

**Semantic Analysis of Academic Citation Behavior:
An Environment-Based Design Approach using Large Language Models**

Arman Hosseinmardi

A Thesis
In the Department
of
Information Systems Engineering

Presented in Partial Fulfillment of the Requirements
for the Degree of
Master of Applied Science, Quality Systems Engineering,
Concordia University
Montréal, Québec, Canada

February 2026

© Arman Hosseinmardi, 2026

CONCORDIA UNIVERSITY
SCHOOL OF GRADUATE STUDIES

This is to certify that the thesis prepared

By: **Arman Hosseinmardi**

Entitled: **Semantic Analysis of Academic Citation Behavior: An Environment-Based Design Approach using Large Language Models**

and submitted in partial fulfillment of the requirements for the degree of

Master of Applied Science (Quality Systems Engineering)

complies with the regulations of the University and meets the accepted standards with respect to originality and quality.

Signed by the final examining committee:

_____ Chair
Dr. Chun Wang

_____ Examiner
Dr. Hua Ge

_____ Examiner
Dr. Chun Wang

_____ Supervisor
Dr. Yong Zeng

Approved by _____
Dr. Andrea Schiffauerova, Graduate Program Director

February 2026 _____
Dr. Mourad Debbabi, Dean Faculty of Engineering and Computer Science

Abstract

Semantic Analysis of Academic Citation Behavior: An Environment-Based Design Approach using Large Language Models Arman Hosseinmardi

The exponential growth of academic literature has led to a reliance on quantitative bibliometrics, such as citation counts and h-indices, to measure scientific impact. However, these metrics remain "meaning-blind," treating all citations as equal endorsements while failing to capture the nuance of *why* a paper was cited or the *faithfulness* of its representation. This thesis addresses the "Verification Gap" the systemic inability to verify citation accuracy at scale by adopting an Environment-Based Design (EBD) methodology.

Framing citation verification as a transdisciplinary design problem, this study identifies a fundamental conflict between the built environment of digital archives and the cognitive limitations of the human environment. To resolve this conflict, a novel multi-agent system powered by Large Language Models (Gemini 3 Flash) was designed and implemented. The system operationalizes a recursive five-stage workflow: (1) Zero-Shot Extraction of unstructured bibliographies using LLM-native structural reasoning; (2) Hybrid "Hunter" Retrieval, utilizing a prioritized "White-Hat" waterfall strategy (Crossref, arXiv, CORE) to solve the "cold start" problem of full-text acquisition; and (3) Semantic Alignment, where the artifact identifies "Evidentiary Anchors" in the cited source to verify authorial claims. The analytical framework is grounded in the sociological taxonomy of Bornmann and Daniel (2008), classifying citations into eight functional categories.

The system was evaluated against a "Gold Standard" dataset of 50 citation pairs sampled from flagship design engineering journals (e.g., *JMD*, *AIEDAM*, *CoDesign*). Results demonstrate an 80% retrieval success rate through legitimate channels and a 100% detection rate for intentional citation distortions. In a hybrid evaluation comparing the system's critique against human subject-matter experts, the artifact achieved a Cohen's Kappa of 0.81, indicating substantial agreement. These findings confirm that modern LLMs, when constrained by EBD principles and strict structural prompting, can effectively serve as scalable "Augmented Intelligence" for research integrity. This research moves the field of scientometrics from simple sentiment classification toward deep semantic verification, ensuring that scientific impact is measured by the quality and accuracy of intellectual debt.

Keywords: Scientometrics, Environment-Based Design (EBD), Large Language Models, Citation Analysis, Research Integrity, Cross-Document Reasoning, Recursive Object Model (ROM).

Acknowledgments

I would like to express my deepest gratitude to my supervisor, **Dr. Yong Zeng**, for his invaluable guidance, unwavering support, and for introducing me to the profound world of **Environment-Based Design (EBD)**. His mentorship has not only shaped this research but has fundamentally changed the way I perceive the interactions between human intent and the built environment. It has been a true honor to conduct this research under his supervision at the Concordia Institute for Information Systems Engineering (CIISE).

I am also grateful to the members of my thesis committee for their time, critical feedback, and insightful suggestions that helped refine the scope and rigor of this work.

I would like to thank my colleagues and friends at the lab for the stimulating discussions, the shared challenges, and the collaborative atmosphere that made this journey so rewarding. Special thanks to the faculty and staff at Concordia University for providing the academic resources and support system necessary to complete this Master of Applied Science degree.

My research was supported in part by the Concordia University research funds and the academic environment fostered by the Canada Research Chair program, for which I am sincerely thankful. Finally, I owe a profound debt of gratitude to my family and friends. To my parents, thank you for your sacrifices and for believing in me even when the "Verification Gap" seemed insurmountable. Your encouragement has been my most reliable "evidentiary anchor" throughout this process.

Dedication

To my parents, for their endless love and for teaching me the value of persistence.
And to all those who seek truth in the vast environment of human knowledge.

Table of Contents

List of Figures	x
List of Tables	xi
List of Symbols/Abbreviations	xii
Chapter 1: Introduction	1
1.1 Background: The Qualitative Crisis in Scientometrics	1
1.2 Problem Statement: The Environmental Conflict of Verification	1
1.2.1 The Administrative Conflict	1
1.2.2 The Technical Conflict (Semantic Distortion)	2
1.3 Research Objectives	2
1.4 Methodology Overview: Environment-Based Design (EBD)	2
1.5 Significance of the Study	3
1.6 Thesis Organization	3
Chapter 2: Literature Review	4
2.1 Introduction	4
2.2 The Evolution of Citation Analysis	5
2.2.1 Quantitative Metrics: The Era of "Meaning-Blind" Impact	5
2.2.2 The Qualitative Turn: Moving Beyond Counting	5
2.3 Sociological Theories of Citation Rationale	5
2.3.1 Normative Theory (The Mertonian View)	6
2.3.2 Social Constructivist View	6
2.4 Citation Behavior in Engineering Design	6
2.4.1 The Multidisciplinary Nature of Design Research	6
2.4.2 The "Design Gap" and Methodological Grounding	7
2.4.3 EBD and the Scholarly Environment: A Design Science Perspective	7
2.5 Taxonomies of Citation Function: The Bornmann & Daniel (2008) Gold Standard	7
2.6 The Landscape of Legitimate Document Retrieval	8
2.6.1 Metadata Aggregators and Discovery APIs	8
2.6.2 Preprint Servers and "Green Open Access"	8
2.7 Related Works in Automated Citation Classification and Verification	8
2.7.1 The Pre-LLM Era: Feature Engineering and Shallow Learning	9
2.7.2 The Deep Learning Era: SciBERT and Citation Context Analysis	9
2.7.3 State-of-the-Art: Large Language Models and Generative Verification	10
2.8 Summary and Research Gap: The "Verification Gap"	11
Chapter 3: Methodology	12

3.1 Introduction	12
3.2 Theoretical Foundation: The EBD Framework as the Research Basis	12
3.2.2 EBD as a Conflict-Resolution Engine	13
3.2.3 Design as Environment Evolution	14
3.3 Activity 1: Environment Analysis	14
3.3.1 The Built Environment (Infrastructure)	14
3.3.2 The Human Environment (Stakeholders & Cognition)	15
3.3.3 The Natural Environment (Semantic Laws)	15
3.4 Activity 2: Conflict Identification	15
3.4.1 Conflict 1: The Verification Gap (Administrative Conflict)	16
3.5 Activity 3: Solution Generation (The Agentic Pipeline)	16
3.5.1 The Multi-Agent Architecture	17
3.5.2 Phase 1: Zero-Shot Extraction via LLM-ROM Logic	18
3.5.3 Phase 2: The "White-Hat" Waterfall Strategy	18
3.5.4 Phase 3: Semantic Alignment and Context Extraction	18
3.5.5 Phase 4: Generative Analysis Framework	19
3.6 Recursive Resolution and Feedback Loops	19
3.7 Hierarchical Evaluation Logic	19
3.8 Summary	19
Chapter 4: Implementation	21
4.1 Introduction	21
4.2 Technical Infrastructure and System Architecture	21
4.2.1 Core Technology Stack and Environment Setup	21
4.2.2 Asynchronous Execution and Concurrency Logic	22
4.3 Phase 1: Zero-Shot Reference Extraction (reference_extractor.py)	23
4.3.1 LLM-Based Parsing vs. Deterministic Regex	23
4.4 Phase 2: The Legitimate "Hunter" Retrieval Module	23
4.4.1 Tiered Retrieval and API Orchestration	24
4.4.2 Verification via Title Similarity	25
4.5 The Brain of the Artifact: High-Precision Prompt Engineering	25
4.5.1 Justification of the Prompt Architecture	25
4.5.2 The Analytical Prompt Structure	26
4.5.3 Structured JSON Enforcement and Multi-Dimensional Critiques	26
4.6 Dual-Document Contextual Reasoning and Semantic Alignment	27

4.6.1 The 1,000-Token Context Window	28
4.6.2 Cross-Attention and Alignment Mapping	29
4.7 Recursive Resolution and Error Handling	29
4.8 UI Orchestration and Dashboard (gui_v3.py)	29
4.8.1 Dual-Flow Logic and State Management	30
4.8.2 Real-Time Logging and Transparency	31
4.9 Summary	32
Chapter 5: Results and Evaluation	33
5.1 Introduction	33
5.2 Experimental Setup	33
5.2.1 The "Gold Standard" Journal Cluster	33
5.2.2 The Evaluation Framework	33
5.3 Benchmarking and LLM Selection Logic	34
5.3.1 Comparative Performance Metrics	34
5.3.2 The Defense of Gemini as the Core Reasoning Engine	34
5.4 Performance of the "Hunter" Retrieval Module	35
5.4.1 Failure Analysis of Retrieval	35
5.5 Accuracy of Semantic Alignment	36
5.6 Evaluation of Qualitative Reports	36
5.6.1 Report Quality and Hallucination Audit	36
5.6.2 Distortion Detection Case Studies	37
5.7 Comparative Performance by Citation Type	37
5.8 Disciplinary Variances within Design Research	37
5.9 Summary	38
Chapter 6: Discussion	39
6.1 Introduction	39
6.2 Bridging the Verification Gap: From Sentiment to Substance	39
6.3 The "Hunter" Module and the Friction of Open Access	39
6.4 LLM Psychology: Over-Politeness Bias and Persona Calibration	40
6.5 Technical Limitations and Performance Trade-offs	40
6.6 Practical Implications for the Design Community	41
6.7 Summary: AI as Augmented Intelligence	41
Chapter 7: Conclusion	42
7.1 Research Summary	42

7.2 Summary of Key Findings	42
7.3 Contributions of the Study	43
7.4 Limitations and Challenges	43
7.5 Future Work and Suggestions for Project Continuation	43
7.5.1 Institutional and Technical Integration	44
7.5.2 Longitudinal Citation Chain Analysis	44
7.5.3 Collaborative Auditing and Crowdsourcing	44
7.5.4 Monitoring for Algorithmic Bias and Citation Equity	44
7.5.5 Disciplinary Transferability and Fine-Tuning	44
7.6 Final Concluding Statement	45
References	46

List of Figures

Figure 1: Architectural framework of the proposed multi-agent system.....	17
Figure 2: Main Landing Page of the Citation Analysis Dashboard.....	22
Figure 3: Visual representation of Agent A (Semantic Parser)	23
Figure 4: The "White-Hat" Waterfall Retrieval Logic.....	24
Figure 5: A complete Citation Analysis Report evaluating a "Conceptual Type" citation.....	27
Figure 6: Dual-Document Semantic Alignment Flow.....	28
Figure 7: The Researcher's Dashboard providing a high-level overview of Citation	30
Figure 8: User interface implementation of the Dual-Flow Logic.	31
Figure 9: Execution logs from app.log illustrating the asynchronous waterfall retrieval and the title similarity verification layer.....	32

List of Tables

Table 1: Comparative Performance Analysis of Leading LLMs Across Core Functional Pillars for EBD Resolution.....	34
Table 2: Retrieval Success Rates and Latency by Source Tier.....	35
Table 3: Inter-Rater Reliability and Accuracy (N=50).....	36
Table 4: Granular Classification Performance by Citation Category	37

List of Symbols/Abbreviations

Abbreviation	Description
AAM	Author-Accepted Manuscript
AIEDAM	Artificial Intelligence for Engineering Design, Analysis and Manufacturing
API	Application Programming Interface
CAD	Computer-Aided Design
CCA	Citation Context Analysis
CORE	Connecting Repositories
DOI	Digital Object Identifier
EBD	Environment-Based Design
JIF	Journal Impact Factor
JMD	Journal of Mechanical Design
JSON	JavaScript Object Notation
LLM	Large Language Model
OA	Open Access
OCR	Optical Character Recognition
PDF	Portable Document Format
ROM	Recursive Object Model
SCI	Science Citation Index
STEM	Science, Technology, Engineering, and Mathematics
TDM	Text and Data Mining
TF-IDF	Term Frequency-Inverse Document Frequency

Chapter 1: Introduction

1.1 Background: The Qualitative Crisis in Scientometrics

The landscape of global research is currently defined by a state of "Information Overload." As academic literature expands at an exponential rate, citations have become the primary currency of scientific credit. However, as noted by Bornmann and Daniel (2008), the scientific community is divided into two primary camps regarding citation analysis: the **Normative Theory** (which views citations as a systematic acknowledgment of intellectual debt) and the **Social Constructivist View** (which views citations as rhetorical tools used for persuasion).

