

ECG-Byte: A Tokenizer for End-to-End Generative Electrocardiogram Language Modeling

William Jongwon Han

Carnegie Mellon University, USA

WJHAN@ANDREW.CMU.EDU

Chaojing Duan

Allegheny Health Network, USA

CHAOJING.DUAN@AHN.ORG

Michael A. Rosenberg

University of Colorado, USA

MICHAEL.A.ROSENBERG@CUANSCHUTZ.EDU

Emerson Liu

Allegheny Health Network, USA

EMERSON.LIU@AHN.ORG

Ding Zhao

Carnegie Mellon University, USA

DINGZHAO@ANDREW.CMU.EDU

Abstract

Large Language Models (LLMs) have shown remarkable adaptability across domains beyond text, specifically electrocardiograms (ECGs). More specifically, there is a growing body of work exploring the task of generating text from a multi-channeled ECG and corresponding textual prompt. Current approaches typically involve pretraining an ECG-specific encoder with a self-supervised learning (SSL) objective and using the features output by the pretrained encoder to finetune a LLM for natural language generation (NLG). However, these methods are limited by 1) inefficiency from two-stage training and 2) interpretability challenges with encoder-generated features. To address these limitations, we introduce **ECG-Byte**, an adapted byte pair encoding (BPE) tokenizer pipeline for autoregressive language modeling of ECGs. This approach compresses and encodes ECG signals into tokens, enabling end-to-end LLM training by combining ECG and text tokens directly, while being much more interpretable since the ECG tokens can be directly mapped back to the original signal. Using **ECG-Byte**, we achieve competitive performance in NLG tasks in only *half* the time and \sim 48% of the data required by two-stage approaches.

Data and Code Availability This paper uses the ECG-Chat pretraining (Zhao et al., 2024) and ECG-QA datasets (Oh et al., 2023), which were both created from the MIMIC-IV ECG (Johnson

et al., 2023) and PTB-XL datasets (Wagner et al., 2020). More details about the datasets are provided in Section 3. The ECG-Chat pretraining, ECG-QA, MIMIC-IV ECG, and PTB-XL datasets are all freely available on <https://github.com/YubaoZhao/ECG-Chat>, <https://github.com/Jwoof5/ecg-qa>, and <https://physionet.org/> respectively. Lastly, we have released the code at the following link: <https://github.com/willxxy/ECG-Byte>.

Institutional Review Board (IRB) All datasets used in this study are directly taken from the publicly available, de-identified MIMIC-IV ECG (Johnson et al., 2023) and PTB-XL (Wagner et al., 2020) datasets, thus not requiring IRB approval.

1. Introduction

Cardiovascular diseases (CVDs) are the leading cause of global mortality, with 17.9 million lives taken each year and increasing (Organization, 2024). Due to their readily available, noninvasive and information dense nature, 12-lead ECGs are first-line diagnostic tools for screening/evaluation of potential CVDs. However, accurate ECG analysis is limited in places where ECG expertise is not accessible, exacerbated by the decline and lack of available cardiac electrophysiologists especially in rural areas (Johnson, 2024).

The aforementioned facts calls attention to the need for accessible, accurate, and efficient automation of ECG analysis through deep learning. Deep

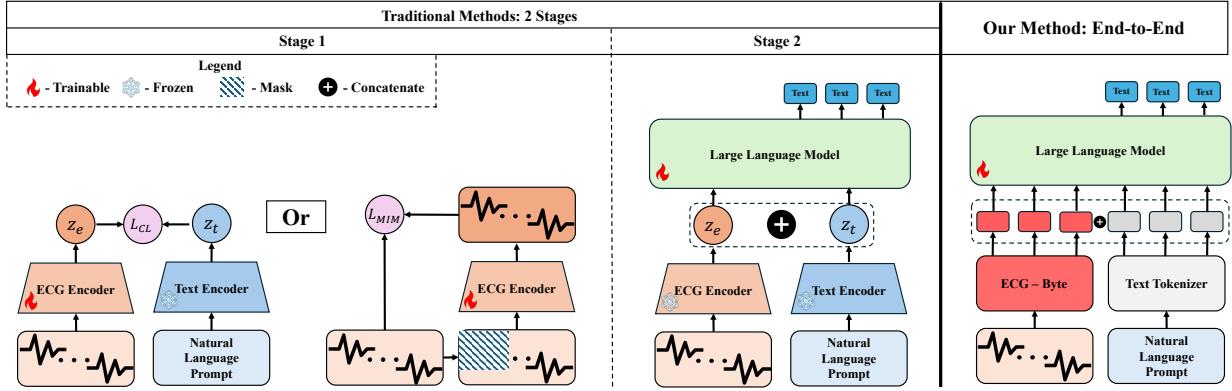


Figure 1: Comparison of traditional methods and our method on ECG language modeling. Traditional methods comprises 2 training stages. The first stage aims to learn a good representation of the 12-lead ECG by training an ECG-specific encoder with a self-supervised learning objective via a combinatorial or individual usage of contrastive learning (L_{CL}) on the hidden states of the ECG z_e and text z_t or masked image modeling (L_{MIM}). The second stage applies the trained encoder to an ECG, which is input alongside a text prompt to a LLM for generation. In contrast, our method is end-to-end and can directly train a LLM for generation utilizing **ECG-Byte**.

learning has reached expert level performance in certain tasks for CVD detection using ECGs (Rajpurkar et al., 2017; Hannun et al., 2019; Qiu et al., 2023b). However, most previous works in this domain have succumbed to a crude classification of hard CVD labels (Choi et al., 2023; Nonaka and Seita, 2020; Martin et al., 2021; Strothoff et al., 2021). A problem with this approach is that ECGs often do not exclusively fall into one diagnostic category, instead, there may be many soft labels annotated by expert physicians and the accumulation of these soft labels allow a more detailed, nuanced, and clinically useful interpretation of the ECG.

The recent onset of Large Language Models (LLMs) provides an opportunity to take a softer, generative, and consequently more flexible approach to ECG analysis. There are some recent works that have taken advantage of this approach to ECG analysis (Qiu et al., 2023a; Tang et al., 2024b; Wan et al., 2024; Fu et al., 2024; Zhao et al., 2024). A commonality among these works is that they treat multi-channel ECGs similarly to images; they first pre-train an encoder specifically for ECGs with some self-supervised learning (SSL) objective, then apply the learned features to finetune a LLM for natural language generation (NLG). However, we observed 2 limitations with this approach: 1) inefficiencies from two-stage training and 2) interpretability challenges

with encoder generated features. Pretraining a good ECG specific encoder can be a significant computational burden due to large datasets, model size, and long training times. Additionally, the latent feature vector output by the ECG encoder cannot be mapped back to the signal, making interpretability difficult when utilizing this feature vector for downstream tasks.

In this study, we introduce **ECG-Byte**, an adapted byte pair encoding (BPE) (Gage, 1994) tokenizer designed for end-to-end training for generative ECG language modeling. Inspired by prior works demonstrating the effectiveness of creating discrete tokens from continuous values (Chen et al., 2022; Han et al., 2024), we leverage quantization to represent amplitude ranges as discrete symbols. Using discrete symbols, we obtain string representations of the ECGs to train **ECG-Byte**. We then apply **ECG-Byte** to directly finetune a LLM for conditional autoregressive NLG, where the text output is conditioned on the text prompt and tokenized signal. We found that our end-to-end training approach is competitive in conditional NLG with only half the time and $\sim 48\%$ of the data required by 2-stage pre-training approaches. Additionally, since the encoded ECG can be reverse transformed back to the original signal, we can interpret token-level attention-based visualizations, whereas it would be impossible to re-

119 verse the hidden latent feature vector outputs by a
 120 pretrained encoder.

121 Our contributions are summarized as the following:

- 122 1. We present **ECG-Byte**, an adapted BPE tok-
 123 enizer for end-to-end training on autoregressive,
 124 conditional NLG.
- 125 2. We empirically show the efficiency of our method
 126 and present competitive performance compared
 127 to conventional 2-stage pretraining approaches,
 128 proposing a new paradigm for conditional NLG
 129 with ECGs.
- 130 3. We conduct an interpretability study on both
ECG-Byte and the LLM by respectively ob-
 131 serving how **ECG-Byte** is merging ECG sig-
 132 nals and by utilizing attention visualizations to
 133 observe how the LLM is processing ECGs and
 134 text.

136 2. Related Works

137 2.1. Deep learning for ECGs

138 There has been a plethora of works utilizing deep
 139 learning for processing ECGs for classification (Ra-
 140 jpurkar et al., 2017; Hannun et al., 2019; Choi et al.,
 141 2023; Nonaka and Seita, 2020; Martin et al., 2021;
 142 Strodtboff et al., 2021). Most of these works utilize
 143 either convolutional neural networks (CNNs) (Ra-
 144 jpurkar et al., 2017) or transformers (Choi et al.,
 145 2023) and exhibit excellent performance at classifica-
 146 tion tasks. There have also been some efforts to frame
 147 classification as a retrieval task in order to recover
 148 cases similar to the given ECG.(Tang et al., 2024b;
 149 Qiu et al., 2023b). However, the retrieval approach
 150 may struggle with rare or unique ECG patterns that
 151 lack good matches in the existing database, poten-
 152 tially leading to missed or inaccurate diagnoses for
 153 unusual cases. Additionally, both classification and
 154 retrieval tasks may be crude formulations of process-
 155 ing ECGs, since ECGs typically exhibit many char-
 156 acteristics of overlapping CVDs.

157 2.2. Large Language Models for ECGs

158 Generative Large Language Models (LLMs) have
 159 given the opportunity to take a softer and more clin-
 160 ically similar approach in processing ECGs by gener-
 161 ating physician-vetted clinical statements (Qiu et al.,
 162 2023a; Tang et al., 2024b; Wan et al., 2024; Fu et al.,
 163 2024; Zhao et al., 2024). The representation of ECG

164 data has been largely considered in previous works
 165 which feed the raw signal into an encoder to obtain
 166 a latent representation of the ECG data, which then
 167 serves as input into a LLM. (Zhao et al., 2024; Wan
 168 et al., 2024; Tang et al., 2024b; Choi et al., 2023). In
 169 order to get good latent representations of the ECGs,
 170 an encoder is first pretrained on a self-supervised
 171 learning (SSL) objective (e.g., contrastive learning,
 172 masked language/image modeling). Although model
 173 performance in terms of label classification has been
 174 excellent when using these approaches, we want to be
 175 able to generate soft labels akin to clinical notation
 176 since ECGs often have overlapping and non-mutually
 177 exclusive descriptors. There have been some efforts
 178 in this direction (Qiu et al., 2023a; Tang et al., 2024b;
 179 Wan et al., 2024; Fu et al., 2024; Zhao et al., 2024),
 180 however, they suffer in efficiency (i.e., requires 2
 181 stages of training). In our work, we challenge this 2-
 182 stage pretraining approach by transforming the ECG
 183 into tokens using **ECG-Byte** and directly training a
 184 LLM for NLG.

185 2.3. Byte Pair Encoding for Domains 186 Outside of Language

187 The Byte Pair Encoding (BPE) algorithm was first
 188 introduced by Gage (1994) for data compression. It
 189 was later adapted to the natural language process-
 190 ing (NLP) domain (Sennrich et al., 2016) and has
 191 been the favored approach to tokenization for the
 192 most popular language models (Grattafiori et al.,
 193 2024; Brown et al., 2020) due to its efficiency and
 194 robustness to rare words. Byte pair encoding has
 195 seen success in representing modalities outside of
 196 language, including molecular graphs (Shen and
 197 Póczos, 2024), electroencephalogram (EEG) (Kly-
 198 menko et al., 2023), and more generally, physiolog-
 199 ical signals.(Tavabi and Lerman, 2021). However, in the-
 200 se cases the byte pair encoded representations are
 201 simply used for classification. Most recently, Tah-
 202 ery et al. (2024) leveraged quantization and BPE to
 203 compress ECG signals and pass as inputs to a BERT
 204 (Devlin et al., 2019) model for SSL. However, in their
 205 work, they only use this representation for classifica-
 206 tion. As previously mentioned, we believe classifica-
 207 tion alone may limit aspects of ECG interpreta-
 208 tion, thus we utilize these representations for genera-
 209 tive diagnosis. Additionally, previous works utilized
 210 a pre-existing BPE tokenizer based on SentencePiece
 211 (Kudo and Richardson, 2018), and do not conduct

212 further analysis of *how* the BPE algorithm is merging
 213 the ECGs.

214 3. Methods

215 This section provides detailed information on the
 216 datasets, preprocessing, ECG signal encoding with
 217 **ECG-Byte**, and LLM training for NLG.

218 3.1. Dataset and Preprocessing

219 **Dataset** In this study, we use variants of the
 220 MIMIC-IV ECG (Gow et al., 2023) and PTB-XL
 221 datasets (Wagner et al., 2020) for NLG. We use
 222 MIMIC-IV ECG pretraining curated by Zhao et al.
 223 (2024) that contains question prompts generated by
 224 GPT-4o alongside the ECG and clinical notes. Ad-
 225 ditionally, we use the ECG-QA dataset (Oh et al.,
 226 2023), a dataset that uses the ChatGPT API to gen-
 227 erate naturalistic, clinically relevant question and an-
 228 swer pairs about the ECG signals from the MIMIC-
 229 IV ECG and PTB-XL datasets. The baselines we
 230 compare our results with all utilize the *single-verify*,
 231 *single-choose*, and *single-query* categorized questions
 232 from the ECG-QA dataset. *single-verify* corresponds
 233 to yes or no questions, *single-choose* corresponds to
 234 where a selected answer is made from two given op-
 235 tions, and *single-query* consists of open-ended ques-
 236 tions. The ECG signals collected from both both
 237 datasets (i.e., MIMIC-IV ECG and PTB-XL) are
 238 sampled at 500 Hz for 10 seconds, resulting in a 5000
 239 length, 12 lead ECG.

240 **Preprocessing** We preprocess all datasets used in
 241 the study in the same manner to maintain consist-
 242 ency. We first convert the ordering of the leads
 243 for the MIMIC-IV ECG dataset (i.e., ['I', 'II', 'III',
 244 ', 'aVR', 'aVF', 'aVL', 'V1', 'V2', 'V3', 'V4', 'V5',
 245 ', 'V6']) to the PTB-XL dataset format (i.e., ['I', 'II',
 246 ', 'III', 'aVL', 'aVR', 'aVF', 'V1', 'V2', 'V3', 'V4', 'V5',
 247 ', 'V6']). We then use a notch filter at 50 Hz and 60 Hz
 248 to remove powerline interference. Each frequency is
 249 targeted with a quality factor of 30 to minimize dis-
 250 tortion, and filtering is applied bidirectionally to pre-
 251 vent phase shifts. Next, a fourth-order Butterworth
 252 bandpass filter with a range of 0.5–100 Hz isolates
 253 relevant ECG components while attenuating high-
 254 frequency noise and low-frequency drift. To address
 255 baseline wander caused by respiratory or movement
 256 artifacts, we bidirectionally apply a fourth-order But-
 257 terworth highpass filter with a cutoff of 0.05 Hz. After
 258 filtering, we apply wavelet denoising to further reduce

259 noise. Using the Daubechies-6 (db6) wavelet at level
 260 4, we decompose each ECG signal into wavelet coef-
 261 ficients. A soft threshold, based on the median ab-
 262 solute deviation of the detail coefficients, is applied
 263 to each coefficient level to suppress noise, ensuring
 264 values near zero are excluded from reconstruction.
 265 Since 250 Hz is a generally accepted sampling fre-
 266 quency adequate for heartbeat analysis (Kwon et al.,
 267 2018), we downsample the 500 Hz sampling frequency
 268 to 250. We then segment the 10 second signal to
 269 non-overlapping windows of 2 seconds, and use each
 270 12 lead 2 second segment of the ECG signal as in-
 271 put to the model. However, for training the tok-
 272 enizer, we did not want to introduce this discontinuity
 273 across the full 10 seconds. Thus, we utilize the un-
 274 segmented, 10 second ECG signal for training **ECG-
 275 Byte**. Lastly, during the unsegmented preprocessing
 276 pipeline, we record the global 1st and 99th percentiles
 277 out of 300,000 samples to utilize in our later steps of
 278 training **ECG-Byte** for normalization.

