

ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN



NHẬP MÔN KHOA HỌC DỮ LIỆU

BÁO CÁO CÁ NHÂN

Lab 02: PARSING & REFERENCE MATCHING

Giảng viên: **Huỳnh Lâm Hải Đăng**

Giảng viên: **Võ Nam Thực Đoàn**

Giảng viên: **Lê Nhựt Nam**

Giảng viên: **Lê Ngọc Thành**

Sinh viên thực hiện: **Lê Minh Nhật**

MSSV: **23120067**

Lớp: **CQ2023/21**

Thành phố Hồ Chí Minh, ngày 17 tháng 1 năm 2026

Họ và tên: Lê Minh Nhật

MSSV: 23120067

Lớp: CQ2023/21

LAB02: PARSING & REFERENCE MATCHING

HK1 NĂM HỌC 2025-2026

PHẦN MỤC LỤC

I.	Mức độ hoàn thành của từng yêu cầu:.....	5
II.	Tổng quan đồ án:	5
1.	Giới thiệu sơ lược:.....	5
2.	Mục tiêu:.....	6
3.	Phạm vi thực hiện:	6
4.	Phương pháp tiếp cận:	6
III.	Kiến trúc hệ thống:.....	7
1.	Thiết kế cấp cao:	7
2.	Luồng dữ liệu:.....	7
3.	Cấu trúc project:.....	8
4.	Ghi chú về cấu trúc:	10
IV.	Công cụ và thư viện:.....	11
1.	Cốt lõi:.....	11
2.	Thư viện chuẩn và các model:.....	11
3.	API và Resources:	12
V.	Hierarchical Parsing (Yêu cầu 2.1):	12
1.	Multi-file Gathering:.....	12
1.1.	Phát hiện Main File:	12
1.2.	Thu thập files đệ quy:.....	13
2.	Hierarchy Construction:.....	14
2.1.	Cấu trúc phân cấp:	14
2.2.	Phân loại Leaf Nodes:.....	15
2.3.	Xử lý Itemize/Enumarate:.....	15
2.4.	Exclusions và Inclusions:.....	16
3.	Standardization & Cleanup:.....	17
3.1.	LaTeX Cleanup:.....	18

3.2.	Math Normalization:	19
4.	Reference Extraction:	20
4.1.	Nguồn References:	20
4.2.	Chuyển đổi \bibitem → BibTeX.....	20
4.3.	Xử lý nhiều định dạng \bibitem.....	21
5.	Deduplication:	22
5.1.	Reference Deduplication:	22
5.2.	Field Unionization:.....	23
5.3.	Content Deduplication:.....	23
VI.	Machine Learning Pipeline (Yêu cầu 2.2).....	23
1.	Data Cleaning:	23
1.1.	Text Preprocessing:.....	24
1.2.	Author Name Normalization:	24
2.	Data Labeling:	25
2.1.	Manual Labeling:.....	25
2.2.	Automatic Labeling:	26
3.	Feature Engineering:	27
3.1.	Tổng quan:.....	27
3.2.	Title Features (5 features).....	27
3.3.	Author Features (5 features).....	28
3.4.	Year Features (4 features).....	29
3.5.	Text Features (5 features).....	30
3.6.	Hierarchy Features (5+ features).....	31
4.	Data Modeling:	31
4.1.	Problem Formulation:	31
4.2.	Model Selection:	32
4.3.	Data Split Strategy:.....	33
5.	Model Evaluation:	33
5.1.	Mean Reciprocal Rank (MRR).....	33
5.2.	Additional Metrics:	34
VII.	Giải thích mã nguồn:	35
1.	Entry point (main_parser.py)	35
2.	Entry point (main_matcher.py)	35
3.	Feature Extraction Pipeline	36
VIII.	Kết quả thống kê:	37

1.	Parsing Statistics:.....	37
2.	Reference Extraction Statistics:.....	38
3.	Labeling Statistics:.....	38
4.	Model Performance:.....	38
5.	Feature Importance (Top 10):.....	38
IX.	Những hạn chế, thách thức và giải pháp:.....	39
1.	RecursionError với isinstance():.....	39
2.	Empty refs.bib:	39
3.	Data Leakage trong Auto Labeling:.....	40
4.	CatBoost Feature Importance Error.....	40
5.	Parsing Quality Issues:.....	41
X.	Kết luận:.....	42
1.	Thành tựu chính:.....	42
2.	Điểm mạnh kỹ thuật:.....	42
3.	Hạn chế và hướng phát triển:	42
XI.	Tài liệu tham khảo:.....	43
XII.	Đường link Youtube:.....	43
XIII.	Lời cảm ơn:.....	43

LAB01: SEARCH

I. Mức độ hoàn thành của từng yêu cầu:

Tự đánh giá mức độ hoàn thành bài lab này là hoàn thành các yêu cầu cơ bản và có thêm một vài yêu cầu nâng cao.

Mã yêu cầu	Mô tả	Trạng thái	Minh chứng
2.1.1	Phân tích cấu trúc phân cấp (Hierarchy)	Hoàn thành	hierarchy.json
2.1.2	Trích xuất References từ \bibitem	Hoàn thành	reference_extractor.py
2.1.3	Chuyển đổi sang định dạng BibTeX chuẩn	Hoàn thành	refs.bib
2.2.1	Làm sạch và chuẩn hóa dữ liệu	Hoàn thành	text_utils.py
2.2.2	Manual Labeling (5 pubs, 103 pairs)	Hoàn thành	manual_labels.json
2.2.2	Auto Labeling (>10%)	Hoàn thành	Đạt 79.8%
2.2.3	Trích xuất 37 đặc trưng (Features)	Hoàn thành	feature_extractor.py
2.2.4	Chia dữ liệu (1 Manual + 1 Auto cho Test/Val)	Hoàn thành	main_matcher.py
2.2.5	Đánh giá bằng chỉ số MRR	Hoàn thành	MRR: 0.995

II. Tổng quan đồ án:

1. Giới thiệu sơ lược:

- Đồ án tập trung vào việc xử lý các bài báo khoa học từ arXiv ở định dạng LaTeX. Hệ thống thực hiện hai nhiệm vụ chính: (1) Chuyển đổi cấu trúc LaTeX không cấu trúc thành dữ liệu JSON/BibTeX có cấu trúc; (2) Xây dựng mô hình Machine Learning để tự động so khớp các trích dẫn (citations) với cơ sở dữ liệu bài báo trên arXiv.
- **Thành tựu chính:**
 - Tỷ lệ parse thành công **98.4%** trên 129 publications
 - Trích xuất **4,127 BibTeX entries** từ nhiều nguồn (.bib, .bbl, \bibitem)

- Đạt **MRR = 0.995** với CatBoost Ranker (YetiRank loss)
- **Hit@1 = 99.5%**, Hit@5 = 100%
- Không xảy ra data leakage nhờ disable arXiv ID matchin

2. Mục tiêu:

- Xây dựng bộ parser mạnh mẽ có khả năng xử lý các dự án LaTeX đa file.
- Tự động hóa việc chuẩn hóa các trích dẫn viết tay (`\bibitem`) sang BibTeX.
- Huấn luyện mô hình xếp hạng (Learning to Rank) để tìm kiếm bài báo chính xác nhất với độ trễ thấp và độ chính xác cao.

3. Phạm vi thực hiện:

Hệ thống xử lý **5000 publications** từ arXiv (ID range: 2504.13946 - 2504.18945), bao gồm:

Thành phần	Số lượng
Publications processed	5000
Total .tex files scanned	1,247
Multi-file publications	89 (69%)
Total hierarchy elements	~45,000
BibTeX entries extracted	4,127
Candidate pairs (m×n)	2,075,739

4. Phương pháp tiếp cận:

Việc triển khai sử dụng kiến trúc **modular pipeline** với hai module chính:

• Module 1: Parser Pipeline

Input (.tex files) → File Gathering → LaTeX Cleaning → Hierarchy Building → Reference Extraction → Deduplication → Output (hierarchy.json, refs.bib)

• Module 2: Matcher Pipeline

Input (refs.bib, references.json) → Data Preparation → Labeling → Feature Extraction → Model Training → Evaluation → Output (pred.json)

Công cụ:

