EMGLLM: Data-to-Text Alignment for Electromyogram Diagnosis Generation with Medical Numerical Data Encoding

Anonymous ACL submission

Abstract

Electromyography (EMG) tables are crucial for diagnosing muscle and nerve disorders, and advancing the automation of EMG diagnostics is significant for improving medical efficiency. EMG tables contain extensive continuous numerical data, which current Large Language Models (LLMs) often struggle to interpret effectively. To address this issue, we propose EMGLLM, a data-to-text model specifically designed for medical inerination tables. EMGLLM employs the EMG Alignment Encoder to simulate the process that doctors compare test values with reference values, aligning the data into word embeddings that reflect health degree. Additionally, we construct ETM, a dataset comprising 17,276 real cases and their corresponding diagnostic results, to support medical data-to-text tasks. Experimental results on ETM demonstrate that EMGLLM outperforms various baseline models in understanding EMG tables and generating highquality diagnoses, which represents an effective paradigm for automatic diagnosis generation from medical examination table.

1 Introduction

004

012

016

017

037

041

Electromyography (EMG) refers to the pattern of electrophysiological signal concomitant with musculations recorded with an electromyograph (Ni et al., 2020), which plays a significant role in evaluating human activities (Cooray et al., 2022; Smedemark-Margulies et al., 2023; Rakhmatulin, 2024). In medicine, the EMG is one of the major diagnostic tools for identifying and characterizing motor unit disorders (Daube, 2002), which is commonly used to examine nerve and muscle excitability and conduction functions, thereby determining the functional status of peripheral nerves, neurons, neuromuscular junctions, and the muscles themselves. After the EMG examination, the physicians perform a two-step analysis based on the records of the electrical signals. They first analyze the waveforms, converting the complex electrical signals into easily interpretable data tables, which contain essential information for medical diagnosis, such as amplitude, conduction velocity, and latency. Subsequently, by completing quantitative analysis, the doctors interpret the converted table data to render their final diagnosis and form a diagnostic report (Boon et al., 2008). In this paper, we focus on the data-to-text task of automatic diagnosis generation from EMG tabular data. 042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

Figure 1 shows an anonymized EMG diagnosis including two parts, Findings and Impression. In the context of an EMG examination, Findings refer to observations of tables, aiming to objectively describe the phenomena reflected by the data, thus facilitating further analysis by the physician. To accurately identify various neuromuscular disorders within tabular data and translate them into Findings, physicians must possess a deep understanding of the distinct patterns associated with neuromuscular junction disorders, radiculopathies, upper motor neuron lesions, and so on. In terms of Impression, it consists of two aspects: a summary and interpretation of the test results, as well as an analysis of the clinical significance of the Findings, which may include diagnostic suggestions or potential issues. Therefore, Impression requires a certain level of clinical experience from doctors. (Katirji, 2002) Basically, EMG diagnosis writing can be error-prone and tedious for underexperienced physicians, and onerous and timeconsuming for experienced physicians. Therefore, considering the powerful reasoning and text generation capabilities of large language models (LLMs) in the medical field (Fan et al., 2024), we are motivated to explore methods for using LLMs to process examination tables and automatically generate medical EMG diagnoses.

The automatic generation of EMG diagnoses involves two major challenges:

084

- 0
- 094 095 096

100

103

104

105

106

107

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

- **Reference value comparison:** This task requires analyzing from the relative size of EMG test values compared to their corresponding reference values to assess the degree of health or abnormality of the test items. Moreover, differences in equipment, environment, and other factors across hospitals may result in varying optimal reference values, which increases the complexity of analysis.
- Intensive numerical data input: It is necessary for EMG diagnosis to quantitatively understand of medical examination tables containing large amounts of continuous numerical data to generate diagnostic results. For LLMs, directly handling numerical data may present certain difficulties (Golkar et al., 2023) since LLMs are not well-adept at comparing numerical values and quantitatively diagnosing the normality of these values, which may lead to errors.

To address this, we proposes EMGLLM, a novel data-to-text framework for automatic diagnosis from medical examination tables, which introduces EMG Alignment Encoder specialized in encoding continuous numerical data in EMG examination tables. The EMG Alignment Encoder can compare the test values with reference values, encoding them into virtual tokens that represent the degree of abnormality, and aligning numerical data to diagnostic text. This allows the LLM to better understand EMG tables, thereby generating more accurate and comprehensive diagnoses. Our main contributions include:

- For automatic diagnosis generation from medical examination tables, we propose a datato-text model, EMGLLM, which includes an EMG Alignment Encoder designed to encode continuous numerical values and enhance data understanding.
- We construct a dataset ETM comprising about 17,000 real EMG tables with their diagnoses annotated by authoritative physicians, which can provide support for researches on automatic diagnosis generation.

127 Compared to all baseline methods, EMG di128 agnoses generated by EMGLLM demonstrates
129 higher quality in all evaluation metrics, fully prov130 ing the effectiveness and robustness of EMGLLM.
131 This method can also be applied to other medical
132 examination tables.

2 Related Work

2.1 Data-to-text Generation

Data-to-text is a significant branch of natural language processing (Sharma et al., 2024). Its goal is to transform complex numerical data and tables into textual descriptions, assisting users in understanding and analyzing data, thereby improving the efficiency of data analysis. Data has the characteristics of complex structure and information density, and many studies have proposed methods to address this challenge. For example, splicing nearest neighbors (Wiseman et al., 2021) is an effective data-to-text policy by inserting or replacing text segments directly from neighbor sourcetarget pairs to construct generations. Search-andlearning method (Jolly et al., 2021) is aimed at enhancing semantic coverage in few-shot data-totext generation. Recently, some research applied LLMs to complete data-to-text. MURMUR (Saha et al., 2022) and TAT-LLM (Zhu et al., 2024) respectively enhanced data-to-text generation capabilities through multi-step and discrete reasoning frameworks. TableLLaMA (Zhang et al., 2024a) and TableLLM (Zhang et al., 2024b) were implemented supervised fine-tuning on table datasets for proficiently handling tabular data.

