme gradient boosting". In: *Entropy* 22.2 (2020), p. 140.
- [109] Wilson E Marcilio and Danilo M Eler. "From explanations to feature selection: assessing SHAP values as feature selection mechanism". In: *2020 33rd SIBGRAPI conference on Graphics, Patterns and Images (SIBGRAPI)*. Ieee. 2020, pp. 340–347.
- [110] Siyuan Qiu, Wenjin Wang, and Hailong Jiao. "LightSeizureNet: A Lightweight Deep Learning Model for Real-Time Epileptic Seizure Detection". In: *IEEE Journal of Biomedical and Health Informatics* (2022).
- [111] Mahmood A Jumaah, Ammar Ibrahim Shihab, and Akeel Abdulkareem Farhan. "Epileptic Seizures Detection Using DCT-II and KNN Classifier in Long-Term EEG Signals". In: *Iraqi Journal of Science* (2020), pp. 2687–2694.

- [112] Rosalind B Marimont and Marvin B Shapiro. “Nearest neighbour searches and the curse of dimensionality”. In: *IMA Journal of Applied Mathematics* 24.1 (1979), pp. 59–70.
- [113] Visar Berisha et al. “Digital medicine and the curse of dimensionality”. In: *NPJ digital medicine* 4.1 (2021), p. 153.
- [114] Thasneem Fathima, P Rahna, and Thanweer Gaffoor. “Wavelet based detection of epileptic seizures using scalp EEG”. In: *2020 International Conference on Futuristic Technologies in Control Systems & Renewable Energy (ICFCR)*. IEEE. 2020, pp. 1–5.
- [115] Atefeh Shariat et al. “Automatic detection of epileptic seizures using Riemannian geometry from scalp EEG recordings”. In: *Medical & Biological Engineering & Computing* 59 (2021), pp. 1431–1445.
- [116] Yuhuan Xiong et al. “Seizure detection algorithm based on fusion of spatio-temporal network constructed with dispersion index”. In: *Biomedical Signal Processing and Control* 79 (2023), p. 104155.
- [117] Vridhi Rohira et al. “Automatic Epilepsy Detection from EEG signals”. In: *Proceedings of the 6th Joint International Conference on Data Science & Management of Data (10th ACM IKDD CODS and 28th COMAD)*. 2023, pp. 272–273.
- [118] Abdul Quaiyum Ansari, Priyanka Sharma, and Manjari Tripathi. “A patient-independent classification system for onset detection of seizures”. In: *Biomedical Engineering/Biomedizinische Technik* 66.3 (2021), pp. 267–274.
- [119] Itaf Ben Slimen et al. “EEG epileptic seizure detection and classification based on dual-tree complex wavelet transform and machine learning algorithms”. In: *Journal of biomedical research* 34.3 (2020), p. 151.
- [120] Hong He, Xinyue Liu, and Yong Hao. “A progressive deep wavelet cascade classification model for epilepsy detection”. In: *Artificial Intelligence in Medicine* 118 (2021), p. 102117.
- [121] Mingyang Li, Wanzhong Chen, and Min Xia. “GNMF-based quadratic feature extraction in SSTFT domain for epileptic EEG detection”. In: *Biomedical Signal Processing and Control* 80 (2023), p. 104274.
- [122] Abdul Quaiyum Ansari, Priyanka Sharma, and Manjari Tripathi. “Automatic seizure detection using neutrosophic classifier”. In: *Physical and Engineering Sciences in Medicine* 43 (2020), pp. 1019–1028.
- [123] Shu-Ling Zhang et al. “A novel EEG-complexity-based feature and its application on the epileptic seizure detection”. In: *International Journal of Machine Learning and Cybernetics* 10 (2019), pp. 3339–3348.
- [124] Shivarudhrappa Raghu et al. “Automated detection of epileptic seizures using successive decomposition index and support vector machine classifier in long-term EEG”. In: *Neural Computing and Applications* 32 (2020), pp. 8965–8984.
- [125] Sumant Kumar Mohapatra and Srikanta Patnaik. “ESA-ASO: An enhanced search ability based atom search optimization algorithm for epileptic seizure detection”. In: *Measurement: Sensors* 24 (2022), p. 100519.