279 3.2. ECG as Bytes

280 **Sampling** Following established practices in NLP
 281 (Dagan et al., 2024), we train **ECG-Byte** on a rep-
 282 resentative subset of the total dataset, selected us-
 283 ing stratified sampling based on morphological clus-
 284 tering. To extract features from each unsegmented
 285 ECG, we compute statistical measures, frequency and
 286 time domain features, morphological characteristics,
 287 and wavelet coefficients. Principal Component Anal-
 288 ysis (PCA) (Wold et al., 1987) is applied for dimen-
 289 sionality reduction, retaining 95% of the variance, fol-
 290 lowed by feature scaling. The optimal number of clus-
 291 ters is determined using the Elbow Method and Sil-
 292 houette Analysis (Rousseeuw, 1987), with the smaller
 293 result chosen. K-means clustering (MacQueen, 1967)
 294 is then applied to the scaled PCA-transformed fea-
 295 tures. If K-means fails to yield distinct clusters,
 296 DBSCAN (Ester et al., 1996) is used as a fallback.
 297 Stratified sampling is performed by randomly select-
 298 ing ECGs from each cluster in proportion to its size,
 299 resulting in a total sample of 200,000 ECGs.

300 **Quantization** To ensure consistency across ECG
 301 signals, we normalize each input by scaling it to a
 302 fixed range and encoding it into a symbolic repres-
 303 entation. Let $X \in \mathbb{R}^{C \times T}$ denote an ECG signal matrix,
 304 where C is the number of ECG leads and T repre-
 305 sents the number of sampled time points per lead.
 306 In this study, $C = 12$ and $T = 500$ unless specified
 307 otherwise. Let p_1 and p_{99} represent the 1st and 99th

percentiles of X across all leads and time points sampled earlier during preprocessing, respectively. The normalization process is defined as follows:

$$X_{\text{norm}} = \frac{X - (p_1 - \epsilon_1)}{(p_{99} + \epsilon_1) - (p_1 - \epsilon_1) + \epsilon_2} \quad (1)$$

where $\epsilon_1 = 0.5$ is a constant to make up for the sampled percentiles and $\epsilon_2 = 10^{-6}$ is a small constant added to prevent division by zero. This transformation shifts and scales X so that the normalized values fall within the range $[0, 1]$. We then apply clipping to ensure that values remain strictly within this range:

$$X_{\text{clipped}} = \text{clip}(X_{\text{norm}}, 0, 1) \quad (2)$$

Inspired by previous works (Klymenko et al., 2023; Tavabi and Lerman, 2021; Chen et al., 2022; Han et al., 2024), we quantize X_{clipped} into discrete levels for symbolic representation. Let \mathcal{A} be the set of 26 symbols, corresponding to the lowercased letters in the English alphabet, $\mathcal{A} = \{a, b, \dots, z\}$. The alphabet size $|\mathcal{A}| = 26$ defines the number of discrete levels. We scale and floor X_{clipped} to integer values, then take the minimum between the floored value and the maximum number of bins as the following:

$$X_{\text{quant}} = \min(\lfloor X_{\text{clipped}} \times |\mathcal{A}| \rfloor, (|\mathcal{A}| - 1)) \quad (3)$$

Finally, each integer value in X_{quant} is mapped to a corresponding symbol in \mathcal{A} to yield the symbolic signal, which serves as a discrete representation of the ECG. After transforming each ECG signal instance into its symbolic form, we first flatten each symbolic ECG instance $X_{\text{symb}}^{(i)}$ into a 1-dimensional sequence of symbols:

$$X_{\text{symb}}^{(i)} = \text{flatten}(X_{\text{quant}}^{(i)}), \quad X_{\text{symb}}^{(i)} \in \mathcal{A}^{CT} \quad (4)$$

where i indexes over all instances in the dataset, and $X_{\text{symb}}^{(i)}$ is the flattened sequence of symbols of length $C \cdot T$. Next, we concatenate all flattened instances $X_{\text{symb}}^{(1)}, X_{\text{symb}}^{(2)}, \dots, X_{\text{symb}}^{(N)}$ across the entire dataset to form a single, long symbolic sequence:

$$\mathbf{X}_{\text{concat}} = X_{\text{symb}}^{(1)} \| X_{\text{symb}}^{(2)} \| \cdots \| X_{\text{symb}}^{(N)}, \quad \mathbf{X}_{\text{concat}} \in \mathcal{A}^{NCT} \quad (5)$$

where $\|$ denotes the concatenation operation, and N is the total number of instances in the dataset. The concatenated symbolic sequence $\mathbf{X}_{\text{concat}}$ of length $N \cdot C \cdot T$ is then used to train **ECG-Byte**.

ECG-Byte Training Process After obtaining the string representation $\mathbf{X}_{\text{concat}}$ of the ECG dataset, we train **ECG-Byte** to compress the discretized ECG signals by iteratively merging the most frequent byte pairs into single tokens, following the BPE algorithm. The process starts by converting $\mathbf{X}_{\text{concat}}$ into an ID vector of 8-bit unsigned integers and initializing a vocabulary map (`vocab`) for string representations of bytes and a `vocab_tokens` map to encode bytes as singleton lists. IDs and `vocab` are initialized to cover the full byte range (0–255), mapping symbols in \mathcal{A} to ASCII values (97–122), while reserving other byte values for unknown bytes. As merging proceeds, new tokens are assigned unique integer IDs starting from 256, acting as abstract labels for progressively larger token units. For each merge iteration, **ECG-Byte** calculates adjacent byte pair frequencies using a parallelized `get_stats` function, efficiently aggregating counts via a fold-and-reduce strategy. The most frequent pair is identified as the "best pair" to merge, and the `merge` function replaces occurrences of this pair in the ID vector with a new token ID, extending the vocabulary and updating `vocab_tokens` accordingly. This process repeats until the specified number of merges is reached or no pairs remain. The output includes the encoded ID vector, the extended vocabulary map, and a history of merge operations. Existing tokenizers, such as SentencePiece (Kudo and Richardson, 2018) or HuggingFace (Wolf et al., 2020), were not used due to their complexity and integration issues, which hindered interpretability. **ECG-Byte**, implemented in Rust for speed, provides a lightweight, flexible framework for representing ECG signals as discrete tokens while drawing inspiration from HuggingFace's tokenizer (Wolf et al., 2020). Detailed pseudocode for the main training pipeline is provided in Algorithm 1, with `merge` and `get_stats` functions detailed in Appendix A.1.

ECG-Byte Encoding Process After training **ECG-Byte**, we encode any quantized ECG signal X_{symb} by first converting each byte in the ECG signal to a 32-bit unsigned integer and building a trie structure, where each node represents a byte or a merged token sequence from prior encoding steps. The trie is initialized with single-byte tokens (0–255) and is extended with custom token sequences from the learned merge history. For each byte sequence in the input, the encoding function traverses the trie to find the longest match, replacing matched sequences with their assigned token IDs. If no match is found,

393 the byte is added to the output as-is. The final en-
 394 coded sequence is returned as `output_ids`, where we
 395 will denote as X_{ID} . We provide the detailed pseu-
 396 docode of the encoding process in Algorithm 2.

Algorithm 1: Training Process for ECG-Byte

Input: Input X_{concat} , Number of merges `num_merges`

Output: Tuple containing final encoded IDs, vocabulary map, and merge history

Function `byte_pair_encoding(X_{concat} , num_merges):`

- | Convert X_{concat} to ID vector `ids` by casting each byte to `u32`;
- | Initialize `vocab` with mappings from IDs 0 to 255 to their string representations;
- | Initialize `vocab_tokens` with mappings from IDs 0 to 255 to singleton lists;
- | Initialize empty list `merges`;
- | **for** $i \leftarrow 0$ **to** `num_merges exclusive` **do**
- | | `pairs` \leftarrow `get_stats(ids)` using parallel processing;
- | | **if** `pairs` is empty **then**
- | | | **break**
- | | **end**
- | | `best_pair` \leftarrow Pair in `pairs` with highest frequency;
- | | **if** `best_pair` is not found **then**
- | | | **break**
- | | **end**
- | | `new_id` $\leftarrow 256 + i$;
- | | `ids` \leftarrow `merge(ids, best_pair, new_id)`;
- | | `vocab[new_id]` \leftarrow concat(`vocab[best_pair.0]`, `vocab[best_pair.1]`);
- | | `new_token` \leftarrow concat(`vocab_tokens[best_pair.0]`, `vocab_tokens[best_pair.1]`);
- | | `vocab_tokens[new_id]` \leftarrow `new_token`;
- | | `merges.append(new_token, new_id)`;
- | **end**
- | **return** (`ids, vocab, merges`);

3.3. Large Language Model

397 In this study, we utilize the Llama-3.2-1B (Grattafiori
 398 et al., 2024) checkpoint through the HuggingFace API
 399 (Wolf et al., 2020) unless specified otherwise. The
 400 Llama 3.2 series model is a variant of the Llama 3
 401 models (Grattafiori et al., 2024) and support context
 402

Algorithm 2: Encoding Process for ECG-Byte

Input: Input X_{symb} , Merge history `merges` containing pairs of token sequences and their token IDs

Output: Vector of encoded IDs

Function `encode(X_{symb} , merges):`

- | `ids` \leftarrow Convert X_{symb} to vector of `u32` by casting each byte;
- | Initialize root `TrieNode trie_root` using `TrieNode::new()`;
- | **for** $b \leftarrow 0$ **to** 255 **do**
- | | `insert(trie_root, [b], b)`;
- | **end**
- | **foreach** (`token_sequence, token_id`) **in** `merges` **do**
- | | `insert(trie_root, token_sequence, token_id)`;
- | **end**
- | Initialize empty list `output_ids`;
- | $i \leftarrow 0$;
- | **while** $i < length of ids$ **do**
- | | `node` \leftarrow `trie_root`;
- | | `match_len` $\leftarrow 0$;
- | | `match_id` \leftarrow `None`;
- | | **for** $j \leftarrow i$ **to** `length of ids` **do**
- | | | `id` \leftarrow `ids[j]`;
- | | | **if** `id` exists in `node.children` **then**
- | | | | `node` \leftarrow `node.children[id]`;
- | | | | **if** `node.token_id` is not `None` **then**
- | | | | | `match_len` $\leftarrow j - i + 1$;
- | | | | | `match_id` \leftarrow `node.token_id`;
- | | | | **end**
- | | | **end**
- | | | **else**
- | | | | **break**
- | | | **end**
- | | **end**
- | | **if** `match_id` is not `None` **then**
- | | | `output_ids.append(match_id)`;
- | | | $i \leftarrow i + match_len$;
- | | **end**
- | | **else**
- | | | `output_ids.append(ids[i])`;
- | | | $i \leftarrow i + 1$;
- | | **end**
- | **end**
- | **return** `output_ids`;

lengths of up to 128k tokens. They are notable for their superior performance despite having 1 billion parameters, making them highly efficient and capable models to test our methodology upon. We also provide an ablation study in subsection 5.4 where we utilize other popular LLMs, such as GPT2 XL 1.5B (Radford et al., 2019), Gemma 2B (Team et al., 2024), and OPT 1.3B (Zhang et al., 2022).

3.4. Learning Objective

The learning objective for training the LLM considers a sequence composed of three parts, $\{X_{ID}, Q, \mathcal{S}\}$, where $X_{ID} \in \mathcal{V}^M$ represents the encoded ECG sequence of length $l_{X_{ID}} = |X_{ID}|$, with each token drawn from the extended vocabulary \mathcal{V} of size M , Q represents the tokenized question, and \mathcal{S} denotes the tokenized answer sequence. The input sequence includes special tokens: - [BOS] as the beginning-of-sequence token, - [SIG_START] and [SIG_END] to indicate the start and end of the encoded ECG sequence, and - [EOS] as the end-of-sequence token for the generated answer. The motivation for adding [SIG_START] and [SIG_END] special tokens is inspired by Liu et al. (2023), where they utilize special tokens indicating the start and end of the image. Thus, the full input sequence is structured as:

$$[\text{BOS}] \parallel [\text{SIG_START}] \parallel X_{ID} \parallel [\text{SIG_END}] \parallel Q \parallel \mathcal{S} \parallel [\text{EOS}],$$

where \parallel denotes concatenation. Let $l_Q = |Q|$, $l_{\mathcal{S}} = |\mathcal{S}|$, and L be the total sequence length, given by:

$$L = 1 + 1 + l_{X_{ID}} + 1 + l_Q + l_{\mathcal{S}} + 1,$$

accounting for the [BOS], [SIG_START], [SIG_END], and [EOS] tokens. The autoregressive objective maximizes the likelihood of each token in $\mathcal{S} \parallel [\text{EOS}]$ conditioned on the preceding context $\text{Context} = \{[\text{BOS}], [\text{SIG_START}], X_{ID}, [\text{SIG_END}], Q\}$ and the previous tokens in \mathcal{S} . The objective is formulated as follows:

$$\mathcal{L}_{NLL} = - \sum_{l'=l_{X_{ID}}+l_Q+4}^L \log P(s_{l'} \mid \text{Context}, s_{<l'}; \theta), \quad (6)$$

where $s_{l'} = \mathcal{S}_{l'-(l_{X_{ID}}+l_Q+4)}$ is the $(l' - (l_{X_{ID}} + l_Q + 4))$ -th token in $\mathcal{S} \parallel [\text{EOS}]$, and $s_{<l'} = \{s_1, s_2, \dots, s_{l'-(l_{X_{ID}}+l_Q+4)-1}\}$ denotes all tokens in $\mathcal{S} \parallel [\text{EOS}]$ preceding $s_{l'}$.

4. Experiments

4.1. Experimental Settings

We fine-tuned the LLM using the AdamW optimizer (Kingma and Ba, 2017) with a learning rate of $1e-4$, weight decay of $1e-2$, and a custom learning rate scheduler. This scheduler applies an initial learning rate init_lr scaled by the model’s hidden dimension ($d_{\text{model}}^{-0.5}$) and dynamically adjusts it based on training steps, with a warm-up phase of 500 steps. The learning rate at step n_{steps} is updated as $\text{lr} = \text{init_lr} \times \min(n_{\text{steps}}^{-0.5}, n_{\text{warmup}}^{-1.5} \times n_{\text{steps}})$. We set $\beta_1 = 0.9$, $\beta_2 = 0.99$, $\epsilon = 1e-8$, batch size 2, and trained for 1 epoch. Additionally, we only train on a randomly sampled subset of 400,000 ECG instances for each respective dataset due to computational resources, unless specified otherwise. We also utilize LoRA (Hu et al., 2021) to finetune the LLM with rank = 16, $\alpha_{LoRA} = 32$, and dropout = 0.05. We conduct our experiments on 4 NVIDIA RTX A6000 48 GB GPUS.

