

**Ministry of Education and Science of Ukraine
Taras Shevchenko National University of Kyiv
Educational and Scientific Institute of Philology
Department of English Philology and Intercultural Communication**

Master's thesis

Expert Artificial Intelligence Terminology Landscape (based on English)

Oleh Hebura

2nd year student of the Education Program
“English Communication Studies and Translation
and Two Western European Languages”
Field of science: 03 “Humanities”
Specialty: 035 “Philology”
Supervised by:
PhD., **Skybytska Nadiia**

«Допущено до захисту»

Протокол засідання кафедри англійської філології

та міжкультурної комунікації

Протокол № 9 від 28.04.2025

Зав. кафедри _____ д. філол. н., проф. Алла БЄЛОВА

ABSTRACT

This research investigates the contemporary terminology landscape of Artificial Intelligence (AI), focusing on the specialized English vocabulary employed by experts. Given AI's dynamic nature and profound impact, understanding the characteristics of its expert lexicon is vital for clear communication, effective knowledge transfer, and navigating the complexities of the domain. The study aims to identify, characterize, and analyze the key linguistic features of modern expert AI terminology using an empirical, corpus-based methodology.

A specialized, synchronic corpus comprising 2,865 titles and abstracts (approx. 544,000 words) from the 2024 AAAI Conference on Artificial Intelligence proceedings was compiled for this purpose. Analysis was conducted employing corpus linguistic tools, primarily within the Sketch Engine environment, alongside initial processing by the Gemini API. The specialized corpus was compared against a large general English reference corpus (English Trends 2014-today) to assess domain specificity.

The findings highlight a terminology heavily concentrated on core methodological themes such as model architectures, learning paradigms, optimization techniques, and data engineering. Structurally, the lexicon is dominated by nominal forms, exhibiting a high prevalence of multi-word terms (MWTs), particularly following N+N and Adj+N patterns, and extensive use of acronyms (e.g., LLM, GNN, RL). Keyword analysis confirmed a high degree of domain specificity, resulting from both technical neologisms unique to AI and the significant semantic specialization of common English words (e.g., *model*, *attention*, *learning*, *training*, *bias*, *hallucination*). Semantic specialization narrows general meanings to precise computational or algorithmic concepts, while metaphorical extension (e.g., mapping cognitive or psychological concepts like learning, attention, or hallucination onto computational processes) serves as a crucial mechanism for term creation and conceptualization.

Key words: artificial intelligence, AI terminology, expert discourse, corpus linguistics, semantic specialization, terminologization, determinologization, collocation.

АНОТАЦІЯ

Це дослідження присвячене аналізу сучасних термінів сфери штучного інтелекту (далі — ШІ) з фокусом на фаховій лексиці, якою послуговуються експерти цієї галузі. Зважаючи на значний вплив ШІ, розуміння особливостей експертного лексикону є критично важливим для чіткої комунікації, ефективної передачі знань та орієнтування у тонкощах сфери. Метою роботи є виявлення, характеристика та аналіз ключових лінгвістичних рис сучасної фахової термінології ШІ за допомогою емпіричної, корпусно-орієнтованої методології.

Для цього було скомпільовано спеціалізований синхронний корпус, який містить 2865 назв та анотацій (близько 544 000 слів) з матеріалів конференції 2024 AAAI Conference on Artificial Intelligence. Аналіз проводився з використанням інструментів корпусної лінгвістики, переважно в середовищі Sketch Engine, а також із залученням Gemini API для первинної обробки. Спеціалізований корпус порівнювався з великим референсним корпусом загальноживаної англійської мови (English Trends 2014-today) для ретельної оцінки доменної специфічності.

Результати дослідження вказують на термінологію, сконцентровану навколо ключових методологічних тем, як-от архітектури моделей ШІ, парадигм навчання, технік оптимізації та інженерії даних. Структурно у лексиконі переважає іменне словотворення, зокрема багатокomпонентні терміни, утворені за моделями N+N та Adj+N, а також акроніми (напр., LLM, GNN, RL). Аналіз ключових слів підтвердив високий ступінь доменної специфічності унаслідок появи технічних неологізмів, унікальних для ШІ, та суттєвої семантичної спеціалізації загальноживаних англійських слів (напр., *model*, *attention*, *learning*, *training*, *bias*, *hallucination*). Семантична спеціалізація звужує загальні значення до точних обчислювальних або алгоритмічних понять, тоді як метафоричне розширення слугує важливим механізмом термінотворення та концептуалізації.

Ключові слова: штучний інтелект, термінологія ШІ, експертний дискурс, корпусна лінгвістика, семантична спеціалізація, термінологізація, детермінологізація, колокація.

TABLE OF CONTENTS

INTRODUCTION	8
Chapter 1: THEORETICAL FOUNDATIONS FOR ANALYSING EXPERT ARTIFICIAL INTELLIGENCE TERMINOLOGY	16
1.1 Terminology Studies: Concepts, Cognition, and Communication	16
1.1.1 Core Concepts and Theoretical Perspectives (Traditional vs. Modern)	16
1.1.2 Cognitive Dimensions of Terminology	18
1.1.3 Communicative Dimensions of Terminology	20
1.2 Corpus Linguistics as a Methodology in Terminology Research.....	22
1.2.1 Suitability of Corpus Linguistics for LSP and Terminology.....	22
1.2.2 Principles of Corpus Design for LSP/Terminology	24
1.2.3 Corpus Analysis Tools and Techniques for Terminology.....	26
1.3 AI Discourse and Terminology	29
1.3.1 Evolution of AI Concepts and Terminology	29
1.3.2 Previous Linguistic/Terminological Studies on AI Vocabulary.....	31
1.4 Defining the Framework for Analysis	34
1.4.1 Foundations and Context	34
1.4.2 Justification of the Corpus-Based Methodology	35
CONCLUSION TO CHAPTER 1	37
CHAPTER 2. METHODOLOGY AND CORPUS FOR THE ANALYSIS OF MODERN EXPERT AI TERMINOLOGY	39
2.1. Corpus Compilation and Characteristics	39
2.1.1. Source Selection Rationale (AAAI-24).....	39
2.1.2. Data Collection and Corpus Size.....	40
2.1.3. Initial Keyword Extraction and Frequency Ranking	40
2.2. Analytical Tools and Reference Data	41
2.2.1. Sketch Engine Environment	41
2.2.2. English Trends Corpus as Reference	42
2.3. Research Procedures.....	42
2.3.1. Phase 1: Thematic Clustering and Structural Analysis Procedure	42

2.3.2. Phase 2: Domain Specificity Analysis Procedure (Keyword Comparison)	43
2.3.3. Phase 3: Semantic and Contextual Analysis Procedure (Word Sketch, Concordance).....	43
CONCLUSION TO CHAPTER 2	44
CHAPTER 3. CHARACTERISTICS OF MODERN EXPERT ARTIFICIAL INTELLIGENCE TERMINOLOGY: A CORPUS-BASED ANALYSIS	46
3.1. Thematic Structure and Linguistic Form of Expert AI Terminology (Phase 1 Results)	46
3.1.1. Dominant Conceptual Categories in the Expert Corpus.....	46
3.1.2. Structural Features: Acronyms and Multi-Word Term Patterns	52
3.2.1. Statistical Keyword Analysis: AI Expert Corpus vs. General English.....	55
3.2.2. Interpretation of Domain-Specific Vocabulary	60
3.3. Semantic Behaviour and Contextual Usage of Key AI Terms (Phase 3 Results)	63
3.3.1.1. Analysis of: model (lemma)	63
3.3.1.2. Collocational Profiles and Grammatical Patterns.....	63
3.3.1.3. Semantic Specialization and Meaning Shifts	65
3.3.1.4. Metaphorical Patterns in Expert AI Terminology	66
3.3.2.1. Analysis of: learning (lemma)	67
3.3.2.2. Collocational Profiles and Grammatical Patterns.....	67
3.3.3.3. Semantic Specialization and Meaning Shifts	69
3.3.3.4. Metaphorical Patterns in Expert AI Terminology	70
3.3.4.1 Analysis of: attention (lemma)	71
3.3.4.2. Collocational Profiles and Grammatical Patterns.....	71
3.3.4.3. Semantic Specialization and Meaning Shifts	73
3.3.4.4. Metaphorical Patterns in Expert AI Terminology	74
3.3.5.1. Analysis of: hallucination (lemma)	75
3.3.5.2. Collocational Profiles and Grammatical Patterns.....	75
3.3.5.3. Semantic Specialization and Meaning Shifts	76
3.3.5.4. Metaphorical Patterns in Expert AI Terminology	77
3.3.6.1. Analysis of: training (lemma).....	78

3.3.6.2. Collocational Profiles and Grammatical Patterns.....	79
3.3.6.3. Semantic Specialization and Meaning Shifts	80
3.3.6.4. Metaphorical Patterns in Expert AI Terminology	81
3.3.7.1. Analysis of: bias (lemma).....	82
3.3.7.2. Collocational Profiles and Grammatical Patterns.....	83
3.3.7.3. Semantic Specialization and Meaning Shifts	84
3.3.7.4. Metaphorical Patterns in Expert AI Terminology	85
CONCLUSION TO CHAPTER 3	86
CONCLUSION	89
List of References.....	93
Appendix.....	105
Summary	151

INTRODUCTION

Relevance of the Study

Artificial Intelligence (AI) is one of the most rapidly developing and transformative fields of the 21st century, profoundly impacting diverse spheres from scientific research and industry to daily life. “Typically, ML algorithms operate by learning models from existing data and generalizing them to unseen data” [88] and are involved in many activities: From fighting the climate crisis, tackling the problems of ageing societies, reducing global poverty, stopping terror, detecting copyright infringements or curing cancer to improving evidence-based politics, improving predictive police work, local transportation, self-driving cars and even waste removal [82]. This active development of AI is linked to the evolution of its specialized language, a complex and constantly expanding terminology. Effective communication, knowledge dissemination, and conceptual clarity within any advanced scientific domain hinge upon a shared and well-understood specialized vocabulary. Indeed, the ability to precisely articulate ideas, replicate findings, and foster collaboration is heavily dependant on the clarity and consistency of the terms employed.

The field of AI, however, presents unique terminological challenges due to its exceptionally rapid pace of advancement. Nowadays, “generative artificial intelligence (AI) and large language models (LLMs) are driving a paradigm shift in terminology work and terminology science” [89]. As new theories emerge and technologies, such as machine learning, deep learning, neural networks, and, more recently, large language models (LLMs), achieve breakthroughs, the terminology used by experts adapts and shifts almost continuously. This flux manifests in the coining of neologisms, the proliferation of acronyms (e.g., *LLMs* for *large language models*), and the assignment of specialized meanings to existing words (e.g., *attention*, *transformer*, *hallucination*).

While tech-speak helps AI specialists communicate complex ideas accurately and efficiently with each other, the rapid growth of this jargon can leave others feeling lost.

When scientists from different backgrounds try to collaborate, this specialized language can get in the way, making it harder for them to connect and share insights. Previous research already suggests that “using jargon significantly disrupts processing fluency, in addition to and separate of comprehension” [38]. The sheer speed of AI development likely makes this communication gap even wider. What this means is, the hurdle to understanding the expert conversation keeps getting higher. This blocks not only public comprehension but also the crucial collaboration needed between different disciplines.

Although previous linguistic research may have examined AI terminology from historical perspectives or based on broader sources like dictionaries, there is a pressing need for a contemporary analysis focused specifically on the *expert discourse* shaping the field today. General analyses or lexicographical resources often lag behind the cutting edge of a field evolving as rapidly as AI. Understanding the characteristics of the terms currently employed by leading researchers—their structure, frequency, semantic behaviour in context, and the conceptual landscape they map—is crucial. This study addresses this gap by undertaking a systematic, corpus-based investigation of the terminology used in recent, high-impact AI research literature, specifically the proceedings of the 2024 AAAI (The Association for the Advancement of Artificial Intelligence) conference [34]. Such an analysis, grounded in empirical evidence drawn directly from expert communication, provides objective insights into the linguistic features defining the forefront of AI discourse and offers value to linguists, terminologists, translators, educators, and AI practitioners navigating this complex domain.

Aim and Objectives

The **aim** of this thesis is to identify, characterize, and analyze the key features of modern expert Artificial Intelligence terminology using a corpus-based methodology.

To achieve this aim, the following **objectives** are set:

1. To compile a specialized corpus representing contemporary expert AI discourse (based on AAAI 2024 proceedings titles and abstracts).
2. To extract and quantify the frequency of key AI-related terms from the specialized corpus, identifying the most salient lexical items.
3. To perform thematic clustering and structural analysis of the identified high-frequency terms, examining patterns in word formation (such as compounding and affixation), the prevalence of multi-word units, and the use of acronyms, reflecting processes of neologism and lexical adaptation.
4. To determine the domain specificity of the expert AI terminology by comparing term frequencies in the specialized corpus against a large reference corpus of general contemporary English (Sketch Engine's English Trends), thereby distinguishing truly specialized vocabulary from general language use.
5. To conduct a semantic and contextual analysis of selected key AI terms using corpus linguistic tools (Sketch Engine), focusing on collocation patterns, grammatical behaviour, semantic specialization or shifts in meaning (investigating phenomena like semantic drift or broadening), and metaphorical usage compared to general English, utilizing functionalities like Word Sketch and concordance analysis.

Object and Subject of Research

The **object** of this research is the terminology employed in modern expert Artificial Intelligence discourse, as manifested in recent academic publications, specifically the titles and abstracts from the Proceedings of the AAAI Conference on Artificial Intelligence, Volume 38 (2024) — 2865 articles.

The **subject** of this research encompasses the lexico-semantic, structural, frequency, and conceptual characteristics of this specialized terminology, identified and analyzed through corpus linguistic methods. This includes investigating patterns of term formation such as compounding, affixation, and acronymization; assessing the domain specificity of terms relative to general English usage; analyzing collocational behaviour and grammatical patterns; exploring semantic nuances, including

specialization and meaning shifts often observed in technical language; and examining the underlying conceptual structures reflected in the language, potentially through the analysis of metaphorical patterns.

Theoretical and Methodological Basis

This study is grounded in the theoretical frameworks of **Terminology Studies** and **Corpus Linguistics**. Terminology studies provide the foundation for understanding the nature, function, and formation of specialized vocabulary within specific domains, viewing terms not merely as labels but as units intrinsically linked to conceptual systems and knowledge structures, and have become “a very much formulized field of research” [61]. Perspectives such as the Communicative Theory of Terminology (CTT), which conceptualizes terminological units as entities with cognitive, linguistic, and communicative dimensions, offer a valuable lens for analyzing how AI terms represent knowledge and function within expert communication. Corpus linguistics, in turn, provides the methodological framework for the empirical study of language use based on extensive, principled collections of authentic texts (corpora). It’s important to mention that “the texts in such a collection are often (but not always) annotated in order to enhance their potential for linguistic analysis” [87]. Its suitability for analyzing Language for Specific Purposes (LSP) and terminology is quite established with many techniques for data-driven linguistic description.

The methodology employed is primarily **corpus-based**, combining quantitative and qualitative approaches. Quantitative methods include frequency analysis to identify key terms (such as *machine learning*, *reinforcement learning*, and *LLMs*) and statistical comparisons between the specialized AI corpus and a general English reference corpus. This comparison, utilizing functions like Sketch Engine’s keyword analysis, serves to establish the domain specificity of the identified terms. Qualitative methods involve the thematic clustering of terms to map the conceptual landscape of expert discourse, structural analysis of term formation (including multi-word units and acronyms, which reflect processes such as compounding and affixation), and detailed contextual and semantic analysis. This latter analysis utilizes corpus tools such as concordancers (for

examining terms in context), and Word Sketches (for summarizing grammatical and collocational behaviour) within the Sketch Engine platform [68]. These tools facilitate the investigation of collocational profiles, grammatical patterns, semantic specialization, meaning shifts, and metaphorical usage.

The initial term extraction and categorization from the corpus were facilitated using the Gemini 2.0 Flash API. This methodology thus adopts a workflow that utilizes AI tools for the efficient initial identification of potential terms across the large dataset, a task that can be time-consuming when performed manually. Then, this initial and broad identification was followed by closer in-depth analysis using established corpus linguistic techniques and tools, such as Sketch Engine, to ensure empirical validity, explore nuanced linguistic features, and provide the necessary contextual and semantic depth. Integration of AI/NLP techniques within a primarily corpus linguistic framework stands in line with contemporary approaches in handling and analyzing specialized linguistic data for terminological research, acknowledging the potential of such tools while maintaining analytical thoroughness through established methods. The outputs from AI tools often require careful validation and refinement, which corpus linguistic methods are well-suited to provide.

Research Material (Corpus Description)

The primary empirical basis for this research is a specialized, synchronic corpus compiled specifically for this study. It consists of the titles and abstracts of 2,865 articles published in the *Proceedings of the AAAI Conference on Artificial Intelligence, Volume 38 (2024)*. This corpus comprises approximately 544,224 words (around 4,000,000 characters). The choice of AAAI proceedings as the source ensures that the corpus represents a snapshot of contemporary *expert* discourse in the AI field, capturing the language used in high-impact, peer-reviewed research at the forefront of the discipline.

From this primary corpus, a keyword list was generated. Initially, AI-related terms were extracted using the Gemini 2.0 Flash API, leveraging its capabilities for

processing large text volumes. As noted by Incelli, “AI’s integration into corpus linguistics is not an entirely new phenomenon” [55], and using LLMs for analysis displays recent trends in Linguistic Research. Subsequently, a frequency analysis was performed on these extracted terms, resulting in a list of approximately 550 unique terms (single and multi-word) that appeared at least 5 times in the corpus. This frequency-ranked list, including prominent terms such as *machine learning*, *large language models*, *neural network*, *reinforcement learning*, *transformer*, and *LLMs*, serves as a key dataset for the structural, thematic, and semantic analyses conducted in Chapter 3.

For comparative analysis, the **English Trends (2014–today)** [69] corpus available via the Sketch Engine platform will be used as a reference corpus. Using this specific reference corpus really strengthens the methodological framework. Given the rapid evolution of AI language and its tendency to repurpose general English words with specialized meanings (e.g., *model*, *network*, *agent*, *attention*, *bias*, *hallucination*), comparing the expert corpus against an up-to-date baseline of general contemporary English allows for a more accurate assessment of domain specificity and semantic divergence than comparison against older or static general language resources would permit.

Novelty and Practical Significance

The **scientific novelty** of this thesis lies in its focused, empirical investigation of contemporary terminology used by experts in the rapidly advancing field of AI. Unlike broader or earlier studies, this research utilizes a large, specialized corpus drawn from highly current (2024) primary sources (AAAI proceedings), ensuring relevance to the present state of the field. It employs a systematic corpus linguistic methodology to analyze not only structural patterns and frequency but also crucial aspects like domain specificity (through comparison with general English) and semantic behaviour in context (exploring collocations, meaning shifts, and grammatical patterns using tools like Word Sketch). Moreover, the use of AI (Gemini API) for initial term extraction and categorization represents a modern approach to handling and analyzing specialized

linguistic data for terminological research, demonstrating how newer computational tools can be incorporated into established linguistic practices.

The **practical significance** of the findings extends across several areas:

- **Contribution to Terminology Studies:** Providing an empirically grounded characterization of terminology dynamics, including neologism and semantic shifts, in a high-impact, fast-paced scientific field.
- **Lexicography:** Offering data that could inform the updating or creation of specialized AI dictionaries and glossaries, ensuring they accurately reflect current expert usage rather than lagging behind.
- **Understanding AI Research:** Revealing dominant themes, key concepts, and emerging trends within current AI research through the objective lens of its specialized language.
- **Applied Linguistics and Communication:** Potentially informing the development of curricula for English for Specific Purposes (ESP) tailored to AI professionals, refining technical writing practices within the field, and improving the accuracy and consistency of translating AI-related texts.

Thesis Structure

This Master's thesis consists of an Introduction, three chapters, Conclusion, References, and Appendices.

The **Introduction** outlines the relevance, aim, objectives, object, subject, theoretical and methodological basis, research material, novelty, practical significance, and structure of the thesis.

Chapter 1, "Theoretical Foundations for Analysing Expert Artificial Intelligence Terminology," reviews key concepts in terminology studies (including traditional and modern perspectives, cognitive and communicative aspects), discusses corpus linguistics as a methodology for specialized language research (covering principles, tools like Sketch Engine), surveys the evolution of AI discourse and its key concepts

alongside previous linguistic research on technical terminology, and defines the specific analytical framework adopted for this study.

Chapter 2, “Methodology and Corpus for the Analysis of Modern Expert AI Terminology,” details the corpus compilation process (source selection rationale, data collection, initial keyword extraction via API and frequency ranking), describes the analytical tools (Sketch Engine environment) and reference data (English Trends corpus), and outlines the specific procedures for the three phases of analysis: Phase 1 (thematic clustering and structural analysis), Phase 2 (domain specificity analysis via keyword comparison), and Phase 3 (semantic and contextual analysis using Word Sketch, and Concordance features).

Chapter 3, “Characteristics of Modern Expert Artificial Intelligence Terminology: A Corpus-Based Analysis,” presents the results of the empirical analysis conducted according to the methodology described in Chapter 2. It discusses the thematic structure and linguistic forms (dominant conceptual categories, patterns in acronyms and multi-word terms) identified in Phase 1, presents the findings on domain-specific terms derived from the statistical keyword comparison against general English in Phase 2, and details the analysis of semantic behaviour and contextual usage (collocational profiles, grammatical patterns, semantic specialization, meaning shifts, metaphorical patterns) for key AI terms investigated in Phase 3.

The **Conclusion** summarizes the main findings of the research, highlights its contribution to terminology studies and the understanding of AI communication, and acknowledges the limitations of the study.

The **References** provide full bibliographic details for all academic sources cited throughout the thesis.

Chapter 1: THEORETICAL FOUNDATIONS FOR ANALYSING EXPERT ARTIFICIAL INTELLIGENCE TERMINOLOGY

This chapter lays the groundwork for the thesis. It introduces some core theoretical concepts from Terminology Studies and Corpus Linguistics, contextualizing the study within the specific domain of Artificial Intelligence (AI) terminology, and defining the analytical framework employed. It establishes the theoretical underpinnings necessary for understanding the nature of specialized vocabulary and justifies the methodological choices made to investigate the characteristics of terms used by experts in the rapidly evolving field of AI.

1.1 Terminology Studies: Concepts, Cognition, and Communication

This section introduces the field of Terminology Studies, its core concerns and contrasting traditional and modern theoretical perspectives. It further explores the cognitive base of terms as knowledge access points and their crucial role in specialized communication.

1.1.1 Core Concepts and Theoretical Perspectives (Traditional vs. Modern)

Terminology Studies, or simply Terminology, is the discipline concerned with the study, description, processing, and management of specialized vocabulary (terms) used within specific subject fields or domains of activity. It is an inherently interdisciplinary field, drawing insights and methods from linguistics, cognitive science, communication studies, documentation science, and computer science. Its central object of study is the “term” – “a word or combination of words that is clear, unambiguous, devoid of expressive features, and is clear in the field of its terminological field” [30] – and the relationship between these terms, the concepts they represent, and the real-world objects or phenomena they refer to.

Historically, the field was dominated by the **General Theory of Terminology (GTT)**, primarily developed by Eugen Wüster in the mid-20th century. Eugen Wüster is known as “the father of modern terminology” [90]. The GTT emerged from the practical needs

of science and technology, particularly the need for clear, unambiguous communication and international standardization. Consequently, the GTT adopted a fundamentally prescriptive approach. It emphasized the primacy of concepts, viewed as clearly defined, objective units of knowledge with universal validity within a given field. The ideal relationship between concept and term, according to GTT, was one-to-one (monosemy), ensuring clarity and avoiding ambiguity. The primary goal of terminological work under the GTT framework was to establish standardized, stable, and internationally harmonized terminologies, often codified in standards documents or specialized dictionaries, thereby facilitating precise technical communication and knowledge transfer.

However, beginning in the latter part of the 20th century, and gaining momentum with the rise of computational linguistics and large-scale text analysis, alternative perspectives emerged, challenging the strict prescriptivism and idealized view of the GTT. **Modern perspectives**, exemplified by approaches such as the **Communicative Theory of Terminology (CTT)** developed by M. Teresa Cabré [39, 40, 41], shift the focus from idealized, standardized terms to the description and analysis of terms *in use* within authentic communicative contexts. CTT conceptualizes the “terminological unit” (TU) not simply as a label for a pre-defined concept, but as a complex entity with distinct but interconnected cognitive (concept), linguistic (term form), and communicative (situation/context) dimensions.

This shift entails a move towards a descriptive methodology. Instead of prescribing how terms *should* be used, modern approaches seek to understand how they *are* actually used by expert communities. This involves acknowledging and analyzing linguistic phenomena often downplayed by GTT, such as variation (different terms for the same concept, or variations in term form), polysemy (a single term having multiple related meanings within the domain), synonymy (multiple terms sharing similar meanings), and the inherent dynamism of terminology, especially in rapidly evolving fields. CTT, for instance, views concepts not as static, objective entities but as culturally constructed and dynamically stabilized through consensus within expert

communities. The meaning and appropriateness of a term are seen as heavily dependent on the specific communicative situation, including the participants, purpose, and textual context.