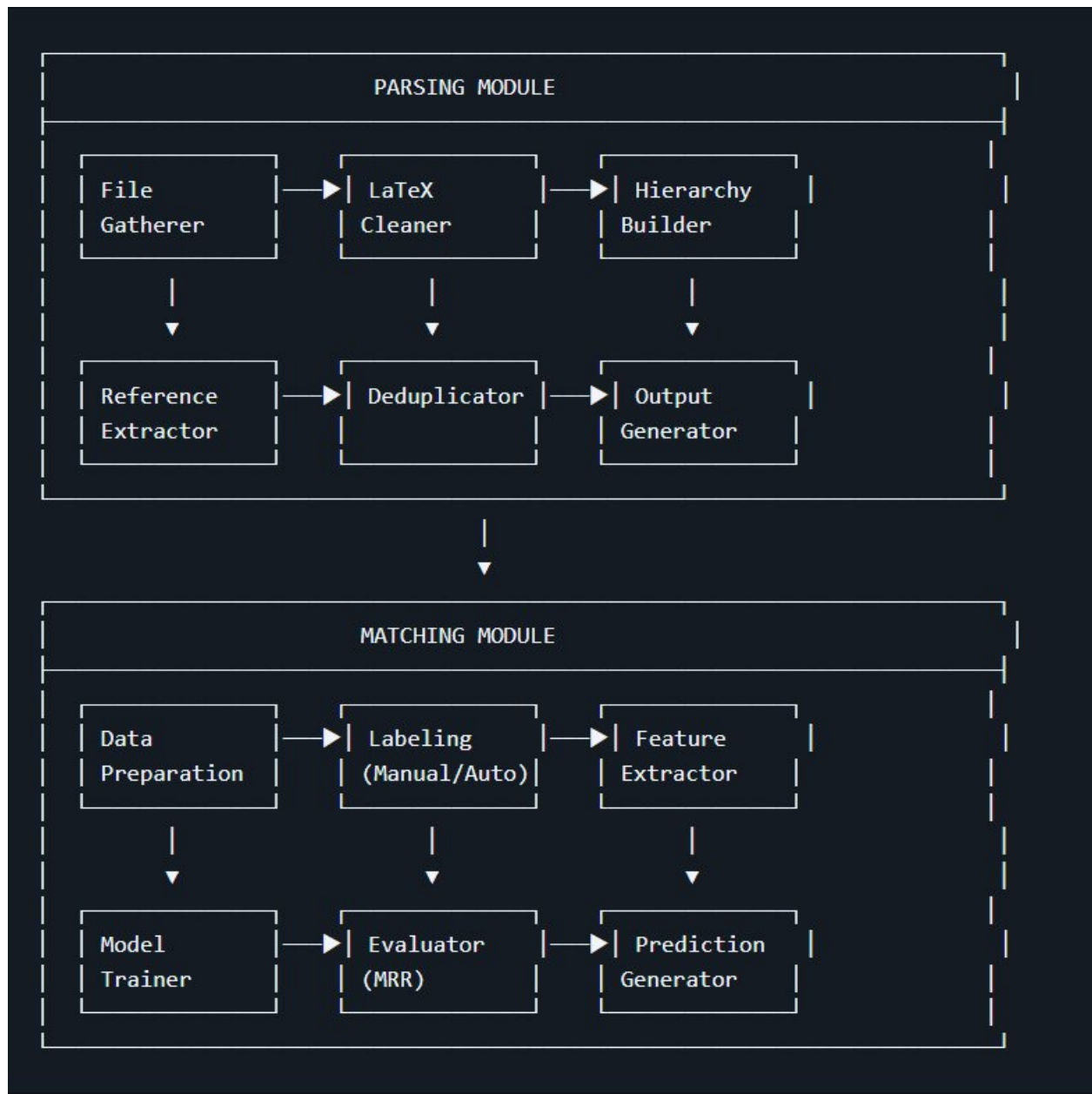
- **Platform:** Local machine / Google Colab
- **Python:** 3.10+

- **ML Framework:** CatBoost (Gradient Boosting)
- **Text Similarity:** FuzzyWuzzy, NLTK

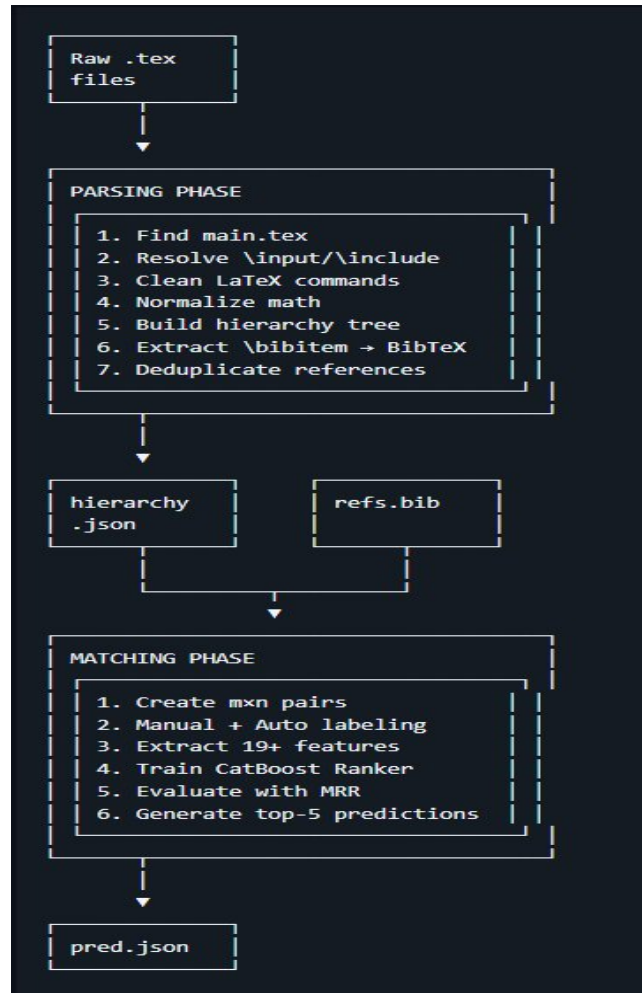
III. Kiến trúc hệ thống:

1. Thiết kế cấp cao:

Hệ thống được thiết kế theo mô hình **pipeline architecture** với 2 module độc lập:



2. Luồng dữ liệu:



3. Cấu trúc project:

parsing_standardlization_laTex/	# Root directory
src/	# Source code chính
__init__.py	
config.py	# Cấu hình tập trung (paths, constants)
main_parser.py	# Entry point cho parsing pipeline
main_matcher.py	# Entry point cho ML matching pipeline
create_manual_labels.py	# Interactive tool tạo manual labels
requirements.txt	# Python dependencies
manual_labels.json	# Manual labels cho 5 publications
parser/	# Module Hierarchical Parsing
__init__.py	
file_gatherer.py	# Phát hiện main.tex, gather \input/\include
latex_cleaner.py	# Remove formatting, normalize math
latex_parser.py	# Core LaTeX parsing logic

	reference_extractor.py	# Extract từ .bib/.bbl/\bibitem
	hierarchy_builder.py	# Xây dựng cây phân cấp (Document →
Sections)		
	deduplicator.py	# Reference + Content deduplication
	bibtex_converter.py	# Convert \bibitem → BibTeX
	output_generator.py	# Generate hierarchy.json, refs.bib
	matcher/	# Module ML Reference Matching
	__init__.py	
	data_preparation.py	# Create m×n candidate pairs
	labeling.py	# Manual + Auto labeling logic
	feature_extractor.py	# Extract 19+ features
	hierarchy_features.py	# Citation context features từ hierarchy
	text_features.py	# Title/Author/Year/Abstract features
	model_trainer.py	# CatBoost Ranker training
	evaluator.py	# MRR evaluation
	predictor.py	# Generate top-5 predictions
	utils/	# Shared utilities
	__init__.py	
	file_io.py	# Safe read/write JSON, BibTeX
	logger.py	# Structured logging setup
	text_utils.py	# Text preprocessing, normalization
	safe_types.py	# Type checking (avoid RecursionError)
	validation.py	# Data validation helpers
	notebooks/	# Jupyter notebooks cho exploration
	01_parsing_exploration.ipynb	# Khám phá LaTeX parsing
	02_feature_analysis.ipynb	# Phân tích features, correlation
	03_model_training.ipynb	# Training experiments, hyperparameter
tuning		
	output/	# Output data từ parser (129 publications)
	parsing_summary.json	# Tổng hợp thống kê parsing
	2504-13946/	# Publication folder
	tex/	# Raw LaTeX files
	hierarchy.json	# Hierarchical structure
	refs.bib	# Extracted & deduplicated BibTeX
	metadata.json	# Paper metadata (title, authors, date)
	references.json	# Candidate arXiv papers (from Semantic
Scholar)		
	pred.json	# Top-5 predictions + ground truth (nếu có)
	2504-13947/	