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

In the medical fields, data-to-text generation also holds vast application prospects. For instance, language models can complete automatic drug description generation from medical information tables (Yermakov et al., 2021) and diagnosis from examination tables (Gu et al., 2020; Guo et al., 2024).

2.2 EMG Diagnosis

EMG has a wide range of applications in medical diagnosis (Gaso et al., 2021; Nguyen et al., 2023; Li et al., 2023). EMG signals can be used to construct an end-to-end sleep stage neural classification model for diagnosing sleep disorders (Chambon et al., 2017). They can also be classified by Markov model (Bureau et al., 2021) for diagnosing potential neuropathies. Specifically, a dataset MIME (Gu et al., 2020) for EMG table tasks is used to train models such as hierarchical transformers. The model and dataset are both closed source and the method is relatively simple, which cannot fully adapt to complex data-to-text tasks. Therefore, this paper aims to explore automatic diagnostic generation based on EMG tables.

				Electro	mvor	rram (EMG	- -				
被检肌肉 Examined Mu	scle	竏籔 Fibs	正儀 PSW	東顧 Fascics	其他	MUPs ars MUP	ッ 体相 Polyph	MUP形态 MUP Form	募集 Rec	6相 urit	
左指总伸肌 L Ext.Dig.Com		- 1+		-	-	-	-		干扰相 Interference		
左 第一背側骨 L Dors.Int.I	向肌	-	-	-	-	-		-	混合 Mix	混合相 Mix	
			Nerv	e Condu	ctior	n Velocity (!	NCV)				
收检神经 xamined Nerv	项目 e Type	刺激 Stim	记录 Rec	潜伏期 Lat (L)	(左)	潜伏期 (右) Lat(R)	波幅 (左) Amp (L)	波幅 (右) Amp (R)	速度 (左) CV (L)	速度 (左) CV (R)	
E中神经 fedian	运动 Motor	腕 Wrist	拇短展肌 APE	4.2			7.0		7.0		
と神经 Unar	感觉 Sensory	中指 Dig III	腕 Wrist	3.4		3.8	15.0	14.2	46.2	44.8	
非总神经 eroncal	F波 F-wave	踝 Ankle	趾短伸肌 FDB	50.5							
Findings: EMG: 被 (No dene NCV: 左 神经传导 (Mildly p left Med sensory n	检肌未 ervation 侧正中 速度和 prolonge ian Ner erves ar erves ar	见明显则 or reinna 神经运动 波幅正行 d motor ve is rev e within 损害	U源性或补 ervation in 力传导潜化 常范围。 distal late vealed. Th the norma	申经源性 the mus 大期正常 ency and ne condu il range.)	损害 cles 上限 sligl	肌电改变. examined.) ,感觉神绪 htly slowed n velocity a	各传导速度 sensory n nd amplif	轻度减慢; erve condua ude of resi	余运动; ction velo idual mot	和感觉 city of or and	

Figure 1: An EMG diagnostic report example, including EMG tables (EMG and NCV) and their corresponding diagnosis (*Findings* and *Impression*). For our automatic diagosis generation task, the input is the EMG tables and the output is the diagnosis.

3 Method

182

183

185

190

191

194

195

196

199

206

3.1 Model

EMGLLM is composed of two fundamental components: the EMG Alignment Encoder and the LLM. The EMG Alignment Encoder is a specialized module tailored for understanding medical examination tables such as EMG tables. As illustrated in Figure 2, when an EMG table is input, text and discrete data are tokenized and vectorized by LLM's tokenizer and embedder directly. For continuous numerical data, they are encoded into virtual tokens using the EMG Alignment Encoder. The model's output is the *Findings* and *Impression* of the EMG tables.

The process by which the EMG Alignment Encoder analyzes continuous numerical table cells is analogous to the approach employed by doctors. In EMG examinations, reference values are the most critical criterion for determining whether a test parameter is within normal ranges. The reference range defines the upper and lower limits of normalcy for a specific examination item. The extent to which the test value exceeds the reference range reflects the degree of pathological alteration in the muscle or nerve. In practice, doctors first assess



Figure 2: EMGLLM Framework. Medical examination tables contain a large amount of continuous data. The numbers marked in blue in the instruction could be encoded by the EMG Alignment Encoder.

the relative magnitude of test values based on reference values, then make a annotation within the table cell to indicate the degree of abnormality. The EMG Alignment Encoder is designed to emulate this process by comparing the continuous test data with multiple reference ranges and encoding the abnormality degree semantically into virtual tokens that are more interpretable by the LLM.

Reference Value Acquisition In practice, doctors rely on their clinical experience to make appropriate adjustments to reference values for certain individual cases. This process involves strong subjectivity. Therefore, we propose a method for mining reference values based on percentiles from the training dataset. We filter out the completely healthy cases without any abnormalities from the training dataset and statistically analyze

207

208

294

295

297

298

299



Figure 3: EMG Alignment Encoder Structure. EMG Alignment Encoder regards the reference matrix of the test values as an image, extracts features using the convolutional layer, and aligns continuous data with text through the Attention mechanism.

the healthy case subset for each item. For a given examination item *i*, we use the *k* upper percentiles $u_i^{p_1}, u_i^{p_2}, \dots, u_i^{p_k}$ and the *k* lower percentiles $l_i^{p_1}, l_i^{p_2}, \dots, l_i^{p_k}$ as multiple reference values at different levels, where p_1, p_2, \dots, p_k represent different percentile thresholds. These percentiles allow us to estimate the boundaries of the reference ranges from data distribution of healthy individuals.