1.1.2 Cognitive Dimensions of Terminology

Beyond their linguistic form, terms are fundamentally linked to cognition. They serve as crucial interfaces between language and knowledge. They function as linguistic labels for **concepts**, which are mental representations of classes of items, objects, processes, or ideas within a specialized knowledge domain. Terms are more than just words; they represent units of knowledge, encapsulating complex information relevant to a specific field, and cognitive linguistics aims to “model the worldview as well as to model the means of linguistic thinking” [32].

This role connects terminology directly to **knowledge representation**, the study of how knowledge is structured, organized, and stored, both mentally and in external systems. The set of terms used within a discipline, and the relationships between them (e.g., hierarchical, associative), reflect the underlying conceptual structure and organization of knowledge in that field. “Once the knowledge is acquired, it should be stored in a knowledge base for the later use in reasoning” [67]. Terms act as access points or pointers to these larger knowledge structures, schemas, or mental models. When experts encounter a familiar term, it triggers the corresponding idea in their head and potentially related knowledge too. This helps them reason effectively within their field. By looking at the terminology, we can get useful insights into how these experts think about their subject, structure their knowledge, and perceive their domain.

The relationship between terms and concepts is also intertwined with fundamental cognitive processes like **categorization**. Concepts allow humans to organize information efficiently by treating distinct entities as members of a class. “For example, apples and lemons can be classified as fruits; therefore, categorization facilitates knowledge organization” [93]. Research in cognitive psychology suggests concepts are often organized hierarchically (e.g., superordinate, basic, subordinate levels) and may

be structured around prototypes or characteristic features rather than strict definitions. While traditional terminology (GTT) aimed for sharply defined concepts, modern views acknowledge that conceptual boundaries can be fuzzy and context-dependent, aligning more closely with cognitive models like Prototype Theory. The terms used in a field reflect these categorization processes, labeling the conceptual distinctions deemed relevant by the expert community.

Thinking about how our minds process information, specialized terms actually have high **efficiency**. They let us condense complex ideas and a whole lot of background knowledge into just a few short words or phrases. This makes communication faster and more precise, and it also makes it easier for experts to mentally juggle information specific to their field.

This close link between terminology and how experts think means that looking at patterns in the language can give us clues about the ideas floating around in a field. When technology changes, or when people understand things differently or make new discoveries, new concepts often pop up or old ones get tweaked. These shifts in thinking naturally show up in the field's terminology – maybe through brand new words being coined, existing words changing their meaning, or terms being used more or less often, or in different ways. Corpus linguistic methods, which allow for the systematic tracking of term frequencies, collocational patterns (words typically used alongside a term), and contextual usage over time or across different contexts, provide a powerful, data-driven approach to observing these linguistic reflexes of cognitive and conceptual change. By analyzing the diachronic or synchronic variations in terminology within a representative corpus, such as the AAI proceedings corpus used in this thesis, it becomes possible to make empirically grounded inferences about shifts in the conceptual frameworks and knowledge organization within the AI expert community.

1.1.3 Communicative Dimensions of Terminology

While terms certainly relate to how we think, their main job really comes alive in communication. Think of terms as the essential tools experts need for effective **specialized communication** within their particular field or area of work. They act like a shared code, a common language, that lets these experts swap complex information precisely, clearly, and quickly. Having this shared understanding is absolutely vital for working together, pushing research forward, and putting specialized knowledge to practical use.

Terminology also plays a huge role in knowledge dissemination – which is really just the process of sharing scientific and technical discoveries. Within groups of experts, terms make it easier to communicate research results, whether that’s in published articles, conference talks, or team projects. For this sharing to work well, it really depends on everyone using the terminology consistently and understanding what it means. As noted by Bianchini et al., “[m]ultidisciplinary collaborations of researchers, clinicians, developers, etc. benefit from a common language and collective understanding of fundamental principles, concepts and techniques in order to strengthen communication among collaborators” [36]. Therefore, a unified language for knowledge exchange could improve the readability of scientific outputs and their impact on society. Beyond expert circles, specialized knowledge often needs to be communicated to broader audiences, including policymakers, practitioners in related fields, students, and the general public. In these contexts, the clarity and accessibility of terminology become even more critical, though often more challenging.

However, the very specialization that makes terminology efficient for insiders can create **potential barriers** for outsiders or even inside communities. Specialized vocabulary, often referred to pejoratively as **jargon**, can be opaque and inaccessible to those who do not share the specific knowledge background of the expert community. As Lucy et al. observe, “scholarly text is often laden with jargon, or specialized language that can facilitate efficient in-group communication within fields but hinder

understanding for out-groups. ... our findings suggest that though multidisciplinary venues intend to cater to more general audiences, some fields' writing norms may act as barriers rather than bridges, and thus impede the dispersion of scholarly ideas" [72]. This can impede communication not only between experts and non-experts but also between experts from different sub-disciplines within a larger field, hindering interdisciplinary understanding and collaboration. Quantitative research suggests that the overuse of jargon, particularly in crucial communicative elements like titles and abstracts of scientific papers, can negatively impact the reach and visibility of research [75], potentially leading to lower citation rates because fewer people can readily grasp the paper's core message.

This underscores the importance of **context and user needs** in terminological practice and analysis, a central tenet of modern perspectives like CTT. The meaning of a term is not fixed and absolute but is actualized and potentially modulated by the specific communicative situation – the text genre, the purpose of the communication, the intended audience, and the surrounding linguistic context. Effective terminological choices, whether in writing, translation, or teaching, require careful consideration of these contextual factors and the specific needs of the users or audience.

This creates a dynamic tension, a kind of "jargon paradox." On one hand, specialized terminology is indispensable for the precision and efficiency required for communication and knowledge advancement *within* an expert group. On the other hand, it can simultaneously function as a significant obstacle to knowledge dissemination and communication *beyond* that immediate group. This tension is particularly relevant in a field like Artificial Intelligence. AI research involves deep internal specialization, generating highly specific terminology within its various subfields (e.g., reinforcement learning, computer vision, natural language processing). At the same time, AI's impact on society is profound and growing quickly. That makes good communication crucial between many different groups: researchers, engineers, policymakers, ethicists, and the public. Therefore, studying expert AI terminology isn't just about understanding the experts themselves (which is our main focus here). It's

also about being aware that their specialized language can create barriers for others. Figuring out exactly how experts use their terms is the necessary starting point for potentially making AI easier for everyone to understand down the road.

1.2 Corpus Linguistics as a Methodology in Terminology Research

This section introduces Corpus Linguistics as a methodology, explains its suitability for studying specialized language and terminology, discusses key principles for designing relevant corpora, and describes essential analysis tools and techniques, particularly those available in Sketch Engine.

1.2.1 Suitability of Corpus Linguistics for LSP and Terminology

Corpus Linguistics (CL) is an approach to the study of language that relies on the analysis of large, principled collections of authentic electronic texts, known as corpora. Its methodology is fundamentally **empirical and data-driven**, drawing conclusions about language structure and use based on observable evidence from real-world communication rather than solely on introspection or prescriptive rules. CL “spurred the development of methods and applications of corpus linguistics and had a notable impact on shifting language analysis” [43].

The focus on **authentic usage** also makes CL particularly well-suited for the study of **Language for Specific Purposes (LSP)** – the language used in specialized domains like science, technology, law, or medicine – and, consequently, for terminology research. Instead of just assuming terms have perfect, dictionary-like meanings (like one word always meaning exactly one thing), Corpus Linguistics lets us dig into how people actually use these terms in their everyday work and communication. This focus on real-world evidence fits perfectly with modern ways of thinking about terminology, like CTT, which highlight how important things like context, slight variations in meaning, and the act of communicating really are.

CL methodologies facilitate both **quantitative and qualitative insights** into language use. Quantitative analysis involves measuring frequencies of words or structures,

calculating the statistical significance of co-occurrence patterns (collocations), and identifying trends. In linguistic research, “it is common to calculate the quantitative relations between parts of speech, considering them as one of the components of the statistical characteristics of the text” [84]. Qualitative analysis typically involves the close examination of concordance lines (examples of a word or phrase in its surrounding context) to understand nuances of meaning, semantic prosody, and specific usage patterns. “It attempts to accurately represent how words, phrases, and grammatical structures are used in various written and spoken discourse types” [78].

Applied specifically to **terminology**, CL provides a powerful toolkit for various tasks. Corpus analysis can help identify potential term candidates based on frequency and statistical measures (like “keywords” specific to a domain – “elements of a text [that] are more important than others in informing readers about the text’s characteristics” [64]). It enables the detailed analysis of a term’s behavior, including its typical collocates (words it frequently combines with, as they are “becoming prominent in our understanding of language learning and use” [48]), its syntactic patterns (the grammatical structures it participates in), and its semantic range as observed across numerous contextual examples. Furthermore, corpora allow researchers “to trace language change and model how languages evolve under various influences” [44] (diachronic analysis) or across different genres or subdomains (synchronic variation). Corpora have thus become indispensable resources for modern terminography (the practice of compiling terminological resources) and terminology research.

The suitability of CL for terminology research goes beyond mere practicality; it represents a methodological operationalization of modern theoretical perspectives. Theories like CTT posit that terms are dynamic linguistic units whose meaning and function are shaped by communicative context and usage within expert communities. These claims about context-dependency, variation, and dynamism require empirical investigation based on substantial amounts of authentic language data. Intuition or the analysis of isolated examples is insufficient to capture the complexities of real-world term usage. CL provides precisely the methods needed to undertake this empirical

investigation. Techniques such as frequency analysis, keyword extraction, collocation analysis, and concordance examination allow researchers to directly observe and measure patterns of usage, variation, co-occurrence, and contextual meaning in large text datasets. Therefore, CL offers the practical means to test, refine, and apply the principles of modern terminology theory, grounding theoretical constructs in observable linguistic evidence drawn from specialized discourse, such as the AI terminology under investigation in this thesis.

1.2.2 Principles of Corpus Design for LSP/Terminology

Whether the findings from research using text collections (corpora) are any good really depends on the collection itself – what’s in it and how well it was put together. Designing a corpus isn’t just about following technical steps; it’s a super important part of the research method that follows some key rules.

The biggest rule is making sure it’s **representative**. This means the collection needs to accurately show what the specific type of language you’re interested in actually looks like. It can be achieved “by balancing the corpus through sampling a wide range of text categories which are defined primarily in terms of external criteria” [76]. For LSP and terminology research, this means the corpus must be representative of the specialized discourse under investigation (e.g., modern expert AI research papers). Factors influencing representativeness include the selection of appropriate text types and genres (e.g., peer-reviewed articles, conference proceedings), the time period covered, the domain scope (e.g., specific subfields of AI), and the characteristics of the authors/speakers (e.g., experts in the field).

Closely related is the principle of **balance**. A well-balanced corpus ensures that different facets or sub-varieties within the target language domain are included in proportions that reflect their actual occurrence or importance. For instance, a corpus of AI research might aim for balance across different subfields (e.g., machine learning, robotics, NLP) or different publication venues, if these variations are relevant to the

research questions. Lack of balance can skew results, over-representing features of certain sub-varieties while under-representing others.

Corpus size is another important consideration. While the axiom “bigger is better” often holds true in CL, providing more data for statistical reliability, the optimal size depends on the research goals and the nature of the language variety. LSP corpora, dealing with potentially denser and more restricted vocabularies, can often yield valuable insights even at sizes considerably smaller than those required for general language research. Generally, “specialized corpora have smaller size than reference corpora representing language for general purpose” [92]. Studies have shown “LSP corpora ranging from tens of thousands to several hundreds of thousands of words to be useful” [66]. The key criterion is whether the corpus is large enough to contain sufficient instances of the phenomena being studied (e.g., occurrences of specific terms and their collocates) to allow for reliable analysis and generalization.

A fundamental distinction exists between **specialized corpora** and **general reference corpora**. Specialized corpora, like the AAI corpus intended for this thesis, contain texts from a specific domain, field, or genre. They are essential for in-depth analysis of LSP and terminology within that domain. General reference corpora (e.g., the British National Corpus (BNC) [37], the Corpus of Contemporary American English (COCA) [45], or large web corpora) aim to represent a language more broadly, encompassing a wide range of genres and text types. General corpora are crucial for comparative purposes in terminology research, particularly for identifying domain-specific keywords by contrasting frequencies in the specialized corpus with those in the general corpus. “In case of a general corpus, the size has to be as large as possible, in order for it to represent the variety of a language. However, in a specialized, domain-specific corpus it is easier to achieve representativeness of this particular language type based on a smaller-sized corpus” [65].

As noted by Incelli, “[c]orpus linguistics traditionally relies on software to analyse extensive language datasets, uncovering patterns and insights. However, the rise of generative AI offers opportunities to enhance linguistic research by automating tasks”

[56]. Building a collection of texts (a corpus) for research isn't just about gathering words; it's deeply connected to what one is trying to figure out and the ideas guiding the study. When deciding which texts to include, how to pick them, how many one needs, and whether they have a good variety, it's not just about making technical choices. These decisions actually show what one thinks defines the specific type of language they research. For example, using papers from a peer-reviewed conference like AAAI essentially means, 'This is what I consider expert AI talk for this time.' Designing a corpus inevitably shapes how we see the language and what we end up finding. That's why it's absolutely crucial to clearly explain the design choices, making sure they line up with the study's goals and core beliefs about the language (e.g., a CTT-informed focus on authentic expert communication).

1.2.3 Corpus Analysis Tools and Techniques for Terminology

Once a suitable corpus is compiled, various computational tools and analytical techniques are employed to extract meaningful information about terminology and language use.

Several **core techniques** are commonly applied in corpus-based terminology research:

- **Frequency Lists:** Generating lists of words or lemmas ordered by their frequency of occurrence in the corpus provides a basic profile of the vocabulary. While high frequency alone doesn't guarantee term status, it can highlight prominent lexical items for further investigation.
- **Keyword Analysis:** This technique statistically compares word frequencies in the specialized (target) corpus against frequencies in a larger general reference corpus. Therefore, it "provides a means for a quantitative linguistic analysis of textual content" [51]. Words that are statistically significantly more frequent in the specialized corpus are identified as "keywords", often corresponding to important domain-specific terms or concepts.
- **Concordance Lines (KWIC):** Concordancers retrieve all occurrences of a specific search word or phrase (the node) and display them with a span of

surrounding context (Key Word In Context), providing insights “into their usage patterns and semantic interpretations” [31]. Examining concordance lines allows for qualitative analysis of a term’s meaning-in-context, its typical usage patterns, grammatical behavior, and semantic prosody (the subtle attitudinal meanings it tends to co-occur with). As noted by Gillings and Mautner, “[c]oncordances are a staple tool in corpus-assisted discourse studies (CADS), as they are in corpus linguistics generally” [49]. Manual analysis of concordance lines is often crucial for disambiguating meaning and understanding function.

- **Collocation Extraction:** This involves identifying words that tend to co-occur with a target word (node) more often than expected by chance. Collocates often reveal significant semantic associations and typical lexical combinations (e.g., verb-object, adjective-noun pairs). Various statistical measures (e.g., Mutual Information, t-score, LogDice) are used to rank the strength or significance of collocations. Analyzing collocations is fundamental to understanding how terms are used and what concepts they relate to. It is an important process in many research branches and “is useful for language learning, dictionary compilation and downstream NLP” [85].
- **N-grams:** These are contiguous sequences of n items (typically words) extracted from the corpus: “a 2-gram ... is a two-word sequence of words like The water, or water of, and a 3-gram ... is a three-word sequence of words like The water of, or water of Walden” [59]. Analyzing frequent n-grams (bigrams, trigrams, etc.) is a common method for identifying potential multi-word terms or formulaic sequences.

While these core techniques can be performed using various software tools, **advanced platforms like Sketch Engine** offer integrated environments specifically designed for sophisticated linguistic analysis of large corpora. Sketch Engine incorporates several functionalities, particularly relevant for the aims of this thesis:

- **Word Sketch:** This is arguably Sketch Engine’s flagship feature. For a given word (lemma), it automatically generates a one-page summary of its significant

grammatical and collocational behavior based on predefined grammatical relations (e.g., “subject_of”, “object_of”, “modifier”, and “and/or”). It groups collocates according to these relations and ranks them by statistical significance (typically using the LogDice score). Word Sketches provide a powerful and efficient way to understand a term’s syntactic roles and semantic preferences within the corpus, moving beyond simple co-occurrence lists. “It can be used as a one-page summary of the word’s grammatical and collocational behavior” [86]. This requires the corpus to be appropriately processed (POS-tagged and lemmatized).

- **Thesaurus (Distributional):** Sketch Engine can automatically generate a distributional thesaurus, listing words that are semantically similar to a target word based on the principle that words occurring in similar contexts tend to have similar meanings. It identifies words that share significant collocates with the target word. “The thesaurus is using word sketches for computing the similarity score: it basically compares word sketch collocates for every pair of words in the corpus and the similarity relates to the fraction of shared collocates between these two words” [58]. This is useful for exploring semantic fields, identifying potential synonyms or related concepts within the AI domain.
- **Keywords:** The platform includes tools for performing keyword analysis, comparing the specialized corpus against various built-in general reference corpora or user-uploaded corpora to identify domain-specific vocabulary.
- **Concordance:** Sketch Engine offers a highly flexible concordance tool that supports complex queries using the Corpus Query Language (CQL), allowing searches based on lemmas, tags, syntactic structures, and combinations thereof.
- **Term Extraction:** While this thesis focuses on characterizing terminology rather than extracting it initially, Sketch Engine does possess automated term extraction capabilities that combine linguistic and statistical information.

The selection of analysis tools and techniques is not merely a matter of convenience; it actively shapes the research process and outcomes. The specific algorithms and functionalities offered by a platform like Sketch Engine, particularly the Word Sketch feature, directly enable certain types of fine-grained analysis of grammatical and

collocational patterns that would be extremely laborious or impossible to conduct manually on a large scale. The capabilities of the chosen tools determine which aspects of terminology can be effectively investigated and influence the nature of the insights derived. Therefore, the decision to utilize Sketch Engine and its specific functionalities (Word Sketch, Keywords, Concordance) is a deliberate methodological choice designed to facilitate the detailed, usage-based characterization of modern expert AI terminology envisioned in this thesis, consistent with the descriptive goals informed by modern terminology theory.

1.3 AI Discourse and Terminology

This section narrows the focus to the specific domain of Artificial Intelligence (AI), examining the interplay between the evolution of its core concepts and the specialized language used to express them. It also reviews prior linguistic or terminological research related to AI or similar technical vocabularies to situate the present study within the broader academic landscape.

1.3.1 Evolution of AI Concepts and Terminology

Over time, AI has seen big changes in approaches and focus. These changes always affected the specialized words used in the field. Tracking how key terms evolved helps us understand how AI's main ideas have grown and shifted.

The **early days of AI**, roughly from the 1950s through the 1970s, were dominated by **symbolic approaches**, often referred to as “Good Old-Fashioned AI” (GOFAI). This era focused on manipulating symbols according to logical rules to achieve intelligent behavior, based on the belief that human cognition could be modeled as symbol processing. Key concepts included logical reasoning, problem-solving, and knowledge representation through explicit rules and structures. The terminology characteristic of this period includes the term *artificial intelligence* itself, along with *symbolic computation*, *logical reasoning*, *heuristic search*, *knowledge representation* (often rule-based), and later, *expert systems* – programs designed to emulate the decision-making ability of a human expert in a narrow domain.

A significant paradigm shift occurred with the rise of **connectionism** and **machine learning**, gaining traction in the 1980s and 1990s, converting to “symbolic AI and the so-called “expert systems” or “knowledge-based systems”” [46] and experiencing explosive growth in the 21st century. Taking inspiration from how our own brains work, connectionist approaches used networks of linked “neurons” to learn patterns straight from data. This was a real change from needing to program specific rules by hand, shifting the focus towards **data-driven methods**. A new lexicon emerged, featuring terms like *machine learning*, *neural network*, *connectionism*, *pattern recognition*, *backpropagation* (a key algorithm for training networks), and *parallel distributed processing*. The adoption of biological metaphors became more prevalent during this period, as “AI makes heavy use of anthropomorphisms” [83].

The early 21st century witnessed the **Deep Learning Revolution**, fueled by the availability of massive datasets (“Big Data”), significant increases in computational power (especially GPUs), and algorithmic advancements. Deep learning involves neural networks with many layers (“deep” architectures), enabling the learning of complex hierarchies of features directly from data. This led to breakthroughs in areas like image recognition, speech recognition, and natural language processing. The terminology expanded rapidly to include *deep learning*, specific architectures like *convolutional neural networks (CNNs)* and *recurrent neural networks (RNNs)*, techniques like *dropout* and *batch normalization*, and approaches such as *reinforcement learning* and *generative adversarial networks (GANs)*.

The **current era** is characterized by the development and application of extremely large-scale models, particularly in natural language processing with “the integration of large language models (LLMs) with agent architectures” seen as “the most significant recent advancement in the field” [62]. The invention of the *transformer* architecture, based on *self-attention mechanisms*, proved highly effective and scalable. This led to the rise of *large language models (LLMs)* like BERT and GPT, trained on vast amounts of text data (*pre-training*) and then adapted for specific tasks (*fine-tuning*). Associated terminology includes *LLM*, *transformer*, *attention*, *pre-training*, *fine-tuning*, *prompt*

engineering (the art of crafting inputs for LLMs), *foundation models* (large models adaptable to many tasks), and *Generative AI (GenAI)* – AI focused on creating new content (text, images, code). Alongside these developments, concerns about transparency and interpretation have spurred research in *Explainable AI (XAI)*.

AI terminology is constantly changing because the field of AI itself changes so rapidly. New methods and inventions need new labels, or old labels get reused in new ways. This ongoing change in language basically creates a timeline, showing how AI ideas have developed. Key terms signal important concepts or breakthroughs from specific times. As AI moves forward, older words might fade, persist, or take on new meanings. Studying the language used today (like in the AAAI texts) gives us a picture of current thinking. Comparing today’s language with the past reveals how the field has grown and changed direction. The words themselves tell the story of AI’s progress and what’s considered important at different times.

1.3.2 Previous Linguistic/Terminological Studies on AI Vocabulary

Although artificial intelligence is a rapidly expanding domain with an ever-evolving specialized vocabulary, dedicated linguistic and terminological studies focusing specifically on the systematic, corpus-based characterization of its modern expert terminology remain relatively scarce compared to the vast number of technical AI publications. This ongoing scarcity is echoed in recent literature, which observes that “AI does not eliminate the corpus itself but reshapes how it is analyzed and interpreted” [57]. Such observations highlight a research gap that this thesis seeks to address.

Nevertheless, relevant background and methodological grounding can be found in several related areas:

- **Studies on Scientific/Technical Vocabulary:** Research on the general characteristics of scientific and technical language continues to inform our understanding of AI terminology. As Liu and Lei note, “the complexity involved in the definition of technical vocabulary” and the challenges of “its identification” remain central issues, especially given the “indispensability [of technical

vocabulary] in professional and vocational communication” [71]. These studies emphasize features such as the prevalence of noun phrases, specialized syntactic structures, and the evolving nature of technical definitions.

- **LSP Research:** The field of LSP offers theoretical frameworks and analytical tools for studying language use in specialized domains, including corpus-based approaches. Volkova highlights that “corpus-related linguistic analysis allows educators to study terms in their natural context to identify typical multi-word units, language patterns, and the frequency of their usage” [91]. This approach is increasingly adopted for both research and pedagogical purposes.
- **Computational Terminology / Automatic Term Extraction (ATE):** Recent advances in computational terminology, particularly in automatic term extraction, have been transformative for terminology research in AI and related fields. As Ville-Ometz et al. (2023) summarize, “deep learning and neural methods have become the state of the art for most NLP applications,” including “automatic term extraction, language mining, and assessment of quality in machine translation”. These methods leverage linguistic, statistical, and machine learning features to identify and manage domain-specific terms efficiently [50].
- **Word Formation Studies:** Technical fields often exhibit characteristic word formation processes. Studies on neologisms in science and technology highlight processes relevant to AI, including:
 - **Compounding:** Combining existing words (e.g., *machine learning, neural network, knowledge graph*).
 - **Affixation:** Adding prefixes or suffixes (e.g., *pre-training, transformer, explainability*).
 - **Acronymization:** Forming terms from initial letters (e.g., *AI, NLP, LLM, CNN, GAN*).
 - **Semantic Shift/Repurposing:** Assigning new, specialized meanings to existing words (e.g., *attention, embedding, layer, agent, hallucinate* in the context of LLMs).

