

LogLLM: 基于日志的大语言模型异常检测

Wei Guan¹, 曹建^{1*}, 钱石友¹, 高建奇¹, 欧阳春²
¹Department of Computer Science and Engineering, SJTU, Shanghai, China
The 昆士兰科技大学信息系统学院, 布里斯班, 澳大利亚 a
{关伟, 曹健, 钱石友, 193139}@sjtu.edu.cn, c.ouyang@qut.edu.au

摘要—— 软件系统通常在日志中记录重要的运行时信息, 以帮助故障排除。基于日志的异常检测已成为一个关键的研究领域, 旨在通过日志数据识别系统问题, 最终提高软件系统的可靠性。传统的深度学习方法往往难以捕捉嵌入在日志数据中的语义信息, 这些数据通常以自然语言组织。在本文中, 我们提出了 LogLLM, 一个基于日志的异常检测框架, 该框架利用大型语言模型 (LLMs)。LogLLM 使用 BERT 从日志消息中提取语义向量, 同时利用基于 transformer 解码器的 Llama 对日志序列进行分类。此外, 我们引入了一个投影器, 以对齐 BERT 和 Llama 的向量表示空间, 确保对日志语义的连贯理解。与需要日志解析器提取模板的传统方法不同, LogLLM 使用正则表达式预处理日志消息, 简化了整个流程。我们的框架通过一个新颖的三阶段程序进行训练, 旨在提高性能和适应性。在四个公开数据集上的实验结果表明, LogLLM 优于最先进的方法。即使处理不稳定的日志, 它也能有效地捕捉日志消息的语义意义并准确检测异常。

关键词: 系统日志, 异常检测, 大语言模型, 深度学习, 日志分析

1. 引言

确保高可用性和可靠性对于大规模软件密集型系统至关重要。随着这些系统变得更加复杂和庞大, 异常的发生变得不可避免。即使是微不足道的问题也可能导致性能下降、数据完整性问题, 以及客户和收入的大量损失。因此, 异常检测对于维护复杂软件密集型系统的健康和稳定性至关重要。

软件密集型系统通常会产生控制台日志, 记录系统状态和关键运行时事件。工程师可以利用这些日志数据来评估系统健康、识别异常和追踪问题的根本原因。然而, 由于日志的潜在大量, 手动分析它们以查找异常既费时又容易出错。因此, 基于日志的异常检测已成为自动日志分析的关键领域, 专注于通过日志数据自动识别系统异常。

近年来, 针对基于日志的异常检测, 已经提出了许多基于深度学习的算法。这些方法通常采用序列深度学习模型, 如 LSTM 和 Transformer。这些方法可以进一步分为基于重建的方法和基于二分类的方法。基于重建的方法涉及设计和训练一个深度神经网络来重建输入日志序列, 异常检测基于重建误差。其基本原理是异常样本无法被准确重建。另一方面, 基于二分类的方法涉及设计一个二分类器来将样本分类为正常或异常。这些方法通常需要标记的异常样本进行训练。众所周知, 系统日志是用自然语言编写的, 并且包含大量的语义信息。然而, 传统的基于深度学习的算法在有效捕获这些信息方面存在困难。

近年来, LLM (大型语言模型) 取得了显著的进展, 如 GPT、Llama 和 ChatGLM。这些模型以其庞大的参数规模为特征, 并在大量数据集上进行了预训练, 数据集大小从几个 GB 到 TB 不等。这种广泛的预训练使它们具备了卓越的语言理解能力, 即使在零样本或少量样本的情况下, 也能在摘要、释义和指令遵循等任务中表现出色。现有的利用 LLM 进行基于日志的异常检测的方法可以分为基于提示工程的方法和基于微调的方法。基于提示工程的方法利用 LLM 的零样本 / 少量样本能力, 仅基于模型内部知识来检测异常。然而, 这些方法往往难以针对特定数据集定制解决方案, 导致检测性能不佳。基于微调的方法将 LLM 集成到深度神经网络中, 并针对用户特定的数据集进行定制。然而, 这些方法面临着诸如语义理解有限、LLM 利用不足 (仅依赖 LLM 进行语义信息提取) 以及未充分考虑输入数据格式等问题, 可能导致内存溢出。

为了应对上述挑战, 我们提出了 LogLLM, 这是一种基于日志的异常检测框架, 它利用了大型语言模型。与传统方法不同, 传统方法依赖于

通讯作者。

LogLLM 通过正则表达式预处理日志消息，从而简化了整个流程。LogLLM 是一种基于微调的方法，利用基于 Transformer 编码器的模型 BERT 从日志消息中提取语义向量。此外，它还采用基于 Transformer 解码器的模型 Llama 对日志序列进行分类。为了确保日志语义的一致性，我们引入了一个投影器，以对齐 BERT 和 Llama 的向量表示空间。我们的框架采用了一种新颖的三阶段训练程序，旨在提高性能和适应性。

如第 V-G 节所述，LLM 由于其庞大的参数规模，经常面临内存不足的问题 [41]。直接将整个日志序列（通过将日志消息连接成一个长字符串）输入到 Llama 中可能导致内存不足问题，并可能使 LLM 困惑，使其难以关注区分异常的关键点。通过采用 BERT 对每个日志消息进行总结，LogLLM 有效地缓解了这些问题。我们在四个公开数据集上进行了实验，结果表明 LogLLM 优于现有方法。即使处理不稳定的日志，由于软件演变，新日志模板频繁出现，它也能有效地捕捉日志消息的语义意义并准确检测异常。消融研究证实了三阶段训练程序的有效性。

我们工作的主要贡献如下：

我们引入了 LogLLM，这是一个利用 LLM 的基于日志的异常检测框架。这项研究是首次尝试同时使用基于 Transformer 编码器和解码器的 LLM，即 BERT 和 Llama，进行基于日志的异常检测。我们提出了一种新颖的三阶段程序，以优化深度模型中不同组件的训练和协调，从而提高性能和适应性。

我们在四个公开可用的真实世界数据集上进行了广泛的实验，证明了 LogLLM 取得了卓越的性能。

II. 相关工作

在本节中，我们探讨了基于日志的异常检测领域的相关工作，特别关注基于深度学习的方法。我们特别关注利用预训练大型语言模型（LLM）的方法。

A. 基于传统深度学习的日志异常检测

针对基于日志的异常检测，已经提出了许多基于传统深度学习的方法。根据训练范式，这些工作可以分为两类：基于重建的方法和基于二分类的方法。

基于重建的方法 [8][15] 涉及设计和训练一个深度神经网络来重建输入日志

基于重建的方法 $\{v_2\}$ - $\{v_4\}$ 包括设计和训练一个深度神经网络来重建输入日志序列。异常检测基于重建误差。正常的日志序列可以以最小的误差重建，而异常的日志序列则无法有效重建，导致显著更高的重建误差。这些方法始终在无异常的正常数据上训练深度模型，这意味着它们是半监督的。

DeepLog [8] 采用 LSTM 预测基于过去日志序列的下一个日志模板 ID。类似地，LogAnomaly [9] 基于序列和定量模式预测下一个日志模板 ID。自编码器（AEs）[10]-[13] 和生成对抗网络（GANs）[14], [15] 在基于重建的方法中得到广泛应用。例如，LogAttn [10] 采用一个结合时间卷积网络（TCN）以捕获时间语义相关性和深度神经网络（DNN）以捕获统计相关性的自编码器。Duan 等人 [14] 使用一个 GAN，其中基于 LSTM 的编码器-解码器框架作为生成器。卷积神经网络（CNNs）用作判别器。重建误差基于输入和生成器输出的差异计算。

基于二分类的方法 [16]-[22] 通常使用输出一个或两个值的深度神经网络。通常，一个值表示样本属于异常类的概率，通过应用阈值将此概率转换为二进制分类来检测异常。当输出两个值时，它们分别表示样本属于正常类和异常类的概率。

