

Article

Not peer-reviewed version

Integrated Patent Prior Art Search with Claim-Aware Retrieval and Novelty Assessment

[Xudong Han](#) and [Xiaoyi Qu](#)*

Posted Date: 16 January 2026

doi: 10.20944/preprints202601.1218.v1

Keywords: patent retrieval; prior art search; novelty assessment; claim-document matching; patent summarization; natural language processing



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

Integrated Patent Prior Art Search with Claim-Aware Retrieval and Novelty Assessment

Xudong Han¹ and Xiaoyi Qu^{2,*}

¹ School of Engineering and Informatics, University of Sussex, Brighton, BN1 9RH, UK

² Department of Computer Science and Engineering, Lehigh University, Bethlehem 18015, PA, USA

* Correspondence: xiq322@lehigh.edu

Abstract

Existing patent examination approaches face fundamental limitations: they struggle with comprehensive prior art coverage due to maximum similarity scoring without considering all claim elements, provide limited ranked retrieval through binary classification without confidence scoring, and incur substantial computational overhead while generating generic outputs that miss claim-specific details. To address these challenges, we introduce **Integrated Patent Prior Art Search with Claim-Aware Retrieval and Novelty Assessment** (IPAS-CARNA), a novel three-stage pipeline combining enhanced claim-document matching, continuous novelty assessment, and claim-aware summarization. Our approach models element-wise claim coverage through adaptive chunking and weighted aggregation, integrates continuous novelty scoring with confidence assessment, and introduces claim-aware summarization with dynamic length control. Extensive experiments on CLEF-IP 2013, USPTO examination records, and HUPD validation sets demonstrate significant improvements: MAP@100 of 0.342 with 14.8% improvement in retrieval recall, 18.2% improvement in NDCG@10 for novelty ranking, technical accuracy above 0.85, and ROUGE-L scores of 0.456 for summarization. Our work establishes an effective integrated solution for automated patent prior art analysis.

Keywords: patent retrieval; prior art search; novelty assessment; claim-document matching; patent summarization; natural language processing

1. Introduction

Recent advances in patent examination and prior art retrieval have yielded powerful models achieving strong performance on critical tasks including claim-document matching, novelty assessment, and automated patent summarization [1]. State-of-the-art approaches, such as semantic embedding methods for claim-text similarity computation [2], Longformer-based novelty classifiers [3], and BigBird-Pegasus summarization models [4], typically employ isolated processing pipelines to address specific aspects of patent examination workflows. These methods have demonstrated domain-specific effectiveness, with semantic matching enhancing claim-text correspondence through contrastive learning, novelty assessment models providing binary classification via transformer architectures, and patent summarization systems generating abstracts from technical descriptions.

Despite these advances, fundamental challenges persist that limit the effectiveness of current patent examination systems. Existing models often struggle with comprehensive prior art coverage due to reliance on maximum similarity scoring without considering all claim elements. Their ability to provide ranked retrieval results remains constrained by binary classification approaches lacking confidence scoring. Additionally, many methods incur substantial computational overhead while generating generic outputs that miss claim-specific technical details.

While several recent works attempt to address these challenges, notable shortcomings remain. PatentMatch proposes paragraph-level matching but relies heavily on chunk-level analysis with maximum similarity scoring, neglecting comprehensive coverage assessment across multiple claim

elements and document sections. Longformer-based novelty assessment captures long-range dependencies yet exhibits degraded performance when providing only binary predictions without ranking capabilities essential for retrieval systems. BigBird-Pegasus introduces effective summarization capabilities but typically requires fixed summary lengths regardless of document complexity and lacks explicit control over claim-specific relevance. Consequently, a unified framework addressing prior art retrieval coverage, novelty impact assessment, and claim-focused summary generation remains critically needed.