- **Historical and Conceptual Analyses of AI:** Works tracing the history of AI concepts and debates have consistently shown how terminology and metaphor shape both the field’s development and its public perception. As Mitchell (2024) observes, “the metaphors we use are not merely rhetorical devices—they actively shape the ethical frameworks we construct to govern these systems” [79]. Recent scholarship has documented “a notable increase in anthropomorphic language over recent years” [54] in LLM research, with terms like “hallucination” and “agent” subtly influencing how capabilities and risks are understood. Kūlis (2024) further notes that “metaphors in artificial intelligence (AI) do more than illustrate; they shape ideas, research aims, and ethical postures, carrying historical and societal biases into complex technological concepts” [63].
- **Early Corpus Studies:** While large-scale, contemporary corpus-based studies of AI terminology are still emerging, earlier research has already illustrated the value of corpus methods for tracing the evolution and contextual use of technical language. As recent reviews emphasize, “terminology today has successfully adopted an approach to collecting lexical data based on corpora” [74]. Corpora have become a cornerstone in the success of machine translation.

The rapid advancement, increasing specialization, and growing societal impact of AI create a compelling need for a more systematic, empirical characterization of its expert terminology. Understanding how experts actually use key terms – their frequency, specificity, typical combinations, syntactic functions, and contextual meanings – is crucial for researchers within the field, for educators training the next generation of AI specialists, for practitioners applying AI techniques, and for facilitating clearer communication across disciplinary and societal boundaries. While existing linguistic research in related areas provides essential theoretical and methodological foundations, a focused, large-scale corpus study dedicated to characterizing *modern expert* AI terminology using contemporary tools like Sketch Engine represents a timely and valuable contribution. This thesis seeks to fill this niche by moving beyond anecdotal

observations, dictionary definitions, or purely technical descriptions to provide a data-driven linguistic profile of the vocabulary shaping contemporary AI discourse.

1.4 Defining the Framework for Analysis

This section explains how the ideas and methods we've covered lead to the specific research plan for this thesis. It gives the reasons for using a corpus approach, describes the text collection, tools, and analysis steps, and connects this whole setup directly to our goal of characterizing today's expert AI terms.

1.4.1 Foundations and Context

The preceding sections have established the key foundations for this research. Firstly, **Terminology Studies** has evolved from prescriptive, concept-centric models (GTT) towards descriptive, usage-based perspectives (like CTT) that recognize terms as dynamic cognitive-linguistic units functioning within specific communicative contexts. These modern theories emphasize the importance of analyzing variation, polysemy, and context in understanding specialized vocabulary. Secondly, **Corpus Linguistics** gives us the practical tools to study how language is actually used in the real world, especially when dealing with large amounts of text. Its methods allow us to systematically look at things like how often words appear, what words typically go together, which terms are particularly important (keywords), common sentence structures, and how context shapes meaning. This way, we can base our ideas about language on solid evidence from real examples. Thirdly, the specific object of study, **modern AI terminology**, exhibits characteristics that demand such an approach: it is rapidly evolving, highly dynamic, interdisciplinary, characterized by a mix of neologisms and repurposed general language words, heavily reliant on acronyms, and prone to semantic shifts and variations in meaning depending on context.

Connecting these leads to a clear conclusion: the inherent nature of modern AI terminology necessitates a theoretical stance aligned with descriptive, context-aware perspectives like CTT, operationalized through the empirical, data-driven methodologies offered by Corpus Linguistics. An approach focused solely on

prescriptive ideals or dictionary definitions would fail to capture the richness, complexity, and dynamism of how terms are actually used by experts in this field.

1.4.2 Justification of the Corpus-Based Methodology

Based on the research background mentioned above, this thesis uses a method focused on analyzing real texts (**corpus-based methodology**) to identify and characterize the key terminological features of modern expert AI discourse. This approach enables us to provide empirical evidence based on authentic language use, aligning with modern terminological theory and directly addressing the dynamic nature of AI language.

The **primary corpus** for this study is a specialized, contemporary corpus compiled from the full-text proceedings of the main conference of the Association for the Advancement of Artificial Intelligence (AAAI) covering recent years. This source is chosen for its **representativeness** of high-impact, peer-reviewed research communication by leading experts across various subfields of AI, ensuring the focus is on *expert* discourse. Its **timeliness** ensures the analysis captures *modern* terminology prevalent in the current research landscape.

For comparative analysis, a large **general English reference corpus** available at SketchEngine (English Trends, 2014–present) will be utilized. This comparison is essential for identifying keywords – terms that are statistically significantly more frequent in the AI expert corpus than in general language, thus highlighting domain specificity.

The **primary analysis platform** will be **Sketch Engine**. This platform is selected due to its comprehensive suite of tools specifically designed for sophisticated corpus linguistic analysis, which are particularly well-suited for terminological characterization. Key functionalities to be employed include:

- **Keywords:** To identify statistically salient, domain-specific terms in the AAAI corpus compared to the general reference corpus.

- **Frequency Lists:** To determine the overall prominence of specific terms within the expert discourse.
- **Concordance (KWIC & CQL):** To examine terms in their immediate linguistic context, facilitating qualitative analysis of meaning, usage variations, and semantic nuances.
- **Word Sketch:** To generate detailed profiles of terms' typical collocational partners and grammatical functions, revealing syntactic patterns and semantic associations preferred in expert usage.

This methodological framework is specifically designed to **address the dynamic characteristics of AI terminology** identified earlier:

- **Neologisms and Domain Specificity:** Keyword analysis is the primary tool for identifying terms that are statistically prominent and characteristic of the AI domain, including newly coined terms or existing words used with highly specialized frequency.
- **Semantic Variation and Shift:** Close analysis of concordance lines will reveal different shades of meaning in context. Examination of collocates and grammatical patterns via Word Sketches can highlight different usage profiles associated with polysemy or semantic nuances.
- **Usage Patterns (Collocations/Syntax):** Word Sketches provide the core data for characterizing the typical lexical environment (collocates) and syntactic roles (grammatical relations) of key AI terms within expert discourse. Frequency data establishes the relative importance or salience of terms.
- **Acronyms:** Frequency lists will identify common acronyms, while concordance analysis will allow examination of their definition, expansion, and contextual use.

Crucially, this framework **directly supports the aims and objectives** of the thesis, which are broadly defined as identifying and characterizing the key terminological features of modern expert AI discourse using corpus-based methods. The chosen corpus provides authentic data reflecting this discourse, the selected tools (Sketch Engine) offer the necessary analytical power, and the specific techniques (keyword

analysis, concordance, Word Sketch, etc.) yield empirical data on term frequency, specificity, collocational behavior, syntactic function, semantic associations, and contextual usage – precisely the features required for a comprehensive characterization.

CONCLUSION TO CHAPTER 1

In conclusion, this chapter has laid the essential theoretical groundwork for the analysis undertaken in this thesis. It began by outlining the evolution of Terminology Studies, moving from traditional, prescriptive approaches towards modern, descriptive perspectives like the Communicative Theory of Terminology (CTT). This shift emphasizes understanding terms as dynamic cognitive-linguistic units functioning within specific communicative contexts, acknowledging phenomena like variation, polysemy, and semantic change – aspects particularly relevant to fast-evolving fields.

Furthermore, the chapter established the crucial role of Corpus Linguistics (CL) as the methodological engine best suited for implementing these modern theoretical viewpoints. We highlighted how CL's focus on authentic language use, combined with its powerful quantitative and qualitative techniques (such as frequency analysis, keyword extraction, collocation analysis, and concordance examination), provides the empirical rigor necessary to investigate specialized language in action. This data-driven approach moves beyond intuition or anecdotal evidence, allowing for systematic observation of terminological patterns.

Crucially, we then contextualized these theoretical and methodological considerations within the specific domain of Artificial Intelligence. The discussion traced the evolution of key AI concepts and their corresponding terminology, noting characteristics like rapid neologism, the prevalence of acronyms, semantic specialization of general English words, and the significant role of metaphor in conceptualizing new technologies. These features underscore why a descriptive, corpus-based approach is not just suitable, but essential for capturing the contemporary reality of expert AI language.

Therefore, the analytical framework adopted for this thesis – grounded in descriptive terminology theory (CTT-aligned) and operationalized through robust empirical methods (CL using a representative, contemporary corpus like the AAAI-24 collection and sophisticated tools like Sketch Engine) – constitutes a logical, coherent, and well-justified strategy. It provides a sound theoretical and methodological foundation for the subsequent chapters, which will delve into the specific procedures (Chapter 2) and empirical findings (Chapter 3) emerging from the application of this framework to the AAAI corpus. The justification of this methodology is thus not merely a procedural description but an integral part of the thesis’s argument, demonstrating that the chosen approach is the most appropriate and effective means to address the research questions posed about the landscape of modern expert AI terminology.

CHAPTER 2. METHODOLOGY AND CORPUS FOR THE ANALYSIS OF MODERN EXPERT AI TERMINOLOGY

This chapter details our research method for looking into the terms AI experts use right now. We'll explain how we put together our collection of texts (our corpus), what tools and data we used to examine it, and the exact procedures we followed. Our strategy involved using computer tools for the initial data collection and processing, combined with established language analysis techniques (corpus linguistics), mainly through the Sketch Engine platform. This allowed us to really characterize the AI language we were studying using solid data.

2.1. Corpus Compilation and Characteristics

A cornerstone of corpus linguistics is the careful design and compilation of the corpus itself, as the nature and quality of the data fundamentally shape the potential findings. For this study, a specialized, synchronic corpus was created to reflect contemporary expert discourse in AI.

2.1.1. Source Selection Rationale (AAAI-24)

The source material chosen for the specialized corpus consists of the titles and abstracts from the proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence (AAAI-24), held in 2024 [10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29]. This selection was driven by the need for **representativeness** and **timeliness**. AAAI is a leading international conference encompassing a wide range of AI sub-disciplines, ensuring the corpus reflects high-impact, peer-reviewed research communication by experts. Focusing exclusively on the 2024 proceedings guarantees that the data captures the most *contemporary* state of expert discourse in this rapidly evolving field, which is crucial for analyzing *modern* AI terminology. Titles and abstracts were selected as they typically contain a high density of key terminology while offering concise summaries of research, a common practice in corpus-based studies focusing on specific discourse elements.

2.1.2. Data Collection and Corpus Size

The titles and abstracts from the 2,865 articles in the AAAI-24 main conference proceedings were collected using a custom Python script developed for this purpose. The script systematically scraped the relevant text data from the official conference proceedings website. The resulting specialized corpus comprises approximately 544,224 words and around 4,000,000 characters. While specialized corpora can vary significantly in size, a corpus of this size is generally considered sufficient to provide valuable insights into Language for Specific Purposes (LSP) and terminology usage, containing enough instances of relevant terms for reliable analysis.

2.1.3. Initial Keyword Extraction and Frequency Ranking

First, we needed a manageable list of the most important terms to analyze in detail. To get this, we extracted keywords and ranked them by frequency. We did this in two main steps:

1. **Term Candidate Identification:** The Gemini API was utilized to process the collected titles and abstracts. It was prompted to identify and extract potential AI-related terms, both single-word and multi-word units. This leverages the capabilities of large language models for initial, broad identification of candidate terms within the large text dataset. The use of advanced models aligns with the trend toward leveraging neural architectures and contextual embeddings for more accurate and semantically meaningful term extraction [33].
2. **Frequency Calculation and Filtering:** Python scripts were then used to calculate the frequency of each identified candidate term across the entire AAAI-24 corpus. As highlighted by Beliga et al., “the most important features for classifying a keyword candidate in these systems are the frequency and location of the term in the document” [35]. Therefore, a minimum frequency threshold of 5 occurrences was applied. This filtering step helps to exclude extremely low-frequency items, which might be errors, highly idiosyncratic usages, or terms not yet established

within the discourse community captured by this specific corpus, thus focusing the analysis on terms with a demonstrable level of usage.

This procedure resulted in a frequency-ranked list of approximately 550 unique candidate terms (e.g., *machine learning*, *reinforcement learning*, *LLMs*, *transformer*). This list served as the primary input for the thematic and structural analysis conducted in Phase 1 of the research.

2.2. Analytical Tools and Reference Data

The core linguistic analysis relies on the Sketch Engine corpus analysis platform and a large general English corpus for comparative purposes.

2.2.1. Sketch Engine Environment

Sketch Engine was chosen as the central platform for the detailed analysis phases (Phases 2 and 3) due to its robust suite of tools specifically designed for corpus linguistic research, lexicography, and terminology work. The AAI-24 corpus was uploaded to the platform for analysis. The key functionalities of Sketch Engine relevant to this thesis include:

- **Wordlist Generation:** Creating frequency-ranked lists of words or lemmas in the corpus.
- **Keywords Tool:** Identifying statistically significant keywords by comparing the specialized corpus against a reference corpus.
- **Concordance:** Allowing detailed examination of terms in their linguistic context (KWIC - Key Word In Context) using simple searches or the advanced Corpus Query Language (CQL).
- **Word Sketch:** Providing automated, one-page summaries of a word's grammatical and collocational behaviour, categorized by grammatical relations.

These combined tools allow for a really thorough look at the terminology – covering its structure, meaning, context, and comparisons as outlined in the research objectives.

2.2.2. English Trends Corpus as Reference

To assess the domain specificity of the AI terms identified in the AAI-24 corpus, comparison with a large general language corpus is necessary. The “English Trends (2014–today)” corpus, integrated within Sketch Engine, was selected as the reference corpus for this study. This multi-billion-word corpus, compiled from web sources, represents a broad range of contemporary general English usage from 2014 and “is updated weekly with new texts and grows by about 70 million words every week” [70]. Comparing the AAI-24 corpus against this up-to-date baseline allows for a statistically grounded identification of terms that are significantly more characteristic of the specialized AI discourse than of general contemporary English, as planned for Phase 2.

2.3. Research Procedures

The empirical investigation was structured into three phases, applying different analytical techniques to the corpus data.

2.3.1. Phase 1: Thematic Clustering and Structural Analysis Procedure

This initial phase focused on characterizing the ~550 high-frequency candidate terms identified in Section 2.1.3. The analysis was performed using the Gemini API, guided by specific prompts designed for thematic and structural classification:

- **Thematic Clustering:** Each term was automatically assigned to one of 10 predefined conceptual categories (*Learning Paradigms, Model Structures & Components, Optimization & Training, Generative Models & Techniques, Data & Feature Engineering, Evaluation & Analysis, Domains & Applications, Challenges & Issues, Other AI Concepts*, as listed in the prompt). This automated clustering provides a high-level overview of the conceptual landscape represented by the most frequent terms in the corpus.
- **Structural Analysis:** Simultaneously, each term underwent automated structural analysis based on a detailed prompt. The API extracted information on: Structure

Type (e.g., Noun Phrase, Acronym), Structure Pattern (e.g., Adj+N, N+N), Number of Compounding Words, and Acronym/Initialism Details. This yielded systematic data on the typical linguistic forms (e.g., single words, multi-word units, acronyms) prevalent in the expert terminology.

This phase demonstrates the integration of Generative AI tools into linguistic analysis. It benefits from their capacity for rapid processing and classification of lexical items based on defined criteria. The results of this phase, including the categorized and structurally annotated term list are presented in Chapter 3.

2.3.2. Phase 2: Domain Specificity Analysis Procedure (Keyword Comparison)

The objective of this phase was to identify terms that are statistically distinctive of the AI expert domain represented by the AAI-24 corpus. While using GenAI in linguistic analysis allows for automating tedious tasks to some extent, “corpora along with corpus tools still preserve their privilege” [60] in reliability and control. The procedure involved using Sketch Engine’s “Keywords” functionality. This tool was used to statistically compare the word frequencies in the AAI-24 corpus against the frequencies in the “English Trends (2014–today)” reference corpus. The output is a list of keywords ranked by a statistical measure (LogDice), indicating terms that are significantly overrepresented in the specialized AI corpus. This analysis provides empirical evidence for the domain specificity of the vocabulary and identifies the core terminology characteristic of this expert discourse, which is discussed in Chapter 3.

2.3.3. Phase 3: Semantic and Contextual Analysis Procedure (Word Sketch, Concordance)

This phase aims to delve deeper into the meaning and usage of selected key AI terms identified through frequency analysis (Phase 1) and keyword analysis (Phase 2). The planned procedure involves:

1. **Term Selection:** A small subset of terms was chosen for detailed analysis. Selection criteria included high frequency, high keyword score, theoretical interest (e.g., terms known for polysemy or semantic shift like *model*, *attention*, *bias*,

hallucination), or representation of key conceptual categories identified in Phase 1.

2. Analysis using Sketch Engine Tools:

- **Word Sketch:** Generating Word Sketches for the selected terms within the AAI-24 corpus to obtain structured summaries of their significant collocates and grammatical patterns. This reveals typical syntactic roles and semantic associations within the expert discourse.
- **Concordance:** Examining numerous concordance lines for each selected term to observe its behavior in diverse linguistic contexts. This qualitative analysis will help identify different senses, semantic nuances, typical phrasings, potential metaphorical uses, and evidence of semantic specialization or shifts compared to general usage.

3. **Comparative Interpretation:** The findings from the specialized corpus were interpreted, and where relevant, compared against the usage patterns observed in the general English reference corpus (English Trends (2014–today)) to highlight specific aspects of semantic specialization or divergence in the AI domain.

The results of this detailed semantic and contextual analysis provide the basis for the discussion in Chapter 3 regarding the nuanced behavior of key terms in modern expert AI communication.

CONCLUSION TO CHAPTER 2

This chapter has detailed the methodological for investigation of modern expert AI terminology presented in this thesis. It laid out the practical steps taken to assemble the necessary linguistic resources and the analytical procedures employed to examine them.

Firstly, the chapter described the process of corpus compilation. The selection of the AAI-24 conference proceedings (titles and abstracts) was justified based on its representativeness of contemporary, high-impact expert discourse and its timeliness.

The data collection method and resulting corpus size (approx. 544,224 words) were specified, establishing the primary dataset.

Secondly, the procedure for initial term identification and filtering was explained, involving a hybrid approach leveraging computational tools for candidate extraction followed by frequency-based filtering to yield a core list of 550 high-frequency candidate terms for analysis.

Furthermore, the chapter identified the core analytical tools and reference data. The Sketch Engine platform was designated as the primary environment for in-depth linguistic analysis, leveraging its robust suite of tools. The English Trends (2014-today) corpus was introduced as the vital reference corpus for comparative analysis.

The overall analytical process, executed systematically, involved initial thematic clustering and structural analysis to gain a broad overview of the lexicon, followed by statistical keyword comparison against general English to determine domain specificity. Finally, the methodology incorporates detailed semantic and contextual investigation of key terms, examining collocational profiles, grammatical patterns, meaning shifts, and metaphorical usage through tools like Word Sketch and concordance analysis.

CHAPTER 3. CHARACTERISTICS OF MODERN EXPERT ARTIFICIAL INTELLIGENCE TERMINOLOGY: A CORPUS-BASED ANALYSIS

This chapter presents the empirical findings derived from the analysis of the specialized AAI-24 corpus, focusing on the characteristics of modern expert Artificial Intelligence terminology. Following the methodology outlined in Chapter 2, the analysis proceeds in stages, beginning with an overview of the dominant conceptual themes and structural patterns identified in Phase 1, followed by analyses of domain specificity (Phase 2) and detailed semantic-contextual behaviour (Phase 3). This initial section, 3.1, details the results of Phase 1, examining the thematic structure and linguistic forms prevalent in the high-frequency terminology of the expert corpus.

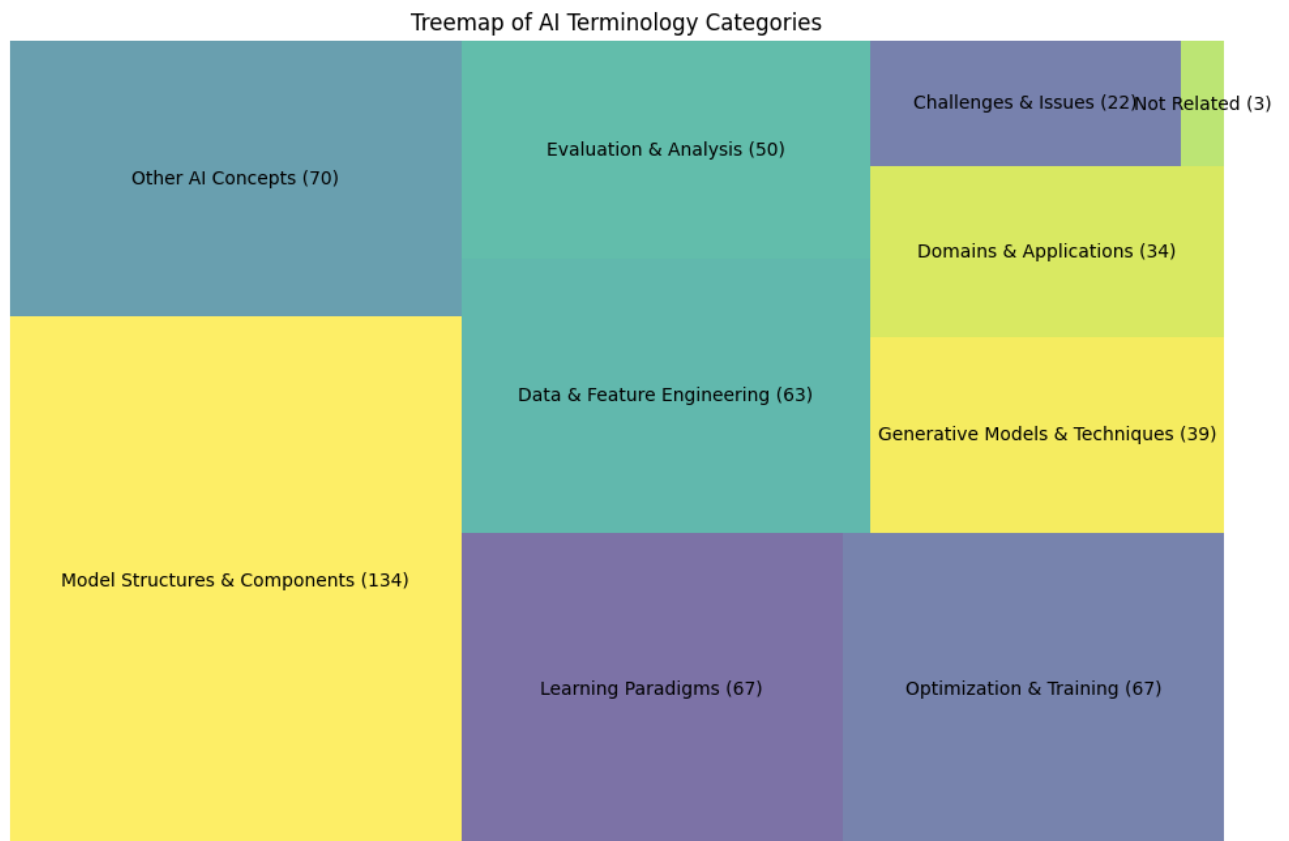
3.1. Thematic Structure and Linguistic Form of Expert AI Terminology (Phase 1 Results)

Phase 1 of the analysis aimed to provide a broad characterization of the terminology landscape within the compiled AAI-24 corpus. This involved thematic clustering and structural analysis of the approximately 550 unique candidate terms identified as occurring five or more times.

3.1.1. Dominant Conceptual Categories in the Expert Corpus

To understand the primary areas of focus within the expert AI discourse represented by the corpus, each of the ~550 high-frequency terms was automatically assigned to one of ten predefined conceptual categories (*Learning Paradigms, Model Structures & Components, Optimization & Training, Generative Models & Techniques, Data & Feature Engineering, Evaluation & Analysis, Domains & Applications, Challenges & Issues, Other AI Concepts*; see Chapter 2 for more methodological details). The resulting distribution provides a quantitative overview of the thematic concentrations within the expert terminology.

Figure 3.1. Distribution of AI Terminology between different categories.



The analysis reveals a distinct concentration of terms within a few core conceptual domains, indicating the primary preoccupations of the research presented in the AAAI-24 proceedings. A treemap visualization (See Figure 3.1) illustrates this distribution graphically, with the area of each rectangle proportional to the number of unique terms assigned to that category. The quantitative counts further confirm the thematic landscape:

1. **Model Structures & Components (134 unique terms):** This category is the largest by a significant margin, highlighting the great importance of the underlying architectures, components, and types of models employed in AI research. The terminology encompasses both general architectural concepts and specific instantiations. High-frequency general terms include *Model* (and its plural *Models*), *Network* (and *Networks*), *Architecture* (and *Architectures*), *Module*, and *Model Parameters*. Alongside these are terms denoting specific influential architectures and components, such as *Neural Network* (along with

variants like *Deep Neural Networks*, *Artificial Neural Networks*, *DNN*, *DNNs*, *NN*, *NNs*), *Transformer* (and *Transformers*, *Vision Transformer*, *ViT*, *Transformer Architecture*), *Graph Neural Networks* (*GNN*, *GNNs*, *GCN*, *GCNs*, *Graph Convolutional Network*), *Convolutional Neural Networks* (*CNN*, *CNNs*), as well as core components like *Encoder*, *Decoder*, *Attention* (and *Self-Attention*, *Cross-Attention*, *Attention Mechanism*, *Attention Module*), and increasingly prominent types like *Diffusion Models*. The sheer volume of terms in this category underscores that a primary focus of the expert discourse is the design, description, and analysis of the systems themselves.