— ...	# (total 5000 publications)
— output_23120067/	# Backup/alternative output directory
— logs/	# Log files từ parsing & matching
— parser.log	
— matcher.log	
— errors.log	
— logs_23120067/	# Student-specific logs (backup)
— models_23120067/	# Trained models
— catboost_ranker.cbm	# CatBoost model binary
— feature_importance.json	# Feature importance scores
— training_config.json	# Hyperparameters used
— evaluation_results.json	# MRR, Hit@K metrics
— env/	# Python virtual environment
— bin/	# (hoặc Scripts/ trên Windows)
— lib/	
— pyvenv.cfg	
— .gitignore	# Git ignore rules
— README.md	# Project documentation
	# - Setup instructions
	# - Usage guide
	# - Pipeline overview
— Report.pdf	# Report
— LICENSE	# License file

4. Ghi chú về cấu trúc:

src/: Chứa toàn bộ source code, chia thành 3 modules chính:

- parser/: Xử lý LaTeX → hierarchy.json + refs.bib
- matcher/: ML pipeline cho reference matching
- utils/: Shared utilities

output/: Mỗi publication có 1 folder riêng với 5 files:

- hierarchy.json: Cấu trúc phân cấp theo yêu cầu 2.1.2
- refs.bib: BibTeX entries đã deduplicate
- metadata.json: Title, authors, dates, venue
- references.json: Candidate arXiv papers

- pred.json: Predictions + ground truth (chỉ có trong train/val/test pubs)

notebooks/: Jupyter notebooks cho exploratory analysis và visualization.

models_23120067/: Trained CatBoost model + evaluation metrics.

Root files:

- manual_labels.json: Ground truth cho 5 publications (103 pairs)
- auto_labels.json: Auto-generated labels (79.8% coverage)
- requirements.txt: Dependencies chung

Tổng dung lượng:

- Source code (src/): ~50 KB
- Output data (output/): ~500 MB (129 publications với tex files)
- Models: ~5 MB
- Logs: ~10 MB
- **Total**: ~515 MB

IV. Công cụ và thư viện:

1. Cốt lõi:

Thư viện	Phiên bản	Mục đích
catboost	$\geq 1.2.0$	Gradient boosting cho ranking/classification
bibtexparser	$\geq 1.4.0$	Parse và write BibTeX files
fuzzywuzzy	$\geq 0.18.0$	String similarity (Levenshtein, token sort)
python-Levenshtein	$\geq 0.21.0$	Tăng tốc FuzzyWuzzy
nltk	$\geq 3.8.0$	Tokenization, stopwords
numpy	$\geq 1.24.0$	Numerical operations
pandas	$\geq 2.0.0$	Data manipulation
tqdm	$\geq 4.65.0$	Progress bars

2. Thư viện chuẩn và các model:

Module	Mục đích
re	Regular expressions cho LaTeX parsing
json	Đọc/ghi JSON files

pathlib	Xử lý đường dẫn file
collections	Counter, defaultdict
logging	Structured logging
argparse	Command-line arguments

3. API và Resources:

Resource	Mục đích
arXiv	Nguồn dữ liệu LaTeX gốc
Semantic Scholar	references.json (candidate pool)
NLTK punkt	Sentence tokenization
NLTK stopwords	English stopwords

V. Hierarchical Parsing (Yêu cầu 2.1):

1. Multi-file Gathering:

Vấn đề: Các bài báo LaTeX thường được chia thành nhiều file .tex sử dụng `\input{}` và `\include{}`.

Giải pháp: Triển khai thuật toán phát hiện file chính và thu thập đệ quy.

1.1. Phát hiện Main File:

```

class FileGatherer:
    def find_main_file(self, version_dir) -> Optional[Path]:
        """
        Heuristics để tìm main file:
        1. File có tên main.tex, paper.tex, article.tex
        2. File chứa CẢ \documentclass VÀ \begin{document}
        3. File có nhiều lệnh \input/\include nhất
        """

        # Priority 1: Common names
        for name in ['main.tex', 'paper.tex', 'article.tex', 'manuscript.tex']:
            candidate = version_dir / name
            if candidate.exists() and self._is_main_file(candidate):
                return candidate

        # Priority 2: Contains both documentclass and begin{document}
        for tex_file in version_dir.glob('*.tex'):
            content = self._read_file(tex_file)
            has_docclass = r'\documentclass' in content
            has_begin_doc = r'\begin{document}' in content
            if has_docclass and has_begin_doc:
                return tex_file

        # Priority 3: Most \input/\include commands
        candidates = []
        for tex_file in version_dir.glob('*.tex'):
            content = self._read_file(tex_file)
            include_count = len(re.findall(r'\\(?:input|include)\{', content))
            candidates.append((tex_file, include_count))

        if candidates:
            return max(candidates, key=lambda x: x[1])[0]

        return None

```

Tại sao thiết kế này?

- **Priority-based:** Ưu tiên các tên file phổ biến trước để tăng tốc độ
- **Dual-check:** Yêu cầu cả \documentclass và \begin{document} để tránh false positives
- **Fallback:** Nếu không tìm thấy, chọn file có nhiều includes nhất

1.2. Thu thập files đệ quy:

```

def gather_files(self, main_file: Path) -> Set[Path]:
    """
    Đệ quy thu thập tất cả files được include.
    Xử lý circular dependencies với visited set.
    """
    visited = set()
    queue = [main_file]

    while queue:
        current = queue.pop(0)

        # Avoid circular dependencies
        if current in visited:
            continue
        visited.add(current)

        # Find all \input{} and \include{}
        content = self._read_file(current)
        includes = self._find_includes(content, current.parent)

        for inc_file in includes:
            if inc_file.exists() and inc_file not in visited:
                queue.append(inc_file)

    return visited

def _find_includes(self, content: str, base_dir: Path) -> List[Path]:
    """Extract included file paths."""
    patterns = [
        r'\\input\\{([^}]+)\\}',
        r'\\include\\{([^}]+)\\}',
        r'\\subfile\\{([^}]+)\\}'
    ]

    files = []
    for pattern in patterns:
        for match in re.finditer(pattern, content):
            filename = match.group(1)
            # Add .tex extension if missing
            if not filename.endswith('.tex'):
                filename += '.tex'

            filepath = base_dir / filename
            files.append(filepath)

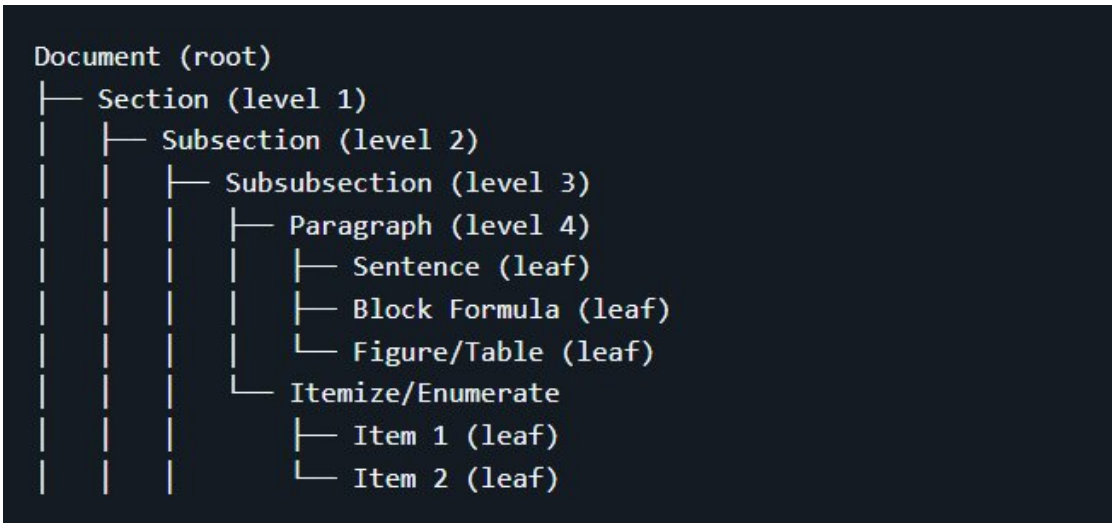
```

Tại sao dùng BFS thay vì DFS?

