

# A synthetic dataset for Time Series Super-Resolution with Deep Learning

Julio Ibarra-Fiallo<sup>1</sup>, Juan A. Lara<sup>2</sup>, and D'hamar Agudelo-Moreno<sup>1</sup>

<sup>1</sup>Colegio de Ciencias e Ingenierías, Universidad San Francisco de Quito, Cumbayá, Ecuador

<sup>2</sup>Universidad de Córdoba, Córdoba, España

\*corresponding author: Julio Ibarra-Fiallo (jibarra@usfq.edu.ec)

## ABSTRACT

The increasing application of time-series analysis in fields like biomedical engineering or telecommunications emphasizes the need for high-quality data to train and evaluate advanced machine learning models. Acquiring real-world temporal data at suitable resolutions is often limited by ethical, economic, or practical constraints. We introduce CoSiBD (Complex Signal Benchmark Dataset for Super-Resolution), a synthetic dataset designed for reproducible time-series super-resolution research. CoSiBD provides 2,500 high-resolution signals ( $N = 5,000$  samples each over a reference domain  $\tau \in [0, 4\pi]$ ) with aligned low-resolution versions at four levels (150, 250, 500, and 1,000 samples) obtained via uniform decimation. Signals are distributed as NumPy arrays, plain text, and JSON, with comprehensive metadata describing segment structure, generation parameters, and seeds for full reproducibility. The generator produces diverse non-stationary behaviors through piecewise frequency modulation and spline-based amplitude envelopes, and provides both clean and noisy variants. Technical validation analyzes spectral properties and reports baseline SR benchmarking and transfer experiments on EEG and speech to illustrate utility on real-world signals.

## Background & Summary

The analysis and simulation of temporal signals are fundamental across science and engineering. These techniques provide critical insights into dynamic processes in multiple domains. For instance, in biomedical research, physiological time series such as electroencephalography (EEG) and electrocardiography (ECG) support the study of neural and cardiac function<sup>1-3</sup>. Another example is telecommunications, which rely on signal processing to ensure data fidelity across noisy media<sup>4</sup>. Developing robust tools for interpreting time-varying data continues to support both scientific discovery and practical applications.

Recent advances in deep learning have contributed significantly to this field by enabling automatic extraction of complex features from raw signals. Deep learning approaches, such as Convolutional Neural Networks (CNNs) or Generative Adversarial Networks (GANs) have demonstrated improved performance over traditional techniques in image, speech, and time-series processing tasks<sup>5,6</sup>. These models support fine-grained signal reconstruction, allowing researchers to explore temporal dynamics in new ways.

Despite this progress, deep learning methods for temporal signal processing often require large quantities of labeled, high-quality data. Access to such data is frequently constrained. For instance, in medicine, data access is constrained by medical privacy regulations such as the General Data Protection Regulation (GDPR) and the Health Insurance Portability and Accountability Act (HIPAA)<sup>7</sup>. In other domains, including telecommunications, data availability is limited by proprietary protocols and the high cost of acquiring large-scale channel sounding datasets for diverse environments<sup>8</sup>. These limitations are particularly relevant in super-resolution (SR) tasks, where models require paired low- and high-resolution signals for effective training. Here, temporal SR refers to reconstructing a higher-sample-count (higher temporal resolution) discrete sequence from a uniformly decimated low-resolution version.

This task has broad potential. In biomedical monitoring and sensing, for instance, SR can help reconstruct higher-resolution physiological time series (e.g., EEG), potentially improving the analysis of subtle physiological irregularities<sup>3</sup>. SR also applies to audio/speech enhancement and telecommunications (e.g., neural audio upsampling and bandwidth extension, relevant to telephony and compression)<sup>9,10</sup>.

Deep learning offers adaptive alternatives to these traditional methods. For instance, CNNs are capable of modeling spatio-temporal structure present in real datasets across these domains. Preliminary work on synthetic time-series generation indicates

38 potential for SR<sup>10,11</sup>, but the lack of accessible, high-quality paired datasets remains a significant barrier to progress.

39  
40 Synthetic datasets offer one solution to this problem, allowing researchers to design reproducible training environments  
41 that reflect the characteristics of real-world signals. Prior studies have used synthetic data in domains such as biomedical  
42 signal analysis<sup>12</sup> and wireless communications<sup>13</sup>, demonstrating that synthetic approaches can help simulate complexity while  
43 avoiding legal and practical restrictions associated with real-world data.

44  
45 To support research in super-resolution for time-series data, we present the Complex Signal Benchmark Dataset (CoSiBD).  
46 CoSiBD is a synthetic dataset composed of time-series signals with variable resolution, frequency characteristics, and noise  
47 levels. As it is designed to resemble real-world signals, our dataset is intended with a double purpose: a) to provide a resource  
48 for training and evaluating deep learning SR models under controlled, reproducible conditions, which can constitute a sort  
49 of benchmark for this problem; and b) a resource for training deep learning models to be used for SR of real-world signals  
50 (either directly or finetuning them with the real data), particularly in scenarios where real signals are scarce and not enough  
51 for a complete training of those models. It includes non-stationary, piecewise-structured signals (via non-uniform interval  
52 partitioning with change-points), multiple levels of resolution and noise, a technical validation suite, and publicly available  
53 Python code to facilitate use. CoSiBD has been previously used in research presented at the International Conference on Signal  
54 Processing and Machine Learning<sup>10</sup>, with good preliminary results for signal reconstruction using deep learning. CoSiBD is  
55 made available to support further development in deep learning approaches for temporal super-resolution.

56 To further position CoSiBD with respect to existing public synthetic time-series resources, we summarize in the next  
57 subsection representative datasets and simulators and highlight the practical gap addressed by our approach.

## 58 Related synthetic time-series resources

59 Publicly available synthetic resources for temporal signals exist, but they are typically designed for tasks other than time-series  
60 SR, or they target a specific domain. In wireless communications, the RadioML family provides large collections of synthetic  
61 complex I/Q sequences with varying SNR (Signal-to-noise ratio) and channel impairments, mainly to benchmark automatic  
62 modulation classification rather than paired SR reconstruction<sup>13–15</sup>. In biomedical signal processing, physiological simulators  
63 such as ECGSYN (electrocardiography) and SEREEGA (EEG) enable controlled generation with tunable morphology, sampling  
64 settings, and noise, supporting method development when real data access is constrained<sup>12,16,17</sup>. In power systems, LoadGAN  
65 provides multi-resolution generation of load time series across sampling rates and time horizons (from sub-second to long-term  
66 scales), but it is not distributed as a standardized paired SR benchmark<sup>18</sup>. Domain-specific paired low-/high-resolution training  
67 data can also be produced via physical forward modeling, e.g., low- and high-resolution 1D seismic traces for learning-based  
68 resolution enhancement<sup>19</sup>.

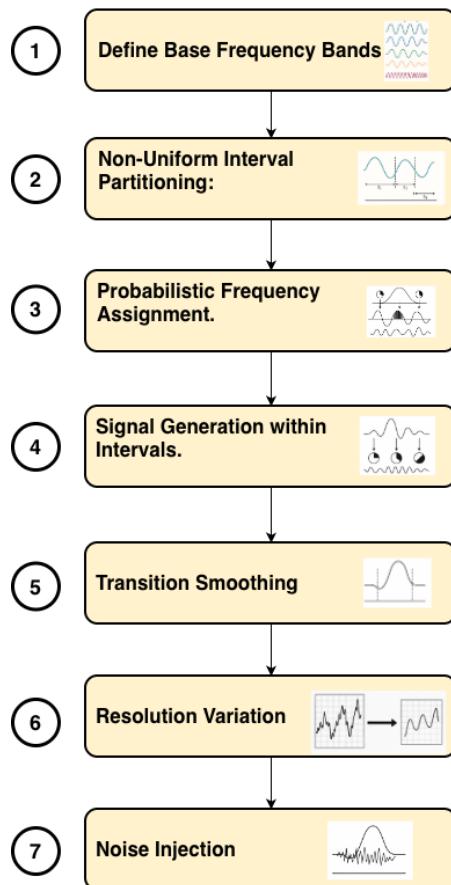
69 Table 1 summarizes these representative resources. Columns indicate the primary \*\*Domain\*\*, the \*\*Form\*\* of distribution  
70 (fixed dataset vs. generator), and the specific \*\*Noise\*\* models included. We explicitly check for \*\*Paired LR–HR SR\*\*  
71 capability (Low Resolution - High Resolution); resources marked as “Configurable” in this column allow users to generate  
72 pairs by running simulators at different settings, but typically do not distribute a standardized, fixed benchmark for SR. The  
73 \*\*Reproducibility granularity\*\* column notes whether exact regeneration is supported at the sample level. CoSiBD is designed  
74 to close the gap by providing a fixed, standardized set of LR–HR pairs with explicit nuisance modeling (noise and structured  
75 interference) and comprehensive metadata, enabling reproducible benchmarking across multiple difficulty levels (defined by  
76 varying downsampling factors and noise intensities).

## 77 Methods

78 The methodology used to generate the synthetic temporal signals that constitute the CoSiBD dataset is illustrated in Figure 1.  
79 The signal generation process is designed to produce time series exhibiting structural properties commonly observed in  
80 real-world temporal data, including variable frequency content, smooth transitions, and intermittent high-frequency activity. A  
81 key aspect of the procedure is the generation of signals at multiple temporal resolutions, enabling the construction of paired  
82 datasets for super-resolution (SR) benchmarking. REV

83  
84  
85 **Design rationale inspired by real signals.** It is important to note that one of the main applications of the proposed dataset will  
86 be the training of deep learning models to be used for SR purposes in other real-world datasets, like, for instance, physiological  
87 or speech signals. Therefore, there is a need for our dataset to resemble real-world data. In particular, real signals exhibit (i)  
88 non-stationary regime changes, (ii) coexisting low- and high-frequency components with intermittent transients, (iii) smooth  
89 amplitude-envelope evolution, and (iv) slow baseline drift and measurement noise. CoSiBD instantiates these properties  
90 via non-uniform interval partitioning with change-points, separate low/high-frequency bands, spline-based envelopes and

## CoSiBD Dataset Generation Process



**Figure 1.** Schematic overview of the CoSiBD signal generation process.

| Resource   | Domain   | Form                       | Paired LR–HR SR               | Multi-resolution                     | Noise / artifacts  | Reproducibility granularity                                       |
|--|--|----------------------------|-------------------------------|--------------------------------------|--|---|
| CoSiBD (this work)                                     | Generic time series (complex-structured signals) | Dataset generator          | Yes (LR → HR targets)         | Yes (150/250/500/1000 → 5000)        | Gaussian + structured interference; primary benchmark uses direct decimation | Per-signal metadata; deterministic regeneration (seed-controlled) |
| RadioML 2016.10A <sup>13, 14</sup>                     | Wireless communications (I/Q)                    | Dataset                    | No (classification benchmark) | N/A (not SR)                         | Variable SNR + channel impairments   | Dataset-level (labels/SNR); not per-sample “recipe”               |
| RadioML 2018.01A <sup>15</sup>                         | Wireless communications (I/Q)                    | Dataset                    | No (classification benchmark) | N/A (not SR)                         | Simulated channel effects + SNR variability                                  | Dataset-level; not SR-paired                                      |
| ECGSYN <sup>12, 16</sup>                               | ECG (physiology)                                 | Simulator/tool             | Configurable                  | Configurable (via sampling settings) | Model-based; supports controlled variability                                 | Configurable via simulator parameters (user-defined)              |
| SEREEGA <sup>17</sup>                                  | EEG (physiology)                                 | Simulator/toolbox          | Configurable                  | Configurable (user-defined)          | Supports noise and event-related components                                  | Configurable via simulator parameters (user-defined)              |
| LoadGAN <sup>18</sup>                                  | Power systems load time series                   | Generator/tool             | No (generation)               | Yes (variable sampling rates)        | Domain-specific variability (load patterns)                                  | Tool-based; generation is configurable                            |
| Synthetic LR–HR seismic traces (example) <sup>19</sup> | Seismic traces (geophysics)                      | Paper-specific paired data | Yes (LR–HR pairs)             | Typically limited (study-specific)   | Study-dependent  | Paired data available for the study; limited generality           |

**Table 1.** Representative publicly available synthetic time-series datasets and simulators related to signal processing and learning. “Form” indicates whether the resource is distributed primarily as a fixed dataset or as a simulator/generator. “Reproducibility granularity” summarizes whether exact per-sample regeneration is supported via documented parameters and seeds.