2. **Other AI Concepts (70 unique terms):** This category ranks second in size and groups fundamental, often abstract, or cross-cutting concepts essential to AI but not fitting exclusively into other specific categories. It includes terminology related to agent-based systems and decision-making (*Agent*, *Agents*, *Policy*, *Policies*, *Markov Decision Process / MDP / MDPs*, *Planning*), foundational ideas (*Framework*, *Reasoning*, *Knowledge*, *Causality*, *Causal Inference*, *Symbolic Reasoning*), information flow (*Knowledge Transfer*, *Message Passing*, *Mutual Information*), and system properties (*Convergence*, *Sparsity*, *Semantics*, *Inductive Bias*). Its prominence suggests that, beyond specific models or learning methods, the discourse heavily relies on a set of core operational and theoretical concepts.
3. **Learning Paradigms (67 unique terms):** This category, equal in size to Optimization & Training, covers the diverse approaches and methodologies by which AI systems learn from data. It features high-level terms like *Machine Learning (ML)* and *Deep Learning (Deep-Learning)*, along with a wide array of specific paradigms reflecting the field's dynamism. Notable examples include *Reinforcement Learning* (and *RL*, *Deep Reinforcement Learning / DRL*, *Multi-Agent Reinforcement Learning / MARL*, *Offline Reinforcement Learning*), *Contrastive Learning* (and *Graph Contrastive Learning / GCL*), *Representation Learning* (and *Graph Representation Learning*), *Federated Learning (FL)*, *Self-*

Supervised Learning (SSL), Transfer Learning, Semi-Supervised Learning, Supervised Learning, Unsupervised Learning, Few-Shot Learning, Zero-Shot Learning, Continual Learning, Active Learning, Meta-Learning, Imitation Learning, and newer approaches like *Prompt Learning*. The breadth and frequency of these terms confirm that *how* models learn remains a central research question.

4. **Optimization & Training (67 unique terms):** Sharing the third rank, this category focuses specifically on the practical processes, algorithms, techniques, and mathematical foundations involved in building, refining, and improving model performance. General terms include *Algorithm* (and *Algorithms*), *Training*, *Optimization*, *Learning Algorithm*, and *Optimization Problem*. Specific techniques are abundant, such as *Knowledge Distillation (Distillation, Self-Distillation)*, *Fine-Tuning (Finetuning, Fine-Tune, Prompt Tuning, Parameter-Efficient Fine-Tuning / Pefit)*, *Pre-Training (Pretraining)*, *Regularization*, *Quantization*, *Sampling*, *Backpropagation*, *Gradient Descent* (and *Stochastic Gradient Descent, Gradient, Gradients*), *Policy Optimization*, *Loss* (and *Loss Function / Loss Functions, Contrastive Loss, Cross-Entropy Loss*), *Convergence Rate*, *Dropout*, and specific optimizers like *Adam*. This highlights the significant focus on the engineering aspects of making AI models effective and efficient.
5. **Data & Feature Engineering (63 unique terms):** Reflecting the data-centric nature of modern AI, this large category includes terms related to the input data itself and its transformation into usable formats. It covers concepts like *Features (Feature, Visual Features, Semantic Features)*, *Data Augmentation (Augmentation, Mixup)*, *Representations (Representation, Node Representations, Feature Representation, Latent Representation)*, *Embeddings (Embedding, Embedding Space)*, specialized data structures like *Knowledge Graphs (Knowledge Graph, KGs, Knowledge Graph Embedding, Knowledge Graph Completion)* and *Point Clouds (Point Cloud)*, data sources (*Training*

Data, Dataset, Datasets, Benchmark Datasets, Imagenet, MNIST, CIFAR-10, CIFAR-100), data properties (*Distribution, Data Distribution, Data Heterogeneity, Noisy Labels, Label Noise*), and processes such as *Feature Extraction, Feature Selection, Feature Fusion, Feature Learning, Feature Matching, Feature Alignment, Sampling, and Encoding*.

6. **Evaluation & Analysis (50 unique terms):** This category comprises terms used to measure, assess, and understand the performance, behaviour, and limitations of AI models. It includes task-specific evaluation terms like *Classification* (and *Classifier/ Classifiers, Node Classification, Image Classification, Text Classification, Binary Classification, Multi-Label Classification*), *Segmentation* (and *Semantic Segmentation, Instance Segmentation*), *Prediction* (and *Predictions, Link Prediction*), *Regression, Detection* (and *Object Detection, Anomaly Detection, OOD Detection, Outlier Detection*), as well as terms related to model properties and analysis methods like *Generalization* (and *Model Generalization, Domain Generalization, Zero-Shot Generalization, Generalization Ability*), *Clustering* (and *Multi-View Clustering, K-Means Clustering*), *Inference* (and *Variational Inference*), *Interpretability, Explainability* (and *Explainable AI / XAI*), *Robustness* (*Model Robustness*), *Accuracy* (*Model Accuracy, Classification Accuracy*), *Transferability, Uncertainty, Fairness, Regret* (*Regret Bound*), and methodologies like *Ablation Studies*. Note that some terms like *Out-Of-Distribution* (*OOD*) function adjectivally but represent evaluation concepts.
7. **Generative Models & Techniques (39 unique terms):** This rapidly growing area focuses on models that create new data. While *Transformer* and *Diffusion Model* also appear under Model Structures, their high relevance here is notable. Key terms specific to or strongly associated with this category include *Large Language Models* (and *LLM, LLMs, LMS, Language Model, Generative Language Models*), *Generative Models* (*Generative Model, Generative Modeling, Deep Generative Models*), *Generative AI*, specific models like

Chatgpt, GPT-4, GPT-3, GPT-3.5, Bert, network types like Generative Adversarial Networks (GAN, GANs) and Variational Autoencoder (VAE), and techniques like Text-To-Image (Text-To-Image Generation, Text-To-Image Models), Text Generation, Image Generation, Image Synthesis, Chain-Of-Thought (CoT), and Stable Diffusion.

8. **Domains & Applications (34 unique terms):** Although relatively smaller in the count of unique high-frequency terms drawn from abstracts, this category indicates key areas where the developed AI methods are applied. Prominent examples include *Natural Language Processing (NLP, Natural Language Understanding, Natural Language Generation), Object Detection, Computer Vision, Image Classification, Autonomous Driving, Recommender Systems (Recommendation), Question Answering (VQA), Sentiment Analysis, Machine Translation, and Named Entity Recognition (NER)*. The focus in abstracts on methodology might explain the lower unique term count here compared to method-centric categories.
9. **Challenges & Issues (22 unique terms):** This category groups terms addressing the problems, limitations, risks, and ethical considerations associated with AI. It includes concepts like *Adversarial Attacks (Adversarial Attack, Adversarial Examples, Adversarial Perturbations, Adversarial Robustness), Overfitting (Over-Fitting, Over-Smoothing), Catastrophic Forgetting, Fairness (Algorithmic Fairness), Bias (Debiasing), Backdoor Attacks (Backdoor), Hallucinations (Hallucination), Label Noise, and Domain Shift*. While crucial, these terms form a smaller proportion of the highest-frequency discourse compared to core methodology terms.
10. **Not Related (3 unique terms):** This category contained a negligible number of terms (*AI itself, Natural Language, PLL*), very broad umbrella terms or artifacts caught by the frequency threshold, but not central to specific technical discussions.

In summary, the thematic analysis of high-frequency terms in the AAI-24 corpus paints a clear picture of the expert discourse. The landscape is dominated by terminology related to the **design and components of AI models** (especially neural networks and transformers), the diverse **learning paradigms** employed, the crucial processes of **optimization and training**, and the essential role of **data and its representation**. While evaluation, specific applications, generative techniques, and challenges are all represented, the core focus, as reflected in the most frequent terminology used in abstracts, remains firmly rooted in the foundational methodologies and architectures for building and training sophisticated AI systems. This provides a crucial map of the conceptual landscape upon which the structural and semantic properties of the terminology are further explored.

3.1.2. Structural Features: Acronyms and Multi-Word Term Patterns

Following the thematic overview, this section delves into the structural characteristics of the high-frequency terminology identified in the AAI-24 corpus. Understanding how terms are formed—whether as single words, complex multi-word units, or acronyms—provides insight into the linguistic mechanisms used to create and designate concepts within the expert AI community. The analysis focuses on the predominant structure types and the recurring patterns observed in term construction.

The terminology landscape is heavily dominated by nominal structures. An analysis of the `structure_type` reveals that Noun Phrases (NPs) constitute the vast majority of the high-frequency terms (457 instances). This aligns with the fundamental function of terminology in scientific discourse: naming the concepts, models, components, processes, data types, and metrics that are central to the field.

Beyond simple Noun Phrase formation, two key structural strategies are particularly prominent: the use of acronyms and the construction of multi-word terms.

Acronyms represent a significant feature of the expert AI lexicon. They were identified as the explicit `structure_type` for 60 unique terms within the high-frequency list. This substantial number underscores the field's reliance on abbreviation as a primary means

of achieving communicative efficiency and brevity among specialists. Acronyms span diverse conceptual categories and include foundational concepts (*AI, ML*), learning paradigms (*RL, SSL, DRL, MARL, GCL*), model types and components (*LLM / LLMs, GNN / GNNs, CNN / CNNs, DNN / DNNs, ViT, BERT, GAN / GANs, VAE, MLP, ResNet*), applications (*NLP, VQA, NER, Lidar*), evaluation concepts (*OOD, XAI, SOTA, DETR*), and specific datasets or optimizers (*MNIST, CIFAR-10, CIFAR-100, Adam, MuJoCo*). Their direct function is usually in compressing longer, descriptive Noun Phrases.

As noted by Wang et al. (2023), “[t]erms may be simple (STs) or multi-word (MWTs)”. The AI lexicon is rich in **MWTs**, predominantly realized as complex Noun Phrases. While single-word nouns form a crucial base layer (105 terms identified with the “N” structure_pattern, e.g., *Model, Algorithm, Network, Training, Attention, Encoder, Loss, Bias, Fairness, Knowledge*), the majority of the terminology involves combinations of two or more words to achieve greater conceptual specificity. Analyzing the structure_pattern for these MWTs reveals highly productive construction templates:

N + N (Noun + Noun): This is the most frequent pattern for constructing MWTs (144 instances). Typically, the first noun modifies the head noun, creating a more specialized concept. This pattern is evidently a highly fertile ground for term creation, yielding examples like *Machine Learning, Object Detection, Knowledge Distillation, Data Augmentation, Language Models, Loss Function, Node Classification, Graph Learning, Feature Extraction, Policy Optimization, Reward Function, Point Cloud, Teacher Model, Student Model, and Search Space*.

Adj + N (Adjective + Noun): Ranking second in frequency (105 instances), this pattern uses an adjectival modifier to qualify or specify a core noun concept. It generates fundamental terms such as *Neural Network, Deep Learning, Contrastive Learning, Federated Learning, Artificial Intelligence, Latent Space, Supervised Learning, Active Learning, Optimal Transport, Adversarial Training, Noisy Labels, Gaussian Process, Semantic Segmentation, Mutual Information, and Inductive Bias*.

Adj + N + N (Adjective + Noun + Noun): This three-word pattern (29 instances) provides further specificity, often building upon Adj+N structures. Key examples include Large Language Models, Graph Neural Networks, Natural Language Processing, Self-Supervised Learning, Few-Shot Learning, Multi-Task Learning, Visual Question Answering, Neural Radiance Fields, Semantic Space, and Partial Label Learning.

Adj + Adj + N (Adjective + Adjective + Noun): Allowing for multiple layers of modification, this pattern (22 instances) is common for descriptive technical terms. Examples include *Deep Neural Networks*, *Semi-Supervised Learning*, *Spiking Neural Networks*, *Artificail Neural Networks*, *Convolutional Neural Networks* (in its full form), and *Multi-Objective Optimization*.

N + N + N (Noun + Noun + Noun): Though less frequent among the top patterns (11 instances), this structure constructs complex concepts by compounding three nouns, seen in terms like *Graph Representation Learning*, *Knowledge Graph Completion*, *Vision-Language Models*, *Markov Decision Processes*, *Knowledge Graph Embedding*, and *Neural Architecture Search*.

In conclusion, the structural analysis of the high-frequency terminology within the AAI-24 corpus highlights a system built primarily on nominal concepts. While single nouns form a foundation, the lexicon expands significantly through the highly productive formation of Multi-Word Terms, chiefly Noun Phrases constructed via recurring patterns like N+N and Adj+N. Concurrently, the extensive use of Acronyms serves as a vital mechanism for abbreviation and efficiency. These structural tendencies — compositionality via MWTs and compression via acronyms — are characteristic linguistic strategies employed by the expert AI community to precisely label and efficiently communicate about the complex and evolving concepts within their domain.

3.2. Identifying Domain-Specific Terms in Expert AI Discourse (Phase 2 Results)

Having established the general thematic landscape and structural patterns in Phase 1, Phase 2 focuses on pinpointing the lexical elements that most distinguish the expert AI

discourse within the AAI-24 corpus from general language use. This is achieved through a comparative statistical analysis, identifying terms whose frequency is significantly higher in the specialized corpus, thereby marking them as characteristic keywords of the domain.

3.2.1. Statistical Keyword Analysis: AI Expert Corpus vs. General English

To quantify the domain specificity of the terminology, a keyword analysis was conducted using Sketch Engine. This process compared the **AI_Expert_AAI24** corpus (focus corpus, consisting of AAI-24 abstracts) against the **English Trends (2014–today)** corpus (reference corpus, representing contemporary general English). The comparison was performed on lemmas (case-insensitive) to group related word forms, using the LogDice statistical measure to rank terms by their relative prominence in the AI corpus. A minimum frequency of 5 occurrences in the AI_Expert_AAI24 corpus was required for a term to be included in the analysis. A higher LogDice score signifies a stronger statistical association of the term with the expert AI domain.

The analysis was performed separately for single-word lemmas and multi-word terms (n-grams) to capture different facets of lexical specificity.

Table 3.1 presents the highest-scoring single-word keywords (lemmas). The list is dominated by terms achieving exceptionally high LogDice scores, indicating their usage is vastly more frequent in the AI expert corpus than in general English.

Table 3.1: Top 30 Single-Word Keywords (Lemmas) in AAI-24 Corpus vs. English Trends (2014-today) ranked by LogDice Score

Keyword (Lemma)	LogDice Score	Frequency (Focus)	Frequency (Reference)	Relative Frequency (Focus)	Relative frequency (reference)
contrastive	3452.352	392	8381	637.68146	0.08474
cross-modal	2188.508	176	3053	286.30594	0.03087

multi-view	1887.73	205	7587	333.48135	0.07671
few-shot	1817.139	194	7292	315.58725	0.07373
vision- language	1453.225	99	1077	161.0471	0.01089
gnns	1423.208	113	2891	183.82144	0.02923
zero-shot	1374.13	142	6743	230.99684	0.06818
ood	1318.708	133	6344	216.3562	0.06414
pre-trained	1226.133	285	27515	463.62042	0.2782
pseudo-label	1208.093	80	772	130.13907	0.00781
self- supervised	1193.265	149	10208	242.38402	0.10321
pre-training	1086.016	157	13378	255.39792	0.13526
gnn	1083.722	124	8528	201.71556	0.08622
semi- supervised	1027.714	119	8749	193.58186	0.08846
out-of- distribution	1021.598	75	1931	122.00538	0.01952
transformer- based	1008.105	88	4164	143.15297	0.0421
multi-agent	994.458	130	11152	211.47598	0.11276
llms	940.609	393	57343	639.30817	0.57978

open-set	931.82	60	480	97.6043	0.00485
sota	915.641	131	13139	213.10272	0.13285
unlabeled	840.68	148	18446	240.75728	0.1865
multi-label	824.404	71	3978	115.49842	0.04022
adaptively	818.704	93	8398	151.28667	0.08491
denoising	808.055	102	10431	165.92731	0.10547
cross-domain	799.73	88	7826	143.15297	0.07913
image-text	785.274	54	1186	87.84387	0.01199
multi-modal	780.491	238	39184	387.16373	0.39618
fine-grained	771.744	200	31818	325.34769	0.3217
interpretability	765.065	109	13045	177.31448	0.13189
discriminative	729.131	104	13072	169.18079	0.13217

Table 3.2 presents the highest-scoring multi-word terms identified through the same process. This list similarly shows very high specificity scores, confirming that particular combinations of words function as highly characteristic terminological units within the AI domain.

Table 3.2: Top 30 Multi-Word Term Keywords in AAI-24 Corpus vs. English Trends (2014-today) ranked by LogDice Score

Keyword (Lemma)	LogDice Score	Frequency (Focus)	Frequency (Reference)	Relative Frequency (Focus)	Relative frequency (reference)
----------------------------	--------------------------	------------------------------	----------------------------------	---	---

extensive experiment	5370.904	648	9523	1054.12646	0.09628
diffusion model	2469.64	260	7052	422.95197	0.0713
state-of-the-art method	2100.773	250	9261	406.6846	0.09364
student abstract	2001.888	123	0	200.08882	0
state-of-the-art performance	1838.093	190	6746	309.08029	0.06821
contrastive learning	1646.199	123	2137	200.08882	0.02161
existing method	1534.136	271	18537	440.8461	0.18742
benchmark dataset	1498.02	151	6334	245.6375	0.06404
downstream task	1478.554	127	3936	206.59578	0.0398
real-world dataset	1409.991	124	4266	201.71556	0.04313
experimental result	1200.127	436	48569	709.25793	0.49107
point cloud	1074.368	203	20519	330.22787	0.20746

graph neural network	1039.269	91	4207	148.03319	0.04254
reinforcement learn	1004.79	197	21664	320.46747	0.21904
generalization ability	959.799	70	1854	113.87169	0.01875
backdoor attack	951.411	65	1112	105.73799	0.01124
proposed method	946.667	274	36688	445.72632	0.37094
knowledge distillation	943.803	72	2394	117.12516	0.02421
target domain	917.254	78	3802	126.88559	0.03844
language models	882.794	84	5430	136.64603	0.0549
adversarial attack	882.645	89	6344	144.77971	0.06414
representation learn	833.083	67	3061	108.99147	0.03095
representation learning	828.105	64	2556	104.11126	0.02584
comprehensive experiment	826.942	56	1017	91.09735	0.01028

novel framework	791.316	62	2728	100.85778	0.02758
pseudo label	788.172	52	737	84.59039	0.00745
domain adaptation	772.37	70	4704	113.87169	0.04756
semantic segmentation	757.855	85	8168	138.27277	0.08258
federated learning	755.214	81	7379	131.76581	0.07461
unlabeled datum	737.869	62	3642	100.85778	0.03682

Together, these ranked lists provide the empirical basis for understanding the vocabulary that defines the lexical distinctiveness of this specialized discourse.

3.2.2. Interpretation of Domain-Specific Vocabulary

The combined results from the single-word and multi-word keyword analyses offer significant insights into the nature of the vocabulary that sets expert AI discourse apart. Several key characteristics emerge from the highest-scoring terms across both lists:

- Prominence of Highly Technical Neologisms and Domain-Specific Formations:** A large proportion of the top keywords consists of terms primarily confined to the AI/ML domain or closely related technical fields. This includes specific modifiers forming characteristic compounds (*Contrastive, Cross-modal, Multi-view, Few-shot, Zero-shot, Pre-trained, Pseudo-label, Self-supervised, Semi-supervised, Out-of-distribution, Transformer-based, Multi-agent*), acronyms representing core concepts (*GNNs, OOD, GNN, LLMs, SOTA, VQA, T2I, SNN, VLMs, DNNS, MAPF, LTLF*), and complete multi-word units

like *Diffusion Model*, *Contrastive Learning*, *Benchmark Dataset*, *Downstream Task*, *Point Cloud*, *Graph Neural Network*, *Backdoor Attack*, *Knowledge Distillation*, *Domain Adaptation*, *Semantic Segmentation*, *Federated Learning*, *Neural Radiance (Fields)*, *Domain Generalization*, *Latent Space*, *Optimal Transport*, *Vision Transformer*, *Data Augmentation*, *Label Noise*, *Style Transfer*, *Prompt Tuning*, *Attention Mechanism*, and *Continual Learning*. The extremely low frequency or absence of many of these in the massive general English reference corpus confirms their high specialization.

- **Specialized Usage of General Vocabulary:** Concurrently, many keywords are words or phrases found in general English but employed with significantly higher frequency and often specialized technical meanings within AI. Single-word examples include *Semantic*, *Generalization*, *Embedding*, *Adversarial*, *Robustness*, *Sparsity*, *Distillation*, *Quantization*, *Multimodal*, *Diffusion*, *Interpretable*, *Trainable*, *Regularization*, and *Dataset*. Multi-word examples containing common words include *State-of-the-art Method / Performance / Result*, *Existing Method*, *Experimental Result*, *Proposed Method*, *Novel Framework*, *Source Domain*, *Target Domain*, *Theoretical Analysis*, *Previous Method*, and *Computational Cost / Complexity*. Their high LogDice scores reflect their terminologization – their crucial role and elevated frequency as labels for specific AI concepts, even if the constituent words themselves are common.
- **Key Role of Compounding and Modification:** The keyword lists underscore the vital role of specific compositional patterns in creating domain-specific terminology. The multi-word keyword list is dominated by NP structures, particularly N+N (*Diffusion Model*, *Point Cloud*, *Knowledge Distillation*, *Data Augmentation*, *Loss Function*, *Node Classification*, *Graph Structure*, *Feature Space*) and Adj+N (*Semantic Segmentation*, *Optimal Transport*, *Theoretical Analysis*, *Computational Cost*, *Public Dataset*). Furthermore, the single-word list highlights the keyword status of specific modifying elements (*Contrastive*,

Multi-view, Few-shot, Pre-trained, Self-supervised, Transformer-based) that combine with core nouns to form a large part of the specialized MWT lexicon.

- **Reflection of Core Research Themes:** The domain-specific vocabulary identified strongly aligns with the primary thematic areas noted in section 3.1.1. Keywords related to **Model Structures & Components** (*Diffusion Model, Graph Neural Network, Neural Radiance, Vision Transformer, Embedding Space, GNN, LLMs, Transformer-based*), **Learning Paradigms** (*Contrastive Learning, Federated Learning, Representation Learning, Semi-supervised Learning, Few-shot, Zero-shot, Self-supervised, Continual Learning*), **Optimization & Training** (*Knowledge Distillation, Domain Adaptation, Optimal Transport, Loss Function, Adversarial Training, Prompt Tuning, Pre-trained, Pre-training, Distillation, Quantization, Regularization*), **Evaluation & Analysis** (*State-of-the-art variants, Benchmark Dataset, Downstream Task, Generalization Ability, Semantic Segmentation, Node Classification, OOD, Interpretability, Robustness*), and **Data & Feature Engineering** (*Point Cloud, Data Augmentation, Pseudo Label, Latent Space, Feature Representation, Knowledge Graph, Embedding, Dataset, Unlabeled*) are heavily represented among the highest-scoring terms. This demonstrates that the lexical core differentiating AI discourse directly corresponds to its central objects of study and methodological concerns.

In summary, the keyword analysis comparing the AAI-24 expert corpus to general English effectively isolates a vocabulary characterized by a mix of highly technical neologisms (including numerous acronyms and specific MWTs), specialized applications of more common words, and recurring patterns of compounding and modification. This distinctive lexicon clearly reflects the field's focus on model architectures, learning techniques, data manipulation, optimization strategies, and evaluation methods, providing a quantitative confirmation of the lexical features that define modern expert AI communication. These findings set the stage for Phase 3, which will explore the semantic and contextual nuances of selected key terms from this domain-specific vocabulary.

3.3. Semantic Behaviour and Contextual Usage of Key AI Terms (Phase 3 Results)

This section initiates the analysis planned in Phase 3, focusing on the semantic behavior and contextual usage of selected key terms. The goal is to move beyond frequency and structure to understand how these terms function meaningfully within the expert AI discourse, particularly in comparison to their usage in general English. We begin with one of the most frequent and foundational terms in the corpus: model.

The analysis below is based on the examination of concordance lines extracted from the AI_Expert_AAAl24 corpus and the English Trends (2014-today) reference corpus using Sketch Engine, focusing on identifying collocational profiles, grammatical patterns, semantic specialization, and potential metaphorical usage.

3.3.1.1. Analysis of: model (lemma)

The term model (including its plural form models) is exceptionally frequent in the AI_Expert_AAAl24 corpus (focus: 5216 instances based on concordance sample size, compared to a vast but proportionally much smaller usage in the general English corpus). Its centrality necessitates a careful examination of its specific usage within the AI domain.

3.3.1.2. Collocational Profiles and Grammatical Patterns

Examining the concordance lines reveals distinct collocational and grammatical preferences for model in the AI corpus compared to general English.

