

Finding Short Signals in Long Irregular Time Series with Continuous-Time Attention Policy Networks

Thomas Hartvigsen¹ Jidapa Thadajarassiri² Xiangnan Kong² Elke Rundensteiner²

Abstract

Irregularly-sampled time series (ITS) are native to high-impact domains like healthcare, where measurements are collected over time at uneven intervals. However, for many classification problems, only small portions of long time series are often relevant to the class label. In this case, existing ITS models often fail to classify long series since they rely on careful imputation, which easily over- or under-samples the relevant regions. Using this insight, we then propose CAT, a model that classifies multivariate ITS by explicitly seeking highly-relevant portions of an input series' timeline. CAT achieves this by integrating three components: (1) A *Moment Network* learns to seek relevant moments in an ITS's continuous timeline using reinforcement learning. (2) A *Receptor Network* models the temporal dynamics of both observations *and* their *timing* localized around predicted moments. (3) A recurrent Transition Model models the sequence of transitions between these moments, cultivating a representation with which the series is classified. Using synthetic and real data, we find that CAT outperforms ten state-of-the-art methods by finding short signals in long irregular time series.

1. Introduction

Background. Irregularly-sampled time series (ITS) have uneven spaces between their observations and are common in impactful domains like healthcare (Hong et al., 2020), environmental science (Cao et al., 2018), and human activity recognition (Singh et al., 2019). Uneven gaps can arise from many sources. For example, in physiological streams, clinicians drive the collection of medical record data by requesting different lab tests and measurements in real time as they investigate the root causes of their patient's conditions (Lipton et al., 2016). *Which* measurements are taken

¹MIT ²WPI. Correspondence to: Thomas Hartvigsen <tomh@mit.edu>.

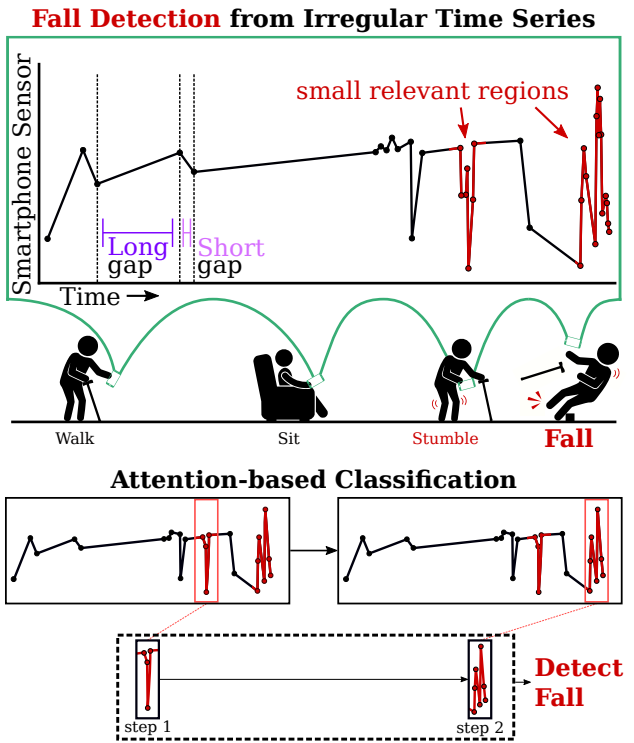


Figure 1: Attention-based classification for ITS.

when differs between patients. When classifying such time series, there are often relationships between *when* observations are made and the class label for the resulting time series. For instance, sicker patients may have more measurements. ITS can also be quite long, while the regions most-relevant to the classification may be quite short, taking up only a small portion of the timeline and creating a small *signal-to-noise* ratio. A successful model must find the best regions in the timeline at which to capture signals in both the values themselves and the patterns in *when* observations were made, or *informative irregularity*, while ignoring irrelevant regions.

Motivating Example. Consider detecting if a person *Fell* using their smartphone's sensors, as illustrated in Figure 1. To extend battery life, a listening probe is used to only collect data when certain conditions are met, for example when the accelerometer changes rapidly. Since the phone

is not always moving, the stored time series are naturally irregularly-sampled. To detect a fall, some regions of the accelerometer’s records are far more relevant than others. Leading up to a fall, for instance, a person may have stumbled earlier in the day. However, there can also be many false positives where the phone moves quickly even though the person is not falling (setting the phone down, for instance). Additionally, *when* observations are made can also be useful: if the phone moves after a long gap, the person may be getting out of bed. Since only some regions are relevant, all a classifier needs are the few most relevant moments in the timeline. Finding these moments is especially important for the long series that naturally exist in many domains.

State-of-the-art. There have been many recent advances in classifying ITS data, though most focus on sparse series with few observations. Many works treat ITS classification as a *missing value imputation* problem (Che et al., 2018; Lipton et al., 2016; Zheng et al., 2017), converting ITS to regular series then performing standard classification. However, to capture short signals, many values need to be imputed to avoid aggregating intricate signals. Plus, this increases the length of the series and imputes values in irrelevant regions of the timeline. As many ITS methods rely on Recurrent Neural Networks, making time series longer will likely decay performance. On the other hand, imputing too few values easily bypasses short signals, aggregating away crucial information. Some works capture informative irregularity by computing statistical features such as *missingness indicators* (Lipton et al., 2016) or the time since last observations (Che et al., 2018) as additional input variables, inflating the feature space.

Some recent works have leaned into learning continuous-time representations *directly* from raw ITS data (Kidger et al., 2020b; Shukla & Marlin, 2019; Li & Marlin, 2016, 2020; Rubanova et al., 2019; De Brouwer et al., 2019; Oh et al., 2018; Shukla & Marlin, 2021). However, they still rely on hand-picking new *reference* timesteps at which to estimate values or compute representations, falling prey to the same challenges of imputation. To-date, these methods do not adapt their reference timesteps to the inputs. Overall, current machine learning methods for classifying ITS data are expected to underperform on long series where the relevant signals are proportionally short.

Problem Definition. We specifically address the problem of *Attention-based ITS Classification (ABC)*, which is to classify long ITS by finding small discriminative signals in the continuous timeline, as illustrated in Figure 1. Given a set of labeled ITS, where each series consists of one sequence of (*timestep*, *value*) pairs per variable, our aim is to produce a classifier that can correctly assign class labels y to previously-unseen instances. For long series, the rele-

vant time window, or the proportion of the timeline needed for classification, may be very small in practice. A successful model should explicitly find these *discriminative moments* with which it can make an accurate classification.

Challenges. Solving the ABC problem is challenging for the three following reasons:

- *Finding Short Signals.* Short, relevant windows of a continuous timeline can be hard to identify, akin to finding a needle in a haystack. For long series, this means that much of the timeline contains effectively irrelevant information, which a model must learn to ignore. Meanwhile, the model must also while avoiding learning spurious correlations found outside relevant regions.
- *Unknown Signal Locations.* Relevant signals may occur anywhere in the continuous timeline. However, rarely are the *true* signal locations labeled, so we assume no prior knowledge of which moments *should* be used for classification. Still, a good model must successfully find these discriminative moments, even without supervision.
- *Informative Irregularity.* Discriminative information often arises in the patterns of *when* observations are made (Rubin, 1976). For instance, rapid measurements may indicate a sicker patient. Learning from such irregularity is often crucial to accurate classification, yet few methods exist for capturing such signals.