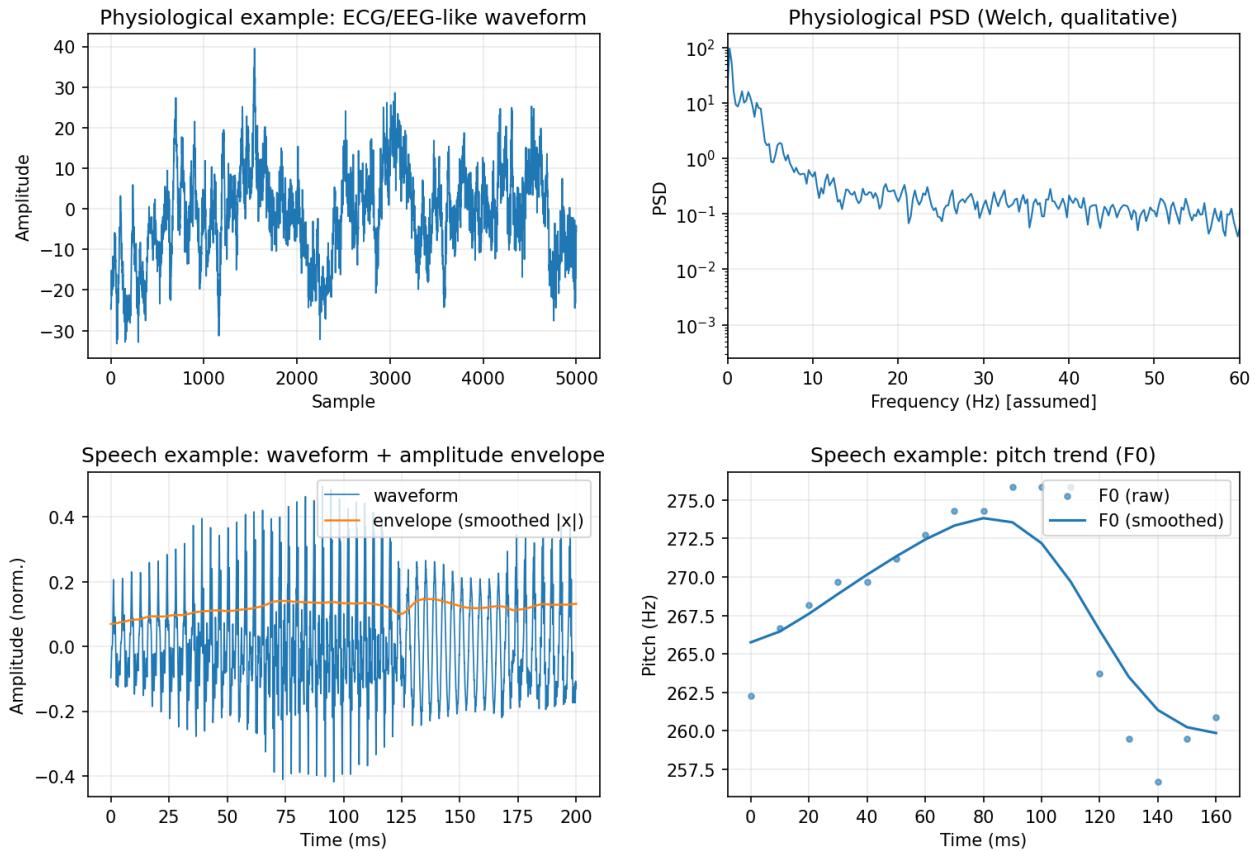
frequency profiles, and explicit offset/noise terms. Figure 2 provides qualitative examples of these motivating properties found in real-world time series; the main goal of our dataset is to be able to capture challenging structure for SR benchmarking rather than match a specific domain distribution.

Figure 2 presents four subfigures labeled A–D to illustrate these concepts. Figure 2A (top left) shows a representative physiological example (ECG/EEG-like waveform; based on a synthetic ECG model<sup>12</sup> and the PhysioNet ECGSYN implementation<sup>16</sup>) illustrating non-stationary activity. Figure 2B (top right) displays the corresponding spectrogram, highlighting structured spectral content. Figure 2C (bottom left) presents a speech signal segment (source: VCTK corpus<sup>20</sup>) showing clear amplitude-envelope dynamics. Figure 2D (bottom right) shows the speech spectrogram, revealing a smoothly varying pitch (F0) trend.

The signal generation pipeline involves the following steps:

1. **Base frequency band definition:** A set of distinct frequency bands is defined to represent the underlying spectral content of the signals. These can be adjusted to reflect application-specific characteristics.
2. **Non-uniform interval partitioning:** The total signal duration is divided into multiple intervals of variable length. The interval lengths are determined probabilistically to introduce variability in the signal structure.
3. **Frequency assignment:** Each interval is assigned a dominant frequency band, sampled according to a predefined probability distribution. This introduces spectral variation over time.
4. **Signal synthesis:** A sinusoidal waveform, or a combination of sinusoids within the assigned frequency band, is generated for each interval. Signal parameters such as amplitude and phase are configurable.
5. **Transition smoothing:** To avoid discontinuities at interval boundaries, a smoothing function is applied to overlapping segments. This ensures gradual transitions between intervals with different frequency content.

Real-signal properties motivating CoSiBD design (qualitative examples)



**Figure 2.** Qualitative real-signal properties motivating the CoSiBD design. The physiological example illustrates non-stationarity in the waveform and structured spectral content; the speech example illustrates amplitude-envelope dynamics and a smoothly varying pitch (F0) trend. These observations motivate CoSiBD mechanisms such as regime partitioning with change-points, low/high-frequency bands, and spline-based envelopes/frequency profiles.

6. **Resolution variation:** All signals are initially synthesized at a high temporal resolution (5,000 samples over the domain [0,  $4\pi$ ]). Lower-resolution versions are created using simple decimation (uniform subsampling). This keeps the SR task aligned with reconstructing the original high-resolution target; the low-resolution observation is obtained by subsampling the original sequence without pre-filtering. Reconstructing low-pass filtered signals is not an objective of CoSiBD. For reproducibility, given a high-resolution sequence  $x_{\text{HR}}[n]$  of length  $N = 5000$  and a target low-resolution length  $M \in \{1000, 500, 250, 150\}$ , we form  $x_{\text{LR}}[i] = x_{\text{HR}}[n_i]$  using the fixed index set  $n_i = \left\lfloor \frac{i(N-1)}{M-1} + 0.5 \right\rfloor$  for  $i = 0, \dots, M-1$  (applied identically to the time array). This reduces to standard stride decimation when  $M$  divides  $N$ .
7. **Noise injection:** Controlled levels of synthetic noise are added to the signals to emulate different data acquisition scenarios. Two noise types are implemented: (1) Additive white Gaussian noise (AWGN) with configurable standard deviation (relative to signal RMS amplitude), representing broadband sensor thermal noise; and (2) structured sinusoidal noise bursts (deterministic sinusoidal components), representing narrow-band interference such as powerline hum (50/60 Hz). Noise is applied probabilistically (in the released dataset, `noise_profile` records `p_has_noise=0.5` per signal, and the realized subset with noise is approximately half). When noise is present, the specific parameters (type, amplitude, frequency for structured noise) are recorded in the metadata, allowing users to benchmark denoising or super-resolution under specific degradation conditions.

**Rationale for structured 50/60 Hz interference and noise.** Real measurement pipelines frequently contain narrow-band interference (e.g., mains hum) superimposed on broadband sensor noise. To reflect this common acquisition artifact, CoSiBD includes an optional structured sinusoidal component in addition to Gaussian noise. CoSiBD signals are generated over a reference domain (by default  $\tau \in [0, 4\pi]$ ); interpreting  $\tau$  as physical time (and therefore reporting frequencies in Hz) requires an explicit time scaling. Throughout this manuscript we adopt an illustrative convention that maps the reference domain to a duration  $T = 4\pi$  seconds, under which the structured component can be interpreted as a 50/60 Hz-like powerline interference term, while the broadband term represents the measurement noise floor.

Figure 3 illustrates this qualitative motivation through four panels (A–D). Panel A displays a synthetic wideband noise floor. Panel B adds the structured narrow-band component. Panel C shows the combined time-domain signal, and Panel D presents the frequency spectrum, clearly showing the powerline-like artifact. This explicitly models the periodic contamination observed in real recordings<sup>21</sup>; the intent is not to reproduce a specific device transfer function but to include realistic nuisance factors that SR models must handle.

**Sampling units and frequency interpretation.** CoSiBD signals are provided as discrete sequences  $x[n]$  (e.g.,  $N = 5,000$  samples) that are directly used as inputs/targets by SR models. The internal generation domain  $\tau \in [0, 4\pi]$  is a reference parameterization; interpreting it as physical time requires choosing a duration  $T$  (in seconds) for the reference interval. Under this convention, the implied sampling rate is  $f_s = N/T$  and all frequencies reported in Hz scale linearly with  $4\pi/T$ . Throughout this manuscript, when reporting example frequencies in Hz we adopt the illustrative convention  $T = 4\pi$  s, yielding  $f_s \approx 5000/(4\pi) \approx 398$  Hz; other equally valid mappings exist depending on application. Consequently, any band-specific interpretation in Hz (e.g., “low/high” frequency ranges) should be understood under the chosen  $T$ ; changing  $T$  rescales all reported Hz values while preserving the underlying discrete sequences, which is a key feature of CoSiBD’s design.

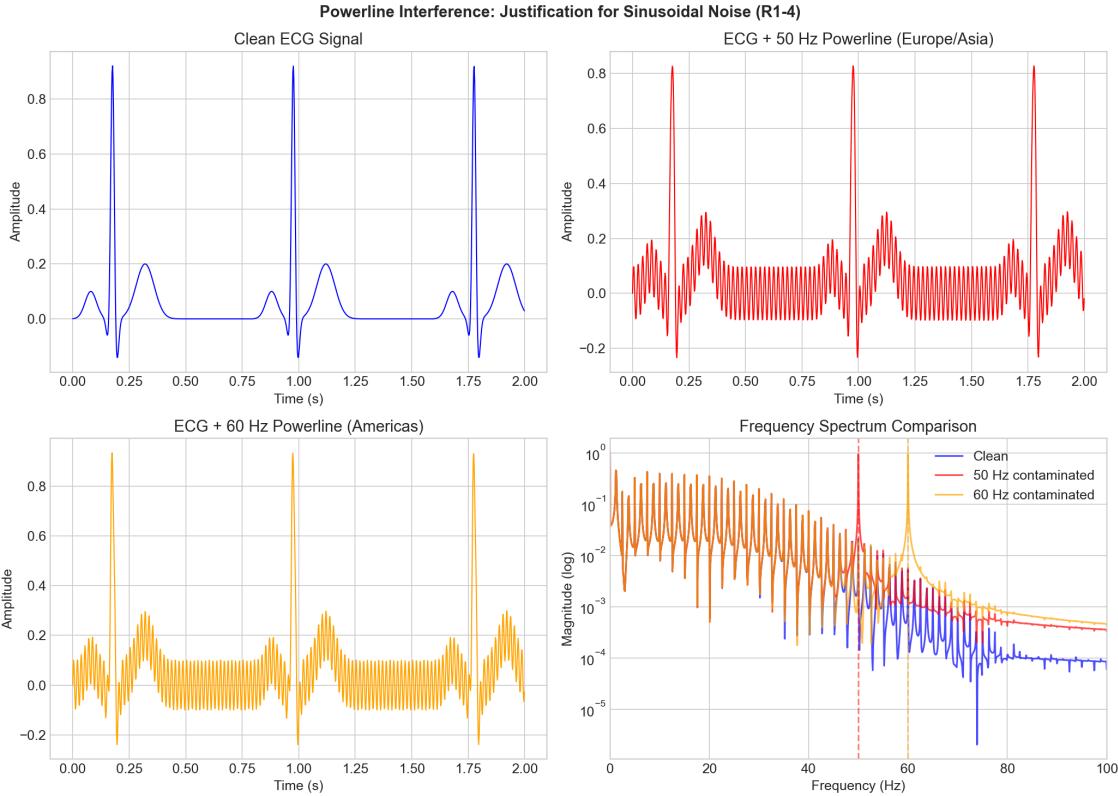
Figure 4 clarifies this convention using six panels (A–F). Panels A–C illustrate the time-domain representation: A shows the raw discrete sequence against sample index; B shows the same sequence mapped to a duration  $T_1$ ; and C mapped to  $T_2$ , demonstrating that the sample values are invariant. Panels D–F show the corresponding frequency-domain effects: D is the normalized spectrum (cycles/sample), while E and F show the Hz mappings for  $T_1$  and  $T_2$  respectively, illustrating how the physical frequency interpretation scales with the assumed duration.