Grammatical Role: In the AI corpus, model functions almost exclusively as a Noun. While the verb form exists in English, instances of “to model” something (meaning to create a representation or simulation) appear rare in the high-frequency abstract context, which favors nominal descriptions of systems and methods.

Common Modifiers (Adjectives preceding model): The AI concordance lines are rich with highly specific technical adjectives modifying the model. Frequent patterns include:

- **Specific Architectures/Types:** Diffusion Model, generative models, Contrastive Diffusion model, language models, neural model, pretrained language models, Heterogeneous-graph-based model, transformer model, probabilistic model.
- **Scale/Scope:** Large language models.
- **Learning Paradigm Association:** Self-supervised models, contrastive models.
- **General Descriptors:** Proposed model, existing models, deep generative models, global model, surrogate model.

In **General English**, modifiers are more varied, referring to product versions (flagship model, 2000 model, Model D), systems/approaches (US model), types (data model), or examples (environmental model). The range is broader and less technically constrained.

Common Noun Compounds (N + model): While model often functions as the head noun, it also appears as the first element in N+N compounds, such as *model performance*, *model design*, *model parameters*, *model accuracy*, *model robustness*, *model optimization*, *model compression*, and *model learning*. This pattern emphasizes aspects of the AI model. This compounding seems less structurally central in the general English samples provided, where model is more often the head noun.

Verbs taking model as Object: Common verbs indicating actions performed on the model include: propose a model, train models, deploy models, evaluate models, fine-tune models, and develop models.

Verbs with model as Subject (Agency): Concordance lines frequently show the model performing actions, indicating a degree of grammatical agency: models have exhibited potential, model learns multifaceted aspects, models distinguish orders, model hampers performance, models fail to explore, model for generating predictions, model performs better. This pattern of attributing actions directly to the model appears more pronounced in the AI context than in general English, where agency is typically reserved for human or physical models acting as exemplars.

Prepositional Phrases: Phrases like model for X (e.g., model for Therapeutic Peptide Generation, model for Fact Checking, model for similarity recognition, model for generating predictions) are common, specifying the model’s purpose.

3.3.1.3. Semantic Specialization and Meaning Shifts

Comparing the usage in the AI corpus with the general English concordance lines reveals significant semantic specialization and a narrowing of meaning for model within the AI domain.

Dominant AI Sense: In the overwhelming majority of AI concordance lines, model refers to a computational construct or mathematical representation designed to learn patterns from data, make predictions, generate outputs, or simulate a process. It signifies an abstract system implemented in software, ranging from statistical algorithms to complex neural networks. This core meaning encompasses various specific types identified by modifiers (language model, diffusion model, generative model, neural network model, etc.).

General English Senses: The general English concordance lines showcase a much broader semantic range for model:

- **Product Version/Type:** flagship model, 2000 model, Model D (specific product designation).
- **System/Approach/Blueprint:** US model, environmental model, conceptual models, RDF data model.
- **Physical Representation/Replica:** (Not strongly present in the small sample, but a common meaning, e.g., scale model).
- **Person Exhibiting Traits (Role Model):** (Not present in sample, but common).
- **Person Employed to Display Fashion (Fashion Model):** Indonesian models.
- **Exemplar/Pattern:** (Implicit in some uses, like model environmental community).

Shift/Specialization: While the general sense of “representation” or “system” exists in both corpora, the AI domain almost exclusively uses model to refer to the specific computational/mathematical construct. The common general English meanings related to product types, physical replicas, or people (role/fashion models) are absent in the expert AI usage observed in the abstracts. This constitutes a clear case of semantic specialization where a general term acquires a precise, dominant technical meaning within a specific field.

3.3.1.4. Metaphorical Patterns in Expert AI Terminology

While model itself isn’t a novel metaphor created by AI (the idea of mathematical/conceptual models predates it), its usage within AI contexts often participates in broader metaphorical frameworks:

- **MODEL AS ARTIFACT/MACHINE:** Concepts like model design, model architecture, model parameters, building models, and deploying models frame the AI model as something constructed, engineered, and having components. Its performance can be measured (model performance, model accuracy), suggesting a machine-like entity evaluated on its operational effectiveness.
- **MODEL AS LEARNER/AGENT:** The frequent collocation with learning and verbs showing agency (model learns, model distinguishes, model predicts, model generates) positions the model as an active entity capable of acquiring knowledge and performing tasks. This anthropomorphic framing is pervasive, though generally understood as shorthand within the expert community.
- **MODEL AS REPRESENTATION:** This is the most direct, non-metaphorical sense, but the kind of representation is specific – an abstract, computational representation of patterns, functions, or knowledge derived from data.

Comparing this to general English, while the “model as representation” idea exists (e.g., data model, conceptual model), the specific metaphorical framings of “model as engineered artifact” and “model as active learner/agent” seem much more central and

pervasive in the AI discourse, reflecting the field's focus on building functional, data-driven systems.

Summary for model

The term *model* in the expert AI corpus is overwhelmingly a noun referring to a specific type of computational/mathematical construct. Its usage is characterized by highly specific technical collocates (modifiers like *diffusion*, *language*, *generative*, *contrastive*, *neural*) and participation in common compounds (*model performance*, *model design*). Grammatically, it frequently appears as the object of verbs related to creation and training (*propose*, *train*) and often functions as the subject of verbs indicating performance or action (*learns*, *predicts*, *generates*). This contrasts with the broader semantic range in general English (*product types*, *exemplars*, *people*). While building on the general idea of a representation, AI usage emphasizes the model as an engineered artifact and an active (if metaphorical) agent capable of learning and performing tasks.

3.3.2.1. Analysis of: learning (lemma)

The lemma *learn* (appearing frequently as the gerund/noun *learning*) is another cornerstone term within the **AI_Expert_AAAI24** corpus. While ubiquitous in general English, its high frequency and central role in defining many AI methodologies warrant a detailed examination of its specific behavior within the expert discourse.

3.3.2.2. Collocational Profiles and Grammatical Patterns

In the AI corpus, *learning* functions overwhelmingly as a **Noun** (often as part of a compound noun phrase) denoting a process or paradigm, rather than the verb *learn* denoting an action performed by a human subject. The Word Sketch data reveals strong collocational patterns distinct from general usage:

Modifiers (Defining Types of Learning): The most striking feature in the AI corpus is the vast array of highly specific technical terms modifying *learning*. The top modifiers identified in the Word Sketch (**AI_Expert_AAAI24**) include: *contrastive*,

representation, machine, reinforcement, deep, continual, semi-supervised, graph, self-supervised, supervised, online, active, prompt, incremental, transfer, feature, policy, and unsupervised. These modifiers are crucial for distinguishing the specific computational paradigm being discussed.

In the English Trends corpus Word Sketch, modifiers are more general and relate primarily to human education or cognition: *remote, distance, in-person, deep* (can overlap, but often refers to depth of understanding), *lifelong, experiential, online, hands-on, continuous, virtual, project-based, language, personalized, student, social-emotional.* While some overlap exists (e.g., *online, deep*), the type and frequency of technical modifiers in the AI corpus are defining characteristics.

Modified Nouns (What learning specifies): In the AI corpus, learning frequently modifies other nouns, indicating a specific kind of framework, technique, or process. Top examples from the Word Sketch include: *learning framework, learning paradigm, learning technique, learning strategy, learning algorithm, learning approach, learning setting, learning environment, learning method, learning model.*

In English Trends, frequently modified nouns relate strongly to education and personal development: *learning curve, learning environment, learning experience, learning difficulty, learning disability, learning outcome, learning opportunity, learning platform, learning style.* Again, the focus is predominantly human-centric.

Verbs with learning as Object: Verbs in the AI corpus often denote enabling, facilitating, or applying specific learning processes: *federate learning, facilitate learning, enable learning, perform learning, supervise learning, accelerate learning, guide learning, integrate learning, employ learning, base learning, use learning, propose learning* (methods).

In English Trends, verbs often relate to the human experience of learning: *facilitate learning* (can overlap), *leverage learning, foster learning, incorporate learning, accelerate learning, enhance learning, personalize learning, promote learning, continue learning.* The focus is often on educational delivery or personal development.

Verbs with learning as Subject: The AI corpus occasionally treats the process of learning as an agent: *learning aims, learning has, learning emerges, learning becomes*.

English Trends shows similar patterns but also verbs reflecting the human experience: *learning happens, learning occurs, learning requires, learning evolves*.

Prepositional Phrases: Patterns like *learning from data, learning of representations, and learning for detection/classification* are common, specifying the input, output, or goal of the computational process.

3.3.3.3. Semantic Specialization and Meaning Shifts

The primary semantic difference lies in the referent of learning.

Dominant AI Sense: In the AI corpus, learning almost invariably refers to the **computational process** whereby an algorithm modifies its parameters or internal structure based on exposure to data, leading to improved performance on a specific task. It is an algorithmic, data-driven process within a machine system. Examples: “multi-modal contrastive learning strategy”, “deep reinforcement learning (DRL)”, “representation learning”, “federated learning”, “self-supervised learning”, “prompt learning”.

Dominant General English Senses: The data from English Trends confirms the primary meaning relates to the **cognitive process of acquiring knowledge, skills, or understanding** by humans (or sometimes animals) through study, experience, or teaching. Examples: “required learning for every child”, “learning how to convert prospects”, “learning our community”, “lifelong learning”, “experiential learning”, “learning a lot of lessons”. A secondary common meaning refers to the body of knowledge or skills acquired.

Shift/Specialization: The AI usage represents a clear **metaphorical extension** and **technical specialization** of the primary cognitive sense. The core concept of “acquiring capability through exposure/experience” is mapped from the biological/cognitive domain onto computational systems processing data. Within AI,

this technical sense has become the default and dominant meaning, largely displacing the direct reference to human cognition found in general English. The proliferation of specific technical modifiers (contrastive, reinforcement, federated, etc.) further solidifies this specialized meaning – these types of learning do not exist in the general human sense.

3.3.3.4. Metaphorical Patterns in Expert AI Terminology

The use of learning in AI is fundamentally metaphorical, drawing heavily on the **AI SYSTEM AS LEARNER** conceptual metaphor.

Core Mapping: The process of algorithms adjusting to data is conceptualized as analogous to humans or animals learning from experience or instruction. Data becomes the “experience” or “teaching material”, improved performance is “knowledge” or “skill”, and the algorithm/model is the “learner”.

Extension to Paradigms: This core metaphor extends to the names of specific paradigms, often borrowing from psychology or education:

- **Supervised Learning:** Metaphor of learning with a “teacher” providing correct answers (labeled data).
- **Unsupervised Learning:** Metaphor of learning patterns without explicit guidance.
- **Reinforcement Learning:** Metaphor of learning through trial and error, guided by “rewards” and “punishments” (feedback signals).
- **Transfer Learning:** Metaphor of applying knowledge learned in one domain to another.
- **Continual/Lifelong Learning:** Metaphor of ongoing learning and adaptation, avoiding “forgetting” (catastrophic forgetting).

Agency: As noted in 3.3.1, attributing actions to the learning process or the model reinforces this agentic, learner-centric metaphor.

Contrast with General English: While the word learning inherently relates to cognition in general English, its application to computational processes in AI is a systematic, technical metaphor that structures much of the field's understanding and discourse about how AI systems develop capabilities.

Summary for learning

In the expert AI corpus, learning predominantly functions as a noun denoting a computational process of algorithmic adaptation to data. It is characterized by collocation with highly specific technical modifiers (e.g., contrastive, reinforcement, federated, self-supervised) that specify the learning paradigm, and it frequently modifies nouns like framework, paradigm, or strategy. This technical usage is a specialized, metaphorical extension of the primary general English sense referring to human/animal cognitive acquisition of knowledge or skill. The AI SYSTEM AS LEARNER metaphor is foundational, structuring the terminology for various learning paradigms (supervised, unsupervised, reinforcement, etc.) and framing the process of algorithmic improvement in cognitive terms.

3.3.4.1 Analysis of: attention (lemma)

The term attention is present in both the AI_Expert_AAAl24 corpus and the English Trends general corpus, but its frequency and usage patterns differ significantly, highlighting a strong case of semantic specialization within the AI domain. While its raw frequency in the AI corpus (512 instances in the concordance sample) is much lower than in the massive general English corpus, its keyword score and specific collocational profile indicate its terminological importance in AI.

3.3.4.2. Collocational Profiles and Grammatical Patterns

Grammatical Role: In both corpora, attention primarily functions as a Noun.

Modifiers (AI Corpus): The Word Sketch for AI_Expert_AAAl24 reveals highly specific technical modifiers: self-attention (though often treated as a compound itself), cross-attention, mask attention, multilevel attention, gated attention, spatial-semantic

attention, multi-dimensional gated attention, content-enhanced mask attention, and dot-product attention. These indicate specific computational mechanisms or types. More general modifiers like significant, considerable, sparse, adaptive, joint, hierarchical, temporal, spatial, global, and local also occur, describing properties *of the mechanism*.

Modifiers (General English): The English Trends Word Sketch shows modifiers related to human cognition, perception, and social interaction: close, medical, immediate, particular, urgent, special, careful, unwanted, national, undivided, widespread, international, meticulous, negative, media, rapt, public, due, prompt, heightened. The focus is on the kind or source of human/public focus.

Modified Nouns (AI Corpus): Attention frequently modifies other nouns, forming crucial compound terms: attention mechanism (very high frequency), attention map, attention module, attention layer, attention matrix, attention head, attention network, attention weight, attention computation, attention score, attention distribution. These compounds refer to specific components or outputs of the attention process in AI models.

Modified Nouns (General English): Common modified nouns relate to human cognitive limits or social phenomena: attention span, attention disorder, attention grabber, attention seeker, attention deficit, attention mechanism (can overlap but less frequent/central), attention shift, attention economy.

Verbs with attention as Object (AI Corpus): Verbs often relate to incorporating or using the mechanism: leverage attention, utilize attention, incorporate attention, calculate attention, share attention, guide attention, introduce attention, propose attention (mechanisms). The collocation pay attention exists but seems less central than in general English, sometimes referring to the *model's* focus (e.g., “agent further considers underlying web-specific content... the existing VLN task that only pays attention to vision and instruction”).

Verbs with attention as Object (General English): The dominant pattern is verbs related to directing or receiving human focus: pay attention (extremely frequent), draw attention, attract attention, catch attention, turn attention, focus attention, garner attention, divert attention, grab attention, gain attention, capture attention, deserve attention, receive attention, get attention, command attention, seek attention, call attention.

3.3.4.3. Semantic Specialization and Meaning Shifts

This is a clear case of **terminologization** where a word with a well-established cognitive meaning is adopted and redefined for a specific technical purpose.

Dominant AI Sense: Attention in the AI corpus almost exclusively refers to a **computational mechanism**, particularly within neural networks (especially Transformers), that allows a model to dynamically weigh the importance of different parts of the input data (e.g., different words in a sentence, different parts of an image) when producing an output. It's a specific, mathematically defined technique involving queries, keys, and values, enabling models to focus on relevant information. Examples: “multi-dimensional gated attention unit”, “self-attention mechanism”, “content-enhanced mask attention learning scheme”, “Gated Attention Coding”, “leverages the multi-dimensional gated attention unit”.

Dominant General English Sense: In English Trends, attention overwhelmingly refers to the **cognitive process of selectively concentrating on one aspect of the environment while ignoring other things**, or the notice, interest, or consideration given to someone or something by people. Examples: “Pay attention during opening and closing hours”, “dispute was already drawing attention”, “seek immediate attention [medical]”, “attract significant attention”, “attracted the most public attention”, “enjoy this type of attention”, “played Christmas music to get the attention of passersby”.

Shift/Specialization: The AI sense is a **metaphorical extension** of the cognitive concept, but has become a distinct technical term. It borrows the idea of “focusing” but applies it to a computational weighting mechanism. The link to the original cognitive

meaning is often metaphorical rather than literal. Crucially, the *specific mechanism* denoted by attention in AI (self-attention, cross-attention, etc.) has no direct equivalent in the general cognitive sense. The technical meaning has become the default within the AI domain, while the general cognitive/social meanings are largely absent in the expert discourse samples.

3.3.4.4. Metaphorical Patterns in Expert AI Terminology

The use of attention in AI is deeply rooted in the **MODEL AS COGNITIVE AGENT** or **MODEL AS PERCEIVER** metaphor.

Core Mapping: The computational mechanism is conceptualized as the model “paying attention” to parts of its input, similar to how a human selectively focuses their cognitive resources. The weighting scores produced by the mechanism are framed as the degree of “attention” the model “gives” to different elements.

Related Metaphors:

- **INFORMATION AS LANDSCAPE/SPACE:** The model attends to different parts of the input (sequence, image), implying the input exists in a space that can be selectively focused upon.
- **MECHANISM AS MODULE/COMPONENT:** Terms like attention mechanism, attention module, and attention unit frame it as a distinct part of the larger model architecture.

In general English, attention is the cognitive process. In AI, attention *names* a computational mechanism *analogous* to that process. The metaphorical distance is key – AI experts understand it refers to a specific set of calculations, not subjective consciousness.

Summary for attention

Within the expert AI corpus, attention is overwhelmingly a technical noun referring to a specific computational mechanism used in neural networks to weigh the relevance of input features. It is characterized by highly specific technical modifiers (self-, cross-,

mask, gated) and forms core compounds like attention mechanism, attention map, and attention module. While derived metaphorically from the human cognitive process of selective focus, the AI sense is distinct and refers to a specific mathematical operation. This contrasts sharply with its general English usage, which predominantly relates to human cognitive focus or social notice. The term exemplifies terminologization through metaphorical extension and specialization within the AI domain, participating in the broader MODEL AS COGNITIVE AGENT metaphor.

3.3.5.1. Analysis of: hallucination (lemma)

The lemma hallucination appears with notable frequency (33 instances in the AI concordance sample) in the **AI_Expert_AAAl24** corpus, particularly given its specific technical application. Its usage is almost entirely confined to the AI domain, contrasting sharply with its established meaning in psychology and general language.

3.3.5.2. Collocational Profiles and Grammatical Patterns

Grammatical Role: Exclusively used as a Noun in the AI corpus samples.

Modifiers (AI Corpus): The Word Sketch and concordance lines show modifiers specifying the *type* or *context* of the AI hallucination: Factual hallucination, visual hallucination (inferred from LVLMs context), object hallucination. Modifiers also relate to its management: mitigating hallucination.

Modifiers (General English): The English Trends Word Sketch reveals modifiers strongly tied to the psychological/medical context: auditory, visual, hypnagogic, drug-induced, olfactory, tactile, vivid, terrifying, psychotic, feverish, sleep-related, grief-induced.

Modified Nouns (AI Corpus): Hallucination modifies other nouns to create specific concepts: hallucination detection, hallucination prevention, hallucination problem, hallucination rates, Domain Information Hallucination module.

Modified Nouns (General English): Similar patterns exist but relate to the psychological experience: hallucination delusion, hallucination paranoia, hallucination experience.

Verbs with hallucination as Object (AI Corpus): Verbs centre on managing or dealing with the phenomenon in AI models: mitigate hallucination (very common), reduce hallucination, prevent hallucination, detect hallucination, address hallucination, correct hallucination, suppress hallucination, facilitate hallucination (in the specific DHU context).

Verbs with hallucination as Object (General English): Verbs relate to experiencing or causing the psychological phenomenon: experience hallucination, induce hallucination, cause hallucination, suffer hallucination, trigger hallucination, provoke hallucination, treat hallucination.

Verbs with hallucination as Subject (AI Corpus): While less common in the noun form, the underlying concept implies the model hallucinates. One concordance line shows the gerund form acting agentively: hallucination elevating instability.

Verbs with hallucination as Subject (General English): The psychological state can be the subject: hallucination torments, hallucination haunts, hallucination plagues, hallucination characterizes, hallucination occurs, hallucination subsides.

3.3.5.3. Semantic Specialization and Meaning Shifts

This term represents a striking case of **neologistic terminologization through metaphor**.

Dominant AI Sense: In the AI corpus, hallucination refers to the tendency of **generative models** (especially LLMs and LVLMs) **to produce output that is nonsensical, factually incorrect, irrelevant to the input prompt**, or unfaithful to the provided source data (e.g., generating descriptions of objects not present in an image). It signifies a specific type of model failure where the output appears plausible but is fabricated or erroneous. Examples: “mitigating the hallucination of large language

models (LLMs)”, “factual hallucination generated by LLMs during its reasoning process”, “Detecting and Preventing Hallucinations in Large Vision Language Models”, “hallucinatory text in the form of non-existent objects”, “effectively reducing the risk of hallucination [in generated data]”.

Dominant General English Sense: In English Trends, hallucination overwhelmingly refers to a **sensory perception** (seeing, hearing, smelling, tasting, or feeling something) **that seems real but does not exist outside the mind**, typically caused by mental illness, neurological conditions, or drugs. Examples: “psychotic symptoms such as hallucinations and delusions”, “suffered vivid, violent hallucinations”, “drug-induced hallucinations”, “auditory hallucinations”.

Shift/Specialization: The AI usage is a direct **metaphorical borrowing** from the psychological sense. The model’s generation of fabricated/unfaithful output is likened to a human experiencing a perception without an external stimulus. This metaphorical usage has rapidly solidified into a standard technical term within AI to describe this specific failure mode of generative models. The original sensory perception meaning is entirely absent in the AI context. The AI sense is purely about fabricated *information* or *output*, not sensory experience.

3.3.5.4. Metaphorical Patterns in Expert AI Terminology

The use of hallucination is itself the primary metaphorical pattern.

Core Metaphor: GENERATIVE MODEL FAILURE IS (LIKE) A MENTAL/PERCEPTUAL DEVIATION:

- **Mapping:** The generation of incorrect/fabricated output by an AI model (the target domain) is conceptualized in terms of a human experiencing a hallucination (the source domain).
- **Implications:** This metaphor frames the model’s error not just as a technical glitch but implicitly likens it to a cognitive malfunction. It highlights the

apparent plausibility but *actual falsity* of the output, much like a real-seeming hallucination.

- **Entailments:** This leads to related concepts like detecting hallucination (identifying the error), preventing hallucination (designing models to avoid it), mitigating hallucination (reducing its effects), and treating it as a “problem” or “risk”.

In general English, hallucination is a psychological phenomenon. In AI, it names a type of model output error by analogy to that phenomenon. The metaphorical nature is central to its technical meaning in AI. While experts understand it refers to output errors, the evocative power of the metaphor strongly shapes discussions around model reliability and trustworthiness.

Summary for hallucination

The term hallucination in the expert AI corpus is a specialized noun referring to the generation of factually incorrect, nonsensical, or ungrounded output by generative AI models. Its collocational profile revolves around concepts of large models (LLMs, LVLMs), detection, prevention, and mitigation. This usage is a direct metaphorical borrowing from the term’s primary meaning in psychology (a false sensory perception). Within AI, it has rapidly become a standard technical term to denote a specific type of model failure, framed implicitly via the MODEL FAILURE AS MENTAL DEVIATION metaphor. The original psychological meaning is absent in the AI context, making this a clear example of neologistic terminologization through metaphor.

3.3.6.1. Analysis of: training (lemma)

The lemma train (most frequently appearing as the noun/gerund training in the AI corpus) is a fundamental term in both general English and the AI domain. However, its specific usage, collocates, and semantic focus differ significantly, marking it as a term with strong domain specialization in AI contexts.

3.3.6.2. Collocational Profiles and Grammatical Patterns

Grammatical Role: In the AI_Expert_AAAl24 corpus, training overwhelmingly functions as a Noun, often as the head of a Noun Phrase or as part of a compound. While the verb train exists (e.g., “training different types of triplets”, “train the masker”), the nominal form describing the *process* or *phase* is much more prominent in the abstract discourse.

Modifiers (AI Corpus): The Word Sketch for AI_Expert_AAAl24 highlights technical modifiers specifying the type or method of training: adversarial training, joint training, end-to-end training, collaborative training, local training, cross-domain training (inferred), pre-training (highly associated, see lemma analysis), fine-tuning (related concept). General modifiers like efficient, stable, additional, and effective also appear, describing the quality of the training process.

Modifiers (General English): In the English Trends corpus, modifiers are vastly more diverse and typically relate to human skill acquisition, professional development, or physical conditioning: spring training, strength training, vocational training, job training, basic training, hands-on training, teacher training, on-the-job training, weight training, military training, staff training, specialized training (can overlap), interval training, resistance training, safety training, pilot training, combat training, potty training. The contexts are overwhelmingly human or animal-centric.

Modified Nouns (AI Corpus): Training frequently modifies other nouns, denoting elements related to the training process: training data, training process, training phase, training sample, training cost, training efficiency, training set, training strategy, training stage, training scheme, training procedure, training example, training time, training dataset, training objective, training loss, training algorithm. These compounds pinpoint specific aspects of the computational training workflow.

Modified Nouns (General English): Nouns modified by training in English Trends typically refer to locations, events, materials, or schedules related to human training: training camp, training session, training ground, training course, training

programme/program, training exercise, training facility, training center, training manual, training schedule, training wheels, training partner.