大多数方法 [16]-[20] 通常以监督方式训练深度模型。例如，张等人 [16] 提出了 LayerLog，该模型整合了词、日志和日志序列层以从日志序列中提取语义特征。CNNs 在 [17], [18] 中被用于开发二分类器。LogRobust [19] 整合了一个预训练的 Word2Vec 模型，特别是 FastText [42]，并将其与 TF-IDF 权重结合起来学习日志模板的表示向量。然后，这些向量被输入到一个基于注意力的 Bi-LSTM 模型中进行异常检测。LogGD [20] 将日志序列转换为图，并利用一个结合图结构和节点语义的图变换神经网络进行基于日志的异常检测。

一些工作 [21], [22] 涉及以半监督方式训练二元分类器。例如，Trine [21] 使用 Transformer 编码器 [24] 将正常日志序列编码成向量表示，并使用生成器生成随机的虚假向量表示。判别器由 Transformer 和多层感知机（MLP）组成，用于区分给定的向量表示是否为正常日志序列，随后用于检测异常。PLELog [22] 通过采用概率标签估计和开发基于注意力的 GRU 神经网络来解决标签不足的挑战，用于异常检测。

公认的是，系统日志以自然的方式记录。

Header	Content
1117838978 2005.06.03 R02-M1-N0-C:J12-U11 2005-06-03-15.49.38.026704 R02-M1-N0-C:J12-U11 RAS KERNEL INFO	instruction cache parity error corrected
1117843015 2005.06.03 R21-M1-N6-C:J08-U11 2005-06-03-16.56.55.309974 R21-M1-N6-C:J08-U11 RAS KERNEL INFO	141 double-hammer alignment exceptions
1117848119 2005.06.03 R16-M1-N2-C:J17-U01 2005-06-03-18.21.59.871925 R16-M1-N2-C:J17-U01 RAS KERNEL INFO	CE sym 2, at 0x0b85eee0, mask 0x05
...	...

图 1: 系统日志的一个示例。

语言并包含大量语义信息。然而，传统的基于深度学习的传统方法在捕捉这些信息方面面临挑战。

B. 基于日志的异常检测的 LLM

现有的 LLM 可以分为基于 transformer 编码器的模型，如 BERT [43], RoBERTa [44], 和 Span-BERT [45], 以及基于 transformer 解码器的模型，包括 GPT-4 [25], Llama 3 [26], 和 ChatGLM [27]。利用 LLM 的两种常见策略是提示工程和微调。

基于提示工程的方法仅通过依赖 LLM 的内部知识来检测异常。这些方法通常采用基于 transformer 解码器的模型。例如，Qi 等人 [7] 使用 ChatGPT 进行零样本和少量样本的基于日志的异常检测，利用将日志序列直接整合到提示模板中的提示模板。然而，当使用大窗口大小对日志消息进行分组时，这种方法变得不切实际。

Egersdoerfer 等人 [30] 通过维护基于摘要的内存来解决此问题，该内存总结了之前的日志消息，从而消除了异常检测中输入整个日志序列的需要。RAGLog [31] 使用检索增强生成（RAG）框架 [46] 通过查询其存储的正常日志条目样本来分析日志条目。他们为 LLM 设计提示模板以确定查询的日志条目是正常还是异常。基于提示工程的方法通常难以针对特定数据集定制解决方案，这可能导致特定数据集中的检测性能不佳。

基于微调的方法将 LLM 纳入深度神经网络并将它们定制到用户的自己的数据集。一些方法 [32]–[35]，虽然采用基于 transformer 编码器的 LLM 进行异常检测，但没有捕捉到日志序列中的语义信息。例如，LogBERT [32] 和 LAnoBERT [33] 使用 BERT 重建日志模板 ID（日志字符串模板的 ID）的输入序列，并根据重建错误检测异常，忽略了语义信息。其他方法 [3], [36]–[39] 使用基于 transformer 编码器的 LLM 仅从日志消息中提取语义信息，然后使用较小的模型 [3], [36]–[38] 或基于距离的比较 [39] 进行分类。例如，NeuralLog [3] 利用 BERT 从原始日志消息中提取语义向量，这些向量随后用于通过基于 transformer 的分类模型检测异常。同样，RAPID [39] 利用基于 transformer 编码器的模型提取语义向量并执行异常检测。

通过将每个查询日志序列与其最近的文档日志序列进行比较来进行检测。Hadadi 等人 [40] 直接将解析自日志序列的模板序列输入到 GPT 模型中，并对其进行微调以准确预测序列标签。然而，这种方法面临两个关键挑战。首先，序列中模板之间的边界不明确，这使得模型难以学习序列依赖关系。其次，每个模板可能被 LLM 的标记器标记成多个标记，并且单个序列可以包含多个日志模板。因此，可能会生成过多的标记，通常超过 LLM 的标记（内存）限制，从而限制了可以处理的序列长度。这两个挑战在第五节 -G 中进一步演示。

LogLLM 是一种基于微调的方法，它使用 BERT 从日志消息中提取语义向量，并使用 Llama（一个基于 transformer 解码器的模型）进行日志序列分类。该方法使用投影仪对 BERT 和 Llama 的向量表示空间进行对齐。BERT 的使用确保了日志消息之间的清晰边界，因为每个消息都由一个独特的嵌入向量表示，从而提高了分类性能。此外，当保持 Llama 的内存和参数大小不变时，这种方法可以处理比直接使用 Llama 的标记器对整个日志序列进行标记更长的序列。

III. 预备知识

为了为后续章节奠定基础，我们介绍了**系统日志**，该日志记录了系统在运行时的事件和内部状态。系统日志包含按时间顺序排列的日志消息列表。

图 1 展示了由 BGL（BlueGene/L 超级计算机系统）生成的原始系统日志片段，每个日志消息都按记录的时间顺序排列。这些原始日志消息是半结构化文本，由一个**标题**和**内容**组成。标题由日志框架确定，包括时间戳、详细程度（例如，WARN/INFO）和组件 [47] 等信息。日志内容包含一个固定部分（揭示日志模板的关键词）和一个可变部分（携带动态运行时信息的参数）。在本文中，我们仅关注每个日志消息的内容。

日志消息可以根据会话或固定 / 滑动窗口（即**记录特定执行流程的一系列日志消息**）进行分组。[48] 会话窗口分区根据日志消息的会话 ID 将日志消息分组，从而生成包含日志的消息序列。

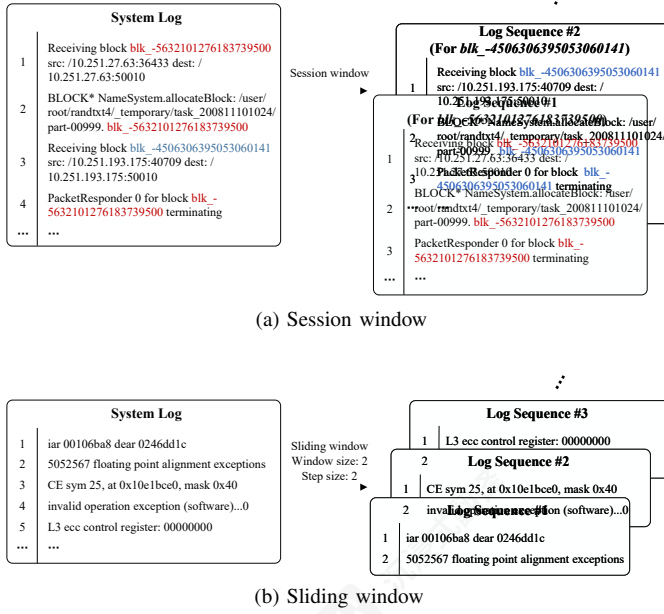


图 2: 日志消息分区示例的说明。

每个会话中的消息。例如，图 2a 展示了 HDFS [49] 日志在会话窗口分组过程中的情况，其中 *block_id* 作为会话 ID。相比之下，*xed/sliding window* partitioning 根据固定大小（窗口大小）对日志消息进行分组，该大小可以通过时间跨度或日志消息数量来定义。这种方法创建了捕获系统日志消息随时间快照的序列。例如，图 2b 展示了 BGL[50] 日志在滑动窗口分组过程中的情况，窗口大小为 2 条消息，步长为 2 条消息。

