



UNDERGRADUATE THESIS

**DEVELOPMENT OF SEMANTIC SEARCH SYSTEM
WITH QUERY EXPANSION FOR THESIS RETRIEVAL
IN THE UPN ACADEMIC REPOSITORY**

**DELA PUSPITA LASMININGRUM
NPM 22081010209**

THESIS ADVISORS

**Eva Yulia Puspaningrum, S.Kom., M.Kom.
Budi Mukhamad Mulyo, S.Kom., M.T.**

**MINISTRY OF HIGHER EDUCATION, SCIENCE, AND TECHNOLOGY
UNIVERSITAS PEMBANGUNAN NASIONAL VETERAN JAWA TIMUR
FACULTY OF COMPUTER SCIENCE
INFORMATICS STUDY PROGRAM
SURABAYA
2026**

APPROVAL SHEET

DEVELOPMENT OF SEMANTIC SEARCH SYSTEM WITH QUERY
EXPANSION FOR THESIS RETRIEVAL IN THE UPN ACADEMIC
REPOSITORY

By :
DELA PUSPITA LASMININGRUM
NPM. 22081010209


Has been defended before, and accepted by, the Board of Assessors of the Thesis Examination of the Informatics Study Program, Faculty of Computer Science, Universitas Pembangunan Nasional Veteran Jawa Timur, on April 15, 2026.

Approved,

Eva Yulia Puspaningrum, S.Kom., M.Kom.
NIP. 19890705 2021212 002


..... (Advisor I)


Budi Mukhamad Mulyo, S.Kom., M.T.
NIP. 19891118 202406 1 003


..... (Advisor II)


Dr. Rizky Parlita, S.Kom., M.Kom.
NIP. 19840518 2021211 003


..... (Head Assessor)

Achmad Junaidi, S.Kom., M.Kom.
NIP. 19781110 2025211 048


..... (Assessor I)

Acknowledge by,
Dean of the Faculty of Computer Science


Prof. Dr. H. Novirina Hendrasarie, MT
NIP. 19681126 199403 2 001

APPROVAL SHEET

**DEVELOPMENT OF SEMANTIC SEARCH SYSTEM WITH QUERY
EXPANSION FOR THESIS RETRIEVAL IN THE UPN ACADEMIC
REPOSITORY**

By:
DELA PUSPITA LASMININGRUM
NPM. 22081010209

Approved to proceed to the Thesis Examination

Approved by,

**Coordinator of Informatics Study Program
Faculty of Computer Science**



Dr. Intan Yuniar Purbasari, S.Kom. MSc.
NIP. 19800602 202521 2 029

STATEMENT OF ORIGINALITY

I am the undersigned:

Student Name : DELA PUSPITA LASMININGRUM
NPM : 22081010209
Degree Program : Bachelor (S1)
Study Program : Informatics
Faculty : Faculty of Computer Science

Hereby declares that this undergraduate thesis contains no part of any other scientific work that has been submitted to obtain an academic degree at any higher education institution. Furthermore, it does not contain any work or opinions previously written or published by others, except for those which are explicitly cited in this thesis and listed completely in references.

And I declare that this scientific document is free from elements of plagiarism. If in the future indications of plagiarism are found in this Thesis, I am willing to accept sanctions in accordance with the applicable laws and regulations.

Thus, I made this statement without any coercion from anyone and to be used as it should.



Surabaya, May 18th, 2026

Declarant,



DELA PUSPITA LASMININGRUM

NPM. 22081010209

ABSTRACT

Student Name / NPM : DELA PUSPITA LASMININGRUM / 22081010209
Thesis Title : DEVELOPMENT OF SEMANTIC SEARCH SYSTEM WITH QUERY EXPANSION FOR THESIS RETRIEVAL IN THE UPN ACADEMIC REPOSITORY
Supervisor : 1. Eva Yulia Puspaningrum, S.Kom., M.Kom.
2. Budi Mukhamad Mulyo, S.Kom., M.T.

Information retrieval systems play an important role in helping users efficiently find relevant documents. In higher education, thesis retrieval has become an essential need for students as a writing reference and to avoid similarity in research topics. The digital repository of UPN “Veteran” Jawa Timur provides access to student theses; however, the existing search system still relies on explicit keyword matching, making it less effective when users submit short, general queries or use different terms with similar meanings.

