

面向高质量数据集的全流程自动化评测平台构建研究

蒋亚军

中电数据产业集团有限公司

DOI:10.32629/acair.v3i4.17909

[摘要] 作为人工智能模型可靠性与可复现性的基础,高质量数据集的现有评测体系大多以人工判定模式为主,或仅覆盖静态模态,不符合现阶段模型的跨模态、可复现与合规化要求。本次研究以高质量数据集为研究对象,构建一种面向建设、格式、分类与质量评测四类内容并支持多模态数据的全流程自动化评测平台,采用算子化可执行指标实现可复现的规则化检测,通过结构化证据包与模板化报告实现评测结论的可流转与可审计化。

[关键词] 高质量数据集; 全流程; 自动化评测; 多模态; 算子化; 可审计性

中图分类号: P415.1+3 **文献标识码:** A

Research on Constructing a Full-Process Automated Evaluation Platform for High-Quality Datasets

Yajun Jiang

China Electronics Data Industry Group Co., Ltd.

[Abstract] As the foundation for the reliability and reproducibility of artificial intelligence models, existing evaluation systems for high-quality datasets predominantly rely on manual assessment or cover only static modalities. This approach fails to meet the current requirements for cross-modal capabilities, reproducibility, and compliance in model development. This study focuses on high-quality datasets to construct a fully automated evaluation platform supporting multimodal data across four domains: construction, format, classification, and quality assessment. It employs operatorised executable metrics for reproducible, rule-based detection, while structured evidence packages and templated reports ensure transferable and auditable evaluation conclusions.

[Key words] high-quality datasets; end-to-end; automated evaluation; multimodal; operatorisation; auditability

引言

伴随着数据驱动方法在视觉、语音与语言等领域的广泛应用,数据集质量已成为模型性能与应用风险的重要影响因素。现阶段监管与行业标准已对数据可追溯性、标注规范与下游适配性提出明确要求。基于对现有标准与解读材料的梳理,本文以测试视角设计评测框架与算子库,以期提供一个兼顾自动检测与专家评审、可实现能量化评分并生成数据集评测报告的技术方案。

1 评测目标与指标体系

1.1 平台目标与定位

对于平台目标来说,本平台以期将标准要求转化为可执行的检测项与算子化工具链,实现可重复、可追溯、可审计的评测过程,最终形成符合TC609系列规范的标准化评测结论与证据包,为数据集的合规签发与复现性验证提供有效支撑。定位方面平台将会是覆盖建设评测、格式评测、分类评测与质量评测四类内容的综合自动化评测系统,可面向数据集开展全生命周期的

结构化判定、证据核验与统一报告生成工作^[1]。

1.2 评测体系与判定规则

在评测框架当中,建设评测聚焦数据需求与建设佐证材料的完整性核验,格式评测则需验证元数据字段与表示方法的规范性,类型划分与特征分析的合理性是分类评测的关注要点,质量评测需通过说明文档、数据质量与模型应用三条路径实施“评估+测试”的复合判定。各类评测可产出结构化结论,其中前三类评测直接输出“符合/不符合”标签,第四类评测则是以分数形式呈现评测结果,所有结构化结论均适配证据清单与可执行整改项。合格规则由前三类评测的符合性检查和第四类评测的质量分数红线共同构成,凡必检项不满足或任一质量维度未达到九十分可直接判定不合格并触发整改流程。平台可在决策层统一收敛四类评测结果,通过全链路快照留存,每一判定均可回溯至具体指标、算子输出或证明材料,形成可量化、可审计的综合评价结论。

2 多模态与分项评测方法

2.1 多模态覆盖原则与统一接口

平台设计应保证多模态数据的完整性、可追溯性和可复现性,支持静态图像、文本的细粒度检测以及长视频、音频流等高成本模态的抽样检测和证据保留。为实现跨模态一致的调用与结果表达,需构建统一的数据接入与预处理接口。数据入库阶段可生成标准化清单与结构化元数据,预处理模块负责按模态执行解码、抽帧、声学预处理与文本规范化等操作,算子执行层通过统一的API接口完成具体的数据处理任务(接收 data_id、modality、metadata 与配置项,返回分数、证据片段、置信度与建议项),评测管线支持批量、流式两种模式并兼容样本化检查与并行加速,确保检测结果既可自动复现又便于人工复核与归档^[2]。