基于日志的异常检测的目的是识别异常日志序列，以便识别系统操作行为中可能存在的问题。

IV. 方法论

在本节中，我们介绍了我们创新的异常检测框架 LogLLM。如图 3 所示，日志序列在输入到深度神经网络之前，使用正则表达式进行预处理，该神经网络集成了 BERT [43]、a 投影器和 Llama [26] 用于日志序列分类。在接下来的章节中，我们将详细介绍日志序列预处理、深度模型架构以及模型训练过程。

A. 预处理

考虑到日志消息内容包含携带动态运行时信息的变量参数，这些参数始终与异常无关，并使深度模型训练复杂化，如第 V-F 节所示，需要一种技术来识别这些参数并将它们替换为常量标记。日志解析器，如 Drain [51] 和 Spell [52]，在基于日志的异常检测方法中得到广泛应用，并似乎是一种有用的技术。然而，正如 Le 等人所指出的。

[3]，现有的日志解析器并不总是对所有日志数据集都表现正确，并且难以处理新日志消息中的词汇表外（OOV）单词，导致语义信息丢失。当日志不稳定时，这些解析器随着时间的推移变得越来越无效，使得支持后续异常检测变得困难。

得益于结构化日志生成过程，可以使用正则表达式 [53] 轻松地识别表示特定对象的参数的文本格式。因此，我们将每个变量参数，如账户、目录路径和 IP 地址，替换为 ‘<*>’。尽管这种方法很简单，但它提供了显著的性能优势。与日志解析器相比，这种预处理技术更有效，且不需要训练。

B. 模型架构

如图 3 所示，我们的深度模型由三个主要组件组成：BERT、投影器和 Llama。BERT 和 Llama 都是预训练的大型语言模型。BERT 用于提取日志消息的向量表示，而 Llama 用于分类日志序列。投影器作为桥梁，将 BERT 和 Llama 的向量表示空间对齐。需要注意的是，我们的模型只包含一个 BERT 实例和一个投影器。

1) BERT: BERT 通过一个线性层后跟一个 *tanh* 激活函数处理分类标记（[CLS]）的语义向量来生成一个语义向量。每个预处理后的日志消息，使用 BERT 标记器和 BERT 模型编码成一个语义向量。对于预处理后的日志序列，BERT 的输出是一个语义向量序列 $C = (c_1, c_2, \dots, c_N) \in \mathbb{R}^{N \times d_{BERT}}$ ，其中 N 代表日志序列的长度（即日志消息的数量）， d_{BERT} 是每个语义向量的维度（即隐藏大小）。

2) 投影仪: 投影仪是一个线性层，它将语义向量 $C \in \mathbb{R}^{N \times d_{BERT}}$ 映射到 Llama 接受的标记嵌入向量，表示为 $E = (e_1, e_2, \dots, e_N) \in \mathbb{R}^{N \times d_{Llama}}$ ，其中 d_{Llama} 是 Llama 的隐藏大小。投影仪的设计是为了使 BERT 和 Llama 的向量表示空间对齐。

3) Llama: 为了对基于 Transformer 解码器的 LLM Llama 进行提示微调，我们根据嵌入的日志序列生成相应的文本查询。具体来说，每个查询由三个部分组成。

第一个组件引入了日志序列，例如“以下是系统日志消息的序列：”。第二个组件包括投影仪输出的标记嵌入 E 。第三个组件查询序列是否异常，例如询问“。这个序列是正常还是异常？”。第一个和第三个组件依次输入到 Llama tokenizer 和 Llama embedding layer，产生 $E_1 \in \mathbb{R}^{A \times d_{Llama}}$ 和 $E_3 \in \mathbb{R}^{Q \times d_{Llama}}$ ，其中 A 和 Q 分别是第一个和第三个组件标记化后产生的标记数量。然后，将三个组件的标记嵌入连接起来，

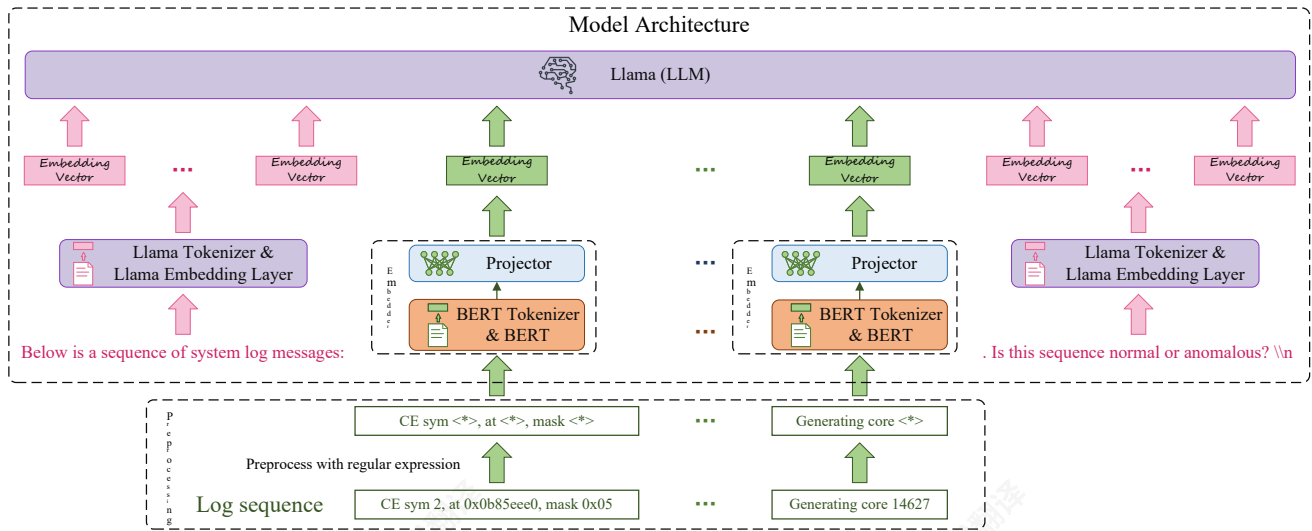


图 3: LogLLM 的框架。值得注意的是, 该模型包含一个 BERT 实例和投影

tor.

被重命名为 $[E_1||E||E_3] \in \mathbb{R}^{(A+N+Q) \times d_{Llama}}$ 并输入到 Llama 中。

C. 训练

1) 少数类过采样: LogLLM 是一种监督式异常检测方法, 这意味着它需要标记的正常和异常样本进行训练。然而, 监督式异常检测方法通常面临数据不平衡的挑战, 这可能导致模型训练偏差。在异常检测任务中, 只有两个类别: 正常和异常, 每个类别的实例数量是不确定的。为了应对数据不平衡, 我们对样本数量较少的类别进行过采样, 确保少数类别的比例不少于 β 。形式上, 设少数类别的比例为 α 和 $\alpha < \beta$, 样本总数为 $Sample_num$ 。为了使少数类别的比例为 β , 需要过采样到以下数量:

$$\frac{\beta(1-\alpha)}{1-\beta} \times Sample_num \quad (1)$$

这种调整将使少数类别的比例等于 β 。

2) 训练目标: 我们的目标是训练深度模型以预测给定的日志序列是正常还是异常。我们微调模型以做出适当的响应: 如果序列是异常的, 它输出 ‘该序列是异常的。’; 如果是正常的, 它输出 ‘该序列是正常的。’。我们使用交叉熵损失 [54] 作为我们的损失函数。