Proposed Method. To address these challenges, we propose the Continuous-time Attention policy network (CAT) as an effective approach to the ABC problem. CAT searches for relevant regions of input series via a reinforcement learning-based *Moment Network*, which learns to find *moments of interest* in the continuous timeline, one by one. At each predicted moment, a *Receptor Network* reads and represents the local temporal dynamics in the measurements along with patterns that exist in the timing of observations through a continuous-time density function. Along the way, a recurrent *Transition Model* constructs a discriminative representation of the *transitions between moments of interest*, which is ultimately used to classify the series. CAT thus presents a novel paradigm for classifying ITS where intricate signals in long series are explicitly sought out and captured. Additionally, CAT generalizes recent ITS classifiers with its flexible *Receptor Network*, which can easily be augmented to leverage components of other recent ITS models.

Contributions. Our contributions are as follows:

- We identify a new, real problem setting for classifying irregularly-sampled time series on which existing

state-of-the-art methods underperform.

- Using insights from this problem, we develop CAT, a novel framework for classifying long irregular time series by finding relevant moments in the *continuous* timeline, generalizing recent work.
- We show that CAT successfully discovers intricate signals in ITS, outperforming the the main competitors on both synthetic and real-world data.

2. Related Work

The ABC problem for ITS relates to both *ITS Classification* and *Input Attention*.

Classifying Irregularly-Sampled Time Series. Classifying irregularly-sampled time series has recently become a popular and impactful problem as it generalizes many prior classification settings. To-date, most approaches (Lipton et al., 2016; Zheng et al., 2017; Che et al., 2018) treat ITS classification as a *missing value imputation* problem: Create a set of evenly-spaced bins, then aggregate multiple values within each bin and estimate one value per empty bin. After imputation, regular time series classification may be performed. Some recent ITS classifiers extend beyond simple imputation options (*e.g.*, mean) approaches by either including auxiliary information such as a *missingness-indicator* (Lipton et al., 2016) or *time-since-last-observation* (Che et al., 2018) as extra features to preserve properties found in the irregularity. Others build more complex value estimators by either learning generative models (Li & Marlin, 2020), using differentiable gaussian kernel adapters (Shukla & Marlin, 2019), or including decay mechanisms in Recurrent Neural Networks (RNN) to encode information-loss when variables go unobserved over long periods of time (Mozer et al., 2017; Che et al., 2018). Many recent works have also begun parameterizing differential equations to serve as time series models (Kidger et al., 2020a; Lechner & Hasani, 2020; Rubanova et al., 2019; Jia & Benson, 2019; Hasani et al., 2021; Schirmer et al., 2022; Salvi et al., 2022), though most still estimate values at hand-picked time steps, then use the estimated values for classification.

Some recent models have also integrated attention mechanisms into ITS classification (Shukla & Marlin, 2021; Chen & Chien, 2021; Tan et al., 2021). However, they still hand-pick reference timesteps for each input time series. Given long ITS with short signals, this decision is hugely impactful, as we show in our experiments. Moreover, by relying on RNNs for classification, these recent methods easily fail to capture signals when the number of estimated values gets too large. This requires the RNN to filter out many irrelevant timesteps in a long series, which is notoriously challenging due to both their slow inference and vanishing

gradients (Hochreiter, 1998).

Input Attention. The goal of *Input Attention* is to discover relevant regions in the *input* space of a given instance and it has recently broken major ground in classifying images (Mnih et al., 2014), graphs (Lee et al., 2019), text (Sood et al., 2020), and regularly-spaced time series (Ismail et al., 2019). We refer to this as *input* attention, as such methods search for relevant regions in the *input space* of each instance. This approach is particularly impactful when inputs are high-dimensional as it explicitly disregards irrelevant regions of the input space. These methods also aid interpretability by clearly displaying which regions of an input were used to make a classification. Input attention has yet to be considered for ITS despite strong implications of successful models.

Input attention differs from *attention mechanisms for recurrent neural networks* (Bahdanau et al., 2014), where attention distributions are predicted over the timesteps in the *latent* space of an RNN. While there is some conceptual overlap, input attention is more data-driven in that it finds regions in the *input* space as opposed to the *latent* space of this specific neural network architecture.

3. Methodology

3.1. Problem Formulation

Given a set of N labeled irregularly-sampled time series $\mathcal{D} = \{(X_i, y_i)\}_{i=0}^N$, consider the D variables of instance $X_i = [X_i^1, \dots, X_i^D]$. To aid readability, all descriptions are provided in terms of one instance and one variable wherever possible. For each variable d , $X^d = [(t_1^d, v_1^d), \dots, (t_{T^d}^d, v_{T^d}^d)]$, where t_i^d is the i -th timestamp of the d -th variable and v_i^d is its corresponding value. Timestamps t may differ between variables and the number of observations T^d may be unique to variable d . We also assume that the inputs X have short signals: Most of the relevant information comes from a small proportion of a series’ timeline. There may still be multiple relevant regions, however. The goal is to learn a function $f : \mathbb{X} \rightarrow \mathcal{Y}$ that accurately maps input $X \in \mathbb{X}$ to its class $y \in \mathcal{Y}$ for previously-unseen time series, where \mathbb{X} is the input space of ITS and $\mathcal{Y} = \{0, \dots, C\}$ is the set of C classes.

3.2. Proposed Method

We propose a **Continuous-time Attention Policy Network (CAT)**, a novel model containing four key steps that work in concert to find short discriminative signals in long ITS:

1. A *Receptor Network* learns to model ITS observations (both the raw values *and* informative irregularity) local to a given *moment of interest*.