224

225

226

229

230

234

236

237

238

241

242

245

246

247

EMG Alignment Encoder Structure The input to the EMG Alignment Encoder for item i is a reference matrix X_i :

$$X_{i} = \begin{pmatrix} u_{i}^{p_{1}} & u_{i}^{p_{2}} & \cdots & u_{i}^{p_{k}} \\ x_{i} & x_{i} & \cdots & x_{i} \\ l_{i}^{p_{1}} & l_{i}^{p_{2}} & \cdots & l_{i}^{p_{k}} \end{pmatrix}$$
(1)

where x_i denotes the continuous test value. The EMG Alignment Encoder views the reference matrix X_i as a form of image, where the pixels represent the arrangement of the examined value and reference ranges, as illustrated in Figure 3. Using a convolutional layer Conv with d_C output channels, the model sequentially compares the test value with the reference values. Subsequently, a linear layer integrates the output vectors of Convto produce the vector \hat{X}_i representing the feature of the test value x_i :

$$\hat{X}_i = f_1(W_1 \text{flatten}(Conv(X_i)) + b_1) \quad (2)$$

249 where f_1 and b_1 are the activation function and bias. 250 When *Conv* outputs *m* vectors of dimension d_C , $W_1 \in \mathbf{R}^{m \times N}$, where N represents the number of vectors output by the linear layer. Consequently, \hat{X}_i contains N vectors of dimension d_C .

These data features are then aligned with the word embeddings in the LLM's vocabulary. As shown in Figure 3, the alignment process first involves learning a set of text prototypes $E' \in$ $\mathbf{R}^{V' \times D}$ from the LLM's vocabulary $E \in \mathbf{R}^{V \times D}$ through $E' = W_2 E$, where V and V' refers to the size of vocabulary and text prototypes respectively subject to $V' \ll V$, D denotes dimension of the LLM embeddings, and $W_2 \in \mathbf{R}^{V \times V'}$. Text prototypes E' serve as a compressed version of the vocabulary capable of semantically implying health or abnormality in medical diagnosis, such as "prolonged", "slowed", and "decreased". The EMG Alignment Encoder then connects the continuous data features X_i in Equation 2 to these text prototypes via a multi-head attention layer:

$$EMGA lignment Encoder(X_i)$$

= MultiHeadAttention(Q_i, K, V, n_{head}) (3)

where $Q_i = \hat{X}_i W_Q$, $K = E'W_K$, $V = E'W_V$, n_{head} is the number of heads, $W_Q \in \mathbf{R}^{d_C \times d}$. $W_K, W_V \in \mathbf{R}^{D \times d}$, $d = \lfloor d_C/n_{head} \rfloor$. The output of EMG Alignment Encoder is N data embeddings of dimension D. In Equation 3, Query is computed from the continuous data in tables, while the Key and Value are derived from the LLM embeddings. The EMG Alignment Encoder leverages this Attention mechanism to associate continuous data with text.

The additional reference value information and reasonable continuous data encoding contribute to enhancing the performance of LLMs in data-totext medical tasks. Another advantage of the EMG Alignment Encoder lies in its continuous function, where similar numeric values are encoded into correspondingly similar virtual tokens. In contrast, standard LLMs tokenize numeric values, a process that discretizes the table's data. For instance, two numerically close values, such as 9.99 and 10.0, may result in significantly different word embeddings in LLMs, which may be not reasonable in data-to-text scenario.

Finally, EMGLLM integrates the EMG Alignment Encoder with the LLM. With the assistance of the EMG Alignment Encoder, the LLM gains better understanding of the EMG table. Combined with the LLM's strong generative capabilities, this



Figure 4: Examples of two pre-training tasks for EMG Alignment Encoder. [Embedding] represents the virtual tokens encoded by the EMG Alignment Encoder from a reference matrix X_i in Equation 1, which is required to enable pre-trained LLMs to complete single data diagnosis without fine-tuning EMG data.

enhancement endows EMGLLM with better automated diagnostic abilities.

3.2 Training

300

307

310

311

312

313

315

316

EMG Alignment Encoder Pre-training Before training on EMG diagnosis task, we pre-train the EMG Alignment Encoder to ensure it can properly perform data understanding in a single continuous test value. The purpose of pre-training is to help the model understand the underlying medical semantics behind the relative size relationship between a test value with its reference values. Freezing the LLM, two types of pre-training tasks based on one test value are applied: (1) Classification of abnormality. (2) Making LLM generate a diagnostic description of a word. The loss function in pre-training is same as the supervised fine-tuning of LLM.

Figure 4 presents examples of the pre-training 317 data. The instructions for pre-training tasks require 318 EMG Alignment Encoder to provide reasonable virtual tokens so that the base LLM can clearly understand their meaning. In pre-training dataset con-321 struction, the test value x_i and the reference val-322 ues u_i and l_i can be obtained by sampling from the 323 diagnosis generation training dataset, and the output labels can be constructed directly from manu-325 ally defined rules. For example, if a test value exceeds the $u_i^{0.02}$ by 20%, the virtual tokens should 327 convey the meaning of "significantly high". This 329 rule-based approach does not rely on any authoritative reference values from hospitals, but can naturally learn an understanding of reference values from the data distribution of healthy individuals, which has good generality. 333

Measurement	Value
# of Samples	17,276
Avg # of Continous Numerical Data	33.14
Avg Length (Findings)	85.04
Avg Length (Impression)	26.82

Table 1: Dataset Statistics

Model Fine-tuning Upon the completion of pretraining, we proceed with supervised training for the EMG data-to-text task. In the fine-tuning phase, we further train both the LLM and the EMG Alignment Encoder on EMG train dataset, where LLM is efficiently trained by the Low-Rank Adaptation (LoRA) (Hu et al., 2021) method. 334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

352

353

354

355

356

357

358

359

361

362

363

364

365

366

367

369

370

371

372

373

Through the aforementioned steps, the EMG Alignment Encoder's representation of continuous numerical data can be enhanced, enabling the EMGLLM to better understand continuous data and diagnose from EMG tables.