2.2 模态关键质量与格式项举例

为确保数据的可用性和下游适配性,平台可针对不同模态定义若干关键质量项。格式规范、元数据完整性与语义一致性是文本数据的关键指标;图像数据则需要关注分辨率、噪声/模糊、标注一致性与去重;音频数据的采样率、信噪比、转录一致性与说话人标注准确性均处于检测范围;考虑到视频数据的特殊性,质量项需覆盖帧完整性、帧率与分辨率一致性、运动模糊/曝光问题、时序标注一致性、关键帧覆盖率与音视频同步等^[3]。对数据标识、关联标识、原始时间、最后修改时间、模态类型、内容字段与标注方式等元数据字段的合规性检查属于格式层面的要求,易量化的格式规范性、安全性、结构完整性与内容一致性等项优先采用自动化算子检测,而需要语义判断的标注规范性、内容真实性与类型一致性则通过专家评估或抽样复核确定。

2.3 建设与分类专用检测项

建设评测的侧重点在于对数据需求佐证、数据规划与建设过程材料的核验,需求范围与可用性证明、数据质量模型与工作量计划等文档证据的完整性与一致性检查均属于检测项。分类评测主要关注类型要素覆盖与类型特征分析的充分性,检验分类方法是否按“行业专识→行业通识→通识”这一顺序合理落位及其特征描述是否支持该判定。为降低工具化实现的难度,将上述建设与分类项映射为两类可执行检查:材料合规性算子用于验证文件存在性、字段一致性与版本控制证据,特征一致性算子用于对类型要素的统计与规则匹配^[4]。在检测流程方面,采用文档审查与自动化抽样相结合的策略,评测报告需列出未通过项的证据清单与可操作的改进路径,确保建设/分类判定既有文本依据又具备可复现的测试痕迹。

3 算子化实现与评测工具箱

3.1 算子定义与注册

从概念上说,算子被定义为面向单一质量评估的最小可执行单元,核心要素涉及输入规范、可配置参数集、输出结构以及版本元数据,注册目录需对算子的分类标签进行登记(建设类、格式类、分类类、质量类),同时为便于调度与治理,还需要登记依

赖声明、资源需求。注册流程需对算子参数化、权限控制与元数据检索提供支持,还可以在算子实例化视图中提供运行示例、输入-输出格式说明与测试用例,确保算子具备可复用性、可审计性与可治理性^[5]。

3.2 执行引擎、采样与组合策略

执行引擎的核心在于插件化调度,执行框架支持批处理与流式两种模式,考虑到适配CPU/GPU异构环境的问题,框架还可内置任务队列、并行调度与资源感知调度器。在算子的运行过程中,可通过沙箱化容器保证依赖隔离与安全边界,为降低重复计算开销,执行引擎还需承担输入预处理、缓存策略与证据片段裁剪职能。对于高成本模态来说,为平衡覆盖与资源消耗,抽样策略可采用分层抽样、异常驱动抽样与代表性片段抽样等方法。平台遵循分层聚合原则实现算子组合,首要任务是在原子级合并多算子输出形成指标得分,后续则是在维度级执行权重加权与冲突解析,其中权重可由专家经验初始化并参考历史评测数据完成校准,最终采用贝叶斯或加权投票方法进行置信度汇总以表达不确定性并驱动人工复核触发规则。

3.3 基准评测与模型适配性协议

基准评测遵循TC609规范提出的双轨测试方案,也就是低成本的外挂知识库(RAG)式验证与基于训练的通用测试。两种方案均采用交叉验证或K-折划分对数据的下游适配性进行评估,单模态与多模态代表性任务均处于基准模型清单的覆盖范围并以公开或行业认可的模型作为参考,可选择召回、精确率、F1等经典指标开展性能评价工作,相关指标可用于量化理解类与生成类任务的表现。评测流程要求明确基准模型选择理由、训练/测试分割策略、目标性能阈值与可复现的运行环境描述,所有模型适配性结论仅作为数据集适用性的测试证据的一部分,决策层仍需结合建设、格式与分类的符合性判定得出最终结论。

4 平台实现、报告生成与评测流程

4.1 平台模块与评测流水线

模块化架构是平台支持可扩展的评测能力与可维护的运维流程的重要基础,主要模块包括数据接入与清单管理、模态预处理与特征抽取、算子执行引擎、聚合评分与决策引擎、专家复核与交互面板、报告生成与证据打包、审计与存档子系统以及运维监控与回归测试模块。多源数据的标准化入库与manifest生成由数据接入模块负责;解码、抽帧、ASR转写、文本规范化等模态特定流程则是在预处理模块执行,此模块需产生可供算子消费的中间表示;算子执行引擎以插件化方式调度检测单元并输出结构化证据;各维度与总体得分由聚合评分引擎进行计算,计算结束后可驱动合格判定;专家复核面板则提供对自动判定的审查、编辑与签发服务;报告生成子系统可根据结构化数据生成可视化报告和机器可读的评测快照,并以证据包形式将所有输入快照、算子版本与运行日志存档,便于日后复查与复现。评测流水线以数据宣告和材料核验为起点,依次经过格式与建设项的符合性检测、分项算子化质量测试、抽样复检与人工审查,决策引擎在合成评定结论后先生成结构化评测报告并触