During inference, we evaluate our model with number of merges $\text{num_merges} = 3500$, sequence length $L = 1024$, and ECG length $T = 500$ unless specified otherwise. We use popular metrics for NLG namely the BLEU-4 (Papineni et al., 2002), Rouge-L (Lin, 2004), Meteor (Banerjee and Lavie, 2005), and BertScore F1 (Zhang et al., 2020) metrics.

5. Results

5.1. Natural Language Generation

We present our main results in Table 1, comparing **ECG-Byte** with prior works and self-implemented two-stage pretraining methods. Notably, Zhao et al. (2024) is not directly comparable due to differing data splits and pretraining datasets, though reported metrics use the same datasets. Zhao et al. (2024) train on the *full* MIMIC-IV ECG Pretrain dataset, finetune on an instruction-tuning dataset for ECG-related conversations, and evaluate on PTB-XL (Wagner et al., 2020) using a unified question: “Could you please help me explain my ECG?” To establish baselines, we implement generic two-stage pretraining methods: L_{CL} , L_{MIM} , and $L_{CL} + L_{MIM}$. Here, L_{CL} employs contrastive learning (Liu et al., 2024; Gopal et al., 2021; Pham et al., 2024; Kiyasseh et al., 2021), L_{MIM} uses Masked Image Modeling (MIM) (Choi et al., 2023; Na et al., 2024; Yang et al., 2022), and $L_{CL} + L_{MIM}$ combines both (Oh et al., 2022;

Table 1: NLG mean results with standard deviations over 5 random seeds comparing against different baselines.

Method	Trained Dataset	Inferenced Dataset	BLEU-4	Rouge-L	Meteor	BertScore F1
ECG-Chat (Zhao et al., 2024)			11.19	29.93	35.10	-
L_{CL}			8.10 ± 0.25	31.36 ± 0.31	27.55 ± 0.36	89.35 ± 0.04
L_{MIM}			6.21 ± 0.22	30.63 ± 0.13	24.91 ± 0.14	90.44 ± 0.04
L_{MERL} (Liu et al., 2024)	MIMIC-IV ECG Pretrain	PTB-XL	10.22 ± 0.25	32.95 ± 0.12	25.60 ± 0.17	89.94 ± 0.01
$L_{CL} + L_{MIM}$			9.33 ± 0.22	30.45 ± 0.21	24.37 ± 0.36	90.29 ± 0.02
ECG-Byte			11.00 ± 0.19	33.41 ± 0.05	24.95 ± 0.09	90.02 ± 0.01
L_{CL}			10.22 ± 0.06	38.41 ± 0.48	24.66 ± 0.23	90.42 ± 0.09
L_{MIM}			7.90 ± 0.23	29.28 ± 0.38	19.03 ± 0.11	67.91 ± 0.17
L_{MERL} (Liu et al., 2024)	ECG-QA MIMIC-IV	ECG-QA MIMIC-IV	10.95 ± 0.24	38.18 ± 0.58	26.24 ± 0.36	90.80 ± 0.06
$L_{CL} + L_{MIM}$			8.57 ± 0.14	34.00 ± 0.25	25.22 ± 0.30	87.72 ± 0.04
ECG-Byte			11.23 ± 0.12	42.49 ± 0.53	27.08 ± 0.15	91.30 ± 0.04
L_{CL}			8.89 ± 0.25	28.63 ± 0.47	18.45 ± 0.31	72.63 ± 0.40
L_{MIM}			15.14 ± 0.28	46.71 ± 0.41	29.64 ± 0.30	92.12 ± 0.10
L_{MERL} (Liu et al., 2024)	ECG-QA PTB-XL	ECG-QA PTB-XL	13.84 ± 0.19	40.14 ± 0.39	26.24 ± 0.35	91.88 ± 0.09
$L_{CL} + L_{MIM}$			14.72 ± 0.27	42.88 ± 0.13	28.25 ± 0.27	89.40 ± 0.01
ECG-Byte			13.93 ± 0.21	47.08 ± 0.56	29.17 ± 0.31	92.53 ± 0.07

McKeen et al., 2024). These approaches utilize pre-trained CLIP (Radford et al., 2021) and ViT (Dosovitskiy et al., 2021), where ECG signals are transformed into three-channel images for finetuning. We adapt Liu et al. (2024)'s state-of-the-art contrastive method (L_{MERL}) for fair comparison. Their most effective model uses a 1D ResNet backbone (He et al., 2015); hence, we employ ResNet101 for direct ECG signal processing. For L_{CL} , L_{MIM} , $L_{CL} + L_{MIM}$, and L_{MERL} , training is conducted on the *full*, pre-processed MIMIC-IV ECG dataset with a batch size of 64 during the first stage. Implementation details for both training stages are in Appendix B. Table 1 demonstrates **ECG-Byte**'s effectiveness, showing competitive or superior performance across all metrics and datasets compared to other methods. Qualitative examples are provided in Appendix C.2.

5.2. Cross Dataset Transferability

We present the results of cross-dataset transferability in Table 2, comparing our approach, **ECG-Byte**, with two-stage pretraining methods. **ECG-Byte** achieves the best zero-shot transfer performance from the ECG-QA PTB-XL dataset to the ECG-QA MIMIC-IV dataset. When transferring from the ECG-QA MIMIC-IV dataset to the ECG-QA PTB-XL dataset, although other 2-stage pre-training methods demonstrate higher performance, **ECG-Byte** maintains competitive results across all metrics.

5.3. Efficiency of ECG-Byte

We compare the efficiency of our end-to-end approach using **ECG-Byte** with 2-stage pretraining methods in Table 3. First, we examine the amount of data required for each method. As previously noted, the first stage of the two-stage pretraining methods is trained on the full MIMIC-IV ECG dataset (Johnson et al., 2023) using segmented ECGs, which are also used as input during the second stage. While **ECG-Byte** is trained on unsegmented ECGs, we convert the number of unsegmented ECGs to an equivalent count of segmented ECGs. Additionally, the reduced data requirement for **ECG-Byte** is due to our sampling approach, where only a subset of ECGs is used to train the tokenizer. Under these settings, our proposed method achieves competitive results using approximately ~48% of the data required for 2-stage pretraining methods. In terms of training time, our approach requires less than *half* the time needed for two-stage pretraining. The training time for the two-stage methods is averaged across our self-implemented approaches (L_{CL} , L_{MIM} , and $L_{CL} + L_{MIM}$) and the L_{MERL} method proposed by Liu et al. (2024).

5.4. Ablation Study

We conduct several ablation studies to show the variability of performance with **ECG-Byte** when we alter the LLM used for finetuning, use different sequence lengths L when inputting to the LLM, training **ECG-Byte** with various number of merges

Table 2: Mean results with standard deviations over 5 random seeds on zero shot cross-dataset transferability.

Method	Trained Dataset	Inferenced Dataset	BLEU-4	Rouge-L	Meteor	BertScore F1
L_{CL}			11.64 ± 0.45	41.48 ± 0.11	25.74 ± 0.13	91.24 ± 0.05
L_{MIM}			$\mathbf{11.70} \pm 0.29$	$\mathbf{42.22} \pm 0.28$	$\mathbf{26.41} \pm 0.10$	91.51 ± 0.03
L_{MERL} (Liu et al., 2024)	ECG-QA MIMIC-IV	ECG-QA PTB-XL	11.53 ± 0.19	39.23 ± 0.40	25.58 ± 0.28	$\mathbf{91.59} \pm 0.03$
$L_{CL} + L_{MIM}$			9.71 ± 0.10	35.10 ± 0.28	24.91 ± 0.19	87.88 ± 0.08
ECG-Byte			8.70 ± 0.04	40.39 ± 0.40	23.29 ± 0.18	91.51 ± 0.03
L_{CL}			5.10 ± 0.04	22.77 ± 0.28	14.63 ± 0.32	77.89 ± 0.13
L_{MIM}			7.68 ± 0.46	$\mathbf{35.77} \pm 0.13$	$\mathbf{22.32} \pm 0.33$	90.28 ± 0.07
L_{MERL} (Liu et al., 2024)	ECG-QA PTB-XL	ECG-QA MIMIC-IV	7.39 ± 0.15	28.33 ± 0.58	18.59 ± 0.35	89.30 ± 0.05
$L_{CL} + L_{MIM}$			7.49 ± 0.21	30.53 ± 0.59	20.25 ± 0.27	86.53 ± 0.11
ECG-Byte			$\mathbf{7.86} \pm 0.13$	35.01 ± 0.41	21.49 ± 0.24	$\mathbf{90.29} \pm 0.07$

Table 3: Efficiency of our method compared against 2-stage pretraining methods.

	1st Stage	2nd Stage	ECG-Byte	end-to-end
# of Data	2,513,435	400,000	1,000,000	400,000
Total # of Data		2,913,435		1,400,000
Time (minutes)	~1258.50	~469.25	~385.12	~420.32
Total Time (minutes)		~1727.75		~805.44

547 num_merges, and varying ECG lengths T . With the
548 exception of the ablating parameter, we fix all other
549 parameters to num_merges = 3500, $L = 1024$, and
550 $T = 500$. We report results on the test set of the
551 PTB-XL variant of ECG-QA (Oh et al., 2023) unless
552 specified otherwise.

Table 4: Ablation study on using different LLMs.

LLM	BLEU-4	Rouge-L	Meteor	BertScore F1
GPT2 XL 1.5B (Radford et al., 2019)	12.30 ± 0.19	41.33 ± 0.57	26.48 ± 0.33	92.00 ± 0.06
Gemma 2B (Team et al., 2024)	13.78 ± 0.18	45.48 ± 0.55	28.32 ± 0.23	92.01 ± 0.02
OPT 1.3B (Zhang et al., 2022)	12.26 ± 0.20	41.84 ± 0.52	26.21 ± 0.29	91.78 ± 0.04
Llama 3.2 1B (Grattafiori et al., 2024)	13.93 ± 0.21	$\mathbf{47.08} \pm 0.56$	29.17 ± 0.31	$\mathbf{92.53} \pm 0.07$

553 **Different LLMs** We show the variability in per-
554 formance of **ECG-Byte** when using different LLMs
555 with similar numbers of parameters in Table 4. While
556 Llama 3.2 1B (Grattafiori et al., 2024) achieves the
557 best results, GPT2 XL 1.5B (Radford et al., 2019),
558 Gemma 2B (Team et al., 2024), and OPT 1.3B
559 (Zhang et al., 2022) also deliver comparable perfor-
560 mances. These findings demonstrate that our method
561 is not limited to Llama 3.2 1B but can achieve similar
562 results across a variety of LLMs.

563 **Sequence Length** Input lengths for LLMs are an
564 important parameter to consider for efficient training

since the calculation of attention is quadratic with
565 respect to the input length (Vaswani et al., 2023).
566 We present results on different sequence lengths L in
567 Table 5. Although the difference is not substantial,
568 we can see that when $L = 1024$ and $L = 2048$ the
569 model yields higher performance than $L = 512$. We
570 attribute this to the rate of truncation and padding
571 for the encoded ECG. Observing Figure 2, most
572 ECGs were being encoded to token sequence lengths
573 of around 500 to 1500. Therefore, we hypothesize
574 that when $L = 512$ a large portion of the ECG to-
575 ken sequence gets truncated, resulting in lower per-
576 formance.
577

Table 5: Ablation study on varying sequence lengths L .

L	BLEU-4	Rouge-L	Meteor	BertScore F1
512	13.61 ± 0.15	$\mathbf{48.15} \pm 0.57$	29.10 ± 0.28	92.41 ± 0.05
1024	$\mathbf{13.93} \pm 0.21$	47.08 ± 0.56	$\mathbf{29.17} \pm 0.31$	$\mathbf{92.53} \pm 0.07$
2048	13.88 ± 0.22	45.21 ± 0.48	28.31 ± 0.27	90.88 ± 0.02

580 **Number of Merges** The number of merges
581 num_merges performed during training **ECG-Byte**
582 corresponds to how much the algorithm compresses
583 the concatenated sequence of quantized ECGS
584 $\mathbf{X}_{\text{concat}}$. More num_merges means more compression,
585 which can affect the expressiveness of the encoded se-
586 quence. In Table 6, we show the performance of our
587 method with different num_merges. The results indi-
588 cate that while performance varies slightly with the
589 number of merges, values of num_merges greater than
590 500 generally yield similar outcomes.

591 **ECG length** Lastly, we show the effect of the
592 length T being considered when encoding the ECG

Table 6: Ablation study on varying number of merges `num_merges`.

<code>num_merges</code>	BLEU-4	Rouge-L	Meteor	BertScore F1
500	13.61 ± 0.53	46.50 ± 0.28	28.49 ± 0.49	92.33 ± 0.02
1750	14.50 ± 0.25	46.74 ± 0.48	30.03 ± 0.25	92.55 ± 0.01
2500	15.10 ± 0.39	46.37 ± 0.28	30.12 ± 0.23	92.53 ± 0.05
3500	13.93 ± 0.21	47.08 ± 0.56	29.17 ± 0.31	92.53 ± 0.07

591 with **ECG-Byte** in Table 7. We want to note that
 592 for the results of $T = 2000$, the full unsegmented
 593 ECG is utilized. Consequently, the number of in-
 594 stances available is less than the targeted dataset size
 595 of 400,000 (i.e., 97,244). Thus, when $T = 2000$, we
 596 use the full dataset to train the model. For shorter
 597 segment lengths, such as $T = 250$ and $T = 500$,
 598 the model demonstrates strong performances indicat-
 599 ing that shorter segments can effectively preserve rel-
 600 evant information for NLG. Interestingly, for $T =$
 601 2000, the model achieves the highest performance
 602 across all metrics. This suggests that when the model
 603 is trained with the full 10 second encoded ECG, it
 604 benefits from richer contextual information present
 605 in the complete ECG waveform.

Table 7: Ablation study on varying lengths T .

T	BLEU-4	Rouge-L	Meteor	BertScore F1
250	12.64 ± 0.20	47.31 ± 0.26	27.97 ± 0.21	92.32 ± 0.06
500	13.93 ± 0.21	47.08 ± 0.56	29.17 ± 0.31	92.53 ± 0.07
1250	11.01 ± 0.19	43.84 ± 0.28	25.49 ± 0.20	93.07 ± 0.03
2000	14.54 ± 0.17	48.03 ± 0.27	32.11 ± 0.22	92.91 ± 0.04

5.5. ECG-Byte Analysis

We analyze **ECG-Byte** by visualizing the usage of merged tokens, length of the encoded ECG, and mapping between the encoded tokens and original ECG. Unless specified otherwise, we analyze **ECG-Byte** when `num_merges` = 3500, L = 1024, and T = 500.