4 Experiments

4.1 ETM Dataset

In this section, we introduce a high quality EMG diagnostic report dataset ETM (Electromyogram Table Mart) derived from Huashan Hospital Affiliated to Fudan University¹ with high authenticity, accuracy, and authority, which contains a total of 17,276 diagnostic reports from 2006 to 2013, and each data includes:

• Basic information of real anonymized patients (age, gender, and height).

• EMG tables (EMG and NCV tests) from the real EMG examination in the hospital.

• Diagnosis (*Findings* and *Impression*) personally written by experienced physicians.

The data format is shown in Figure 1. The full dataset is further proportionally divided into training, validation, and testing set, with data volumes of 13820, 1728, and 1728 respectively, which can effectively support medical data-to-text research.

Some statistical information of the ETM dataset is displayed in Table 1 basic statistics for our Some statistics information ETM dataset. The average number of continuous numerical data in tables is 33.14, indicating that the model's input contains dense numerical information. The automatic diagnostic task requires the model to have a sufficient understanding of continuous test values.

¹https://www.huashan.org.cn/

	Model	Automatic							Model	
	model	ROUGE-1	ROUGE-2	ROUGE-L	BLEU-1	BLEU-2	BLEU-3	Correctness	Completeness	Human
	DeepSeek-V2 (zero-shot)	60.19(0.88)	42.91(1.25)	51.53(1.10)	49.67(0.87)	42.22(1.04)	37.16(1.16)	3.44	3.61	3.41
Overall	Lattice	71.59(0.94)	56.82(1.25)	65.25(0.96)	55.97(0.82)	46.57(0.85)	39.41(0.83)	3.46	3.28	3.50
Overall	TableLLM-7B	74.41(0.42)	58.03(0.80)	67.37(0.55)	65.48(0.56)	58.31(0.67)	53.22(0.79)	3.72	3.70	3.88
	Chinese-Alpaca-2-7B-16K	79.24(0.33)	65.15(0.70)	73.18(0.45)	71.26(0.50)	65.18(0.50)	60.67(0.74)	4.02	3.92	4.21
	EMGLLM (Ours)	80.44(0.23)	66.26(0.50)	74.24(0.26)	72.86(0.51)	66.70(0.58)	62.14(0.64)	4.11	4.09	4.38
Findings	DeepSeek-V2 (zero-shot)	60.29(0.70)	42.20(0.90)	52.00(0.80)	45.21(0.76)	38.66(0.82)	34.15(0.85)	3.53	3.74	3.53
	Lattice	71.83(0.73)	56.80(0.71)	65.67(0.71)	54.73(0.56)	46.33(0.56)	39.69(0.53)	3.63	3.41	3.56
	TableLLM-7B	74.19(0.58)	57.45(0.81)	66.93(0.67)	64.55(0.65)	57.35(0.74)	51.92(0.81)	3.85	3.86	3.90
	Chinese-Alpaca-2-7B-16K	79.02(0.66)	64.35(0.77)	72.43(0.66)	70.11(0.73)	63.91(0.78)	59.16(0.80)	4.03	4.00	4.36
	EMGLLM (Ours)	80.36(0.52)	66.03(0.69)	73.92(0.53)	71.83(0.54)	65.75(0.61)	61.05(0.66)	4.10	4.13	4.40
	DeepSeek-V2 (zero-shot)	51.83(0.72)	33.10(0.80)	48.54(0.81)	40.89(1.02)	33.10(0.90)	28.02(0.83)	3.36	3.48	3.30
T	Lattice	65.06(0.56)	46.51(0.74)	63.04(0.59)	50.77(0.59)	39.60(0.60)	30.19(0.63)	3.29	3.14	3.43
impressions	TableLLM-7B	70.36(0.54)	51.53(0.79)	67.91(0.60)	62.69(0.66)	53.60(0.72)	46.87(0.81)	3.59	3.54	3.86
	Chinese-Alpaca-2-7B-16K	76.68(0.34)	61.04(0.55)	74.72(0.40)	70.41(0.50)	62.85(0.59)	57.14(0.62)	4.01	3.85	4.06
	EMGLLM (Ours)	77.21(0.41)	61.49(0.86)	75.26(0.43)	70.91(0.36)	63.29(0.64)	57.38(0.86)	4.13	4.05	4.36

Table 2: **Main Results.** Average results (standard deviation) of EMGLLM and baseline models on the ETM test set. All automatic evaluations are tested with 5 random seeds.