## 151 Data Records

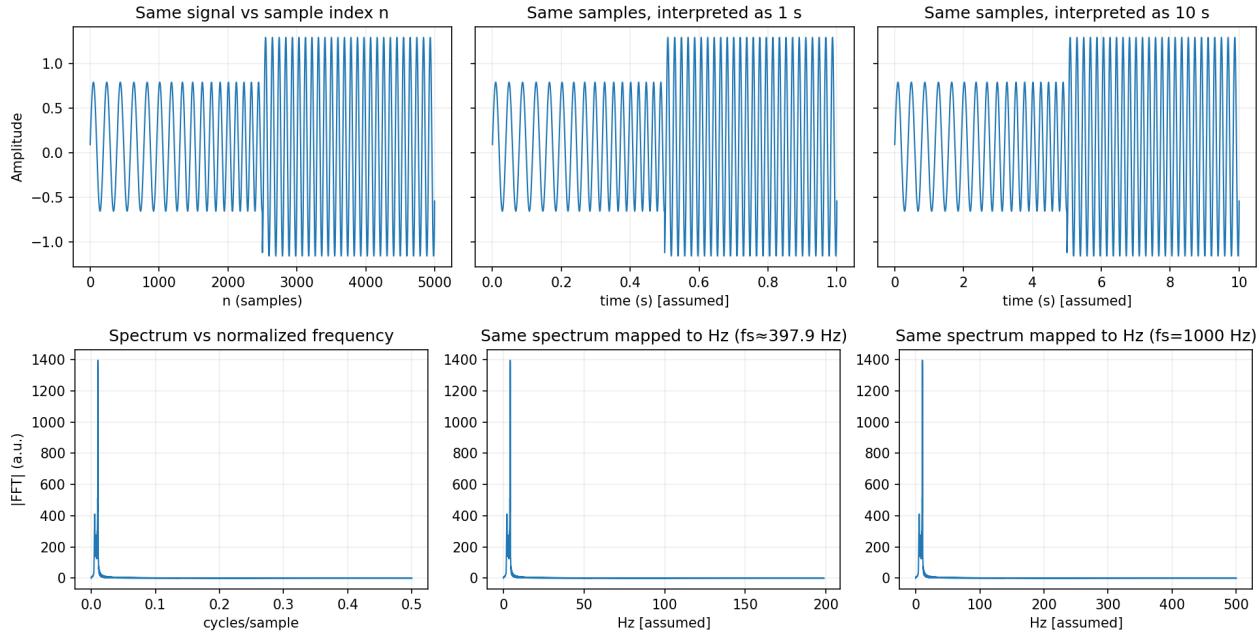
152 The Complex Signal Benchmark Dataset (CoSiBD) is publicly available on Zenodo<sup>22</sup> and consists of synthetic temporal signals, 153 mainly created to support the development and evaluation of temporal super-resolution (SR) algorithms, and also to train deep 154 learning models that can be used for SR in real-world signals. This section provides an overview of the dataset structure, 155 content, and storage format, as well as the parameters that rule the generation of data and the metadata that enrich our dataset.

156 The dataset comprises 2,500 high-resolution signals, each with corresponding subsampled versions at four resolution levels, 157 organized into three main categories:

- **High-resolution signals:** 2,500 signals with 5,000 samples each, spanning the domain  $T = [0, 4\pi]$  (s), which, under the 158 illustrative convention used, corresponds to  $f_s = 5000/(4\pi) \approx 398$  Hz. Each signal is stored in three formats: NumPy 159 compressed format (.npz), plain text (.txt), and JSON (.json). Per-signal metadata (frequency profiles with explicit 160 change-points (`base_points` and `high_freq_points`) and segment labels (`variation_type`), amplitude 161



**Figure 3.** Qualitative motivation for the structured interference term used in CoSiBD. An illustrative example shows how adding a narrow-band sinusoidal component (interpretable as 50/60 Hz under the illustrative convention  $T = 4\pi$  s) produces the characteristic periodic contamination observed in real recordings<sup>21</sup>, while broadband noise captures the measurement floor.



**Figure 4.** Sampling/unit convention in CoSiBD. Top: the same discrete sequence  $x[n]$  can be plotted against the sample index or under different assumed time scalings. Bottom: the intrinsic frequency axis is normalized (cycles/sample); mapping to “Hz” depends on the assumed sampling rate  $f_s$  (two example mappings shown).

envelopes, spline parameters, vertical offsets, noise configurations, and seeds) is provided in a consolidated JSON file (`signals_metadata.json`) with one entry per signal, enabling exact regeneration.

- **Simple subsampled signals:** Uniform decimation (uniform subsampling) of each signal to four target resolutions: 150 (illustrative  $f_s \approx 11.9$  Hz for  $T = 4\pi$  s), 250 (illustrative  $f_s \approx 19.9$  Hz for  $T = 4\pi$  s), 500 (illustrative  $f_s \approx 39.8$  Hz for  $T = 4\pi$  s), and 1,000 samples (illustrative  $f_s \approx 79.6$  Hz for  $T = 4\pi$  s). These low-resolution versions serve as inputs for SR benchmarking against the original 5,000-sample target. Stored in `.npz`, `.txt`, and `.json` formats.

The dataset is provided as consolidated files under `SignalBuilderC/data/`. High-resolution signals are stored as `signals_high_resolution_5000.[npz|txt|json]`. Simple subsampled (decimated) signals are stored as `signals_subsampled_simple_{150,250,500,1000}.[npz|txt|json]`. Dataset-level metadata (described later in detail) and configuration are stored in `signals_metadata.json` (per-signal metadata, one entry per signal), `signals_metadata Consolidated_2500.json`, and `dataset_summary.json`.

Regarding the three formats used for both high-resolution and subsampled signals, we provide here some additional information for each format: (1) NumPy compressed format (`.npz`) containing the signal array, time array, and (for high-resolution only) clean signal without noise; (2) consolidated plain text format (`.txt`) with one signal per row (samples separated by whitespace) for maximum portability; and (3) JSON format (`.json`) with both time and signal arrays for web-based applications and interoperability.

Reproducibility is ensured through documented seeds: each high-resolution signal is generated using a unique seed (ranging from 10,000 to 12,499), enabling exact regeneration of individual signals or the entire dataset. All generation parameters (described later in detail) are stored in metadata JSON files, including: (1) frequency profile parameters—`tau_frequency` values from uniform distribution [1, 2] with 0.05 step; (2) amplitude envelope parameters—`amplitude_spline_type` is predominantly `zero_order` (approximately 70%), otherwise `tension` with `tau_amplitude` drawn from {1, 3, 5, 8, 10, 12, 15, 20}; (3) vertical offsets—normally distributed (`mean=0`, `SD=3.0`); and (4) noise configurations—per-signal `noise_profile` with `p_has_noise=0.5` (and `p_gaussian=0.5` when noise is present).

## Metadata schema and example

CoSiBD provides per-signal metadata to support (i) deterministic regeneration, (ii) principled partitioning (e.g., by noise type/level or segment labels), and (iii) analysis of the piecewise structure induced by change-points. Table 2 summarizes representative fields contained in `signals_metadata.json`. A minimal example entry is shown below (one signal; values truncated for brevity).

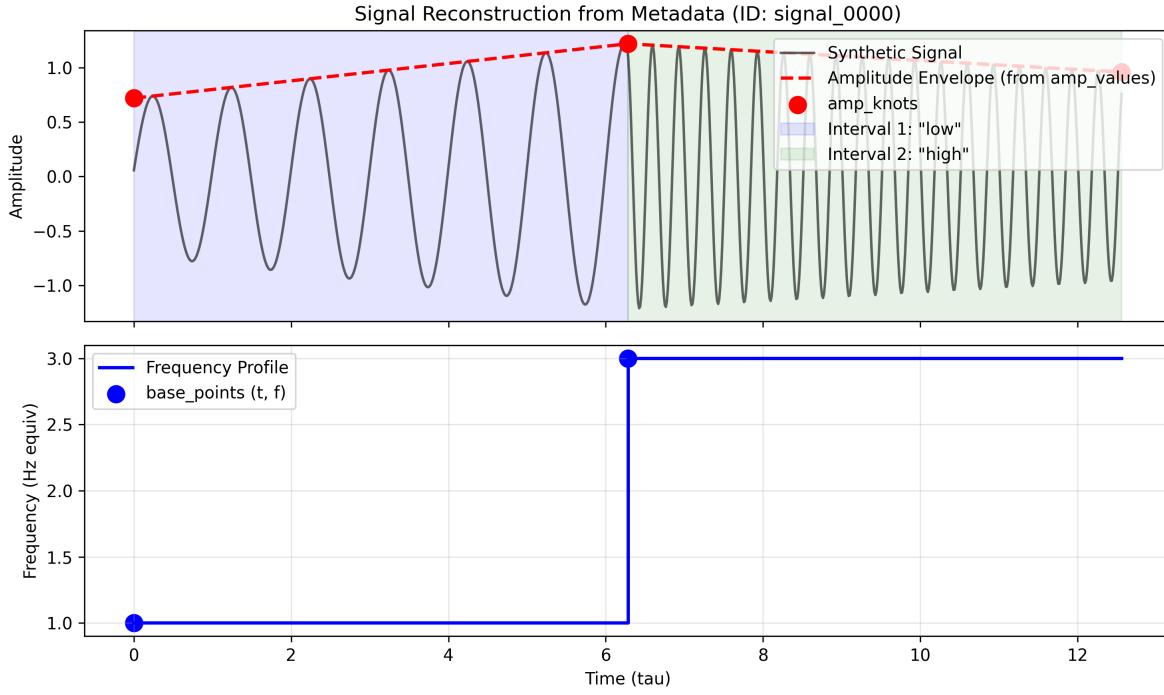
| Field                 | Type           | Example                 | Meaning  |
|-----------------------|----------------|-------------------------|--|
| signal_id             | String         | "signal_0000"           | Unique identifier for the generated signal.  |
| index                 | Integer        | 0                       | Numeric index of the signal (0 to N-1).  |
| seed                  | Integer        | 10000                   | Seed used for deterministic random <sup>REV</sup> generation of this specific signal.  |
| t_start, t_end        | Float          | 0.0, 12.56              | Start and end time of the signal domain.   |
| base_points           | Array          | [[0.0, 2.07]...]        | Control points $(t, f)$ for the frequency spline.  |
| high_freq_points      | Array          | [[0.0, 0.0]...]         | Control points $(t, f)$ for an optional high-frequency component.  |
| variation_type        | List[Str]      | ["low", "high", "high"] | Categorical labels aligned with base_points anchors (same length), describing the intended band at each change-point (e.g., "low", "high", "no_change"); interval labels are derived between successive anchors. |
| tau_frequency         | Float          | 1.15                    | Tension parameter for the frequency spline.  |
| amplitude_spline_type | String         | "zero_order"            | Interpolation family for the amplitude envelope (e.g., step-wise or tension spline).   |
| tau_amplitude         | Float / String | 10.0                    | Tension parameter for the amplitude spline (or "N/A" for zero-order).  |
| amp_knots             | Array          | [0.0, 6.28, 12.56]      | Knot times defining the amplitude envelope segments.   |
| amp_values            | Array          | [0.72, 1.22, 0.96]      | Amplitude values at knots (piecewise-constant or spline-interpolated).   |
| noise_profile         | Dict           | {...}                   | Parameters for added noise (if any).   |

**Table 2.** Metadata schema describing the JSON structure for each signal. Key fields define the frequency trajectory (base\_points, variation\_type) and amplitude envelope, allowing precise reconstruction or filtering.

190 Table 2 details the key fields describing each signal, specifically the variation\_type list which provides categorical  
 191 labels aligned with the base\_points anchors (e.g., "low", "high", or "no\_change"), enabling users to filter for signals  
 192 containing specific dynamic behaviors. The following JSON snippet shows a simplified illustrative metadata record (Example  
 193 1; values truncated for readability). Figure 5 provides a visual counterpart of the same record: it highlights how consecutive  
 194 base\_points (and, when present, high\_freq\_points) define temporal intervals and frequency targets, while  
 195 variation\_type provides the corresponding per-anchor labels (here, yielding two intervals: "low" then "high"), and  
 196 amp\_knots/amp\_values define the amplitude envelope.

```

197 {
198   "t_start": 0.0,
199   "t_end": 12.566370614359172,
200   "fs_high": 397.88735772973837,
201   "tau_frequency": 1.15,
202   "amplitude_spline_type": "zero_order",
203   "tau_amplitude": "N/A",
204   "vertical_offset": 0.06905161748158965,
205   "base_points": [[0.0, 1.0], [6.28, 3.0], [12.56, 3.0]],
206   "high_freq_points": [[0.0, 0.0], [6.28, 0.0], [12.56, 0.0]],
207   "variation_type": ["low", "high", "high"],
208   "amp_knots": [0.0, 6.28, 12.56],
209   "amp_values": [0.72, 1.22, 0.96],
210   "noise_profile": {"has_noise": true, "noise_type": "gaussian", "p_has_noise": 0.5, ...},
211   "seed": 10000,
212   "signal_id": "signal_0000",
213   "index": 0
214 }
```



**Figure 5.** Visual representation of an illustrative metadata record for signal\_0000 (Example 1). The plot illustrates how `base_points` define temporal intervals and frequency targets, while `amp_values` control the amplitude envelope segments, corresponding directly to the JSON structure.

## 215 Parameters for signal generation

216 As we anticipated before, our dataset is generated from a set of configuration parameters summarized in Table 3. Each  
 217 high-resolution signal was generated with a unique seed (10,000–12,499) and sampled parameter values within the defined  
 218 ranges, supporting diversity while maintaining reproducibility. Parameters marked as *Sampled (seed-controlled)* are drawn  
 219 per-signal using the recorded seed, while fixed defaults are held constant across the dataset. Unless otherwise stated, these  
 220 defaults correspond to the settings used in the Zenodo release, and the CLI follows the same behavior unless a user explicitly  
 221 overrides a parameter.