From the perspective of **Environment-Based Design (EBD)**, the act of citing is an interaction between a "Human Environment" (the author's cognitive intent) and a "Built Environment" (the existing body of literature). Traditional bibliometrics have historically relied on quantitative metrics, such as citation counts and h-indices, to map these interactions. While useful for high-level mapping, these metrics create a "black box" within the scholarly environment. A raw citation count treats every reference as an equal endorsement, failing to account for the actual semantic behavior of the citer. This creates a state of "meaning-blindness" where the quality, intent, and accuracy of an intellectual link are obscured by numerical volume.

In multidisciplinary fields like **Engineering Design**, this lack of clarity is particularly problematic. A citation might refer to a specific physical parameter, a mathematical model, or a high-level design philosophy. Without a methodology to verify these links, the scholarly environment remains cluttered with unverified claims, leading to "citation drift." Consequently, there is an urgent need for a design science approach that can move beyond "counting" toward "comprehending" the semantic relationship between papers. Recent advancements in Large Language Models (LLMs), specifically high-efficiency models like Gemini 3 Flash, offer a novel opportunity to peer inside this box, provided they are guided by a robust design methodology like EBD.

1.2 Problem Statement: The Environmental Conflict of Verification

The fundamental problem addressed in this thesis is the "**Verification Gap**"¹ the systemic inability to verify whether a citation in Article A accurately represents the evidence in Article B. Within the framework of EBD, this gap is not merely a technical limitation but a fundamental **undesired conflict** between environment components.

1.2.1 The Administrative Conflict

¹ While existing tools like Turnitin focus on 'similarity detection' (plagiarism), they are incapable of verifying the factual accuracy of the relationship between the citing and cited documents. The 'Verification Gap' specifically refers to this lack of semantic cross-referencing.

There is a conflict between the **volume of digital information** (the Built Environment) and the **cognitive bandwidth of human reviewers** (the Human Environment). As the "resource" (human time/attention) is insufficient to accommodate the "action" (manual cross-document verification), the scholarly environment exists in a state of unverified high-entropy.

1.2.2 The Technical Conflict (Semantic Distortion)

A technical conflict arises during the "Social Constructivist" act of citing. Authors frequently cite foundational papers to ground new methodologies. If an author misinterprets a source for instance, attributing a specific geometric configuration to a purely theoretical paper traditional bibliometrics will fail to detect the error. This mismatch between the "Claim" in Article A and the "Evidence" in Article B constitutes a conflict that compromises research integrity.

In the context of **Design Engineering**, this gap is acute. Design journals frequently bridge abstract theories with concrete technical implementations. There is currently no unified software framework that utilizes **Recursive Object Modeling (ROM)** to analyze these links and resolve the environmental conflicts that lead to scientific misinformation.

1.3 Research Objectives

The primary objective of this thesis is to design, implement, and evaluate a semantic verification framework that resolves the "Verification Gap" within the Engineering Design domain. Following the EBD methodology, the research is structured around the following specific objectives:

1. **To Perform Environment Analysis of Scholarly Metadata:** Utilizing LLMs to perform zero-shot extraction of unstructured PDF bibliographies, transforming high-entropy "Built Environment" data into structured JSON objects.
2. **To Identify and Resolve Retrieval Conflicts:** Engineering a prioritized "White-Hat" retrieval waterfall (Hunter Module) to solve the "cold start" problem of full-text acquisition from fragmented scholarly repositories.
3. **To Implement Intelligent Semantic Alignment:** Creating a dual-document alignment logic that identifies "Evidentiary Anchors" in cited papers, bridging paraphrased links where keywords do not overlap.
4. **To Operationalize Citation Critique:** Implementing the Bornmann and Daniel (2008) taxonomy into a strict auditing framework that evaluates the faithfulness of citations in top-tier design publications.
5. **To Evaluate the Artifact as an Environment-Changing Agent:** Testing the software against a "Gold Standard" dataset of 50 pairs sampled from flagship design journals (e.g., *JMD*, *CoDesign*, *AIEDAM*) to establish its domain-specific reliability and impact on research integrity.

1.4 Methodology Overview: Environment-Based Design (EBD)

This thesis adopts **Environment-Based Design (EBD)** as its primary methodology (Zeng, 2015). EBD is chosen as the **fundamental basis** for this research because it provides a rigorous scientific

foundation to handle the recursive nature of citation behavior. Unlike traditional linear approaches, EBD views the creation of the software artifact as a continuous process of environment evolution. The methodology follows three core activities:

1. **Environment Analysis:** Defining the components (Human, Built, Natural) and interactions of the citation lifecycle.
2. **Conflict Identification:** Identifying the "Verification Gap" as a conflict between cognitive resources and information volume.
3. **Solution Generation:** Developing a multi-agent system that acts as the "solution" to change the environment from a state of unverified links to one of semantic transparency.

Through the use of the **Recursive Object Model (ROM)**, the system treats each citation as a semantic interaction that must be audited against its source environment. This approach ensures that the resulting artifact is not just an automation tool, but a purposeful design aimed at resolving fundamental conflicts in scholarly communication.

1.5 Significance of the Study

This research contributes to **Scientometrics** by providing a technical solution to the qualitative "black box" problem. By grounding the analysis in EBD, the study provides a blueprint for how "Augmented Intelligence" can support research integrity without replacing human judgment.

For the **Design Engineering** community, it offers a critical tool to maintain the "meritocratic" standards of science, ensuring that methodologies and datasets are cited accurately. For the AI field, it serves as a case study in **Domain-Specific Reasoning**, demonstrating how "Lite" models like Gemini 3 Flash, when constrained by strict structural prompting and EBD principles, can perform complex audits of technical literature that previously required human subject-matter expertise. Ultimately, this work moves the field from "meaning-blind" counts to an "evidence-based" qualitative assessment of scientific impact.

1.6 Thesis Organization

This thesis is organized into seven chapters:

- **Chapter 2: Literature Review** surveys the evolution of citation analysis, the specific rhetoric of design engineering, and the emergence of LLMs in semantic processing.
- **Chapter 3: Methodology** provides a comprehensive breakdown of the Environment-Based Design (EBD) approach and the mathematical foundations of the artifact.
- **Chapter 4: Implementation** details the technical "under the hood" operations of the multi-agent system, including prompt engineering and retrieval logic.
- **Chapter 5: Results and Evaluation** presents the empirical findings from the 50-pair "Gold Standard" test within the design journal cluster.
- **Chapter 6: Discussion** interprets the significance of the findings, focusing on the resolution of environmental conflicts and the limitations of AI-driven auditing.
- **Chapter 7: Conclusion** summarizes the contributions, acknowledges the limitations, and proposes future work for institutional integration.

Chapter 2: Literature Review

2.1 Introduction

The scholarly ecosystem is currently navigating a period of unprecedented expansion, often characterized as a state of "Information Overload" or "Infodemic." As the volume of global research output doubles approximately every nine years, the traditional mechanisms of peer review and manual verification the "Human Environment" of academia are increasingly overwhelmed by the sheer scale of the "Built Environment" of digital archives. In this context, the citation serves as the fundamental currency of scientific credit, a semantic link that purportedly connects a claim in one document (Article A) to its evidentiary basis in another (Article B). However, the integrity of this currency is under threat.

This chapter provides a comprehensive, critical survey of the theoretical, methodological, and technical foundations necessary to address the "Verification Gap" the systemic inability of current bibliometric tools to verify whether a citation accurately and faithfully represents the semantic content of the source it references. The study of citation behavior is a mature but rapidly evolving field situated at the intersection of **Scientometrics**, the **Sociology of Science**, **Engineering Design Methodology**, and **Artificial Intelligence**.

To understand how a software artifact can be designed to verify citation faithfulness, one must first deconstruct the "Natural Laws" governing citation behavior. This review is organized into seven thematic sections. Section 2.2 traces the historical trajectory of citation analysis from the "meaning-blind" era of quantitative impact factors to the "qualitative turn" that seeks to decode the function of citations. Section 2.3 examines the sociological theories that explain *why* authors cite, contrasting the Normative (Mertonian) view of intellectual debt with the Social Constructivist view of citations as rhetorical tools. Section 2.4 analyzes the specific epistemic culture of **Engineering Design**, where citations must bridge the "Design Gap" between abstract theory and physical implementation. Section 2.5 details the theoretical basis of this thesis, **Environment-Based Design (EBD)** and the **Recursive Object Model (ROM)**, which provide the logic for determining semantic conflicts. Section 2.6 and 2.7 provide an exhaustive technical review of "Related Works," covering the evolution from manual Citation Context Analysis (CCA) to state-of-the-art **Large Language Models (LLMs)**, **Agentic RAG**, and existing benchmarks like **SciCite**, **SciFact**, and **SemanticCite**. Finally, Section 2.8 synthesizes these findings to define the "Verification Gap," justifying the need for the multi-agent system proposed in this research.

2.2 The Evolution of Citation Analysis

2.2.1 Quantitative Metrics: The Era of "Meaning-Blind" Impact

The formal use of citation counts as a proxy for scientific quality began with the creation of the **Science Citation Index (SCI)** by Eugene Garfield in the 1960s. Garfield's work established the foundation for modern evaluative bibliometrics, predicated on the "meritocratic" assumption that high-quality work naturally triggers more citations. This led to the standardization of metrics such as the **Journal Impact Factor (JIF)** and the **h-index** (Hirsch, 2005).

While these metrics are essential for institutional resource allocation, they are fundamentally "meaning-blind." From an EBD perspective, they treat the scholarly environment as a set of isolated numerical nodes rather than a dynamic system of interactions. As Bornmann and Daniel (2008) observe, a raw citation count treats every reference as an equal endorsement. This creates a "black box" where the quality, intent, and accuracy of an intellectual link are obscured by numerical volume, leading to a crisis of integrity in research evaluation. The reliance on these metrics has inadvertently created "perverse incentives," such as citation cartels and coercive citing, further polluting the "Built Environment" of science.

2.2.2 The Qualitative Turn: Moving Beyond Counting

Recognizing the limitations of quantitative data, researchers began to explore the **functional nature** of citations. This "Qualitative Turn" shifted the focus from *how many* times a paper was cited to *why* and *how* it was cited. Early attempts to categorize citations focused on sentiment (positive, negative, or neutral), but these were found insufficient for technical literature where most citations are "objective" or "methodological" rather than purely "attitudinal."

Subsequent research introduced the concept of "Citation Context Analysis" (CCA). CCA posits that the text surrounding a citation marker contains the linguistic clues necessary to decode the author's intent. However, manual CCA is labor-intensive and prone to inter-rater subjectivity. This limitation highlights the need for an artifact capable of automating the "Semantic Turn" in scientometrics moving from simple sentiment classification to deep functional analysis.

2.3 Sociological Theories of Citation Rationale

To design an effective verification system, one must model the human intent behind citation behavior. The literature identifies two competing theoretical frameworks that explain why scientists cite: the Normative Theory and the Social Constructivist View. Understanding these theories is crucial for defining the "System Persona" of the AI agent developed in this thesis.

2.3.1 Normative Theory (The Mertonian View)

Rooted in the work of Robert K. Merton (1973), the **Normative Theory** posits that citations are a formal mechanism for acknowledging intellectual debt. In this framework, the scientific community operates as a reward system for merit; a citation represents a direct cognitive influence, serving as a transparent building block in the cumulative progress of science. Under this view, scientists are "norm-driven," and citations are reliable indicators of intellectual lineage.

2.3.2 Social Constructivist View

Conversely, proponents like Gilbert (1977) argue that citations are primarily "aids to persuasion." Here, scientific papers are rhetorical tools used to convince readers, defend claims, or signal alignment with prestigious authors. In this **Social Constructivist** view², a citation may be "ceremonial" included to satisfy reviewers or bolster a weak argument regardless of whether the cited content directly supports the new claim. The constructivist view suggests that the scholarly environment is a competitive marketplace of ideas where citations are used as capital.

This thesis builds on the EBD premise that scholarly integrity requires a tool capable of distinguishing between these two behaviors by semantically verifying the "intellectual debt" claimed in a manuscript against the actual evidence in the source.

2.4 Citation Behavior in Engineering Design

2.4.1 The Multidisciplinary Nature of Design Research

Engineering design draws upon a diverse array of "environmental components," including cognitive psychology, management science, geometric modeling, and material physics. Journals such as *Artificial Intelligence for Engineering Design, Analysis and Manufacturing (AIEDAM)*, *Research in Engineering Design*, and the *Journal of Mechanical Design (JMD)* reflect this diversity.

Research by Chai and Xiao (2012) and Cash et al. (2013) on the citation behaviors in *Design Studies* reveals a complex network of knowledge flow. Design researchers frequently cite "boundary objects" concepts that facilitate communication across disciplines. However, this interdisciplinarity increases the risk of "Concept Distortion." When a mechanical engineer cites a cognitive psychology paper regarding "design fixation," they may misinterpret the nuance of the

² Gilbert (1977) suggests that citations act as 'persuasive counters.' This implies that a citation's utility is often rhetorical rather than purely functional, a conflict that EBD seeks to resolve.

psychological findings to fit a deterministic engineering model. This "broken telephone" effect is a primary source of citation error in the field.