Model				Model Evaluation					
		ROUGE-1	ROUGE-2	ROUGE-L	BLEU-1	BLEU-2	BLEU-3	Correctness	Completeness
Overall	Chinese-Alpaca-2-7B-16K	77.48(0.71)	62.82(0.71)	70.86(0.62)	68.25(0.80)	61.97(0.78)	57.48(0.78)	3.88	3.85
	EMGLLM (Ours)	79.48(0.65)	65.73(0.76)	73.51(0.70)	71.64(0.72)	65.63(0.77)	61.25(0.80)	3.93	3.92
Findings	Chinese-Alpaca-2-7B-16K	77.63(0.65)	62.81(0.72)	70.48(0.59)	67.72(0,81)	61.54(0.81)	56.83(0.82)	3.88	3.91
	EMGLLM (Ours)	79.30(0.48)	65.36(0.53)	72.92(0.46)	70.64(0.54)	64.71(0.55)	60.14(0.56)	3.90	3.93
Impressions	Chinese-Alpaca-2-7B-16K	73.36(0.71)	55.17(0.88)	71.16(0.66)	65.82(0.74)	57.07(0.78)	50.92(0.87)	3.87	3.78
	EMGLLM (Ours)	76.91(1.10)	60.56(1.85)	74.63(1.20)	70.02(1.31)	62.16(1.76)	56.38(2.10)	3.95	3.92

Table 3: **Data-intensive Results.** Average results (standard deviation) of experiment on a subset with larger average amount of continuous values. EMGLLM demonstrates more significant advantages. All automatic evaluations are tested with 5 random seeds.

4.2 Setup

375

376

377

382

388

389

391

4.2.1 Baseline Methods

We select various baseline models capable of performing automatic EMG diagnosis, including both general text-to-text generation models and models specifically designed for data-to-text tasks.

Chinese-Alpaca-2-7B-16K Chinese-Alpaca-

2-7B-16K (Cui et al., 2023) is a widely used LLM. It also serves as the base LLM module for EMGLLM. The prompt template for Chinese-Alpaca-2-7B-16K is fully consistent with that of EMGLLM, with the only difference being that Chinese-Alpaca-2-7B-16K directly process the continuous data in textual form. Besides, this model is similarly fine-tuned using the LoRA method, with training hyperparameters consistent with those of EMGLLM. The comparison with EMGLLM can clearly demonstrate the effect of the EMG Alignment Encoder on the automatic generation of diagnostic results.

394**TableLLM-7B**TableLLM (Zhang et al., 2024b)395is an LLM specifically designed for tabular data396inputs, fine-tuned on a large dataset of table397tasks. Since the base model used by TableLLM,

CodeLlaMA-7B (Rozière et al., 2023), does not support Chinese, we replicate the training using the official code on Chinese-CodeLlaMA-7B to develop a Chinese version TableLLM , and subsequently fine-tune it on ETM dataset.

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

Lattice Lattice (Wang et al., 2022) is a data-totext generation model with a structure-aware selfattention mechanism and a tranformation-invariant positional encoding mechanism improved from T5base.

DeepSeek-V2 (zero-shot) We conduct experiments on DeepSeek-V2 (DeepSeek-AI, 2024), a powerful general Chinese model with 236B parameters, in a zero-shot setting. In addition to the basic prompt template, we provide several output examples to guide the model in generating *Findings* and *Impressions* in the correct format.

4.2.2 Implementation Details

For the implementation of EMGLLM, we first obtain reference values from the ETM training set. From a total of 13,820 samples, we filter out 7,166 (52%) completely healthy samples based on text rules and perform quantile statistics on each examination item i to determine reference values

	Model	ROUGE-1	ROUGE-2	ROUGE-L	BLEU-1	BLEU-2	BLEU-3
	w/o Reference Value	79.02(0.38)	64.32(0.53)	72.79(0.47)	70.81(0.58)	64.44(0.64)	59.85(0.71)
Orranall	w/o Encoder Pre-training	79.24(0.18)	64.98(0.38)	73.02(0.12)	71.44(0.36)	65.26(0.40)	60.77(0.45)
Overall	EMGLLM (rule-based)	79.86(0.28)	65.50(0.62)	73.69(0.30)	72.38(0.43)	66.09(0.54)	61.45(0.62)
	w/o Encoder Fine-tuning	80.07(0.12)	65.92(0.33)	74.09(0.12)	72.66(0.27)	66.47(0.35)	61.91(0.40)
	EMGLLM	80.44(0.23)	66.26(0.50)	74.24(0.26)	72.86(0.51)	66.70(0.58)	62.14(0.64)
Findings	w/o Reference Value	79.03(0.36)	64.26(0.58)	72.36(0.36)	70.03(0.58)	63.81(0.66)	59.05(0.72)
	w/o Encoder Pre-training	79.00(0.45)	64.20(0.59)	72.27(0.46)	70.21(0.52)	63.96(0.57)	59.19(0.63)
	EMGLLM (rule-based)	79.45(0.35)	64.60(0.52)	72.76(0.32)	70.89(0.32)	64.56(0.39)	59.67(0.45)
	w/o Encoder Fine-tuning	79.83(0.26)	65.18(0.31)	73.29(0.28)	71.42(0.36)	65.15(0.37)	60.34(0.41)
	EMGLLM	80.36(0.52)	66.03(0.69)	73.92(0.53)	71.83(0.54)	65.75(0.61)	61.05(0.66)
	w/o Reference Value	74.88(0.29)	57.68(0.54)	72.62(0.24)	67.93(0.29)	59.56(0.40)	53.37(0.47)
Immunationa	w/o Encoder Pre-training	76.27(0.29)	60.19(0.48)	74.15(0.38)	69.94(0.47)	62.16(0.53)	56.31(0.59)
mpressions	EMGLLM (rule-based)	76.93(0.49)	60.79(0.81)	74.83(0.54)	70.56(0.69)	62.69(0.83)	56.71(0.95)
	w/o Encoder Fine-tuning	77.45(0.55)	61.72(0.70)	75.49(0.42)	71.15(0.48)	63.52(0.60)	57.58(0.66)
	EMGLLM	77.21(0.41)	61.49(0.86)	75.26(0.43)	70.91(0.36)	63.29(0.64)	57.38(0.86)

Table 4: Ablation Study Results. All automatic evaluations are tested with 5 random seeds.

 u_i and l_i . Subsequently, 7 quantile thresholds $\{p_1, p_2, p_3, ..., p_7\} = \{0.02, 0.05, 0.08, ..., 0.2\}$ are set to construct the reference matrix \hat{X}_i .

