

---

# Contrastive Pre-Training for Multimodal Medical Time Series

---

**Aniruddh Raghu**  
Massachusetts Institute of Technology  
araghu@mit.edu

**Payal Chandak**  
Massachusetts Institute of Technology

**Ridwan Alam**  
Massachusetts Institute of Technology

**John Guttag**  
Massachusetts Institute of Technology

**Collin M. Stultz**  
Massachusetts Institute of Technology

## Abstract

Clinical time series data are highly rich and provide significant information about a patient’s physiological state. However, these time series can be complex to model, particularly when they consist of multimodal data measured at different resolutions. Most existing methods to learn representations of these data consider only tabular time series (e.g., lab measurements and vitals signs), and do not naturally extend to modelling a full, multimodal time series. In this work, we propose a contrastive pre-training strategy to learn representations of multimodal time series. We consider a setting where the time series contains sequences of (1) high-frequency electrocardiograms and (2) structured data from labs and vitals. We outline a strategy to generate augmentations of these data for contrastive learning, building on recent work in representation learning for medical data. We evaluate our method on a real-world dataset, finding it obtains improved or competitive performance when compared to baselines on two downstream tasks.

## 1 Introduction

In various clinical settings, such as the intensive care unit (ICU), patients are closely monitored, resulting in time series data being generated for each patient. This highly rich data contains significant physiological information about a patient’s state and their disease progression over time [Johnson et al., 2016]. As a result, there have been efforts to conduct representation learning on these data, with the goal being to use these representations for various downstream predictive tasks [McDermott et al., 2021, Weatherhead et al., 2022, Tonekaboni et al., 2021, Yèche et al., 2021, Tipirneni and Reddy, 2022].

Although these works develop effective strategies to model clinical time-series data, they focus on learning representations of a time-series of unimodal, structured data alone (e.g. lab values). However, in reality, data originating from a patient’s encounter is significantly more complex, containing several modalities of time-series measured at different resolutions. For example, a given patient may have two distinct, complementary types of time-series recorded at periodic intervals: (1) high-frequency physiological signals (e.g., a 10 second electrocardiogram recorded at 240 Hz); and (2) structured labs and vitals signs. Extending existing methods to learn representations of these complex, hierarchical time-series is challenging, since they neither model sequences of high-frequency waveforms nor contend with the multimodal nature of the data stream.

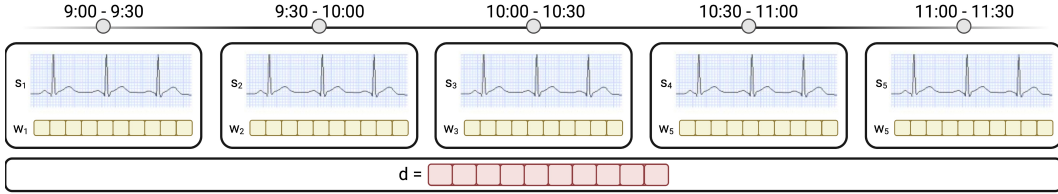


Figure 1: **An example of a multimodal clinical time-series or ‘trajectory’.** The trajectory  $\tau_n$  is characterized by an overall static vector  $d$  consisting of measurements that remain constant, and a time series of signals  $s_t$  and structured data  $w_t$  measured at each time step. Here, each time step is a 30 minute window.

In this paper, we take steps towards addressing this gap and outline an approach for learning representations of multimodal clinical time series. We describe a contrastive pre-training strategy where augmentations of the multimodal time series are generated by combining methods from recent work on representation learning for medical data [Yèche et al., 2021, Kiyasseh et al., 2021]. We evaluate our pre-training strategy on a real-world dataset of heart failure patients from a large tertiary hospital – Massachusetts General Hospital (MGH). Our initial results suggest that our approach yields pretrained models that outperform or perform competitively with baselines on two downstream tasks of clinical interest.

**Related work.** Contrastive objectives encourage the learned representations of different views coming from the same sample to be similar, and from different samples to be dissimilar [Chen et al., 2020a,b, He et al., 2020, Grill et al., 2020, Bardes et al., 2022]. Multiview contrastive learning objectives can generalize to any modality of data given a procedure to generate different views of each sample. Here, we adopt the contrastive framework from SimCLR [Chen et al., 2020a], and define procedures for generating augmentations of multimodal time-series.