This study aims to design and implement a thesis retrieval system based on semantic search capable of understanding the contextual meaning of queries and documents. The proposed approach utilizes IndoSBERT to generate semantic representations (embeddings) of thesis titles, abstracts, and user queries. To improve sensitivity to term variations and limitations of short queries, an embedding-based query expansion technique is applied. The embedding retrieval process is conducted using FAISS with cosine similarity to maintain search efficiency on a large-scale dataset. The study includes system design, embedding generation, query expansion integration, as well as implementation and functional testing of the retrieval system.

The experimental results show that the system is capable of retrieving relevant documents based on semantic meaning despite vocabulary differences between queries and documents. The semantic fine-tuning model with query expansion achieved an nDCG@15 score of 0.8470, indicating that relevant documents tend to appear at the top of the search results. The application of fine-tuning and query expansion proved effective in improving the quality of semantic representations and enriching query context, enabling the system to capture semantic relationships better than the baseline semantic approach. In terms of performance, the system was able to process searches over 15,326 documents with an average retrieval time of approximately 0.5 seconds using FAISS. Therefore, the IndoSBERT-based semantic search approach combined with fine-tuning and query expansion can serve as an alternative approach for thesis retrieval systems, particularly for queries with vocabulary variations.

Keywords: *semantic search*, IndoSBERT, *query expansion*, FAISS, thesis retrieval.

ACKNOWLEDGEMENTS

Praise and gratitude are rendered unto Allah SWT for His abundant grace and guidance, which have enabled the author to complete this thesis entitled "Development of Semantic Search System with Query Expansion for Thesis Retrieval in the UPN Academic Repository". This thesis is prepared as one of the academic requirements for obtaining the degree of Bachelor of Computer Science in the Informatics Study Program, Faculty of Computer Science, Universitas Pembangunan Nasional "Veteran" Jawa Timur.

The process of preparing this research was not without its challenges; however, through the guidance, support, and prayers of various parties, all obstacles were overcome. Therefore, with sincere humility, the author wishes to express deep appreciation and heartfelt gratitude to:

1. Prof. Dr. Ir. Novirina Hendrasarie, MT, as Dean of the Faculty of Computer Science, Universitas Pembangunan Nasional "Veteran" Jawa Timur.
2. Dr. Intan Yuniar Purbasari, S.Kom., MSc., as Head of the Informatics Study Program, Faculty of Computer Science Universitas Pembangunan Nasional "Veteran" Jawa Timur.
3. Eva Yulia Puspaningrum, S.Kom., M.Kom. and Budi Mukhamad Mulyo, S.Kom., M.T., as academic supervisors who generously dedicated their time to providing guidance, serving as critical brainstorming partners, and assisting the author in identifying effective solutions to every technical challenge encountered throughout the thesis process.
4. Dr. Rizky Parluka, S.Kom., M.Kom. and Achmad Junaidi, S.Kom., M.Kom., as thesis examiners who provided invaluable feedback in evaluating and improving the various shortcomings of this research in pursuit of a higher quality thesis.
5. Retno Mumpuni, S.Kom., M.Sc., as academic advisor who has provided academic direction and guidance to the author from the beginning of the study period through to the final stage of completion.
6. All Lecturers of the Informatics Study Program, who have equipped the author with knowledge, insight, and inspiration throughout the course of study.

7. The Coordinator and staff of the UPT Library of UPN "Veteran" Jawa Timur, who granted research permission and facilitated the author in the data collection process essential to this research.
8. All my respondents from various majors in UPN Jatim, who voluntarily participated in the testing process. Thank you for generously dedicating your time and providing the valuable feedback necessary for the evaluation of this research.
9. My beloved Mother and Father, who have ceaselessly provided material and spiritual support, along with my Sister as my greatest inspiration. My deepest gratitude is extended for all the prayers, affection, and unwavering material and spiritual support that sustained the author through to the successful completion of this thesis.
10. Kevin Ricky Pradana, my greatest support system who has consistently stood by the author throughout the challenges of completing this thesis. Thank you for being an exceptional discussion partner, for accompanying exploring various coffeeshops to meet our daily targets, and for the endless mutual motivation as we strive to graduate and achieve our dreams together.
11. Izza, Firli, Dilla, Aileen, Rahma, and Dwi, as the author's closest friends and fellow thesis fighters. Thank you for sharing the burdens and joys of this journey together.
12. My Zfordda Mojokerto friend circle, who, despite our occasional gatherings, have always provided deep conversations and valuable insights that helped the author gain new perspectives on this thesis.
13. Members of HIMATIFA, particularly the Executive Board (BPH) and Heads of Departments for the 2023 and 2024 periods. Thank you for the extensive insights and knowledge shared across all matters, which significantly contributed to the author's personal growth and the completion of this research.
14. All other friends and parties who cannot be mentioned individually, but have provided support and assistance in any form during the author's academic journey.