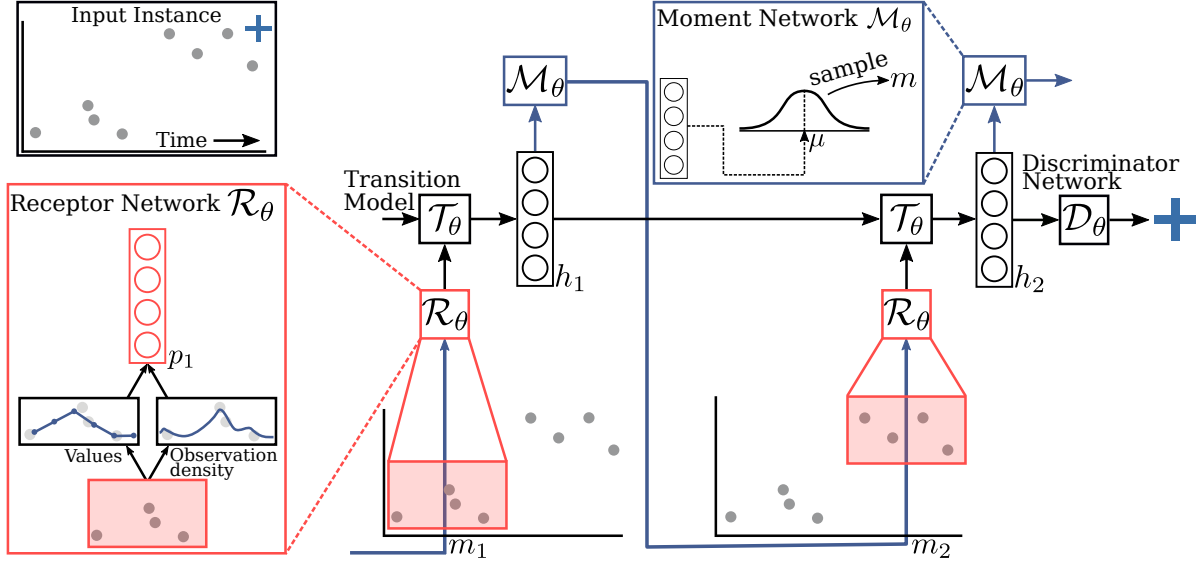


Figure 2: Overview of CAT. The *Receptor Network* models input values and irregularity around a moment m_i in the continuous timeline: $\hat{x}_i = \mathcal{R}(X, m)$. The *Transition Model* then updates its hidden state $h_i = \mathcal{T}(\hat{x}_i)$, modeling the transitions between *moments*. Then, the *Moment Network* parameterizes a Normal distribution from which it samples the next moment $m_{i+1} = \mathcal{M}(h_i)$. After iterating K times, the *Discriminator Network* predicts the final class: $y = \mathcal{D}(h_K)$, classifying the entire series.

2. A recurrent *Transition Model* represents the Receptor Network’s findings across multiple moments.
3. A reinforcement learning *Moment Network* predicts *moments of interest* based on the Transition Model. A *moment of interest* is timestamp around which relevant information may exist.
4. After k repetitions of Steps 1-3, a *Discriminator Network* classifies X using all steps.

Beyond state-of-the-art performance, a clear benefit of CAT is its novel framework for classifying time series: CAT decomposes a time series into a sequence of local representations that is discriminative in (1) which subsequences are modeled, and (2) their relative order. This approach is more flexible than rigidly reading a time series from either left-to-right, right-to-left, or all-at-once; our model adapts the processing order to the input. Since ordering is discrete, Reinforcement Learning is a natural fit. That is, we let the model pick the order, then reward or penalize based on the final classification. This technical novelty helps CAT stand out from alternative ITS models.

3.2.1. RECEPTOR NETWORK

First, the Receptor Network \mathcal{R}_θ creates a vector representation of values and irregularity *local* to a given moment of interest m_i . So given $m_i \in [0, \max T]$, \mathcal{R}_θ predicts a vector \hat{x}_i , representing the local *values* and *informative irreg-*

ularity within a width- δ window of X centered on moment m_i , where $\max T$ is the largest timestamp in X . Thus \mathcal{R}_θ can be placed anywhere in the continuous timeline, where it will proceed to model local signals. In our experiments, the first moment, m_0 , is sampled from a uniform distribution across the timeline. To compute local representations of both values *and* irregularity, we compute two w -dimensional vectors per variable: \mathbf{p} represents X ’s values and \mathbf{q} represents informative irregularity, which are then encoded into a shared representation $\hat{\mathbf{x}}$. For readability, we describe \mathcal{R} for one variable, omitting superscripts d since all variables are processed the same way and in parallel.

To compute vector $\hat{\mathbf{x}}$, all timestamps and values within this window are first extracted into two vectors: τ is a sequence of timestamps in the window $[m_i - \frac{\delta}{2}, m_i + \frac{\delta}{2}]$, and ν contains their corresponding values. We use $[]$ to denote a range in the timeline beginning at the real time $m_i - \frac{\delta}{2}$ and ending at time $m_i + \frac{\delta}{2}$.

To compute \mathbf{p} , the representation of the values within the window surrounding m_i , we linearly interpolate values ν to estimate w values at a set of new timestamps. The j -th element of \mathbf{p} can be interpolated with respect to timestamp $t' = m_i - \frac{\delta}{2} + \frac{j\delta}{w}$ for $j = \{1, \dots, w\}$ as

$$\mathbf{p}_j = \frac{(\text{SG}(t', \tau) - t') \nu_{\text{LL}(t', \tau)} + (t' - \text{LL}(t', \tau)) \nu_{\text{SG}(t', \tau)}}{\text{SG}(t', \tau) - \text{LL}(t', \tau)},$$

where $\text{LL}(t', \tau)$ is the largest timestamp less than t' and $\nu_{\text{LL}(t', \tau)}$ is its corresponding value. Similarly, $\text{SG}(t', \tau)$ is

the smallest timestamp greater than t' and $\nu_{\text{SG}(t',\tau)}$ is its corresponding value. By iterating j across integers from 1 to w , we compute w evenly-spaced values representing the local observations. If a timestamp $t' > \max T$ or $t' < \min T$, the nearest value in the window is returned, flattening the edges of the window. If no observations occur in the window, we set $\hat{x} = \{0\}^w$.

To compute \mathbf{q} , which represents *informative irregularity* within the window, we learn a function to represent the *timing* of observations, quantifying irregularity through the squared exponential kernel, inspired by (Li & Marlin, 2016). Thus the j -th element of \mathbf{q} as computed with respect to each $t' = m_i - \frac{\delta}{2} + \frac{j\delta}{w}$ for $j = \{1, \dots, w\}$ is

$$\mathbf{q}_j = \sum_{k=1}^{|\tau|} e^{-\alpha(t' - \tau_k)}, \quad (1)$$

where τ_K is the K -th element of sequence τ . Thus the *timing* of the observations is converted to a sequence of densities, which often change by class (Lipton et al., 2016). α controls the kernel distance between t' and τ_K and can be picked or learned during training (Shukla & Marlin, 2019).

Since the output of the Receptor Network will eventually be used by the Moment Network to predict the next moment m_{i+1} , we also compute \mathbf{p} and \mathbf{q} for each variable at two granularities: One for fine-grained local information, one for coarse-grained representation of the entire series that is useful for both capturing long-term trends and for finding the next moments of interest, inspired by (Mnih et al., 2014). After computing \mathbf{p} and \mathbf{q} , a neural network predicts a L -dimensional representation $\hat{\mathbf{x}}_i$, creating a dense, vector representation of the width- δ window surrounding moment m_i :

$$\hat{\mathbf{x}}_i = \psi(\mathbf{W}[\mathbf{F}(\{\mathbf{p}^d\}_{d=1}^D), \mathbf{F}(\{\mathbf{q}^d\}_{d=1}^D)] + \mathbf{b}), \quad (2)$$

where $\mathbf{F}(\cdot)$ and $[\cdot]$ denote flattening and concatenation, respectively. \mathbf{W} and \mathbf{b} are a matrix and vector of learnable parameters of shape $L \times 4w$ and $4w$. ψ is the rectified linear unit. To also incorporate *where* the collected data come from in the timeline, we concatenate m_i with $\hat{\mathbf{x}}_i$ before passing it to the Transition Model.