Contrastive representation learning has previously been used to improve predictive performance on downstream clinical tasks [Weatherhead et al., 2022, Tonekaboni et al., 2021]. Prior work has focused on applying contrastive approaches to unimodal time-series, such as electrocardiograms [Kiyasseh et al., 2021, Gopal et al., 2021, Diamant et al., 2021, Oh et al., 2022] and sequences of tabular data (labs and vitals) [Yèche et al., 2021, Li et al., 2021]. Here, we study contrastive learning on multimodal clinical data, building on existing approaches to enable their application to our time-series data with ECGs, labs, and vitals.

## 2 Methods

### 2.1 Problem Setup

We use the term *trajectory* to refer to a sequence of physiological signals and structured data collected over time for a patient. This definition is motivated by our use-case in cardiovascular medicine, where patients may be monitored with telemetry devices that continuously record waveforms in addition to having lab tests and vitals signs periodically measured. The concept of a trajectory could be readily expanded to include other information, such as imaging or medications, depending on the context.

We frame the problem of representation learning as one of pre-training (PT). We assume access to a (large) PT dataset of patient trajectories, and a set of (likely smaller) fine-tuning (FT) datasets that have labels for downstream tasks. We first pre-train a model  $f_\theta$  on the PT dataset, and then evaluate the model based on its performance after FT on the downstream tasks.

More formally, we denote the PT dataset as  $\mathcal{D}_{PT} = \{\tau_n\}_{n=1}^{N_{PT}}$ , which contains  $N_{PT}$  trajectories (denoted  $\tau$ ). A trajectory  $\tau_n$  is of length  $T_n$ , and has the following structure:  $\tau_n = (d, \{(w, s)_t\}_{t=1}^{T_n})$ . Here,  $d \in \mathbb{R}^L$  represents a set of static features that do not change over the trajectory (such as demographic information or infrequently measured lab values). The sequence  $\{(w, s)_t\}_{t=1}^{T_n}$  contains a vector of structured data  $w \in \mathbb{R}^M$ , and a high-dimensional signal  $s \in \mathbb{R}^{C \times Z}$ , where  $C$  is the number of signal channels and  $Z$  is the number of samples in the signal (typically on the order of a few thousand). A visualization is shown in Figure 1.

During FT we have a collection of  $K$  tasks indexed by  $k$ , each with an accompanying FT dataset of trajectories  $\mathcal{D}_{\text{FT}}^{(k)} = \{(\tau_n, y_n)\}_{n=1}^{N_{\text{FT}}^{(k)}}$ . Each FT trajectory has an accompanying label  $y_n$  that the model is fine-tuned to predict. ‘24 hour mortality’ is an example of a fine-tuning task, where the model predicts if a patient will die in the next day.

## 2.2 Contrastive Learning for Trajectories

We use a contrastive learning approach to pre-train the model  $f_\theta$ . This is because contrastive approaches have previously been demonstrated to be (independently) effective on clinical signals data [Kiyasseh et al., 2021, Gopal et al., 2021] and structured data [Yèche et al., 2021]. In particular, we adopt a multiview contrastive approach inspired by SimCLR [Chen et al., 2020a]:

1. Sample a batch of trajectories from the PT dataset.
2. For each trajectory  $\tau_n$ , generate two augmented views of it.
3. Pass the augmented views through the representation model  $f_\theta$  and a projection head  $g_\phi$ .
4. Minimize the normalized temperature-scaled cross-entropy loss (NT-Xent) on these projections. The loss is averaged over all positive pairs that arise from the same underlying trajectory, as in Chen et al. [2020a]. The form of the loss function is detailed in the appendix.