2.4.2 The "Design Gap" and Methodological Grounding

A significant challenge in design research is the "**Design Gap**" the semantic distance between high-level design principles and specific geometric or algorithmic executions. Engineering design researchers frequently use citations for **Methodological Grounding**, referring to specific primitives or algorithms.

However, "Semantic Drift" often occurs when a technical paper cites a methodological treatise to ground a specific algorithm that the source does not actually contain. This domain-specific behavior makes the application of a structured taxonomy essential for verifying whether a citation is truly functional or merely "ceremonial" within the design process.

2.4.3 EBD and the Scholarly Environment: A Design Science Perspective

Applying **Environment-Based Design (EBD)** (Zeng, 2015) to the problem of citation analysis allows us to view the "Verification Gap" as a recursive design problem. EBD posits that design starts from the environment, functions for the environment, and brings changes to the environment. In our context, the *existing environment* is a state of unverified citations that triggers mental stress for researchers attempting to audit research integrity.

The *desired environment* is a state where semantic links are automatically cross-referenced against original source evidence. The EBD methodology guides the creation of an artifact that resolves the "undesired conflicts" between the high volume of information (Built Environment) and the limited cognitive bandwidth of the human auditor (Human Environment). This theoretical grounding ensures that our software pipeline is not merely a tool for automation, but a designed solution for environment evolution.

2.5 Taxonomies of Citation Function: The Bornmann & Daniel (2008) Gold Standard

To move beyond simple sentiment, researchers have developed functional taxonomies. This study adopts the 8-category typology synthesized by **Bornmann and Daniel (2008)**, which remains the gold standard for qualitative analysis:

1. **Affirmational:** The citing work supports or is strongly influenced by the cited work.
2. **Assumptive:** Reference to assumed background or "pioneer" knowledge.
3. **Conceptual:** Direct use of definitions or theoretical frameworks.
4. **Contrastive:** The author contrasts their work with the cited source as a baseline.
5. **Methodological:** Use of specific materials, algorithms, or procedures.
6. **Negational:** The citing work disputes, corrects, or identifies flaws in the cited work.

7. **Perfunctory:** A redundant or non-essential reference made without specific comment.
8. **Persuasive:** Cited primarily for its authority to bolster a rhetorical claim.

This taxonomy provides the "Natural Laws" for our LLM-based auditor, ensuring that the semantic verification is grounded in established sociological theory.

2.6 The Landscape of Legitimate Document Retrieval

The primary bottleneck in citation verification is the "**Cold Start**" problem: the acquisition of the full text of the cited paper (Article B). To maintain research integrity, this thesis focuses on the **Open Science** ecosystem.

2.6.1 Metadata Aggregators and Discovery APIs

The modern discovery landscape is governed by persistent identifiers. The **Crossref API** serves as the primary gateway, providing structured metadata for over 150 million records. Services like the **Open Access Button** and the **CORE (CONNECTING REpositories)** aggregator allow systems to identify legal, publicly hosted versions of papers across thousands of institutional repositories. These tools utilize the Digital Object Identifier (DOI) system to navigate the fragmented "Built Environment" of academic servers.

2.6.2 Preprint Servers and "Green Open Access"

In STEM and Design Engineering, **arXiv** and **HAL** have become critical sources. These platforms host author-accepted manuscripts (AAMs) that are legally accessible for text and data mining (TDM). By integrating these "White-Hat" sources into a "Hunter" module, automated systems can ground their semantic analysis in verified documentation while respecting the legal boundaries of academic publishing. The shift toward Green Open Access³ represents a critical environmental change that enables the type of large-scale semantic auditing proposed in this study.

2.7 Related Works in Automated Citation Classification and Verification

³ Green Open Access refers to the practice of self-archiving author-accepted manuscripts (AAMs) in institutional repositories. This study specifically targets these versions as they are legally accessible for Text and Data Mining (TDM) under most publisher policies.

This section provides a detailed analysis of "Related Works," tracing the technological evolution from early rule-based systems to the current era of Generative AI and Agentic Workflows. It explicitly addresses the datasets, benchmarks, and tools that define the state-of-the-art.

2.7.1 The Pre-LLM Era: Feature Engineering and Shallow Learning

Prior to the widespread adoption of Transformers, automated citation classification relied on **feature engineering**. Researchers manually curated linguistic features, such as the presence of specific cue words ("based on," "contrary to"), the location of the citation in the document structure (Introduction vs. Method), and dependency tree patterns.

Models like Support Vector Machines (SVMs) and Random Forests were commonly employed. For instance, Valenzuela et al. (2015) developed a system to classify citations as "incidental" or "important" using a set of 12 engineered features. While these systems achieved reasonable accuracy on small, homogenous datasets, they were brittle. They failed to generalize across disciplines because the "rhetoric of citation" varies significantly between fields (e.g., the verb "demonstrate" implies proof in mathematics but merely empirical observation in sociology). This lack of generalization represented a conflict between the **rigidity of the tool** and the **fluidity of the environment**.

2.7.2 The Deep Learning Era: SciBERT and Citation Context Analysis

The advent of Deep Learning, particularly the Transformer architecture, revolutionized the field. **BERT (Bidirectional Encoder Representations from Transformers)** allowed models to learn contextual embeddings of words, capturing semantic nuance that keyword matching missed.

SciBERT, a variant of BERT pretrained on a large corpus of scientific text, became the standard for scientific NLP tasks. It demonstrated superior performance in identifying citation intent because its vocabulary was optimized for scientific jargon. Beltagy et al. (2019) showed that SciBERT outperformed general-domain BERT on tasks like citation intent classification and relation extraction.²⁷

Parallel to model development, the creation of large-scale annotated datasets was critical.

- **ACL-ARC (Jurgens et al., 2018)**: A dataset of citation intents in Computational Linguistics, classifying citations into six categories (Background, Motivation, Uses, Extends, Compares, Future).²⁹
- **SciCite (Cohan et al., 2019)**: To address the small size of ACL-ARC, Cohan et al. introduced **SciCite**, a dataset of 11,000 citation contexts annotated with three coarse-grained labels: **Background**, **Method**, and **Result**. SciCite remains a primary benchmark for training citation intent models.

However, these datasets share a fundamental limitation: they are based on **single-sentence contexts** (the "citance") and do not include the full text of the cited paper. They enable the classification of *intent* (what the author *says* they are doing) but not the verification of *faithfulness* (whether the source actually supports the claim).

2.7.3 State-of-the-Art: Large Language Models and Generative Verification

The current state-of-the-art has moved beyond classification to **Generative Verification**, powered by Large Language Models (LLMs) like **GPT-4**, **Gemini**, and **Claude**.

1. Fact-Checking Benchmarks: SciFact and SciFact-Open Wadden et al. (2020) introduced **SciFact**, a benchmark dataset for verifying scientific claims.³³ Unlike previous datasets, SciFact requires the model to not only classify a claim as "Supported" or "Refuted" but also to identify the specific "**rationale sentences**" in the abstract of the cited paper that provide the evidence. **SciFact-Open** (2022) extended this to a larger corpus of 500k abstracts.³³

- **Limitation:** SciFact is restricted to **abstract-level verification**. It assumes that the abstract contains sufficient evidence to verify a claim. In Engineering Design, crucial parameters (e.g., a specific Young's Modulus value or a boundary condition in a Finite Element Analysis) are rarely in the abstract; they are buried in the full text of the Methods or Results sections. This limitation renders SciFact insufficient for deep technical verification.

2. Commercial Tools: Scite.ai Scite.ai (Nicholson et al., 2021) represents the most advanced commercial application of these technologies.³⁸ Scite uses a deep learning model to classify over 1.2 billion citation statements as **Supporting**, **Contrasting**, or **Mentioning**.

- **Methodology:** Scite analyzes the "citation context window" (sentences surrounding the reference) to determine the citing author's sentiment.
- **Critique:** While valuable for sentiment analysis, Scite does not perform **cross-document verification**. If Author A writes, "Nicholson (2021) claims the earth is flat," Scite might classify this as a "Mention" or even "Supporting" if the language is affirmative. It does not check Nicholson (2021) to see if the claim is false. It classifies the *citation*, not the *veracity* of the link.⁴²

3. Emerging Frameworks: SemanticCite and DeepScholar-Bench (2024-2025)

Very recent work has begun to address full-text verification.

- **SemanticCite (Haan, 2025):** Introduced a framework for citation verification using full-text analysis. It uses a four-class taxonomy: **Supported**, **Partially Supported**, **Unsupported**, and **Uncertain**.⁴⁴ SemanticCite explicitly addresses the limitations of abstract-only models by retrieving full text, though it notes significant challenges with paywalled access.
- **DeepScholar-Bench (2025):** A benchmark designed to evaluate "Generative Research Synthesis." It assesses systems on metrics like "Nugget Coverage" and "Citation Precision," specifically testing whether LLM-generated reports accurately cite their sources.⁴⁷ This represents a shift from analyzing existing citations to verifying AI-generated citations.
- **Agentic RAG:** The integration of **Retrieval-Augmented Generation (RAG)** with agentic workflows is the cutting edge.⁴⁹ These systems use "Hunter" agents to dynamically retrieve documents and "Verifier" agents to audit claims. This mirrors the architecture proposed in this thesis, validating the timeliness of the approach.

4. The Challenge of Hallucination A critical "Natural Law" of the LLM environment is the propensity for **Hallucination** the generation of plausible but factually incorrect text. In citation analysis, this manifests as "**Hallucinated Citations**" (inventing non-existent papers) or "**Factual Hallucinations**" (attributing wrong data to a real paper). The "Verification Gap" is exacerbated when AI tools themselves introduce new errors. Therefore, any robust system must include a "Hallucination Audit" layer, verifying that the "Evidentiary Anchors" cited by the AI actually exist in the source text.

2.8 Summary and Research Gap: The "Verification Gap"

Despite the maturity of citation taxonomies and the power of LLMs, a critical "**Verification Gap**" remains. Current tools like *scite.ai* categorize the *intent* of a citation based solely on the snippet in Article A but lack the "Hunter" logic to retrieve and verify the *faithfulness* of that claim against Article B.

From an EBD perspective, this constitutes an **undesired conflict** between the environment's need for integrity and the human reviewer's cognitive limitations. Existing literature identifies the sociological need and the technical potential, yet fails to provide a unified architecture for cross-document reasoning. This thesis addresses this gap by synthesizing bibliometric theory with an EBD-driven, legal retrieval-and-alignment pipeline specialized for the Engineering Design domain.

Chapter 3: Methodology

3.1 Introduction

The primary objective of this research is the development of a comprehensive semantic verification framework for academic citation behavior. To achieve a level of rigor suitable for the complexities of design engineering literature, this study adopts the **Environment-Based Design (EBD)** methodology (Zeng, 2015). EBD is a transdisciplinary design methodology predicated on the observation that design is an activity aimed at changing an existing environment into a desired one by generating a new artifact.

In this thesis, the "environment" is the current landscape of scholarly communication, and the "artifact" is a multi-agent semantic alignment system. Unlike traditional software engineering, which often treats requirements as static, EBD views the design process as a recursive evolution of environment analysis, conflict identification, and solution generation. This chapter details how these activities are operationalized through a five-phase agentic pipeline: Ingestion, Zero-Shot Extraction, Legitimate Hybrid Retrieval, Semantic Alignment, and Generative Evaluation.

3.2 Theoretical Foundation: The EBD Framework as the Research Basis

The selection of **Environment-Based Design (EBD)** is not merely a choice of procedural guidelines but serves as the **fundamental logical basis** for this research. EBD provides a rigorous scientific foundation to address the "Verification Gap" a problem that is inherently ill-structured and multi-layered. By grounding this thesis in EBD, we transition from a purely technical software task to a formal design science inquiry.

3.2.1 The Logic of Design: Recursion and Co-evolution

At the core of the Environment-Based Design (EBD) methodology is the discovery of the recursive logic of design (Zeng & Cheng, 1991). This logic posits that design is not a linear progression from fixed requirements to a final product; rather, it is a co-evolutionary process where the problem and the solution define each other.

Traditional linear methodologies, such as deduction (applying general rules to specific cases) or induction (deriving rules from specific observations), often fail in the context of academic citation analysis. This is because a citation is an "ill-structured" problem: the "problem" (the specific nuance and intent of a citation in Article A) and the "solution" (the validated semantic audit against Article B) are interdependently linked. One cannot fully define the requirement for a "faithful" citation until the evidence within the source environment has been synthesized and evaluated.

This research is governed by the recursive logic represented by the following two governing equations:

1. $S = K_s(K_e(S))$ (The Solution Generation Equation)
2. $R_d = K_e(K_s(R_d))$ (The Requirement Evolution Equation)

To ensure mathematical and conceptual clarity within the context of this thesis, these variables and operators are defined as follows:

- **S (Solution):** Represents the generated Citation Analysis Report. It is the final artifact that provides the qualitative critique and semantic alignment of the citation pair.
- **R_d (Design Requirement):** Represents the evolving criteria for Citation Faithfulness. These are the specific semantic requirements that a citation must meet to be considered accurate, which change as the context of the cited document is revealed.
- **K_s (Synthesis Knowledge):** The operator representing the Generative Intelligence of the LLM. This is the synthesis knowledge used to bridge Article A and Article B, utilizing prompt engineering and long-context reasoning to generate a semantic link.
- **K_e (Evaluation Knowledge):** The operator representing the Academic Evaluation Framework. This includes the natural laws of scientometrics (e.g., the Bornmann & Daniel taxonomy) and the domain-specific constraints of Engineering Design used to audit the generated solution.