3.2.2. TRANSITION MODEL

Next, the Transition Model \mathcal{T}_θ represents the *transitions* between information gathered at each moment of interest. We follow the state-of-the-art for a vast array of sequential learning tasks and implement this component as an RNN, creating one H -dimensional vector representation \mathbf{h}_i per moment-of-interest. To avoid vanishing gradients, we use a Gated Recurrent Unit (GRU) (Cho et al., 2014) to compute the hidden state \mathbf{h}_i .

This recurrent component takes only K steps and K is typically kept very low ($K = 3$ in our experiments). In contrast, most recent models instead step through a large number of imputed timestamps T (typically $T \gg K$) creating slow models that are hard to optimize.

3.2.3. MOMENT NETWORK

Next, the *Moment Network* \mathcal{M}_θ uses the hidden state \mathbf{h}_i and predicts the next moment-of-interest m_{i+1} . There are no ground-truth moments, so we frame this component as a Partially-Observable Markov Decision Process (POMDP), similar to (Mnih et al., 2014). We follow the standard approach and solve this POMDP using on-policy reinforcement learning. In this way, the hidden state \mathbf{h}_i from the *Transition Model* serves as an observation from the environment (representing the data collected at all prior moments of interest). The possible actions include all real-valued timestamps between 0 and $\max T$, and we define the reward to be the final classification success. The goal is to learn a policy $\pi(\mathbf{h}_i)$ that predicts the next moment m_{i+1} .

Since there are infinitely-many moments in the continuous timeline, we parameterize the mean μ_i of a Normal distribution with fixed variance from which we *sample* real-valued m_{i+1} . To acquire good samples, \mathbf{h}_i is first projected into a one-dimensional probabilistic space by a neural network: $\mu_i = \sigma(\mathbf{W}\mathbf{h}_i + \mathbf{b})$. We then scale μ_i by multiplying with $\max T$ and sample moment $m_{i+1} \sim \mathcal{N}(\mu_i, \sigma)$, with tunable σ . If $m_{i+1} > \max T$, we re-assign $m_{i+1} := \max T$, and if $m_{i+1} < 0$, $m_{i+1} := 0$. To train \mathcal{M}_θ , we set reward $r_i = 1$ if the final classification is accurate, otherwise $r_i = -1$. The Moment Network thus seeks *discriminative* regions, which lead to the highest rewards.

CAT predicts K moments of interest, iteratively cycling between the Receptor Network, Moment Network, and Transition Model K times. This packs information from K steps into the final hidden state \mathbf{h}_K .

3.2.4. DISCRIMINATOR NETWORK

The final component of CAT is a *Discriminator Network* \mathcal{D}_θ , which learns to project the Transition Model's final hidden state \mathbf{h}_K into a C -dimensional probabilistic space in which it predicts \hat{y} to be X 's class label. This final classification is made via a single linear layer: $\hat{y} = \text{softmax}(\mathbf{W}\mathbf{h}_K + \mathbf{b})$. The discriminator is naturally connected to the transition model, so is easily expandable according to the required complexity of a task.

3.2.5. CAT TRAINING

The Receptor Network, Transition Model, and Discriminator are optimized together to predict \hat{y} accurately by mini-

mizing cross entropy:

$$\mathcal{L}_s(\theta_s) = - \sum_{c=0}^C y_c \log \hat{y}_c, \quad (3)$$

where y_c is 1 if X is in class c and \hat{y}_c is the corresponding prediction. θ_s denotes these networks’ parameters.

The Moment Network, on the other hand, samples the moments, so its learning objective is the maximization of the expected reward: $R = \sum_{i=0}^K r_i$, so $\theta_{rl}^* = \arg \max_{\theta_{rl}} \mathbb{E}[R]$, where θ_{rl}^* is the optimal parameters for the *Moment Network*. However, this is not differentiable.

To maximize $\mathbb{E}[R]$ using backpropagation, we follow the standard protocol for on-policy reinforcement learning and optimize the Moment Network’s policy using the REINFORCE algorithm (Williams, 1992). Thus, we use a well-justified surrogate loss function that is differentiable, allowing for optimization by taking steps in the direction of $\mathbb{E}[\nabla \log \pi(h_{0:k}, \mu_{0:k}, r_{0:k})R]$. Thus the gradient can then be approximated for the predicted moments. Thus learning progresses, but there may be high variance in the policy updates since this is not the *true* gradient for maximizing $\mathbb{E}[R]$. To reduce variance, we employ the commonly-used *baseline* approach to approximate the expected reward, with which we may adjust the raw reward values, as shown in Equation 4. Here, b_j is a baseline predicted by a two-layer neural network and its predictions approximate the mean R by reducing the mean squared error between b_j and the average R . The weights θ_{rl} are thus updated by how much better than average are the outcomes.

$$\mathcal{L}_{rl}(\theta_{rl}) = -\mathbb{E} \left[\sum_{i=0}^k \log \pi(m_i | h_i) \left[\sum_{j=i}^k (R - b_j) \right] \right] \quad (4)$$

Finally, the entire network can be optimized jointly via gradient descent on the sum of Equations 3 and 4: $\mathcal{L}(\theta) = \mathcal{L}_s(\theta_s) + \mathcal{L}_{rl}(\theta_{rl})$, where θ denotes CAT’s parameters.

4. Experiments

4.1. Datasets

We evaluate CAT using one synthetic dataset and five real-world publicly-available datasets.

MII: We develop a synthetic binary classification dataset to demonstrate that CAT indeed finds short signals in long ITS data. To add signals for different classes, we center a width- Δ discriminative region around a random moment in the timeline for each time series. The values for the timestamps within the width- Δ window take one of two forms, depending on the class. One class is characterized by the values $\{1, 1, 1\}$ (“II”-shaped), and the other by the values $\{1, 0, 1\}$ (“M”-shaped). The timestamps corresponding

to these values are evenly-spaced in the width- Δ window. All timestamps *not* in the discriminative region are sampled uniformly across the timeline and values are sampled from a Normal distribution $\mathcal{N}(0, 1)$. In selecting Δ , we determine the signal-to-noise ratio of the data: A small Δ means that the “II” or “M” signals happen in a short period of time, so overlooking the signal is punished more. We generate 5000 time series, each with 500 timestamps, and have an equal number of instances for each class.