**Trajectory-level augmentation.** This learning framework requires generating augmented views of each trajectory. We use the following approach for each type of data (further details in the appendix):

- **High-frequency signal  $s$ :** For each signal in the trajectory of length  $T_n$ , we form a pair of augmented views by first splitting the signal into two disjoint segments (e.g., taking the first 10 seconds as one view, and the second 10 seconds as the second view) and then applying random masking and noise addition as augmentations to each view independently, similar to Kiyasseh et al. [2021]. The intuition is that two segments of a signal that are close in time should encode similar physiological state, and can therefore be considered paired views.
- **Structured-time series data  $w$ :** The tabular data sequence forms a  $T_n \times M$  matrix over all timesteps of the trajectory. Following prior work [Yèche et al., 2021], we apply two data augmentation strategies to this matrix: history cutout and channel dropout.
- **Static features  $d$ :** Following Bahri et al. [2021], we obtain two views of by (1) randomly selecting features to corrupt; (2) corrupting them by replacing them values drawn from a uniform distribution over the values that feature takes in the training set.

We note that the individual aspects of our approach (contrastive loss function and augmentations for signals/tabular data) have been used in prior work; however, combining these together in a single pipeline for pre-training jointly on a time-series of signals and structured data is a new direction. Investigating alternative augmentation strategies for the multimodal data is an important direction of future work.

## 3 Experiments

We evaluate our method on a dataset from a large tertiary hospital. We assess whether our pre-training scheme on multimodal trajectories can improve predictive performance on downstream tasks over (1) no pre-training and (2) pre-training on structured data alone.

**Dataset.** Our cohort consists of a total of 51,921 trajectories from 10,990 unique patients with a diagnosis of heart failure from the Massachusetts General Hospital (MGH). Each element in a trajectory corresponds to data from 30 minutes of a patient’s stay, and for simplicity, we fix each trajectory to be of length 10. Within each trajectory,  $d \in \mathbb{R}^9$  contains information from a basic metabolic panel and blood pressures. These values are measured infrequently in our dataset (every 6-12 hrs) so we treat them as constant within a trajectory.  $w_t \in \mathbb{R}^9$  represents summary statistics related to heart rate and SpO2 within each 30 minute window.  $s_t \in \mathbb{R}^{4 \times 7200}$  is a 4-channel ECG signal (leads I, II, III, and V1) measured for 30 seconds at 240 Hz within that 30 minute window. Missing structured data are forward-fill imputed where possible (for example, if part of a time series) and otherwise imputed with the mean over the training dataset. Missing signals are represented with zeros. Further dataset details are in the appendix.

**Downstream tasks.** We evaluate models on two downstream tasks (Table 1):

Table 1: Statistics for fine-tuning tasks.

Task	# Patients	# Trajectories	Prevalence
Elevated mPAP	1,484	5,225	77.8%
24 hour mortality	10,990	15,835	1.3%

Table 2: **Multimodal contrastive pre-training obtains improved or competitive performance when compared to baselines.** Using our multimodal pre-training approach improves performance over (1) not pre-training and (2) training with only the structured data on the elevated mPAP detection task. On mortality prediction, our method performs comparably with baselines that do not use ECGs, and improves over the RandInit baseline that uses the multimodal trajectory as input, suggesting the pre-trained model better utilizes the rich feature space.

Task (Metric)	RandInit	RandInit (No ECG)	Contrastive	Contrastive (No ECG)
Elevated mPAP (AUROC)	0.70 ± 0.01	0.67 ± 0.01	<b>0.72 ± 0.01</b>	0.66 ± 0.01
24 hour mortality (AUPRC)	0.07 ± 0.03	<b>0.12 ± 0.01</b>	<b>0.12 ± 0.01</b>	<b>0.12 ± 0.02</b>

- **Elevated mPAP:** Predict whether a patient’s mean Pulmonary Arterial Pressure (mPAP) is abnormally high (greater than 20 mmHg) at a time marker, given the trajectory data until that point. This task is of clinical interest since the mPAP is typically only measured via an invasive study, and so inferring whether it is abnormal using minimally invasive signals (i.e., the ECG, labs, and vitals signs) is valuable. Prior work [Schlesinger et al., 2022] has studied a similar task from the 12-lead ECG, but not in the context of trajectory data, as we do here.
- **24 hour mortality:** Predict whether the patient is likely to die in the next day, given a trajectory from the start of a patient’s admission. This task is commonly used to evaluate pre-training strategies on clinical time series data from the ICU [McDermott et al., 2021, Yèche et al., 2021].