Token Usage and Length Distribution We examine the token usage and length distributions for **ECG-Byte** with `num_merges` = 3500 on a subsample of 277,840 ECGs from the PTB-XL dataset. The left panel of Figure 2 displays the token usage distribution, showing token frequency (y-axis) ranked in descending order (x-axis). A small subset of tokens dominates the occurrences, while the rest are infrequently used—a typical characteristic of BPE-based

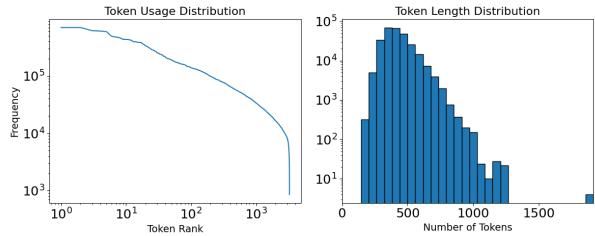


Figure 2: Plots of the token usage and length distributions for **ECG-Byte** where `num_merges` = 3500. More examples with varying `num_merges` are provided in Appendix A.3.

tokenization, where common patterns are compressed into frequent tokens and rare patterns into infrequent ones. The right panel of Figure 2 illustrates the token length distribution of the encoded ECGs, with most falling between 500 and 1000 tokens, demonstrating **ECG-Byte**'s effective compression of the original signal. Additional examples of these distributions are provided in Appendix A.3.

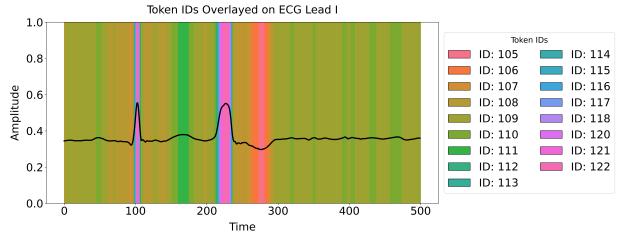


Figure 3: A mapping between tokens used for a given ECG Lead I. More examples are provided in Appendix A.2.

Token to ECG Mapping To illustrate how **ECG-Byte** encodes ECG signals, we analyze the mapping between tokens and signal features. Figure 3 shows an example Lead I ECG signal with unique token IDs (represented by different colors) overlaid. The P wave, QRS complex, and T wave are distinctly captured by different tokens, though this precision varies across instances. As demonstrated, **ECG-Byte** effectively merges key regions of the signal. Additional examples are provided in Appendix A.2 due to page limitations.

640 **Attention Visualizations** Figure 4 visualizes attention weights across a selected ECG lead and text portions of the input after training. We focus on one lead due to the uniformity of attention across encoded signal tokens. For interpretability, the reversed ECG signal is overlayed on the encoded ECG. The model primarily attends to the textual portion of the input sequence, as shown in Figure 4. Previous studies have debated whether attention visualizations are inherently explainable (Jain and Wallace, 2019; Wiegrefe and Pinter, 2019) and explored their role in vision-language models (Aflalo et al., 2022; Woo et al., 2024; Arif et al., 2024; Cui et al., 2024). These works often observe minimal attention to visual input, with models relying primarily on text. We hypothesize that a similar phenomenon occurs in Figure 4, as the ECG tokens, though represented like text, are 1) newly introduced and 2) perceived as a different modality (e.g., vision). Additional examples are provided in Appendix A.4.

6. Discussion and Conclusion

661 In this study, we introduce **ECG-Byte**, a custom
 662 BPE algorithm to encode ECGs into a discrete sequence of tokens for conditional autoregressive NLG.
 663 **ECG-Byte** introduces a paradigm shift in generative
 664 ECG language modeling by enabling efficient end-to-end training, compared to traditional two-stage
 665 pretraining approaches. Our pipeline demonstrates
 666 strong performances, achieving results comparable
 667 to two-stage methods while requiring only *half* the
 668 training time and approximately 48% of the data. In
 669 addition to its efficiency, **ECG-Byte** enhances inter-
 670 pretability. By analyzing its underlying mechanism,
 671 we observe that critical ECG regions, such as the P
 672 wave, the QRS complex, and the T wave, are effec-
 673 tively grouped during tokenization, as illustrated in
 674 Figure 3. Furthermore, the reversibility of the com-
 675 pressed token sequence allows us to trace each token
 676 back to its original ECG signal segment, providing in-
 677 sight into the specific portions of the signal attended
 678 to by the model. However, as shown in Figure 4,
 679 the model’s attention weight distribution resembles
 680 that of vision language models, focusing primarily on
 681 the textual components of the input sequence during
 682 generation.

683 This work is in its early stages and needs fur-
 684 ther exploration. Future directions include: (1) re-
 685 fining BPE merging rules to better capture ECG-
 686 specific features, (2) adopting more advanced quan-

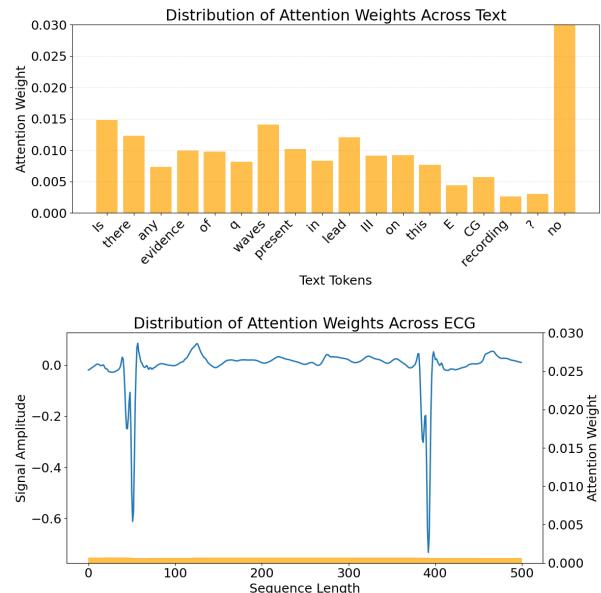


Figure 4: The attention weight overlaid on top of both the text (top) and ECG (bottom). More examples are provided in Appendix A.4.

687 tization techniques that preserve time-series character-
 688 istics (Carson et al., 2024b; Elsworth and Güttel,
 689 2020; Carson et al., 2024a), (3) introducing stronger
 690 modality-specific distinctions, such as embeddings
 691 beyond [SIG_START] and [SIG_END] (Gui et al.,
 692 2023), and (4) extending **ECG-Byte** for conversa-
 693 tional tasks through instruction tuning.
 694

Limitations One of the main limitations of this
 695 work is the scale in terms of computing and data.
 696 Since we only used a subset of the data to train and
 697 test our method, we were unable to train the model
 698 to its full potential. However, even with only using
 699 a small subset of the data, we are able to see ex-
 700 tremely promising results compared to other 2-stage
 701 SSL pretraining methods. Therefore, we do not view
 702 this limitation as a major bottleneck.
 703

Acknowledgments

705 This work is done in collaboration with the Mario
 706 Lemieux Center for Heart Rhythm Care at Allegheny
 707 General Hospital. We thank Wenhao Ding, Haohong
 708 Lin, Shiqi Liu, and Hyoeun Kang for the valuable
 709 discussions.
 710

711 References

- 712 Estelle Aflalo, Meng Du, Shao-Yen Tseng, Yongfei
 713 Liu, Chenfei Wu, Nan Duan, and Vasudev Lal.
 714 Vl-interpret: An interactive visualization tool for
 715 interpreting vision-language transformers, 2022.
 716 URL <https://arxiv.org/abs/2203.17247>.
- 717 Kazi Hasan Ibn Arif, JinYi Yoon, Dimitrios S.
 718 Nikolopoulos, Hans Vandierendonck, Deepu John,
 719 and Bo Ji. Hired: Attention-guided token drop-
 720 ping for efficient inference of high-resolution vision-
 721 language models in resource-constrained environ-
 722 ments, 2024. URL <https://arxiv.org/abs/2408.10945>.
- 723 Satanjeev Banerjee and Alon Lavie. Meteor: An au-
 724 tomatic metric for mt evaluation with improved
 725 correlation with human judgments. In *IEEvalu-
 726 ation@ACL*, 2005.
- 727 Tom B. Brown, Benjamin Mann, Nick Ryder,
 728 Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
 729 Arvind Neelakantan, Pranav Shyam, Girish Sastry,
 730 Amanda Askell, Sandhini Agarwal, Ariel Herbert-
 731 Voss, Gretchen Krueger, Tom Henighan, Rewon
 732 Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey
 733 Wu, Clemens Winter, Christopher Hesse, Mark
 734 Chen, Eric Sigler, Mateusz Litwin, Scott Gray,
 735 Benjamin Chess, Jack Clark, Christopher Berner,
 736 Sam McCandlish, Alec Radford, Ilya Sutskever,
 737 and Dario Amodei. Language models are few-
 738 shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>.
- 739
- 740 Erin Carson, Xinye Chen, and Cheng Kang. Llm-
 741 abba: Understand time series via symbolic approx-
 742 imation, 2024a. URL <https://arxiv.org/abs/2411.18506>.
- 743
- 744 Erin Carson, Xinye Chen, and Cheng Kang. Quan-
 745 tized symbolic time series approximation, 2024b.
 746 URL <https://arxiv.org/abs/2411.15209>.
- 747
- 748 Ting Chen, Saurabh Saxena, Lala Li, David J. Fleet,
 749 and Geoffrey Hinton. Pix2seq: A language model-
 750 ing framework for object detection, 2022.
- 751 Seokmin Choi, Sajad Mousavi, Phillip Si, Haben G.
 752 Yhdego, Fatemeh Khadem, and Fatemeh Afghah.
 753 Ecgbert: Understanding hidden language of ecgs
 754 with self-supervised representation learning, 2023.
- Chenhang Cui, Jiabing Yang, Yiyang Zhou, Peng
 755 Xia, Ying Wei, and Huaxiu Yao. Fading focus:
 756 Mitigating visual attention degradation in
 757 large vision-language models, 2024. URL <https://openreview.net/forum?id=gam5LiMPKT>.
- Gautier Dagan, Gabriel Synnaeve, and Baptiste
 760 Rozière. Getting the most out of your tokenizer for
 761 pre-training and domain adaptation, 2024. URL
 762 <https://arxiv.org/abs/2402.01035>.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
 764 Kristina Toutanova. Bert: Pre-training of deep
 765 bidirectional transformers for language under-
 766 standing, 2019.
- Alexey Dosovitskiy, Lucas Beyer, Alexander
 768 Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
 769 Thomas Unterthiner, Mostafa Dehghani, Matthias
 770 Minderer, Georg Heigold, Sylvain Gelly, Jakob
 771 Uszkoreit, and Neil Houlsby. An image is worth
 772 16x16 words: Transformers for image recogni-
 773 tion at scale, 2021.
- 774
- Steven Elsworth and Stefan Güttel. Abba: Adaptive
 775 brownian bridge-based symbolic aggregation
 776 of time series, 2020. URL <https://arxiv.org/abs/2003.12469>.
- Martin Ester, Hans-Peter Kriegel, Jörg Sander, and
 779 Xiaowei Xu. A density-based algorithm for discov-
 780 ering clusters in large spatial databases with noise.
 781 In *Proceedings of the Second International Confer-
 782 ence on Knowledge Discovery and Data Mining*,
 783 KDD'96, page 226–231. AAAI Press, 1996.
- 784
- Guohua Fu, Jianwei Zheng, Islam Abudayyeh,
 785 Chizobam Ani, Cyril Rakowski, Louis Ehwerhe-
 786 muepha, Hanna Lu, Yongjuan Guo, Shenglin Liu,
 787 Huimin Chu, and Bing Yang. Cardiogpt: An ecg
 788 interpretation generation model. *IEEE Access*, PP:
 789 1–1, 01 2024. doi: 10.1109/ACCESS.2024.3384349.
- 790
- Philip Gage. A new algorithm for data compres-
 791 sion. *The C Users Journal archive*, 12:23–38,
 792 1994. URL <https://api.semanticscholar.org/CorpusID:59804030>.
- 793
- Bryan Gopal, Ryan W. Han, Gautham Raghupathi,
 795 Andrew Y. Ng, Geoffrey H. Tison, and Pranav Ra-
 796 jpurkar. 3kg: Contrastive learning of 12-lead elec-
 797 trocardiograms using physiologically-inspired aug-
 798 mentations, 2021. URL <https://arxiv.org/abs/2106.04452>.
- 799
- 800

801	Brian Gow, Tom Pollard, Larry A Nathanson,	de Oliveira, Madeline Muzzi, Mahesh Pasupuleti,	851
802	Alistair Johnson, Benjamin Moody, Chrystinne	Mannat Singh, Manohar Paluri, Marcin Kardas,	852
803	Fernandes, Nathaniel Greenbaum, Jonathan W	Maria Tsimpoukelli, Mathew Oldham, Mathieu	853
804	Waks, Parastou Eslami, Tanner Carbonati,	Rita, Maya Pavlova, Melanie Kambadur, Mike	854
805	Ashish Chaudhari, Elizabeth Herbst, Dana	Lewis, Min Si, Mitesh Kumar Singh, Mona Has-	855
806	Moukheiber, Seth Berkowitz, Roger Mark, and	san, Naman Goyal, Narjes Torabi, Nikolay Bash-	856
807	Steven Horng. Mimic-iv-ecg: Diagnostic electro-	lykov, Nikolay Bogoychev, Niladri Chatterji, Ning	857
808	cardiogram matched subset, 2023. URL https://physionet.org/content/mimic-iv-ecg/1.0/ .	Zhang, Olivier Duchenne, Onur Çelebi, Patrick Al-	858
809		rassy, Pengchuan Zhang, Pengwei Li, Petar Va-	859
810	Aaron Grattafiori, Abhimanyu Dubey, Abhinav	sic, Peter Weng, Prajjwal Bhargava, Pratik Dubal,	860
811	Jauhri, Abhinav Pandey, Abhishek Kadian, Ah-	Praveen Krishnan, Punit Singh Koura, Puxin Xu,	861
812	mad Al-Dahle, Aiesha Letman, Akhil Mathur,	Qing He, Qingxiao Dong, Ragavan Srinivasan,	862
813	Alan Schelten, Alex Vaughan, Amy Yang, An-	Raj Ganapathy, Ramon Calderer, Ricardo Silveira	863
814	gela Fan, Anirudh Goyal, Anthony Hartshorn,	Cabral, Robert Stojnic, Roberta Raileanu, Ro-	864
815	Aobo Yang, Archi Mitra, Archie Sravankumar,	han Maheswari, Rohit Girdhar, Rohit Patel, Ro-	865
816	Artem Korenev, Arthur Hinsvark, Arun Rao, As-	main Sauvestre, Ronnie Polidoro, Roshan Sum-	866
817	ton Zhang, Aurelien Rodriguez, Austen Gregerson,	baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang,	867
818	Ava Spataru, Baptiste Roziere, Bethany Biron,	Saghar Hosseini, Sahana Chennabasappa, Sanjay	868
819	Binh Tang, Bobbie Chern, Charlotte Caucheteux,	Singh, Sean Bell, Seohyun Sonia Kim, Sergey	869
820	Chaya Nayak, Chloe Bi, Chris Marra, Chris	Edunov, Shaoliang Nie, Sharan Narang, Sharath	870
821	McConnell, Christian Keller, Christophe Touret,	Raparthy, Sheng Shen, Shengye Wan, Shruti Bhos-	871
822	Chunyang Wu, Corinne Wong, Cristian Canton	ale, Shun Zhang, Simon Vandenhende, Soumya	872
823	Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel	Batra, Spencer Whitman, Sten Sootla, Stephane	873
824	Song, Danielle Pintz, Danny Livshits, Danny Wy-	Collot, Suchin Gururangan, Sydney Borodinsky,	874
825	att, David Esiobu, Dhruv Choudhary, Dhruv Ma-	Tamar Herman, Tara Fowler, Tarek Sheasha,	875
826	hajan, Diego Garcia-Olano, Diego Perino, Dieuwke	Thomas Georgiou, Thomas Scialom, Tobias Speck-	876
827	Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina	bacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn,	877
828	Lobanova, Emily Dinan, Eric Michael Smith,	Vedanuj Goswami, Vibhor Gupta, Vignesh Ra-	878
829	Filip Radenovic, Francisco Guzmán, Frank Zhang,	manathan, Viktor Kerkez, Vincent Gonguet, Vir-	879
830	Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis	ginie Do, Vish Vojeti, Vítor Albiero, Vladan	880
831	Anderson, Govind Thattai, Graeme Nail, Gre-	Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin	881
832	goire Mialon, Guan Pang, Guillem Cucurell, Hai-	Fu, Whitney Meers, Xavier Martinet, Xiaodong	882
833	ley Nguyen, Hannah Korevaar, Hu Xu, Hugo	Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide	883
834	Touvron, Iliyan Zarov, Imanol Arrieta Ibarra,	Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang,	884
835	Isabel Kloumann, Ishan Misra, Ivan Evtimov,	Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei,	885
836	Jack Zhang, Jade Copet, Jaewon Lee, Jan Gef-	Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yun-	886
837	fert, Jana Vranes, Jason Park, Jay Mahadeokar,	ing Mao, Zacharie Delpierre Coudert, Zheng Yan,	887
838	Jeet Shah, Jelmer van der Linde, Jennifer Bil-	Zhengxing Chen, Zoe Papakipos, Aaditya Singh,	888
839	lock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng	Aayushi Srivastava, Abha Jain, Adam Kelsey,	889
840	Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao	Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria,	890
841	Yu, Joanna Bitton, Joe Spisak, Jongsoo Park,	Ahuva Goldstand, Ajay Menon, Ajay Sharma,	891
842	Joseph Rocca, Joshua Johnstun, Joshua Saxe,	Alex Boesenberg, Alexei Baevski, Allie Feinstein,	892
843	Junteng Jia, Kalyan Vasuden Alwala, Karthik	Amanda Kallet, Amit Sangani, Amos Teo, Anam	893
844	Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li,	Yunus, Andrei Lupu, Andres Alvarado, Andrew	894
845	Kenneth Heafield, Kevin Stone, Khalid El-Arini,	Caples, Andrew Gu, Andrew Ho, Andrew Poulton,	895
846	Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal	Andrew Ryan, Ankit Ramchandani, Annie	896
847	Bhalla, Kushal Lakhota, Lauren Rantala-Yeary,	Dong, Annie Franco, Anuj Goyal, Aparajita Saraf,	897
848	Laurens van der Maaten, Lawrence Chen, Liang	Arkabandhu Chowdhury, Ashley Gabriel, Ash-	898
849	Tan, Liz Jenkins, Louis Martin, Lovish Madaan,	win Bharambe, Assaf Eisenman, Azadeh Yazdan,	899
850	Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke	Beau James, Ben Maurer, Benjamin Leonhardi,	900
		Bernie Huang, Beth Loyd, Beto De Paola, Bhar-	901