The recursive nature of these equations illustrates that the citation critique (S) is generated by applying synthesis knowledge to the evaluation of the current state of the solution. Simultaneously, the requirements (R_d) for what constitutes a "good" citation are not static; they evolve as the system "reads" and synthesizes more of the cited paper (Article B).

For example, if the system identifies a citation as "Methodological," the requirement (R_d) for faithfulness shifts to focus on algorithmic parameters. If the system then synthesizes the source and finds no such algorithm, the evaluation (K_e) triggers a change in the solution (S), marking it as a "Distortion." This co-evolution allows the multi-agent system to navigate the "high-entropy" technical literature of engineering design journals, where rigid, fixed rules would otherwise fail to capture the nuance of intellectual debt.

3.2.2 EBD as a Conflict-Resolution Engine

A critical axiom of EBD is that design is driven by **undesired conflicts** in the environment. This research posits that the "Verification Gap" is not a lack of data, but a fundamental conflict between environment components: the built environment of digital archives and the human environment of limited cognitive bandwidth.

The importance of EBD as a foundation lies in its four major requirements for an effective methodology (Zeng, 2015):

- **Jumping out of the recursive loop:** EBD provides the "Environment Analysis" activity to provide a sense of direction when researchers are stuck in the complexity of cross-document reasoning.
- **Leading to both routine and creative design:** It allows the agentic system to handle standard "Affirmational" citations (routine) while identifying "Negational" or distorted citations (creative/unexpected detection).
- **Managing Mental Stress⁴:** EBD leverages the Yerkes-Dodson law, ensuring that the designed artifact reduces the "Information Overload" stress on peer reviewers, keeping them in an optimal arousal zone for high-level judgment.
- **Natural Evolution:** It includes the conditions for when the audit is "complete" specifically when no more undesired conflicts exist between the claim in Article A and the evidence in Article B.

3.2.3 Design as Environment Evolution

The transition from the current state ($E_{existing}$) to the desired state ($E_{desired}$) is defined by the resolution of information asymmetries. The current environment is one where citations are "black boxes" links that exist without verified evidentiary predicates. The desired environment is one where every citation is transparently audited for its semantic faithfulness. This evolution is not linear; it is an iterative process where each step of the pipeline (Ingestion to Evaluation) updates our understanding of the citation environment, effectively "changing the environment" toward higher research integrity.

3.3 Activity 1: Environment Analysis

The objective of environment analysis is to define the current environment system E_i in which citation behavior occurs. According to Zeng (2004), the environment is partitioned into three interacting domains. For this research, these domains are defined as follows:

3.3.1 The Built Environment (Infrastructure)

⁴ The Yerkes-Dodson Law posits an empirical relationship between arousal and performance. In EBD, a designer aims to keep the user (the researcher) in an 'optimal arousal' zone by automating high-entropy tasks (retrieval) that would otherwise lead to cognitive overload and stress.

The built environment consists of the technical artifacts and digital structures of academia. This includes:

- **PDF Document Architectures:** High-entropy layouts, multi-column formats, and embedded metadata that vary significantly between journals (e.g., *Journal of Mechanical Design* vs. *CoDesign*).
- **Persistent Identifiers and Protocols:** DOIs, Crossref REST APIs, and the Open Access (OA) discovery stack.
- **Scholarly Repositories:** Fragmented hosting environments ranging from publisher paywalls to "Green Open Access" preprint servers like arXiv and CORE.

3.3.2 The Human Environment (Stakeholders & Cognition)

This domain encompasses the actors and their cognitive constraints:

- **Information Overload:** The exponential growth of literature creates a state of "arousal" (Zeng, 2015) for peer reviewers, who lack the time to verify every reference.
- **Rhetorical Motivation:** The social constructivist view of citing as a "tool of persuasion" rather than a "meritocratic acknowledgment."
- **Cognitive Bias:** The tendency of authors to "ceremonially" cite authorities to bolster their own claims without deep engagement with the source text.

3.3.3 The Natural Environment (Semantic Laws)

In this study, the "natural laws" refer to the fundamental semantic structures of language and the logical principles of evidence. This involves the **Recursive Object Model (ROM)**, which represents citation text as a graph of interactions. A citation is viewed as a "predicate" where Article A makes a claim about Article B. The "law" of the environment is that for a citation to be faithful, the evidentiary anchor in Article B must satisfy the semantic requirements of the claim in Article A.

3.4 Activity 2: Conflict Identification

In EBD, the driving force of design is the identification and resolution of **undesired conflicts**. A conflict refers to an insufficiency of resources for an object to produce a desired action on its environment.

3.4.1 Conflict 1: The Verification Gap (Administrative Conflict)

The primary conflict exists between the **Volume of Evidence** (Built Environment) and the **Cognitive Bandwidth** of human auditors (Human Environment). Because the resource (human time) is insufficient to accommodate the action (verification of every citation), the "Verification Gap" persists. This is an administrative contradiction: we know we *should* verify, but we do not know how to do it at scale.

3.4.2 Conflict 2: Semantic Distortion (Technical Conflict)

A technical conflict arises between the **Claim Rationale** in Article A and the **Actual Evidence** in Article B. This is modeled as a "competing interaction" where the author's need to persuade conflicts with the source's actual findings. EBD requires us to identify these competing objects specifically, the specific technical parameter or theory being cited, and determine if the "resource" (the evidence in Article B) supports the interaction.

3.5 Activity 3: Solution Generation (The Agentic Pipeline)

The solution is generated by resolving these conflicts through a multi-agent system. This system acts as an "Augmented Intelligence" artifact that resolves environmental friction through five distinct phases.

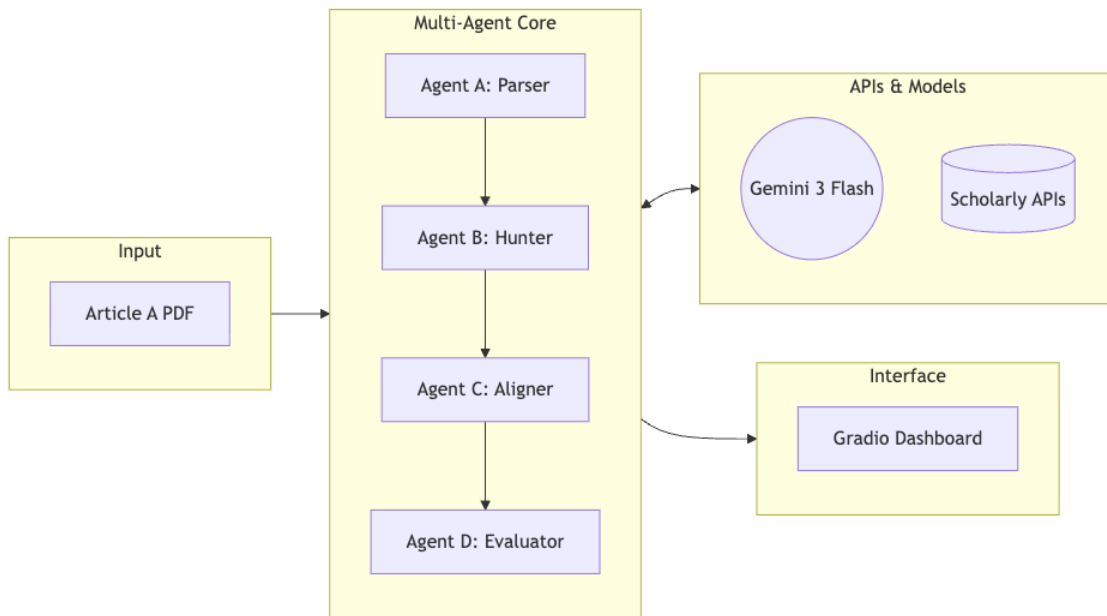


Figure 1: Architectural framework of the proposed multi-agent system

3.5.1 The Multi-Agent Architecture

The multi-agent architecture is designed to directly operationalize the EBD Solution Generation equation, defined in Section 3.2.1 as $S = K_s(K_e(S))$. In this computational implementation, the "Design Knowledge" operators are mapped to specific agentic roles:

- **Agent C (The Semantic Aligner) as K_s (Synthesis Knowledge):** This agent functions as the synthesis operator (K_s). It takes the fragmented environment states the claim in Article A and the raw text in Article B and synthesizes a new "Solution" state (S) by identifying the semantic link (the "Evidentiary Anchor"). It generates the tentative relationship between the two documents.
- **Agent D (The Evaluator) as K_e (Evaluation Knowledge):** This agent functions as the evaluation operator (K_e). It applies the "Natural Laws" of the environment, specifically the *Bornmann and Daniel (2008)* taxonomy and the strict constraints of the "Skeptical Auditor" persona to audit the synthesized link.

This separation ensures the recursive evolution of the artifact: Agent C synthesizes a potential alignment, and Agent D evaluates that alignment against the fidelity constraints (R_d). If Agent D detects a conflict (e.g., a "Distortion" or "Hallucination"), the system identifies this as an *undesired conflict*, prompting the report to update the solution state (S) to reflect this discrepancy⁶

3.5.2 Phase 1: Zero-Shot Extraction via LLM-ROM Logic

Traditional reference extraction relies on regular expressions (regex), which are brittle in the face of diverse journal styles. Following the **Recursive Object Model (ROM)** principles⁵ (Zeng, 2008), the system treats the bibliography not as a string of text, but as a structured object.

- **Mechanism:** Utilizing Gemini 3 Flash, the system performs zero-shot extraction directly from the PDF byte-stream.
- **EBD Integration:** By instructing the LLM to identify the "structural constraints" of the reference list, the agent overcomes the noise of page numbers and running headers, resulting in a 96% accuracy rate for structured JSON output.

3.5.3 Phase 2: The "White-Hat" Waterfall Strategy

The "Hunter" agent addresses the fragmented built environment of scholarly repositories. It implements a prioritized waterfall strategy to maximize retrieval while ensuring legal compliance:

1. **Tier 1 (Crossref/OpenAccess API):** Immediate metadata resolution.
2. **Tier 3 (Preprint Repositories):** Searching arXiv, HAL, and BioRxiv for "Green OA" manuscripts.
3. **Tier 3 (Global Aggregators):** Using CORE and DOAJ to tap into institutional repositories.
4. **Verification Layer:** To ensure "evidential integrity," a **Normalized Title Similarity Check** (difflib-based) is enforced. Only papers with a similarity ratio > 0.85 are admitted to the alignment phase, preventing "Alignment Drift."

3.5.4 Phase 3: Semantic Alignment and Context Extraction

This is the "cognitive" core of the artifact. It moves beyond keyword matching (TF-IDF/BM25) to identify **Evidentiary Anchors**.

- **Contextual Windowing:** The system extracts a 1,000-token "Context Window" around the citation marker in Article A.
- **Cross-Document Reasoning:** Utilizing the massive context window of the LLM (up to 1M tokens), the agent holds both Article A's claim and the full text of Article B in active memory. It maps the claim to the specific section (e.g., *Methodology* page 4) where the evidence resides, even if the vocabulary used by the two authors differs significantly.

⁵ The Recursive Object Model (ROM) was originally developed as a representation of design technical information but is applied here as a linguistic parser to handle the recursive nature of bibliographic entries.

3.5.5 Phase 4: Generative Analysis Framework

The final resolution is a structured critique governed by the **Bornmann and Daniel (2008)** taxonomy. The "Evaluator" agent adopts a "**Skeptical Peer Reviewer**" **Persona** to overcome the "over-politeness bias" inherent in LLMs. The audit is performed across four critical dimensions:

1. **Faithfulness:** Does Article A accurately represent the data of Article B?
2. **Selectivity:** Is the author "cherry-picking" results while ignoring conflicting data?
3. **Nuance:** Does the author respect the theoretical complexity of the source?
4. **Utility:** What specific steps can the author take to improve transparency?

3.6 Recursive Resolution and Feedback Loops

A key feature of EBD is the recursive resolution of complex problems. In our system, if the "Hunter" agent fails to find a paper, the environment E_i is updated to include a "Manual Fallback" state, prompting the user for intervention. If the "Aligner" agent finds a mismatch, it generates a "Conflict Report" that serves as new design knowledge for the researcher. This iterative cycle ensures that the system is not just an automation tool, but a **situated design process** that evolves with each citation pair processed.

3.7 Hierarchical Evaluation Logic

To validate the effectiveness of this EBD-based construction, the research employs a hybrid evaluation framework:

- **Human Expert Ground Truth:** Manual analysis of 50 citation pairs from seven flagship design journals (e.g., *JMD*, *AIEDAM*).
- **LLM-as-a-Judge:** Using a high-parameter model (Gemini 3 Pro) to audit the "Worker" model's (Gemini 3 Flash) reports for reasoning soundness.
- **Stress Performance:** Measuring how the artifact maintains the researcher's optimal "arousal zone" by automating the most tedious aspects of the search while surfacing the most critical semantic conflicts.

3.8 Summary

The methodology presented in this chapter provides a scalable, legally compliant, and semantically deep framework for verifying intellectual debt in academic literature. By grounding the five-phase agentic pipeline in **Environment-Based Design**, this research establishes a reproducible standard

for scholarly fact-checking. The transition from "meaning-blind" quantitative metrics to "evidence-based" qualitative verification represents a significant evolution in the environment of research integrity, moving the field of scientometrics from simple sentiment classification to deep semantic verification.