222 Concretely, generation proceeds as a deterministic pipeline conditioned on the seed: (i) the seed initializes all pseudo-random  
 223 draws; (ii) a set of frequency-trajectory and envelope parameters is sampled within the reported ranges; (iii) a continuous-time  
 224 signal is synthesized over the fixed domain ( $T = 4\pi$  s convention) by combining a time-varying frequency trajectory with  
 225 an amplitude envelope; and (iv) optional nuisance terms (offset and additive noise) are applied. Therefore, fixing the seed  
 226 fully determines the sampled parameters and the resulting signal instance, which supports exact reproduction and controlled  
 227 ablations.

228 **Frequency trajectory.** The overall (instantaneous) frequency evolution is controlled by the sampled low/high frequency bands  
 229 (“Low Frequency” and “High Frequency”), the number of transitions (“Change Points”), and their temporal placement (“Change  
 230 Locations”). The categorical “Variation Type” ({low, high, no\_change}) provides labels aligned with the `base_points`  
 231 anchors, indicating whether the trajectory targets the low band, the high band, or remains stable around each change-point;  
 232 this supports controlled balancing and downstream filtering via metadata. The smoothness/shape of the resulting trajectory is  
 233 governed by the spline settings, including `tau_frequency` (“Tension (freq)”).

234 **Amplitude envelope.** Independently of the frequency trajectory, the signal amplitude is modulated by an envelope whose  
 235 overall scale is constrained by “Amplitude Range”. The envelope shape is generated via the specified interpolation scheme  
 236 (“Spline Type”) and, when applicable, by the sampled “Tension (amp)” parameter. This yields signals with comparable  
 237 frequency content but different amplitude dynamics.

238 **Offsets and noise.** To emulate common acquisition artifacts without tying the generator to a specific instrument, we add a  
 239 per-signal “Vertical Offset” and inject additive noise stochastically according to “Noise Prob.”. When noise is enabled for a  
 240 given signal, its parameters are recorded in metadata to preserve reproducibility.

241 **Resolution and sampling.** Each instance is first generated at the target high-resolution length (e.g., 5,000 samples in the  
 242 released dataset), and lower-resolution views are derived by subsampling, enabling consistent multi-scale analysis and SR

243 benchmarking (Figure 6).

244 **Mathematical formulation.** Let  $t \in [t_{\text{start}}, t_{\text{end}}]$  with the paper's convention  $T = t_{\text{end}} - t_{\text{start}} = 4\pi$  s. A generated signal can be  
245 written as

$$x(t) = v + a(t) \sin(\phi(t)) + \varepsilon(t),$$

246 where  $v$  is a vertical offset,  $a(t) \geq 0$  is the amplitude envelope, and  $\varepsilon(t)$  is an optional additive noise term. The instantaneous  
247 frequency trajectory is encoded through the phase derivative

$$f(t) = \frac{1}{2\pi} \frac{d\phi(t)}{dt},$$

248 and is constructed from a spline defined by a set of control points and a tension parameter (Table 3). This separation (frequency  
249 trajectory vs. amplitude envelope) allows generating signals that share the same transition structure but differ in amplitude  
250 dynamics, or vice versa.

251 **Mapping Table 3 to metadata fields.** The parameters in Table 3 correspond to concrete entries in each per-signal metadata  
252 record (Figure 5 shows an illustrative example):

253 • *Seed* → `seed` (and derived identifiers `signal_id`, `index`). Fixing `seed` reproduces all sampled values.

254 • *Domain T = 4π convention* → `t_start`, `t_end`.

255 • *Change points / locations / variation type* → `base_points` (frequency control points), `variation_type` (interval  
256 labels), and auxiliary arrays such as `high_freq_points` used to represent high-frequency targets during transitions.

257 • *Frequency tension* → `tau_frequency`.

258 • *Amplitude range / spline settings* → `amplitude_spline_type`, `tau_amplitude`, and the envelope knot repre-  
259 sentation `amp_knots` and `amp_values`.

260 • *Vertical offset* → `vertical_offset`.

261 • *Noise probability and settings* → `noise_profile` (including whether noise was applied and the sampled noise  
262 parameters for that signal).

263 Some fields (e.g., `fs_high`) are reported as part of the complete record to make downstream processing explicit; they are  
264 derived from the chosen resolution and the domain convention.

| Parameter        | Range / Options                 | Default  | Description   |
|------------------|---------------------------------|--|---|
| Low Frequency    | 1–5 Hz                          | Sampled (seed-controlled)Random <sup>REV</sup> | Low-frequency component ( $T = 4\pi$ s convention).             |
| High Frequency   | 20–100 Hz                       | Sampled (seed-controlled)Random <sup>REV</sup> | High-frequency variations for transitions.                      |
| Change Points    | 2–11                            | Sampled (seed-controlled)Random <sup>REV</sup> | Number of frequency transitions per signal.                     |
| Change Locations | Continuous                      | Sampled (seed-controlled)Random <sup>REV</sup> | Time locations where transitions occur.                         |
| Variation Type   | {low, high, no_change}          | Balanced                                       | Category of frequency change per segment.                       |
| Amplitude Range  | 3–16                            | Sampled (seed-controlled)Random <sup>REV</sup> | Bounds for amplitude envelope generation.                       |
| Vertical Offset  | $N(0, 3.0)$                     | 0.0  | Normally distributed offset added to signals.                   |
| Spline Type      | Zero-Order (70%), Tension (30%) | -  | Interpolation method for envelopes.                             |
| Tension (freq)   | [1, 2]                          | 1.5  | Tension parameter for frequency splines.                        |
| Tension (amp)    | {1, 3, 5, 8, 10, 12, 15, 20}    | Sampled (seed-controlled)Random <sup>REV</sup> | Tension parameter for amplitude splines ("N/A" for zero-order). |
| Noise Prob.      | 0.0–1.0                         | 0.5  | Probability of adding noise to a signal.                        |
| Seed             | Integer                         | -  | Unique initialization seed per signal.                          |

**Table 3.** Configuration parameters used to generate the dataset. Ranges define the sampling space for the 2,500 signals, ensuring diversity. When not explicitly set by the user, parameters are either sampled per signal (seed-controlled) or take fixed defaults as listed (defaults used for the Zenodo release). Configuration parameters used to generate the dataset. Ranges define the sampling space for the 2,500 signals, ensuring diversity. Default values apply when a parameter is not randomized.<sup>REV</sup>

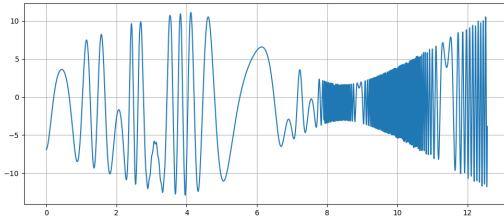
To explicitly characterize dataset diversity and complexity, CoSiBD spans multiple controlled axes of variation (Table 3), including the number and location of change points, categorical transition types, low/high frequency bands, and amplitude envelope configurations. The resulting variability is visible in representative realizations (Figures 6 and 7) and is quantified in Technical Validation via the distribution of dominant frequencies (Figure 9 and Table 4) and PSD behavior under different resolutions and noise settings (Figures 11 and 12). While the dataset is synthetic and not fitted to match a single domain-specific distribution, these controlled variations provide reproducible coverage of common real-world time-series phenomena such as non-stationarity, transient high-frequency events, and additive noise.

Figure 6 shows a representative signal from the dataset sampled at different resolution levels, as well as a version with added noise. This illustrates the variety of sampling and noise conditions included in CoSiBD.

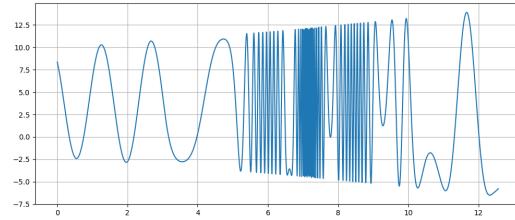
Figure 7 displays four additional synthetic signals generated using different configuration parameters. These examples demonstrate the variability in temporal structure across instances in the dataset.

## Custom Dataset Generation

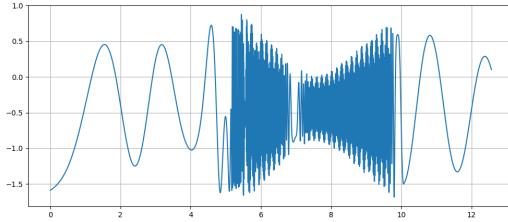
In addition to the pre-generated dataset, the CoSiBD package includes a command-line interface (CLI) that allows users to generate custom datasets with their own parameter distributions. Figure 8 demonstrates the usage of this tool. The command specifies the number of signals (`-n_signals`), the target high-resolution length in samples (`-resolution`, e.g., 5000), and the probability of applying noise to each signal (`-noise_prob`, e.g., 0.5, meaning noise is injected into roughly half of the signals). The output log confirms the resolved configuration (signal count, domain, and output directory), reports progress during generation, and summarizes the created artifacts (per-signal files and the consolidated `signals_metadata.json`).



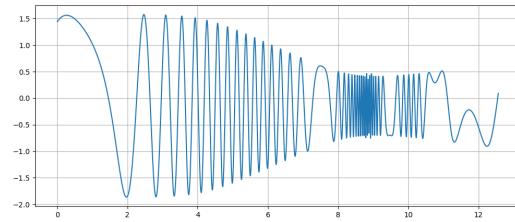
**(a)** High-resolution signal (5000 samples).



**(b)** Medium-resolution signal (500 samples).

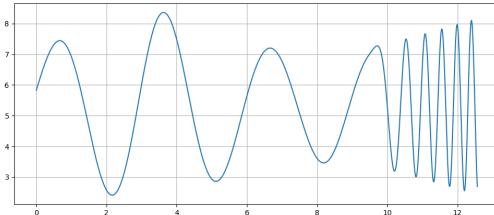


**(c)** Low-resolution signal (250 samples).

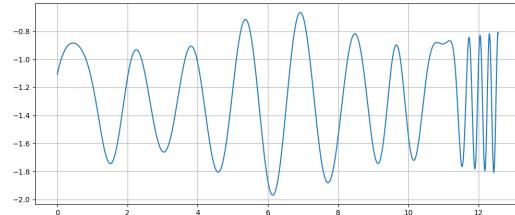


**(d)** Signal with added noise.

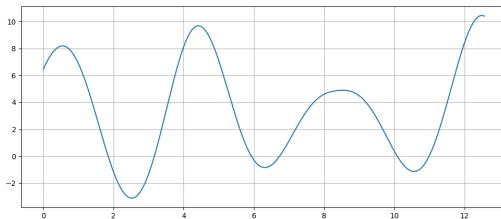
**Figure 6.** A synthetic signal sampled at different resolutions: (a) high (5000 samples), (b) medium (500 samples), (c) low (250 samples), and (d) with added noise. These examples reflect the multi-resolution and noise conditions present in the dataset.



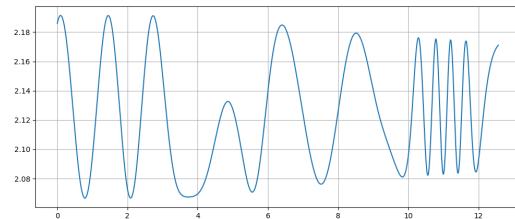
**(a)** Signal with increasing frequency over time.



**(b)** Signal with localized frequency variation.



**(c)** Signal with smooth oscillations and broad amplitude cycles.



**(d)** Signal with irregular peak spacing.

**Figure 7.** Examples of synthetic signals in the dataset generated with different parameter configurations. Each signal presents a distinct temporal profile.

```

$ python generate_dataset.py --n_signals 2500 --resolution 5000 --noise_prob 0.5
[INFO] Initializing Signal Generator...
[INFO] Configuration:
- Signals: 2500
- Resolution: 5000 samples
- Domain: [0, 4pi]
- Noise Probability: 0.5
[INFO] Output Directory: ./dataset_output/
Processing: 100%|██████████| 2500/2500 [00:45<00:00, 55.20it/s]

[SUCCESS] Dataset generation complete.
[INFO] Generated 2500 .npz files.
[INFO] Metadata saved to ./dataset_output/signals_metadata.json

$ ls -lh ./dataset_output/ | head -n 5
total 420M
-rw-r--r-- 1 user staff 120K Jan 10 10:00 signal_0000.npz
-rw-r--r-- 1 user staff 120K Jan 10 10:00 signal_0001.npz
-rw-r--r-- 1 user staff 120K Jan 10 10:00 signal_0002.npz

```

**Figure 8.** Demonstration of the CoSiBD Command Line Interface (CLI). Users can generate new dataset instances by specifying parameters such as signal count, sampling resolution, and noise probability directly from the terminal.

## 283 Technical Validation

284 This section evaluates the signal generation procedure by analyzing spectral properties under different conditions, including the  
 285 distribution of dominant frequencies, spectral stability across sampling rates, and the effect of noise. Additionally, we provide a  
 286 multi-scale super-resolution benchmark to demonstrate the dataset's utility in training deep learning models, and illustrative  
 287 transfer learning experiments using real-world EEG and speech data. These analyses aim to assess variability and stability  
 288 under the reported settings, and to document the dataset's behavior for reproducible use. Below, the methodologies and results  
 289 are described in detail.