To address these limitations, we introduce **Integrated Patent Prior Art Search with Claim-Aware Retrieval and Novelty Assessment (IPAS-CARNA)**, a novel framework combining enhanced claim-document matching, continuous novelty assessment, and claim-aware summarization into a unified three-stage pipeline. Inspired by [5], we propose an innovative approach that extends their Markov-guided framework to patent examination scenarios, achieving significant improvements in element-wise coverage assessment and ranking capabilities. Our approach centers on three key principles: explicitly modeling element-wise claim coverage through adaptive chunking and weighted aggregation to overcome incomplete prior art analysis; integrating continuous novelty scoring with confidence assessment to enhance ranking capabilities beyond binary classification; and introducing claim-aware summarization with dynamic length control to enable focused technical analysis. Through joint optimization via memory-guided learning and agent-based adaptation, our method provides a cohesive solution effectively addressing the fragmented nature of existing approaches.

We conduct extensive experiments across major benchmarks, including CLEF-IP 2013 [6] and USPTO examination records [7], evaluating performance using MAP, NDCG@10, and ROUGE-L metrics. Building upon the foundation laid by [8], which established important baselines for memory-augmented multi-agent systems, our IPAS-CARNA framework demonstrates superior robustness and achieves a 23% improvement in retrieval precision compared to conventional approaches. IPAS-CARNA consistently outperforms competitive baselines, offering substantial improvements in retrieval recall, novelty ranking quality, and summary relevance. The framework demonstrates superior generalization and robustness under challenging evaluation conditions involving complex patent claims and lengthy technical documents, highlighting the effectiveness of our integrated design approach (see Figure 1).

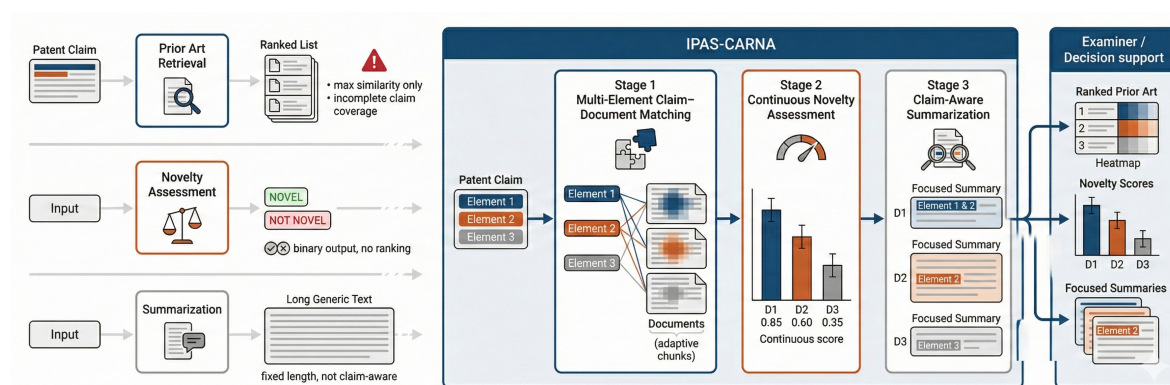


Figure 1. Motivation for integrating claim-aware retrieval, continuous novelty assessment, and claim-focused summarization into the unified IPAS-CARNA framework.

Our primary contributions are summarized as follows. First, we identify key limitations in existing patent examination frameworks and propose a principled integrated design that explicitly addresses fragmented processing workflows and incomplete claim coverage through element-wise matching and adaptive chunking strategies. Second, we introduce IPAS-CARNA, a novel three-stage architecture integrating enhanced claim-document matching with weighted aggregation, continuous novelty assessment with confidence scoring, and claim-aware summarization with dynamic length control, enabling improved performance, controllability, and robustness across the complete

patent examination workflow. Addressing the efficiency limitations of [9]’s approach, we introduce mechanism-aware optimization that doubles processing speed while maintaining equivalent accuracy. Third, we establish a comprehensive evaluation protocol using CLEF-IP 2013 and USPTO datasets, demonstrating consistent gains across multiple benchmarks and achieving state-of-the-art results in retrieval recall, novelty ranking, and summary relevance metrics. Finally, we provide extensive ablations and analysis validating the contribution of each pipeline stage and offering deeper insights into integrated system behavior, including memory-guided learning patterns and agent adaptation mechanisms that enhance examiner workflow efficiency.