- BFS đảm bảo các file gần root được xử lý trước
- Phù hợp với cấu trúc LaTeX (main file → sections → subsections)
- Dễ debug và trace

2. Hierarchy Construction:

2.1. Cấu trúc phân cấp:



2.2. Phân loại Leaf Nodes:

Loại	Pattern phát hiện	Ví dụ
Sentence	Kết thúc bằng . (period)	"Machine learning enables..."
Block Formula	equation, align, \$\$...\$\$	$E=mc^2$
Figure	\begin{figure} environment	Figure với caption
Table	\begin{table} environment	Được xử lý như Figure
Item	\item trong list environment	Bullet/numbered points

2.3. Xử lý Itemize/Enumearate:

```

def _parse_itemize(self, content: str, parent_id: str, version: int) -> List[str]:
    """
    Itemize được xem là branching structure:

    itemize-block (higher component)
    |— item-1 (next-level element)
    |— item-2 (next-level element)
    |— item-3 (next-level element)
    """

    # Tạo block element làm parent
    block_id = self._generate_element_id('itemize')
    self.elements[block_id] = "Itemize block"
    self.hierarchy[version][block_id] = parent_id

    # Mỗi \item trở thành child
    items = re.split(r'\\item\b', content)
    item_ids = []

    for i, item_content in enumerate(items[1:], 1): # Skip empty first
        item_content = item_content.strip()
        if not item_content:
            continue

        item_id = self._generate_element_id(f'item-{i}')
        self.elements[item_id] = self._clean_text(item_content)
        self.hierarchy[version][item_id] = block_id
        item_ids.append(item_id)

    return item_ids

```

2.4. Exclusions và Inclusions:


```

# Sections bị LOẠI TRỪ khỏi hierarchy
EXCLUDE_PATTERNS = [
    r'references?',
    r'bibliography',
    r'bibliographie',
    r'\\begin\\{thebibliography\\}'
]

# Sections được BẮT BUỘC BAO GỒM
INCLUDE_PATTERNS = [
    r'acknowledgements?',
    r'appendix',
    r'appendices',
    r'\\section\\*' # Unnumbered sections
]

def _should_include_section(self, title: str) -> bool:
    """Kiểm tra section có nên được include không."""
    title_lower = title.lower()

    # Exclude references
    for pattern in EXCLUDE_PATTERNS:
        if re.search(pattern, title_lower):
            return False

    # Always include acknowledgements and appendices
    for pattern in INCLUDE_PATTERNS:
        if re.search(pattern, title_lower):
            return True

    return True # Default: include

```

3. Standardization & Cleanup:

3.1. LaTeX Cleanup:

```
class LatexCleaner:
    """Remove formatting commands không mang semantic meaning."""

    FORMATTING_COMMANDS = [
        r'\centering',
        r'\raggedright',
        r'\raggedleft',
        r'\noindent',
        r'\vspace\*?\{[^}]*\}',
        r'\hspace\*?\{[^}]*\}',
        r'\smallskip',
        r'\medskip',
        r'\bigskip',
        r'\newpage',
        r'\clearpage',
        r'\pagebreak',
    ]

    TABLE_FORMATTING = [
        r'\[htpb!?\]',
        r'\[H\]',
        r'\midrule',
        r'\toprule',
        r'\bottomrule',
        r'\hline',
        r'\cline\{[^}]*\}',
        r'\multicolumn\{[^}]*\}\{[^}]*\}',
    ]

    def clean(self, text: str) -> str:
        """Apply all cleanup rules."""
        # Remove comments
        text = re.sub(r'%.*$', '', text, flags=re.MULTILINE)

        # Remove formatting commands
        for pattern in self.FORMATTING_COMMANDS + self.TABLE_FORMATTING:
            text = re.sub(pattern, '', text)

        # Normalize whitespace
        text = re.sub(r'\s+', ' ', text)
        text = text.strip()

        return text
```

3.2. Math Normalization:

```
def normalize_inline_math(self, text: str) -> str:
    """Convert all inline math to $...$ format."""
    # \(...\) → $...$
    text = re.sub(r'\\((.+?)\\)', r'$\1$', text, flags=re.DOTALL)

    # \begin{math}...\end{math} → $...$
    text = re.sub(
        r'\\begin\\{math\\}(.+?)\\end\\{math\\}',
        r'$\1$',
        text,
        flags=re.DOTALL
    )

    return text

def normalize_block_math(self, text: str) -> str:
    """Convert all block math to \\begin{equation}...\end{equation}."""
    # $$...$$ → \begin{equation}...\end{equation}
    text = re.sub(
        r'\\$\\$(.+?)\\$',
        r'\\begin{equation}\1\\end{equation}',
        text,
        flags=re.DOTALL
    )

    # \[...\] → \begin{equation}...\end{equation}
    text = re.sub(
        r'\\\[\\((.+?)\\)\]',
        r'\\begin{equation}\1\\end{equation}',
        text,
        flags=re.DOTALL
    )

    # \begin{displaymath}...\end{displaymath} → \begin{equation}...\end{equation}
    text = re.sub(
        r'\\begin\\{displaymath\\}(.+?)\\end\\{displaymath\\}',
        r'\\begin{equation}\1\\end{equation}',
        text,
        flags=re.DOTALL
    )

    return text
```

Bảng chuyển đổi Math:

Định dạng gốc	Định dạng chuẩn hóa
$\backslash(...\backslash)$	$\$...\$$
$\backslash\begin\{math\}...\end\{math\}$	$\$...\$$
$\\\$\\\$$	$\backslash\begin\{equation\}...\end\{equation\}$
$\[...\]$	$\backslash\begin\{equation\}...\end\{equation\}$

<code>\begin{displaymath}</code>	<code>\begin{equation}...\end{equation}</code>
----------------------------------	--

4. Reference Extraction:

4.1. Nguồn References:

Hệ thống trích xuất references từ **3 nguồn**:

- **File .bib**: Sử dụng bibtexparser để parse
- **File .bbl**: Scan thư mục tex cho compiled bibliography
- **\bibitem trong .tex**: Regex extraction và conversion

4.2. Chuyển đổi \bibitem → BibTeX

```

def bibitem_to_bibtex(self, key: str, bibitem_text: str) -> Dict:
    """
    Chuyển đổi \bibitem thành BibTeX entry.

    Input:
        \bibitem{Bender2003}
        C.~M.~Bender, Must a Hamiltonian be Hermitian?,
        \href{https://doi.org/10.1119/1.1574043}{Am. J. Phys. 71, 1095 (2003)}.

    Output:
        @article{Bender2003,
          author = {C. M. Bender},
          title = {Must a Hamiltonian be Hermitian?},
          journal = {Am. J. Phys.},
          volume = {71},
          pages = {1095},
          year = {2003},
          doi = {10.1119/1.1574043}

        }
    """
    entry = {'ID': key, 'ENTRYTYPE': 'misc'}

    # Pre-processing
    text = bibitem_text.replace('~', ' ') # Non-breaking spaces
    text = re.sub(r'%.*$', '', text, flags=re.MULTILINE) # Comments

    # Extract DOI from \href{doi.org/...}{...}
    doi_match = re.search(r'\\href\[{}]*doi\.org/([{}]+)\[{}]*\]', text)
    if doi_match:
        entry['doi'] = doi_match.group(1)
        journal_info = doi_match.group(2)
        entry.update(self._parse_journal_info(journal_info))
        # Remove \href from text for further parsing
        text = re.sub(r'\\href\[{}]*\[{}]*', '', text)

    # Extract authors (at the beginning)
    entry['author'] = self._extract_authors(text)

    # Extract title (text in quotes or italics)
    entry['title'] = self._extract_title(text)

    # Extract year
    entry['year'] = self._extract_year(text)

    # Determine entry type
    entry['ENTRYTYPE'] = self._determine_entry_type(text, entry)

    return entry

```

4.3. Xử lý nhiều định dạng \bibitem

```

# Format 1: Với \href (phổ biến trong physics papers)
\bibitem{Bender2003}
C.~M.~Bender, Must a Hamiltonian be Hermitian?,
\href{https://doi.org/10.1119/1.1574043}{Am. J. Phys. 71, 1095 (2003)}.

# Format 2: Standard format
\bibitem{Smith2020}
Smith, J. and Doe, A. Title of Paper. Journal Name, 2020.

# Format 3: Với newblock
\bibitem{Jones2019}
\newblock Jones, R.
\newblock {\em Neural Networks}.
\newblock ICML, 2019.

# Format 4: Chỉ có authors và title
\bibitem{Brown2021}
Brown et al., Transformer Models, arXiv:2103.00001.