The author acknowledges that this thesis contains many shortcomings. Therefore, constructive criticism and suggestions from all parties are greatly welcomed in the interest of improving this work. Finally, with all the limitations the author possesses, it is hoped that this report may be of benefit to all parties in general, and to the author in particular.

Surabaya, May, 18th 2026

Penulis

TABLE OF CONTENTS

COVER PAGE	i
APPROVAL SHEET	ii
APPROVAL SHEET	iv
STATEMENT OF ORIGINALITY	vi
ABSTRACT	i
ACKNOWLEDGEMENTS.....	iii
TABLE OF CONTENTS	vi
LIST OF FIGURE.....	xi
LIST OF TABLES	xiii
LIST OF FORMULA	xv
CHAPTER I INTRODUCTION	1
1.1 Background.....	1
1.2 Problem Formulation.....	3
1.3 Objectives	3
1.4 Benefits.....	4
1.5 Scope and Limitations	4
CHAPTER II LITERATURE REVIEW	5
2.1 Previous Research	5
2.2 Information Retrieval (IR).....	7
2.3 Search System	8
2.4 Natural Language Processing (NLP).....	9
2.5 Text Preprocessing	9
2.6 Word Embedding / Sentence Embedding.....	11
2.7 Semantic Representation	12
2.7.1 BERT (Bidirectional Encoder Representations from Transformers)	13

2.7.2	SBERT (<i>Sentence Bidirectional Encoder Representations from Transformers</i>).....	13
2.7.3	Indo-Sentence BERT	14
2.8	Semantic search	15
2.9	Fine Tuning with Multiple Negative Ranking Loss	16
2.10	Cosine Similarity	17
2.11	Query Expansion	18
2.12	FAISS (Facebook AI Similarity Search)	19
2.13	Evaluation Metrics.....	21
2.14	Hugging Face.....	24
CHAPTER III METHODOLOGY		25
3.1	Research Stages	25
3.2	Problem Analysis.....	26
3.3	Data Collection.....	27
3.3.1	Data Source.....	27
3.3.2	Data Collection Process	27
3.3.3	Data Storage.....	29
3.4	Data Preprocessing	30
3.4.1	Data Translation.....	30
3.4.2	Normalization and Case Folding.....	31
3.4.3	Tokenization	32
3.4.4	Stopword Removal	33
3.4.5	Text Representation Column Formation.....	35
3.5	Query expansion	36
3.5.1	Loading the Pre-trained FastText Model	37
3.5.2	Fine Tuning Word2Vec	38

3.5.3	Query Tokenization.....	40
3.5.4	Expansion Candidate Retrieval Using Word2Vec and FastText	40
3.5.5	Expansion Term Filtering	42
3.5.6	Query Expansion Formation	42
3.6	Embedding Model Formation (IndoSBERT)	43
3.6.1	Loading the Baseline (Pre-trained) IndoSBERT Model.....	43
3.6.2	Document Embedding Formation Using the Baseline Model	44
3.6.3	Fine-Tuning Data Preparation (<i>Title-Abstract</i>)	45
3.6.4	Fine-Tuning IndoSBERT.....	45
3.6.5	Document Embedding Formation (Fine-Tuned).....	46
3.7	Semantic Retrieval with FAISS.....	47
3.7.1	FAISS Index Construction for IndoSBERT (Baseline)	48
3.7.2	FAISS Index Construction for Fine-Tuned IndoSBERT	49
3.7.3	Semantic Retrieval Process via FAISS	50
3.8	System Testing	51
3.8.1	Test Query Compilation.....	51
3.8.2	System Testing Scenarios	53
3.8.3	Document Pooling and Relevance Assessment	54
3.9	Result Evaluation.....	56
3.10	System Visualization	57
CHAPTER IV RESULT AND DISCUSSION.....		61
4.1	Results of Data Collection and Preprocessing.....	61
4.1.1	Data Collection Results.....	61
4.1.2	Data Preprocessing Results.....	61
4.2	Query Expansion Result	66
4.2.1	Tokenization and Document Frequency Computation.....	66