In the EMG Alignment Encoder, the output channel number $d_C = 64$, the number of output embeddings N = 2, the size of text prototypes V' = 192, and the number of heads $n_{head} = 8$.

For the LLM component of EMGLLM, we select the widely-used Chinese-Alpaca-2-7B-16K as base model. Pre-training of EMG Alignment Encoder is conducted for 2000 steps, followed by 5 epochs of fine-tuning, with a batch size of 1 and a gradient accumulation step of 16. Optimization is performed using the Adam optimizer, with a learning rate of 5e-5. The LLM is trained using the LoRA method, with a rank of 8, an alpha value of 16, and the training target set to ['q_proj', 'v_proj'].

For the training of baseline models, we preprocess the dataset according to the input and output formats required by the model and employ the recommended hyperparameters of the projects.

4.2.3 Metrics

422

423

494

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

To comprehensively evaluate the quality of EMG automatic diagnosis, we use multidimensional metrics. The automatic metrics include:

• ROUGE (Recall-Oriented Understudy for Gisting Evaluation) (Lin, 2004): It measures the generation quality by comparing the overlap between texts. ROUGE-1, ROUGE-2, and ROUGE-L are selected as metrics.

• BLEU (Bilingual Evaluation Understudy) (Papineni et al., 2002): It compares the n-gram match between texts. We use BLEU-1, BLEU-2, and BLEU-3 to evaluate the model's capabilities.

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

In addition, we introduce model evaluation, using GPT-40 (OpenAI, 2023) as a judge to assess the quality of the model-generated diagnoses. We provide authoritative doctors' ground truth diagnoses as a reference for GPT-40, simultaneously inputting the model-generated results, allowing GPT-40 to analyze and compare the differences between the two and provide a multidimensional evaluation. GPT-40's scoring criteria include:

• **Correctness:** evaluate whether a diagnosis falsely reports non-existent abnormalities. A higher score indicates fewer false positives.

• **Completeness:** evaluate whether a diagnosis has missed reporting existing abnormalities. A higher score indicates fewer missed abnormalities.

In model evaluation, GPT-40 evaluates the *Find-ings* and *Impressions* separately, and we use the average of these two evaluations as the overall score for the diagnosis. The evaluation template for GPT-40 is presented in Appendix A.

Finally, we conduct the human evaluation and provide human scores. We sample 50 examples from the test set and rate from 1 to 5 to the generated outputs of each model. These human experts are graduate students responsible for research and development projects in the medical technology field. The criteria for the human evaluation can be found in Appendix A. We ask human experts to score the Findings and Impression separately based on the following scoring criteria, and the average of these two scores is then taken as the Overall score.

4.3 Results

489

490

491

492

493

494

495

496

497

498

499

502

503

504

508

509

510

512

513

514

515

516

518

519

521

522

525

527

529

530

531

533

534

535

538

4.3.1 Main Results

Table 2 presents the main results of the EMG automatic diagnosis generation. In automatic, model and human evaluations, it can be observed that EMGLLM outperforms all baseline methods, including data-to-text models such as Lattice and TableLLM-7B. In particular, the comparison between EMGLLM and Chinese-Alpaca-2-7B-16K clearly illustrates the improvement brought by the EMG Alignment Encoder to the LLM in EMG automatic diagnosis. This demonstrates that the EMGLLM framework effectively utilizes test values and reference values to reasonably encode numerical data in medical tables, resulting in higherquality diagnosis generation.

Additionally, zero-shot DeepSeek-V2 shows lower performance, indicating that a general LLM without specific fine-tuning lacks the knowledge of medical data. This underscores the importance of datasets for medical tables and highlights the contribution of the ETM.

We also observe that for all models in the experiment, the rankings of the model evaluation metrics are generally consistent with those of the human evaluation scores, indicating that GPT-4 can serve as a substitute for human evaluation in our task.

4.3.2 Effectiveness on Data-intensive Input Scenario

To further validate the effectiveness of data encoding method, we extract samples with a relatively large number of continuous numerical data from the ETM dataset, resulting in a data-intensive subset. This subset contains 5,000 training samples and 600 test samples, with an average of 43.49 continuous numerical values per sample, higher than 33.14 shown in Table 1. As shown in Table 3, compared to the results from training and testing on the full dataset in Table 2, the performance gap between EMGLLM and Chinese-Alpaca-2-7B-16K widens, exceeding 3 in overall diagnoses, 2 in Findings, and 6 in Impressions in terms of ROUGE-2. Therefore, as the data amounts in the tables increase and the task becomes more challenging, EMGLLM demonstrates greater robustness.

4.3.3 Ablation Study

In Section 3.1, we propose a method for obtaining reference values and attempt to compare test values with them using the EMG Alignment Encoder. A natural question arises: once reference values are obtained, is it effective to directly convert continuous numerical data into categorical terms such as "high", "normal", or "low" based on rules without the EMG Alignment Encoder? Therefore, we conduct experiment on a rule-based approach for processing data input. Specifically, if the test value for item *i* exceeds $u_i^{0.05}$, it is denoted as "high"; if it is below $l_i^{0.05}$, it is denoted as "low"; otherwise, it is denoted as "normal". The LLM trained by this rule-based method is denoted as EMGLLM (rule-based) in Table 4. It is shown that replacing the EMG Alignment Encoder with rules leads to a certain degradation in performance. This indicates that the rule-based method is relatively inflexible in handling medical examination tables. Besides, to verify the necessity of introducing reference values, we evaluate EMGLLM without reference values by replacing each detection item's reference values with random numbers from a standard normal distribution during model fine-tuning phase. As shown in Table 4, this leads to a significant performance drop.