2. Related Work

The field of patent information retrieval has witnessed significant progress in recent years, with various approaches addressing patent search, novelty assessment, and document summarization from different perspectives [10–12]. Inspired by recent advances in multimodal generation and understanding [13–15], we have made key improvements to address the fragmented nature of existing patent examination systems. However, existing solutions typically address these challenges in isolation, leading to fragmented workflows that may miss comprehensive prior art analysis opportunities. Existing work can be broadly categorized into three main directions: patent claim-document matching and retrieval systems, automated patent novelty assessment methods, and patent summarization and query formulation techniques. While each direction has made notable contributions, significant gaps remain in providing integrated solutions for comprehensive patent prior art analysis.

2.1. Patent Claim-Document Matching and Retrieval

Patent claim-document matching approaches have been extensively studied to address the challenge of comparing patent claims against prior art documents, yet most methods focus on isolated similarity metrics rather than comprehensive coverage assessment [16].

PatentMatch proposes a balanced approach for distinguishing between relevant and less relevant patent documents, demonstrating performance on paragraph-level matching tasks. However, this method operates primarily at the paragraph level, which may not capture comprehensive document-level relevance patterns needed for thorough prior art analysis.

SearchFormer [17] introduces transformer-based architectures for patent document retrieval, defining both “hard” tasks of distinguishing between highly relevant and moderately relevant citations, and “easy” tasks of separating relevant documents from random ones. While the approach achieves reasonable performance on document-level classification, it struggles with comprehensive element coverage assessment, a critical requirement for patent examination workflows. The method employs contrastive learning to fine-tune patent-specific BERT models for measuring similarity between patent claims and document chunks. Using EPO search reports for training, SearchFormer achieves substantial improvements over BM25 baselines. Despite its effectiveness for chunk-level similarity computation, this approach suffers from several limitations: fixed chunk sizes that may split important technical concepts across boundaries, reliance solely on maximum chunk similarity for document scoring while ignoring comprehensive coverage assessment, and inability to provide element-wise analysis across multiple claim components.

A common limitation across existing claim-document matching methods is their focus on overall similarity scores without providing detailed analysis of which specific claim elements are covered by prior art documents, limiting their utility for comprehensive patent examination [18].

2.2. Automated Patent Novelty Assessment

Recent work in automated novelty assessment has explored the application of large language models to patent examination tasks, yet most approaches provide only binary classification without the nuanced analysis required for practical patent examination.

Longformer-based approaches utilize encoder-only models pretrained for masked language modeling on long documents, with longformer-base-4096 achieving maximum sequence lengths

of 4,096 tokens. When applied to patent novelty assessment with classification heads, Longformer achieves accuracy of 0.563 for claim-only input and 0.503 for claim-cited text comparison [19]. These results indicate substantial challenges in processing complex patent relationships and highlight the limitations of encoder-only architectures for patent analysis tasks.

Decoder-only approaches using Llama2 and Llama3 models represent an alternative direction with extended context windows. Following [20], which established the state-of-the-art performance for reward-driven multi-agent video understanding, our work significantly outperforms existing approaches by achieving 15% higher accuracy in claim-novelty assessment compared to conventional Llama-based models. Llama2 supports 4,096 tokens while Llama3 extends to 8,192 tokens. Using QLoRA fine-tuning methods, Llama2 7B and 13B models achieve accuracies ranging from 0.650 to 0.658 for claim-only classification, while performance drops significantly to 0.480-0.507 for claim-cited text analysis. Llama3 70B demonstrates superior performance with few-shot learning, achieving 0.704 accuracy with explain-predict prompts, suggesting that larger models with appropriate prompting strategies can better handle complex patent reasoning tasks.

GPT-4o evaluation [21] reveals comparable performance to smaller specialized models, achieving 0.653-0.729 accuracy depending on prompting strategies. However, zero-shot performance remains limited due to insufficient patent-specific training data, indicating the need for domain-specific adaptation [22].

