

## Overview: A Two-Stage Augmentation Pipeline

The core idea is to create a data-loading pipeline that, for each clean ECG signal from the public dataset, applies a series of random transformations before feeding it to your model. This process should mimic the two primary sources of error you identified: lead displacement and noise.

The key is to perform this **"on-the-fly"** during training. Instead of creating a massive, static augmented dataset that lives on your hard drive, you'll generate unique augmented samples for each training epoch. This exposes your model to a much wider variety of realistic variations.

### 1. Simulating Lead Displacement

A slight physical shift of an ECG electrode means the signal it records will be a mixture of the potential at its ideal location and the potentials of the surrounding areas. We can simulate this by creating new, **augmented leads that are linear combinations of the original, physically adjacent leads.**

#### Method: Lead Mixing

For each lead, you'll create a new signal by blending it with its neighbors, controlled by small, random coefficients.

Limb Leads (Frontal Plane):

The limb leads (I, II, III, aVR, aVL, aVF) are mathematically related through Einthoven's triangle and the Goldberger terminals. Augmenting any two of the primary leads (I and II) will naturally augment the others.

- $II_{aug} = (1 - \alpha) * II + \alpha * I$
- $I_{aug} = (1 - \beta) * I + \beta * II$
- Then, recalculate the others based on these new, slightly shifted signals:
  - $III_{aug} = II_{aug} - I_{aug}$
  - $aVR_{aug} = -(I_{aug} + II_{aug}) / 2$
  - $aVL_{aug} = I_{aug} - II_{aug} / 2$
  - $aVF_{aug} = II_{aug} - I_{aug} / 2$

Here,  $\alpha$  and  $\beta$  are small random numbers, for example, drawn from a uniform distribution  $U(-0.1, 0.1)$ . This simulates a slight rotation/shift in the frontal plane axis.

Precordial Leads (Transverse Plane):

The chest leads (V1-V6) are arranged physically across the chest. A shift in one lead will cause it to pick up signals from its immediate neighbors.

Lead to Augment	Physically Adjacent Leads	Augmentation Formula
V1	V2	$V1_{aug} = (1 - \alpha) * V1 + \alpha * V2$
V2	V1, V3	$V2_{aug} = (1 - \alpha - \beta) * V2 + \alpha * V1 + \beta * V3$
V3	V2, V4	$V3_{aug} = (1 - \alpha - \beta) * V3 + \alpha * V2 + \beta * V4$
V4	V3, V5	$V4_{aug} = (1 - \alpha - \beta) * V4 + \alpha * V3 + \beta * V5$
V5	V4, V6	$V5_{aug} = (1 - \alpha - \beta) * V5 + \alpha * V4 + \beta * V6$
V6	V5	$V6_{aug} = (1 - \alpha) * V6 + \alpha * V5$

- For each step,  $\alpha$  and  $\beta$  should be new, small random numbers (e.g., from  $U(0, 0.15)$ ), representing a pull towards one neighbor or the other.
- The key is that the coefficients for a given lead sum to 1 to preserve the overall signal power ( $(1 - \alpha - \beta) + \alpha + \beta = 1$ ).

## 2. Simulating Noise

After simulating the lead displacement, you should add various types of noise that are characteristic of wearable, non-gelled sensors. You can apply these in combination.

- **Baseline Wander:** This low-frequency drift is often caused by breathing.
  - **Simulation:** Generate a very low-frequency sine wave (or a combination of a few) and add it to the signal.
  - $ECG_{aug} = ECG_{displaced} + A * \sin(2 * \pi * f * t + \phi)$
  - Where amplitude  $A$  is a percentage of the signal's peak-to-peak voltage, frequency  $f$  is low (e.g., 0.05 Hz to 0.5 Hz for breathing), and phase  $\phi$  is random.
- **Powerline Interference:** Noise from electrical mains (50 or 60 Hz).
  - **Simulation:** Add a sine wave at the powerline frequency. Its amplitude should be small and randomized for each sample.
  - $ECG_{aug} = ECG_{aug} + A_{pl} * \sin(2 * \pi * f_{pl} * t + \phi)$
  - Where  $f_{pl}$  is 50 or 60 Hz.
- **Muscle Artifacts (EMG):** Higher-frequency noise from muscle contractions.
  - **Simulation:** The simplest method is to add Gaussian noise with a randomized standard deviation.

- $ECG\_aug = ECG\_aug + N(0, \sigma)$
- A more advanced method is to download actual EMG signal samples from a database (like PhysioNet) and add random segments of this real noise to your ECG data.
- **Motion Artifacts:** Sudden, sharp spikes or baseline shifts from body or shirt movement.
  - **Simulation:** Randomly add step functions or sharp Gaussian spikes to the signal at random locations. You can also randomly select a segment of the signal and abruptly shift its baseline up or down.

### 3. Dataset Size and Validation Strategy

**Answering your question directly: Yes, your augmented dataset will be conceptually massive.**

If you have 100,000 ECGs and you implement this "on-the-fly" augmentation pipeline, your model will see a slightly different, unique version of each ECG every time it trains on it. Over 100 epochs, you've effectively shown the model 10 million unique (though related) examples ( $100,000 * 100$ ). This is the primary strength of this approach.

#### **Crucial Validation Step:**

Your goal is to perform well on data *that looks like it came from the shirt*. Therefore, you must create **two separate test sets**:

1. **Test Set A (Clean):** A held-out portion of the original, clean public data. This tells you how well your model performs on "perfect" data and serves as a baseline.
2. **Test Set B (Augmented):** Take another held-out portion of the clean public data and apply the *exact same random augmentation pipeline* to it. This test set simulates your prototype's data and provides the most realistic estimate of your model's real-world performance.

Your primary metric for success should be the model's performance on **Test Set B**. **Comparing performance on A and B** will tell you how robust your model is to the noise and displacement you expect.