UWave (Liu et al., 2009): The popular UWave dataset contains 4478 length-945 gesture pattern time series collected from a handheld device. Each series is a member of one of eight classes. We follow the preprocessing procedure outlined by (Li & Marlin, 2016), randomly downsampling to 10% of the original values to create irregularity.

ExtraSensory (Vaizman et al., 2017): Following (Hartvigsen et al., 2022), we augment existing human activity data by simulating listening probes on smartphone data. Listening probes collect data from devices only when certain conditions are met, creating realistic ITS. For example, consider detecting hand tremors for digital health (García-Magariño et al., 2016). A listening probe on a smartphone’s accelerometer will collect data only when the phone moves rapidly, capturing hand tremors while the phone is carried. However, false positives are common: when the phone is set down or dropped, data are *also* collected, resulting in irrelevant regions. Our sampling is more realistic than prior works, which randomly downsample without encoding meaning into the irregularity of samples.

We extract four disjoint, non-overlapping datasets from the challenging ExtraSensory human activity recognition database (Vaizman et al., 2017) via a simulated listening probe on the 3-dimensional (x, y, and z axes) accelerometer records. When the norm of the difference between consecutive records surpasses a threshold $\gamma = 0.001$, the corresponding accelerometer data are collected. We collect four datasets, one for each of four human activities: WALKING (2636 time series), RUNNING (1066 time series), LYING-DOWN (7426 time series), and SLEEPING (9276 time series). For each class, we extract data for the person who performed the activity the most since people’s activity patterns are often incomparable. We then break each series into windows of 200 timestamps, then apply the listening probe. The task is to detect whether the person performed the activity within this window. We finally balance each dataset to have an equal number of positive and negative series and ensure no extracted segments overlap.

4.2. Compared Methods

We compare CAT to ten recent ITS classifiers. The first four methods are use imputation and feature expansions: linear interpolation (GRU-interp), mean imputa-

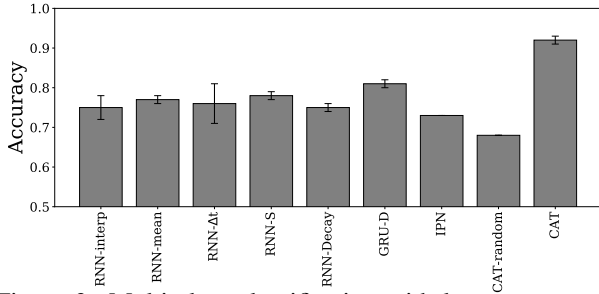


Figure 3: Multi-class classification with long UWave series. CAT outperforms recent methods when sampling many points from the timeline of each UWave time series.

tion (GRU-mean), mean imputation with extra time-since-last-observation features (GRU- Δt), and mean imputation with a missingness indicator (GRU-S) (Lipton et al., 2016). The second group contains state-of-the-art ITS classifiers: GRU-Decay (Mozer et al., 2017), GRU-D (Che et al., 2018), IPN (Shukla & Marlin, 2019), mTAN (Shukla & Marlin, 2021), and NCDE (Kidger et al., 2020a). We also ablate CAT by replacing the Moment Network with randomly-selected moments of interest, which we refer to as CAT w/o Moment.

4.3. Implementation Details

For the UWave dataset, we use a standard 80% training, 10% validation, and 10% testing split. The ExtraSensory datasets contain instances taken from different windows along a single timeline. To avoid cross-contamination, we split instances *in time*, aiming for 80% training and 20% testing splits. The training/testing process is repeated five times and we report the average and standard deviation for all experiments. All methods use 64-dimensional hidden states for their respective RNNs. For CAT, we set $k = 3$, use a 50-dimensional representation for the Receptor Network, and set $\alpha = 100$ in Equation 1. All models are optimized using Adam with a learning rate of $1e^{-3}$ and weight decay of $1e^{-5}$ and all methods are run until their losses converge, taking around 200 epochs. Each model is implemented in PyTorch in our public code.

4.4. Experimental Results

4.4.1. EXPERIMENTS ON REAL-WORLD DATA.

First, we demonstrate that CAT indeed handles long series better than the state-of-the-art methods. To achieve this, we impute the UWave data with 200 timestamps, which is much higher than prior experiments (Shukla & Marlin, 2019). For ease of comparison, we also have CAT observe the data at the same “resolution” by setting $w = \delta * 200$ where δ is the receptor-width hyperparameter. This resolution can be tuned within CAT. Our results are reported in

Table 1: Accuracy with *infrequently*-imputed values for the Human Activity datasets. **Bold** indicates best performance.

Methods	Datasets			
	Walking	Running	Lying Down	Sleeping
GRU-interp	0.65 (0.02)	0.52 (0.04)	0.83 (0.05)	0.76 (0.05)
GRU-mean	0.64 (0.02)	0.52 (0.09)	0.78 (0.04)	0.79 (0.01)
GRU- Δt	0.64 (0.01)	0.52 (0.07)	0.83 (0.04)	0.79 (0.03)
GRU-S	0.64 (0.01)	0.46 (0.03)	0.82 (0.08)	0.76 (0.03)
GRU-Decay	0.64 (0.02)	0.44 (0.03)	0.77 (0.04)	0.78 (0.03)
GRU-D	0.65 (0.01)	0.43 (0.01)	0.78 (0.03)	0.76 (0.03)
IPN	0.65 (0.01)	0.46 (0.03)	0.85 (0.04)	0.77 (0.03)
mTAN	0.68 (0.02)	0.56 (0.01)	0.73 (0.00)	0.76 (0.02)
NCDE	0.66 (0.00)	0.53 (0.00)	0.75 (0.00)	0.71 (0.00)
CAT w/o Moment	0.60 (0.02)	0.47 (0.01)	0.76 (0.01)	0.73 (0.01)
CAT (ours)	0.81 (0.03)	0.62 (0.01)	0.87 (0.04)	0.91 (0.02)

Table 2: Accuracy with *frequently*-imputed values.

Methods	Datasets			
	Walking	Running	Lying Down	Sleeping
GRU-interp	0.62 (0.04)	0.52 (0.05)	0.78 (0.07)	0.76 (0.05)
GRU-mean	0.59 (0.02)	0.51 (0.09)	0.8 (0.04)	0.79 (0.02)
GRU- Δt	0.56 (0.02)	0.48 (0.02)	0.83 (0.04)	0.80 (0.03)
GRU-S	0.61 (0.04)	0.51 (0.07)	0.89 (0.02)	0.80 (0.02)
GRU-Decay	0.59 (0.02)	0.51 (0.05)	0.88 (0.05)	0.77 (0.02)
GRU-D	0.63 (0.02)	0.45 (0.01)	0.82 (0.03)	0.73 (0.01)
IPN	0.62 (0.01)	0.46 (0.01)	0.85 (0.05)	0.78 (0.04)
mTAN	0.64 (0.01)	0.55 (0.01)	0.74 (0.02)	0.76 (0.02)
NCDE	0.66 (0.00)	0.53 (0.00)	0.75 (0.00)	0.71 (0.00)
CAT w/o Moment	0.61 (0.01)	0.47 (0.00)	0.77 (0.00)	0.74 (0.01)
CAT (ours)	0.83 (0.03)	0.65 (0.04)	0.89 (0.03)	0.90 (0.01)

Figure 3. As expected, CAT achieves state-of-the-art accuracy on these data while the compared methods underperform their accuracy with roughly 100 imputed values. This indicates that CAT is far more robust to longer series than the state-of-the-art ITS classifiers.