3) 训练流程: 为了训练我们的深度模型, 我们遵循三个主要阶段。

第一阶段: 微调 Llama 以捕捉答案模板: 第一阶段涉及微调 Llama 以捕捉答案模板。具体来说, 我们训练 Llama 对提示 “这个序列正常还是异常?” 做出 “这个序列是异常/正常的。” 的回应。这一阶段只需要少量数据样本。

阶段 2: 训练日志消息的嵌入器: 第二阶段涉及训练日志消息的嵌入器, 特别是 BERT 和投影器。这一阶段的目的是将每个日志消息投影到 Llama 中最合适的标记的嵌入中, 使 Llama 能够判断给定的日志序列是正常还是异常。

第 3 阶段. 微调整个模型: 最后, 我们将整个模型进行微调, 以确保所有组件的性能一致且准确。

4) 在 LLMs 上的高效微调: 为了降低使用大量参数对 LLMs (BERT 和 Llama) 进行微调的成本, 我们利用 QLoRA [55] 来最小化内存使用。QLoRA 通过将梯度反向传播到冻结的 4 位量化模型中, 同时保持在全 16 位微调过程中达到的性能水平。

V. 实验

在本节中, 我们对四个真实日志进行了广泛的实验, 以研究以下研究问题 (RQs):

RQ1: LogLLM 在基于日志的异常检测中效果如何?

RQ2: 不同的预处理技术如何影响 LogLLM 的性能?

RQ3: Llama 的嵌入器效果如何?

RQ4: Llama 模型的大小如何影响 LogLLM 的性能?

RQ5: 三个阶段训练过程中的每个阶段如何影响 LogLLM 的性能?

RQ6: 由超参数 β 确定的少数类过采样不同级别如何影响 LogLLM 的性能?

LogLLM 是用 Python 编写的, 源代码可在 <https://github.com/guanwei49/LogLLM> 上找到。

A. 基准方法

为了验证所提方法的优越性，我们将 LogLLM 与五种最先进的半监督方法进行了比较：DeepLog [8], LogAnomaly [9], PLELog [22], Fast-LogAD [34], 和 LogBERT [32]。我们还将其与三种监督方法进行了比较：LogRobust [19], CNN [18] 和 NeuralLog [3], 以及一种不需要训练深度模型但需要一些正常样本用于检索的方法：RAPID [39]。

值得注意的是，FastLogAD、LogBERT、NeuralLog 和 RAPID 采用了 LLM 进行异常检测。

B. 实验设置

我们在配备 Intel Xeon Gold 6330 CPU（38 核）、256GB 内存和 NVIDIA A40 GPU（48GB 内存）的服务器上进行了所有实验。

在我们的实验中，我们利用 BERT-base 模型¹ 和 Llama-3-8B 模型² 作为骨干。第 IV-C1 节中描述的超参数 β 设置为 30%。我们使用 AdamW 优化器 [56] 以 16 个样本的小批量大小训练模型。除非另有说明，否则训练过程配置如下：在第一阶段，只涉及 1,000 个样本，学习率为 $5e-4$ 。第二和第三阶段各包含两个 epoch，学习率为 $5e-5$ 。

为了进行公平的比较，我们将所有比较方法的超参数配置为它们原始文章中提供的值。

C. 评估指标

我们使用广泛采用的 *Precision*、*Recall* 和 *F1-score* 来评估这些方法的性能。这些指标的计算方法如下：

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

$$F1-score = \frac{2 * Precision * Recall}{Precision + Recall} \quad (4)$$

, where TP , FN , FP represent true positives, false negatives and false positives respectively.

精确率指的是模型识别出的所有异常中正确检测到的异常所占的百分比，而召回率表示从所有真实异常中正确识别出的异常所占的百分比。F1 分数将这两个指标结合成一个单一指标，为模型在检测异常方面的性能提供平衡评估。

D. 数据集

为了评估我们的基于日志的异常检测方法，我们选择了四个公开数据集 [57]: HDFS、BGL、Liberty 和 Thunderbird。每个数据集的详细信息如下：

HDFS（Hadoop 分布式文件系统）数据集是通过在超过

200 个 Amazon EC2 节点上运行基于 Hadoop 的 mapreduce 作业生成的。这些日志消息根据它们的块 *_id* 分组，这些 id 反映了 HDFS 中的程序执行，从而产生了 575,061 个块。在这些块中，有 16,838 个块（2.93%）表示系统异常。

BGL（Blue Gene/L）数据集是从劳伦斯利弗莫尔国家实验室（LLNL）的 BlueGene/L 超级计算机系统收集的超级计算机系统日志数据集。该数据集包含 4,747,963 条日志消息，每条消息都已被人工标记为正常或异常。其中有 348,460 条日志消息（7.34%）被标记为异常。

Thunderbird 数据集是从桑迪亚国家实验室（SNL）的 Thunderbird 超级计算机获取的公开可访问的日志数据集。该数据集包含正常和异常消息，每条消息都已被人工分类。尽管该数据集包含超过 2 亿条日志消息，但我们关注的是 1000 万条连续日志消息的子集，以提高计算效率。这个子集包括 4,937 条异常日志消息，占总数的约 0.049%。

Liberty 数据集由位于阿尔伯克基的桑迪亚国家实验室（SNL）的 Liberty 超级计算机的系统日志组成。这台超级计算机具有 512 个处理器和 944 GB 的内存，数据集包含超过 2 亿条日志消息。为了提高计算效率，我们采样了 500 万条连续的日志消息，其中 1,600,525 条被识别为异常，占总采样消息的约 32.01%。

在 HDFS 的背景下，我们采用会话窗口策略，该策略涉及根据每个日志消息中存在的块 *_id* 将日志消息分组为序列。每个会话都使用真实标签进行标记。对于其他数据集，包括 BGL、Thunderbird 和 Liberty，我们使用滑动窗口策略来分组日志消息，窗口大小为 100 条消息，步长为 100 条消息。如果一个日志序列包含至少一条根据真实标签标记为异常的日志消息，则认为该日志序列是异常的。

与现有工作 [8], [9], [19], [22], [34], [39]，类似，我们将每个数据集分为训练集和测试集，比例为 8:2，以评估基于日志的异常检测方法的性能。对于 HDFS 数据集，我们随机将日志序列分为训练数据和测试数据。相比之下，对于 BGL、Thunderbird 和 Liberty 数据集，我们遵循时间顺序分割 [6]。这种策略确保训练集中的所有日志序列都早于测试集中的日志序列，反映了现实世界的情况，并减轻了不稳定日志数据可能引起的数据泄露。

表 I 总结了实验中使用的各个数据集的统计数据。

E. 性能评估（RQ1）

表 II 展示了各种基于日志的异常检测方法在 HDFS、BGL、Liberty 和 Thunderbird 数据集上的实验结果。最佳结果以粗体显示。我们观察到以下几点：

¹<https://huggingface.co/google-bert/bert-base-uncased>

²<https://huggingface.co/meta-llama/Meta-Llama-3-8B>

表 I: 实验中使用的数据集的统计信息。

	日志消息	日志序列	训练数据			测试数据		
			# Log sequences	异常	异常比率	# Log sequences	# 异常	Anomaly ratio
HDFS	11,175,629	575,061	460,048	13,497	2.93%	115,013	3,341	2.90%
BGL	4747963	47135	37708	4009	10.63%	9427	817	8.67%
自由	500,000	50,000	40,000	34,144	85.36%	10,000	651	6.51%
Thunderbird	10,000,000	99,997	79,997	837	1.05%	20,000	29	0.15%

T表 2: 在 HDFS、BGL、Liberty 和 Thunderbird 数据集上的实验结果。最佳结果已突出显示 n bold.