- [126] SR Ashokkumar et al. “Application of Multi-Domain Feature for Automated Seizure Detection from EEG Signal”. In: *2022 3rd International Conference on Smart Electronics and Communication (ICOSEC)*. IEEE. 2022, pp. 280–285.
- [127] S Ramakrishnan and AS Muthanantha Murugavel. “Epileptic seizure detection using fuzzy-rules-based sub-band specific features and layered multi-class SVM”. In: *Pattern Analysis and Applications* 22.3 (2019), pp. 1161–1176.
- [128] Yuhuan Xiong et al. “Seizure Detection Based on Improved Genetic Algorithm Optimized Multilayer Network”. In: *IEEE Access* 10 (2022), pp. 81343–81354.

- [129] Duanpo Wu et al. “Automatic epileptic seizures joint detection algorithm based on improved multi-domain feature of cEEG and spike feature of aEEG”. In: *IEEE Access* 7 (2019), pp. 41551–41564.
- [130] Yannis Misirlis et al. “Pediatric Epilepsy Assessment Based on EEG Analysis”. In: *2022 45th International Conference on Telecommunications and Signal Processing (TSP)*. IEEE. 2022, pp. 377–380.
- [131] Fang Dong et al. “Novel seizure detection algorithm based on multi-dimension feature selection”. In: *Biomedical Signal Processing and Control* 84 (2023), p. 104747.
- [132] Kashif Ahmad Khan et al. “A hybrid Local Binary Pattern and wavelets based approach for EEG classification for diagnosing epilepsy”. In: *Expert Systems with Applications* 140 (2020), p. 112895.
- [133] Georgepeter Jaffino, Muniasamy Sundaram, and J Prabin Jose. “Weighted 1D-local binary pattern features and Taylor-Henry gas solubility optimization based Deep Maxout network for discovering epileptic seizure using EEG”. In: *Digital Signal Processing* 122 (2022), p. 103349.
- [134] Ruofan Wang et al. “Epileptic Seizure Detection Using Geometric Features Extracted from SODP Shape of EEG Signals and AsyLnCPSO-GA”. In: *Entropy* 24.11 (2022), p. 1540.
- [135] Scot Davidson et al. “Epileptic Seizure Classification Using Combined Labels and a Genetic Algorithm”. In: *2022 IEEE 21st Mediterranean Electrotechnical Conference (MELECON)*. IEEE. 2022, pp. 430–435.
- [136] Deba Prasad Dash, Maheshkumar H Kolekar, and Kamlesh Jha. “Multi-channel EEG based automatic epileptic seizure detection using iterative filtering decomposition and Hidden Markov Model”. In: *Computers in biology and medicine* 116 (2020), p. 103571.
- [137] Gulshan Kumar, Subhash Chander, and Ahmad Almadhor. “An intelligent epilepsy seizure detection system using adaptive mode decomposition of EEG signals”. In: *Physical and Engineering Sciences in Medicine* 45.1 (2022), pp. 261–272.
- [138] Fatih Altindiş and Bülent Yilmaz. “Detection of Epileptic Seizures with Tangent Space Mapping Features of EEG Signals”. In: *2021 Medical Technologies Congress (TIPTEKNO)*. IEEE. 2021, pp. 1–4.
- [139] Amer A Sallam et al. “Epilepsy detection from EEG signals using artificial neural network”. In: *Intelligent Computing & Optimization 1*. Springer. 2019, pp. 320–327.
- [140] Chengbin Huang, Weiting Chen, and Guitao Cao. “Automatic epileptic seizure detection via attention-based cnn-birnn”. In: *2019 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. IEEE. 2019, pp. 660–663.
- [141] Sudesh Kumar, Rekh Ram Janghel, and Satya Prakash Sahu. “Convolutional Neural Network for Seizure Detection Using Scalp Electroencephalogram (EEG)”. In: *Proceedings of 6th International Conference on Recent Trends in Computing: ICRTC 2020*. Springer. 2021, pp. 431–442.
- [142] Roneel V Sharan and Shlomo Berkovsky. “Epileptic seizure detection using multi-channel EEG wavelet power spectra and 1-D convolutional neural networks”. In: *2020 42nd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*. IEEE. 2020, pp. 545–548.
- [143] Mingyang Li and Wanzhong Chen. “FFT-based deep feature learning method for EEG classification”. In: *Biomedical Signal Processing and Control* 66 (2021), p. 102492.