Verbs with training as Object (AI Corpus): In the AI context, verbs often relate to initiating, enabling, or optimizing the computational process: accelerate training, set training (parameters), guide training, facilitate training, enable training, require training.

Verbs with training as Object (General English): Verbs relate to participating in, providing, or completing human training: undergo training, complete training, resume training, impart training, receive training, provide training, conduct training, begin training, attend training, require training, offer training.

Prepositional Phrases (AI Corpus): Phrases often specify the context or goal: training from scratch, training for keypoints description, training on data/datasets, training under the state-of-the-art network, training during the stage, training between domains.

3.3.6.3. Semantic Specialization and Meaning Shifts

Similar to learning, the term training undergoes significant semantic specialization in the AI domain.

Dominant AI Sense: Training in AI refers specifically to the **computational process of optimizing the parameters of a model (typically a machine learning model) by exposing it to data**. The goal is to minimize a loss function and improve the model's ability to perform a specific task (e.g., classification, generation). It is an automated, algorithmic procedure. Examples: “validating the effectiveness of the model design and pre-training strategies”, “naively selecting more triplets for training under the state-of-the-art network”, “enhance both the keypoint detector and descriptor training”, “regularize the model training”, “unseen observed during the training stage”.

Dominant General English Senses: The primary meaning in English Trends relates to the **process of teaching a person or animal a particular skill or type of behavior**,

often through practice and instruction over a period. It also refers to the **process of undertaking a course of exercise** and diet in preparation for a **sporting event**. Examples: “pilot training”, “next CNA training program”, “lost control during a training run”, “dedication to training”, “driver training program”, “Petroleum Training Institute”, “training new bereavement volunteers”, “vocational training”, “spring training”.

Shift/Specialization: The AI usage is a **metaphorical extension** of the sense of teaching or preparing someone/something for a task. The algorithm/model is implicitly framed as the entity being “trained”, and the data serves as the “instruction” or “practice material”. While the core idea of improving capability through a process is preserved, the nature of the process (computational optimization vs. human instruction/practice) and the entity being trained (algorithm vs. person/animal) are fundamentally different. Within AI, this specialized computational meaning has become the default.

3.3.6.4. Metaphorical Patterns in Expert AI Terminology

The use of training heavily relies on the **MODEL AS TRAINEE/STUDENT** or **MODEL AS ATHLETE** metaphor, closely related to the **MODEL AS LEARNER** metaphor associated with learning.

Core Mapping: The process of optimizing model parameters using data is conceptualized as “training” the model, analogous to how a human or animal is trained for a specific skill or task.

Metaphorical Entailments:

- **Data as Curriculum/Exercise:** The training data or training set functions as the material used for instruction or practice.
- **Process Stages:** Concepts like pre-training and fine-tuning map onto ideas of foundational learning, followed by specialized refinement.

- **Effort/Cost:** Terms like training cost, training time, and training efficiency frame the computational process in terms of resource expenditure, similar to human training efforts.
- **Optimization as Guidance:** Training strategy, training objective, guide training frame the optimization process as a directed effort to shape the model's capabilities.
- **Adversarial Training:** This term specifically evokes the metaphor of training against an opponent or challenging conditions to improve resilience.

Contrast with General English: While general English uses training for instructing humans/animals or physical conditioning, the AI application maps these concepts onto an automated, algorithmic process of parameter optimization driven by data exposure. The metaphorical framing allows experts to talk about this complex computational process using familiar concepts of instruction, practice, and preparation.

Summary for training:

In the expert AI corpus, training primarily functions as a noun referring to the computational process of optimizing a model's parameters using data. Its collocational profile is dominated by technical terms specifying the type of training (adversarial, joint, end-to-end) or elements of the process (data, cost, phase, strategy, set, sample). This contrasts with its general English usage centered on human/animal skill acquisition or physical preparation. The AI sense represents a specific metaphorical extension (MODEL AS TRAINEE/ATHLETE), framing algorithmic optimization in terms of teaching, practice, and preparation. This technical meaning is paramount within the AI discourse.

3.3.7.1. Analysis of: bias (lemma)

The lemma bias occurs frequently in both the AI_Expert_AAAl24 corpus (360 instances in the concordance sample) and the English Trends general corpus (over 1.2 million instances), but its specific applications and collocations reveal distinct patterns and a strong terminological sense within the field of AI.

3.3.7.2. Collocational Profiles and Grammatical Patterns

Grammatical Role: Primarily used as a Noun in both corpora. The adjective biased also appears frequently in both.

Modifiers (AI Corpus): Modifiers in the AI_Expert_AAAI24 corpus often specify the source or type of bias within the AI/ML context: inductive bias, dataset bias, distribution bias, label bias, domain bias, visual bias, semantic bias, prediction bias, systematic bias, estimation bias. Some overlap with general modifiers exists but is less frequent (e.g., potential bias, strong bias).

Modifiers (General English): The English Trends Word Sketch shows a very wide range of modifiers, often related to social categories, attitudes, or statistical contexts: unconscious bias, implicit bias, racial bias, confirmation bias, gender bias, cognitive bias, inherent bias, systemic bias, liberal bias, political bias, negativity bias, optimism bias, selection bias, publication bias, survivorship bias. While some statistical concepts overlap, the focus on social and cognitive biases is much stronger in general English.

Modified Nouns (AI Corpus): Bias modifies other nouns, typically referring to management techniques: bias mitigation, bias metric, bias learning (less common).

Modified Nouns (General English): Frequent compounds include bias training (often anti-bias training), bias incident, bias crime, bias binding, bias ply (tire), bias tape (sewing), bias correction, bias audit, bias circuit. The range covers social issues, technical (non-AI) contexts, and statistics.

Verbs with bias as Object (AI Corpus): Verbs strongly relate to *addressing or measuring* bias in AI systems: mitigate bias, eliminate bias, alleviate bias, reduce bias, measure bias, evaluate bias, address bias, incorporate bias (e.g., inductive bias), introduce bias, amplify bias, cause bias, exhibit bias.

Verbs with bias as Object (General English): Verbs cover a broader range, including experiencing, perceiving, or introducing bias: perceive bias, perpetuate bias, eliminate bias, harbor bias, minimize bias, reinforce bias, correct bias, combat bias, reflect bias,

counteract bias, avoid bias, address bias, expose bias, show bias, confirm bias, have bias.

Verbs with bias as Subject (AI Corpus): Less frequent for the noun form, but the concept implies the model or data exhibits bias: bias exists, bias causes.

Verbs with bias as Subject (General English): Similar to AI, but also includes verbs reflecting its pervasive nature in human contexts: bias taints, bias motivates, bias creeps in, bias influences, bias clouds, bias exists, bias skews, bias distorts, bias affects.

Prepositional Phrases: bias towards X (e.g., bias towards classes, bias towards seen objects) is a key pattern indicating the direction of the skew. Bias in Y (e.g., bias in dataset, bias in model, bias in RL) specifies the location.

3.3.7.3. Semantic Specialization and Meaning Shifts

The term bias has multiple related senses, but its prominence and specific focus differ between the domains.

Dominant AI Sense: In AI, bias most frequently refers to systematic error or unfairness introduced into a model or its outputs due to various factors. This includes:

- **Data Bias:** Skews present in the training data that reflect societal biases or unrepresentative sampling (e.g., gender bias, racial bias in facial recognition datasets).
- **Algorithmic Bias:** Systematic errors introduced by the model architecture or learning algorithm itself (e.g., favoring certain outcomes).
- **Inductive Bias:** The set of assumptions a learning algorithm uses to make predictions on unseen data. This is a more neutral technical sense, but still refers to inherent model properties influencing outcomes.
- **Estimation Bias:** Statistical bias in parameter estimates.
- **Examples:** “mitigate the prediction bias”, “hidden bias issue”, “human bias”, “visual bias caused by compositions”, “potential space biases introduced by semantics”.

Dominant General English Senses: While the statistical sense exists, the most frequent meanings in English Trends relate to:

- **Cognitive Bias:** Systematic patterns of deviation from norm or rationality in judgment (e.g., confirmation bias, negativity bias, unconscious bias).
- **Social Bias:** Prejudice for or against one person or group, especially in a way considered to be unfair (e.g., racial bias, gender bias, political bias, anti-Israel bias).
- **Statistical Bias:** Systematic error in estimation.
- **Technical Bias:** In electronics, a bias voltage. In materials, a diagonal cut/direction (bias tape, bias ply).

Shift/Specialization: The AI usage heavily borrows from both the statistical sense (systematic error) and the social sense (unfairness related to groups). The critical aspect in AI is that **cognitive/social biases present in human data or design choices can become encoded as statistical biases in the model's behavior**, leading to unfair or inaccurate outcomes. While inductive bias is a more neutral technical term, much of the discourse around bias in AI explicitly concerns the negative consequences of mirroring societal prejudices or leading to systemic errors. The electronic or material senses are irrelevant in the AI context.

3.3.7.4. Metaphorical Patterns in Expert AI Terminology

While not always strictly metaphorical, the discourse around bias often employs frames related to deviation and contamination:

BIAS AS SKEW/TILT: The core idea is a deviation from a neutral or fair baseline. Phrases like bias towards classes or verbs like skew and tilt reflect this spatial metaphor.

BIAS AS CONTAMINATION/FLAW: Verbs like taint, creep in, infect, permeate, corrupt, and flaw frame bias as an undesirable element that degrades the model's

integrity or fairness. Addressing bias involves mitigation, elimination, correction, rooting out, weeding out.

BIAS AS INHERITED TRAIT: The idea that models inherit bias from data connects to biological metaphors.

Summary for bias

In the expert AI corpus, bias primarily refers to systematic error or unfairness in models or data, encompassing data-driven reflections of societal prejudice, algorithmic tendencies, and statistical properties like inductive bias. Its collocational profile includes modifiers specifying the type of bias (inductive, dataset, visual) and verbs focused on its mitigation (mitigate, eliminate, alleviate, reduce). While drawing heavily on both the statistical and social meanings from general English, the AI context uniquely fuses these, focusing on how data/algorithmic biases lead to unfair or inaccurate outcomes, often mirroring real-world prejudices. Metaphorically, bias is often framed as a skew, a contamination, or an inherited flaw.

CONCLUSION TO CHAPTER 3

This chapter embarked on a corpus-based characterization of modern expert Artificial Intelligence terminology, drawing empirical evidence from a specialized corpus compiled from AAAI-24 conference abstracts. Through a multi-faceted analysis encompassing thematic structure, linguistic form, domain specificity, and detailed semantic-contextual usage, several key characteristics of this specialized vocabulary have been identified.

Firstly, the **thematic analysis (Section 3.1.1)** revealed a strong concentration of high-frequency terminology within core methodological areas. The discourse, as reflected in the abstracts, is overwhelmingly focused on **Model Structures & Components** (e.g., neural network, transformer, GNN), **Learning Paradigms** (e.g., reinforcement learning, contrastive learning, self-supervised learning), **Optimization & Training** (e.g., algorithm, fine-tuning, loss function), and **Data & Feature Engineering** (e.g.,

embedding, representation, data augmentation). While crucial areas like Evaluation, Applications, Generative Models, and Challenges & Issues are certainly present, their associated high-frequency terminology appears less voluminous in this specific snapshot of expert communication, suggesting a primary focus on the foundational “how-to” of building and refining AI systems in this venue.

Secondly, the **structural analysis (Section 3.1.2)** highlighted specific linguistic forms favoured in AI terminology. The lexicon is predominantly nominal, built upon Noun Phrases. Conceptual precision is largely achieved through two complementary strategies: the extensive use of **Multi-Word Terms (MWTs)**, particularly compounds following recurring patterns like N+N (e.g., machine learning, knowledge distillation) and Adj+N (e.g., neural network, deep learning), and the significant reliance on **Acronyms** (e.g., LLM, GNN, RL, SSL, OOD). The identification of 60 unique acronym types among the high-frequency terms underscores their critical role in achieving communicative efficiency within the expert community.

Thirdly, the **keyword analysis comparing the AI corpus against general English (Section 3.2)** statistically confirmed the highly specialized nature of the expert lexicon. The distinctiveness of AI terminology stems from a combination of highly technical neologisms and acronyms often absent in general language (e.g., contrastive, gnn, ood, diffusion model), the assignment of specific, narrowed technical senses to common words (e.g., semantic, generalization, embedding, adversarial, robustness, attention, model), and the prevalence of characteristic compounding patterns (e.g., state-of-the-art method, benchmark dataset, point cloud). These domain-specific keywords strongly align with the core thematic areas identified earlier, reinforcing the link between the field’s conceptual focus and its unique linguistic encoding.

Finally, the **detailed semantic and contextual analysis of selected key terms (Section 3.3)** — model, learning, attention, hallucination, training, and bias — provided deeper insights into their specific behaviour. A consistent pattern observed was **semantic specialization**, where terms with broader meanings in general English acquire precise, often computational or algorithmic, definitions within AI (e.g., model

as computational construct, attention as weighting mechanism, learning/training as algorithmic optimization, bias as systemic error/unfairness). Furthermore, **metaphorical extension** proved to be a vital mechanism for term creation and conceptualization (e.g., learning based on MODEL AS LEARNER, attention on MODEL AS PERCEIVER, hallucination on FAILURE AS MENTAL DEVIATION, bias on BIAS AS SKEW/CONTAMINATION). These metaphors structure understanding but refer to specific technical phenomena. Distinct **collocational profiles** and **grammatical patterns** (e.g., the prevalence of technical modifiers for learning, the agentive use of model, the specific verbs associated with mitigating bias or hallucination) further differentiate the expert usage in the field of AI from general language patterns.

This chapter has demonstrated through corpus-based evidence that modern expert AI terminology, as represented in the AAI-24 corpus, is a dynamic and highly specialized system. It is characterized by its focus on core methodologies, its reliance on nominal structures, multi-word compounding, and acronyms for precision and efficiency, its distinctiveness from general English vocabulary through neologism and semantic specialization, and its conceptual richness derived from systematic metaphorical extensions. These findings provide an empirically grounded understanding of the linguistic features shaping communication at the forefront of Artificial Intelligence research.

CONCLUSION

This Master's thesis aimed to identify, characterize, and analyze key features of modern expert AI terminology. Using an empirical, corpus-based method, the study focused on abstracts from the 2024 AAAI Conference. The collected corpus was analyzed with automated techniques and corpus linguistic tools in Sketch Engine, and compared to a general English reference corpus to determine domain specificity.

The investigation yielded several key findings regarding the nature of expert AI terminology:

1. **Thematic Concentration:** The analysis confirmed a strong thematic focus within the expert discourse on the foundational elements of AI systems. Terminology related to **Model Structures & Components** (e.g., neural network, transformer, GNN), **Learning Paradigms** (e.g., reinforcement learning, contrastive learning), **Optimization & Training** procedures (e.g., algorithm, fine-tuning, loss), and **Data & Feature Engineering** (e.g., embedding, representation) dominated the high-frequency lexicon of the corpus. While other crucial areas like Evaluation, Applications, and Challenges were represented, the core terminological landscape highlighted a primary preoccupation with the methodologies for building and refining AI models.
2. **Dominant Structural Patterns:** Expert AI terminology relies heavily on nominal structures. While single-word nouns form a base, conceptual precision is largely achieved through the highly productive formation of **Multi-Word Terms (MWTs)**, predominantly Noun Phrases constructed using recurring patterns like N+N (e.g., machine learning) and Adj+N (e.g., deep learning). Concurrently, **Acronyms** (e.g., LLM, GNN, RL, OOD) were found to be exceptionally prevalent (approx. 60 unique types among high-frequency terms), serving as a critical mechanism for communicative efficiency and compression of complex MWTs within the expert community.
3. **High Degree of Domain Specificity:** The comparative keyword analysis demonstrated that the expert AI lexicon is markedly distinct from general

contemporary English. This specificity arises from a blend of: (a) **highly technical neologisms** and acronyms largely unique to the domain (e.g., GNN, OOD, diffusion model), and (b) the **semantic specialization** of words common in general English, which acquire precise, often computational or algorithmic meanings within AI (e.g., model, network, attention, learning, training, bias, embedding, robustness).

4. **Key Semantic Mechanisms:** The detailed analysis of selected terms (model, learning, attention, hallucination, training, bias) revealed characteristic semantic behaviours. **Semantic specialization and narrowing** were consistently observed, restricting broader general meanings to specific technical contexts (e.g., attention as a computational mechanism). **Metaphorical extension** proved a vital process for terminologization, mapping concepts from cognitive science, engineering, or even psychiatry onto computational phenomena (e.g., learning, attention, hallucination, training). These terms exhibited distinct **collocational profiles** and **grammatical patterns** within the AI corpus compared to general English, further solidifying their specialized status.

Contributions and Implications:

This research contributes to several fields:

- **Terminology Studies & LSP Research:** It provides an empirically grounded characterization of terminology dynamics—including neologism, semantic specialization, metaphor, acronymization, and compounding—within a rapidly evolving, high-impact scientific domain. It demonstrates the utility of corpus linguistic methods, enhanced by preliminary automated analysis, for investigating contemporary specialized language use aligned with descriptive theoretical frameworks like CTT.
- **AI Understanding and Communication:** The findings offer insights into how AI experts structure their conceptual world and communicate key ideas. This characterization can aid newcomers, educators, technical writers, and translators

in navigating the complexities of AI discourse, potentially mitigating communication barriers within and beyond the field.

- **Lexicography:** The study provides valuable data reflecting *current expert usage*, which can inform the creation or updating of specialized AI dictionaries, glossaries, and terminological databases, ensuring they accurately represent the contemporary state of the field's language.

Limitations of the Study:

While providing valuable insights, this study has several limitations:

- **Corpus Scope:** The analysis was based solely on abstracts from a single conference (AAAI-24). Abstracts possess specific genre conventions and may not fully represent the terminological range or nuances found in full research papers, textbooks, patents, or spoken expert communication.
- **Synchronic Focus:** The study provides a snapshot of terminology at a specific point in time (early 2024). It does not capture the diachronic evolution of terms, although the methodology could be adapted for such analysis.
- **Limited Scope of Phase 3:** The detailed semantic and contextual analysis (Phase 3) was performed only on a limited subset of key terms. While indicative of broader patterns like specialization and metaphor, the specific findings for these terms may not universally apply to all AI vocabulary.
- **Methodological Dependencies:** The results are contingent on the specific corpus compiled, the reference corpus used (English Trends), the analytical tools (Sketch Engine, initial processing via Gemini API), and the chosen statistical measures (LogDice). Different choices could yield slightly different results.

In conclusion, this thesis has provided a systematic, data-driven profile of modern expert Artificial Intelligence terminology based on a contemporary corpus. By analyzing thematic concentrations, structural preferences (including MWTs and acronyms), domain specificity, and the semantic behaviour of key terms, the study

offers valuable insights into the linguistic characteristics shaping communication at the forefront of this critical and rapidly advancing scientific field.

List of References

1. Клименко, Н. (2017). Колокації, терміни та їхня прагматика в суспільно-політичних і науково-публіцистичних текстах. Рідне слово в етнокультурному вимірі, 84–117. http://nbuv.gov.ua/UJRN/rsev_2017_2017_10
2. Клименко, Н. (2009). Термінування і детермінування в процесах інтелектуалізації сучасної української мови. *Studia Linguistica*, (III), 100–108. <https://studia-linguistica.knu.ua/2009-3-100-107-klimenko-n-f-terminuvannja-i-determinuvannja-v-procesah-intelektualizacii-suchasnoi-ukrainskoi-movi/>
3. Кочан, І. (2018). Українське термінознавство сьогодні. *Studia Philologica*, (11), 93–100. <https://studiap.kubg.edu.ua/index.php/journal/article/download/235/219>
4. Кочан, І. (2023). Особливості опрацювання українських терміносистем на сучасному етапі. *Науковий вісник Ужгородського університету. Серія: Філологія*, 2(50), 257–263. <https://dspace.uzhnu.edu.ua/jspui/bitstream/lib/58430/1/%D0%9E%D0%A1%D0%9E%D0%91%D0%9B%D0%98%D0%92%D0%9E%D0%A1%D0%A2%D0%86%20%D0%9E%D0%9F%D0%A0%D0%90%D0%A6%D0%AE%D0%92%D0%90%D0%9D%D0%9D%D0%AF%20%D0%A3%D0%9A%D0%A0%D0%90%D0%87%D0%9D%D0%A1%D0%AC%D0%9A%D0%98%D0%A5.pdf>
5. Куньч, З. (2023). Варіантність термінів (процеси термінування та детермінування). *Термінологічний вісник*, (7), 84–91. <https://termvisnyk.iul-nasu.org.ua/wp-content/uploads/sites/11/2023/06/7.pdf>
6. Петрова, Т. (2021). Розвиток теорії терміна в українській та зарубіжних лінгвістичних традиціях. *Studia Warمیńskie*, 58, 83–98. <https://wuwr.pl/swr/article/download/11836/10761>
7. Туровська, Л. (2019). *Термінологічний вісник: Збірник наукових праць* (Вип. 5). Київ: Інститут української мови НАНУ. https://iul-nasu.org.ua/pdf/termvisnik/terv_2019_5.pdf
8. Тодор, О. (2007). Детермінологізація. *Українська мова: Енциклопедія* (3-є вид., с. 139–140). Київ.