方法	数据集	Log parser	HDFS			BGL			自由			Thunderbird			Avg. F ₁
			精度	Rec.	F ₁	精度	Rec.	F ₁	精度	Rec.	F ₁	精度	Rec.	F ₁	
DeepLog		✓	0.835	0.994	0.908	0.166	0.988	0.285	0.751	0.855	0.800	0.017	0.966	0.033	0.506
日志异常		✓	0.886	0.893	0.966	0.176	0.985	0.299	0.684	0.876	0.768	0.025	0.966	0.050	0.521
PLELog		✓	0.893	0.979	0.934	0.595	0.880	0.710	0.795	0.874	0.832	0.808	0.724	0.764	0.810
FastLogAD		✓	0.721	0.893	0.798	0.167	1.000	0.287	0.151	0.999	0.263	0.008	0.931	0.017	0.341
LogBERT		✓	0.989	0.614	0.758	0.165	0.989	0.283	0.902	0.633	0.744	0.022	0.172	0.039	0.456
LogRobust		✓	0.961	1.000	0.980	0.696	0.968	0.810	0.695	0.979	0.813	0.318	1.000	0.482	0.771
CNN		✓	0.966	1.000	0.982	0.698	0.965	0.810	0.580	0.914	0.709	0.870	0.690	0.769	0.818
NeuralLog		✗	0.971	0.988	0.979	0.792	0.884	0.835	0.875	0.926	0.900	0.794	0.931	0.857	0.893
RAPID		✗	1.000	0.859	0.924	0.874	0.399	0.548	0.911	0.611	0.732	0.200	0.207	0.203	0.602
LogLLM		✗	0.994	1.000	0.997	0.861	0.979	0.916	0.992	0.926	0.958	0.966	0.966	0.966	0.959

所提出的 LogLLM 在所有数据集上实现了最高的 F₁-分数。平均而言, LogLLM 的 F₁-分数比现有最佳方法 NeuralLog 高出 6.6%, 证明了其在基于日志的异常检测中的有效性。尽管 FastLogAD、LogBERT、NeuralLog 和 RAPID 等方法采用了 LLM 进行异常检测, 但它们的性能仍然不尽人意。FastLogAD 和 LogBERT 利用基于 BERT 的 transformer 编码器模型, 通过日志序列重建错误来检测异常。它们的输入由日志解析器从日志消息中提取的日志模板 ID (日志字符串模板的 ID) 序列组成, 缺乏语义信息。相比之下, NeuralLog 和 RAPID 利用基于 transformer 编码器的模型从日志消息中提取语义向量。然而, NeuralLog 使用较小的模型, 而 RAPID 则依赖于基于距离的比较进行异常序列分类。另一方面, LogLLM 利用 BERT 提取语义向量, 并利用 Llama (一种基于 transformer 解码器的 LLM) 进行异常检测。通过投影器对齐 BERT 和 Llama 的表示空间, 充分利用 LLM 在基于日志的异常检测中的潜力。

此外, LogLLM 在精确度和召回率之间取得了平衡, 表明它保持了低误报率并最大限度地减少了漏报。相比之下, 像 Fast-LogAD 这样的方法对异常过于敏感, 往往会导致许多误报。例如, 在 BGL 数据集上, 尽管 FastLogAD 的召回率为 1, 但其精确率仅为 0.167, 使其在实际应用中不切实际。类似地, DeepLog、LogAnomaly 和 LogBERT 等方法也存在类似问题。另一方面, RAPID 对异常不够敏感, 导致许多异常未被检测到。例如, 在 BGL 数据集上, RAPID 的精确率为 0.874, 但召回率仅为 0.399。

标签异常的影响: 如图表 II 所示, 在

表 III: 计算成本。

	训练时间 (分钟)	Testing time (Minutes)
DeepLog	72.17	3.42
LogAnomaly	156.16	7.25
PLELog	315.47	33.59
LogRobust	108.42	2.48
CNN	98.16	2.16
FastLogAD	254.17	0.29
LogBERT	429.04	43.77
NeuralLog	267.46	21.44
RAPID	63.98	38.43
LogLLM	1,065.15	64.48

与 DeepLog、LogAnomaly、FastLogAD、LogBERT 和 RAPID 等方法相比, 这些方法需要无异常的干净数据集来构建异常检测模型, 而像 PLELog、LogRobust、CNN、NeuralLog 和 LogLLM 这样的方法则表现出更优越的性能。这些模型不仅使用正常样本进行训练, 还使用标记的异常样本进行训练。例如, 这五种方法在四个数据集上实现了平均 F1-score 超过 0.771, 而那些没有利用标记异常的方法在四个数据集上的平均 F3-score 低于 0.602。这表明, 结合标记的异常可以为异常检测方法提供显著的优点。

计算成本: 每种方法的耗时在表 III 中给出。这些结果是在所有数据集上平均得出的。

尽管 RAPID 不需要训练深度模型, 但向量和表示的提取和检索仍然耗时。与其他方法相比, FastLogAD 需要相对较高的训练时间, 因为它在测试时只使用模型的判别器, 所以测试时间最短。正如预期的那样, 虽然我们提出的 LogLLM 表现出最佳性能, 但它也产生了最高的

表 IV: 不同预处理技术对 HDFS、BGL、Liberty 和 Thunderbird 数据集的影响。最佳结果以粗体显示。

	HDFS				BGL			Liberty			Thunderbird			Avg. F ₁
	准确度	Rec.	F ₁	精度	Rec.	F ₁	精度	Rec.	F ₁	精度	Rec.	F ₁		
Raw	0.994	0.991	0.993	0.943	0.767	0.846	0.911	0.908	0.909	0.806	0.862	0.833		0.895
模板 ID	0.995	0.945	0.969	0.775	0.286	0.418	0.994	0.270	0.425	1.000	0.379	0.550		0.591
模板	0.991	1.000	0.995	0.861	0.919	0.889	0.968	0.931	0.949	0.950	0.655	0.776		0.902
RE (LogLLM)	0.994	1.000	0.997	0.861	0.979	0.916	0.992	0.926	0.958	0.966	0.966	0.966		0.959

TABLEV: 嵌入器 (BERT & adapter) 和 LLaMA 模型大小的影响, 其中 ‘Mem.’ 表示 GPU 内存使用量 (GB), ‘Tim.’ 表示训练时间 (分钟)。(‘-’ 表示内存不足 (OOM) 错误)。

	HDFS					BGL					Liberty					雷鸟				
	Prec.	Rec.	F ₁	Mem.	Tim.	精度	Rec.	F ₁	Mem.	Tim.	精度	Rec.	F ₁	Mem.	Tim.	精度	Rec.	F ₁	Mem.	Tim.
L.-1B	0.986	0.995	0.991	16.5	1022.1	-	-	-	-	-	0.960	0.699	0.809	42.6	443.2	1.000	0.724	0.840	44.5	1732.1
Emb. & L.-1B	0.996	0.996	0.996	8.0	1412.2	0.734	0.944	0.825	32.4	187.1	0.950	0.905	0.927	29.3	173.2	0.875	0.966	0.918	32.4	715.1
L.-8B	0.988	0.997	0.992	43.0	4712.1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
Emb. & L.-8B	0.994	1.000	0.997	16.6	2168.2	0.861	0.979	0.916	38.0	396.2	0.992	0.926	0.958	36.1	412.1	0.966	0.966	0.966	38.2	1284.2

由于其参数数量庞大, 计算成本较高。然而, 与使用 LLM 的其他方法 (如 LogBERT、NeuralLog 和 RAPID) 相比, LogLLM 的测试时间仍然可接受。

F. 不同预处理技术 (RQ2)

我们评估了不同预处理技术的有效性。结果如表 IV 所示。在此表中, ‘Raw’ 表示日志消息的内容未经预处理, 直接输入到所提出的深度模型中。‘Template’ 表示由 Drain (一个日志解析器) 生成的日志模板序列被用作所提出的深度模型的输入。‘TemplateID’ 表示由 Drain 获得的日志模板 ID 被简单地编码为数值向量, 而不是使用 BERT。预处理技术 ‘Template ID’ 使模型无法捕获日志消息中的语义信息。值得注意的是, Drain 解析器应用于整个数据集, 而不仅仅是训练数据集, 以避免由于 OOV 问题导致的性能下降。‘RE’ 表示使用正则表达式进行日志消息的预处理。