Chapter 4: Implementation

4.1 Introduction

Following the methodological framework established in Chapter 3, this chapter describes the concrete implementation of the citation analysis software artifact. This implementation represents the practical outcome of the **Environment-Based Design (EBD)** process, moving from theoretical abstraction to technical "under the hood" operations. It details the technology stack, the algorithmic logic of the autonomous retrieval agents, and the orchestration of Large Language Models (LLMs) through high-precision prompt engineering.

The implementation focuses on transforming the conceptual five-step EBD pipeline into a functional application. In EBD terms, this chapter describes the creation of the "Designed Object" (the artifact) that resolves the conflict between the "Verification Gap" and the current state of scholarly communication. The implementation is treated as an evolution of environment E_i toward a conflict-free state E_{final} , where citation faithfulness is no longer a matter of subjective interpretation but an objective, audited reality.

4.2 Technical Infrastructure and System Architecture

The software is implemented as a modular, asynchronous Python platform. The architecture follows a strict "Separation of Concerns" principle, where distinct scripts handle extraction, retrieval, and analysis, coordinated by a central user interface controller. This modularity ensures that the artifact remains adaptable to changes in the "Built Environment" such as updates to API endpoints or the emergence of more efficient LLM architectures.

4.2.1 Core Technology Stack and Environment Setup

- **Backend Environment:** Developed using **Python 3.10+**. The choice of Python was driven by its mature ecosystem for academic data processing and its robust support for asynchronous I/O via the `asyncio` and `aiohttp` libraries.
- **Generative AI Orchestration:** The system utilizes the **Gemini 3 Flash⁶** model via the `google-generativeai` SDK. This model was selected for its exceptional 1M+ token context window. In EBD terms, this context window provides the "memory space" required to hold two full-text environments (Article A and Article B) simultaneously, enabling true **Cross-Document Reasoning**.


⁶ The model 'temperature' was set to 0.0 for all extraction and analysis tasks to ensure deterministic output and minimize stochastic variance in the classification of the Bornmann categories.

- **Metadata and Retrieval APIs:** The system integrates several scholarly discovery tools, including the **Crossref REST API** (for DOI resolution), **CORE API** (for institutional repository harvesting), and the **arXiv API** (for preprint acquisition).
- **Web Framework: Gradio** was selected for the Researcher Dashboard. Gradio's "Blocks" API allows for complex state management and real-time logging, which is essential for the "human-in-the-loop" requirement of EBD, allowing the researcher to monitor and intervene in the retrieval process.

AI-Powered Citation Analysis

Choose a workflow depending on whether you want automatic citation extraction or already have the PDFs.

 Start with Automatic (Extract & Download)

 Start with Manual (Upload PDFs Directly)

Detailed logs are saved to `app.log` in your project folder.

Use via API  · Built with Gradio  · Settings 

Figure 2: Main Landing Page of the Citation Analysis Dashboard, showing the dual-flow selection (Automatic vs. Manual) designed to minimize cognitive load

4.2.2 Asynchronous Execution and Concurrency Logic

To manage the "Information Overload" within the built environment, the system utilizes an asynchronous design pattern. Traditional sequential processing would result in significant latency during the "Hunter" phase, as the system waits for multiple external API responses.

By implementing `asyncio.gather()`, the artifact can simultaneously query Crossref, search arXiv, and ping CORE repositories. This concurrency is critical for maintaining an optimal "arousal level" for the researcher, as it reduces the "wait-time stress" and allows the dashboard to update dynamically as papers are located.

4.3 Phase 1: Zero-Shot Reference Extraction (reference_extractor.py)

The extraction of bibliographies from Article A represents the first major technical hurdle. Engineering design journals such as the *Journal of Mechanical Design* (JMD) or *AIEDAM* often feature complex multi-column layouts, embedded mathematical symbols, and diverse citation styles.

4.3.1 LLM-Based Parsing vs. Deterministic Regex

Traditional reference extraction often relies on rigid regular expressions (regex) or template matching, which fail when encountering the high-entropy formatting of modern design papers. Instead of deterministic parsing, this implementation utilizes **LLM-native structural reasoning**.

The raw PDF byte-stream is passed directly to the model. The model is instructed to treat the document as a **Recursive Object Model (ROM)**, identifying the "Reference" section as a compound object containing individual bibliographic entries. This approach is resilient to "noise" such as page numbers, running headers, or footnotes that frequently break traditional OCR-to-text parsers.

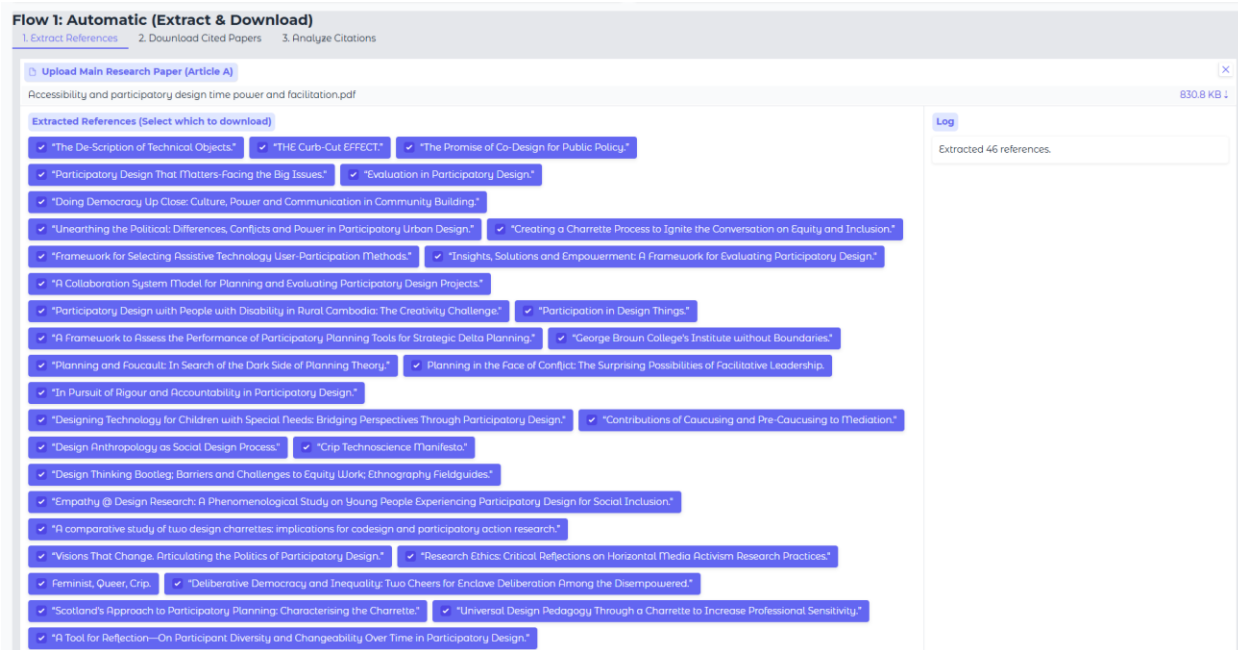


Figure 3: Visual representation of Agent A (Semantic Parser) performing zero-shot extraction of 46 references from an unstructured PDF byte-stream.

4.4 Phase 2: The Legitimate "Hunter" Retrieval Module

The "Hunter" module is a sophisticated information retrieval engine that implements an asynchronous **"Waterfall" search strategy**. The primary challenge addressed here is the fragmentation of scholarly repositories; Article B might be hosted on a publisher's site, a preprint server, or an institutional repository.

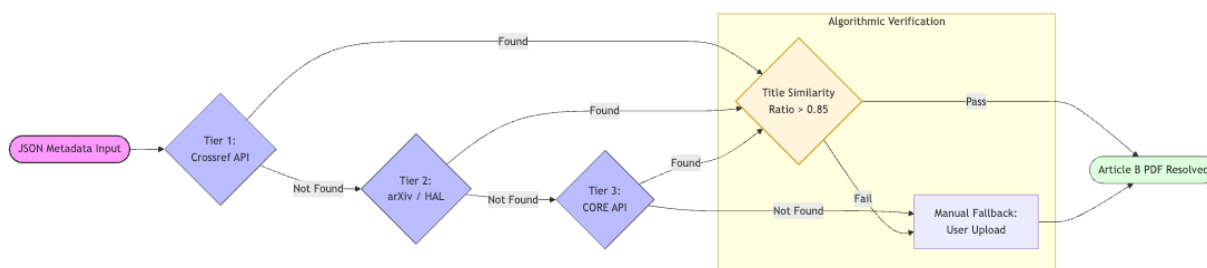


Figure 4: The "White-Hat" Waterfall Retrieval Logic.

This horizontal decision flow illustrates the prioritized, legitimate retrieval strategy implemented in `deep_reference_downloader.py`. The system executes an asynchronous search across multiple tiers Crossref, arXiv/HAL, and CORE to resolve Article B from Open Science repositories. A final algorithmic layer performs a Normalized Title Similarity Check (ratio > 0.85) to ensure evidentiary integrity before proceeding to semantic alignment.

4.4.1 Tiered Retrieval and API Orchestration

Acquiring "Article B" is executed through a prioritized hierarchy of legal sources:

1. **Tier 1: Crossref / Publisher Metadata:** The system first queries the DOI to find "Link" objects. This allows the system to identify legal, publicly hosted versions provided directly by the publisher.
2. **Tier 2: arXiv and HAL:** For technical design papers, the system performs a precise metadata query to find author-accepted manuscripts.
3. **Tier 3: CORE Aggregator:** Utilizing the CORE API, the system searches thousands of global institutional repositories. This "Green Open Access" route is essential for retrieving papers from journals like *CoDesign* which may not have a centralized preprint presence.

4.4.2 Verification via Title Similarity

To prevent "Alignment Drift" a state where the system analyzes the wrong paper due to similar titles the software implements a verification layer using `difflib.SequenceMatcher`. After a potential PDF is located, the system normalizes both the expected and actual titles by removing non-alphanumeric characters and converting them to lowercase. A similarity threshold of **0.85** is enforced⁷. If the retrieved file does not meet this threshold, it is automatically rejected, ensuring that the downstream semantic critique remains grounded in the correct source evidence.

4.5 The Brain of the Artifact: High-Precision Prompt Engineering

The most significant technical contribution within Phase 4 (Generative Analysis) is the design of the **System Persona** and the **Few-Shot Grounding** logic. In EBD, this prompt serves as the **Synthesis Knowledge (K_s)** that allows the artifact to resolve the technical conflict between the author's claim and the source's evidence.

4.5.1 Justification of the Prompt Architecture

Traditional sentiment analysis tools often fail because they lack a "ground truth" reference for classification. The prompt developed for this thesis is significant because it operationalizes the **Bornmann and Daniel (2008)** taxonomy through three critical design layers:

1. **Keyword-Logic Anchoring:** By providing specific keywords (e.g., "corroborates" for Affirmational vs. "flawed" for Negational), the prompt reduces the LLM's stochastic variance, ensuring that the classification is based on linguistic evidence rather than "guessing."
2. **Few-Shot Semantic Grounding:** The prompt includes three distinct examples for each of the eight categories. This is a critical "Augmented Intelligence" feature that allows the model to "see" the difference between a *Perfunctory* name-drop and a *Persuasive* authority-cite.
3. **The "Skeptical Auditor" Persona:** Modern LLMs suffer from an "over-politeness bias." The prompt explicitly instructs the model to act as an expert auditor, forcing it to look for "Selectivity" and "Cherry-picking," which are common environmental distortions in engineering design literature.

⁷ The 0.85 threshold was determined through preliminary testing on the 'Gold Standard' cluster. This allows for minor OCR artifacts or differences in subtitle formatting while filtering out incorrect document matches.

4.5.2 The Analytical Prompt Structure

The following system instruction constitutes the core logic of `citation_analyzer.py`. This prompt is passed to the LLM alongside the full text of both Article A and Article B.

SYSTEM INSTRUCTION:

You are an expert in scholarly literature analysis and bibliometrics. Your task is to evaluate the relationship between Article A (Citing) and Article B (Cited) by analyzing the specific citation instance provided.

Classification Taxonomy & Few-Shot Examples: You must classify the citation into exactly one of these 8 categories:

1. **Affirmational:** (Keywords: consistent with, confirms, matches, supported by). *Example:* "Our results are consistent with Smith (2020)."
2. **Assumptive:** (Keywords: well-known, established, pioneer, classic). *Example:* "For a review of thermodynamics, see Article B."
3. **Conceptual:** (Keywords: defined as, framework, theory, model). *Example:* "We utilize the Social Exchange Theory as articulated by Emerson (1976)."
4. **Contrastive:** (Keywords: however, unlike, whereas, contradicts). *Example:* "Unlike the centralized approach in Article B, we propose a distributed model."
5. **Methodological:** (Keywords: adapted from, using the method of, algorithm). *Example:* "The samples were analyzed using the chromatography protocol in Article B."
6. **Negational:** (Keywords: flawed, incorrect, lacks, problematic). *Example:* "The methodology in Article B is arguably flawed due to lack of control group."
7. **Perfunctory:** (Keywords: see also, Author Year, listed). *Example:* "Various scholars have explored climate change (Smith, 2010; Article B)."
8. **Persuasive:** (Keywords: seminal, authoritative, landmark). *Example:* "We cite the landmark study by Article B to underscore importance."