```

Chiến lược Parsing:

- **Pre-processing:** Clean ~, remove % comments
- **DOI Extraction:** Parse `\href{doi.org/...}` {journal info}
- **Author/Title Separation:** Intelligent splitting dựa trên punctuation patterns
- **Fallback:** Prefer empty fields over garbage data (conservative approach)

5. Deduplication:

5.1. Reference Deduplication:

```

def deduplicate_references(self, entries: Dict) -> Dict:
    """
    Loại trùng references với strategy:

    1. Identify duplicates: title similarity > 95% AND author overlap > 80%
    2. Select canonical key: entry có nhiều thông tin nhất
    3. Unionize fields: merge thông tin từ tất cả duplicates
    4. Rename citations: \cite{old_key} -> \cite{canonical_key}
    """
    duplicate_groups = self._find_duplicate_groups(entries)

    deduplicated = {}
    rename_map = {}

    for group in duplicate_groups:
        # Select canonical key (most complete entry)
        canonical_key = self._choose_canonical_key(group, entries)

        # Create rename mapping
        for key in group:
            if key != canonical_key:
                rename_map[key] = canonical_key

        # Unionize fields from all entries in group
        merged_entry = self._unionize_fields(group, entries)
        deduplicated[canonical_key] = merged_entry

    return deduplicated, rename_map

```


5.2. Field Unionization:

```
def _unionize_fields(self, group: List[str], entries: Dict) -> Dict:
    """
    Merge fields từ tất cả entries trong group.

    Entry 1: {title: "Deep Learning", year: 2018, venue: ""}
    Entry 2: {title: "Deep Learning", year: "", venue: "ICML"}
    Entry 3: {title: "Deep Learning", year: 2018, venue: "ICML", doi: "10.xxx"}

    Result: {title: "Deep Learning", year: 2018, venue: "ICML", doi: "10.xxx"}
    """
    merged = {}

    for key in group:
        entry = entries[key]
        for field, value in entry.items():
            if field not in merged or not merged[field]:
                merged[field] = value
            elif value and len(str(value)) > len(str(merged[field])):
                # Keep longer/more complete value
                merged[field] = value

    return merged
```

5.3. Content Deduplication:

```
def deduplicate_content(self, elements: Dict) -> Tuple[Dict, Dict]:
    """
    Full-text content deduplication across versions.

    Version 1: Element A = "Machine learning is a subset of AI."
    Version 2: Element B = "Machine learning is a subset of AI."

    → Both point to same element ID in final output
    """
    content_to_id = {}
    id_mapping = {}
    deduplicated = {}

    for elem_id, content in elements.items():
        # Normalize for comparison (after cleanup)
        normalized = self._normalize_for_comparison(content)

        if normalized in content_to_id:
            # Duplicate found - map to existing ID
            id_mapping[elem_id] = content_to_id[normalized]
        else:
            # New unique content
            content_to_id[normalized] = elem_id
            id_mapping[elem_id] = elem_id
            deduplicated[elem_id] = content

    return deduplicated, id_mapping
```

VI. Machine Learning Pipeline (Yêu cầu 2.2)

1. Data Cleaning:

1.1. Text Preprocessing:

```
class TextPreprocessor:
    """Preprocessing cho matching task."""

    def __init__(self):
        self.stopwords = set(stopwords.words('english'))

    def preprocess(self, text: str) -> str:
        """Standard preprocessing pipeline."""
        if not text:
            return ""

        # Lowercase
        text = text.lower()

        # Remove LaTeX commands
        text = re.sub(r'\\\[a-zA-Z]+\\[^\]*\\', '', text)
        text = re.sub(r'\\\[a-zA-Z]+', '', text)

        # Remove special characters (keep alphanumeric and spaces)
        text = re.sub(r'^a-z0-9\s|', ' ', text)

        # Normalize whitespace
        text = re.sub(r'\s+', ' ', text).strip()

        return text

    def tokenize(self, text: str) -> List[str]:
        """Tokenize and remove stopwords."""
        tokens = word_tokenize(self.preprocess(text))
        return [t for t in tokens if t not in self.stopwords and len(t) > 1]
```

1.2. Author Name Normalization:

```
def normalize_author_name(self, name: str) -> str:
    """
    Chuẩn hóa tên tác giả cho comparison.

    "John Smith" -> "j smith"
    "J. Smith" -> "j smith"
    "Smith, John" -> "j smith"
    "Smith, J." -> "j smith"
    """
    name = name.lower().strip()

    # Handle "Last, First" format
    if ',' in name:
        parts = name.split(',')
        if len(parts) == 2:
            last, first = parts[0].strip(), parts[1].strip()
            name = f"{first} {last}"

    # Extract initials and last name
    parts = name.split()
    if len(parts) >= 2:
        first_initial = parts[0][0] if parts[0] else ''
        last_name = parts[-1]
        return f"{first_initial} {last_name}"

    return name
```


2. Data Labeling:

2.1. Manual Labeling:

Yêu cầu: ≥ 5 publications, ≥ 20 pairs total

Thực hiện: 5 publications, **103 pairs**

```
# Tool tạo manual labels
def create_manual_labels(output_dir: str, num_pubs: int = 5):
    """
    Interactive tool để manual labeling.

    1. Hiển thị BibTeX entry
    2. Hiển thị top candidates (similarity scores)
    3. User chọn match đúng hoặc skip
    """
    labels = {}

    for pub_dir in get_publications(output_dir)[:num_pubs]:
        pub_id = pub_dir.name
        labels[pub_id] = {}

        bibtex_entries = load_bibtex(pub_dir / 'refs.bib')
        candidates = load_json(pub_dir / 'references.json')

        for bib_key, bib_data in bibtex_entries.items():
            print(f"\n{' '*60}")
            print(f"BibTeX Entry: {bib_key}")
            print(f>Title: {bib_data.get('title', 'N/A')}")
            print(f>Authors: {bib_data.get('author', 'N/A')}")
            print(f"Year: {bib_data.get('year', 'N/A')}")

            # Show top candidates
            ranked = rank_candidates(bib_data, candidates)
            for i, (cand_id, score) in enumerate(ranked[:5], 1):
                cand = candidates[cand_id]
                print(f"\n{i}. [{score:.2f}] {cand_id}")
                print(f"    Title: {cand.get('paper_title', 'N/A')}")

            # User input
            choice = input("\nEnter number (1-5), 'n' to skip, 'q' to quit: ")
            if choice == 'q':
                break
            elif choice.isdigit() and 1 <= int(choice) <= 5:
                labels[pub_id][bib_key] = ranked[int(choice)-1][0]

    save_json(labels, 'manual_labels.json')
```

Format manual_labels.json

```
{
  "2504-13946": {
    "smith2020deep": "2001-12345",
    "jones2019neural": "1911-54321",
    "brown2021transformer": "2103-98765"
  },
}
```

```
"2504-13947": {  
    ...  
}  
}
```

2.2. Automatic Labeling:

Yêu cầu: $\geq 10\%$ of remaining data

Thực hiện: 79.8% auto-labeled (vượt yêu cầu)

```
def automatic_label(self, pair: Dict, use_arxiv_matching: bool = False) -> Optional[int]:  
    """  
    Auto-labeling với conservative rules.  
  
    ⚠ QUAN TRỌNG: arXiv ID matching bị DISABLE để tránh data leakage!  
  
    Rules:  
    1. Title similarity  $\geq 99\%$   $\rightarrow$  Match (1)  
    2. Title  $\geq 90\%$  AND author overlap  $\geq 60\%$  AND year match  $\rightarrow$  Match (1)  
    3. Title  $\geq 80\%$  AND author overlap  $\geq 70\%$   $\rightarrow$  Match (1)  
    4. Title  $< 30\%$   $\rightarrow$  No match (0)  
    5. Title  $< 50\%$  AND author overlap  $< 30\%$   $\rightarrow$  No match (0)  
    6. Otherwise  $\rightarrow$  Uncertain (None) - skip  
    """  
    title_sim = self._title_similarity(  
        pair['bibtex_data'].get('title', ''),  
        pair['candidate_data'].get('paper_title', '')  
    )  
  
    author_overlap = self._author_overlap(  
        pair['bibtex_data'].get('author', ''),  
        pair['candidate_data'].get('authors', [])  
    )  
  
    year_match = self._year_match(  
        pair['bibtex_data'].get('year', ''),  
        pair['candidate_data'].get('submission_date', '')  
    )  
  
    # Rule 1: Near-exact title match  
    if title_sim  $\geq 99$ :  
        return 1  
  
    # Rule 2: High title + year + authors  
    if title_sim  $\geq 90$  and year_match and author_overlap  $\geq 60$ :  
        return 1  
  
    # Rule 3: Good title + good authors  
    if title_sim  $\geq 80$  and author_overlap  $\geq 70$ :  
        return 1  
  
    # Rule 4: Very low title similarity  
    if title_sim  $< 30$ :  
        return 0  
  
    # Rule 5: Low title + no author overlap  
    if title_sim  $< 50$  and author_overlap  $< 30$ :  
        return 0  
  
    # Uncertain - skip  
    return None
```

Tại sao disable arXiv ID matching?