正如预期的那样, 预处理技术 ‘RE’ 在所有数据集中产生了最高的 F₁ 分数。相反, 预处理技术

‘Template ID’ 在所有数据集中始终产生最低的 F₁ 分数, 平均比 ‘RE’ 低 36.8%。这可以归因于

‘Template ID’ 阻碍了模型捕获日志消息中语义信息的能力, 从而损害了其从自然语言角度检测异常的能力。预处理技术 ‘Raw’ 和 ‘Template’ 的结果相对较好, 但它们的 F₁ 分数仍然比 ‘RE’ 低 6.4% 和 5.7%。对于预处理技术 ‘Raw’, 每个日志消息内容中的可变部分 (携带动态运行时信息的参数) 对异常检测的影响很小。然而, 由于它们的高度随机性, 它们可能会混淆。

模型, 使得难以辨别异常。对于预处理技术 ‘模板’, 解析器并不总是可靠的, 有时会错误地删除常数部分或保留变量部分, 这可能导致信息丢失或模型混淆, 使得难以辨别异常。

G. Embedder 的影响 (RQ3)

我们研究了嵌入器 (BERT 和适配器) 对 LogLLM 是否必要。结果如表 V 所示。‘L.-1B’ 指的是直接将日志序列 (通过将日志消息以分号 (;) 作为分隔符连接成一个长字符串) 输入到 ‘Llama-3.2-1B’ 模型³中。‘Emb. & L.-1B’ 表示基于 ‘Llama-3.2-1B’ 的 LogLLM。

不出所料, 在嵌入器的帮助下, 模型需要的 GPU 内存更少, 从而避免了内存不足 (OOM) 错误。此外, 它通过阐明序列中消息之间的边界来提高模型性能。这种改进的表示使得 LLM 能够更好地捕捉序列依赖关系。

H. Llama 模型大小的影响 (RQ4)

如表 V 所示, 更大的 LLaMA 模型尺寸会导致更好的性能, 但代价是增加了 GPU 内存使用量和更长的训练时间。

平均而言, 与使用 Llama-3.2-1B 相比, 采用 Llama-3-8B 将 F₁ 分数提高 4.3%, 但增加了 7.7 GB 的 GPU 内存使用量, 并将训练时间延长了 443.2 分钟。

I. 训练过程消融研究 (RQ5)

我们通过消融研究调查了每个训练过程的影响。结果如表 VI 所示, 其中 ‘W/O’ 表示 ‘without’。我们的观察如下:

跳过任何训练阶段都会导致所有数据集上的 F₁ 分数下降, 证明了我们方法的有效性。

³<https://huggingface.co/meta-llama/Llama-3.2-1B>

表 VI: 在 HDFS、BGL、Liberty 和 Thunderbird 数据集上对训练过程的消融研究。最佳结果以粗体显示。

	HDFS				BGL			Liberty			Thunderbird			Avg. F ₁
	准确率	Rec.	F ₁	精确度	Rec.	F ₁	精确度	Rec.	F ₁	精确度	Rec.	F ₁		
无第 1 阶段	0.991	1.000	0.995	0.578	0.971	0.725	0.685	0.290	0.408	0.381	0.828	0.522	0.662	
不包含阶段 2	0.994	1.000	0.997	0.858	0.920	0.888	0.995	0.906	0.949	0.848	0.966	0.903	0.934	
无阶段 1&2	0.992	1.000	0.996	0.853	0.882	0.868	0.995	0.906	0.949	0.897	0.897	0.897	0.927	
无阶段 3	0.993	0.999	0.996	0.704	0.776	0.738	1.000	0.684	0.812	0.958	0.793	0.868	0.854	
LogLLM	0.994	1.000	0.997	0.861	0.979	0.916	0.992	0.926	0.958	0.966	0.966	0.966	0.959	

三阶段训练流程。值得注意的是，没有第一阶段训练会导致最差的表现，所有数据集上平均的 F₁ 分数下降了高达 29.7%。然而，没有第一阶段和第二阶段训练（仅采用第三阶段：对整个模型进行 ne-tuning）仍然可以获得可接受的表现，平均 F₁ 分数仅下降了 3.2%。这表明在训练日志消息的嵌入器（BERT 和 projector）（第二阶段）之前，先对 Llama 进行 ne-tuning 以捕获答案模板（第一阶段）是至关重要的。如果没有第一阶段（即直接训练嵌入器），嵌入器可能会被误导，导致对日志消息的语义捕获不正确，从而导致模型失败。没有第三阶段训练会导致相对较差的表现，平均 F₁ 分数下降了 10.5%。这表明仅仅对 Llama 进行 ne-tuning 和单独训练嵌入器不足以让模型捕获异常模式；对整个模型进行连贯的 ne-tuning 是必要的。没有第二阶段和第一阶段和第二阶段训练也会导致性能下降，平均 F₁ 分数分别下降了 2.5% 和 3.2%。这表明在 ne-tuning 整个模型之前单独训练嵌入器也可以提高性能。这一阶段允许嵌入器为 Llama 生成更好的日志消息语义向量，以便其能够识别异常。

总之，我们提出的基于日志的深度模型三阶段训练流程非常适合我们的异常检测。

J. Impact of Minority Class Oversampling (RQ6)

请注意，训练数据集中的正常样本和异常样本是不平衡的，如表 I 所示。对于 HDFS、BGL 和 Thunderbird 数据集，正常样本的数量多于异常样本。相反，在 Liberty 数据集中，异常样本的数量超过正常样本。如第 IV-C1 节所述，超参数 β 通过过采样来控制少数类的比例，以解决数据不平衡问题。在本节中，我们通过改变其值来研究 β 的影响。图 4 展示了在不同 β 幅度下 LogLLM 在四个数据集上的性能。当 $\beta = 0$ 时，样本不进行过采样；而是直接使用原始数据集进行训练。

如图 4b 所示，对于 HDFS、BGL 和 Thunderbird 数据集，召回率始终在增加，而对于 Liberty 数据集，随着 β 的增加，召回率在下降。这可以归因于，对于 HDFS、BGL 和 Thunderbird 数据集，当 β 增加时，异常样本被过采样，使得模型更容易将样本识别为异常。

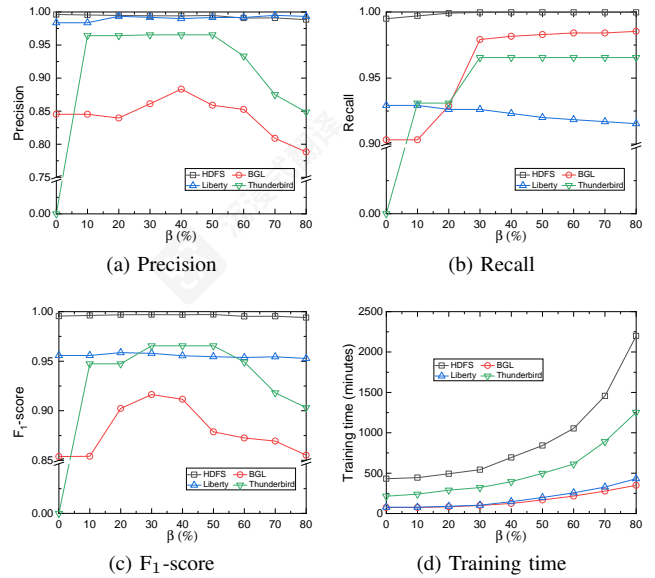


图 4: 少数类过采样的影响。

相比之下，对于 Liberty 数据集，当 β 增加时，正常样本被过采样，使得模型更容易将样本识别为正常。

如图 4c 所示，F₁ 分数的趋势在所有数据集中基本上是相同的。F₁ 分数随着 β 的增加先上升后下降。然而，LogLLM 似乎对 β 不敏感；当 β 在 10% 到 80% 之间时，F₁ 分数的变化不超过 0.07。得益于 LLM 中嵌入的丰富语义知识，训练好的模型可以有效地学习异常模式并检测异常，即使少数类仅占数据集的 10%。然而，LogLLM 似乎无法有效地处理极端不平衡的场景。例如，在 Thunderbird 数据集中，异常样本仅占样本的 1.05%，导致训练好的模型产生偏差，将所有样本分类为正常。因此，精确度、召回率和 F₁ 分数都等于 0。