### 290 Validation Context

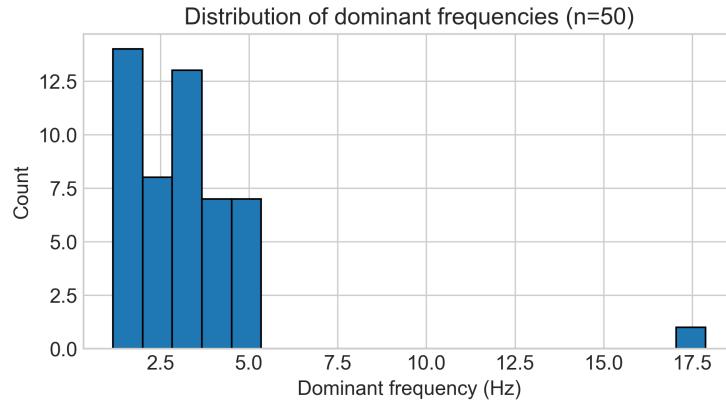
291 Experimental parameters were selected to support reproducibility and to illustrate representative behaviors of the generator  
 292 under the reported settings. The number of signals ( $n = 50$  for spectral analysis,  $n = 2500$  for benchmarks) provides a compact  
 293 but informative sample to summarize variability. Sampling resolutions (ranging from 150 to 5000 samples) reflect scenarios  
 294 requiring different levels of detail, aligning with typical signal processing use cases. Noise amplitudes (Gaussian noise with  
 295  $\sigma \in [0.0, 0.2]$ ) were motivated by common acquisition artifacts, with the goal of providing a controllable benchmark rather  
 296 than an exhaustive model of any specific measurement pipeline. Unless otherwise stated, signal-generation settings follow the  
 297 configuration in Table 3.

### 298 Analysis of Dominant Frequency Distribution

299 To assess the stability and variability of the primary spectral components, we analyzed the distribution of dominant frequencies  
 300 across multiple generated signals. A total of fifty independent signals were synthesized using identical input parameters. To  
 301 examine their spectral characteristics, we computed the power spectral density (PSD) of each signal, which quantifies how  
 302 signal power is distributed across different frequencies.

303 The PSD was estimated using Welch's method, selected for its ability to reduce noise and provide a smoother spectral  
 304 representation<sup>23</sup>. This method stabilizes spectral estimation by dividing the signal into overlapping segments, computing  
 305 their individual spectra, and averaging them. This reduces variance from random fluctuations and yields a smoother estimate.  
 306 For each signal, the dominant frequency was identified as the frequency at which the PSD reaches its maximum value. This  
 307 corresponds to the most prominent spectral component, indicating where the signal concentrates most of its energy.

309 By analyzing the distribution of dominant frequencies across the dataset, we evaluate whether the generated signals ex-  
 310 hibit consistent spectral patterns or if there is significant variation. High consistency would indicate stability in the data  
 311 generation process, whereas high variability could suggest the influence of stochastic sampling effects (seed-controlled)random  
 312 factors<sup>REV</sup>.



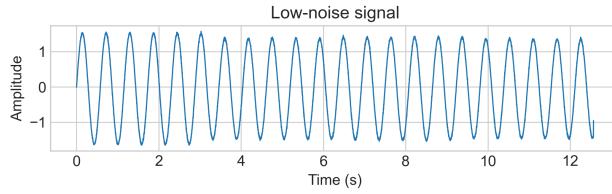
**Figure 9.** Distribution of dominant frequencies in 50 independently generated signals (reported in Hz, assuming the illustrative convention  $T = 4\pi$  s; for other time domains, the axis rescales by  $4\pi/T$ ).

| Statistic                  | Value (Hz; illustrative $T = 4\pi$ s) |
|----------------------------|---------------------------------------|
| Average Dominant Frequency | 0.508                                 |
| Standard Deviation         | 0.195                                 |
| Minimum Dominant Frequency | 0.390                                 |
| Maximum Dominant Frequency | 1.171                                 |

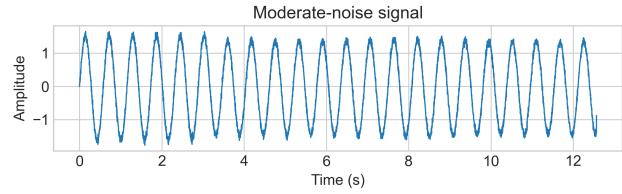
**Table 4.** Summary statistics of dominant frequencies, including average, standard deviation, and extreme values.

314 The results, shown in Figure 9 and Table 4, indicate that the dominant frequency (reported in Hz under the illustrative convention  
 315  $T = 4\pi$  s) is predominantly low under the reported settings. Quantitatively, the mean dominant frequency is 0.508 Hz with a  
 316 standard deviation of 0.195 Hz, while the observed range spans 0.390–1.171 Hz (Table 4). Figure 9 shows a strong concentration  
 317 in the low-frequency region, with a thinner tail toward higher values, corresponding to a smaller subset of realizations where  
 318 the strongest spectral peak shifts upward. This pattern is consistent with the stochastic sampling of frequency profiles (e.g.,  
 319 random control points and segment-wise variations): signals share the same parameter ranges, but some draws produce higher  
 320 instantaneous-frequency segments that become dominant in the PSD. Overall, the concentration around low frequencies  
 321 combined with controlled spread supports the goal of stable primary structure with deliberate diversity, which is desirable for  
 322 training models that must generalize across plausible spectral configurations<sup>24</sup>.

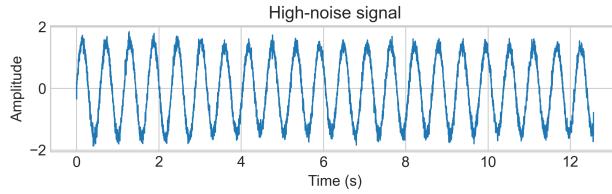
323 Figure 10 presents examples of signals from the CoSiBD dataset with increasing levels of added noise, illustrating how  
 324 amplitude fluctuations progressively obscure the underlying temporal structure.



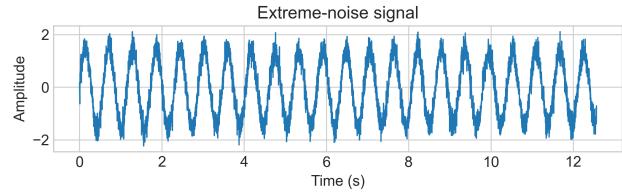
**(a)** Low-noise signal, where amplitude variations are present but minimally distorted.



**(b)** Moderate-noise signal, with irregular peaks and troughs beginning to distort the oscillatory pattern.



**(c)** High-noise signal, where significant distortion leads to unpredictable fluctuations.



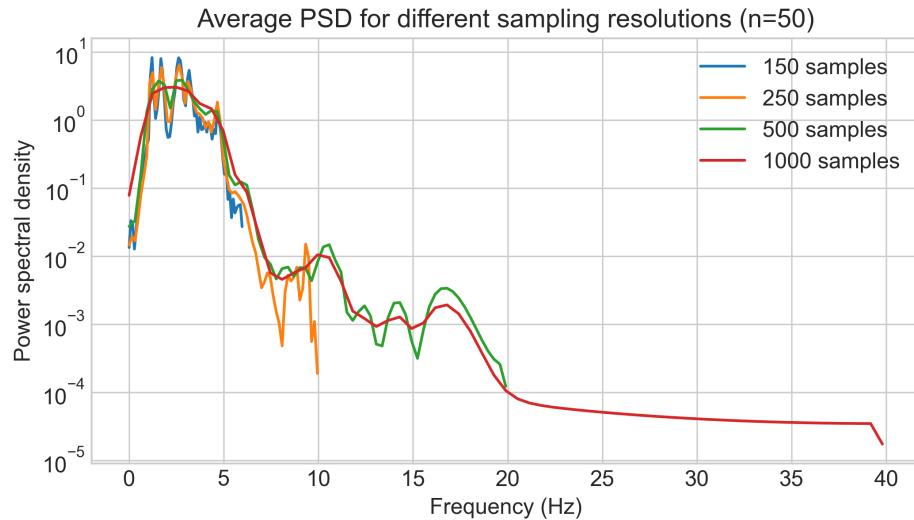
**(d)** Extreme-noise signal, where the original oscillatory structure is almost entirely masked by chaotic interference.

**Figure 10.** Visualization of signals under increasing noise conditions, showing how added noise progressively masks the original temporal patterns. From low (a) to extreme noise levels (d), this degradation highlights reconstruction challenges for super-resolution models.

### 325 Spectral Stability Across Sampling Resolutions

326 This analysis aims to investigate the influence of sampling resolution (number of samples) on the robustness of spectral estimates  
 327 under varying frequency content. When frequency axes are reported in Hz, they follow the illustrative convention  $T = 4\pi$  s; for  
 328 other choices of  $T$ , the Hz axis rescales by  $4\pi/T$ . At lower resolutions, reduced sampling density and coarser frequency grids  
 329 can obscure or merge spectral peaks, compromising the ability to distinguish closely spaced spectral components<sup>25</sup>. Conversely,  
 330 higher resolutions improve the granularity of the frequency axis, allowing for better separation of spectral features and reducing  
 331 the risk of misrepresenting the signal's underlying structure<sup>26</sup>.

332 This evaluation documents how spectral summaries vary with sampling resolution under the reported settings. The intent is to  
 333 provide descriptive context for using CoSiBD at different resolutions (and computational budgets) in benchmark protocols,  
 334 rather than to prescribe a universal sampling rate.



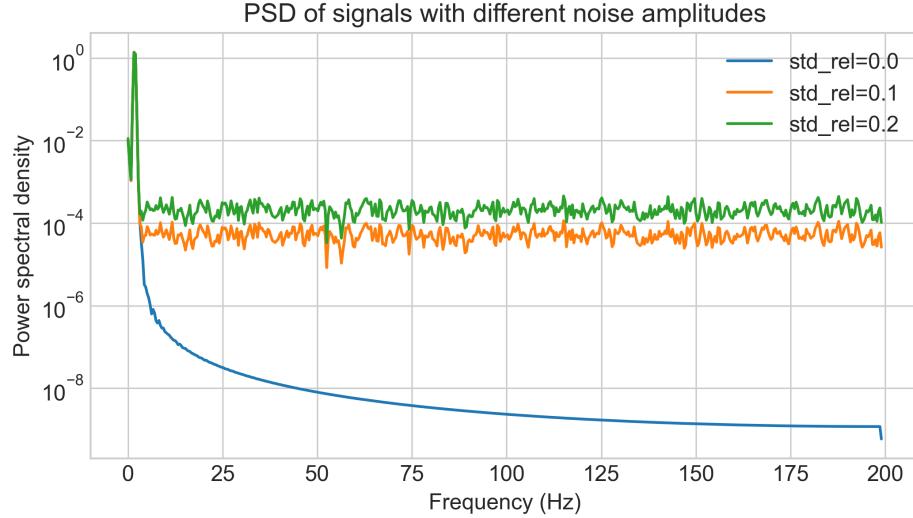
**Figure 11.** Average power spectral density (PSD) for different sampling resolutions based on 50 independent runs (Hz axis under the illustrative convention  $T = 4\pi$  s).

335 As shown in Figure 11, lower sampling resolutions, specifically the blue curve (150 samples) and the orange curve (250

samples), exhibit a noticeable reduction in detail within the higher-frequency range. These lower-resolution curves display greater fluctuations, consistent with the theoretical effects of subsampling and aliasing<sup>27</sup>. In contrast, the higher sampling resolutions (500, 1000 samples) demonstrate a smoother and more stable spectral profile. This analysis confirms that while lower sampling rates introduce aliasing artifacts, the dataset provides spectral fidelity comparable to theoretical expectations when sufficient resolution is employed.

### Impact of Noise on Frequency Characteristics

We analyze how varying the noise amplitude affects the power spectral density (PSD), with particular attention to differences between low- and high-frequency regions. Figure 12 illustrates this effect under the reported settings.



**Figure 12.** Power spectral density (PSD) of signals generated with different noise amplitudes (Hz axis under the illustrative convention  $T = 4\pi s$ ).