- [144] Muhammet Varlı and Hakan Yilmaz. “Multiple classification of EEG signals and epileptic seizure diagnosis with combined deep learning”. In: *Journal of Computational Science* 67 (2023), p. 101943.



- [145] Yanna Zhao et al. "Graph attention network with focal loss for seizure detection on electroencephalography signals". In: *International journal of neural systems* 31.07 (2021), p. 2150027.
- [146] Poojan Dalal, Chirag N Paunwala, and Neeta Chapatwala. "Statistical feature rich Deep learning based Epileptic Seizure detection". In: *2022 IEEE Region 10 Symposium (TEN-SYMP)*. IEEE. 2022, pp. 1–6.
- [147] Peter Z Yan et al. "Automated spectrographic seizure detection using convolutional neural networks". In: *Seizure* 71 (2019), pp. 124–131.
- [148] Dalibor Cimr et al. "Automatic seizure detection by convolutional neural networks with computational complexity analysis". In: *Computer Methods and Programs in Biomedicine* 229 (2023), p. 107277.
- [149] S Ramakrishnan et al. "Seizure detection with local binary pattern and CNN classifier". In: *Journal of Physics: Conference Series*. Vol. 1767. 1. IOP Publishing. 2021, p. 012029.
- [150] Zuochen Wei et al. "Automatic epileptic EEG detection using convolutional neural network with improvements in time-domain". In: *Biomedical Signal Processing and Control* 53 (2019), p. 101551.
- [151] Yongxin Sun and Xiaojuan Chen. "Automatic Detection of Epilepsy Based on Entropy Feature Fusion and Convolutional Neural Network". In: *Oxidative Medicine and Cellular Longevity* 2022 (2022).
- [152] AJ Tallón-Ballesteros. "An Effective Deep Neural Network Architecture for Cross-Subject Epileptic Seizure Detection in EEG Data". In: *Proceedings of CECNet 2021: The 11th International Conference on Electronics, Communications and Networks (CECNet), November 18-21, 2021*. Vol. 345. IOS Press. 2022, p. 54.
- [153] William Sukaria et al. "Epileptic Seizure Detection Using Convolution Neural Networks". In: *2022 IEEE International Symposium on Medical Measurements and Applications (MeMeA)*. IEEE. 2022, pp. 1–5.
- [154] Chengbin Huang et al. "A Feature Fusion Framework and Its Application to Automatic Seizure Detection". In: *IEEE Signal Processing Letters* 28 (2021), pp. 753–757.
- [155] Xiang Zhang et al. "Adversarial representation learning for robust patient-independent epileptic seizure detection". In: *IEEE journal of biomedical and health informatics* 24.10 (2020), pp. 2852–2859.
- [156] Busra Mutlu İpek, Hüseyin Oktay Altun, and Kasım Öztoprak. "Towards fully automated detection of epileptic disorders: a novel CNSVM approach with Clough–Tocher interpolation". In: *Biomedical Engineering/Biomedizinische Technik* 67.3 (2022), pp. 151–159.
- [157] Mingkan Shen et al. "Real-time epilepsy seizure detection based on EEG using tunable-Q wavelet transform and convolutional neural network". In: *Biomedical Signal Processing and Control* 82 (2023), p. 104566.
- [158] Aref Einizade et al. "Explainable automated seizure detection using attentive deep multi-view networks". In: *Biomedical Signal Processing and Control* 79 (2023), p. 104076.
- [159] Dhousha Sagga et al. "Epileptic Seizures Detection on EEG Signal Using Deep Learning Techniques". In: *2022 6th International Conference on Advanced Technologies for Signal and Image Processing (ATSIP)*. IEEE. 2022, pp. 1–6.
- [160] Susanta Kumar Rout et al. "An efficient epileptic seizure classification system using empirical wavelet transform and multi-fuse reduced deep convolutional neural network with digital implementation". In: *Biomedical Signal Processing and Control* 72 (2022), p. 103281.
- [161] Weidong Dang et al. "Studying multi-frequency multilayer brain network via deep learning for EEG-based epilepsy detection". In: *IEEE sensors journal* 21.24 (2021), pp. 27651–27658.