与 BGL 和 Thunderbird 数据集相比，HDFS 和 Liberty 数据集的精确度、召回率和 F₁ 分数在 β 方面表现出最小变化。这种一致性源于 HDFS 和 Liberty 数据集中异常样本和正常样本之间更明显的模式，使得 LogLLM 能够轻松区分它们，无论比例如何。

正常和异常样本。

正如预期的那样,随着 β 的增加,训练时间也随之增加,如图4d所示。这种关系产生的原因是更高的 β 导致更多的过采样数据样本,如方程(1)所示,从而扩大了训练数据集。

总结来说,少数类过采样是必要的;然而,超参数 β 的值对LogLLM的性能没有显著影响,因此无需仔细选择。此外,过大的 β 值是不理想的,因为它们会导致训练时间延长。30%到50%之间的值被认为是可接受的。

VI. 结论

在本文中,我们提出了LogLLM,这是一种基于日志的、利用LLMs的异常检测新框架。LogLLM使用基于transformer编码器和解码器的LLMs,特别是BERT和Llama,进行基于日志的异常检测。BERT用于从日志消息中提取语义向量,而Llama用于分类日志序列。为了确保日志语义的一致性,我们引入了一个投影器,以对齐BERT和Llama的向量表示空间。LogLLM使用一个创新的三阶段程序进行训练,旨在提高性能和适应性。在四个公开的真实世界数据集上进行的广泛实验表明,LogLLM取得了显著的性能。随后的消融研究进一步证实了我们的三阶段训练程序的有效性。

参考文献

- [1] R. S. Kazemzadeh and H.-A. Jacobsen, ‘可靠且高度可用的分布式发布/订阅服务’, 载于2009年第28届IEEE国际可靠分布式系统研讨会, IEEE, 2009, 第41–50页。
- [2] E. Bauer and R. Adams, 云计算的可靠性和可用性. 约翰·威利与Sons, 2012。
- [3] V.-H. Le and H. Zhang, ‘基于日志的无日志解析异常检测’, 载于2021年第36届IEEE/ACM国际自动软件工程会议(ASE). IEEE, 2021, 第492–504页。
- [4] W. Guan, J. Cao, H. Zhao, Y. Gu and S. Qian, ‘业务流程异常检测综述与基准测试’, IEEE知识和数据工程杂志, 第1–23页, 2024。
- [5] S. Zhang, Y. Ji, J. Luan, X. Nie, Z. Chen, M. Ma, Y. Sun and D. Pei, ‘面向无监督日志异常检测的端到端自动机器学习’, 自动软件工程(ASE’24), 2024。
- [6] V.-H. Le and H. Zhang, ‘基于深度学习的日志异常检测: 我们还有多远?’, 载于第44届国际软件工程会议论文集, 2022, 第1356–1367页。
- [7] J. Qi, S. Huang, Z. Luan, S. Yang, C. Fung, H. Yang, D. Qian, J. Shang, Z. Xiao, and Z. Wu, “Loggpt: Exploring chatgpt for log-based anomaly 基于日志的大语言模型异常检测”, 在2023 IEEE国际高性能计算与通信、数据科学及系统、智能城市与传感器、云及大数据系统与应用(HPCC/DSS/SmartCity/DependSys)会议论文集中, IEEE, 2023, 第273–280页。
- [8] M. Du, F. Li, G. Zheng, and V. Srikumar, 在2017年ACM SIGSAC计算机与通信安全会议论文集中, ‘Deeplog: 通过深度学习从系统日志中进行异常检测与诊断’, 2017, 第1285–1298页。
- [9] W. Meng, Y. Liu, Y. Zhu, S. Zhang, D. Pei, Y. Liu, Y. Chen, R. Zhang, S. Tao, P. Sun 等人, ‘Loganomaly: 无监督检测非结构化日志中的顺序和定量异常’, 在IJCAI会议论文集中, 第19卷第7期, 2019, 第4739–4745页。
- [10] L. Zhang, W. Li, Z. Zhang, Q. Lu, C. Hou, P. Hu, T. Gui, and S. Lu, ‘Logattn: 基于自动编码器的注意力机制的日志无监督异常检测’, 在知识科学、工程与管理国际会议上, Springer, 2021, 第222–235页。
- [11] M. Catillo, A. Pecchia, and U. Villano, “通过系统日志的深度自动编码进行异常检测的Autolog”, 在Expert Systems with Applications期刊中, 第191卷, p. 116263, 2022。
- [12] Y. Xie and K. Yang, “通过图特征融合的对抗性自动编码器进行日志异常检测”, 在IEEE Transactions on Reliability期刊中, 2023。
- [13] 张X, 蔡X, 余M, 邱D, “基于LSTM网络和变分自编码器的日志异常检测模型”, 载于2023年第4届信息科学、并行与分布式系统国际会议(ISPDS), IEEE, 2023, 第239–244页。
- [14] 段X, 英S, 袁W, 程H, 尹X, “用于日志异常检测的生成对抗网络。” 计算机系统科学与工程, 第37卷, 第1期, 第135–148页, 2021年。
- [15] 何Z, 唐Y, 赵K, 刘J, 陈W, “基于对抗训练的图日志异常检测”, 载于国际可靠软件工程研讨会: 理论、工具和应用, Springer, 2023, 第55–71页。
- [16] 张C, 王X, 张H, 张J, 张H, 刘C, 韩P, “基于层次语义的日志序列异常检测: Layerlog”, 应用软件计算, 第132卷, 第109860页, 2023年。
- [17] 哈希米S, 马恩蒂拉M, “Onelog: 面向端到端软件日志异常检测”, 自动化软件工程, 第31卷, 第2期, 第37页, 2024年。
- [18] 卢S, 魏X, 李Y, 王L, “使用卷积神经网络在大数据系统日志中检测异常”, 载于2018年第16届IEEE可靠、自主和安全的计算国际会议, 第16届国际普适智能计算会议, 第4届大数据智能计算和网络安全技术大会(DASC/PiCom/DataCom/CyberSciTech), IEEE, 2018, 第151–158页。
- [19] 张X, 徐Y, 林Q, 乔B, 张H, 当Y, 谢C, 杨X, 程Q, 李Z等. (v1) “基于日志的鲁棒不稳定日志数据异常检测”, 载于v3第27届ACM欧洲软件工程会议和软件工程基础研讨会联合会议论文集(v4), 2019年, 第807–817页。
- [20] 谢Y, 张H, 巴巴尔M.A., “Loggd: 利用图神经网络从系统日志中检测异常”, 载于v1第22届IEEE国际软件质量、可靠性和安全性会议(v2), IEEE, 2022年, 第299–310页。
- [21] 赵Z, 牛W, 张X, 张R, 余Z, 黄C, “Trine: 在生成对抗网络中使用三个Transformer编码器进行syslog异常检测”, 载于v1《应用智能》(v2), 第52卷第8期, 第8810–8819页, 2022年。
- [22] 杨L, 陈J, 王Z, 王W, 江J, 董X, 张W, “通过概率标签估计进行半监督日志异常检测”, 载于v1第43届IEEE/ACM国际软件工程会议(ICSE)(v2), IEEE, 2021年, 第1448–1460页。
- [23] 霍克瑞特S, “长短期记忆”, 载于v1《神经计算》(MIT-Press)(v2), 1997年。
- [24] 瓦斯瓦尼A, “注意力即是所需”, 载于v2《神经信息处理系统进展》(v3), 2017年。
- [25] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat 等人, “Gpt-4技术报告”, arXiv预印本arXiv:2303.08774, 2023年。
- [26] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan 等人, “Llama 3模型群”, arXiv预印本arXiv:2407.21783, 2024年。
- [27] T. GLM, A. Zeng, B. Xu, B. Wang, C. Zhang, D. Yin, D. Rojas, G. Feng, H. Zhao, H. Lai 等人, “ChatGLM: 从glm-130b到glm-4的全系列大语言模型工具”, arXiv预印本arXiv:2406.12793, 2024年。
- [28] W. Guan, J. Cao, J. Gao, H. Zhao, and S. Qian, “Dabl: 使用大语言模型检测业务流程中的语义异常”, arXiv预印本arXiv:2406.15781, 2024年。
- [29] Y. Liu, S. Tao, W. Meng, F. Yao, X. Zhao, and H. Yang, “Logprompt: 面向零样本和可解释日志分析的提示工程”, 在2024 IEEE/ACM第46届国际软件工程会议: 补充会议论文集, 2024年, 第364–365页。
- [30] C. Egersdoerfer, D. Zhang, and D. Dai, “早期探索使用ChatGPT对并行文件系统日志进行基于日志的异常检测”, 在第32届高性能并行和分布式计算国际研讨会, 2023年, 第315–316页。
- [31] 潘杰, 梁伟胜, 以及叶一笛, “Raglog: 基于检索增强生成的日志异常检测”, 载于2024年IEEE公共安全技术世界论坛(WFPST)。IEEE, 2024年, 第169–174页。