902	gavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden	953
903	Hancock, Bram Wasti, Brandon Spence, Brani Sto-	954
904	jkovic, Brian Gamido, Britt Montalvo, Carl Parker,	955
905	Carly Burton, Catalina Mejia, Ce Liu, Chang-	956
906	han Wang, Changkyu Kim, Chao Zhou, Chester	957
907	Hu, Ching-Hsiang Chu, Chris Cai, Chris Tin-	958
908	dal, Christoph Feichtenhofer, Cynthia Gao, Da-	959
909	mon Civin, Dana Beaty, Daniel Kreymer, Daniel	960
910	Li, David Adkins, David Xu, Davide Testuggine,	961
911	Delia David, Devi Parikh, Diana Liskovich, Di-	962
912	dem Foss, Dingkang Wang, Duc Le, Dustin Hol-	963
913	land, Edward Dowling, Eissa Jamil, Elaine Mont-	964
914	gomery, Eleonora Presani, Emily Hahn, Emily	965
915	Wood, Eric-Tuan Le, Erik Brinkman, Esteban Ar-	966
916	caute, Evan Dunbar, Evan Smothers, Fei Sun,	967
917	Felix Kreuk, Feng Tian, Filippos Kokkinos, Fi-	968
918	rat Ozgenel, Francesco Caggioni, Frank Kanayet,	969
919	Frank Seide, Gabriela Medina Florez, Gabriella	970
920	Schwarz, Gada Badeer, Georgia Swee, Gil Halpern,	971
921	Grant Herman, Grigory Sizov, Guangyi, Zhang,	972
922	Guna Lakshminarayanan, Hakan Inan, Hamid	973
923	Shojanazeri, Han Zou, Hannah Wang, Hanwen	974
924	Zha, Haroun Habeeb, Harrison Rudolph, Helen	975
925	Suk, Henry Aspegren, Hunter Goldman, Hongyuan	976
926	Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tu-	977
927	fanov, Ilias Leontiadis, Irina-Elena Veliche, Itai	978
928	Gat, Jake Weissman, James Geboski, James Kohli,	979
929	Janice Lam, Japhet Asher, Jean-Baptiste Gaya,	980
930	Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen,	981
931	Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong,	982
932	Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill,	983
933	Jon Shepard, Jonathan McPhie, Jonathan Torres,	984
934	Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou	985
935	U, Karan Saxena, Kartikay Khandelwal, Katayoun	986
936	Zand, Kathy Matosich, Kaushik Veeraraghavan,	987
937	Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun	988
938	Huang, Kunal Chawla, Kyle Huang, Lailin Chen,	989
939	Lakshya Garg, Lavender A, Leandro Silva, Lee	990
940	Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron	991
941	Moshkovich, Luca Wehrstedt, Madian Khabsa,	992
942	Manav Avalani, Manish Bhatt, Martynas Mankus,	
943	Matan Hasson, Matthew Lennie, Matthias Reso,	
944	Maxim Groshev, Maxim Naumov, Maya Lathi,	
945	Meghan Keneally, Miao Liu, Michael L. Seltzer,	
946	Michal Valko, Michelle Restrepo, Mihir Patel,	
947	Mik Vyatskov, Mikayel Samvelyan, Mike Clark,	
948	Mike Macey, Mike Wang, Miquel Jubert Her-	
949	moso, Mo Metanat, Mohammad Rastegari, Mun-	
950	ish Bansal, Nandhini Santhanam, Natascha Parks,	
951	Natalsha White, Navyata Bawa, Nayan Singh, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Niko-	
952	lay Pavlovich Laptev, Ning Dong, Norman Cheng,	
	Oleg Chernoguz, Olivia Hart, Omkar Salpekar,	
	Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul	
	Saab, Pavan Balaji, Pedro Rittner, Philip Bon-	
	trager, Pierre Roux, Piotr Dollar, Polina Zvyag-	
	ina, Prashant Ratanchandani, Pritish Yuvraj, Qian	
	Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub,	
	Raghatham Murthy, Raghu Nayani, Rahul Mi-	
	tra, Rangaprabhu Parthasarathy, Raymond Li, Re-	
	bekkah Hogan, Robin Battey, Rocky Wang, Russ	
	Howes, Ruty Rinott, Sachin Mehta, Sachin Siby,	
	Sai Jayesh Bondu, Samyak Datta, Sara Chugh,	
	Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru	
	Pan, Saurabh Mahajan, Saurabh Verma, Seiji	
	Yamamoto, Sharadh Ramaswamy, Shaun Lind-	
	say, Shaun Lindsay, Sheng Feng, Shenghao Lin,	
	Shengxin Cindy Zha, Shishir Patil, Shiva Shankar,	
	Shuqiang Zhang, Shuqiang Zhang, Sinong Wang,	
	Sneha Agarwal, Soji Sajuyigbe, Soumith Chin-	
	tala, Stephanie Max, Stephen Chen, Steve Ke-	
	hoe, Steve Satterfield, Sudarshan Govindaprasad,	
	Sumit Gupta, Summer Deng, Sungmin Cho, Sunny	
	Virk, Suraj Subramanian, Sy Choudhury, Sydney	
	Goldman, Tal Remez, Tamar Glaser, Tamara Best,	
	Thilo Koehler, Thomas Robinson, Tianhe Li, Tian-	
	jun Zhang, Tim Matthews, Timothy Chou, Tzook	
	Shaked, Varun Vontimitta, Victoria Ajayi, Vic-	
	toria Montanez, Vijai Mohan, Vinay Satish Ku-	
	mar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru,	
	Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li,	
	Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will	
	Constable, Xiaocheng Tang, Xiaoqian Wu, Xiaolan	
	Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman,	
	Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin	
	Zhang, Ying Zhang, Yossi Adi, Youngjin Nam,	
	Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian,	
	Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito,	
	Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei	
	Zhao, and Zhiyu Ma. The llama 3 herd of models,	
	2024. URL https://arxiv.org/abs/2407.21783 .	
	Liangke Gui, Yingshan Chang, Qiuyuan Huang, Sub-	993
	hojit Som, Alex Hauptmann, Jianfeng Gao, and	994
	Yonatan Bisk. Training vision-language transfor-	995
	mers from captions, 2023. URL https://arxiv.org/abs/2205.09256 .	996
		997
	William Jongwon Han, Diana Gomez, Avi Alok,	998
	Chaojing Duan, Michael A. Rosenberg, Douglas	999
	Weber, Emerson Liu, and Ding Zhao. Interpreta-	1000
	tion of intracardiac electrograms through textual	1001

- 1002 representations, 2024. URL <https://arxiv.org/abs/2402.01115>. 1048
- 1003
- 1004 Awni Y. Hannun, Pranav Rajpurkar, Masoumeh 1049
1005 Haghpanahi, Geoffrey H. Tison, Codie Bourn, 1049
1006 Mintu P. Turakhia, and Andrew Y. Ng. 1050
1007 Cardiologist-level arrhythmia detection and 1050
1008 classification in ambulatory electrocardiograms 1051
1009 using a deep neural network. *Nature Medicine*, 25: 1051
1010 65–69, 01 2019. doi: 10.1038/s41591-018-0268-3. 1051
1011 URL <https://www.nature.com/articles/s41591-018-0268-3>. 1051
- 1012
- 1013 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian 1052
1014 Sun. Deep residual learning for image recognition, 1052
1015 2015. URL <https://arxiv.org/abs/1512.03385>. 1053
- 1016 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan 1054
1017 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and 1054
1018 Weizhu Chen. Lora: Low-rank adaptation of large 1055
1019 language models, 2021. URL <https://arxiv.org/abs/2106.09685>. 1056
- 1020
- 1021 Sarthak Jain and Byron C. Wallace. Attention is not 1057
1022 explanation, 2019. 1057
- 1023 Alistair E. W. Johnson, Lucas Bulgarelli, Lu Shen, 1058
1024 Alvin Gayles, Ayad Shammout, Steven Horng, 1058
1025 Tom J. Pollard, Benjamin Moody, Brian Gow, 1059
1026 Li-wei H. Lehman, Leo A. Celi, and Roger G. 1060
1027 Mark. Mimic-iv, a freely accessible electronic 1061
1028 health record dataset. *Scientific Data*, 10, 01 2023. 1062
1029 doi: 10.1038/s41597-022-01899-x. 1063
- 1030 Mark Johnson. Counties most in need 1064
1031 of cardiologists are the most likely to 1064
1032 have none, 07 2024. URL <https://www.washingtonpost.com/science/2024/07/29/cardiologists-rural-counties-shortage/>. 1065
- 1033
- 1034
- 1035 Diederik P. Kingma and Jimmy Ba. Adam: A 1066
1036 method for stochastic optimization, 2017. 1066
- 1037 Dani Kiyasseh, Tingting Zhu, and David A. Clifton. 1067
1038 CloCS: Contrastive learning of cardiac signals 1067
1039 across space, time, and patients, 2021. URL 1068
1040 <https://arxiv.org/abs/2005.13249>. 1068
- 1041 Mykola Klymenko, Sam M Doesburg, George 1069
1042 Medvedev, Pengcheng Xi, Urs Ribary, and Vasily A 1069
1043 Vakorin. Byte-pair encoding for classifying 1070
1044 routine clinical electroencephalograms in adults over 1070
1045 the lifespan. *IEEE Journal of Biomedical and 1071
1046 Health Informatics*, pages 1–11, 01 2023. doi: 10.1109/jbhi.2023.3236264. 1071
- 1047
- Taku Kudo and John Richardson. Sentencepiece: 1048
A simple and language independent subword 1049
tokenizer and detokenizer for neural text processing, 1049
2018. URL <https://arxiv.org/abs/1808.06226>. 1050
- Ohhwan Kwon, Jinwoo Jeong, Hyung Bin Kim, In Ho 1051
Kwon, Song Yi Park, Ji Eun Kim, and Yuri Choi. 1051
Electrocardiogram sampling frequency range ac- 1052
ceptable for heart rate variability analysis. *Health- 1052
care Informatics Research*, 24:198, 2018. doi: 10. 1053
4258/hir.2018.24.3.198. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6085204/>. 1053
- Chin-Yew Lin. Rouge: A package for automatic eval- 1054
uation of summaries. In *ACL 2004*, 2004. 1054
- Che Liu, Zhongwei Wan, Cheng Ouyang, Anand 1055
Shah, Wenjia Bai, and Rossella Arcucci. Zero-shot 1055
ecg classification with multimodal learning and 1056
test-time clinical knowledge enhancement, 2024. 1056
URL <https://arxiv.org/abs/2403.06659>. 1057
- Haotian Liu, Chunyuan Li, Qingyang Wu, and 1058
Yong Jae Lee. Visual instruction tuning, 2023. 1058
- URL <https://arxiv.org/abs/2304.08485>. 1059
- J MacQueen. Some methods for classification and 1060
analysis of multivariate observations. *Project Eu- 1060
clid*, 5.1:281–298, 1967. 1061
- Harold Martin, Ulyana Morar, Walter Izquierdo, 1062
Mercedes Cabrerizo, Anastasio Cabrera, and 1062
Malek Adjouadi. Real-time frequency-independent 1063
single-lead and single-beat myocardial infarction 1063
detection. *Artificial intelligence in medicine*, 121: 1064
102179, 2021. 1065
- Kaden McKeen, Laura Oliva, Sameer Masood, Au- 1066
gustin Toma, Barry Rubin, and Bo Wang. Ecg- 1066
fm: An open electrocardiogram foundation model, 1067
2024. URL <https://arxiv.org/abs/2408.05178>. 1068
- Yeongyeon Na, Minje Park, Yunwon Tae, and 1069
Sunghoon Joo. Guiding masked representation 1069
learning to capture spatio-temporal relationship of 1070
electrocardiogram, 2024. URL <https://arxiv.org/abs/2402.09450>. 1070
- Naoki Nonaka and Jun Seita. Electrocardiogram 1071
classification by modified efficientnet with data 1071
augmentation. In *2020 Computing in Cardiology*, 1072
pages 1–4. IEEE, 2020. 1072