4.2.2	Fine-Tuning Word2Vec with FastText Initialization.....	66
4.2.3	Expansion Term Filtering	68
4.2.4	Hybrid Query Expansion Results.....	69
4.3	IndoSBERT Embedding Model Construction Results	72
4.3.1	Loading the Baseline Model	72
4.3.2	Fine-Tuning Data Preparation.....	73
4.3.3	IndoSBERT Fine Tuning	74
4.3.4	Comparison of Baseline vs. Fine-Tuned Retrieval Results	75
4.4	System Testing and Evaluation Results.....	76
4.4.1	Document Pooling	76
4.4.2	Relevance Assessment	77
4.4.3	Retrieval Results by Query Across Various Test Scenarios	79
4.4.4	Evaluation Results Using Metrics (Precision, Recall, MAP, and nDCG)	95
4.5	Discussion.....	100
4.5.1	System Capability in Handling Vocabulary Differences	100
4.5.2	Application of IndoSBERT and Query Expansion	101
4.5.3	System Performance on Large-Scale Collections	102
4.5.4	Comparative Analysis of Lexical and Semantic Matching Characteristics	102
4.5.5	Research Limitations and Future Development.....	104
4.6	Interface Implementation.....	105
4.6.1	Interface Code Structure	105
4.6.2	Interface Functionality	107
CHAPTER V CONCLUSION AND RECOMENDATION		109
5.1	Conclusion.....	109

5.2 Recommendations110
BIBLIOGRAPHY111

LIST OF FIGURE

Figure 3. 1 Research Stages.....	25
Figure 3. 2 Example of Data Translation in Spreadsheet	30
Figure 3. 3 Normalization and Case Folding Process.....	31
Figure 3. 4 Illustration of the Tokenization Process	32
Figure 3. 5 Stopword Removal Process.....	34
Figure 3. 6 Query Expansion Stages	37
Figure 3. 7 Loading the FastText Model	38
Figure 3. 8 Vocabulary Formation Based on Thesis Data.....	38
Figure 3. 9 Embedding Initialization from FastText	39
Figure 3. 10 Word2Vec Fine-Tuning Process	40
Figure 3. 11 Expansion Word Candidate Retrieval.....	41
Figure 3. 12 Expansion Word Candidate Filtering.....	42
Figure 3. 13 IndoSBERT Embedding Formation	44
Figure 3. 14 Fine-Tuning Data Preparation for IndoSBERT	45
Figure 3. 15 IndoSBERT Fine-Tuning Process with Title and Abstract.....	46
Figure 3. 16 Baseline FAISS Index Construction.....	49
Figure 3. 17 End-to-End FAISS Retrieval Flow	50
Figure 3. 18 Document Pooling and Relevance Assessment.....	56
Figure 3. 19 System Evaluation Process	57
Figure 3. 20 System Wireframe.....	58
Figure 4. 1 Visual Evaluation of Short Queries Across All Scenarios	98
Figure 4. 2 Visual Evaluation of Long Queries Across All Scenarios	100
Figure 4. 3 System Interface Display	107

LIST OF TABLES

Table 2. 1 Previous Research	5
Table 3. 1 Sample of Scraped Data	29
Table 3. 2 Data After Normalization and Case Folding	32
Table 3. 3 Sample of Tokenized Data	33
Table 3. 4 Sample of Stopword Removal Results	35
Table 3. 5 Text Representation Columns	36
Table 3. 6 Test Queries.....	52
Table 4. 1 Data Distribution per Faculty	61
Table 4. 2 Sample Query Expansion Results	70
Table 4. 3 IndoSBERT Fine-Tuning Parameters.....	75
Table 4. 4 Expansion Output from the Pooling Process.....	76
Table 4. 5 Retrieval Results Across Five Model Scenarios.....	80
Table 4. 6 Short Query Evaluation Results by Metric.....	96
Table 4. 7 Long Query Evaluation Results by Metric	98

LIST OF FORMULA

(2. 1) Semantic Representation Construction	15
(2. 2) Fine Tuning with <i>Multiple Negative Ranking Loss</i>	16
(2. 3) <i>Cosine Similarity</i>	17
(2. 4) Vector based on Euclidean Distance.....	20
(2. 5) ANN Process	20
(2. 6) Formula of <i>Precision@K</i>	22
(2. 7) Formula of <i>Recall @K</i>	22
(2. 8) Formula of AP	23
(2. 9) Formula of MAP	23
(2. 10) Formula of <i>DCG@K</i>	23
(2. 11) Formula of <i>nDCG@K</i>	23