9. Томіленко, Л. (2015). Термінологічна лексика в сучасній тлумачній лексикографії української літературної мови. Івано-Франківськ: Фоліант.
- 10.AAAI-24 Technical Tracks 1
 AAAI Press. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(1). <https://ojs.aaai.org/index.php/AAAI/issue/view/576>
- 11.AAAI-24 Technical Tracks 2
 AAAI Press. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(2). <https://ojs.aaai.org/index.php/AAAI/issue/view/577>
- 12.AAAI-24 Technical Tracks 3
 AAAI Press. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(3). <https://ojs.aaai.org/index.php/AAAI/issue/view/578>
- 13.AAAI-24 Technical Tracks 4
 AAAI Press. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(4). <https://ojs.aaai.org/index.php/AAAI/issue/view/579>
- 14.AAAI-24 Technical Tracks 5
 AAAI Press. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(5). <https://ojs.aaai.org/index.php/AAAI/issue/view/580>
- 15.AAAI-24 Technical Tracks 6
 AAAI Press. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(6). <https://ojs.aaai.org/index.php/AAAI/issue/view/581>
- 16.AAAI-24 Technical Tracks 7
 AAAI Press. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(7). <https://ojs.aaai.org/index.php/AAAI/issue/view/582>
- 17.AAAI-24 Technical Tracks 8
 AAAI Press. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(8). <https://ojs.aaai.org/index.php/AAAI/issue/view/583>

18.AAAI-24	Technical	Tracks	9
AAAI Press. (2024). <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 38(9). https://ojs.aaai.org/index.php/AAAI/issue/view/584			
19.AAAI-24	Technical	Tracks	10
AAAI Press. (2024). <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 38(10). https://ojs.aaai.org/index.php/AAAI/issue/view/585			
20.AAAI-24	Technical	Tracks	11
AAAI Press. (2024). <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 38(11). https://ojs.aaai.org/index.php/AAAI/issue/view/586			
21.AAAI-24	Technical	Tracks	12
AAAI Press. (2024). <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 38(12). https://ojs.aaai.org/index.php/AAAI/issue/view/587			
22.AAAI-24	Technical	Tracks	13
AAAI Press. (2024). <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 38(13). https://ojs.aaai.org/index.php/AAAI/issue/view/588			
23.AAAI-24	Technical	Tracks	14
AAAI Press. (2024). <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 38(14). https://ojs.aaai.org/index.php/AAAI/issue/view/589			
24.AAAI-24	Technical	Tracks	15
AAAI Press. (2024). <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 38(15). https://ojs.aaai.org/index.php/AAAI/issue/view/590			
25.AAAI-24	Technical	Tracks	16
AAAI Press. (2024). <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 38(16). https://ojs.aaai.org/index.php/AAAI/issue/view/591			

- 26.AAAI-24 Technical Tracks 17
 AAAI Press. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(17). <https://ojs.aaai.org/index.php/AAAI/issue/view/592>
- 27.AAAI-24 Technical Tracks 18
 AAAI Press. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(18). <https://ojs.aaai.org/index.php/AAAI/issue/view/593>
- 28.AAAI-24 Special Track Safe, Robust and Responsible AI Track
 AAAI Press. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(19). <https://ojs.aaai.org/index.php/AAAI/issue/view/594>
- 29.AAAI-24 Special Track AI for Social Impact, Senior Member Presentations, New Faculty Highlights, Journal Track
 AAAI Press. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(20). <https://ojs.aaai.org/index.php/AAAI/issue/view/595>
- 30.Abduvokhidova, H. (2023). Terminology as the basis of linguistics. *European Journal of Research Development and Sustainability*, 4(07), 99–101. <https://scholarzest.com/index.php/ejrds/article/download/3754/2998/6809>
- 31.Ahmad, S. (2025). The vocabulary of thinking machines: A linguistic inquiry into AI terminology. *International Journal of Innovative Research in Technology*, 11(10), 2858–2884. https://ijirt.org/publishedpaper/IJIRT174038_PAPER.pdf
- 32.Allanazarova, M. A. (2020). Basic concepts and principles of cognitive linguistics. *The American Journal of Social Science and Education Innovations*, 2(09), 399–404. <https://doi.org/10.37547/tajssei/Volume02Issue09-61>
- 33.Altuncu, E., Nurse, J. R. C., Xu, Y., Guo, J., & Li, S. (2025). Improving performance of automatic keyword extraction (AKE) methods using PoS-tagging and enhanced semantic-awareness. *arXiv preprint arXiv:2211.05031*. <https://arxiv.org/pdf/2211.05031.pdf>

34. Association for the Advancement of Artificial Intelligence. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 38). <https://aaai.org/proceeding/aaai-38-2024/>
35. Beliga, S. (2014). Keyword extraction: A review of methods and approaches. University of Rijeka, Department of Informatics. https://langnet.uniri.hr/papers/beliga/Beliga_KeywordExtraction_a_review_of_methods_and_approaches.pdf
36. Bianchini, E., Climie, R. E., Mayer, C. C., Martina, M. R., Nandi, M., Schmidt-Trucksäss, A., Segers, P., Park, C., Pucci, G., Terentes-Printzios, D., & Charlton, P. H. (2024). Unified Language for Knowledge Dissemination: The Vascular Ageing Glossary, an Initiative by VascAgeNet. *Artery Research*, 30(Suppl 1), 1. <https://doi.org/10.1007/s44200-023-00041-5>
37. British National Corpus. (n.d.). *About the BNC*. Oxford University Computing Services. <http://www.natcorp.ox.ac.uk/>
38. Bullock, O. M., Colón Amill, D., Shulman, H. C., & Dixon, G. N. (2019). Jargon as a barrier to effective science communication: Evidence from metacognition. *Public Understanding of Science*, 28(7), 845–853. <https://comm.osu.edu/sites/comm.osu.edu/files/PUS%202019-%20Bullock%20et%20al..pdf>
39. Cabré Castellví, M. T. (2003). Theories of terminology: Their description, prescription and explanation. *Terminology. International Journal of Theoretical and Applied Issues in Specialized Communication*, 9(2), 163–199. John Benjamins Publishing Company.
40. Cabré, M. T. (1999). *Terminology: Theory, methods, and applications* (J. C. Sager, Ed.; J. A. DeCesaris, Trans.). John Benjamins Publishing Company. <http://www.dawsonera.com/depp/reader/protected/external/AbstractView/S9789027298652>

41. Cabré, M. T. (2023). *Terminology: Cognition, language and communication*. John Benjamins Publishing Company. <https://benjamins.com/catalog/ivitra.36>
42. Cambridge University Press. (n.d.). In Cambridge Dictionary. Retrieved May 1, 2025, from <https://dictionary.cambridge.org/>
43. Curry, N., & McEnery, T. (2025). Corpus linguistics for language teaching and learning: A research agenda. *Language Teaching*, 1–20. doi:10.1017/S0261444824000430
44. Dalieva, M. (2024). Diachronic corpora and language evolution over time. *International Journal of Language and Literary Studies*, 2(10), 58–60. <https://webofjournals.com/index.php/1/article/download/1892/1875/3711>
45. Davies, M. (2008–). *The Corpus of Contemporary American English (COCA)*. Available online at <https://www.english-corpora.org/coca/>
46. Delipetrev, B., Tsinaraki, C., & Kostić, U. (2020). *Historical evolution of artificial intelligence* (EUR 30221 EN). Publications Office of the European Union. <https://doi.org/10.2760/801580>
47. Dictionary.com, LLC. (n.d.). In Dictionary.com. Retrieved May 1, 2025, from <https://www.dictionary.com/>
48. Gablasova, D., Brezina, V. and McEnery, T. (2017), Collocations in Corpus-Based Language Learning Research: Identifying, Comparing, and Interpreting the Evidence. *Language Learning*, 67: 155-179. <https://doi.org/10.1111/lang.12225>
49. Gilquin, G., & Marchi, A. (2024). Concordancing for CADS: Practical challenges and theoretical implications. *International Journal of Corpus Linguistics*, 29(1), 34–58. <https://www.jbe-platform.com/content/journals/10.1075/ijcl.21168.gil>
50. Haddad Haddad, A., Rigouts Terryn, A., & Mitkov, R. (Eds.). (2023). *Computational Terminology in NLP and Translation Studies (ConTeNTS): Proceedings of the First ConTeNTS Workshop and the 16th BUCC Workshop, co-located with the 14th*

International Conference on Recent Advances in Natural Language Processing (RANLP 2023), Varna, Bulgaria, 7 September 2023 (Online ISBN 978-954-452-090-8). INCOMA Ltd. <https://aclanthology.org/2023.contents-1.pdf>

51. Hagen, L., & Balahur, A. (2023). Using machine learning to explain document characteristics. *Frontiers in Artificial Intelligence*, 5, 975729. <https://doi.org/10.3389/frai.2022.975729>
52. HarperCollins Publishers. (n.d.). In Collins English Dictionary. Retrieved May 1, 2025, from <https://www.collinsdictionary.com/>
53. IAAI-24, EAAI-24, AAAI-24 Student Abstracts, Undergraduate Consortium and Demonstrations
AAAI Press. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(21). <https://ojs.aaai.org/index.php/AAAI/issue/view/596>
54. Ibrahim, L., & Cheng, M. (2025). Thinking beyond the anthropomorphic paradigm benefits LLM research. *ArXiv*. <https://arxiv.org/abs/2502.09192>
55. Incelli, E. . (2025). Exploring the Future of Corpus Linguistics: Innovations in AI and Social Impact. *International Journal of Mass Communication*, 3, 1–10. <https://doi.org/10.6000/2818-3401.2025.03.01>
56. Incelli, E. (2025). Exploring the future of corpus linguistics: Innovations in AI and social impact. *International Journal of Mass Communication*, 3, 1–10. <https://doi.org/10.6000/2818-3401.2025.03.01>
57. Incelli, E. (2025). Exploring the future of corpus linguistics: Innovations in AI and social impact. *International Journal of Mass Communication*, 3(1), 1–10. <https://www.lifescienceglobal.com/pms/index.php/IJMC/article/download/10111/5215/24025>
58. Jakubíček, M., & Rychlý, P. (2019). A distributional multi-word thesaurus in Sketch Engine. In A. Horák, P. Rychlý, & A. Rambousek (Eds.), *Proceedings of Recent*

- Advances in Slavonic Natural Language Processing (RASLAN 2019)* (pp. 143–147).
Tribun EU. https://www.sketchengine.eu/wp-content/uploads/Distributional_multi-word_thesaurus.pdf
59. Jurafsky, D., & Martin, J. H. (2025). N-gram language models. In *Speech and language processing* (3rd ed. draft, Chapter 3). Stanford University. <https://web.stanford.edu/~jurafsky/slp3/3.pdf>
60. Kalaš, F. (2025). Bridging tradition and innovation: Analysing language data with ChatGPT-4 in corpus linguistics. *SSRN Electronic Journal*. <https://ssrn.com/abstract=5126316>
61. Khan, S. A. (2016). The Distinction between Term and Word: A Translator and Interpreter Problem and the Role of Teaching Terminology. *Procedia - Social and Behavioral Sciences*, 232, 696-704. <https://doi.org/10.1016/j.sbspro.2016.10.095>
62. Krishnan, N. (2025). AI Agents: Evolution, Architecture, and Real-World Applications. *ArXiv*. <https://arxiv.org/abs/2503.12687>
63. Kūlis, M. (2023). Lost in Translation: Artificial Intelligence and the Burden of Bad Metaphors. *PhilArchive*. Institute of Philosophy and Sociology, University of Latvia. <https://philarchive.org/archive/KLILIT>
64. Kyröläinen, A., & Laippala, V. (2023). Predictive keywords: Using machine learning to explain document characteristics. *Frontiers in Artificial Intelligence*, 5, 975729. <https://doi.org/10.3389/frai.2022.975729>
65. Laktionova, A. (2021). Compiling a specialised corpus for translation research in the environmental domain. In *Proceedings of the RANLP 2021 Student Research Workshop* (pp. 94–98). INCOMA Ltd. https://doi.org/10.26615/issn.2603-2821.2021_014

- 66.Lattanzi, M. (2016). Corpora and culture: Theoretical and methodological issues. *Culture e Corpora*, 4(1), 7–26. <http://sibaese.unisalento.it/index.php/culturecorpora/article/download/12434/11073>
- 67.Lenko, V. S., Pasichnyk, V. V., & Shcherbyna, Y. M. (2017). Knowledge representation models. *Computer Science and Information Technologies*, 864, 158–168. Lviv Polytechnic National University. <https://science.lpnu.ua/sites/default/files/journal-paper/2018/jul/13788/21.pdf>
- 68.Lexical Computing CZ s.r.o. (2024). *Sketch Engine* [Computer software]. <https://www.sketchengine.eu/>
- 69.Lexical Computing CZ s.r.o. (2025). *English Trends corpus*. Sketch Engine. <https://www.sketchengine.eu/english-trends-corpus/>
- 70.Lexical Computing CZ s.r.o. (2025). *English Trends corpus*. Sketch Engine. <https://www.sketchengine.eu/english-trends-corpus/>
- 71.Liu, D., & Lei, L. (2020). Technical vocabulary. In S. Webb (Ed.), *The Routledge handbook of vocabulary studies* (pp. 111–124). Routledge. <https://doi.org/10.4324/9780429291586-8>
- 72.Lucy, L., Dodge, J., Bamman, D., & Keith, K. A. (2022). Words as Gatekeepers: Measuring Discipline-specific Terms and Meanings in Scholarly Publications. *ArXiv*. <https://arxiv.org/abs/2212.09676>
- 73.Macmillan Education. (n.d.). In Macmillan Dictionary. Retrieved May 1, 2025, from <https://www.macmillandictionary.com/>
- 74.Madina Anafinova. (2021). Corpus based research in terminology. Habarşy - Ál-Farabi atyndagy Qazaq Memlekettik ultiq Universiteti. Filologiya Seriâsy. https://www.academia.edu/116361940/Corpus_based_research_in_terminology

75. Martínez, A., & Mammola, S. (2021). Specialized terminology reduces the number of citations of scientific papers. *Proceedings of the Royal Society B: Biological Sciences*, 288(1948), 20202581. <https://doi.org/10.1098/rspb.2020.2581>
76. McEnery, T., Xiao, R., & Tono, Y. (2006). Unit 11: Corpus representativeness and balance. In *Corpus-based language studies: An advanced resource book* (pp. [insert page numbers if available]). Routledge. <https://www.lancaster.ac.uk/fass/projects/corpus/ZJU/xCBLS/chapters/B01.pdf>
77. Merriam-Webster. (n.d.). In Merriam-Webster.com dictionary. Retrieved May 1, 2025, from <https://www.merriam-webster.com/>
78. Misnawati, M., Nur, S., & Tahir, S. Z. (2024). Corpus linguistics today: A qualitative approach. *Research and Innovation in Applied Linguistics (RIAL)*, 2(1), 45–62. <https://doi.org/10.31963/rial.v2i1.4486>
79. Mitchell, M. (2024). The metaphors of artificial intelligence. *Science*. <https://doi.org/adt6140>
80. Oxford University Press. (n.d.). In Oxford English Dictionary. Retrieved May 1, 2025, from <https://www.oed.com/>
81. Pearson Education. (n.d.). In Longman Dictionary of Contemporary English Online. Retrieved May 1, 2025, from <https://www.ldoceonline.com/>
82. Rehak, R. (2021). The Language Labyrinth: Constructive Critique on the Terminology Used in the AI Discourse. *ArXiv*, *abs/2307.10292*. <https://www.semanticscholar.org/paper/The-Language-Labyrinth%3A-Constructive-Critique-on-in-Rehak/6a0f29bdcd68c6cb0f0b1be569968d0d8e8265dd>
83. Rehak, R. (2021). The language labyrinth: Constructive critique on the terminology used in the AI discourse. In P. Verdegem (Ed.), *AI for everyone? Critical perspectives* (pp. 87–102). University of Westminster Press. <https://doi.org/10.16997/book55.f>

84. Shestakevych, T., Shyika, Y., & Tsiokh, L. (2024). Quantitative characteristics of the author's idiostyle. In *Proceedings of the Computational Linguistics Workshop at the 8th International Conference on Computational Linguistics and Intelligent Systems (CoLInS-2024)*. CEUR Workshop Proceedings, Vol-3722. <http://ceur-ws.org/Vol-3722/paper26.pdf>
85. Shvets, A., Mohammadshahi, A., Henderson, J., & Wanner, L. (2022). Multilingual Extraction and Categorization of Lexical Collocations with Graph-aware Transformers. *ArXiv*. <https://arxiv.org/abs/2205.11456>
86. Sketch Engine. (2024). *Word sketch – collocations and word combinations*. <https://www.sketchengine.eu/guide/word-sketch-collocations-and-word-combinations/>
87. Stefanowitsch, A. (2020). *Corpus linguistics: A guide to the methodology*. Language Science Press. <https://library.oapen.org/handle/20.500.12657/43768>
88. Suresh, H., & Guttag, J. V. (2021). A framework for understanding sources of harm throughout the machine learning life cycle. arXiv preprint *arXiv:1901.10002*. <https://arxiv.org/abs/1901.10002>
89. Terminology in the age of AI: The transformation of terminology theory and practice. (2024). *Journal of Translation Studies*, 4(1), Article 5. <https://www.ingentaconnect.com/content/plg/jts/2024/00000004/00000001/art00005>
90. Trojar, M. (2017). Wüster's view of terminology. Inštitut za slovenski jezik Frana Ramovša ZRC SAZU, Ljubljana. https://www.academia.edu/34943245/W%C3%BCsters_View_of_Terminology
91. Volkova, O. (2021). Corpus-based ESP learning and teaching. In *1st Ukrainian Conference on Applied Linguistics: Corpora and Discourse* (pp. 130–133). National Technical University of Ukraine “Igor Sikorsky Kyiv Polytechnic Institute”. <https://corpora.kamts1.kpi.ua/cad-2021/paper/download/25165/13900>

92. Wissik, T. (2010). Development of comparable specialized corpora of national German varieties: The UNI-Corpus. In R. Xiao, S. Bernardini, & K. Wang (Eds.), *Using Corpora in Contrastive and Translation Studies (UCCTS 2010) Proceedings* (pp. 9–18). Lancaster University.
<https://www.lancaster.ac.uk/fass/projects/corpus/UCCTS2010Proceedings/papers/Wissik.pdf>
93. Yang, J., & Long, C. (2020). Common and distinctive cognitive processes between categorization and category-based induction: Evidence from event-related potentials. *Brain Research*, 1749, 147134. <https://doi.org/10.1016/j.brainres.2020.147134>

Appendix

List of the 550 terms analyzed in the research, sorted by frequency

lemma	frequency	capitalization	category
machine learning	444	Machine Learning	Learning Paradigms
model	290	Model	Model Structures & Components
algorithm	195	Algorithm	Optimization & Training
reinforcement learning	164	Reinforcement Learning	Learning Paradigms
neural network	161	Neural Network	Model Structures & Components
large language models	148	Large Language Models	Generative Models & Techniques
deep learning	143	Deep Learning	Learning Paradigms
transformer	121	Transformer	Model Structures & Components
models	114	Models	Model Structures & Components
algorithms	112	Algorithms	Optimization & Training

llms	112	LLMs	Generative Models & Techniques
learning	103	Learning	Learning Paradigms
contrastive learning	103	Contrastive Learning	Learning Paradigms
training	96	Training	Optimization & Training
classification	94	Classification	Evaluation & Analysis
neural networks	89	Neural Networks	Model Structures & Components
network	88	Network	Model Structures & Components
graph neural networks	77	Graph Neural Networks	Model Structures & Components
rl	76	RL	Learning Paradigms
ai	75	AI	Not Related
deep neural networks	66	Deep Neural Networks	Model Structures & Components
diffusion models	66	Diffusion Models	Generative Models & Techniques

generalization	65	Generalization	Evaluation & Analysis
representation learning	64	Representation Learning	Learning Paradigms
federated learning	63	Federated Learning	Learning Paradigms
clustering	55	Clustering	Evaluation & Analysis
artificial intelligence	55	Artificial Intelligence	Other AI Concepts
diffusion model	54	Diffusion Model	Generative Models & Techniques
gnns	53	GNNs	Model Structures & Components
knowledge distillation	51	Knowledge Distillation	Optimization & Training
natural language processing	48	Natural Language Processing	Domains & Applications
object detection	47	Object Detection	Domains & Applications
agent	47	Agent	Other AI Concepts
policy	47	Policy	Other AI Concepts
optimization	47	Optimization	Optimization & Training

classifier	46	Classifier	Evaluation & Analysis
adversarial attacks	46	Adversarial Attacks	Challenges & Issues
encoder	45	Encoder	Model Structures & Components
attention	45	Attention	Model Structures & Components
features	43	Features	Data & Feature Engineering
computer vision	43	Computer Vision	Domains & Applications
self-attention	41	Self-Attention	Model Structures & Components
fine-tuning	41	Fine-Tuning	Optimization & Training
clip	40	Clip	Model Structures & Components
pre-training	39	Pre-Training	Optimization & Training
inference	39	Inference	Evaluation & Analysis

generative models	38	Generative Models	Generative Models & Techniques
data augmentation	38	Data Augmentation	Data & Feature Engineering
language models	38	Language Models	Model Structures & Components
self-supervised learning	37	Self-Supervised Learning	Learning Paradigms
zero-shot	36	Zero-Shot	Learning Paradigms
transfer learning	36	Transfer Learning	Learning Paradigms
semi-supervised learning	35	Semi-Supervised Learning	Learning Paradigms
image classification	34	Image Classification	Domains & Applications
chatgpt	34	Chatgpt	Generative Models & Techniques
semantic segmentation	34	Semantic Segmentation	Evaluation & Analysis
self-supervised	32	Self-Supervised	Learning Paradigms

attention mechanism	30	Attention Mechanism	Model Structures & Components
representations	30	Representations	Data & Feature Engineering
graph neural network	30	Graph Neural Network	Model Structures & Components
latent space	30	Latent Space	Data & Feature Engineering
gnn	30	GNN	Model Structures & Components
unsupervised	29	Unsupervised	Learning Paradigms
model training	29	Model Training	Optimization & Training
deep reinforcement learning	28	Deep Reinforcement Learning	Learning Paradigms
few-shot learning	28	Few-Shot Learning	Learning Paradigms
fl	28	FL	Other AI Concepts
adversarial training	28	Adversarial Training	Optimization & Training
segmentation	27	Segmentation	Evaluation & Analysis

adversarial examples	26	Adversarial Examples	Challenges & Issues
domain adaptation	26	Domain Adaptation	Learning Paradigms
embeddings	26	Embeddings	Data & Feature Engineering
ml	26	ML	Learning Paradigms
pre-trained models	25	Pre-Trained Models	Model Structures & Components
agents	25	Agents	Other AI Concepts
vision transformer	25	Vision Transformer	Model Structures & Components
node classification	24	Node Classification	Evaluation & Analysis
framework	24	Framework	Other AI Concepts
reasoning	24	Reasoning	Other AI Concepts
continual learning	24	Continual Learning	Learning Paradigms
prediction	24	Prediction	Evaluation & Analysis
knowledge transfer	24	Knowledge Transfer	Other AI Concepts

multi-modal	23	Multi-Modal	Other AI Concepts
cnn	23	CNN	Model Structures & Components
planning	23	Planning	Other AI Concepts
generative model	23	Generative Model	Generative Models & Techniques
overfitting	23	Overfitting	Challenges & Issues
cnn	23	CNNs	Model Structures & Components
classifiers	23	Classifiers	Model Structures & Components
embedding	22	Embedding	Data & Feature Engineering
nlp	22	NLP	Domains & Applications
contrastive loss	22	Contrastive Loss	Optimization & Training
decoder	22	Decoder	Model Structures & Components
neural radiance fields	22	Neural Radiance Fields	Model Structures & Components

llm	22	LLM	Generative Models & Techniques
transformers	22	Transformers	Model Structures & Components
meta-learning	22	Meta-Learning	Learning Paradigms
gpt-4	21	GPT-4	Generative Models & Techniques
large language model	21	Large Language Model	Generative Models & Techniques
convolutional neural networks	21	Convolutional Neural Networks	Model Structures & Components
loss function	21	Loss Function	Optimization & Training
multi-agent reinforcement learning	21	Multi-Agent Reinforcement Learning	Learning Paradigms
nerf	21	Nerf	Generative Models & Techniques
interpretability	21	Interpretability	Evaluation & Analysis
gan	21	GAN	Generative Models & Techniques

predictions	20	Predictions	Evaluation & Analysis
policies	20	Policies	Other AI Concepts
training data	20	Training Data	Data & Feature Engineering
autoencoder	20	Autoencoder	Model Structures & Components
few-shot	20	Few-Shot	Learning Paradigms
anomaly detection	20	Anomaly Detection	Evaluation & Analysis
dnns	20	DNNs	Model Structures & Components
multimodal	20	Multimodal	Data & Feature Engineering
domain generalization	19	Domain Generalization	Learning Paradigms
graph learning	19	Graph Learning	Learning Paradigms
generator	19	Generator	Model Structures & Components
catastrophic forgetting	19	Catastrophic Forgetting	Challenges & Issues

vision-language models	19	Vision-Language Models	Model Structures & Components
regression	19	Regression	Evaluation & Analysis
feature extraction	19	Feature Extraction	Data & Feature Engineering
loss	18	Loss	Optimization & Training
knowledge graphs	18	Knowledge Graphs	Data & Feature Engineering
robustness	18	Robustness	Evaluation & Analysis
pseudo-labels	18	Pseudo-Labels	Data & Feature Engineering
exploration	18	Exploration	Other AI Concepts
autonomous driving	18	Autonomous Driving	Domains & Applications
model-agnostic	18	Model-Agnostic	Other AI Concepts
model performance	17	Model Performance	Evaluation & Analysis
markov decision processes	17	Markov Decision Processes	Other AI Concepts

supervised learning	17	Supervised Learning	Learning Paradigms
marl	17	MARL	Learning Paradigms
feature extractor	17	Feature Extractor	Data & Feature Engineering
gradients	17	Gradients	Optimization & Training
prompt learning	17	Prompt Learning	Learning Paradigms
ssl	16	SSL	Learning Paradigms
gcn	16	GCN	Model Structures & Components
pre-trained language models	16	Pre-Trained Language Models	Model Structures & Components
regularization	16	Regularization	Optimization & Training
cross-attention	16	Cross-Attention	Model Structures & Components
bert	16	Bert	Model Structures & Components

distillation	16	Distillation	Optimization & Training
generative adversarial networks	15	Generative Adversarial Networks	Generative Models & Techniques
image generation	15	Image Generation	Domains & Applications
semantic	15	Semantic	Other AI Concepts
unsupervised domain adaptation	15	Unsupervised Domain Adaptation	Learning Paradigms
knowledge graph	15	Knowledge Graph	Data & Feature Engineering
language model	15	Language Model	Model Structures & Components
architecture	15	Architecture	Model Structures & Components
style transfer	15	Style Transfer	Domains & Applications
ai systems	15	AI Systems	Other AI Concepts
unsupervised learning	14	Unsupervised Learning	Learning Paradigms
prompt tuning	14	Prompt Tuning	Optimization & Training

deep neural network	14	Deep Neural Network	Model Structures & Components
feature space	14	Feature Space	Data & Feature Engineering
representation	14	Representation	Data & Feature Engineering
active learning	14	Active Learning	Learning Paradigms
message passing	14	Message Passing	Other AI Concepts
optimal transport	14	Optimal Transport	Optimization & Training
vit	14	ViT	Model Structures & Components
out-of-distribution	14	Out-Of-Distribution	Evaluation & Analysis
graph representation learning	14	Graph Representation Learning	Learning Paradigms
supervised	13	Supervised	Learning Paradigms
feature fusion	13	Feature Fusion	Data & Feature Engineering
generative	13	Generative	Generative Models & Techniques