- 1091 Jungwoo Oh, Hyunseung Chung, Joon myoung 1138
 1092 Kwon, Dong gyun Hong, and Edward Choi. Lead- 1139
 1093 agnostic self-supervised learning for local and 1140
 1094 global representations of electrocardiogram, 2022. 1141
 1095 URL <https://arxiv.org/abs/2203.06889>. 1142
 1096 1143
- 1096 Jungwoo Oh, Gyubok Lee, Seongsu Bae, Joon myoung 1144
 1097 Kwon, and Edward Choi. Ecg-qa: A comprehensive 1145
 1098 question answering dataset combined with 1146
 1099 electrocardiogram, 2023. URL <https://arxiv.org/abs/2306.15681>. 1147
 1100 1148
- 1101 World Health Organization. Cardiovascular diseases, 2024. URL https://www.who.int/health-topics/cardiovascular-diseases#tab=tab_1. 1149
- 1105 Kishore Papineni, Salim Roukos, Todd Ward, and 1150
 1106 Wei-Jing Zhu. Bleu: a method for automatic eval- 1151
 1107 uation of machine translation. In *ACL*, 2002. 1152
- 1108 Manh Pham, Aaqib Saeed, and Dong Ma. C-melt: 1153
 1109 Contrastive enhanced masked auto-encoders for 1154
 1110 ecg-language pre-training, 2024. URL <https://arxiv.org/abs/2410.02131>.
- 1112 Jielin Qiu, William Han, Jiacheng Zhu, Mengdi Xu, 1155
 1113 Michael Rosenberg, Emerson Liu, Douglas Weber, 1156
 1114 and Ding Zhao. Transfer knowledge from natu- 1157
 1115 ral language to electrocardiography: Can we de- 1158
 1116 tect cardiovascular disease through language mod- 1159
 1117 els? In Andreas Vlachos and Isabelle Augen- 1160
 1118 stein, editors, *Findings of the Association for Com- 1161
 1119 putational Linguistics: EACL 2023*, pages 442– 1162
 1120 453, Dubrovnik, Croatia, May 2023a. Associa- 1163
 1121 tion for Computational Linguistics. doi: 10. 1164
 1122 18653/v1/2023.findings-eacl.33. URL <https://aclanthology.org/2023.findings-eacl.33>. 1165
- 1124 Jielin Qiu, Jiacheng Zhu, Shiqi Liu, William Han, 1166
 1125 Jingqi Zhang, Chaojing Duan, Michael A. Rosen- 1167
 1126 berg, Emerson Liu, Douglas Weber, and Ding 1168
 1127 Zhao. Automated cardiovascular record retrieval 1169
 1128 by multimodal learning between electrocardiogram 1170
 1129 and clinical report. In Stefan Hegselmann, Anto- 1171
 1130 nio Parziale, Divya Shanmugam, Shengpu Tang, 1172
 1131 Mercy Nyamewaa Asiedu, Serina Chang, Tom 1173
 1132 Hartvigsen, and Harvineet Singh, editors, *Pro- 1174
 1133 ceedings of the 3rd Machine Learning for Health 1175
 1134 Symposium*, volume 225 of *Proceedings of Machine 1176
 1135 Learning Research*, pages 480–497. PMLR, 10 Dec 1177
 1136 2023b. URL <https://proceedings.mlr.press/v225/qiu23a.html>. 1178
- 1137 1179
- Alec Radford, Jeffrey Wu, Rewon Child, David 1138
 Luan, Dario Amodei, and Ilya Sutskever. 1139
 Language models are unsupervised multitask 1140
 learners, 2019. URL https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf. 1141
 1142 1143
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya 1144
 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish 1145
 Sastry, Amanda Askell, Pamela Mishkin, Jack 1146
 Clark, Gretchen Krueger, and Ilya Sutskever. 1147
 Learning transferable visual models from natural 1148
 language supervision. In *ICML*, 2021. 1149
- Pranav Rajpurkar, Awni Y. Hannun, Masoumeh 1150
 Haghpanahi, Codie Bourn, and Andrew Y. Ng. 1151
 Cardiologist-level arrhythmia detection with con- 1152
 volutional neural networks, 2017. URL <https://arxiv.org/abs/1707.01836>. 1153
 1154
- Peter J. Rousseeuw. Silhouettes: A graphical 1155
 aid to the interpretation and validation of clus- 1156
 ter analysis. *Journal of Computational and Ap- 1157
 plied Mathematics*, 20:53–65, 1987. ISSN 0377- 1158
 0427. doi: [https://doi.org/10.1016/0377-0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7). 1159
 URL <https://www.sciencedirect.com/science/article/pii/0377042787901257>. 1160
 1161
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 1162
 Neural machine translation of rare words with sub- 1163
 word units, 2016. URL <https://arxiv.org/abs/1508.07909>. 1164
 1165
- Yuchen Shen and Barnabás Póczos. Graphbpe: 1166
 Molecular graphs meet byte-pair encoding, 2024. 1167
 URL <https://arxiv.org/abs/2407.19039>. 1168
- Nils Strodtthoff, Patrick Wagner, Tobias Schaeffter, 1169
 and Wojciech Samek. Deep learning for ecg anal- 1170
 ysis: Benchmarks and insights from ptb-xl. *IEEE 1171
 Journal of Biomedical and Health Informatics*, 25: 1172
 1519–1528, 2021. 1173
- Saeedeh Tahery, Fatemeh Hamid Akhlaghi, Termeh 1174
 Amirsoleimani, and Saeed Farzi. Heartbert: A 1175
 self-supervised ecg embedding model for efficient 1176
 and effective medical signal analysis, 2024. URL 1177
<https://arxiv.org/abs/2411.11896>. 1178
- Jialu Tang, Tong Xia, Yuan Lu, Cecilia Mascolo, and 1179
 Aaqib Saeed. Electrocardiogram-language model 1180
 for few-shot question answering with meta learn- 1181
 ing, 2024a. URL <https://arxiv.org/abs/2410.14464>. 1182
 1183

1184	Jialu Tang, Tong Xia, Yuan Lu, Cecilia Mascolo, and Aaqib Saeed. Electrocardiogram report generation and question answering via retrieval-augmented self-supervised modeling, 2024b. URL https://arxiv.org/abs/2409.08788 .	1234
1185		1235
1186		1236
1187		1237
1188		1238
1189	Nazgol Tavabi and Kristina Lerman. Pattern discov- ery in time series with byte pair encoding, 2021. URL https://arxiv.org/abs/2106.00614 .	1239
1190		1240
1191		1241
1192	Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussonot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Cas- bon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieil- lard, Piotr Stanczyk, Sertan Girgin, Nikola Mom- chev, Matt Hoffman, Shantanu Thakoor, Jean- Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Co- enen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijayku- mar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Mor- eira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak- Plucińska, Harleen Batra, Harsh Dhand, Ivan Nar- dini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiran- bir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lau- ren Usui, Laurent Sifre, Lena Heuermann, Leti- cia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofiz Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan,	1242
1193		1243
1194		1244
1195		1245
1196		1246
1197		1247
1198		1248
1199		1249
1200		1250
1201		1251
1202		1252
1203		1253
1204		1254
1205		1255
1206		1256
1207		1257
1208		1258
1209	URL https://arxiv.org/abs/2408.00118 .	1259
1210		
1211		
1212		
1213		
1214		
1215		
1216		
1217		
1218		
1219		
1220		
1221		
1222		
1223		
1224		
1225		
1226		
1227		
1228		
1229		
1230		
1231		
1232		
1233		
1234	Akhil Vaid, Joy Jiang, Ashwin Sawant, Stamatios Lerakis, Edgar Argulian, Yuri Ahuja, Joshua Lam- pert, Alexander Charney, Hayit Greenspan, Ben- jamin Glicksberg, Jagat Narula, and Girish Nad- karni. Heartbeit: Vision transformer for electro- cardiogram data improves diagnostic performance at low sample sizes, 2022.	1260
1235		1261
1236		1262
1237		1263
1238		1264
1239		1265
1240		1266
1241		
1242		
1243		
1244		
1245		
1246		
1247		
1248		
1249		
1250		
1251		
1252		
1253		
1254		
1255		
1256		
1257		
1258		
1259		
1260	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023.	1267
1261		1268
1262		1269
1263		1270
1264		
1265		
1266		
1267		
1268		
1269		
1270		
1271	Patrick Wagner, Nils Strodthoff, Ralf-Dieter Bous- seljot, Dieter Kreiseler, Fatima I. Lunze, Wo- jciech Samek, and Tobias Schaeffter. PTB- XL, a large publicly available electrocardiogra- phy dataset. <i>Scientific Data</i> , 7(1):154, May 2020. ISSN 2052-4463. doi: 10.1038/ s41597-020-0495-6. URL https://www.nature.com/articles/s41597-020-0495-6 . Number: 1	1271
1272		1272
1273		1273
1274		1274
1275		1275
1276		1276
1277		1277
1278		1278
1279	Publisher: Nature Publishing Group.	1279
1280	Zhongwei Wan, Che Liu, Xin Wang, Chaofan Tao, Hui Shen, Zhenwu Peng, Jie Fu, Rossella Arcucci, Huaxiu Yao, and Mi Zhang. Meit: Multi-modal	1280
1281		1281
1282		1282

- 1283 electrocardiogram instruction tuning on large lan- 1329
 1284 guage models for report generation, 2024. URL 1330
 1285 <https://arxiv.org/abs/2403.04945>. 1331
 1286 Sarah Wiegreffe and Yuval Pinter. Attention is not 1332
 1287 not explanation, 2019. 1333
- 1288 Svante Wold, Kim Esbensen, and Paul Geladi. 1334
 1289 Principal component analysis. *Chemomet- 1335
 1290 rics and Intelligent Laboratory Systems*, 2
 1291 (1):37–52, 1987. ISSN 0169-7439. doi:
 1292 [https://doi.org/10.1016/0169-7439\(87\)80084-9](https://doi.org/10.1016/0169-7439(87)80084-9).
 1293 URL [https://www.sciencedirect.com/
 1294 science/article/pii/0169743987800849](https://www.sciencedirect.com/science/article/pii/0169743987800849). Proceedings
 1295 of the Multivariate Statistical Workshop
 1296 for Geologists and Geochemists.
- 1297 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien 1334
 1298 Chaumond, Clement Delangue, Anthony Moi,
 1299 Pierrick Cistac, Tim Rault, Rémi Louf, Morgan
 1300 Funtowicz, Joe Davison, Sam Shleifer, Patrick von
 1301 Platen, Clara Ma, Yacine Jernite, Julien Plu, Can-
 1302 wen Xu, Teven Le Scao, Sylvain Gugger, Mariama
 1303 Drame, Quentin Lhoest, and Alexander M. Rush.
 1304 Huggingface’s transformers: State-of-the-art natu-
 1305 ral language processing, 2020.
- 1306 Sangmin Woo, Donguk Kim, Jaehyuk Jang, Yubin 1334
 1307 Choi, and Changick Kim. Don’t miss the for-
 1308 est for the trees: Attentional vision calibration for
 1309 large vision language models, 2024. URL <https://arxiv.org/abs/2405.17820>.
- 1311 Shunxiang Yang, Cheng Lian, and Zhigang Zeng. 1334
 1312 Masked autoencoder for ecg representation learn-
 1313 ing. In *2022 12th International Conference on In-*
1314 formation Science and Technology (ICIST), pages
 1315 95–98, 2022. doi: 10.1109/ICIST55546.2022.
 1316 9926900.
- 1317 Susan Zhang, Stephen Roller, Naman Goyal, Mikel 1334
 1318 Artetxe, Moya Chen, Shuhui Chen, Christopher
 1319 Dewan, Mona Diab, Xian Li, Xi Victoria Lin,
 1320 Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt
 1321 Shuster, Daniel Simig, Punit Singh Koura, Anjali
 1322 Sridhar, Tianlu Wang, and Luke Zettlemoyer. Opt:
 1323 Open pre-trained transformer language models,
 1324 2022. URL <https://arxiv.org/abs/2205.01068>.
- 1325 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. 1334
 1326 Weinberger, and Yoav Artzi. Bertscore: Evaluating
 1327 text generation with bert. *ArXiv*, abs/1904.09675,
 1328 2020.
- Yubao Zhao, Tian Zhang, Xu Wang, Puyu Han,
 Tong Chen, Linlin Huang, Youzhu Jin, and Ji-
 aju Kang. Ecg-chat: A large ecg-language model
 for cardiac disease diagnosis, 2024. URL <https://arxiv.org/abs/2408.08849>.

Appendix A. ECG-Byte

A.1. Additional Pseudocode for ECG-Byte

Algorithm 3: Merging a pair in an ID array

Input: Array of IDs *ids*, Pair to merge *pair* as (u_1, u_2) , New ID *new_id*

Output: Merged vector of IDs

Function *merge(ids, pair, new_id):*

```

  Initialize empty vector new_ids with capacity
  of ids;
  i ← 0;
  while i < length of ids do
    if i < length of ids - 1 and
      (ids[i], ids[i+1]) = pair then
        new_ids.append(new_id);
        i ← i + 2;
    end
    else
      new_ids.append(ids[i]);
      i ← i + 1;
    end
  end
  return new_ids;
```

Algorithm 4: Calculating Frequency of Byte Pairs in an Array

Input: Array of IDs *ids*

Output: HashMap of pairs and their frequencies

Function *get_stats(ids):*

```

  pair_counts ← Parallel fold operation;;
  foreach window of size 2 in ids do
    Let  $(u_1, u_2) \leftarrow$  elements of the
    window;
    Increment the count of  $(u_1, u_2)$  in
    local pair_count;
  end
  pair_counts ← Parallel reduce operation to
  combine local pair_count HashMaps;
  return pair_counts;
```

1336 **A.2. Mapping between Token and ECG**

1337 We add more examples of the mapping between the
 1338 ECG signal and the encoded tokens for **ECG-Byte**
 1339 in Figures 5 and 6.