Second, we show that CAT successfully captures *informative irregularity* in long series, as indicated by our results on the human activity recognition datasets (WALKING, RUNNING, LYINGDOWN, and SLEEPING). We compare all models using two settings: infrequent imputation (200 values) and frequent imputation (500 values). Intuitively, *frequent* imputation leads to clearer signals, as there are more values imputed on the signal, while *infrequent* imputation leads to unclear signals. To successfully classify these data given infrequent imputation, finding the relevant regions of the data is more important. On the other hand, frequent imputations provide clear signals but come with the added noise, requiring explicit discovery of the relevant regions. Again, to compare with other methods, we set $w = \delta * 200$ and $w = \delta * 500$ for each respective frequency.

Our results for this experiment, shown in Tables 1 and 2, show that, as expected, CAT outperforms all compared methods in both the *infrequent* and the *frequent* settings for all datasets by an average of over 8%. The baselines also mainly perform their best with *infrequent* imputation, while CAT performs its best at *frequent* imputation as it adapts to different resolutions. Also as expected, the recent GRU-

D , IPN , and $mTAN$ models are generally CAT’s strongest competitors. As expected, methods that model irregularity ($GRU-D$, IPN , $GRU-S$, and CAT) largely beat the methods that disregard irregularity. $GRU-interp$ ’s poor performance indicates that the benefits of CAT do not come from the linear interpolation used by the Receptor Network.

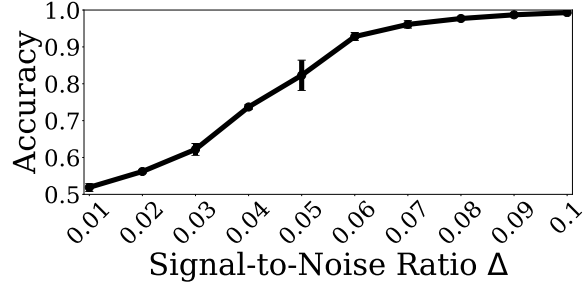
For all datasets, CAT outperforms CAT *w/o* *Moment*, the *policy-free* version of CAT that places the *Receptor Network* at random moments in the timeline. In fact, CAT *w/o* *Moment* is overall the *worst*-performing method, indicating that CAT ’s strong performance comes from a successfully-trained *Moment Network*. Therefore CAT indeed succeeds to *learn* the discriminative regions of the given time series. However, it is possible that CAT *w/o* *Moment* could still perform well with enough moments. To determine if this is the case, we vary the number of moments for both CAT and CAT *w/o* *Moment*. As our results in Figure 5 in the Appendix show, the moment network is effective.

4.4.2. EXPERIMENTS ON LONG SYNTHETIC DATA.

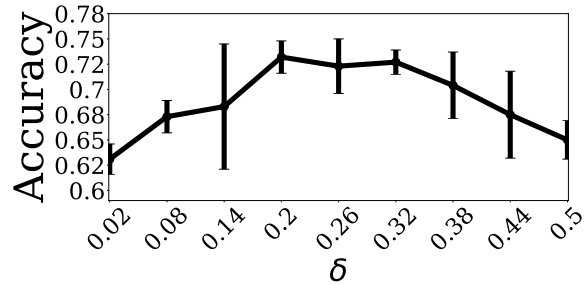
We finally evaluate CAT ’s robustness to signal length using the synthetic MII dataset. We use long, 500 timestep time series for all experiments. Therefore, for short signals, there is a huge amount of noise with very tiny relevant regions. Our results are shown in Figure 4.

First, as shown in Figure 4a, we vary the signal-to-noise ratio in MII, as defined by the length of the length of the relevant signal for each class. Intuitively, as this ratio increases, the signal becomes easier to identify. By updating the *receptor width* δ to match the signal-to-noise ratio as it is increases, we find that the *Moment Network* indeed succeeds in finding the discriminative moments in the timeline, achieving nearly-perfect accuracy even when the signal only takes up 6% of the timeline. Once the signal takes up 10% of the timeline, CAT consistently achieves 100% testing accuracy. We also find that the compared methods fail when the signal-to-noise ratio is lower than 0.1, achieving roughly 50% testing accuracy. This is expected as RNNs are classically hard to train on long series, especially with such noisy inputs.

Second, as shown in Figure 4b, we vary the *receptor width* parameter δ for a signal-to-noise ratio Δ of 0.04 to understand CAT ’s sensitivity to the proper selection of δ . We investigate the signal-to-noise ratio of 0.04 where CAT achieves only 75% accuracy, indicating potential sensitivity to hyperparameters (see Figure 4a). As expected, accuracy suffers both when δ is either too small (0.02) or too large (0.5). The optimal δ lies somewhere between 0.2 and 0.32 for this experiment. Quite interestingly, this is much larger than the data’s signal-to-noise (0.04). While a larger receptor width δ should capture signals more easily, suggesting that the receptor still filters out the noisy regions when they



(a) Effect of signal width Δ .



(b) Effect of δ with $\Delta = 0.04$.

Figure 4: CAT ’s performance on Synthetic MII dataset.

overlap with the receptor’s window. These results also indicate that CAT can be robust to overestimating δ .

5. Conclusions

In this work, we identify the open Attention-Based Classification problem for long and irregularly-sampled time series, which is a challenging and impactful setting common to many important domains. The Attention-Based Classification problem is to classify long irregularly-sampled time series based on small discriminative signals in the continuous timeline while learning to ignore irrelevant regions. Since prior methods rely on good selection of a set of timesteps at which to impute values, they struggle to classify time series in this setting, which we demonstrate experimentally. Using insights from prior methods, we then propose the Continuous-time Attention Policy Network (CAT), which generalizes previous works by learning to searching for short signals in a time series’ potentially-long timeline. CAT includes a reinforcement learning-based *Moment Network* that seeks discriminative moments in the timeline, positioning a novel *Receptor Network* that represents signals from *both* the values themselves and the patterns existing in the timing of the observations. Using a core *Transition Model* that learns to model the transition between moments, a *Discriminator Network* finally classifies the entire series. This approach can intuitively be extended to match the modeling paradigms proposed by other recent methods, like differential equation models and time-representational encodings. We validate our method on a wide range of experiments featuring four real datasets, ab-

lation studies, impacts of hyperparameter selection, a synthetic data highlighting CAT’s strengths, and timing experiments. Across the board, CAT consistently outperforms recent alternatives by successfully finding short signals in long time series.