- [162] Samaneh Nasiri and Gari D Clifford. “Generalizable seizure detection model using generating transferable adversarial features”. In: *IEEE Signal Processing Letters* 28 (2021), pp. 568–572.
- [163] Jian Zhang et al. “Automatic epileptic EEG classification based on differential entropy and attention model”. In: *Engineering Applications of Artificial Intelligence* 96 (2020), p. 103975.
- [164] Fatima Hassan, Syed Fawad Hussain, Saeed Mian Qaisar, et al. “Epileptic Seizure Detection Using a Hybrid 1D CNN-Machine Learning Approach from EEG Data”. In: *Journal of Healthcare Engineering* 2022 (2022).
- [165] Nan Ke et al. “Convolutional transformer networks for epileptic seizure detection”. In: *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 2022, pp. 4109–4113.
- [166] Tian-li Tao et al. “Seizure detection by brain-connectivity analysis using dynamic graph isomorphism network”. In: *2022 44th Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*. IEEE. 2022, pp. 2302–2305.
- [167] Jiatong He et al. “Spatial-temporal seizure detection with graph attention network and bi-directional LSTM architecture”. In: *Biomedical Signal Processing and Control* 78 (2022), p. 103908.
- [168] Lavinia Ferariu and Adela Țucaș. “Using Hand-Crafted and Learned EEG Features for the Detection of Epileptic Seizures”. In: *2021 International Conference on e-Health and Bioengineering (EHB)*. IEEE. 2021, pp. 1–4.
- [169] Waseem Ahmad Mir, Mohd Anjum, and Sana Shahab. “Deep-EEG: an optimized and robust framework and method for EEG-based diagnosis of epileptic seizure”. In: *Diagnostics* 13.4 (2023), p. 773.
- [170] J Prabin Jose, Muniasamy Sundaram, and Georgepeter Jaffino. “Adaptive rag-bull rider: A modified self-adaptive optimization algorithm for epileptic seizure detection with deep stacked autoencoder using electroencephalogram”. In: *Biomedical Signal Processing and Control* 64 (2021), p. 102322.
- [171] Peizhen Peng et al. “Domain adaptation for epileptic EEG classification using adversarial learning and Riemannian manifold”. In: *Biomedical Signal Processing and Control* 75 (2022), p. 103555.
- [172] T Jhansi Rani and D Kavitha. “Effective Epileptic Seizure Detection Using Enhanced Salp Swarm Algorithm-based Long Short-Term Memory Network”. In: *IETE Journal of Research* (2022), pp. 1–18.
- [173] Prabin Jose Johnrose, Sundaram Muniasamy, and Jaffino Georgepeter. “Rag-bull rider optimisation with deep recurrent neural network for epileptic seizure detection using electroencephalogram”. In: *IET Signal Processing* 15.2 (2021), pp. 122–140.
- [174] Anviti Pandey et al. “Epileptic Seizure Classification Using Battle Royale Search and Rescue Optimization-Based Deep LSTM”. In: *IEEE Journal of Biomedical and Health Informatics* 26.11 (2022), pp. 5494–5505.
- [175] Md Nurul Ahad Tawhid, Siuly Siuly, and Tianning Li. “A convolutional long short-term memory-based neural network for epilepsy detection from EEG”. In: *IEEE Transactions on Instrumentation and Measurement* 71 (2022), pp. 1–11.
- [176] Haole Xi et al. “Two-Stage Multi-task Learning for Automatic Epilepsy Detection”. In: *Advances in Natural Computation, Fuzzy Systems and Knowledge Discovery: Proceedings of the ICNC-FSKD 2021 17*. Springer. 2022, pp. 866–873.

- [177] Lavinia Ferariu and Angela Mihalachi. “Detection of Epileptic Seizures With Autoencoders Working With Multiple EEG Frames”. In: *2022 E-Health and Bioengineering Conference (EHB)*. IEEE. 2022, pp. 1–4.
- [178] Katarzyna Lewandowska et al. “Impact of alarm fatigue on the work of nurses in an intensive care environment—a systematic review”. In: *International journal of environmental research and public health* 17.22 (2020), p. 8409.
- [179] Vinit Shah et al. “Objective evaluation metrics for automatic classification of EEG events”. In: *Biomedical Signal Processing: Innovation and Applications* (2021), pp. 223–255.