1340 **A.3. Token usage and length distribution for
 1341 varying num_merges**

1342 We add more examples of the token usage and length
 1343 distributions for varying `num_merges` in Figure 7.

1344 **A.4. Attention Visualizations**

1345 We add more visualizations of the attention weights
 1346 in Figures 8, 9, 10, 11, 12, 13, 14, 15.

1347 **Appendix B. 2-stage Pretraining
 1348 Approaches**

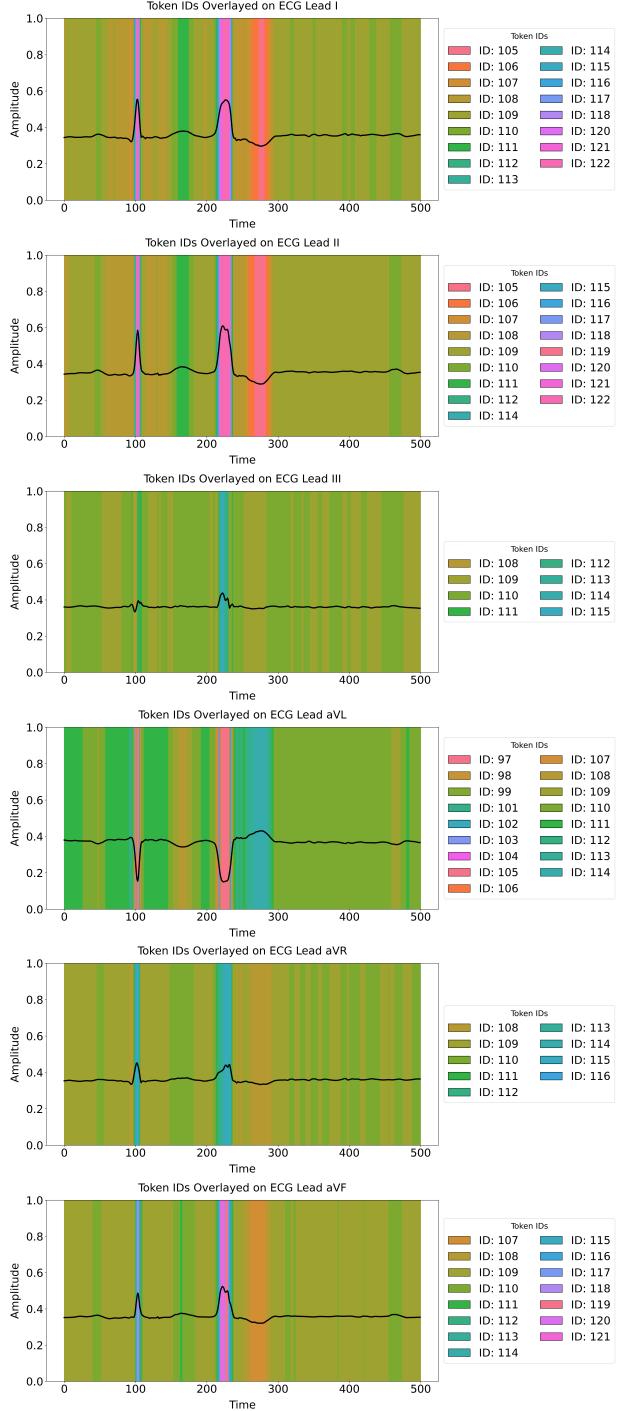
1349 To be consistent, we normalize each ECG in the
 1350 same manner as described in subsection 3.2. Consider a dataset of N ECG-image and clinical note
 1351 pairs, denoted as $\{(I_i, O_i)\}_{i=1}^N$, where: $I_i \in \mathbb{R}^{3 \times C \times T}$
 1352 is the i -th normalized and replicated ECG im-
 1353 age, obtained by stacking the clipped ECG sig-
 1354 na- X_{clipped} along the channel dimension: $I_i =$
 1355 $\text{stack}(X_{\text{clipped}}, X_{\text{clipped}}, X_{\text{clipped}})$. The reason we do
 1356 this is because we need to create RGB images to use
 1357 pretrained image models like ViT (Dosovitskiy et al.,
 1358 2021) and CLIP (Radford et al., 2021).
 1359

1360 O_i is the corresponding clinical note for the i -th
 1361 ECG, serving as the textual description. Note that
 1362 O_i differs from S in the autoregressive setup, where
 1363 S represents the tokenized answer sequence provided
 1364 by either ECG-QA (Oh et al., 2023) or MIMIC-IV
 1365 ECG pretraining (Zhao et al., 2024).

1366 Given these two features I and O we then de-
 1367 scribe the contrastive, masked, and dual approaches
 1368 implemented for our baselines that are derived from
 1369 commonly used techniques used throughout previous
 1370 works (Oh et al., 2022; Choi et al., 2023; McKeen
 1371 et al., 2024; Pham et al., 2024; Tang et al., 2024a,b;
 1372 Vaid et al., 2022).

1373 **B.1. Contrastive learning approaches**

1374 We utilize a pretrained CLIP Radford et al. (2021)
 1375 checkpoint, namely ‘openai/clip-vit-base-patch32’,
 1376 provided by HuggingFace (Wolf et al., 2020) to en-
 1377 code ECG signals I and text labels O into a shared
 1378 embedding space. Let $f_{\text{img}} : \mathbb{R}^{3 \times C \times T} \rightarrow \mathbb{R}^d$ and



1375 Figure 5: A mapping between tokens used for a given
 1376 ECG Leads I, II, III, aVL, aVR, aVF.

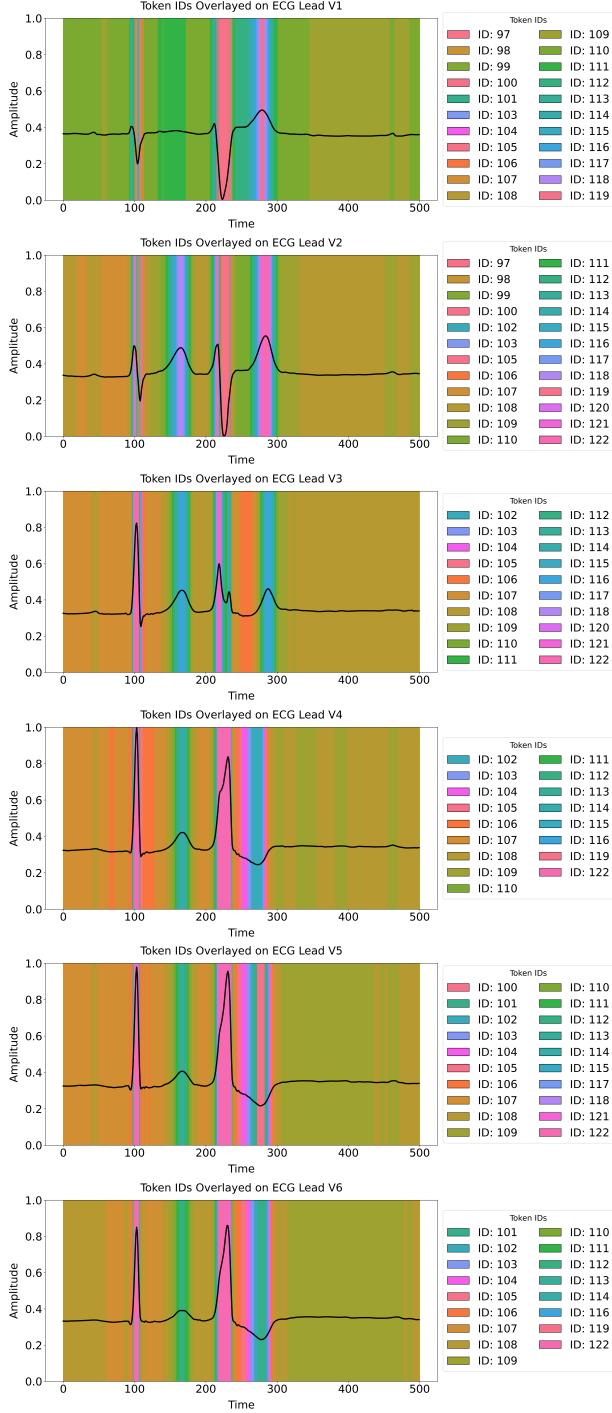


Figure 6: A mapping between tokens used for a given ECG Leads V1, V2, V3, V4, V5, V6.

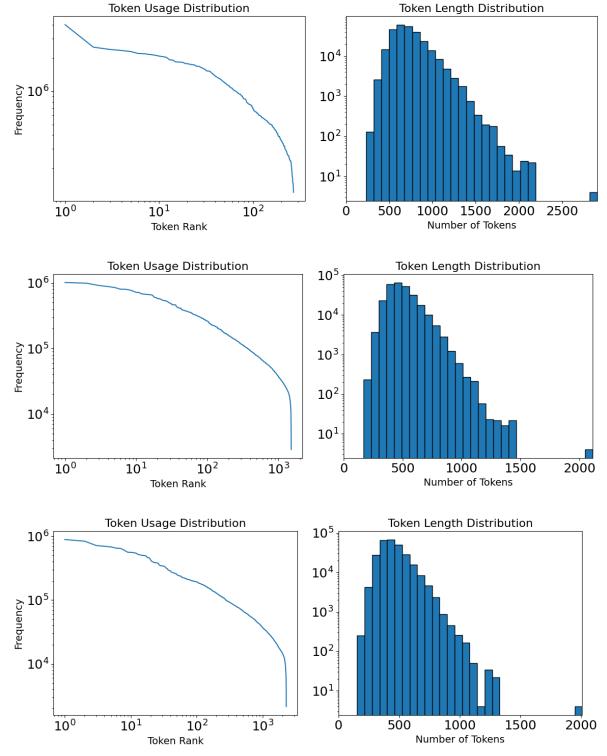


Figure 7: Plots of the token usage and length distributions for **ECG-BYTE** where `num_merges` is 500, 1750, and 2500 from top to bottom.

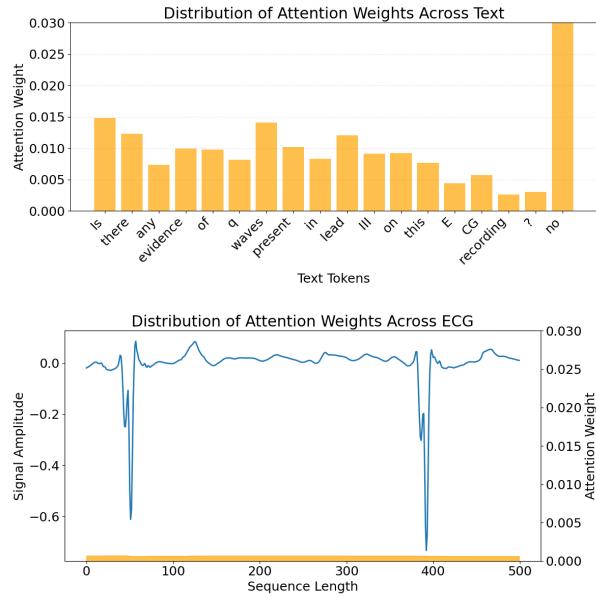


Figure 8: The attention weight overlayed on both text (top) and ECG (bottom).

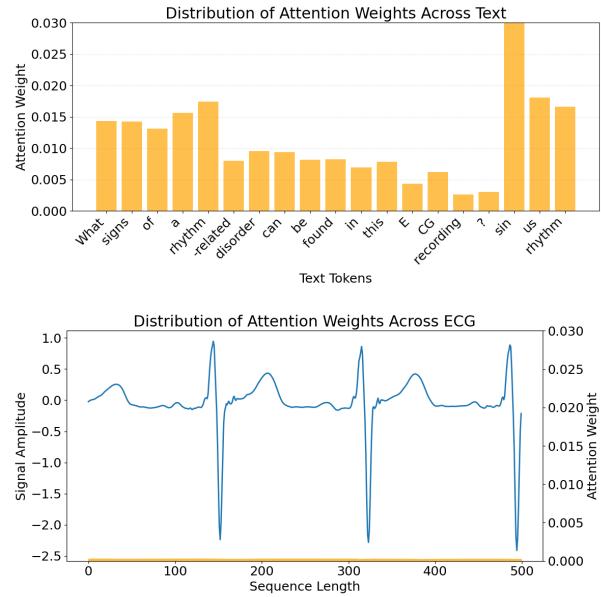


Figure 10: The attention weight overlayed on both text (top) and ECG (bottom).

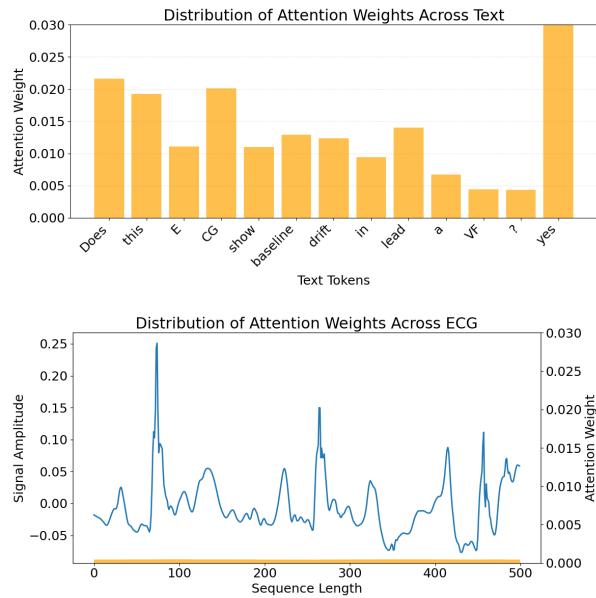


Figure 9: The attention weight overlayed on both text (top) and ECG (bottom).

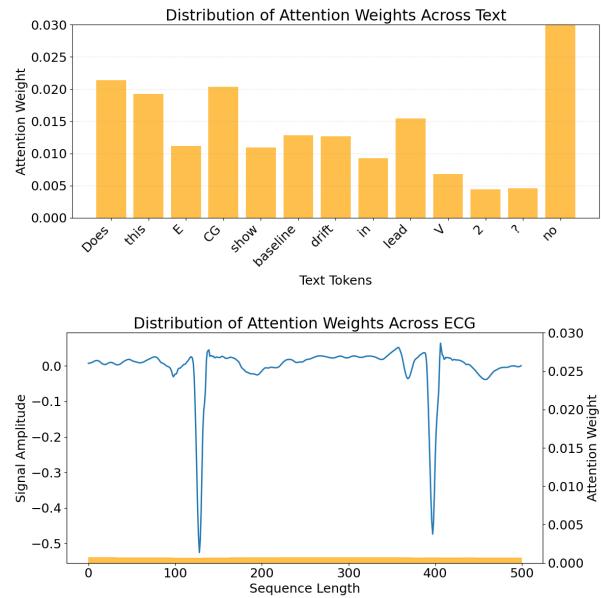


Figure 11: The attention weight overlayed on both text (top) and ECG (bottom).

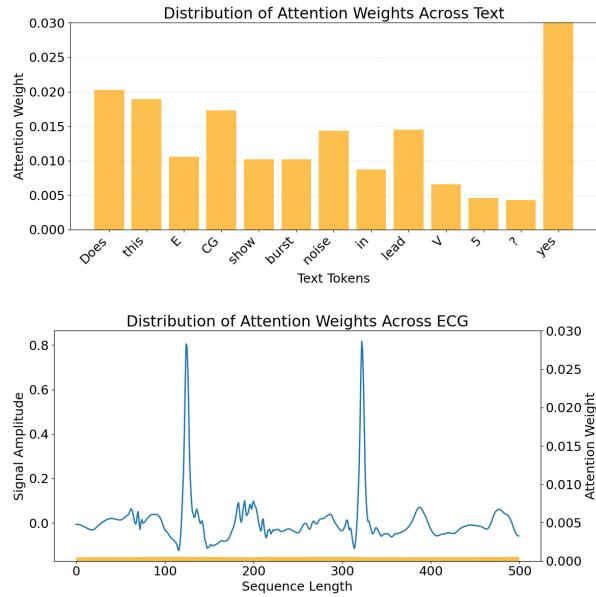


Figure 12: The attention weight overlayed on both text (top) and ECG (bottom).

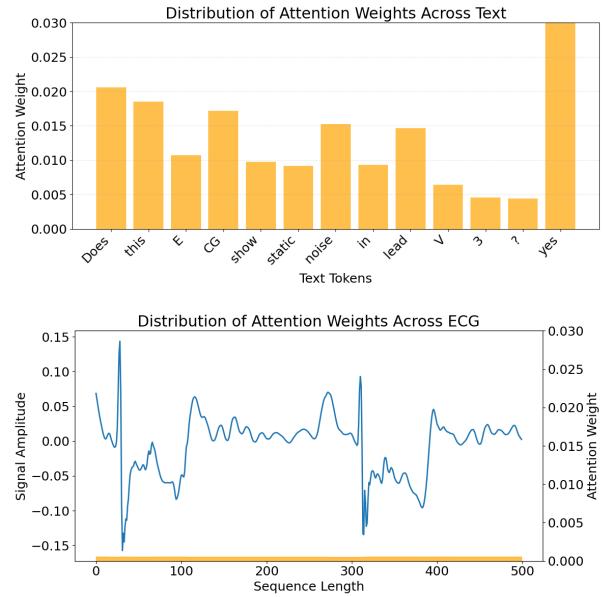


Figure 14: The attention weight overlayed on both text (top) and ECG (bottom).