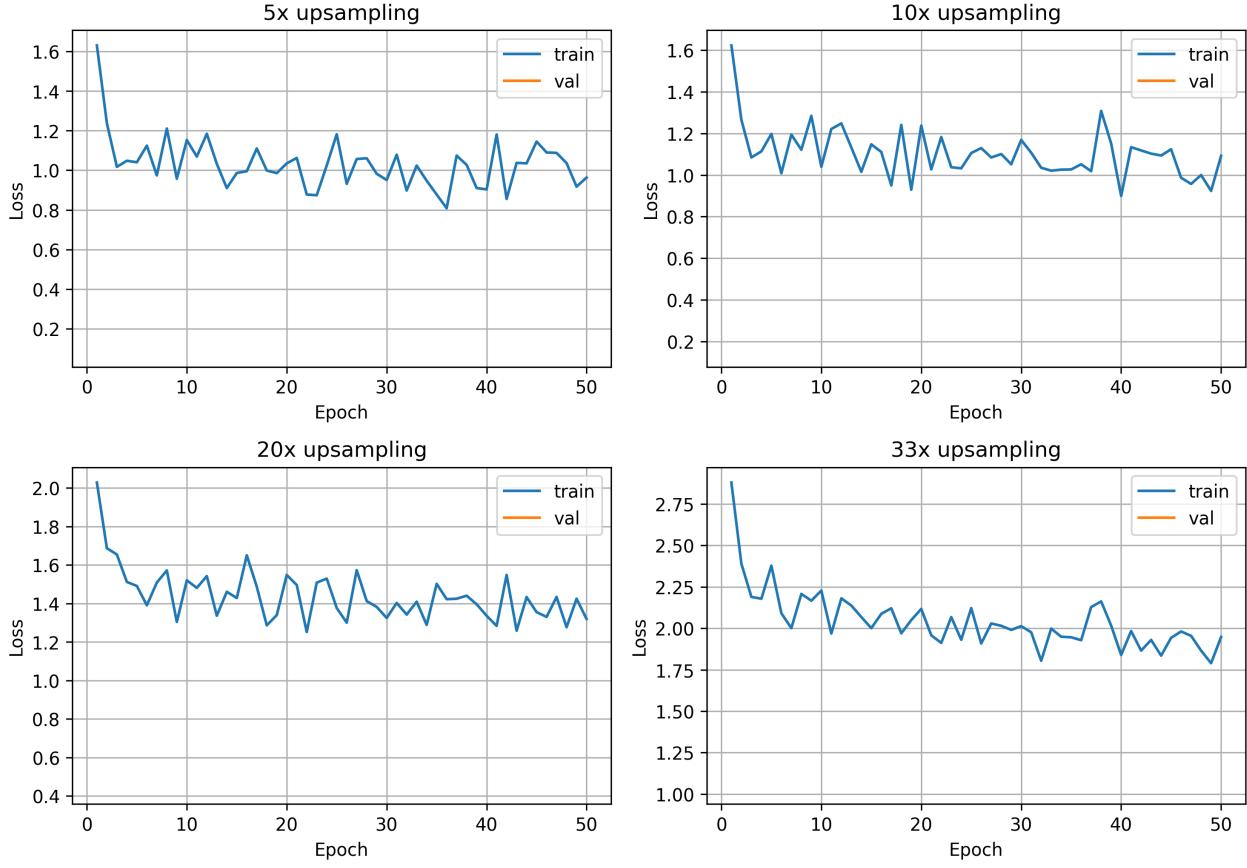
Across these settings, as the noise amplitude increases—from 0.0 (blue curve) to 0.2 (green curve)—the estimated PSD exhibits increased variability at higher frequencies, while the low-frequency region remains comparatively stable. This stability suggests that CoSiBD signals retain their primary structural characteristics even under significant noise, a critical property for robust representation learning<sup>28</sup>.

### Multi-Scale Super-Resolution Benchmark

To illustrate a baseline use case of CoSiBD, we trained a series of convolutional neural network (CNN) models for time series super-resolution at four different scaling factors:  $5 \times$  ( $1000 \rightarrow 5000$ ),  $10 \times$  ( $500 \rightarrow 5000$ ),  $20 \times$  ( $250 \rightarrow 5000$ ), and  $33 \times$  ( $150 \rightarrow 5000$ ). All models employed the TimeSeriesSRNet architecture—a five-layer encoder-decoder network with 1D convolutional layers and bilinear upsampling, inspired by deep residual architectures for audio generation<sup>9</sup>. For this benchmark, the 2,500 high-resolution signals were partitioned into an experiment-specific split of 2,000 paired signals for training and 500 held-out signals for validation.

We selected this architecture as a simple 1D convolutional encoder–decoder baseline: it captures local temporal structure while providing a deterministic upsampling mechanism, enabling consistent comparisons across scaling factors. Each model was trained using mean squared error (MSE) loss, a standard objective for regression tasks requiring broad mode coverage<sup>29</sup>, using the Adam optimizer (learning rate 0.001) and early stopping. A batch size of 16 was used as a practical compromise between optimization stability and MPS memory constraints. Training was conducted on Apple Silicon GPU (MPS backend) to accelerate convergence.

Table 5 summarizes the validation performance. In these runs, validation loss increased systematically with upsampling factor, reflecting the inherent difficulty of reconstructing fine temporal details from severely undersampled inputs (Table 5, Figure 13). Figure 13 illustrates the training and validation loss evolution for all four upsampling factors: curves track each other without pronounced divergence, and final validation loss increases monotonically with the factor, consistent with the expected rise in reconstruction difficulty.

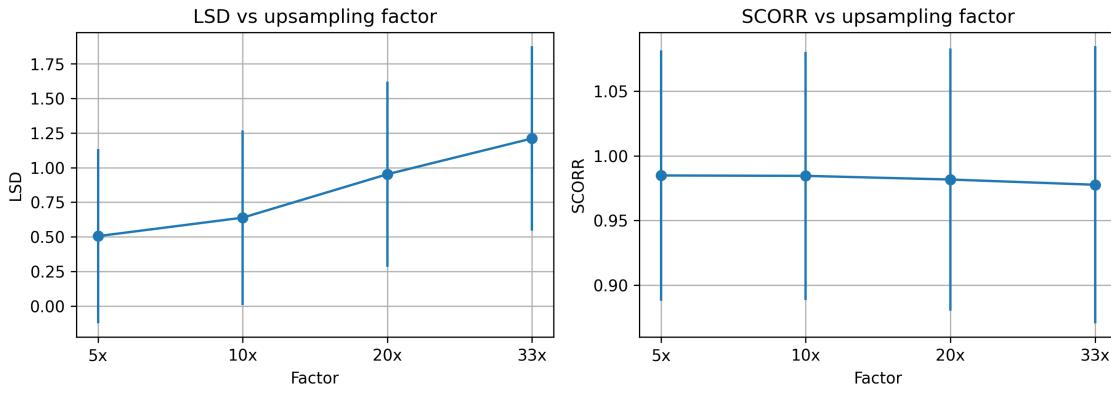


**Figure 13.** Training and validation loss evolution across all four upsampling factors ( $5\times$ ,  $10\times$ ,  $20\times$ ,  $33\times$ ). Each panel shows loss curves during training; in these runs, training and validation curves follow similar trends without pronounced divergence. The systematic increase in final validation loss with upsampling factor reflects the inherent difficulty of reconstructing fine temporal details from severely undersampled inputs.

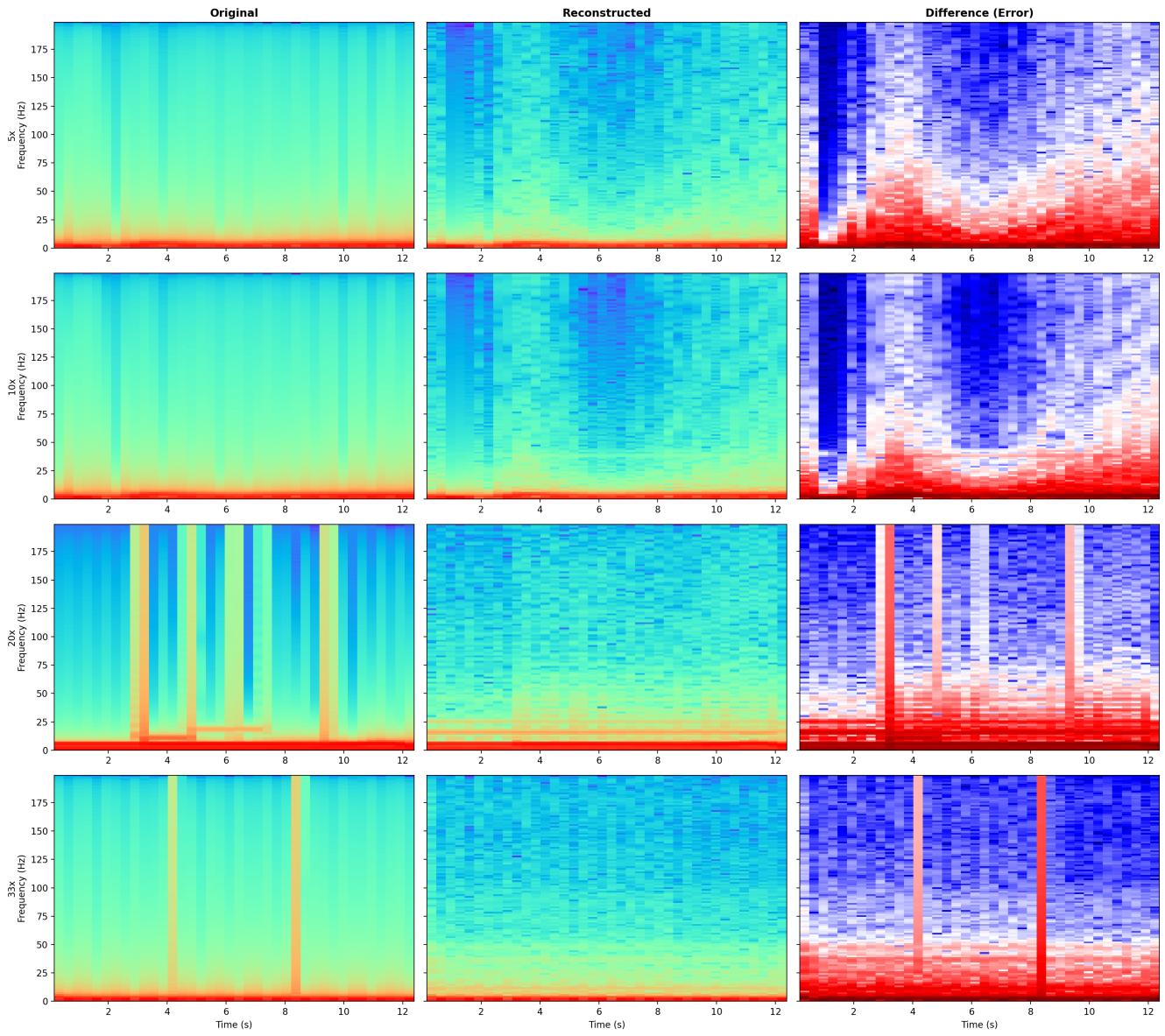
| Input Size   | Factor     | Val Loss | Epochs | Early Stop | LSD            | SCORR          |
|--------------|------------|----------|--------|------------|----------------|----------------|
| 1000 samples | $5\times$  | 0.0845   | 50     | No         | $0.51\pm 0.63$ | $0.98\pm 0.10$ |
| 500 samples  | $10\times$ | 0.1524   | 50     | No         | $0.64\pm 0.63$ | $0.98\pm 0.10$ |
| 250 samples  | $20\times$ | 0.4376   | 50     | No         | $0.95\pm 0.67$ | $0.98\pm 0.10$ |
| 150 samples  | $33\times$ | 1.0326   | 50     | No         | $1.21\pm 0.67$ | $0.98\pm 0.11$ |

**Table 5.** Multi-scale super-resolution benchmark results. Validation loss measured as mean squared error on 500 independent validation signals. LSD (Log Spectral Distance) quantifies spectral content deviation, while SCORR (Spectral Correlation) measures frequency-domain similarity. All models completed the full 50-epoch training without early termination, showing stable convergence.

To complement amplitude-based validation, we computed spectral fidelity metrics. Log Spectral Distance (LSD) increased from 0.51 ( $5\times$ ) to 1.21 ( $33\times$ ), indicating progressively larger deviations in spectral magnitude as reconstruction becomes more challenging, while Spectral Correlation (SCORR) remained consistently high (Table 5, Figure 14), suggesting that broad spectral structure is still largely preserved across factors<sup>26</sup>. Figure 15 provides representative spectrogram comparisons across all upsampling factors; the rightmost column (spectral difference) highlights where reconstruction introduces localized discrepancies (warm colors) that become more pronounced at higher factors, typically around rapid transitions and high-frequency content.