**Methods evaluated.** Motivated by our experimental goals, we evaluate the following methods:

- **RandInit:** We train a model from random initialization on each FT task on the full trajectory (statics, structured data timeseries, and signals timeseries).
- **RandInit (No ECG):** We train a model from random initialization on only the tabular data (statics and structured data timeseries).
- **Contrastive:** we follow the pre-training strategy outlined in Section 2.2 on the PT dataset, before fine-tuning this model on the FT tasks.
- **Contrastive (No ECG):** we pre-train on the tabular data alone, following the same augmentation process as described above, before fine-tuning the model on the tabular data alone for each FT task.

**Experimental setup.** We split our dataset on a per-patient level into an 80/20 development/test set split and use 20% of the development set as a validation set. We standardize the encoder  $f_\theta$  to be a CNN-MLP-GRU. We use a 2-layer MLP for the projection head  $g_\phi$ . The model architecture is shown in the appendix, Figure 2. We conduct pre-training for 10 epochs, using a batch size of 128. We set augmentation strength hyperparameters to values used in prior work [Kiyasseh et al., 2021, Yèche et al., 2021, Bahri et al., 2021], and set the temperature of the NT-Xent loss to 0.1, following [Yèche et al., 2021]. Fine-tuning is for 10 epochs, using early stopping based on validation loss. All models are trained with the Adam optimizer with a learning rate of 1e-3. We run 5 random seeds at fine-tuning time, and report the median/half IQR of test AUROC for elevated mPAP and test AUPRC for 24 hour mortality (since the label has very low prevalence). Further experimental details are in the appendix.

**Results.** Results are shown in Table 2. We highlight two key findings:

- **Multimodal contrastive pre-training improves or obtains competitive performance.** The full contrastive strategy improves performance over all baselines on Elevated mPAP, and performs comparably to the best baselines on the 24 hour mortality task. This result suggests that the benefit of multimodal data in pre-training may be task-dependent. For the mortality task, with the amount of data available, the structured data alone are likely highly predictive, and so also including the ECGs does not improve performance. Improved performance may be observed at lower data regimes, as explored in prior work [McDermott et al., 2021].
- **A more complex feature space can worsen performance without pre-training.** On the 24 hour mortality task, including the ECGs in a model without pre-training worsens performance noticeably,

which is likely a result of overfitting with a more complex feature space. The pre-trained model does not show this behaviour, suggesting that the pre-training process helps mitigate this issue.

## 4 Conclusion

In this work, we outlined a contrastive pre-training strategy for multimodal clinical time series. Our method models trajectories of electrocardiogram signals and structured measurements from lab tests and vitals signs. In an experimental evaluation on a dataset of patients with heart failure, we find that our strategy obtains competitive or improved performance when compared to baselines on two downstream tasks.

**Scope and Limitations.** Given that this is an early-stage investigation, we only evaluate methods on two downstream tasks, and did not conduct extensive search over various model and training hyperparameters. Addressing these, and comparing to other pre-training schemes (including an investigation into different methods of augmenting the multimodal trajectories) is a central direction of future work. In addition, examining other evaluation modes for the pre-trained models (e.g, linear evaluation, training a multilayer perceptron on frozen representations) are important to consider.

**Social impact.** Our contribution in this work is mostly methodological. However, given that our application domain is in medicine, a high-risk setting, it is important that our method is more thoroughly validated in larger retrospective and prospective studies before any real-world use. This is to understand any potential risks with its use in practice.