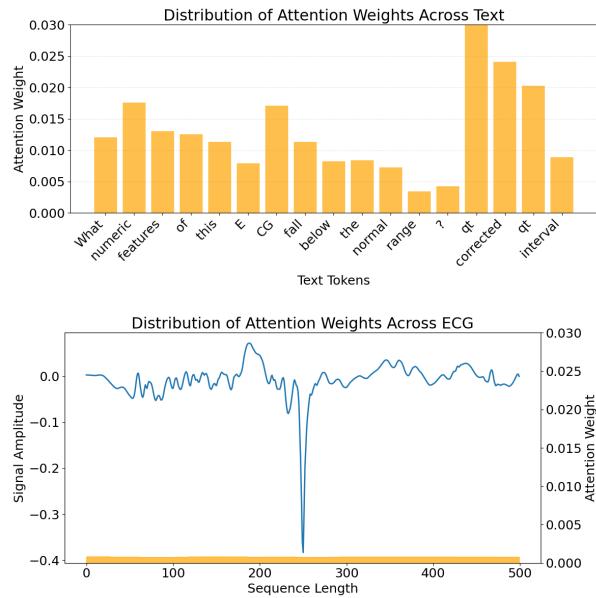


Figure 13: The attention weight overlayed on both text (top) and ECG (bottom).

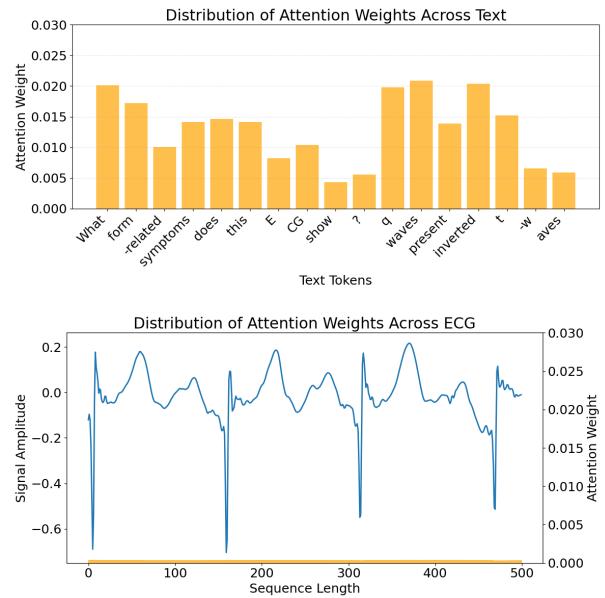


Figure 15: The attention weight overlayed on both text (top) and ECG (bottom).

¹³⁷⁹ f_{txt} : Text $\rightarrow \mathbb{R}^d$ be the image and text encoders
¹³⁸⁰ of the pretrained CLIP model, respectively. The em-
¹³⁸¹ beddings for the i -th pair are computed as:

$$z_i^{\text{img}} = f_{\text{img}}(I_i), \quad z_i^{\text{txt}} = f_{\text{txt}}(O_i),$$

¹³⁸² where $z_i^{\text{img}}, z_i^{\text{txt}} \in \mathbb{R}^d$. The CLIP loss function $\mathcal{L}_{\text{CLIP}}$
¹³⁸³ aligns the embeddings of corresponding ECG signals
¹³⁸⁴ and text labels while contrasting them with non-
¹³⁸⁵ matching pairs. This is formulated as:

$$\begin{aligned} \mathcal{L}_{\text{CL}} = & -\frac{1}{N} \sum_{i=1}^N \left[\log \frac{\exp(\text{sim}(z_i^{\text{img}}, z_i^{\text{txt}})/\tau)}{\sum_{j=1}^N \exp(\text{sim}(z_i^{\text{img}}, z_j^{\text{txt}})/\tau)} \right. \\ & \left. + \log \frac{\exp(\text{sim}(z_i^{\text{txt}}, z_i^{\text{img}})/\tau)}{\sum_{j=1}^N \exp(\text{sim}(z_i^{\text{txt}}, z_j^{\text{img}})/\tau)} \right] \end{aligned}$$

¹³⁸⁶ where $\text{sim}(\cdot, \cdot)$ denotes cosine similarity, and τ is a
¹³⁸⁷ learnable temperature parameter.

¹³⁸⁸ To integrate the pretrained CLIP model into our
¹³⁸⁹ language model for joint reasoning over ECG signals
¹³⁹⁰ and text, we project the frozen image embeddings
¹³⁹¹ z_i^{img} into the language model's hidden space. Let
¹³⁹² $W \in \mathbb{R}^{h \times d}$ be a learnable projection matrix, where h
¹³⁹³ is the hidden dimension of the language model. The
¹³⁹⁴ projected embeddings are:

$$z_i^{\text{clip}} = W z_i^{\text{img}}.$$

¹³⁹⁵ These projected embeddings z_i^{clip} are then
¹³⁹⁶ prepended to the token embeddings of the
¹³⁹⁷ language model, where we get $\text{Context} =$
¹³⁹⁸ $\{[\text{BOS}], [\text{SIG_START}], z_i^{\text{clip}}, [\text{SIG_END}], Q\}$ to train
¹³⁹⁹ the same autoregressive objective, L_{NLL} .

¹⁴⁰⁰ B.2. Masked image modeling approaches

¹⁴⁰¹ Consider the normalized ECG image $I \in \mathbb{R}^{3 \times C \times T}$
¹⁴⁰² obtained as previously described. We utilize a pre-
¹⁴⁰³ trained Vision Transformer (ViT) model ([Dosovitskiy et al., 2021](#)), specifically the ‘google/vit-base-
¹⁴⁰⁴ patch16-224-in21k’ checkpoint provided by Hugging-
¹⁴⁰⁵ Face ([Wolf et al., 2020](#)).

¹⁴⁰⁷ The image I is partitioned into P non-overlapping
¹⁴⁰⁸ patches. Let N be the number of images in our
¹⁴⁰⁹ dataset, and I_i denote the i -th image. The ViT en-
¹⁴¹⁰ coder f_{vit} projects these patches into latent embed-
¹⁴¹¹ dings:

$$z_i^{\text{patch}} = f_{\text{vit}}(I_i) \in \mathbb{R}^{P \times d},$$

¹⁴¹² where d is the embedding dimension of the ViT
¹⁴¹³ model.

During training, we randomly mask a subset of
¹⁴¹⁴ patches for each image I_i , creating a binary mask
¹⁴¹⁵ $M_i \in \{0, 1\}^P$, where $M_{i,j} = 1$ if patch j is masked
¹⁴¹⁶ and $M_{i,j} = 0$ otherwise. The masked embeddings
¹⁴¹⁷ z_i^{masked} are formed by replacing the embeddings of
¹⁴¹⁸ masked patches with a mask token. A reconstruction
¹⁴¹⁹ head f_{rec} is then applied to predict the pixel-level
¹⁴²⁰ content of the masked patches:
¹⁴²¹

$$\hat{I}_i = f_{\text{rec}}(z_i^{\text{masked}}) \in \mathbb{R}^{P \times d}.$$

The masked image modeling loss \mathcal{L}_{MIM} is computed
¹⁴²² as the mean squared error (MSE) between the recon-
¹⁴²³ structed embeddings \hat{I}_i and the original embeddings
¹⁴²⁴ z_i^{patch} at the masked positions:
¹⁴²⁵

$$\mathcal{L}_{\text{MIM}} = \frac{1}{N} \sum_{i=1}^N \frac{1}{\sum_{j=1}^P M_{i,j}} \sum_{j=1}^P M_{i,j} \left\| \hat{I}_i[j] - z_i^{\text{patch}}[j] \right\|_2^2. \quad (7)$$

To integrate the MIM representations into the lan-
¹⁴²⁶ guage model for joint reasoning over ECG signals and
¹⁴²⁷ textual questions, we project the frozen ViT embed-
¹⁴²⁸ dings $z_i^{\text{img}} \in \mathbb{R}^d$ into the language model's hidden
¹⁴²⁹ space. Let $W \in \mathbb{R}^{h \times d}$ be a learnable projection ma-
¹⁴³⁰ trix, where h is the hidden dimension of the language
¹⁴³¹ model. The projected embeddings are given by:
¹⁴³²

$$z_i^{\text{vit}} = W z_i^{\text{img}}.$$

These projected embeddings z_i^{vit} are then prepended
¹⁴³³ to the language model's token embeddings, to get
¹⁴³⁴ $\text{Context} = \{[\text{BOS}], [\text{SIG_START}], z_i^{\text{vit}}, [\text{SIG_END}], Q\}$
¹⁴³⁵ to train the same autoregressive objective, L_{NLL} ,
¹⁴³⁶ mentioned previously.
¹⁴³⁷

¹⁴³⁸ B.3. Dual approaches

The dual approach follows the previous two con-
¹⁴³⁹ trastive and masked image modeling approaches for
¹⁴⁴⁰ pretraining the ECG encoder but simply just com-
¹⁴⁴¹ bines the losses like so:
¹⁴⁴²

$$\mathcal{L}_{\text{Dual}} = \lambda_1 \mathcal{L}_{\text{MIM}} + \lambda_2 \mathcal{L}_{\text{CL}}$$

where $\lambda_1 = \lambda_2 = 1$ in our study.
¹⁴⁴³

However, when training the autoregressive LLM,
¹⁴⁴⁴ we project both embeddings, z_i^{vit} and z_i^{clip} , outputted
¹⁴⁴⁵ by their respective frozen encoders via a learnable
¹⁴⁴⁶ projection matrix into the language model's hidden
¹⁴⁴⁷ space of dimension h . We then concatenate the
¹⁴⁴⁸ projected embeddings and pass them through a fu-
¹⁴⁴⁹ sion network to obtain the fused visual embedding
¹⁴⁵⁰

1451 $z_i^{\text{fused}} \in \mathbb{R}^h$:

$$z_i^{\text{fused}} = f_{\text{fusion}}(\text{concat}(z_i^{\text{vit}}; z_i^{\text{clip}})),$$

1452 where f_{fusion} is a trainable feedforward net-
 1453 work. The fused visual embedding z_i^{fused}
 1454 is prepended to the token embeddings of
 1455 the language model, forming $\text{Context} =$
 1456 $\{[\text{BOS}], [\text{SIG_START}], z_i^{\text{fused}}, [\text{SIG_END}], Q\}$ to
 1457 train the autoregressive objective, L_{NLL} .

1458 Appendix C. Additional Results

1459 C.1. Does Larger LLMs Yield Higher 1460 Performance?

1461 We present the results of ablating the size of the LLM
 1462 in Table 8. Interestingly, the performance across the
 1463 three different model sizes (1B, 3B, 8B) remains fairly
 1464 similar. We believe that the limited dataset size pre-
 1465 vents the larger models from realizing their full per-
 1466 formance potential. We hypothesize that increasing
 1467 the amount of training data would enable the larger
 1468 models to leverage their greater capacity, resulting in
 1469 observable performance improvements.

Table 8: Ablation study on how larger LLMs perform for NLG.

LLM	BLEU-4	Rouge-L	Meteor	BertScore F1
Llama 3.2 1B (Grattafiori et al., 2024)	13.93 ± 0.21	47.08 ± 0.56	29.17 ± 0.31	92.53 ± 0.07
Llama 3.2 3B (Grattafiori et al., 2024)	14.80 ± 0.17	46.55 ± 0.21	29.53 ± 0.16	92.42 ± 0.01
Llama 3.1 8B (Grattafiori et al., 2024)	13.80 ± 0.16	46.29 ± 0.25	28.56 ± 0.11	92.44 ± 0.05

1470 C.2. Qualitative NLG Examples

1471 We provide qualitative NLG examples of successful
 1472 (Figure 17) and unsuccessful generations (Figure 16).

Ground Truth Question	Which diagnostic symptom does this ECG show, incomplete left bundle branch block or incomplete right bundle branch block, excluding uncertain symptoms?	Does the qrs duration shown on this ECG fall within the normal range?	What form-related traits are exhibited by this ECG in lead I?	In lead V2, what form-related features does this ECG display?	What direction is this ECG deviated to?
Ground Truth Answer	incomplete right bundle branch block	yes	low amplitude t-wave	q waves present inverted t-waves	extreme axis deviation
Generated Answer	incomplete left bundle branch block	no	non-specific st depression	none	left axis deviation
Ground Truth Question	Within which numeric range does the qt interval of this ECG fall, above the normal range or within the normal range	Which diagnostic symptom does this ECG show, subendocardial injury in anterolateral leads or subendocardial injury in inferolateral leads, excluding uncertain symptoms?	Which diagnostic symptom does this ECG show, myocardial infarction in inferoposterolateral leads or myocardial infarction in anterolateral leads, excluding uncertain symptoms?	What form-related symptoms does this ECG show in lead II?	What diagnostic symptoms does this ECG show, excluding uncertain symptoms?
Ground Truth Answer	none	none	myocardial infarction in anterolateral leads	high qrs voltage	myocardial infarction in anteroseptal leads non-diagnostic t abnormalities
Generated Answer	qt interval	subendocardial injury in anterolateral leads	myocardial infarction in inferoposterolateral leads	non-specific st depression	myocardial infarction in anteroseptal leads

Figure 16: Randomly sampled NLG results of unsuccessful generations on the PTB-XL test set from ECG-QA.

Ground Truth Question	Is atrial fibrillation detectable from this ECG?	Which diagnostic symptom does this ECG show, left posterior fascicular block or subendocardial injury in lateral leads, including uncertain symptoms?	Which diagnostic symptom does this ECG show, subendocardial injury in lateral leads or incomplete left bundle branch block, including uncertain symptoms?	What is the diagnostic symptom that can be identified from this ECG, excluding any symptoms that are unclear, left atrial overload/enlargement or myocardial infarction in anterolateral leads?	What are the leads on the ECG that are manifesting static noise?
Ground Truth Answer	yes	none	subendocardial injury in lateral leads	left atrial overload/enlargement	lead I lead II lead III lead aVR lead aVL lead aVF lead V1 lead V2 lead V3 lead V4 lead V5 lead V6
Generated Answer	yes	none	subendocardial injury in lateral leads	left atrial overload/enlargement	lead I lead II lead III lead aVR lead aVL lead aVF lead V1 lead V2 lead V3 lead V4 lead V5 lead V6
Ground Truth Question	Does this ECG reveal any signs of sinus bradycardia?	Are there any noises detected in lead aVF on this ECG?	What numeric features of this ECG fall below the normal range?	What types of noises are displayed in lead aVL in this ECG waveform?	By excluding uncertain symptoms, which diagnostic symptom is apparent in this ECG, ischemic in inferior leads or left anterior fascicular block?
Ground Truth Answer	no	no	pr interval qt corrected qt interval	baseline drift static noise	left anterior fascicular block
Generated Answer	no	no	pr interval qt corrected qt interval	baseline drift static noise	left anterior fascicular block

Figure 17: Randomly sampled NLG results of successful generations on the PTB-XL test set from ECG-QA.