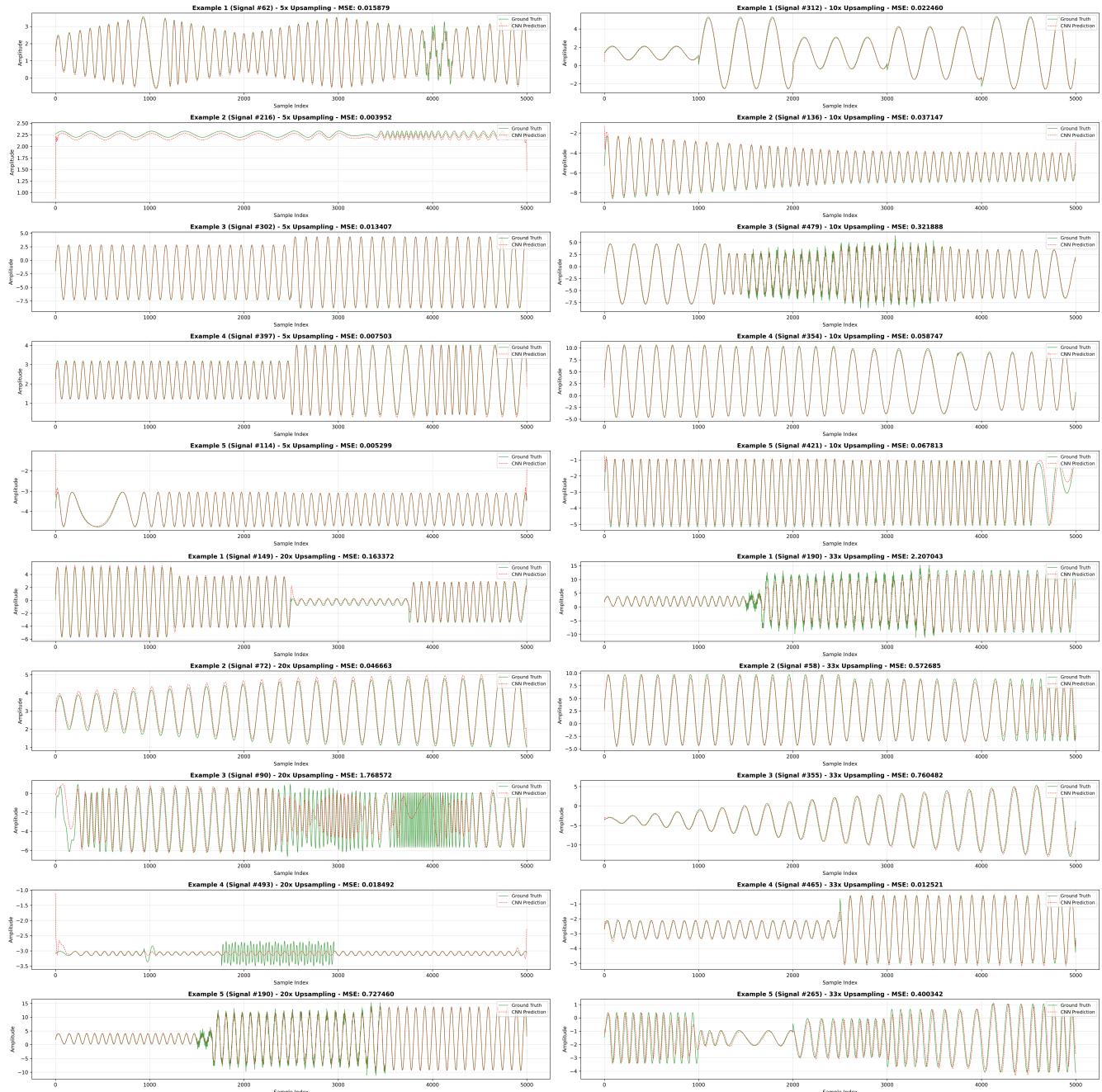


**Figure 14.** Spectral quality metrics vs upsampling factor. Left: Log Spectral Distance (LSD) increases systematically with upsampling factor, from 0.51 (5×) to 1.21 (33×). Right: Spectral Correlation (SCORR) maintains consistently high values (>0.97) across all factors. Error bars represent standard deviation over 500 validation signals per factor.



**Figure 15.** Spectrogram comparison across all upsampling factors. Each row represents a different upsampling factor ( $5\times$ ,  $10\times$ ,  $20\times$ ,  $33\times$ ), showing original signal (left), CNN-reconstructed signal (center), and spectral difference (right). Reconstruction artifacts become more visible at higher upsampling rates. Representative signals selected based on median Log Spectral Distance (LSD) for each factor.

Representative prediction examples (Figure 16) provide qualitative comparisons of reconstructed outputs against ground truth across scaling factors. As the upsampling factor increases, reconstructions tend to preserve the overall waveform trend while increasingly smoothing fast variations and sharp events, which is consistent with the information loss introduced by severe undersampling.



**Figure 16.** Representative prediction examples across all upsampling factors. Each quadrant shows prediction comparisons for a different scaling factor ( $5\times$ ,  $10\times$ ,  $20\times$ ,  $33\times$ ), displaying low-resolution inputs, ground-truth high-resolution signals, and CNN-reconstructed outputs.

These multi-scale experiments provide quantitative baseline results for future benchmarking studies. In addition to providing a reproducible baseline under a fixed architecture and protocol, the systematic increase in task difficulty—from moderate  $5\times$  upsampling to extreme  $33\times$  reconstruction—can serve as a practical reference for comparing architectures, loss functions, and training strategies in the time series super-resolution domain. More broadly, CoSiBD enables controlled benchmarking where methods can be compared on identical LR–HR pairs and nuisance settings, supporting fair ablation studies and future community protocols.

### 383 Illustrative Transfer Learning Experiments

384 To demonstrate the practical utility of CoSiBD for training deep learning models, we conducted transfer learning experiments  
 385 using convolutional neural networks (CNNs) for time-series super-resolution<sup>9,30</sup>. A TimeSeriesSRNet model (encoder-decoder  
 386 architecture with 1D convolutions: 1→64→128→256, bilinear upsampling, decoder 256→128→64→1) was evaluated on two  
 387 distinct real-world domains: EEG clinical signals<sup>1</sup> (500 training, 690 validation samples) and VCTK speech recordings<sup>20</sup> (44  
 388 hours from 109 speakers).

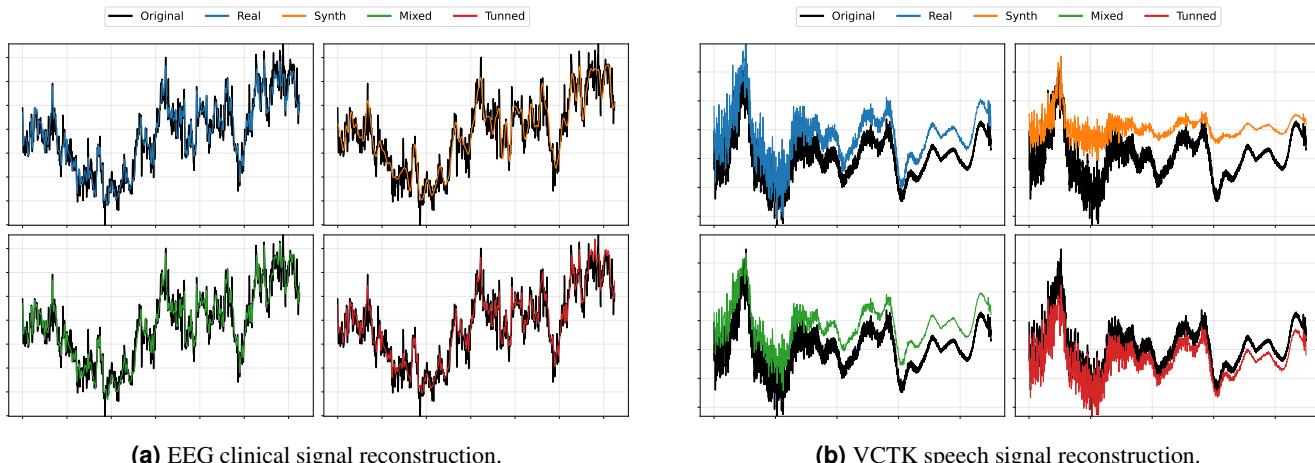
389 Four training strategies were systematically evaluated: (1) **Real-only**: trained exclusively on domain-specific real data  
 390 (baseline); (2) **Synth-only**: trained exclusively on CoSiBD synthetic signals; (3) **Mixed**: trained on a combined dataset of  
 391 synthetic and real data; and (4) **Tuned**: pre-trained on CoSiBD synthetic data, then fine-tuned on the real-world training set.  
 392 Performance was measured using Mean Absolute Error (MAE) between predicted and ground-truth high-resolution signals.

393 The results (Table 6, Figure 17) demonstrate that integrating CoSiBD improves model performance compared to using  
 394 limited real data alone. While models trained only on synthetic data (Synth-only) yielded higher errors due to domain shift, the  
 395 **Mixed** and **Tuned** strategies consistently outperformed the Real-only baseline. In particular, the best EEG result is an MAE of  
 396 **9.73** achieved by the **Mixed** strategy, while the best VCTK result is an MAE of **4.41** achieved by **Tuned** (Table 6). Specifically,  
 397 fine-tuning a CoSiBD-pretrained model reduced error significantly on the VCTK dataset, suggesting that the synthetic dataset  
 398 effectively captures universal temporal structures relevant to super-resolution tasks.

399 Beyond the aggregate MAE values, Figure 17 provides a qualitative interpretation of these results. In the EEG example,  
 400 Mixed/Tuned reconstructions better follow fast transients and reduce amplitude under/overshoot compared to the baseline. In  
 401 the speech example, CoSiBD-enhanced strategies recover sharper oscillatory detail and reduce residual smoothing, visually  
 402 aligning with the lower MAE reported in Table 6.

| Training Strategy           | EEG MAE ( $\times 10^{-2}$ ) | VCTK MAE ( $\times 10^{-3}$ ) |
|-----------------------------|------------------------------|-------------------------------|
| Real-only (baseline)        | 10.77                        | 5.92                          |
| Synth-only                  | 12.11                        | 8.79                          |
| Mixed (synth + real)        | <b>9.73</b>                  | 5.59                          |
| Tuned (pretrain + finetune) | 10.68                        | <b>4.41</b>                   |

**Table 6.** Mean Absolute Error (MAE) for CNN-based super-resolution models. Embedding CoSiBD data (Mixed and Tuned strategies) improves reconstruction accuracy compared to training on limited real data alone.



**Figure 17.** Visual comparison of super-resolution predictions. CoSiBD-enhanced models (Mixed/Tuned) recover finer details in both (a) EEG transients and (b) speech waveforms compared to the baseline. In this figure, the horizontal axis is sample index and the vertical axis is amplitude.

405 As a further validation of structural transferability, we applied the TimeSeriesSRNet trained solely on CoSiBD (5× upsampling)  
 406 to reconstruct full 2-second audio clips from the VCTK corpus without any fine-tuning. Despite the domain gap, the model

407 successfully recovered intelligible speech with preserved harmonic structure (mean Pearson correlation  $r = 0.928$ ). This  
408 indicates that CoSiBD's diverse frequency and envelope parameters generalize well to complex, non-stationary signals like  
409 human speech. Reconstructed examples (degraded vs. reconstructed signals under different strategies) are provided in the  
410 accompanying repository under `AudioModels/results/audio_samples/`.

## 411 Usage Notes

412 The CoSiBD dataset contains high-resolution signals and corresponding subsampled versions at multiple resolutions. Signals  
413 are provided in consolidated `.txt`, `.npz`, and `.json` formats. Pairing between low- and high-resolution versions is  
414 performed by row index: row  $i$  in a subsampled file corresponds to row  $i$  in the high-resolution file, with per-signal parameters  
415 available in `signals_metadata.json`. The dataset is distributed as a single, unified collection without a predefined  
416 train/validation/test split. Users should create partitions appropriate to their objectives (e.g., random splits, stratified splits by  
417 noise type/level or signal characteristics, cross-validation, or scenario-specific test sets), using the provided metadata to support  
418 principled partitioning.

### 419 Reading the Data

420 The signals are stored as consolidated plain text (`.txt`) files, with one signal per row (samples separated by whitespace). Each  
421 file contains multiple time series stacked vertically, where each row corresponds to a single signal. The dataset can be accessed  
422 using standard Python tools:

```
423 import numpy as np
424
425 # Load subsampled (simple decimation) and high-resolution signals
426 # Each .txt file is consolidated: one signal per row
427 x_valid = np.loadtxt('SignalBuilderC/data/signals_subsampled_simple_250.txt')
428 y_valid = np.loadtxt('SignalBuilderC/data/signals_high_resolution_5000.txt')
429
430 # Optional: convert to PyTorch tensors
431 # import torch
432 # x_valid = torch.tensor(x_valid, dtype=torch.float32)
433 # y_valid = torch.tensor(y_valid, dtype=torch.float32)
```

434 These commands return NumPy arrays (each row corresponds to one signal). Users can optionally convert them to PyTorch  
435 tensors.