## References

- D. Bahri, H. Jiang, Y. Tay, and D. Metzler. Scarf: Self-supervised contrastive learning using random feature corruption. In *International Conference on Learning Representations*, 2021.
- A. Bardes, J. Ponce, and Y. LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. In *ICLR*, 2022.
- T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020a.
- T. Chen, S. Kornblith, K. Swersky, M. Norouzi, and G. Hinton. Big self-supervised models are strong semi-supervised learners. *arXiv preprint arXiv:2006.10029*, 2020b.
- N. Diamant, E. Reinertsen, S. Song, A. Aguirre, C. Stultz, and P. Batra. Patient contrastive learning: a performant, expressive, and practical approach to ecg modeling. 2021.
- B. Gopal, R. W. Han, G. Raghupathi, A. Y. Ng, G. H. Tison, and P. Rajpurkar. 3kg: Contrastive learning of 12-lead electrocardiograms using physiologically-inspired augmentations. 2021.
- J.-B. Grill, F. Strub, F. Alché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar, et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020.
- K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020.
- A. E. Johnson, T. J. Pollard, L. Shen, L.-w. H. Lehman, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. Anthony Celi, and R. G. Mark. MIMIC-III, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016.
- D. Kiyasseh, T. Zhu, and D. A. Clifton. Clocs: Contrastive learning of cardiac signals across space, time, and patients. In *International Conference on Machine Learning*, pages 5606–5615. PMLR, 2021.
- Y. Li, M. Mamouei, G. Salimi-Khorshidi, S. Rao, A. Hassaine, D. Canoy, T. Lukasiewicz, and K. Rahimi. Hi-behrt: Hierarchical transformer-based model for accurate prediction of clinical events using multimodal longitudinal electronic health records. *arXiv preprint arXiv:2106.11360*, 2021.
- M. McDermott, B. Nestor, E. Kim, W. Zhang, A. Goldenberg, P. Szolovits, and M. Ghassemi. A comprehensive ehr timeseries pre-training benchmark. In *Proceedings of the Conference on Health, Inference, and Learning*, pages 257–278, 2021.
- J. Oh, H. Chung, J.-m. Kwon, D.-g. Hong, and E. Choi. Lead-agnostic self-supervised learning for local and global representations of electrocardiogram. In *Conference on Health, Inference, and Learning*, pages 338–353. PMLR, 2022.
- D. E. Schlesinger, N. Diamant, A. Raghu, E. Reinertsen, K. Young, P. Batra, E. Pomerantsev, and C. M. Stultz. A deep learning model for inferring elevated pulmonary capillary wedge pressures from the 12-lead electrocardiogram. *JACC: Advances*, 1(1):100003, 2022.
- S. Tipirneni and C. K. Reddy. Self-supervised transformer for sparse and irregularly sampled multivariate clinical time-series. *ACM Trans. Knowl. Discov. Data*, 1(1), 2022.
- S. Tonekaboni, D. Eytan, and A. Goldenberg. Unsupervised representation learning for time series with temporal neighborhood coding. *arXiv preprint arXiv:2106.00750*, 2021.
- A. Weatherhead, R. Greer, M.-A. Moga, M. Mazwi, D. Eytan, A. Goldenberg, and S. Tonekaboni. Learning unsupervised representations for ICU timeseries. In *Proceedings of the Conference on Health, Inference, and Learning*, volume 174 of *Proceedings of Machine Learning Research*, pages 152–168. PMLR, 07–08 Apr 2022.
- H. Yèche, G. Dresdner, F. Locatello, M. Hüser, and G. Rätsch. Neighborhood contrastive learning applied to online patient monitoring. In *International Conference on Machine Learning*, pages 11964–11974. PMLR, 2021.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
  - (b) Did you describe the limitations of your work? [Yes] See Section 4.
  - (c) Did you discuss any potential negative societal impacts of your work? [Yes] See Section 4.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [N/A]
  - (b) Did you include complete proofs of all theoretical results? [N/A]
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [No] The dataset is private so we cannot include it here. We describe key aspects of our experimental setup in the main text and appendix.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes]
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] We describe these details in the appendix.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [N/A]
  - (b) Did you mention the license of the assets? [N/A]
  - (c) Did you include any new assets either in the supplemental material or as a URL? [No]
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [Yes] We describe these details in the appendix.
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [Yes] We describe these details in the appendix.
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [Yes] We have IRB approval for this study – details are provided in the appendix.
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

## A Further Methods Details

**Loss function.** Given a positive pair of projections  $z_i$  and  $z_j$  that originate from different augmentations of the same underlying sample, the NT-Xent loss function is computed as follows:

$$\mathcal{L}(z_i, z_j) = -\log \frac{\exp(\text{sim}(z_i, z_j)/\tau)}{\sum_{k=1}^{2B} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(z_i, z_k)/\tau)}, \quad (1)$$

where  $\text{sim}(a, b) = a^T b / (\|a\| \|b\|)$  is cosine similarity and  $\tau$  is the temperature hyperparameter. This is the exact form used in prior work [Chen et al., 2020a].

