

Domain-Adaptive Pre-Training vs. Instruction Fine-Tuning for Cross-Lingual Retrieval Robustness

Assignee Research

June 11, 2026

Abstract

Large language models (LLMs) are increasingly used to access legal information. Yet, their deployment in multilingual legal settings is constrained by unreliable retrieval and the lack of domain-adapted, open-embedding models. In particular, existing multilingual legal corpora are not designed for semantic retrieval, and PDF-based legislative sources introduce substantial noise due to imperfect text extraction. To address these challenges, we introduce LEMUR, a large-scale multilingual corpus of EU environmental legislation constructed from 24,953 official EUR-Lex PDF documents covering 25 lan

1 Introduction

This paper examines: LEMUR: A Corpus for Robust Fine-Tuning of Multilingual Law Embedding Models for Retrieval. Research question: How does domain-adaptive pre-training on legal corpora compare to instruction fine-tuning for improving cross-lingual retrieval robustness against adversarial perturbations in multilingual embedding models?.

2 Methodology

Systematic literature search across multiple databases yielded 10 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.5/10.

3 Results

10 papers retrieved. 10 claims extracted; 7 independently verified. Quality review score: 7.5/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
CaseHOLD is a multiple-choice benchmark comprising more than 53,000 U.S. case-law holdings.	✓	0.17
LeCaRD is a large-scale case-retrieval dataset for the Chinese criminal law system with expert-designed relevance criterion	✓	0.26
The LEMUR corpus is constructed from official legislative PDFs and targets downstream embedding-model fine-tuning across	×	0.15
SAILER and DELTA are structure-aware models that capture section-level or structural dependencies to improve legal case	✓	0.17
SM-BERT-CR and REAKASE-8B incorporate supporting-relation modeling and reasoning-driven representations.	✓	0.18
LEXCLIPR enables paragraph-level retrieval across ECtHR judgments, showing that off-the-shelf multilingual encoders struggle	✓	0.23
Domain-specific pretraining consistently improves legal NLP tasks.	×	0.14
The LEMUR corpus comprises 1,174 distinct legal acts from 1961–2025, available in 25 official EU languages.	✓	0.21
The LEMUR corpus is converted from original PDF files into a structured and machine-readable text format.	✓	0.17
The LEMUR corpus includes high-quality query–document pairs used in experiments.	×	0.14

References

- <http://arxiv.org/abs/2509.22472v1>

- <http://arxiv.org/abs/2602.09570v1>
- <http://arxiv.org/abs/2310.05276v1>