### 436 Visualizing Signal Pairs

437 To inspect the alignment between low- and high-resolution versions, users can visualize paired signals indexed by the same  
438 row:[To explore the resolution differences, users can visualize aligned pairs of signals: REV](#)

```
439 import matplotlib.pyplot as plt
440 import numpy as np
441
442 # Load subsampled (simple decimation) and high-resolution signals
443 # Each .txt file is consolidated: one signal per row
444 x_lr = np.loadtxt('SignalBuilderC/data/signals_subsampled_simple_250.txt')
445 x_hr = np.loadtxt('SignalBuilderC/data/signals_high_resolution_5000.txt')
446
447 # Access a paired signal by row index
448 i = 0
449 low_res_signal = x_lr[i]
450 high_res_signal = x_hr[i]
451
452 # Visualize a paired low- and high-resolution signal
453 plt.figure(figsize=(10, 4))
454
455 # High-resolution signal
456 plt.plot(high_res_signal, label='High-resolution (5000 samples)', alpha=0.8)
```

```

457
458 # Low-resolution signal (aligned to HR index range for visualization)
459 lr_x = np.linspace(0, len(high_res_signal), len(low_res_signal))
460 plt.scatter(lr_x, low_res_signal, color='red', s=12,
461             label='Low-resolution (250 samples)')
462
463 plt.xlabel('Sample index')
464 plt.ylabel('Amplitude')
465 plt.title('Paired Low- and High-Resolution Signal')
466 plt.legend()
467 plt.grid(True)
468 plt.tight_layout()
469 plt.show()

```

470 This visualization highlights how the same underlying temporal structure is represented at different resolutions while  
471 preserving alignment between paired signals. Additional signal characteristics (e.g., change-points, frequency profiles, or noise  
472 configuration) can be retrieved from `signals_metadata.json` using the same row index.<sup>REV</sup>

### 473 Training a baseline model (synthetic-only)

474 The following example illustrates a minimal synthetic-only training loop for time-series super-resolution using CoSiBD pairs  
475 (LR input from simple uniform decimation, HR target). The intent is to provide a compact, reproducible starting point; full  
476 training scripts and additional configurations are available in the accompanying repository.

```

477 import numpy as np
478 import torch
479 import torch.nn as nn
480 from torch.utils.data import DataLoader, TensorDataset
481
482 # Load paired signals (rows align by index)
483 x = np.loadtxt('SignalBuilderC/data/signals_subsampled_simple_250.txt')    # (2500, 250)
484 y = np.loadtxt('SignalBuilderC/data/signals_high_resolution_5000.txt')      # (2500, 5000)
485
486 # Train/val split (example protocol)
487 x_train, y_train = x[:2000], y[:2000]
488 x_val, y_val = x[2000:2500], y[2000:2500]
489
490 device = torch.device('mps' if torch.backends.mps.is_available() else 'cpu')
491
492 def to_tensor(a):
493     # Convert NumPy array (B, L) into torch tensor (B, 1, L)
494     # The extra channel dimension matches Conv1d input format
495     return torch.tensor(a, dtype=torch.float32).unsqueeze(1)    # (B, 1, L)
496
497 # Create dataloaders for batched training
498 train_loader = DataLoader(TensorDataset(to_tensor(x_train), to_tensor(y_train)),
499                           batch_size=16, shuffle=True)
500 val_loader = DataLoader(TensorDataset(to_tensor(x_val), to_tensor(y_val)),
501                         batch_size=16, shuffle=False)
502
503 class TinySRNet(nn.Module):
504     def __init__(self, out_len=5000):
505         super().__init__()
506         # Encoder: extract local features in the LR domain
507         self.enc = nn.Sequential(
508             nn.Conv1d(1, 64, kernel_size=5, padding=2), nn.ReLU(),
509             nn.Conv1d(64, 128, kernel_size=5, padding=2), nn.ReLU(),

```

```

510         nn.Conv1d(128, 256, kernel_size=5, padding=2), nn.ReLU(),
511     )
512     # Upsampling: map LR features to the HR length
513     self.up = nn.Upsample(size=out_len, mode='linear', align_corners=False)
514     # Decoder: project features back to a 1-channel HR signal
515     self.dec = nn.Sequential(
516         nn.Conv1d(256, 128, kernel_size=5, padding=2), nn.ReLU(),
517         nn.Conv1d(128, 64, kernel_size=5, padding=2), nn.ReLU(),
518         nn.Conv1d(64, 1, kernel_size=5, padding=2),
519     )
520
521     def forward(self, x):
522         z = self.enc(x)
523         z = self.up(z)
524         return self.dec(z)
525
526 model = TinySRNet(out_len=5000).to(device)
527 # Optimizer and loss (MSE is a standard regression objective)
528 opt = torch.optim.Adam(model.parameters(), lr=1e-3, weight_decay=1e-5)
529 loss_fn = nn.MSELoss()
530
531 for epoch in range(1, 11):
532     model.train()
533     for xb, yb in train_loader:
534         xb, yb = xb.to(device), yb.to(device)
535         # Forward + loss + backward + update
536         opt.zero_grad()
537         pred = model(xb)
538         loss = loss_fn(pred, yb)
539         loss.backward()
540         opt.step()
541
542     model.eval()
543     with torch.no_grad():
544         # Validation loop: average loss over batches
545         val_loss = 0.0
546         for xb, yb in val_loader:
547             xb, yb = xb.to(device), yb.to(device)
548             val_loss += loss_fn(model(xb), yb).item()
549         val_loss /= len(val_loader)
550         print(f"epoch={epoch:02d} val_mse={val_loss:.4f}")

```

## 551 Code availability

552 The complete signal generation pipeline, including modules for frequency profile generation, amplitude envelope construction,  
553 spline interpolation, noise application, and data export in multiple formats, is available at: [CoSiBD scripts on GitHub](#). The  
554 repository includes SignalBuilderC, a modular Python package with documented functions for: (1) generating high-resolution  
555 signals with configurable parameters, (2) creating subsampled versions via simple decimation (uniform subsampling), (3)  
556 exporting signals in NumPy, text, and JSON formats, and (4) comprehensive metadata generation. All code is provided with  
557 example notebooks demonstrating dataset regeneration and usage. These scripts are distributed under the MIT License.

558  
559 The dataset itself is published separately at: Zenodo<sup>22</sup> (DOI: [10.5281/zenodo.15138853](https://doi.org/10.5281/zenodo.15138853)). The Zenodo record distributes the  
560 dataset under the Creative Commons Attribution 4.0 International (CC BY 4.0) license.

## 561 References

- 562 1. Luciw, M. D., Jarocka, E. & Edin, B. B. Multi-channel eeg recordings during 3,936 grasp and lift trials with varying  
563 weight and friction. *Sci. Data* **1**, 140047, [10.1038/sdata.2014.47](https://doi.org/10.1038/sdata.2014.47) (2014).
- 564 2. Nayak, S. K. *et al.* A review of methods and applications for a heart rate variability analysis. *Algorithms* **16**, 433,  
565 [10.3390/a16090433](https://doi.org/10.3390/a16090433) (2023).
- 566 3. Shaffer, F. & Ginsberg, J. P. An overview of heart rate variability metrics and norms. *Front. Public Heal.* **5**, 258,  
567 [10.3389/fpubh.2017.00258](https://doi.org/10.3389/fpubh.2017.00258) (2017).
- 568 4. Chen, S.-W. Non-uniform sampling data converters: A journey to uncharted circuits and systems. In *2022 International  
569 Symposium on VLSI Design, Automation and Test (VLSI-DAT)*, 1–1, [10.1109/VLSI-DAT54769.2022.9768053](https://doi.org/10.1109/VLSI-DAT54769.2022.9768053) (2022).
- 570 5. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444, [10.1038/nature14539](https://doi.org/10.1038/nature14539) (2015).
- 571 6. Goodfellow, I. J. *et al.* Generative adversarial networks. *arXiv preprint arXiv:1406.2661* [10.48550/arXiv.1406.2661](https://doi.org/10.48550/arXiv.1406.2661)  
572 (2014).
- 573 7. Isasa, I. *et al.* Comparative assessment of synthetic time series generation approaches in healthcare: leveraging patient  
574 metadata for accurate data synthesis. *BMC Med. Informatics Decis. Mak.* **24**, Article 27 (2024).
- 575 8. Zhang, C., Bengio, S., Hardt, M., Recht, B. & Vinyals, O. Understanding deep learning requires rethinking generalization.  
576 *arXiv preprint arXiv:1611.03530* [10.48550/arXiv.1611.03530](https://doi.org/10.48550/arXiv.1611.03530) (2016).
- 577 9. Kuleshov, V., Enam, S. Z. & Ermon, S. Audio super resolution using neural networks. *arXiv preprint arXiv:1708.00853*  
578 [10.48550/arXiv.1708.00853](https://doi.org/10.48550/arXiv.1708.00853) (2017).
- 579 10. Ibarra-Fiallo, J. & Lara, J. A. Contextual deep learning approaches for time series reconstruction. In *2024 IEEE  
580 International Conference on Omni-Layer Intelligent Systems, COINS 2024* (Institute of Electrical and Electronics Engineers  
581 Inc., London, United Kingdom, 2024).
- 582 11. Brophy, E., Wang, Z., She, Q. & Ward, T. Generative adversarial networks in time series: A systematic literature review.  
583 *ACM Comput. Surv.* **55**, Article 199, [10.1145/3559540](https://doi.org/10.1145/3559540) (2023).
- 584 12. McSharry, P. E., Clifford, G. D., Tarassenko, L. & Smith, L. A. A dynamical model for generating synthetic electrocardio-  
585 gram signals. *IEEE Transactions on Biomed. Eng.* **50**, 289–294, [10.1109/TBME.2003.808805](https://doi.org/10.1109/TBME.2003.808805) (2003).
- 586 13. O’Shea, T. J. & West, N. Radio machine learning dataset generation with GNU radio. In *Proceedings of the GNU Radio  
587 Conference*, vol. 1 (2016).
- 588 14. DeepSig. Datasets (including radioml 2016.10a). <https://www.deepsig.ai/datasets/>. Accessed 2026-01-13.
- 589 15. DeepSig. Radioml 2018.01a dataset. <https://www.deepsig.ai/datasets/>. Accessed 2026-01-13.
- 590 16. McSharry, P. & Clifford, G. D. ECGSYN: A realistic ecg waveform generator (physionet). <https://physionet.org/physiotools/ecgsyn/>. Accessed 2026-01-13.
- 592 17. Krol, L. R., Pawlitzki, J., Lotte, F., Gramann, K. & Zander, T. O. Sereega: Simulating event-related eeg activity. *J.  
593 Neurosci. Methods* **309**, 13–24, [10.1016/j.jneumeth.2018.08.001](https://doi.org/10.1016/j.jneumeth.2018.08.001) (2018).
- 594 18. Pinceti, A., Sankar, L. & Kosut, O. Generation of synthetic multi-resolution time series load data. *arXiv:2107.03547*  
595 (2021).
- 596 19. Yuan, Z., Jiang, Y., An, Z., Ma, W. & Wang, Y. Seismic resolution improving by a sequential convolutional neural network.  
597 *PLOS ONE* **19**, e0304981, [10.1371/journal.pone.0304981](https://doi.org/10.1371/journal.pone.0304981) (2024).
- 598 20. Yamagishi, J., Veaux, C. & MacDonald, K. CSTR VCTK Corpus: English multi-speaker corpus for CSTR voice cloning  
599 toolkit (version 0.92), [10.7488/ds/2645](https://doi.org/10.7488/ds/2645) (2019).
- 600 21. Sörnmo, L. & Laguna, P. *Bioelectrical Signal Processing in Cardiac and Neurological Applications* (Elsevier Academic  
601 Press, 2005).
- 602 22. Ibarra-Fiallo, J., Lara, J. A. & Agudelo Moreno, D. Cosibd, [10.5281/zenodo.15138853](https://doi.org/10.5281/zenodo.15138853) (2025). Version v1. Dataset.
- 603 23. Welch, P. D. The use of fast fourier transform for the estimation of power spectra: A method based on time averaging  
604 over short, modified periodograms. *IEEE Transactions on Audio Electroacoustics* **15**, 70–73, [10.1109/TAU.1967.1161901](https://doi.org/10.1109/TAU.1967.1161901)  
605 (1967).
- 606 24. Bengio, Y., Courville, A. & Vincent, P. Representation learning: A review and new perspectives. *IEEE Transactions on  
607 Pattern Analysis Mach. Intell.* **35**, 1798–1828, [10.1109/TPAMI.2013.50](https://doi.org/10.1109/TPAMI.2013.50) (2013).

- 608 25. Rabiner, L. R. & Gold, B. *Theory and Application of Digital Signal Processing* (Prentice Hall, 1975).
- 609 26. Marple, S. L., Jr. *Digital Spectral Analysis with Applications* (Prentice Hall, 1987).
- 610 27. Shannon, C. E. & Weaver, W. *The Mathematical Theory of Communication* (University of Illinois Press, 1949).
- 611 28. Bishop, C. M. *Pattern Recognition and Machine Learning* (Springer, 2006).
- 612 29. Goodfellow, I., Bengio, Y. & Courville, A. *Deep Learning* (MIT Press, 2016).
- 613 30. Kaniraja, C. P., M, V. D. & Mishra, D. A deep learning framework for electrocardiogram (ecg) super resolution and arrhythmia classification. *Res. on Biomed. Eng.* **40**, 199–211, [10.1007/s42600-023-00320-x](https://doi.org/10.1007/s42600-023-00320-x) (2024).

615 **Acknowledgments**

616 This research was supported by Dean's Office of the Polytechnic College of the San Francisco de Quito University and partially  
617 by ProyExcel-0069 project of the Andalusian University, Research and Innovation Department.

618 **Author Contributions**

619 J. I. F. handled the methodological design for artificial data creation, probabilistic analysis, spline-based variations, noise  
620 distributions, and *noderandom-node<sup>REV</sup>* selection. J. A. L. was responsible for methodology (time series design) and supervision.  
621 D. A. M. performed data processing and validation analysis. All of the authors have contributed to writing the manuscript.

622 **Competing Interests**

623 The authors declare no competing interests