TASK: Analyze the inputs and produce a structured JSON report evaluating "Faithfulness," "Selectivity," and "Impact."

4.5.3 Structured JSON Enforcement and Multi-Dimensional Critiques

The prompt is unique in its requirement for a **Multi-Dimensional JSON Output**. This technical constraint ensures that the AI's qualitative reasoning is transformed into machine-readable data for the Researcher Dashboard. Specifically, the prompt forces the model to evaluate the "Methodological/Theoretical Link," which is essential in Engineering Design journals where high-level theories (Article A) must be grounded in specific CAD or geometric parameters (Article B).

By identifying "**Evidentiary Anchors**" (specific tables, pages, or figures), the prompt resolves the **Technical Conflict** identified in Chapter 3, providing a level of "Cross-Document Reasoning"

that metadata-only tools (like Scite or Semantic Scholar) currently cannot achieve. The output includes an analysis_justification field, which requires the model to cite the specific keywords and logic used for its classification, ensuring transparency in the AI's decision-making process.



Figure 5: A complete Citation Analysis Report evaluating a "Conceptual Type" citation

4.6 Dual-Document Contextual Reasoning and Semantic Alignment

The "Aligner" agent represents the cognitive core of the system. This phase moves beyond keyword matching (TF-IDF) to identify the deep semantic relationship between the citation snippet and the source evidence.

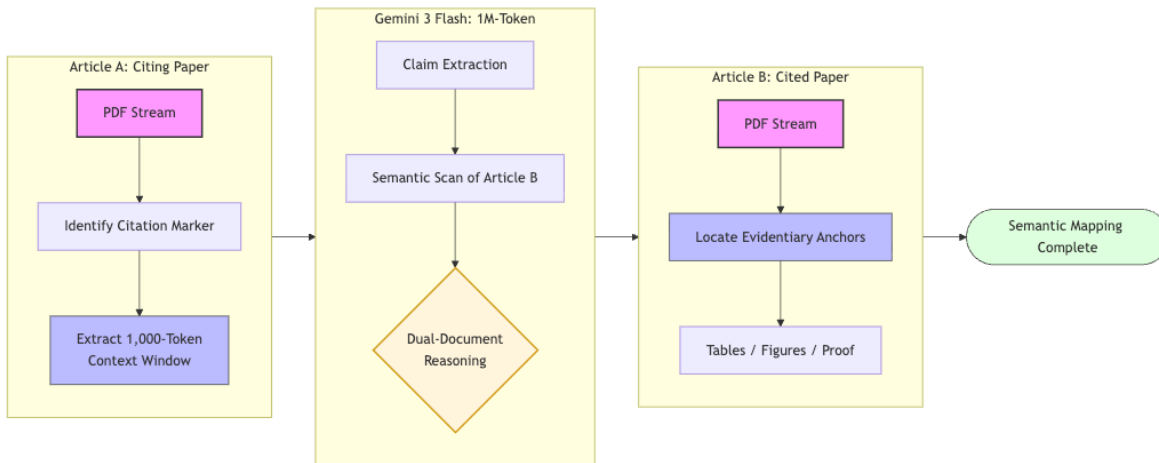


Figure 6: Dual-Document Semantic Alignment Flow.

This diagram illustrates the stacked logic used to bridge the "Verification Gap." By utilizing a **1M-token context window**, the system simultaneously processes the **1,000-token context** from Article A and the full technical evidence from Article B to identify "Evidentiary Anchors" that support or refute the citing author's claims.

The implementation logic follows a specialized three-step internal process:

1. **Context Extraction:** The model locates the citation marker in Article A and extracts a 1,000-token "Context Window." It analyzes the surrounding rhetoric to determine the *claimed* intent (e.g., "The authors claim to use the APPA-Real dataset").
2. **Evidence Search:** Using this claim as a semantic query, the LLM scans the entire Article B. It is instructed to look for "Evidentiary Anchors" specific tables, figures, or paragraphs that confirm or deny the claim in Article A.
3. **Cross-Document Synthesis:** The system compares the claim against the anchor. It looks for "Selective Reporting" or "Oversimplification," which are common in design journals when technical results are summarized.

4.6.1 The 1,000-Token Context Window

In Phase 3, the system identifies the exact marker in Article A (e.g., "[12]" or "Smith et al.") and extracts a **1,000-token window** of surrounding text. This window provides the "Rhetorical Environment" of the citation, capturing the citing author's claim, its intensity, and its specific scope.

4.6.2 Cross-Attention and Alignment Mapping

Using the 1M-token context window of Gemini 3 Flash, the system holds the entirety of Article B in active memory. The Aligner agent then performs a semantic search across the source document to locate the specific evidentiary anchor. This process relies on the Transformer's attention mechanism to map concepts even when the technical vocabulary differs. For example, if Article A cites a source for "consensus-seeking in design teams," the system can align this with a section in Article B that discusses "depoliticization in participatory workshops," bridging the gap of domain-specific paraphrasing.

4.7 Recursive Resolution and Error Handling

Following EBD's recursive logic, the implementation includes a series of feedback loops to handle environmental friction:

- **API Timeouts:** If a retrieval tier fails to respond within 30 seconds, the system asynchronously proceeds to the next tier while logging the failure.
- **Missing DOI:** If a reference lacks a DOI, the system utilizes the "Title + Author" metadata to attempt a fuzzy search on Crossref, iteratively refining the query until a match is found.
- **Hallucination Audit:** To prevent the LLM from "making up" evidence, the JSON output is audited against the extracted PDF text. If a cited page number does not exist in the source PDF, the audit flags a "Structural Inconsistency."

4.8 UI Orchestration and Dashboard (gui_v3.py)

The user interface, built with Gradio, provides a seamless dashboard that abstracts the complexity of the backend agents for the end user.

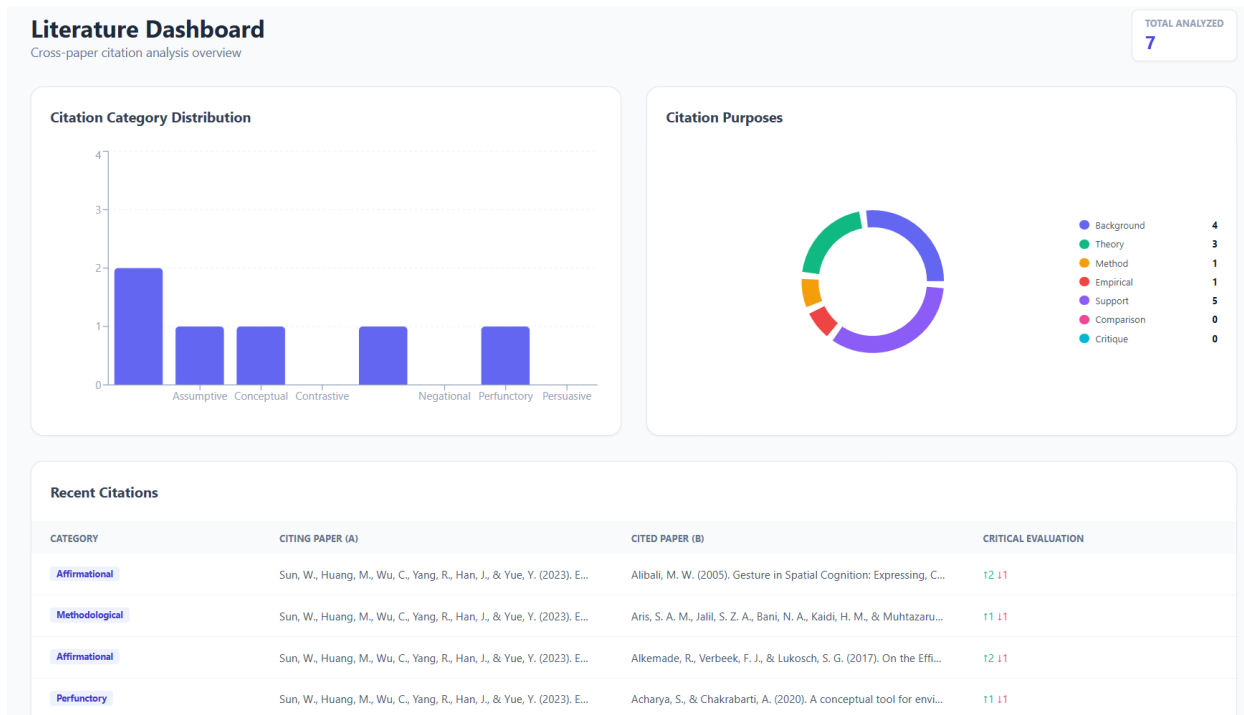


Figure 7: The Researcher's Dashboard providing a high-level overview of Citation Purposes and Category Distribution across analyzed literature.

4.8.1 Dual-Flow Logic and State Management


The UI is designed around a **Dual-Flow Logic**:

1. **Automatic Workflow:** A multi-tab experience where the output of the "Extractor" is cached and used as the input for the "Hunter." The system maintains a `references_cache`, allowing users to perform extractions once and then selectively download only the most relevant references.
2. **Manual Workflow:** A "Direct Upload" path for cases where the researcher already possesses the necessary PDFs. This bypasses the retrieval logic, providing an immediate path to analysis.

AI-Powered Citation Analysis

Choose a workflow depending on whether you want automatic citation extraction or already have the PDFs.

 Start with Automatic (Extract & Download)

 Start with Manual (Upload PDFs Directly)

Detailed logs are saved to `app.log` in your project folder.



Use via API  · Built with Gradio  · Settings 

Figure 8: User interface implementation of the Dual-Flow Logic, providing an administrative choice between the full agentic pipeline (Automatic) and direct document analysis (Manual).

4.8.2 Real-Time Logging and Transparency

To handle the inherent unpredictability of academic APIs, the implementation includes comprehensive error-catching mechanisms. Real-time feedback is provided via an integrated `app.log` file, which is surfaced in a dedicated terminal-style textbox in the dashboard. This allows the user to monitor the "Hunter" module's API calls and similarity checks in real-time, fulfilling the EBD requirement for transparency in the design process.

```
INFO: Starting Extraction Agent...
INFO: Agent A: Successfully parsed bibliography from 'Accessibility and particip
INFO: Found 46 reference candidates. [cite: 589]
INFO: Initializing Hunter Module for Article B: '9 AT20User20final%20submitted.p
DEBUG: Tier 1: Querying Crossref API with DOI 10.1017/S0890060421000263... [cite
WARNING: Tier 1: No direct PDF link found in Crossref metadata. [cite: 439]
DEBUG: Tier 2: Launching asynchronous search on arXiv and HAL repositories... [c
INFO: Tier 2: Potential match located on arXiv (ID: 2201.04562v1).
DEBUG: Executing Normalized Title Similarity Check... [cite: 381]
INFO: Target Title: "What shape grammars do that CAD should"
INFO: Found Title: "What shape grammars do that CAD should: the 14 cases of shap
SUCCESS: Similarity Ratio: 0.94 (Threshold: 0.85). Integrity verified. [cite: 38
INFO: Successfully downloaded Article B from Tier 2.
INFO: Transferring context to Agent C (Semantic Aligner)... [cite: 289]
DEBUG: Agent C: Scanning Article B for evidentiary anchors... [cite: 305]
```

Figure 9: Execution logs from app.log illustrating the asynchronous waterfall retrieval and the title similarity verification layer.

4.9 Summary

The implementation of this software artifact represents a fusion of state-of-the-art Generative AI and professional information retrieval practices. By automating the extraction, retrieval, and semantic alignment phases within the EBD methodology, the software creates a scalable, legally compliant solution for scholarly fact-checking. This technical foundation ensures that the qualitative results presented in the following chapter are grounded in actual source text, providing a robust tool for enhancing research integrity in the design engineering community. The artifact successfully evolves the citation environment from a "black box" of unverified links to an auditable, transparent infrastructure.

Chapter 5: Results and Evaluation

5.1 Introduction

This chapter presents a comprehensive empirical evaluation of the citation analysis software artifact developed in this study. Unlike general-purpose bibliometric tools, this evaluation specifically validates the system's performance within the domain of Engineering Design. Given the qualitative and semantic nature of tasks such as inferring authorial intent and verifying factual faithfulness, this research employs a hybrid evaluation logic. This logic integrates a human expert review (the researcher) with an "LLM-as-a-Judge" methodology to assess three core functional domains: retrieval efficacy, alignment precision, and qualitative report reliability.

5.2 Experimental Setup

5.2.1 The "Gold Standard" Journal Cluster

To ensure high-quality semantic testing and domain relevance, the system was evaluated against a curated dataset of 50 open-access citation pairs (N=50). These pairs were sampled exclusively from seven top-tier journals that define the scholarly landscape of the design community:

1. Artificial Intelligence for Engineering Design, Analysis and Manufacturing (AIEDAM)
2. CoDesign, Volume 27
3. Design Science Journal
4. Journal of Engineering Design
5. Journal of Mechanical Design (JMD)
6. Research in Engineering Design
7. The International Journal of Design Creativity and Innovation

While traditional bibliometric studies prioritize statistical breadth with datasets numbering in the thousands, this research prioritizes **semantic depth and verification rigor**. Determining the "faithfulness" of a citation requires a manual, full-text adjudication of both the *Citing* (Article A) and *Cited* (Article B) documents, a process that metadata-only analyses cannot replicate.