- [180] Jesse Davis and Mark Goadrich. “The relationship between Precision-Recall and ROC curves”. In: *Proceedings of the 23rd international conference on Machine learning*. 2006, pp. 233–240.
- [181] Maria M Cvach. “Managing hospital alarms”. In: *Nursing2020 Critical Care* 9.3 (2014), pp. 13–27.
- [182] Tita Alissa Bach, Lars-Martin Berglund, and Eva Turk. “Managing alarm systems for quality and safety in the hospital setting”. In: *BMJ open quality* 7.3 (2018), e000202.
- [183] Adriana Carla Bridi, Thiago Quinellato Louro, and Roberto Carlos Lyra da Silva. “Clinical Alarms in intensive care: implications of alarm fatigue for the safety of patients”. In: *Revista latino-americana de enfermagem* 22 (2014), pp. 1034–1040.
- [184] Lihan Tang, Menglian Zhao, and Xiaobo Wu. “Accurate classification of epilepsy seizure types using wavelet packet decomposition and local detrended fluctuation analysis”. In: *Electronics Letters* 56.17 (2020), pp. 861–863.
- [185] Jin Jing et al. “Interrater reliability of expert electroencephalographers identifying seizures and rhythmic and periodic patterns in EEGs”. In: *Neurology* 100.17 (2023), e1737–e1749.
- [186] Jin Jing et al. “Development of expert-level classification of seizures and rhythmic and periodic patterns during eeg interpretation”. In: *Neurology* 100.17 (2023), e1750–e1762.
- [187] Lawrence J Hirsch et al. “American Clinical Neurophysiology Society’s standardized critical care EEG terminology: 2021 version”. In: *Journal of clinical neurophysiology: official publication of the American Electroencephalographic Society* 38.1 (2021), p. 1.
- [188] Baharan Kamousi et al. “Monitoring the burden of seizures and highly epileptiform patterns in critical care with a novel machine learning method”. In: *Neurocritical care* 34 (2021), pp. 908–917.
- [189] Anahit Babayan et al. “A mind-brain-body dataset of MRI, EEG, cognition, emotion, and peripheral physiology in young and old adults”. In: *Scientific data* 6.1 (2019), pp. 1–21.
- [190] Paula Sanz Leon et al. “The Virtual Brain: a simulator of primate brain network dynamics”. In: *Frontiers in neuroinformatics* 7 (2013), p. 10.
- [191] Viktor K Jirsa et al. “The virtual epileptic patient: individualized whole-brain models of epilepsy spread”. In: *Neuroimage* 145 (2017), pp. 377–388.
- [192] Viktor Jirsa et al. “Personalised virtual brain models in epilepsy”. In: *The Lancet Neurology* 22.5 (2023), pp. 443–454.
- [193] Micheal Strein et al. “Prevention, treatment, and monitoring of seizures in the intensive care unit”. In: *Journal of clinical medicine* 8.8 (2019), p. 1177.
- [194] M Brandon Westover et al. “The probability of seizures during EEG monitoring in critically ill adults”. In: *Clinical Neurophysiology* 126.3 (2015), pp. 463–471.
- [195] J Claassen et al. “Detection of electrographic seizures with continuous EEG monitoring in critically ill patients”. In: *Neurology* 62.10 (2004), pp. 1743–1748.
- [196] Ajay Rajshekar et al. *Automated Seizure Detection and Management using CERIBELL: A single center cohort study (P10-1.006)*. 2023.

- [197] Irfan Al-Hussaini and Cassie S Mitchell. “Towards Interpretable Seizure Detection Using Wearables”. In: *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE. 2023, pp. 1–2.
- [198] Charles Begley et al. “The global cost of epilepsy: A systematic review and extrapolation”. In: *Epilepsia* 63.4 (2022), pp. 892–903.
- [199] Ben D Fulcher, Max A Little, and Nick S Jones. “Highly comparative time-series analysis: the empirical structure of time series and their methods”. In: *Journal of the Royal Society Interface* 10.83 (2013), p. 20130048.
- [200] Farhad Maleki et al. “Generalizability of machine learning models: Quantitative evaluation of three methodological pitfalls”. In: *Radiology: Artificial Intelligence* 5.1 (2022), e220028.