发专家复核与整改流程,在所有必检项满足且经专家签发或复测确认后,平台发布正式证书,可回溯的评测链路离不开运行环境、算子版本与输入快照的完整记录的大力支持^[6]。

4.2 评测报告模板与证据包

评测报告模板需要以兼顾机器解析与人工可读为设计目标,包含封面与元信息、总体评分摘要、四类评测逐项结果、算子级诊断表与证据索引、问题清单与整改建议清单以及合格性结论与签发记录。对于证据包来说,可按照模块化格式对示例帧/时间段、原始文档副本、元数据快照与运行日志进行打包,同时还需生成机器可读的JSON快照便于再验证,支持多种导出格式并包含可验证签名与版本信息的报告与证据包能够在认证与交换场景中扮演可证明的测试凭证这一重要角色^[7]。

4.3 报告自动化与可审计性

在规则化映射表与模板化文字生成组件的帮助下,算子诊断能够被映射为可执行整改项并填充至报告字段。为提高报告的可用性,专家可依托系统提供的编辑草稿与人工修订通道将修订结果写回元数据并触发算子/权重的校准流程,最终实现自动化输出与人工判断的闭环优化。报告的可审计性离不开端到端快照(包含数据切片、算子版本、配置参数、执行日志与专家批注)与不可篡改的审计链的大力支持,评测结果支持差异化回溯与再评估,报告需明确披露抽样策略、置信区间与触发人工复核的阈值,确保相关结论在法务或监管场景下具有可解释性、可复查性。

4.4 合规、隐私与对外交换

相关人员可围绕数据最小暴露、可审计性与法规映射构建安全与合规体系,为将评测结论与规范条款对齐,体系需包含敏感信息检测与脱敏算子、基于角色的访问控制、端到端传输与存储加密、审计日志与访问溯源以及合规映射模块。合规对接模块可通过内置的规范条目与判定逻辑确保操作的合规性,报告需明确陈述与规范对应的证据项与合格断言,面对涉及隐私或法律风险的数据集可以自动触发合规评估流程并限制报告的

公开粒度。为满足监管要求与审计需求,可采取最小保留策略管理运行环境和证据包,还需要记录保留期限与销毁日志。对于跨机构的评测报告交换行为来说,平台可借助数字签名和权限控制等安全措施把控敏感数据与合规性^[8]。

5 结语

总之,本文构建的算子驱动评测体系与自动化流水线将建设、格式、分类与质量四类评测纳入统一判定框架,可实现从材料核验到下游适配性测试的闭环处理,希望能够为相关行业人员提供参考。

[参考文献]

[1]王诚文,董青秀,穗志方,等.自然语言处理评测数据集质量评估研究[J].中文信息学报,2023,37(2):15.

[2]丁浩,张畅,顾乐,等.面向AI全生命周期的高质量数据集评测体系研究[J].数字化转型,2025,2(8):87-97.

[3]代婕.数字经济背景下数据质量评估模型研究[J].电子商务评论,2025,14(6):2589-2599.

[4]蔡莉,朱扬勇.从数据质量到数据产品质量[J].大数据,2022,8(3):26-39.

[5]姜磊,张涛.GSBERT:一种基于可视化解释的数据标注自动检测方法实证[J].图书情报工作,2025,69(11):111-122.

[6]赵志君,庄馨予.中国人工智能高质量发展:现状,问题与方略[J].改革,2023(9):11-20.

[7]燕江依,李荪,樊威,等.新一代数据标注产业对“人工智能+”范式创新的作用机理与实践路径研究[J].信息技术与政策,2025,51(8):26.

[8]张何灿,易成岐,郭鹏,等.高质量AI数据体系面临的数据版权困境,应对策略解析与实施路径研究[J].农业图书情报学报,2025,36(9):32-43.

作者简介:

蒋亚军(1988--),女,汉族,河南人,中级工程师,研究生,单位:中电数据产业集团有限公司,研究方向:高质量数据集评测。