- [32] H. Guo, S. Yuan, 和 X. Wu, “Logbert: 基于 BERT 的日志异常检测”, 载于 2021 年国际神经网络联合会议 (IJCNN) . IEEE, 2021, 第 1-8 页。
- [33] Y. Lee, J. Kim, 和 P. Kang, “Lanobert: 基于 BERT 掩码语言模型的系统日志异常检测”, 应用软计算, 第 146 卷, 第 110689 页, 2023 年。
- [34] Y. Lin, H. Deng, 和 X. Li, “Fastlogad: 基于掩码引导的伪异常生成和判别的日志异常检测”, *arXiv* 预印本 *arXiv:2404.08750*, 2024 年。
- [35] C. Almodovar, F. Sabrina, S. Karimi, 和 S. Azad, “Log t: 使用 NE 调优的语言模型的日志异常检测”, *IEEE 网络和服务管理杂志*, 2024 年。
- [36] S. Chen 和 H. Liao, “Bert-log: 基于预训练语言模型的系统日志异常检测”, 应用人工智能, 第 36 卷, 第 1 期, 第 2145642 页, 2022 年。
- [37] J. L. Adebba, D.-H. Kim, 和 J. Kwak, “Sarlog: 通过 BERT 增强的对比学习实现的语义感知鲁棒日志异常检测”, *IEEE 物联网杂志*, 2024 年。
- [38] Y. Fu, K. Liang, 和 J. Xu, “Mlog: 基于语义表示的 Mogri er LSTM 日志异常检测方法”, *IEEE Transactions on Services Computing*, 第 16 卷, 第 5 期, 第 3537-3549 页, 2023 年。
- [39] G. No, Y. Lee, H. Kang, 和 P. Kang, “考虑 token 级信息的预训练语言模型的无监督检索日志异常检测”, *Engineering Applications of Artificial Intelligence*, 第 133 卷, 第 108613 页, 2024 年。
- [40] F. Hadadi, Q. Xu, D. Bianculli, 和 L. Briand, “使用 gptmodels 在不稳定日志上的异常检测”, *arXiv* 预印本 *arXiv:2406.07467*, 2024 年。
- [41] M. Burtsev, M. Reeves, 和 A. Job, “大型语言模型的工作局限性”, *MIT Sloan Management Review*, 第 65 卷, 第 2 期, 第 8-10 页, 2024 年。
- [42] A. Joulin, “Fasttext. zip: 压缩文本分类模型”, *arXiv* 预印本 *arXiv:1612.03651*, 2016 年。
- [43] J. Devlin, “Bert: 用于语言理解的深度双向变换器预训练”, *arXiv* 预印本 *arXiv:1810.04805*, 2018 年。
- [44] 刘宇, “Roberta: 一种鲁棒优化的 BERT 预训练方法”, 预印本 *arXiv:1907.11692*, 2019 年。
- [45] M. Joshi, D. Chen, Y. Liu, D. S. Weld, L. Zettlemoyer, 和 O. Levy, “通过表示和预测跨度来改进预训练: SpanBERT”, 计算语言学协会会刊, 第 8 卷, 第 64-77 页, 2020 年。
- [46] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel 等人, “用于知识密集型 NLP 任务的检索增强生成”, 神经信息处理系统进展, 第 33 卷, 第 9459-9474 页, 2020 年。
- [47] 朱珠, 何思, 刘杰, 何平, 谢奇, 郑振, 和吕明仁, “用于自动日志解析的工具和基准”, 在 2019 IEEE/ACM 第 41 届国际软件工程会议: 软件工程实践 (ICSE-SEIP) 中, IEEE, 2019 年, 第 121-130 页。
- [48] 何平, 朱珠, 何思, 李杰, 和吕明仁, “关于日志解析及其在日志挖掘中的应用的评估研究”, 在 2016 第 46 届 IEEE/IFIP 国际可靠系统与网络会议 (DSN) 中, IEEE, 2016 年, 第 654-661 页。
- [49] Xu W, Huang L, Fox A, Patterson D, 和 Jordan M, “通过挖掘控制台日志模式进行在线系统问题检测”, 在 2009 第 9 届 IEEE 国际数据挖掘会议中, IEEE, 2009 年, 第 588-597 页。
- [50] A. Oliner 和 J. Stearley, “超级计算机说了什么: 对五个系统日志的研究”, 载于 第 37 届 IEEE/IFIP 国际可靠系统与网络会议 (DSN'07) . IEEE, 2007 年, 第 575-584 页。
- [51] P. He, J. Zhu, Z. Zheng 和 M. R. Lyu, “Drain: 一种具有固定深度的在线日志解析方法”, 载于 2017 年 IEEE 国际 Web 服务会议 (ICWS) . IEEE, 2017 年, 第 33-40 页。
- [52] M. Du 和 F. Li, “Spell: 系统事件日志的流式解析”, 载于 2016 年 IEEE 第 16 届国际数据挖掘会议 (ICDM) . IEEE, 2016 年, 第 859-864 页。
- [53] V.-H. Le 和 H. Zhang, “基于提示的少样本学习的日志解析”, 载于 2023 年 IEEE/ACM 第 45 届国际软件工程会议 (ICSE) . IEEE, 2023 年, 第 2438-2449 页。
- [54] J. S. Bridle, “前馈分类网络输出的概率解释, 与统计模式识别的关系”, 载于《神经计算: 算法、架构和应用》. Springer, 1990 年, 第 227-236 页。
- [55] T. Dettmers, A. Pagnoni, A. Holtzman 和 L. Zettlemoyer, “Qlora: 量化 LLM 的有效微调”, 《神经信息处理系统进展》, 第 36 卷, 2024 年。6] I. Loshchilov, “解耦权重衰减正则化”, *arXiv* 预印本 *arXiv:1711.05101*, 2017 年。7] J. Zhu, S. He, P. He, [5] J. Liu 和 M. R. Lyu, “Loghub: 用于 AI 驱动日志分析的大型系统日志数据集”, 载于 2023 年 IEEE 第 34 届国际软件可靠性工程研讨会 (ISSRE) . IEEE, 2023 年, 第 355-366 页。