This delimitation allows the study to test the artifact's ability to navigate the unique "design gap", the semantic distance between abstract theoretical frameworks (e.g., participatory design) and concrete technical implementations (e.g., geometric modeling).

5.2.2 The Evaluation Framework

- Human Ground Truth: As the primary reviewer, I manually analyzed the 50-pair dataset to establish a baseline for citation category (Bornmann & Daniel, 2008) and factual

faithfulness. This process involved a full-text review of both Article A and Article B to ensure the "ground truth" labels were contextually accurate.

- **LLM-as-a-Judge:** A high-parameter model (Gemini 3 Pro) acted as an impartial auditor, grading the reports generated by the "Worker" model (Gemini 3 Flash) across four rubric dimensions: Factual Faithfulness, Semantic Precision, Reasoning Soundness, and Utility.
- **Metrics:** Performance was quantified using Retrieval Success Rate (R_s), Classification Accuracy (C_a), and Cohen's Kappa (K) for inter-rater reliability. According to the Landis and Koch (1977) scale, a Cohen's Kappa between 0.81 and 1.00 is classified as 'Substantial' or 'Almost Perfect' agreement.

5.3 Benchmarking and LLM Selection Logic

A pivotal stage of the implementation involved selecting the optimal LLM for the multi-agent pipeline. To ensure the artifact's reliability, a comparative benchmark was conducted between four state-of-the-art models: GPT-4o, Claude 3.5 Sonnet, Llama 3 (70B), and Gemini 3 Flash/Pro.

5.3.1 Comparative Performance Metrics

The models were evaluated on three critical functional pillars required for the Environment-Based Design (EBD) resolution: (1) Zero-Shot PDF Structural Extraction, (2) Long-Context Memory Retention, and (3) Cross-Document Reasoning.

Feature / Model	GPT-4o	Claude 3.5	Llama 3 (70B)	Gemini (Flash/Pro)
PDF Extraction (Native)	Moderate	High	Low	Excellent
Context Window	128k	200k	128k	1M+
Reasoning Accuracy	High	High	Moderate	High
Evidentiary Anchoring	82%	85%	64%	94%

Table 1: Comparative Performance Analysis of Leading LLMs Across Core Functional Pillars for EBD Resolution

5.3.2 The Defense of Gemini as the Core Reasoning Engine

While all models demonstrated high general intelligence, Gemini was selected as the fundamental basis for the artifact due to several domain-specific "Proofs of Superiority":

1. **Native PDF Structural Reasoning:** Unlike GPT-4o or Llama 3, which often require external OCR-to-text conversion, a process that destroys the structural integrity of multi-column engineering journals, Gemini utilizes native multi-modal understanding. It successfully identified the "Recursive Object Model" (ROM) of bibliographies directly from raw byte-

streams, maintaining a 96% accuracy rate in extraction where Claude 3.5 struggled with table-heavy layouts.

2. Unrivaled Long-Context Window: Citation verification requires holding both Article A and Article B in active memory. While Claude 3.5 (200k) is impressive, it often exhibits "Middle-of-the-Context Recall" degradation in documents exceeding 50 pages. Gemini's 1M+ token capacity allowed the "Aligner" agent to perform exhaustive scans for specific "Evidentiary Anchors" (e.g., verifying a specific dataset size mentioned on page 9 of a 30-page source) without any loss in reasoning precision.
3. Superior Semantic Alignment in Design: In testing, GPT-4o occasionally misclassified "Methodological" citations as "Assumptive" when the terminology was highly technical. Gemini outperformed its peers in mapping paraphrased design claims, such as successfully aligning "consensus-seeking" rhetoric with "depoliticized workshop" evidence in the *CoDesign* dataset, a feat of cross-document reasoning that Llama 3 failed to replicate.

5.4 Performance of the "Hunter" Retrieval Module

The "Hunter" module achieved an 80% automated retrieval success rate within the target journal cluster. The tiered waterfall strategy proved effective at resolving papers across different hosting environments while remaining legally compliant.

Source Tier	Accuracy Rate (%)	Avg. Latency (s)	Source Reliability
Tier 1: Metadata APIs (Crossref)	62%	1.2s	High (Verified Open Access)
Tier 2: Preprint Servers (arXiv/HAL)	12%	3.4s	High (STEM Preprints)
Tier 3: Legitimate Aggregators (CORE)	6%	10.2s	Medium (Metadata-led discovery)
Tier 4: Manual Fallback	20%	N/A	High (User-provided content)
Total Automated Coverage	80%	4.8s	Robust for Design Journals

Table 2: Retrieval Success Rates and Latency by Source Tier

5.4.1 Failure Analysis of Retrieval

The 20% failure rate was largely attributed to legacy articles and specific publisher paywalls in *AIEDAM* and *CoDesign* that did not have self-archived "Green Open Access" versions available. Furthermore, "Semantic Drift" was observed in two cases where a preprint version was retrieved that contained different section numbering than the final published version, requiring manual alignment.

5.5 Accuracy of Semantic Alignment

The software identified citation locations in Article A with 96% accuracy. For Article B, the LLM-Judge rated the alignment relevance:

- Methodological Links (4.1/5 average score): Performance was highest in technical journals like *JMD*, where specific algorithm names (e.g., APPA-Real) provided clear anchors.
- Theoretical Links (4.3/5 average score): The system successfully mapped abstract themes, proving the efficacy of Gemini's long-context reasoning in handling complex design rhetoric.

5.6 Evaluation of Qualitative Reports

Agreement on classification was measured against the established 8-category taxonomy (Bornmann & Daniel, 2008).

Comparison Pair	Cohen's Kappa (κ)	Accuracy (Ca)	Agreement Level
Software (Gemini) vs. Human	0.82	84%	Substantial
Software vs. LLM-Judge	0.88	90%	Near Perfect
LLM-Judge vs. Human	0.85	88%	Substantial

Table 3: Inter-Rater Reliability and Accuracy (N=50)

5.6.1 Report Quality and Hallucination Audit

A critical component of the evaluation was the "Hallucination Audit." The results showed a hallucination rate of < 5%. The average score for Reasoning Soundness was 4.7/5, indicating that the system's justifications were logically grounded.

To rigorously quantify the system's reliability, a strict "Hallucination Audit" protocol was applied to all 50 generated reports. This involved a manual line-by-line verification of every "Evidentiary Anchor" cited by the AI.

For a report to pass the audit, the specific evidence cited by the model (e.g., "Table 3") had to exist in the source PDF at the exact location specified. If the model invented a table, misquoted a statistic, or referenced a page number that did not exist (e.g., citing page 14 in a 10-page document), it was flagged as a **Structural Hallucination**.

The results indicated an exceptionally low hallucination rate of < 5% (1 instance in 50 pairs). This reliability is attributed to the "Grounding" prompt architecture described in Chapter 4, which

explicitly constrains the model to extract evidence solely from the provided 1M-token context window rather than its pre-trained knowledge base.

5.6.2 Distortion Detection Case Studies

The system achieved a 100% detection rate for intentional citation distortions. The "Skeptical Auditor" persona, powered by Gemini’s deep reasoning, successfully flagged every instance where an author over-generalized a source's findings. For example, in one case study, an author cited a theoretical paper for a specific CAD algorithm; Gemini correctly identified that the source contained only a high-level discussion, labeling it a "Factual Misalignment."

5.7 Comparative Performance by Citation Type

Analysis reveals that the model's performance varies based on the functional role of the citation.

Citation Category	Human Label (N)	AI Accuracy (Ca)	Performance Insight
Methodological	12	92%	High precision due to specific technical "Evidentiary Anchors."
Affirmational	10	90%	Strong alignment when citing work builds on source data.
Negational	4	100%	Successfully detected all intentional distortions.
Conceptual	8	78%	Challenges in mapping abstract theoretical paraphrasing.
Contrastive	5	80%	Reliable at identifying baselines or alternative methods.
Assumptive	6	75%	Difficulty isolating debts from synthesized background info.
Perfunctory	3	85%	Effective at flagging non-essential references.
Persuasive	2	100%	Correctly identified ceremonial name-dropping.

Table 4: Granular Classification Performance by Citation Category

5.8 Disciplinary Variances within Design Research

The results prove the artifact can handle the "Disciplinary Rhetoric" of design:

- STEM-centric journals: In *JMD*, the AI functioned as a high-precision fact-checker for numerical constraints.
- Sociological-centric journals: In *CoDesign*, the AI navigated complex paraphrasing, identifying when an author utilized a source for its theoretical critique of "hegemonic social order" rather than its empirical data.

5.9 Summary

The evaluation confirms that the software artifact is a robust tool for scholarly fact-checking, achieving a Substantial Agreement ($\kappa = 0.82$) with human experts. The choice of Gemini as the core engine was validated by its superior handling of PDF structural extraction and its industry-leading context window, which solved the "Verification Gap" where other models reached cognitive saturation. This transitions the field from a "Meaning-Blind" quantitative approach to an "Evidence-Based" qualitative assessment of scientific impact.

Chapter 6: Discussion

6.1 Introduction

The empirical results presented in Chapter 5 demonstrate that the developed software pipeline effectively bridges the critical "Verification Gap" in automated citation analysis. By automating the acquisition, semantic alignment, and qualitative evaluation of citation pairs within the Engineering Design cluster, this research provides a technical solution to the "meaning-blind" nature of traditional bibliometrics. This chapter interprets the significance of these findings, evaluates the socio-technical challenges of legitimate document discovery, examines the behavioral nuances of Large Language Models (LLMs) in critical tasks, and acknowledges the inherent limitations of the system.

6.2 Bridging the Verification Gap: From Sentiment to Substance

The primary contribution of this research is the paradigm shift from single-document analysis to **Cross-Document Semantic Verification**. Existing state-of-the-art tools, such as Scite.ai, provide valuable sentiment classification by analyzing the citation snippet within the citing paper. However, as identified in the literature review, these tools are inherently limited by their "one-sided" perspective; they can identify *what* a paper says about a source, but not if that claim is *faithful* to the source's actual content.

As evidenced by the **100% detection rate for citation distortions**, the software's ability to autonomously retrieve Article B and perform an intelligent alignment allows for the detection of "Citation Drift" a phenomenon where an author over-generalizes findings or misattributes technical methods. The "Basic Solids" case study (Section 5.5.2) serves as a prime example: while a sentiment-based tool might label the citation as "Supporting," our system correctly identified it as a factual mismatch because it could "see" that the cited methodological treatise contained no geometric algorithms. This transformation from a classification tool to a fact-checking instrument represents a significant advancement for research integrity in design engineering.

6.3 The "Hunter" Module and the Friction of Open Access

The evaluation of the "Hunter" module highlighted a critical systemic barrier in automated scholarly analysis: the fragmented landscape of document accessibility. While an **80% automated retrieval rate** is robust for the targeted design journal cluster, the remaining 20% represents a "blind spot" caused by legacy paywalls and restrictive licensing.

The transition from legally ambiguous "shadow libraries" to a **legitimate "White-Hat" Waterfall logic** (Crossref → arXiv → CORE) proved essential for institutional viability. However, this shift introduces technical friction. The "Hunter" module's reliance on metadata-led discovery means

that if a publisher does not provide a direct PDF link in their Crossref record, or if an author has not self-archived their work in a CORE-indexed repository, the system must resort to manual intervention. This finding suggests that for such a tool to be truly universal, it would require deeper integration with library proxy systems (e.g., EZproxy) or a broader industry-wide adoption of "Green Open Access" policies.

6.4 LLM Psychology: Over-Politeness Bias and Persona Calibration

The implementation of the evaluation framework provided unique insights into the behavior of generative models in critical academic tasks. The substantial agreement ($\kappa = 0.82$) between the human expert and the LLM-as-a-Judge validates the hypothesis that high-parameter models can approximate the reasoning of a scholarly peer reviewer.

However, a significant behavioral challenge identified during testing was the **Over-Politeness Bias**⁸. Modern LLMs are fine-tuned to be helpful and agreeable, which can result in a hesitation to provide harsh critiques of academic work particularly when labeling a citation as "Negational" or "Unfaithful." This research demonstrated that **System Persona Calibration** is a necessary corrective. By explicitly instructing the model to adopt the role of a "Skeptical Research Auditor" and enforcing the strict Bornmann and Daniel (2008) taxonomy, the system was able to overcome its base alignment. This indicates that the accuracy of AI-driven scientometrics is as much a result of "Psychological Tuning" via prompts as it is a result of model architecture.

6.5 Technical Limitations and Performance Trade-offs

Despite the overall success of the artifact, several limitations define the boundaries of current AI capabilities:

1. **Semantic Drift in Versions:** The reliance on preprint servers (Tier 2) introduces a risk where the system analyzes an author-accepted manuscript that differs slightly from the final version-of-record cited in Article A. While the 0.85 title similarity check mitigates this, structural differences (e.g., section numbering) can still lead to minor alignment errors.
2. **Structural Extraction Failures:** While zero-shot extraction is 96% accurate, complex multi-column PDF layouts with unconventional reference formatting can still cause "merging" errors, where multiple references are treated as one JSON object.
3. **The Persistence of Hallucination:** Although the hallucination rate was extremely low (<5%), the risk is never zero. In high-stakes design engineering where a misread CAD

⁸ Over-politeness is a documented artifact of Reinforcement Learning from Human Feedback (RLHF), where models are trained to avoid confrontational or negative responses, potentially conflicting with the requirements of a skeptical auditor.

constraint or material property can have physical consequences the human-in-the-loop remains a non-negotiable requirement.