## B Further Experimental Details

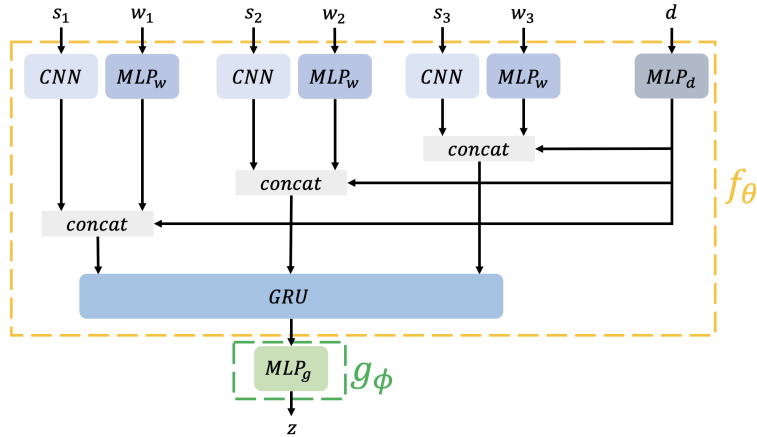


Figure 2: **Model architecture used in our experiments.** We show the architecture used to model trajectories in a scenario where the input trajectory has 3 timesteps.

**Dataset.** Our dataset is derived from the electronic medical record at the Massachusetts General Hospital (MGH) and was obtained with IRB approval, protocol number 2020P003053. Since the dataset has some identifiable information, all computations are performed on a server that sits behind the hospital firewall.

**Architectural details.** We use the following architectures for the encoder and projection head:

- **Encoder ( $f_\theta$ ):** Each signal in the trajectory is passed through a ResNet-styled 1-D CNN encoder with global average pooling. We base our CNN encoder model off a ResNet-18 architecture with kernel size of 15. Following global average pooling over the temporal dimension, each signal is projected into 128 dimensions with a linear layer. The structured data at each timestep is embedded with a 2-layer fully-connected network with 128 hidden units and ReLU activation at each layer, and this embedding is then concatenated with the signal embedding. The static features are passed through a different 2-layer fully-connected network with 128 hidden units and ReLU activation at each layer, and then concatenated with the embeddings of the signal and structured data timeseries at each timestep. The resulting sequence of vectors is passed into a 4-hidden layer GRU with hidden size of 384, with the last hidden state of the GRU being used as the overall trajectory embedding vector.
- **Projection head ( $g_\phi$ ):** This is a 2-layer fully connected network with batch normalization and ReLU activation with 128 hidden units. It takes the last hidden state of the GRU as input. The resulting projection is normalized (following Chen et al. [2020a]) before computing the NT-Xent loss over the batch.

Figure 2 shows the model architecture in a scenario in which the input trajectory has 3 timesteps.

**Augmentations and hyperparameters.** We use the following augmentations in the contrastive learning pipeline. Unless otherwise mentioned, the strength of the transformations were chosen based on the values reported in the original works.



- Random signal masking [Kiyasseh et al., 2021]: Randomly set 50% of the signal to zero.
- Random signal noise: Add Gaussian noise with variance 0.1 to the signal.
- History cutout [Yèche et al., 2021]: With probability 0.8, randomly set 50% of the timesteps in each feature dimension to zero (i.e., impute with the mean). The masking fraction of 50% is different to the value of 17% used in the original paper, and was determined based on validation loss on the downstream tasks.
- Channel dropout [Yèche et al., 2021]: For each feature dimension, with probability 0.2, set all the values over time to zero (i.e., impute with the mean).
- Feature corruption [Bahri et al., 2021]: For each feature dimension, with probability 0.6, replace the value that feature takes with a random value drawn from the uniform distribution over the values that feature takes in the training data.

**Other hyperparameters.** We did not search over architectural choices, such as embedding dimensions for the different modalities, and instead opted to standardize them to be the same (128). The learning rate of  $1e-3$  was deduced based on analyzing learning curves and seeing that this value resulted in stable training across all models.

**Compute details.** All models were trained on a single NVIDIA Quadro RTX 8000 GPU. Pre-training takes about 10 hours, and fine-tuning on an individual task takes about 1 hour. Pre-training uses approximately 20 GB of GPU memory, and fine-tuning uses approximately 10 GB of GPU memory.