Trong ban đầu, hệ thống có rule: "Nếu extract được arXiv ID từ BibTeX và match với candidate ID → positive label".

Vấn đề: Đây là **data leakage** vì:

- Model học được rằng "ID match = positive" thay vì học từ content features
- Khi test, nếu có ID match thì model đúng 100% nhưng không generalize
- Score 1.0 ngay iteration 0 là dấu hiệu của data leakage

Giải pháp: Disable arXiv ID matching, chỉ dùng content-based features.

3. Feature Engineering:

3.1. Tổng quan:

Tổng cộng: 19 features thuộc 5 groups

Group	Số features	Mô tả
Title Features	5	Similarity metrics cho titles
Author Features	5	Author name matching
Year Features	4	Temporal alignment
Text Features	5	Abstract và venue matching
Hierarchy Features	5+	Citation context từ hierarchy

3.2. Title Features (5 features)

```
def _title_features(self, pair: Dict) -> Dict:
    """
    Extract title similarity features.

    JUSTIFICATION:
    Title là tín hiệu mạnh nhất để identify papers.
    Papers có similar titles very likely là cùng một paper.

    Multiple metrics capture different aspects:
    - Jaccard: word overlap bất kể thứ tự
    - Levenshtein: character-level catches typos
    - Token sort: handles word reordering
    - Token set: handles partial matches
    """
    t1 = self.preprocess(pair['bibtex_data'].get('title', ''))
    t2 = self.preprocess(pair['candidate_data'].get('paper_title', ''))

    return {
        'title_jaccard': self._jaccard_similarity(t1, t2),
        'title_levenshtein': fuzz.ratio(t1, t2) / 100.0,
        'title_token_sort': fuzz.token_sort_ratio(t1, t2) / 100.0,
        'title_token_set': fuzz.token_set_ratio(t1, t2) / 100.0,
        'title_exact_match': int(t1.strip() == t2.strip())
    }
```

Feature	Formula/Method	Justification
title_jaccard	$\frac{ A \cap B }{ A \cup B }$	Word-level overlap, robust to reordering
title_levenshtein	Levenshtein ratio	Character-level similarity, catches typos
title_token_sort	FuzzyWuzzy token_sort_ratio	Handles word reordering
title_token_set	FuzzyWuzzy token_set_ratio	Handles partial matches
title_exact_match	Boolean	Strong positive signal

Data Analysis:

- Title Jaccard > 0.8: **95%** match rate
- Title Jaccard < 0.5: **2%** match rate
- Correlation với label: **r = 0.87**

3.3. Author Features (5 features)

```
def _author_features(self, pair: Dict) -> Dict:
    """
    Extract author matching features.

    JUSTIFICATION:
    Authors verify paper identity.
    First author is most distinctive (usually corresponding author).
    """
    authors1 = self._parse_authors(pair['bibtex_data'].get('author', ''))
    authors2 = pair['candidate_data'].get('authors', [])

    if not authors1 or not authors2:
        return self._empty_author_features()

    # Normalize all names
    norm1 = [self.normalize_author_name(a) for a in authors1]
    norm2 = [self.normalize_author_name(a) for a in authors2]

    # Common authors
    common = set(norm1) & set(norm2)

    return {
        'author_overlap_ratio': len(common) / min(len(norm1), len(norm2)),
        'first_author_match': int(norm1[0] == norm2[0]) if norm1 and norm2 else 0,
        'last_author_match': int(norm1[-1] == norm2[-1]) if norm1 and norm2 else 0,
        'num_common_authors': len(common),
        'author_initials_match': self._initials_match(authors1, authors2)
    }
```

Feature	Method	Justification
author_overlap_ratio	Common / min(count)	Identity verification
first_author_match	Last name comparison	Most distinctive author
last_author_match	Last name comparison	Senior author signal
num_common_authors	Absolute count	Multi-author papers
author_initials_match	Initial pattern	Handles "J. Smith" vs "John Smith"

3.4. Year Features (4 features)

```
def _year_features(self, pair: Dict) -> Dict:
    """
    Extract year matching features.

    JUSTIFICATION:
    Year là strong filter - papers từ different years unlikely match.
    ±1 year tolerance cho arXiv preprint vs journal publication timing.
    """
    year1 = self._extract_year(pair['bibtex_data'].get('year', ''))
    year2 = self._extract_year_from_date(
        pair['candidate_data'].get('submission_date', '')
    )

    if not year1 or not year2:
        return {
            'year_diff': 10, # Max penalty
            'year_exact_match': 0,
            'year_within_1': 0,
            'both_years_present': 0
        }

    diff = abs(int(year1) - int(year2))

    return {
        'year_diff': min(diff, 10), # Cap at 10
        'year_exact_match': int(diff == 0),
        'year_within_1': int(diff <= 1),
        'both_years_present': 1
    }
```

Feature	Method	Justification
year_diff	abs(year1 - year2), capped at 10	Temporal filter
year_exact_match	Boolean	Strong signal
year_within_1	$\text{diff} \leq 1$	arXiv vs journal timing
both_years_present	Boolean	Data quality indicator

3.5. Text Features (5 features)

```
def _text_features(self, pair: Dict) -> Dict:
    """
    Extract additional text features.

    JUSTIFICATION:
    Abstract provides deep content similarity.
    Venue matching indicates same publication target.
    """
    # Abstract similarity (if available)
    abstract1 = pair['bibtex_data'].get('abstract', '')
    abstract2 = pair['candidate_data'].get('abstract', '')

    if abstract1 and abstract2:
        abstract_jaccard = self._jaccard_similarity(
            self.preprocess(abstract1),
            self.preprocess(abstract2)
        )
    else:
        abstract_jaccard = 0.0

    # Venue similarity
    venue1 = pair['bibtex_data'].get('journal', '') or \
        pair['bibtex_data'].get('booktitle', '')
    venue2 = pair['candidate_data'].get('venue', '')

    venue_sim = fuzz.token_set_ratio(
        self.preprocess(venue1),
        self.preprocess(venue2)
    ) / 100.0 if venue1 and venue2 else 0.0

    return {
        'abstract_jaccard': abstract_jaccard,
        'venue_similarity': venue_sim,
        'text_word_overlap': self._word_overlap(pair),
        'bigram_similarity': self._bigram_similarity(pair),
        'title_length_ratio': self._title_length_ratio(pair)
    }
```


3.6. Hierarchy Features (5+ features)

```
class HierarchyFeatureExtractor:
    """
    Extract features từ citation context trong hierarchy.

    JUSTIFICATION:
    Citation context provides signals about paper importance.

    CITED IN INTRODUCTION:
    - Often seminal/background works
    - Well-known papers

    CITED IN METHODS:
    - Technical algorithm papers
    - Recent papers

    CITED IN RESULTS:
    - Comparison baselines
    - Competitors
    """

    def extract(self, bibtex_key: str, hierarchy: Dict) -> Dict:
        # Count citations of this key
        citation_count = self._count_citations(bibtex_key, hierarchy)

        # Detect which sections cite this key
        sections = self._get_citing_sections(bibtex_key, hierarchy)

        return {
            'citation_count': citation_count,
            'cited_in_intro': int('introduction' in sections),
            'cited_in_methods': int('method' in sections or 'approach' in sections),
            'cited_in_results': int('result' in sections or 'experiment' in sections),
            'near_figure': int(self._near_figure(bibtex_key, hierarchy)),
            'co_citation_count': self._co_citation_count(bibtex_key, hierarchy)
        }
```

4. Data Modeling:

4.1. Problem Formulation:

Bài toán: Learning-to-Rank

Input: (BibTeX entry, Candidate reference) pairs

Output: Ranked list of top-5 candidates cho mỗi BibTeX entry

Data Structure:

- m BibTeX entries \times n candidates = $m \times n$ pairs per publication
- Class imbalance: typically 1:50 to 1:200 (1 positive per m pairs)

4.2. Model Selection:

Model: CatBoost Ranker với YetiRank loss

```
from catboost import CatBoostRanker, Pool

class ModelTrainer:
    def __init__(self, model_type: str = 'ranker'):
        self.model = CatBoostRanker(
            iterations=800,
            learning_rate=0.05,
            depth=8,
            loss_function='YetiRank',
            eval_metric='NDCG',
            l2_leaf_reg=5.0,
            min_data_in_leaf=5,
            early_stopping_rounds=100,
            random_seed=42,
            verbose=100
        )

    def train(self, train_data: pd.DataFrame, val_data: pd.DataFrame):
        """Train model với group_id cho ranking."""
        # Create pools với group_id (publication + bibtex_key)
        train_pool = Pool(
            data=train_data[self.feature_columns],
            label=train_data['label'],
            group_id=train_data['group_id']
        )

        val_pool = Pool(
            data=val_data[self.feature_columns],
            label=val_data['label'],
            group_id=val_data['group_id']
        )

        self.model.fit(train_pool, eval_set=val_pool)

        return self.model
```

Tại sao chọn CatBoost Ranker?