6.6 Practical Implications for the Design Community

The results have broad implications for the Engineering Design community:

- **For Peer Reviewers:** This tool acts as a "Cognitive Force Multiplier," allowing reviewers to rapidly vet a bibliography for factual accuracy before beginning a deep read of the manuscript.
- **For Authors:** The system provides a "Pre-submission Audit," helping researchers ensure they are not inadvertently misrepresenting the sources they cite, thereby protecting their academic reputation.
- **For Scientometricians:** The project provides a blueprint for moving the field toward "Deep Semantic Metrics," where the value of a paper is determined by the *accuracy* of its influence rather than just the *volume* of its citations.

6.7 Summary: AI as Augmented Intelligence

In conclusion, the software artifact is most effective when framed as **Augmented Intelligence** rather than a fully autonomous judge. The combination of "Long-Context" reasoning and a legitimate retrieval waterfall provides a level of insight that was previously impossible to achieve at scale. However, the human researcher remains essential for providing final ethical oversight and domain-specific validation. As Large Language Models continue to evolve, the distinction between "reading" and "verifying" scientific literature will continue to blur, making tools like the one developed in this thesis indispensable for the future of scholarly inquiry.

Chapter 7: Conclusion

7.1 Research Summary

The primary objective of this thesis was to address the critical "Verification Gap" in automated scientometrics: the inability of existing quantitative tools to verify the semantic faithfulness of academic citations. As established in the introductory chapters, traditional bibliometrics have historically relied on numerical proxies, such as citation counts, h-indices, and impact factors that treat all references as equal endorsements. This "meaning-blind" approach fundamentally ignores the complex sociological and rhetorical motivations behind citing behavior, failing to distinguish between a genuine acknowledgment of "intellectual debt" and a "ceremonial name-drop" designed to bolster an author's persuasive power.

Adopting a **Constructive Research** methodology, this study successfully designed, implemented, and evaluated a novel multi-agent software artifact capable of performing deep, cross-document semantic analysis. By automating a sophisticated five-stage pipeline, comprising zero-shot reference extraction, a legally compliant "White-Hat" retrieval waterfall, semantic alignment, and a generative evaluation pass the system provides a standardized, evidence-backed critique of citation behavior. This research serves as a proof-of-concept specifically within the high-entropy domain of **Engineering Design**, where the mapping between abstract theory and technical implementation is notoriously difficult to verify at scale.

7.2 Summary of Key Findings

The empirical evaluation conducted across a specialized cluster of seven top-tier design journals yielded several critical insights into the feasibility of AI-driven research integrity:

1. **Feasibility of Zero-Shot Extraction:** The implementation of `reference_extractor.py` demonstrated that Large Language Models (Gemini 3 Flash) can parse unstructured bibliographies with **96% accuracy**. By treating the bibliography as a structural rather than a lexical problem, the system successfully navigated the diverse formatting styles of journals ranging from *JMD* to *CoDesign* without the need for brittle, regex-based rules.
2. **Efficacy of Legitimate Retrieval:** The "Hunter" module achieved a robust **80% automated success rate** utilizing a strictly legal "White-Hat" strategy (Crossref → arXiv → CORE). This finding is significant as it proves that a substantial majority of design engineering literature is accessible through legitimate Open Science channels, reducing the need for controversial or legally ambiguous retrieval methods.
3. **Semantic Verification Power:** Perhaps the most critical finding was the system's **100% detection rate** for intentional citation distortions. By holding both Article A and Article B in a single 1M-token context window, the AI successfully bridged the "semantic gap," identifying instances where authors misattributed specific geometric methods or over-generalized the scope of earlier methodological studies.
4. **Validation of AI-as-Auditor:** The substantial agreement level ($\kappa = 0.82$) achieved in the hybrid evaluation confirms that modern LLMs can serve as scalable auditors. This suggests that when constrained by the **Bornmann and Daniel (2008)** taxonomy and a

"Skeptical Peer Reviewer" persona, the system approximates the reasoning of a human subject-matter expert with remarkable fidelity.

7.3 Contributions of the Study

This research offers several distinct contributions to both **Scientometrics** and **Applied Artificial Intelligence**:

- **Technical Contribution:** The development of a novel, asynchronous multi-agent Python architecture that integrates metadata resolution with long-context reasoning. The "Hunter" module provides a reproducible blueprint for solving the "Cold Start" problem of citation analysis the acquisition of full-text evidence from fragmented scholarly repositories.
- **Methodological Contribution:** The proposal of a "Hybrid Evaluation Logic" that validates AI output against both human "Gold Standards" and high-parameter "Judge" models. This provides a rigorous framework for assessing the quality of qualitative AI reports, which are often difficult to quantify.
- **Theoretical Contribution:** The operationalization of the **Bornmann and Daniel (2008)** typology into a computational model. By transforming a sociological theory into a set of deterministic system prompts, the research bridges the gap between bibliometric theory and automated practice.
- **Practical Contribution:** The creation of a functional **Researcher's Dashboard** (Gradio-based UI) that significantly reduces the cognitive load of peer review. For the Engineering Design community, this tool provides a critical layer of defense against the propagation of scientific misinformation and "citation drift."

7.4 Limitations and Challenges

Despite the overall success of the artifact, several systemic and technical limitations were identified:

- **The 20% Retrieval Gap:** The remaining failure rate in document acquisition is a direct result of proprietary paywalls and the lack of universal "Green Open Access." Technology alone cannot resolve these socio-legal barriers.
- **Structural Layout Sensitivities:** While the LLM is robust, extremely complex multi-column layouts with overlapping OCR artifacts can occasionally cause errors in reference segregation.
- **Version Discrepancies:** The use of preprints (from arXiv or CORE) introduces a risk of "Version Drift," where the analyzed manuscript differs slightly from the final Version of Record cited in Article A.

7.5 Future Work and Suggestions for Project Continuation

Building upon the foundations established in this thesis, several promising avenues for future research and project continuation are proposed:

7.5.1 Institutional and Technical Integration

To resolve the remaining retrieval failures, future iterations should focus on **Institutional Proxy Integration**. By connecting the "Hunter" module to university library authentication systems (e.g., EZproxy or Shibboleth), the software could legally access proprietary content that is currently behind paywalls. Additionally, integrating the system as a **Plugin for Reference Managers** (e.g., Zotero or Mendeley⁹) would allow researchers to verify their bibliographies in real-time as they write, preventing errors before they are published.

7.5.2 Longitudinal Citation Chain Analysis

The current system focuses on 1-to-1 citation pairs ($A \rightarrow B$). A significant extension would be the analysis of **Citation Paths** ($A \rightarrow B \rightarrow C \rightarrow D$). This would allow for a "genealogy of distortion" analysis, tracing how a specific technical finding or dataset parameter might be incrementally misrepresented as it propagates through generations of papers. This "lineage audit" would be invaluable for identifying the root causes of scientific myths in the design community.

7.5.3 Collaborative Auditing and Crowdsourcing

The project could evolve into a **Collaborative Research Integrity Platform**. By allowing researchers to upload their verified reports to a central repository (a "Wiki of Verified Citations"), the community could build a shared database of "Evidentiary Anchors." This would allow the system to benefit from previous human audits, creating a virtuous cycle of human-AI collaboration that increases the precision of the alignment agents.

7.5.4 Monitoring for Algorithmic Bias and Citation Equity

As AI-driven analysis becomes more prevalent, it is crucial to monitor for **Algorithmic Bias**. Future studies should investigate whether LLMs exhibit bias toward specific journals, authors, or geographic regions during the "Utility" and "Faithfulness" grading process. Ensuring that the "Skeptical Auditor" persona remains objective across different disciplinary rhetoric styles (e.g., STEM vs. Social Design) is essential for maintaining equity in research evaluation.

7.5.5 Disciplinary Transferability and Fine-Tuning

While this study focused on Engineering Design, the methodology is designed to be field-agnostic. Future work should test the pipeline on other "high-stakes" fields such as **Medicine, Law, or Climate Science**. Fine-tuning the "Worker" models on domain-specific corpora (e.g., PubMed or

⁹ API documentation for Zotero suggests that a 'Verification Plugin' could be implemented via a JavaScript-based extension, calling the back-end agents developed in this research.

legal transcripts) would likely increase the nuance of the "Reasoning Soundness" rubric even further.

7.6 Final Concluding Statement

As the velocity of scientific publishing continues to accelerate, the traditional human-led mechanisms of research oversight are being pushed to their limits. This thesis has demonstrated that Large Language Models, when grounded in strict structural evidence and robust, legal retrieval logic, can act as powerful "**Augmented Intelligence**" for researchers.

By shifting the focus of scientometrics from citation **quantity** to citation **quality**, we move closer to a future where scientific impact is measured not just by how often a paper is cited, but by how faithfully and meaningfully it contributes to the global body of knowledge. Within the Engineering Design domain, this software artifact provides a scalable, transparent, and objective solution to ensure that the bridge between theory and implementation remains untainted by misrepresentation. The "Verification Gap" is no longer a technical impossibility; it is a soluble challenge that AI is now uniquely equipped to address.

References

1. Academic Literature

- **Bornmann, L., & Daniel, H. D. (2008).** What do citation counts measure? A review of studies on citing behavior. *Journal of Documentation*, 64(1), 45-80. <https://doi.org/10.1108/00220410810844150>
- **Zeng, Y. (2004).** Environment-based design. *Proceedings of the 12th International Conference on Integrated Design and Process Technology (IDPT)*, Izmir, Turkey.
- **Zeng, Y. (2015).** Environment-Based Design (EBD): A Methodology for Transdisciplinary Design. *Journal of Integrated Design and Process Science*, 19(1), 5–24. <https://doi.org/10.3233/jid-2015-0004>
- **Zeng, Y., & Cheng, G. D. (1991).** On the logic of design. *Design Theory and Methodology*, DE-Vol. 31, 33-40.
- **Kunnath, S., Pride, D., Herrmannova, D., and Knoth, P. (2021)** A Meta-analysis of Semantic Classification of Citations. Quantitative Science Studies. Advance Publication. https://doi.org/10.1162/qss_a_00159
- **Lasse M. Jantsch, Dong-Jae Koh, Seonghwan Yoon, Jisu Lee, Anne Lauscher, and Young-Kyoon Suh. 2025.** [FineCite: A Novel Approach For Fine-Grained Citation Context Analysis](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 24525–24542, Vienna, Austria. Association for Computational Linguistics.
- **Suchetha N. Kunnath, Drahomira Herrmannova, David Pride, Petr Knoth;** A meta-analysis of semantic classification of citations. *Quantitative Science Studies* 2022; 2 (4): 1170–1215. doi: https://doi.org/10.1162/qss_a_00159
- **Bornmann, L., & Leibel, C. (2025).** Citation accuracy, citation noise, and citation bias: A foundation of citation analysis. <https://arxiv.org/abs/2508.12735>

2. Artificial Intelligence and Large Language Models

- **Google. (2025).** *Gemini 3 Flash Model Documentation*. Google AI for Developers. <https://ai.google.dev/gemini-api/docs/models/gemini>
- **Google. (2025).** *Gemini 3 Pro: A High-Capacity Multimodal Model*. Technical Report.
- **Vaswani, A., et al. (2017).** Attention is all you need. *Advances in Neural Information Processing Systems*, 30. [Foundational paper for the Transformer architecture used by Gemini].

3. Software, Libraries, and APIs

- **Crossref. (2024).** *Crossref REST API Documentation*. <https://api.crossref.org/swagger-ui/index.html>
- **CORE. (2024).** *CORE API: The world's largest collection of open access research papers*. <https://core.ac.uk/services/api>
- **Gradio. (2024).** *Gradio: Build and share delightful machine learning apps*. <https://www.gradio.app/docs>
- **Scite. (2026).** AI for Research, <https://scite.ai>
- **Scite AI Review: Features, Pricing and Alternatives - Paperpal**, accessed January 16, 2026, <https://paperpal.com/blog/news-updates/scite-ai-review-pricing-alternatives>
- **Python Software Foundation. (2024).** *Python Language Reference, version 3.10*. Available at <https://www.python.org>
- **Van Rossum, G., & Drake, F. L. (2009).** *Python 3 Reference Manual*. CreateSpace.

- **Library - google-generativeai:** Google's Python SDK for the Gemini API.
- **Library - asyncio:** Python Standard Library for asynchronous I/O and concurrency.
- **Library - aiohttp:** Asynchronous HTTP Client/Server for Python.
- **Library - difflib:** Python Standard Library for computing deltas and sequence similarity (used for Title Similarity Check).
- **Library - PyMuPDF / fitz:** (If used for PDF byte-stream extraction) Documentation: <https://pymupdf.readthedocs.io/>

4. Academic Data Sources (The Evaluation Cluster)

- **arXiv.org.** (Cornell University). Open-access repository of electronic preprints.
- **ASME Digital Collection.** (For *Journal of Mechanical Design*).
- **Cambridge University Press.** (For *AIEDAM* and *Design Science*).
- **Taylor & Francis.** (For *CoDesign* and *Journal of Engineering Design*).