References

- Bahdanau, D., Cho, K., and Bengio, Y. Neural machine translation by jointly learning to align and translate. *arXiv:1409.0473*, 2014.
- Cao, W., Wang, D., Li, J., Zhou, H., Li, L., and Li, Y. Brits: Bidirectional recurrent imputation for time series. In *NeurIPS*, 2018.
- Che, Z., Purushotham, S., Cho, K., Sontag, D., and Liu, Y. Recurrent neural networks for multivariate time series with missing values. *Scientific reports*, 8(1):6085, 2018.
- Chen, Y. and Chien, J. Continuous-time attention for sequential learning. In *AAAI*, 2021.
- Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. Learning phrase representations using rnn encoder–decoder for statistical machine translation. In *EMNLP*, pp. 1724–1734, 2014.
- De Brouwer, E., Simm, J., Arany, A., and Moreau, Y. Gru-ode-bayes: Continuous modeling of sporadically-observed time series. In *NeurIPS*, pp. 7379–7390, 2019.
- García-Magariño, I., Medrano, C., Plaza, I., and Oliván, B. A smartphone-based system for detecting hand tremors in unconstrained environments. *Personal and Ubiquitous Computing*, 20(6):959–971, 2016.
- Hartvigsen, T., Gerych, W., Thadajarassiri, J., Kong, X., and Rundensteiner, E. Stop&hop: Early classification of irregular time series. In *CIKM*, 2022.
- Hasani, R., Lechner, M., Amini, A., Rus, D., and Grosu, R. Liquid time-constant networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 7657–7666, 2021.
- Hochreiter, S. The vanishing gradient problem during learning recurrent neural nets and problem solutions. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 6(02):107–116, 1998.
- Hong, S., Xu, Y., Khare, A., Priambada, S., Maher, K., Aljiffry, A., Sun, J., and Tumanov, A. Holmes: Health online model ensemble serving for deep learning models in intensive care units. In *KDD*, pp. 1614–1624. ACM, 2020.
- Ismail, A. A., Gunady, M., Pessoa, L., Bravo, H. C., and Feizi, S. Input-cell attention reduces vanishing saliency of recurrent neural networks. In *NeurIPS*, pp. 10814–10824, 2019.
- Jia, J. and Benson, A. R. Neural jump stochastic differential equations. In *NeurIPS*, 2019.
- Kidger, P., Morrill, J., Foster, J., and Lyons, T. Neural controlled differential equations for irregular time series. In *NeurIPS*, 2020a.
- Kidger, P., Morrill, J., and Lyons, T. Generalised interpretable shapelets for irregular time series. *arXiv:2005.13948*, 2020b.
- Lechner, M. and Hasani, R. Learning long-term dependencies in irregularly-sampled time series. In *NeurIPS*, 2020.
- Lee, J. B., Rossi, R. A., Kim, S., Ahmed, N. K., and Koh, E. Attention models in graphs: A survey. *TKDD*, 13(6): 1–25, 2019.
- Li, S. C.-X. and Marlin, B. M. A scalable end-to-end gaussian process adapter for irregularly sampled time series classification. In *NeurIPS*, pp. 1804–1812, 2016.
- Li, S. C.-X. and Marlin, B. M. Learning from irregularly-sampled time series: a missing data perspective. In *ICML*, 2020.
- Lipton, Z. C., Kale, D., and Wetzel, R. Directly modeling missing data in sequences with rnns: Improved classification of clinical time series. In *MLHC*, pp. 253–270, 2016.
- Liu, J., Zhong, L., Wickramasuriya, J., and Vasudevan, V. uwave: Accelerometer-based personalized gesture recognition and its applications. *Pervasive and Mobile Computing*, 5(6):657–675, 2009.
- Mnih, V., Heess, N., Graves, A., et al. Recurrent models of visual attention. In *NeurIPS*, pp. 2204–2212, 2014.
- Mozer, M. C., Kazakov, D., and Lindsey, R. V. Discrete event, continuous time rnns. *arXiv:1710.04110*, 2017.
- Oh, J., Wang, J., and Wiens, J. Learning to exploit invariances in clinical time-series data using sequence transformer networks. In *MLHC*, pp. 332–347, 2018.
- Rubanov, Y., Chen, T. Q., and Duvenaud, D. K. Latent ordinary differential equations for irregularly-sampled time series. In *NeurIPS*, pp. 5321–5331, 2019.
- Rubin, D. B. Inference and missing data. *Biometrika*, 63(3):581–592, 1976.
- Salvi, C., Lemercier, M., and Gerasimovics, A. Neural stochastic pdes: Resolution-invariant learning of continuous spatiotemporal dynamics. In *Advances in Neural Information Processing Systems*, 2022.
- Schirmer, M., Eltayeb, M., Lessmann, S., and Rudolph, M. Modeling irregular time series with continuous recurrent units. In *International Conference on Machine Learning*, pp. 19388–19405. PMLR, 2022.
- Shukla, S. N. and Marlin, B. Interpolation-prediction networks for irregularly sampled time series. In *ICLR*, 2019.
- Shukla, S. N. and Marlin, B. M. Multi-time attention networks for irregularly sampled time series. In *ICLR*, 2021.
- Singh, B. P., Deznabi, I., Narasimhan, B., Kucharski, B., Uppaal, R., Josyula, A., and Fiterau, M. Multi-

- resolution networks for flexible irregular time series modeling (multi-fit). *arXiv:1905.00125*, 2019.
- Sood, E., Tannert, S., Müller, P., and Bulling, A. Improving natural language processing tasks with human gaze-guided neural attention. In *NeurIPS*, 2020.
- Tan, Q., Ye, M., Wong, G. L.-H., and Yuen, P. Cooperative joint attentive network for patient outcome prediction on irregular multi-rate multivariate health data. In *International Joint Conference on Artificial Intelligence*, 2021.
- Vaizman, Y., Ellis, K., and Lanckriet, G. Recognizing de-tailed human context in the wild from smartphones and smartwatches. *IEEE Pervasive Computing*, 16(4):62–74, 2017.
- Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- Zheng, K., Gao, J., Ngiam, K. Y., Ooi, B. C., and Yip, W. L. J. Resolving the bias in electronic medical records. In *KDD*, pp. 2171–2180. ACM, 2017.

Appendix

A. Dataset Descriptions

All datasets statistics are shown in Table 3. For each of the four Human Activity Recognition datasets from the EXTRASENSORY dataset (<http://extrasensory.ucsd.edu/>) (Vaizman et al., 2017) (WALKING, RUNNING, LYINGDOWN, and SLEEPING), we aim for an 80% training and 20% testing split *in time*, though this is challenging to control in practice. Thus, the exact ratio differs between the series. For each dataset, we further split off 10% of the training set for validation.

Table 3: Dataset Statistics. N_{train} denotes the number of training instances, N_{test} is the number of testing instances, Avg. T is the average number of observations per series, and C is the number of classes.