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

588

We also conduct ablation study over training methods in Section 3.2. As shown in Table 4, pretraining the EMG Alignment Encoder is essential, resulting in a well-calibrated initialization. Finetuning the EMG Alignment Encoder in conjunction with the LLM on real EMG diagnostic datasets can further enhance the capabilities.

5 Conclusion

In this paper, we propose EMGLLM, a medical data-to-text model, for the automatic diagnosis generation of Electromyography (EMG) tables. The model framework with the EMG Alignment Encoder can enhance the encoding of continuous numerical data, enabling the model to simulate the process by which physicians compare test values to reference values during diagnosis. This approach facilitates a better model understanding of the degree of health and abnormality reflected by the data. In addition, we construct the ETM dataset, which comprises 17,276 real case examination EMG tables and diagnoses from authoritative doctors, to support the advancement of medical data-to-text Finally, experimental results demonresearch. strate that EMGLLM outperforms baseline methods in all automatic, model and human evaluations for EMG diagnosis generation, confirming the effectiveness of the EMGLLM approach in handling medical examination data for automatic diagnosis.

641 642 643 644 645 646 647 648 649 650 651 652 653 654 655 656 657 658 659 660 661 662 663 664 665 666 667 668 669 670 671 672 673 674 675 676 677 678 679 680 681

682

683

684

685

686

687

688

689

690

639

640

589 Limitation

EMGLLM is designed to augment the model understanding of continuous numerical data in medical examination tables, without addressing other
elements of the tables. At present, our experiments
have been conducted solely on the EMG task. We
will extend our model to other types of medical examinations, such as complete blood counts and urinalysis tables in future works.

Ethics Statement

We acknowledge the limitations of current LLMs 599 and the ethical implications of their use, including the potential for inaccurate or misleading responses in diagnosis. However, our research has shown an improvement in reliability compared to baseline methods. Future research may explore 604 more robust methods to address these challenges. Our dataset is constructed from real medical diagnostic reports and contains basic information about patients. However, the dataset we publish is completely anonymous, and we will only disclose age, gender, and height information without revealing 610 any other private information. The release of this 611 dataset has been approved for use by Huashan Hospital Affiliated with Fudan University.

> All human participants involved in the evaluation of this research were compensated at or above the average local wage rate, ensuring fair remuneration for their time and contributions.

References

614

615

616

617

618

619

625

629

630

631

- Andrea J Boon, Kais I Alsharif, C Michel Harper, and Jay Smith. 2008. Ultrasound-guided needle emg of the diaphragm: technique description and case report. *Muscle & Nerve: Official Journal* of the American Association of Electrodiagnostic Medicine, 38(6):1623–1626.
- Arthur Bureau, Jean-Maxime Le Carpentier, Eric Le Carpentier, Yannick Aoustin, and Yann Péréon.
 2021. Décomposition et analyse de tracés emg pour aider au diagnostic des maladies neuromusculaires. *Preprint*, arXiv:2109.14922.
- Stanislas Chambon, Mathieu Galtier, Pierrick Arnal, Gilles Wainrib, and Alexandre Gramfort. 2017. A deep learning architecture for temporal sleep stage classification using multivariate and multimodal time series. *Preprint*, arXiv:1707.03321.
- Navin Cooray, Zhenglin Li, Jinzhuo Wang, Christine Lo, Mahnaz Arvaneh, Mkael Symmonds, Michele Hu, Maarten De Vos, and Lyudmila S Mihaylova. 2022. Automated movement detection with dirichlet

process mixture models and electromyography. In 2022 25th International Conference on Information Fusion (FUSION), volume 31, page 01–08. IEEE.

- Yiming Cui, Ziqing Yang, and Xin Yao. 2023. Efficient and effective text encoding for chinese llama and alpaca. *arXiv preprint arXiv:2304.08177*.
- Jasper R Daube. 2002. Assessing the motor unit with needle electromyography. *CONTEMPORARY NEU-ROLOGY SERIES*, 66:293–323.
- DeepSeek-AI. 2024. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *Preprint*, arXiv:2405.04434.
- Zhihao Fan, Jialong Tang, Wei Chen, Siyuan Wang, Zhongyu Wei, Jun Xi, Fei Huang, and Jingren Zhou. 2024. Ai hospital: Benchmarking large language models in a multi-agent medical interaction simulator. *Preprint*, arXiv:2402.09742.
- Mekia Shigute Gaso, Selcuk Cankurt, and Abdulhamit Subasi. 2021. Electromyography signal classification using deep learning. In 2021 16th International Conference on Electronics Computer and Computation (ICECCO). IEEE.
- Siavash Golkar, Mariel Pettee, Michael Eickenberg, Alberto Bietti, Miles Cranmer, Geraud Krawezik, Francois Lanusse, Michael McCabe, Ruben Ohana, Liam Parker, Bruno Régaldo-Saint Blancard, Tiberiu Tesileanu, Kyunghyun Cho, and Shirley Ho. 2023. xval: A continuous number encoding for large language models. *Preprint*, arXiv:2310.02989.
- Qizheng Gu, Cong Nie, Ruixiang Zou, Wei Chen, and Dong Tian. 2020. Automatic generation of electromyogram diagnosis report. *IEEE*.
- YiQiu Guo, Yuchen Yang, Ya Zhang, Yu Wang, and Yanfeng Wang. 2024. Dictllm: Harnessing key-value data structures with large language models for enhanced medical diagnostics. *Preprint*, arXiv:2402.11481.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *Preprint*, arXiv:2106.09685.
- Shailza Jolly, Zi Xuan Zhang, Andreas Dengel, and Lili Mou. 2021. Search and learn: Improving semantic coverage for data-to-text generation. *Preprint*, arXiv:2112.02770.
- Bashar Katirji. 2002. The clinical electromyography examination: An overview. *Neurologic clinics*, 20(2):291–303.
- X. Li, X. Zhang, X. Yi, D. Liu, H. Wang, B. Zhang, B. Zhang, D. Zhao, and L. Wang. 2023. Review of medical data analysis based on spiking neural networks. *Preprint*, arXiv:2212.02234.