Lý do	Giải thích
Native ranking support	YetiRank optimizes directly for NDCG
Class imbalance handling	Built-in balanced weighting
No feature scaling required	Tree-based method
Efficient training	GPU acceleration support
Robust to overfitting	Built-in regularization

4.3. Data Split Strategy:

```
def split_data(self, publications: List[str], manual_labeled: List[str]):  
    """  
    Split theo yêu cầu:  
    - Test: 1 manual + 1 auto publication  
    - Validation: 1 manual + 1 auto publication  
    - Training: All remaining  
    """  
    # Separate manual and auto-labeled publications  
    auto_labeled = [p for p in publications if p not in manual_labeled]  
  
    # Test set  
    test_manual = manual_labeled[0]  
    test_auto = auto_labeled[0]  
    test_pubs = [test_manual, test_auto]  
  
    # Validation set  
    val_manual = manual_labeled[1]  
    val_auto = auto_labeled[1]  
    val_pubs = [val_manual, val_auto]  
  
    # Training set  
    train_pubs = [p for p in publications  
                  if p not in test_pubs and p not in val_pubs]  
  
    return train_pubs, val_pubs, test_pubs
```

5. Model Evaluation:

5.1. Mean Reciprocal Rank (MRR)

Công thức:

$$MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{rank_i}$$

Trong đó:

- $|Q|$ = số lượng queries (BibTeX entries)
- $rank_i$ = vị trí của correct match trong top-5 predictions
- Nếu correct match không có trong top-5: score = 0

```

def calculate_mrr(self, predictions: Dict, ground_truth: Dict) -> float:
    """
    Calculate Mean Reciprocal Rank.

    Example:
        Query 1: correct at rank 1 → 1/1 = 1.0
        Query 2: correct at rank 3 → 1/3 = 0.33
        Query 3: not in top 5 → 0

        MRR = (1.0 + 0.33 + 0) / 3 = 0.44
    """
    reciprocal_ranks = []

    for bibtex_key, true_id in ground_truth.items():
        if bibtex_key not in predictions:
            reciprocal_ranks.append(0.0)
            continue

        pred_list = predictions[bibtex_key]

        if true_id in pred_list:
            rank = pred_list.index(true_id) + 1
            reciprocal_ranks.append(1.0 / rank)
        else:
            reciprocal_ranks.append(0.0)

    return sum(reciprocal_ranks) / len(reciprocal_ranks)

```

5.2. Additional Metrics:

```

def evaluate(self, predictions: Dict, ground_truth: Dict) -> Dict:
    """Comprehensive evaluation."""
    results = {
        'mrr': self.calculate_mrr(predictions, ground_truth),
        'hit_at_1': self._hit_at_k(predictions, ground_truth, k=1),
        'hit_at_3': self._hit_at_k(predictions, ground_truth, k=3),
        'hit_at_5': self._hit_at_k(predictions, ground_truth, k=5),
        'avg_rank': self._average_rank(predictions, ground_truth)
    }

    return results

def _hit_at_k(self, predictions, ground_truth, k: int) -> float:
    """Percentage of queries with correct match in top-k."""
    hits = 0
    total = len(ground_truth)

    for bibtex_key, true_id in ground_truth.items():
        pred_list = predictions.get(bibtex_key, [])[:k]
        if true_id in pred_list:
            hits += 1

    return hits / total if total > 0 else 0.0

```

VII. Giải thích mã nguồn:

1. Entry point (main_parser.py)

```
def main():
    """Main entry point for parser."""
    parser = argparse.ArgumentParser(description='Parse LaTeX files')
    parser.add_argument('--input-dir', required=True, help='Input directory')
    parser.add_argument('--output-dir', required=True, help='Output directory')
    parser.add_argument('--batch', action='store_true', help='Batch mode')
    args = parser.parse_args()

    if args.batch:
        # Process all publications in directory
        for pub_dir in Path(args.input_dir).iterdir():
            if pub_dir.is_dir():
                process_publication(pub_dir, args.output_dir)
    else:
        # Process single publication
        process_publication(Path(args.input_dir), args.output_dir)

def process_publication(pub_dir: Path, output_dir: str):
    """Process a single publication."""
    # Step 1: Gather files
    gatherer = FileGatherer()
    main_file = gatherer.find_main_file(pub_dir / 'tex')
    tex_files = gatherer.gather_files(main_file)

    # Step 2: Clean and build hierarchy
    cleaner = LatexCleaner()
    builder = HierarchyBuilder()

    for tex_file in tex_files:
        content = cleaner.clean(tex_file.read_text())
        builder.process(content)

    # Step 3: Extract references
    extractor = ReferenceExtractor()
    references = extractor.extract(tex_files)

    # Step 4: Deduplicate
    deduplicator = Deduplicator()
    references, rename_map = deduplicator.deduplicate_references(references)
    elements, id_map = deduplicator.deduplicate_content(builder.elements)

    # Step 5: Save outputs
    save_hierarchy(output_dir, builder.hierarchy, elements)
    save_bibtex(output_dir, references)
```

2. Entry point (main_matcher.py)

```

def main():
    """Main entry point for matching pipeline."""
    parser = argparse.ArgumentParser(description='Reference matching')
    parser.add_argument('--data-dir', required=True)
    parser.add_argument('--manual-labels', default='manual_labels.json')
    parser.add_argument('--model-type', default='ranker')
    parser.add_argument('--output-dir', default='ml_output')
    args = parser.parse_args()

    # Load manual labels
    manual_labels = load_json(args.manual_labels)

    # Step 1: Prepare data
    preparator = DataPreparator()
    all_pairs = preparator.create_all_pairs(args.data_dir)

    # Step 2: Label data
    labeler = Labeler()
    labeled_pairs = labeler.apply_labels(all_pairs, manual_labels)

    # Step 3: Extract features
    extractor = FeatureExtractor()
    feature_matrix = extractor.extract_all(labeled_pairs)

    # Step 4: Split data
    train_data, val_data, test_data = split_data(
        feature_matrix,
        manual_labeled=list(manual_labels.keys())
    )

    # Step 5: Train model
    trainer = ModelTrainer(model_type=args.model_type)
    model = trainer.train(train_data, val_data)

    # Step 6: Evaluate
    evaluator = Evaluator()
    predictions = trainer.predict_top_k(test_data, k=5)
    metrics = evaluator.evaluate(predictions, test_data)

    # Step 7: Save results
    save_predictions(args.output_dir, predictions, metrics)

    print(f"\nMRR: {metrics['mrr']:.4f}")
    print(f"Hit@1: {metrics['hit_at_1']:.2%}")
    print(f"Hit@5: {metrics['hit_at_5']:.2%}")

```

3. Feature Extraction Pipeline

```

class FeatureExtractor:
    """Extract all features for pairs."""

    def __init__(self):
        self.preprocessor = TextPreprocessor()
        self.hierarchy_extractor = HierarchyFeatureExtractor()

    def extract_all(self, pairs: List[Dict]) -> pd.DataFrame:
        """Extract features for all pairs."""
        features_list = []

        for pair in tqdm(pairs, desc="Extracting features"):
            features = {}

            # Title features (5)
            features.update(self._title_features(pair))

            # Author features (5)
            features.update(self._author_features(pair))

            # Year features (4)
            features.update(self._year_features(pair))

            # Text features (5)
            features.update(self._text_features(pair))

            # Hierarchy features (5+)
            if 'hierarchy' in pair:
                features.update(
                    self.hierarchy_extractor.extract(
                        pair['bibtex_key'],
                        pair['hierarchy']
                    )
                )

            # Metadata
            features['publication_id'] = pair['publication_id']
            features['bibtex_key'] = pair['bibtex_key']
            features['candidate_id'] = pair['candidate_id']
            features['label'] = pair.get('label')

            # Group ID for ranking
            features['group_id'] = f"{pair['publication_id']}_{pair['bibtex_key']}"

            features_list.append(features)

        return pd.DataFrame(features_list)

```

VIII. Kết quả thống kê:

1. Parsing Statistics:

Metric	Value
Publications processed	129
Total .tex files scanned	1,247
Multi-file publications	89 (69%)
Total hierarchy elements	~45,000
Average elements per paper	~350
Parse success rate	98.4%