Dataset	N_{train}	N_{test}	Avg. T	C
MII	4000	1000	500	2
UWAVE	4030	448	94	8
WALKING	1616	1020	99	2
RUNNING	666	400	85	2
LYINGDOWN	6186	1240	80	2
SLEEPING	6462	2814	80	2

B. Further MII Experiments

Expanding on the synthetic experiment discussed in the Experiments section of our main paper, we also run all compared methods for each of the signal-to-noise ratios, the results of which are shown in Figure 7. Again, each series has 500 timesteps, only 3 of which are relevant to the classification task. The 3 relevant timesteps are evenly-spaced in a randomly-placed width- Δ window in the continuous timeline. In this experiment, all compared methods fail to classify these series, even when performing imputation with 500 timesteps, which does not delete the signal. Instead, they fail to *focus* on the the discriminative region and so cannot perform classification. On the contrary, CAT achieves nearly-perfect accuracy with a signal-to-noise ratio as low as .06, indicating that it indeed does find the relevant regions.

C. Timing Experiments

CAT’s *Transition Model* uses an RNN to model the transitions between *moments*, as opposed to the timestamps themselves. This hints that CAT should naturally be much faster than the compared methods. We confirm this by timing the training of all methods on the WALKING dataset with frequent imputation—see Figure 6. As expected, CAT

runs over seven times faster than the next slowest method while achieving much higher testing accuracy. This is particularly meaningful for long series in time-sensitive domains such as healthcare where a model’s inference time is hugely important (Hong et al., 2020). Our reported timing comparisons between compared methods is also largely consistent with prior works’ timing experiments (Shukla & Marlin, 2019). Also as expected, the GRU-D (Che et al., 2018), mTAN (Shukla & Marlin, 2021), and NCDE (Kidger et al., 2020a) run significantly slower than the other compared methods and so omit their results from this figure. Their Accuracies are much lower than CAT’s—see Table 2 in the main paper. All models were trained and evaluated on Intel Xeon Gold 6148 CPUs.

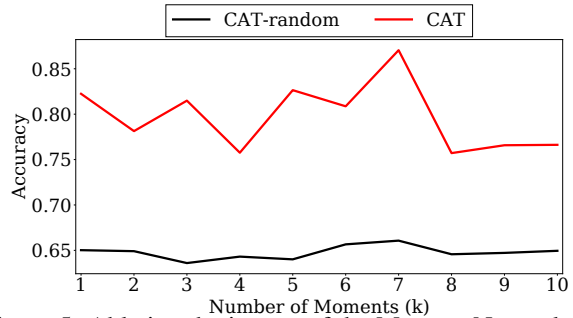


Figure 5: Ablating the impact of the Moment Network using the WALKING dataset. CAT outperforms CAT-random so the number of hops does not increase accuracy alone.

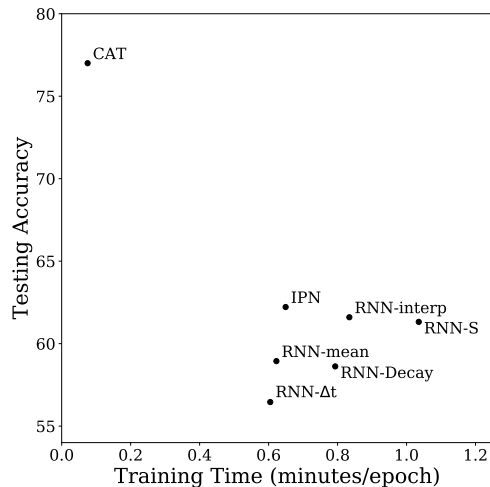


Figure 6: Timing performance for the high-resolution WALKING dataset. GRU-D and NCDE take over 3x longer than the next-slowest RNN-S and so is omitted from this figure.

D. CAT Hyperparameters

We experiment with three key hyperparameters of CAT for each dataset: The receptor-width δ , the hidden dimension

Table 4: Best hyperparameter settings for CAT.

Dataset	δ	Hidden Dim. of \mathcal{R}	Density
UWAVE	0.05	50	Off
Infrequent WALKING	0.2	50	Off
Infrequent RUNNING	0.05	50	On
Infrequent LYINGDOWN	0.05	50	Off
Infrequent SLEEPING	0.05	50	On
Frequent WALKING	0.2	50	Off
Frequent RUNNING	0.2	50	On
Frequent LYINGDOWN	0.2	50	On
Frequent SLEEPING	0.1	50	Off

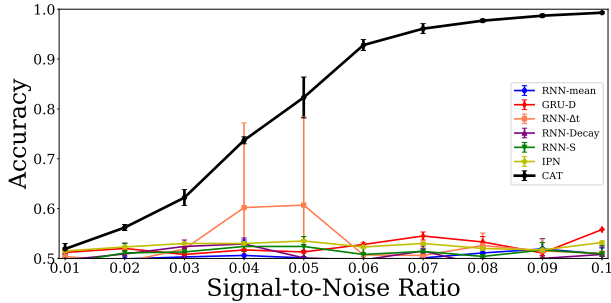


Figure 7: Effect of changing the signal width on accuracy.

of the Receptor Network \mathcal{R} , and whether or not to use the informative irregularity feature of CAT in the Receptor Network. Interestingly, we found that for \mathcal{R} , a hidden dimension of 50 seemed to consistently produce the best results. This hidden dimension largely controls the number of parameters in CAT and influences the timing experiments for which we also use a 50-dimensional representation. Our selections for δ values for different datasets are shown in Table 4.

We tune δ between three values: 0.05, 0.1, and 0.2. For UWAVE, $\delta = 0.05$ was best. $\delta = 0.05$ was also best for all infrequent EXTRASENSORY datasets except for WALKING, which used 0.2. $\delta = 0.2$ was chosen for all frequent EXTRASENSORY datasets except for SLEEPING, for which $\delta = 0.1$.

For δ , we observe that for the *infrequent* experiments, a smaller receptor width is largely the best option while a larger width is beneficial for the *frequent* experiments. This may be due to the fact that with the infrequent representation of the input series, closer focus on the comparatively-fuzzier signals is required. We also found that setting the number of steps $k = 3$ consistently outperformed larger and smaller values. While large values of k conceptually should still learn to classify effectively, in practice the more steps taken by a reinforcement learning agent per episode can make it more challenging to optimize effectively due to the credit assignment problem.

We also find that there are cases where it is not essential to use both channels—Values and Irregularity—in the receptor network. While using the irregularity channel (computed via the squared exponential kernel in the main paper) always leads to state-of-the-art performance by CAT, its omission can sometimes improve CAT’s performance slightly. When irregularity is an essential feature, however, this information cannot be removed. We show for which datasets this is true in Table 4. This may be a feature of (1) how the irregularity is represented—there are other approaches—and (2) how essential it is to the task. We recommend always using the irregularity channel as the potential downside of ignoring irregularity outweighs the minor benefits of omission in some cases.