- 711 712 713 714 715 716 717 719 721
- 725 726 727 728 729
- 731 732 733
- 734
- 737
- 739 740 741
- 742
- 743
- 744 745

- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In Text summarization branches out, pages 74-81.
- Van Hieu Nguyen, Gia Thien Luu, Thien Van Luong, Mai Xuan Trang, Philippe Ravier, and Olivier Buttelli. 2023. After-fatigue condition: A novel analysis based on surface emg signals. Preprint, arXiv:2309.04770.
- Jianmo Ni, Chun-Nan Hsu, Amilcare Gentili, and Julian McAuley. 2020. Learning visual-semantic embeddings for reporting abnormal findings on chest xrays. arXiv preprint arXiv:2010.02467.
- OpenAI. 2023. Gpt-4 technical report. Technical report, OpenAI.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th annual meeting of the Association for Computational Linguistics, pages 311–318.
- Ildar Rakhmatulin. 2024. Pieeg-16 to measure 16 eeg channels with raspberry pi for brain-computer interfaces and eeg devices. Preprint, arXiv:2409.07491.
- Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Ellen Tan, Yossef Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Defossez, Jade Copet, Faisal Azhar, Hugo Touvron, Gabriel Synnaeve, Louis Martin, Nicolas Usunier, and Thomas Scialom. 2023. Code llama: Open foundation models for code. Technical report, MetaAI.
- Swarnadeep Saha, Xinyan Velocity Yu, Mohit Bansal, Ramakanth Pasunuru, and Asli Celikyilmaz. 2022. Murmur: Modular multi-step reasoning for semi-structured data-to-text generation. Preprint, arXiv:2212.08607.
- Mandar Sharma, Ajay Gogineni, and Naren Ramakrishnan. 2024. Innovations in neural data-to-text generation: A survey. Preprint, arXiv:2207.12571.
- Niklas Smedemark-Margulies, Yunus Bicer, Elifnur Sunger, Tales Imbiriba, Eugene Tunik, Deniz Erdogmus, Mathew Yarossi, and Robin Walters. 2023. Fast and expressive gesture recognition using a combination-homomorphic electromyogram encoder. Preprint, arXiv:2311.14675.
- Fei Wang, Zhewei Xu, Pedro Szekely, and Muhao Chen. 2022. Robust (controlled) table-to-text generation with structure-aware equivariance learning. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies.
- Sam Wiseman, Arturs Backurs, and Karl Stratos. 2021. Data-to-text generation by splicing together nearest neighbors. Preprint, arXiv:2101.08248.

Ruslan Yermakov, Nicholas Drago, and Angelo Ziletti. 2021. Biomedical data-to-text generation via fine-tuning transformers. arXiv preprint arXiv:2109.01518.

746

747

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

768

769

771

773

774

- Tianshu Zhang, Xiang Yue, Yifei Li, and Huan Sun. 2024a. Tablellama: Towards open large generalist models for tables. *Preprint*, arXiv:2311.09206.
- Xiaokang Zhang, Jing Zhang, Zeyao Ma, Yang Li, Bohan Zhang, Guanlin Li, Zijun Yao, Kangli Xu, Jinchang Zhou, Daniel Zhang-Li, Jifan Yu, Shu Zhao, Juanzi Li, and Jie Tang. 2024b. Tablellm: Enabling tabular data manipulation by llms in real office usage scenarios. Preprint, arXiv:2403.19318.
- Fengbin Zhu, Ziyang Liu, Fuli Feng, Chao Wang, Moxin Li, and Tat-Seng Chua. 2024. Tat-llm: A specialized language model for discrete reasoning over tabular and textual data. Preprint, arXiv:2401.13223.

Model and Human Evaluation Details Α

In Sections 4.3.1 and 4.3.2, we introduce a GPT-40-based model evaluation with 0.1 model temperature for stability. GPT-40 is tasked with scoring both the Findings and Impression sections, and we use the average of these two scores as the overall score of a diagnosis. The template is shown in Figure 5, where we provide the full scoring criteria for Correctness and Completeness, allowing GPT-40 to reference the authoritative doctor's diagnosis when assigning scores.



Figure 5: Template for GPT-40 evaluation of EMG diagnosis generation

775 776	In the human evaluation, we ask human experts to score based on the following scoring criteria.
777	5 - The generated diagnosis is completely iden-
778	tical to the real diagnosis. Not only the con-
779	clusion but also the detailed descriptions are
780	fully consistent.
781	4 - The generated diagnosis and the real diag-
782	nosis have identical conclusions, and most of
783	the detailed descriptions are accurate. There
784	may be minor omissions or incomplete de-
785	scriptions in certain details, but these discrep-
786	ancies do not affect the overall diagnostic con-
787	clusion.
788	3 - The generated diagnosis and the real diagno-
789	sis have the same direction, and the conclu-
790	sions are generally consistent, but there are 1
791	to 2 notable discrepancies and slight inaccu-
792	racies in details.
793	2 - The generated diagnosis is largely inconsis-
794	tent with the real diagnosis, with only a few
795	minor details agreeing.
796	1 - The generated diagnosis is completely oppo-
797	site to the real diagnosis. The conclusion is
798	significantly erroneous, with a fundamentally
799	incorrect assessment of the condition, which
800	does not meet medical standards.