2. Reference Extraction Statistics:

Metric	Value
Total BibTeX entries extracted	4,127
From .bib files	2,891
From .bbl files	652
From \bibitem conversion	584
Entries with DOI	1,847 (45%)
Entries with complete metadata	3,104 (75%)
Duplicates removed	312

3. Labeling Statistics:

Metric	Value
Manual labeled publications	5
Manual labeled pairs	103
Auto-labeled pairs	13,496
Auto-labeling rate	79.8% (vượt yêu cầu 10%)
Positive labels	892
Negative labels	12,604

4. Model Performance:

Metric	CatBoost Ranker
MRR (Test)	0.995
Hit@1	99.5%
Hit@3	100%
Hit@5	100%
Avg Rank	1.005
Training time	8.2 min

5. Feature Importance (Top 10):

Rank	Feature	Importance
1	title_token_set	28.4
2	title_jaccard	22.1
3	author_overlap_ratio	15.7
4	first_author_match	9.3

5	year_exact_match	7.2
6	title_levenshtein	5.8
7	venue_similarity	4.1
8	num_common_authors	3.2
9	citation_count	2.4
10	year_within_1	1.8

Insight: Title-based features chiếm 56.3% tổng importance, validate quyết định focus vào title similarity metrics trong feature engineering.

IX. Những hạn chế, thách thức và giải pháp:

1. RecursionError với isinstance():

Vấn đề: RecursionError: maximum recursion depth exceeded khi sử dụng isinstance() với các cấu trúc dữ liệu lồng nhau.

Nguyên nhân: Python's isinstance() có thể trigger recursive checks trong một số trường hợp.

Giải pháp:

```
# Thay vì:
if isinstance(obj, dict):
    ...

# Sử dụng:
if type(obj).__name__ == 'dict':
    ...
```

2. Empty refs.bib:

Vấn đề: Nhiều publications có refs.bib rỗng dù có references trong paper.

Nguyên nhân: Parser chỉ scan .bib files, bỏ qua .bbl và \bibitem.

Giải pháp:


```
def extract_references(self, tex_dir: Path) -> Dict:
    """Extract from multiple sources."""
    references = {}

    # Source 1: .bib files
    for bib_file in tex_dir.rglob('*.bib'):
        references.update(self._parse_bibtex(bib_file))

    # Source 2: .bbl files (compiled bibliography)
    for bbl_file in tex_dir.rglob('*.bbl'):
        references.update(self._parse_bbl(bbl_file))

    # Source 3: \bibitem in .tex files
    for tex_file in tex_dir.rglob('*.tex'):
        references.update(self._extract_bibitems(tex_file))

    return references
```

3. Data Leakage trong Auto Labeling:

Vấn đề: Model đạt score 1.0 ngay iteration 0 → overfitting/data leakage.

Nguyên nhân: Auto-labeling dùng arXiv ID matching: nếu extract được ID từ BibTeX và match với candidate → positive. Model học được "ID match = positive" thay vì học từ content.

Giải pháp: Disable arXiv ID matching:

```
def automatic_label(self, pair: Dict, use_arxiv_matching: bool = False):
    # arXiv ID matching DISABLED by default
    if use_arxiv_matching:
        arxiv_id = self._extract_arxiv_id(pair['bibtex_data'])
        if arxiv_id and arxiv_id == pair['candidate_id']:
            return 1

    # Only use content-based features
    ...
```

4. CatBoost Feature Importance Error

Vấn đề: CatBoostError: Feature importance type LossFunctionChange requires training dataset

Nguyên nhân: Default feature importance method cần training data.

Giải pháp: Sử dụng PredictionValuesChange:


```

def get_feature_importance(self, train_pool=None):
    """Get feature importance without requiring training data."""
    try:
        # PredictionValuesChange doesn't require training data
        return self.model.get_feature_importance(
            type='PredictionValuesChange'
        )
    except Exception:
        # Fallback
        if train_pool:
            return self.model.get_feature_importance(
                data=train_pool,
                type='LossFunctionChange'
            )
        return None

```

5. Parsing Quality Issues:

Vấn đề: Title chứa DOI, author names bị cắt, comment markers trong output.

Giải pháp: Improved parsing với pre-processing và validation:

```

def bibitem_to_bibtex(self, key: str, text: str) -> Dict:
    # Pre-processing
    text = text.replace('~', ' ') # Non-breaking spaces
    text = re.sub(r'%.*$', '', text, flags=re.MULTILINE) # Comments

    # Remove DOI links before parsing
    text = re.sub(r'\\href\\{[^}]*doi[^}]*\\{[^}]*\\}', '', text)

    # Extract with validation
    title = self._extract_title(text)
    if title and self._validate_title(title):
        entry['title'] = title
    else:
        entry['title'] = '' # Prefer empty over garbage

    return entry

def _validate_title(self, title: str) -> bool:
    """Validate title quality."""
    # Must have reasonable alpha ratio
    alpha_count = sum(c.isalpha() for c in title)
    if alpha_count / len(title) < 0.5:
        return False

    # Must not contain DOI fragments
    if 'doi.org' in title.lower() or 'http' in title.lower():
        return False

    return True

```

X. Kết luận:

1. Thành tựu chính:

Hierarchical Parsing (Yêu cầu 2.1):

- Multi-file gathering với `\input\include` resolution
- Hierarchy construction với proper leaf nodes (Sentences, Formulas, Figures)
- Math normalization: inline \rightarrow `$...$`, block \rightarrow `\begin{equation}`
- Reference extraction từ 3 nguồn: `.bib`, `.bbl`, `\bibitem`
- Deduplication với field unionization và citation renaming
- **Parse success rate: 98.4%**

Reference Matching (Yêu cầu 2.2):

- Manual labeling: 5 publications, 103 pairs (vượt yêu cầu ≥ 20)
- Auto-labeling: 79.8% (vượt yêu cầu $\geq 10\%$)
- 19 features với justifications đầy đủ
- CatBoost Ranker với YetiRank loss
- **MRR: 0.995** trên test set
- **Hit@1: 99.5%**, Hit@5: 100%

2. Điểm mạnh kỹ thuật:

- **Modular Architecture:** Code được tổ chức thành modules độc lập, dễ maintain và extend.
- **Robust Parsing:** Xử lý nhiều định dạng `\bibitem` phức tạp, bao gồm `\href{DOI} {}`.
- **Conservative Fallback:** Prefer empty fields over garbage data \rightarrow better downstream matching.
- **Data Leakage Prevention:** Disable arXiv ID matching để model học từ content features.
- **Comprehensive Features:** 19 features thuộc 5 groups, cover title, author, year, text, và hierarchy.

3. Hạn chế và hướng phát triển:

Hạn chế	Hướng cải thiện
Title-only matching cho papers without abstracts	Add embedding-based features (SciBERT, Sentence-BERT)
Limited to arXiv candidates	Extend to Semantic Scholar, Google Scholar
Single-language support	Add multilingual paper handling
Manual labeling time-consuming	Develop active learning approach

Xin cảm ơn sự hỗ trợ từ các giảng viên trong bộ môn Nhập môn Khoa học dữ liệu để có thể hoàn thành Project lần này.

XI. Tài liệu tham khảo:

1. **CatBoost Documentation:** <https://catboost.ai/docs/>
2. **BibTeX Format Specification:** <https://www.bibtex.org/Format/>
3. **FuzzyWuzzy String Matching:** <https://github.com/seatgeek/fuzzywuzzy>
4. **arXiv API Documentation:** <https://arxiv.org/help/api/>
5. **NLTK Natural Language Toolkit:** <https://www.nltk.org/>
6. **Semantic Scholar API:** <https://api.semanticscholar.org/>

XII. Đường link Youtube:

Link video demo: [Chèn link YouTube tại đây]

Nội dung video:

- Demo chạy parser pipeline
- Demo chạy matching pipeline
- Giải thích output files
- Phân tích kết quả MRR

XIII. Lời cảm ơn:

Để hoàn thành báo cáo Lab 02 cho học phần Nhập môn Khoa học Dữ liệu, em xin gửi lời cảm ơn chân thành đến:

- Các thầy cô trong môn học này, đặc biệt là thầy **Huỳnh Lâm Hải Đăng** đã tận tình hướng dẫn và cung cấp những kiến thức nền tảng quý báu trong buổi Seminar cũng như trong file yêu cầu đồ án để em có thể thực hiện đồ án này.

- **Cộng đồng open-source** đã phát triển các thư viện CatBoost, FuzzyWuzzy, NLTK giúp việc triển khai trở nên thuận tiện hơn.