3d object detection	13	3D Object Detection	Domains & Applications
gans	13	GANs	Generative Models & Techniques
prototypes	13	Prototypes	Other AI Concepts
learning-based	13	Learning-Based	Learning Paradigms
feature learning	13	Feature Learning	Data & Feature Engineering
spiking neural networks	13	Spiking Neural Networks	Model Structures & Components
bayesian optimization	13	Bayesian Optimization	Optimization & Training
multi-task learning	13	Multi-Task Learning	Learning Paradigms
denoising	13	Denoising	Other AI Concepts
pre-trained model	13	Pre-Trained Model	Model Structures & Components
visual question answering	13	Visual Question Answering	Domains & Applications
quantization	13	Quantization	Optimization & Training

sampling	13	Sampling	Data & Feature Engineering
imagenet	13	Imagenet	Data & Feature Engineering
embedding space	13	Embedding Space	Data & Feature Engineering
feature representation	13	Feature Representation	Data & Feature Engineering
mlp	12	MLP	Model Structures & Components
recommender systems	12	Recommender Systems	Domains & Applications
global model	12	Global Model	Model Structures & Components
feature representations	12	Feature Representations	Data & Feature Engineering
fairness	12	Fairness	Challenges & Issues
teacher model	12	Teacher Model	Model Structures & Components
student model	12	Student Model	Model Structures & Components

loss functions	12	Loss Functions	Optimization & Training
language modeling	12	Language Modeling	Domains & Applications
architectures	12	Architectures	Model Structures & Components
deep learning models	12	Deep Learning Models	Model Structures & Components
machine learning models	12	Machine Learning Models	Model Structures & Components
attention mechanisms	12	Attention Mechanisms	Model Structures & Components
stochastic	12	Stochastic	Optimization & Training
graph contrastive learning	12	Graph Contrastive Learning	Learning Paradigms
imitation learning	11	Imitation Learning	Learning Paradigms
vqa	11	VQA	Domains & Applications
black-box models	11	Black-Box Models	Model Structures & Components

pretrained models	11	Pretrained Models	Model Structures & Components
zero-shot learning	11	Zero-Shot Learning	Learning Paradigms
dnn	11	DNN	Model Structures & Components
detector	11	Detector	Model Structures & Components
decision trees	11	Decision Trees	Model Structures & Components
adversarial attack	11	Adversarial Attack	Challenges & Issues
node representations	11	Node Representations	Data & Feature Engineering
text-to-image	11	Text-To-Image	Generative Models & Techniques
cross-modal	11	Cross-Modal	Other AI Concepts
transformer architecture	11	Transformer Architecture	Model Structures & Components
multi-label classification	11	Multi-Label Classification	Evaluation & Analysis
multi-view clustering	11	Multi-View Clustering	Evaluation & Analysis

drl	11	DRL	Learning Paradigms
gaussian process	11	Gaussian Process	Model Structures & Components
self-training	11	Self-Training	Learning Paradigms
convergence	11	Convergence	Other AI Concepts
mutual information	11	Mutual Information	Data & Feature Engineering
pseudo labels	11	Pseudo Labels	Data & Feature Engineering
neural architecture search	10	Neural Architecture Search	Optimization & Training
transformer-based	10	Transformer-Based	Model Structures & Components
text-to-image models	10	Text-To-Image Models	Generative Models & Techniques
transformer models	10	Transformer Models	Model Structures & Components
graph convolutional network	10	Graph Convolutional Network	Model Structures & Components
prompts	10	Prompts	Other AI Concepts

vlms	10	VLMs	Model Structures & Components
networks	10	Networks	Model Structures & Components
accuracy	10	Accuracy	Evaluation & Analysis
downstream tasks	10	Downstream Tasks	Evaluation & Analysis
ensemble	10	Ensemble	Model Structures & Components
adversarial robustness	10	Adversarial Robustness	Challenges & Issues
decision-making	10	Decision-Making	Other AI Concepts
gradient descent	10	Gradient Descent	Optimization & Training
backdoor attacks	10	Backdoor Attacks	Challenges & Issues
variational autoencoder	10	Variational Autoencoder	Model Structures & Components
feature maps	10	Feature Maps	Data & Feature Engineering
self-distillation	10	Self-Distillation	Optimization & Training

policy optimization	10	Policy Optimization	Optimization & Training
multi-armed bandit	10	Multi-Armed Bandit	Learning Paradigms
adversarial	10	Adversarial	Challenges & Issues
modeling	10	Modeling	Other AI Concepts
optimization problem	10	Optimization Problem	Optimization & Training
stable diffusion	10	Stable Diffusion	Generative Models & Techniques
feature	10	Feature	Data & Feature Engineering
snn	10	SNNs	Model Structures & Components
fine-tune	10	Fine-Tune	Optimization & Training
policy learning	10	Policy Learning	Learning Paradigms
link prediction	9	Link Prediction	Evaluation & Analysis
convolution	9	Convolution	Model Structures & Components

finetuning	9	Finetuning	Optimization & Training
incremental learning	9	Incremental Learning	Learning Paradigms
detectors	9	Detectors	Model Structures & Components
attention module	9	Attention Module	Model Structures & Components
graph convolutional networks	9	Graph Convolutional Networks	Model Structures & Components
adversarial perturbations	9	Adversarial Perturbations	Challenges & Issues
convolutional neural network	9	Convolutional Neural Network	Model Structures & Components
ood detection	9	OOD Detection	Evaluation & Analysis
ai concepts	9	AI Concepts	Other AI Concepts
noisy labels	9	Noisy Labels	Data & Feature Engineering
optimization objective	9	Optimization Objective	Optimization & Training
cifar-10	9	CIFAR-10	Data & Feature Engineering

instance segmentation	9	Instance Segmentation	Evaluation & Analysis
state-of-the-art	9	State-Of-The-Art	Other AI Concepts
learning framework	9	Learning Framework	Learning Paradigms
vision transformers	9	Vision Transformers	Model Structures & Components
question answering	9	Question Answering	Domains & Applications
sam	9	SAM	Domains & Applications
machine learning algorithms	9	Machine Learning Algorithms	Optimization & Training
learning algorithm	9	Learning Algorithm	Optimization & Training
multi-agent systems	9	Multi-Agent Systems	Domains & Applications
pretraining	9	Pretraining	Optimization & Training
foundation models	9	Foundation Models	Model Structures & Components
machine translation	9	Machine Translation	Domains & Applications

text generation	9	Text Generation	Generative Models & Techniques
named entity recognition	9	Named Entity Recognition	Domains & Applications
chain-of-thought	9	Chain-Of-Thought	Generative Models & Techniques
backpropagation	9	Backpropagation	Optimization & Training
generative ai	9	Generative AI	Generative Models & Techniques
lms	9	LMS	Generative Models & Techniques
image captioning	8	Image Captioning	Domains & Applications
beam search	8	Beam Search	Optimization & Training
resnet	8	ResNet	Model Structures & Components
bi-level optimization	8	Bi-Level Optimization	Optimization & Training
visual features	8	Visual Features	Data & Feature Engineering

distributed learning	8	Distributed Learning	Learning Paradigms
adapters	8	Adapters	Model Structures & Components
transformer-based models	8	Transformer-Based Models	Model Structures & Components
ablation studies	8	Ablation Studies	Evaluation & Analysis
attention maps	8	Attention Maps	Model Structures & Components
neural representation	8	Neural Representation	Model Structures & Components
graph convolution	8	Graph Convolution	Model Structures & Components
in-context learning	8	In-Context Learning	Learning Paradigms
data distribution	8	Data Distribution	Data & Feature Engineering
pre-trained	8	Pre-Trained	Other AI Concepts
graph	8	Graph	Other AI Concepts
large language models (llms)	8	Large Language Models (LLMs)	Generative Models & Techniques

surrogate model	8	Surrogate Model	Model Structures & Components
generalization ability	8	Generalization Ability	Evaluation & Analysis
network architecture	8	Network Architecture	Model Structures & Components
online learning	8	Online Learning	Learning Paradigms
recommendation	8	Recommendation	Evaluation & Analysis
cot	8	CoT	Other AI Concepts
datasets	8	Datasets	Data & Feature Engineering
variational inference	8	Variational Inference	Optimization & Training
nn	8	NN	Model Structures & Components
benchmarks	8	Benchmarks	Other AI Concepts
encoding	8	Encoding	Data & Feature Engineering
transferability	8	Transferability	Evaluation & Analysis

encoders	8	Encoders	Model Structures & Components
gpt-3	8	GPT-3	Generative Models & Techniques
feature selection	8	Feature Selection	Data & Feature Engineering
gradient	8	Gradient	Optimization & Training
self-attention mechanism	8	Self-Attention Mechanism	Model Structures & Components
explainability	7	Explainability	Evaluation & Analysis
mnist	7	MNIST	Data & Feature Engineering
natural language understanding	7	Natural Language Understanding	Domains & Applications
artificial neural networks	7	Artificial Neural Networks	Model Structures & Components
text prompts	7	Text Prompts	Data & Feature Engineering
value function	7	Value Function	Other AI Concepts
multi-objective optimization	7	Multi-Objective Optimization	Optimization & Training

continuous control	7	Continuous Control	Domains & Applications
bayesian inference	7	Bayesian Inference	Other AI Concepts
prompt	7	Prompt	Other AI Concepts
multimodal learning	7	Multimodal Learning	Learning Paradigms
plms	7	PLMs	Generative Models & Techniques
autoregressive	7	Autoregressive	Model Structures & Components
debiasing	7	Debiasing	Challenges & Issues
prompting	7	Prompting	Other AI Concepts
partial label learning	7	Partial Label Learning	Learning Paradigms
relation extraction	7	Relation Extraction	Domains & Applications
feature alignment	7	Feature Alignment	Data & Feature Engineering
model generalization	7	Model Generalization	Evaluation & Analysis
benchmark datasets	7	Benchmark Datasets	Data & Feature Engineering

neural	7	Neural	Model Structures & Components
ner	7	NER	Domains & Applications
gcns	7	GCNs	Model Structures & Components
deep models	7	Deep Models	Model Structures & Components
classification accuracy	7	Classification Accuracy	Evaluation & Analysis
offline reinforcement learning	7	Offline Reinforcement Learning	Learning Paradigms
super-resolution	7	Super-Resolution	Domains & Applications
domain shift	7	Domain Shift	Challenges & Issues
convergence rate	7	Convergence Rate	Optimization & Training
explainable ai	7	Explainable AI	Evaluation & Analysis
algorithm names	7	Algorithm Names	Optimization & Training

policy network	7	Policy Network	Model Structures & Components
mdps	7	MDPs	Other AI Concepts
model parameters	7	Model Parameters	Model Structures & Components
image-to-image translation	7	Image-To-Image Translation	Domains & Applications
diffusion	7	Diffusion	Generative Models & Techniques
discriminator	7	Discriminator	Model Structures & Components
dataset	7	Dataset	Data & Feature Engineering
model optimization	7	Model Optimization	Optimization & Training
ai-driven	7	AI-driven	Other AI Concepts
fine-tuned	7	Fine-Tuned	Optimization & Training
semantic information	7	Semantic Information	Data & Feature Engineering
message-passing	7	Message-Passing	Model Structures & Components
end-to-end	7	End-To-End	Other AI Concepts

image encoder	7	Image Encoder	Model Structures & Components
text-to-image generation	7	Text-To-Image Generation	Generative Models & Techniques
xai	7	XAI	Evaluation & Analysis
generation	7	Generation	Generative Models & Techniques
regret bound	7	Regret Bound	Evaluation & Analysis
recurrent neural networks	7	Recurrent Neural Networks	Model Structures & Components
point cloud	7	Point Cloud	Data & Feature Engineering
markov decision process	7	Markov Decision Process	Other AI Concepts
ai planning	6	AI Planning	Other AI Concepts
representation space	6	Representation Space	Data & Feature Engineering
aggregation	6	Aggregation	Other AI Concepts
ai system	6	AI System	Other AI Concepts
kernel	6	Kernel	Model Structures & Components

recommendation systems	6	Recommendation Systems	Domains & Applications
gpt-3.5	6	GPT-3.5	Generative Models & Techniques
model accuracy	6	Model Accuracy	Evaluation & Analysis
vision-language	6	Vision-Language	Domains & Applications
gcl	6	GCL	Learning Paradigms
knowledge graph completion	6	Knowledge Graph Completion	Data & Feature Engineering
auto-encoder	6	Auto-Encoder	Model Structures & Components
backdoor attack	6	Backdoor Attack	Challenges & Issues
label noise	6	Label Noise	Challenges & Issues
symbolic reasoning	6	Symbolic Reasoning	Other AI Concepts
pareto front	6	Pareto Front	Optimization & Training
bias	6	Bias	Challenges & Issues

algorithmic fairness	6	Algorithmic Fairness	Challenges & Issues
model robustness	6	Model Robustness	Evaluation & Analysis
sentiment analysis	6	Sentiment Analysis	Domains & Applications
mdp	6	MDP	Other AI Concepts
autonomous systems	6	Autonomous Systems	Domains & Applications
predictor	6	Predictor	Model Structures & Components
causal inference	6	Causal Inference	Other AI Concepts
pll	6	PLL	Not Related
variational autoencoders	6	Variational Autoencoders	Model Structures & Components
deep generative models	6	Deep Generative Models	Generative Models & Techniques
image features	6	Image Features	Data & Feature Engineering
ddpm	6	DDPM	Generative Models & Techniques
foundation model	6	Foundation Model	Model Structures & Components

deep-learning	6	Deep-Learning	Learning Paradigms
out-of-distribution detection	6	Out-Of-Distribution Detection	Evaluation & Analysis
object localization	6	Object Localization	Domains & Applications
explainable artificial intelligence	6	Explainable Artificial Intelligence	Evaluation & Analysis
segmentation model	6	Segmentation Model	Model Structures & Components
ai-based	6	AI-Based	Other AI Concepts
neural model	6	Neural Model	Model Structures & Components
adversarial learning	6	Adversarial Learning	Learning Paradigms
semi-supervised	6	Semi-Supervised	Learning Paradigms
train	6	Train	Optimization & Training
learning paradigm	6	Learning Paradigm	Learning Paradigms
module	6	Module	Model Structures & Components

encoder-decoder	6	Encoder-Decoder	Model Structures & Components
sample efficiency	6	Sample Efficiency	Optimization & Training
network architectures	6	Network Architectures	Model Structures & Components
model predictions	6	Model Predictions	Evaluation & Analysis
weakly supervised	6	Weakly Supervised	Learning Paradigms
latent spaces	6	Latent Spaces	Data & Feature Engineering
ood	6	OOD	Evaluation & Analysis
stochastic gradient descent	6	Stochastic Gradient Descent	Optimization & Training
generative adversarial network	6	Generative Adversarial Network	Generative Models & Techniques
inductive bias	6	Inductive Bias	Other AI Concepts
natural language	6	Natural Language	Not Related
adaptation	6	Adaptation	Other AI Concepts
augmentation	6	Augmentation	Data & Feature Engineering

segment anything model	6	Segment Anything Model	Model Structures & Components
backdoor	6	Backdoor	Challenges & Issues
vae	6	VAE	Model Structures & Components
semantics	6	Semantics	Other AI Concepts
network parameters	6	Network Parameters	Model Structures & Components
graph data	6	Graph Data	Data & Feature Engineering
semantic space	6	Semantic Space	Data & Feature Engineering
pseudo-labeling	6	Pseudo-Labeling	Learning Paradigms
recurrent neural network	5	Recurrent Neural Network	Model Structures & Components
rendering	5	Rendering	Other AI Concepts
causal discovery	5	Causal Discovery	Other AI Concepts
state-of-the-art models	5	State-Of-The-Art Models	Other AI Concepts
benchmark	5	Benchmark	Other AI Concepts

neurons	5	Neurons	Model Structures & Components
offline rl	5	Offline RL	Learning Paradigms
detr	5	DETR	Model Structures & Components
prompt engineering	5	Prompt Engineering	Other AI Concepts
data heterogeneity	5	Data Heterogeneity	Data & Feature Engineering
pruning	5	Pruning	Optimization & Training
bayesian networks	5	Bayesian Networks	Model Structures & Components
black-box	5	Black-Box	Evaluation & Analysis
causality	5	Causality	Other AI Concepts
model architecture	5	Model Architecture	Model Structures & Components
object detectors	5	Object Detectors	Model Structures & Components
u-net	5	U-Net	Model Structures & Components

mujoco	5	MuJoCo	Domains & Applications
model learning	5	Model Learning	Learning Paradigms
open-set	5	Open-Set	Evaluation & Analysis
activation functions	5	Activation Functions	Model Structures & Components
anns	5	ANNs	Model Structures & Components
distribution shifts	5	Distribution Shifts	Challenges & Issues
non-autoregressive	5	Non-Autoregressive	Model Structures & Components
hallucinations	5	Hallucinations	Challenges & Issues
commonsense reasoning	5	Commonsense Reasoning	Other AI Concepts
encoder-decoder architecture	5	Encoder-Decoder Architecture	Model Structures & Components
trainable parameters	5	Trainable Parameters	Optimization & Training

learning models	5	Learning Models	Learning Paradigms
learn	5	Learn	Learning Paradigms
natural language generation	5	Natural Language Generation	Generative Models & Techniques
over-fitting	5	Over-Fitting	Challenges & Issues
negative sampling	5	Negative Sampling	Optimization & Training
rnn	5	RNN	Model Structures & Components
graph encoder	5	Graph Encoder	Model Structures & Components
mcts	5	MCTs	Other AI Concepts
classification model	5	Classification Model	Model Structures & Components
reward function	5	Reward Function	Optimization & Training
decoding	5	Decoding	Model Structures & Components
text classification	5	Text Classification	Evaluation & Analysis

thompson sampling	5	Thompson Sampling	Optimization & Training
monte carlo tree search	5	Monte Carlo Tree Search	Other AI Concepts
rewards	5	Rewards	Other AI Concepts
adam	5	Adam	Optimization & Training
model compression	5	Model Compression	Optimization & Training
knowledge	5	Knowledge	Other AI Concepts
search space	5	Search Space	Optimization & Training
image synthesis	5	Image Synthesis	Generative Models & Techniques
point clouds	5	Point Clouds	Data & Feature Engineering
nns	5	NNs	Model Structures & Components
local search	5	Local Search	Optimization & Training
sparsity	5	Sparsity	Other AI Concepts
local models	5	Local Models	Model Structures & Components

text encoder	5	Text Encoder	Model Structures & Components
consistency regularization	5	Consistency Regularization	Optimization & Training
meta-reinforcement learning	5	Meta-Reinforcement Learning	Learning Paradigms
contrastive language-image pre-training	5	Contrastive Language-Image Pre-Training	Learning Paradigms
classification models	5	Classification Models	Model Structures & Components
probabilistic model	5	Probabilistic Model	Model Structures & Components
semantic features	5	Semantic Features	Data & Feature Engineering
binary classification	5	Binary Classification	Evaluation & Analysis
sota	5	SOTA	Other AI Concepts
diffusion probabilistic model	5	Diffusion Probabilistic Model	Generative Models & Techniques
dropout	5	Dropout	Optimization & Training
neural radiance field	5	Neural Radiance Field	Model Structures & Components

detection	5	Detection	Evaluation & Analysis
differentiable	5	Differentiable	Optimization & Training
instruction-following	5	Instruction-Following	Other AI Concepts
evolutionary algorithm	5	Evolutionary Algorithm	Optimization & Training
image recognition	5	Image Recognition	Domains & Applications
unet	5	UNET	Model Structures & Components
feature aggregation	5	Feature Aggregation	Data & Feature Engineering
cifar-100	5	CIFAR-100	Data & Feature Engineering
hallucination	5	Hallucination	Challenges & Issues
multi-layer perceptron	5	Multi-Layer Perceptron	Model Structures & Components
mixup	5	Mixup	Data & Feature Engineering
peft	5	Peft	Optimization & Training

causal graph	5	Causal Graph	Data & Feature Engineering
zero-shot generalization	5	Zero-Shot Generalization	Evaluation & Analysis
data-driven	5	Data-Driven	Other AI Concepts
class-incremental learning	5	Class-Incremental Learning	Learning Paradigms
surrogate models	5	Surrogate Models	Model Structures & Components
softmax	5	Softmax	Model Structures & Components
retraining	5	Retraining	Optimization & Training
denoising diffusion probabilistic models	5	Denoising Diffusion Probabilistic Models	Generative Models & Techniques
parameter-efficient fine-tuning	5	Parameter-Efficient Fine-Tuning	Optimization & Training
uncertainty	5	Uncertainty	Evaluation & Analysis
feature disentanglement	5	Feature Disentanglement	Data & Feature Engineering
generative modeling	5	Generative Modeling	Generative Models & Techniques

heuristic	5	Heuristic	Other AI Concepts
lidar	5	LIDAR	Domains & Applications
learning algorithms	5	Learning Algorithms	Optimization & Training
distribution	5	Distribution	Data & Feature Engineering
graph-based	5	Graph-Based	Data & Feature Engineering
over-smoothing	5	Over-Smoothing	Challenges & Issues
feature matching	5	Feature Matching	Data & Feature Engineering
convolutional network	5	Convolutional Network	Model Structures & Components
state-of-the-art methods	5	State-Of-The-Art Methods	Other AI Concepts
clustering algorithms	5	Clustering Algorithms	Evaluation & Analysis
knowledge graph embedding	5	Knowledge Graph Embedding	Data & Feature Engineering
regret	5	Regret	Other AI Concepts

approximation algorithms	5	Approximation Algorithms	Optimization & Training
outlier detection	5	Outlier Detection	Evaluation & Analysis
reward	5	Reward	Other AI Concepts
rlhf	5	RLHF	Learning Paradigms
reward functions	5	Reward Functions	Optimization & Training
bayesian	5	Bayesian	Other AI Concepts
recognition	5	Recognition	Evaluation & Analysis
latent representation	5	Latent Representation	Data & Feature Engineering
graph attention networks	5	Graph Attention Networks	Model Structures & Components
kgs	5	KGs	Data & Feature Engineering
self-supervised methods	5	Self-Supervised Methods	Learning Paradigms
k-means clustering	5	K-Means Clustering	Evaluation & Analysis
disentanglement	5	Disentanglement	Other AI Concepts

node representation	5	Node Representation	Data & Feature Engineering
generative language models	5	Generative Language Models	Generative Models & Techniques
latent diffusion model	5	Latent Diffusion Model	Generative Models & Techniques
vision-language pre-training	5	Vision-Language Pre-Training	Learning Paradigms
optimization algorithm	5	Optimization Algorithm	Optimization & Training
cross-entropy loss	5	Cross-Entropy Loss	Optimization & Training
neural models	5	Neural Models	Model Structures & Components
image translation	5	Image Translation	Domains & Applications
filters	5	Filters	Model Structures & Components
autonomous agents	5	Autonomous Agents	Other AI Concepts

Summary

Магістерська робота присвячена ідентифікації, характеристиці та аналізу ключових особливостей сучасної експертної термінології сфери штучного інтелекту (ШІ) на матеріалі англійської мови. Зважаючи на стрімку еволюцію цієї галузі, дослідження ґрунтується на емпіричному, корпусному підході.

Об'єктом дослідження є термінологія, що використовується в сучасному експертному дискурсі ШІ, як вона представлена в провідних наукових публікаціях. Предметом дослідження виступають лексико-семантичні, структурні, частотні та концептуальні характеристики цієї спеціалізованої лексики.

Методологія дослідження включала укладання спеціалізованого синхронного корпусу (анотації конференції AAAI-2024, ~544 тис. слів) та його багатоаспектний аналіз. Початкова обробка даних (екстракція, категоризація) здійснювалася за допомогою Python та Gemini API. Подальший детальний аналіз проводився на платформі Sketch Engine з використанням її інструментів (частотний аналіз, виявлення ключових слів у порівнянні з референтним корпусом English Trends, аналіз колокацій та граматичних патернів через Word Sketch, контекстуальний аналіз за допомогою конкордансу).

Основні результати дослідження дозволили виявити такі визначальні риси сучасної експертної термінології ШІ:

Тематична сфокусованість: Аналіз підтвердив значну концентрацію термінології навколо фундаментальних методологічних аспектів розробки ШІ-систем, зокрема, структур моделей та їх компонентів, парадигм навчання, процедур оптимізації та тренування, а також інженерії даних та ознак.

Переважаючі структурні патерни: Лексикон характеризується домінуванням номінативних одиниць. Концептуальна точність досягається за рахунок активного використання багатослівних термінів (БСТ), утворених

переважно за моделями N+N та Adj+N, та широкого застосування акронімів (напр., LLM, GNN, RL) як засобу ефективної комунікації.

Високий ступінь доменної специфічності: Порівняльний аналіз ключових слів статистично підтвердив унікальність експертного лексикону ШІ порівняно із загальноживаною англійською мовою. Ця специфічність виникає завдяки поєднанню технічних неологізмів та акронімів, відсутніх у загальному вжитку, та семантичної спеціалізації загальноживаних слів (model, attention, learning, bias), які набувають вузьких технічних значень.

Значущість семантичних механізмів: Детальний аналіз ключових термінів виявив послідовні процеси семантичної спеціалізації (звуження значення) та активне використання метафоричного перенесення для термінологізації понять з інших сфер (learning, hallucination, training, attention). Ці терміни також демонструють специфічні колокаційні та граматичні патерни в експертному дискурсі.

Загалом, дослідження представляє сучасну експертну термінологію ШІ як динамічну, високоспеціалізовану систему. Вона інтенсивно використовує механізми композиційності (багатослівні терміни), компресії (акроніми), семантичної спеціалізації та метафори для точного й ефективного позначення складних та швидкозмінних концепцій галузі.