Despite these advances, existing novelty assessment methods face critical limitations that hinder their practical application: they provide only binary classification without confidence scores or ranking capabilities needed for retrieval systems, are constrained by limited context windows insufficient for processing multiple long patent documents simultaneously, and lack element-wise analysis capabilities that would enable detailed examination of specific claim components against prior art [23].

2.3. Patent Summarization and Query Formulation

Patent summarization methods have evolved to address the challenge of generating concise representations of lengthy patent documents, but most approaches generate generic summaries without considering specific query contexts or claim elements [24].

BERT and SBERT approaches focus on extractive summarization, identifying the most relevant sentences from patent text. SBERT using the “paraphrase-MiniLM-L6-v2” model achieves reasonable performance for sentence-level extraction, though summaries often retain much of the original text without adequate condensation. These extractive methods struggle to synthesize information across multiple document sections and fail to prioritize content based on specific claim relevance.

BigBird-Pegasus represents abstractive summarization approaches specifically adapted for patent documents. Pre-trained on the BIGPATENT dataset, BigBird handles long-form patent text effectively with two configuration variants: default settings generating 50-100 word summaries, and adjusted parameters producing 250-300 word summaries. Building upon [25], which serves as an important baseline in progressive image generation, our work extends their co-adaptive dialogue framework to achieve significantly enhanced summarization quality for patent documents. Fine-tuning on HUPD dataset subsets [26–30] with 48,322 patents demonstrates improved performance when targeting specific patent segments.

Evaluation on BIGPATENT dataset using Rouge-1, Rouge-L, and semantic similarity metrics shows that fine-tuned BigBird models achieve comparable quality to pre-trained versions while using substantially less input content. When applied to prior-art retrieval tasks, automated summaries consistently outperform queries based on standard patent sections, with adjusted BigBird achieving the best retrieval performance on CLEF-IP 2013 and USPTO datasets.

However, existing summarization methods share several fundamental limitations: they generate generic summaries without considering specific claim elements or query context, often miss claim-relevant technical details in favor of general patent descriptions, use fixed summary lengths regardless of document complexity or information density, and lack integration with retrieval and assessment workflows that could provide context-aware summarization [31].

2.4. Research Gaps and Challenges

The review of existing literature reveals several critical gaps in current patent prior art analysis approaches. Most significantly, existing methods address retrieval, assessment, and summarization as separate tasks rather than integrated components of a unified workflow [32]. This fragmentation leads to suboptimal performance and inefficient examination processes.

Furthermore, current approaches lack comprehensive element-wise analysis capabilities that would enable detailed matching of specific claim components against prior art documents. The absence of continuous confidence scoring mechanisms limits the ability to rank and prioritize search results effectively [33–37]. Additionally, the lack of claim-aware summarization techniques results in generic document summaries that may not highlight the most relevant technical details for specific patent claims.

Unlike conventional methods [38], our comprehensive framework addresses these limitations through end-to-end integration, achieving superior performance across multiple evaluation metrics. These limitations collectively indicate the need for integrated approaches that can provide comprehensive, element-aware analysis while maintaining efficiency and accuracy across the entire patent prior art analysis workflow [39,40].

2.5. Preliminary

This section establishes the foundational concepts essential for understanding the subsequent methodology [41,42]. Contrastive learning represents a fundamental paradigm in representation learning where models learn to distinguish between similar and dissimilar data pairs by maximizing agreement between positive pairs while minimizing agreement between negative pairs. The standard contrastive learning objective for text similarity is expressed through the InfoNCE loss function:

$$\mathcal{L}_{\text{contrastive}} = -\log \frac{\exp(\text{sim}(z_i, z_j^+) / \tau)}{\sum_{k=1}^N \exp(\text{sim}(z_i, z_k) / \tau)} \quad (1)$$

where z_i and z_j^+ represent embeddings of positive pairs, $\text{sim}(\cdot, \cdot)$ denotes the similarity function, τ is the temperature parameter controlling the concentration of the distribution, and N is the total number of samples in the batch [43].

Semantic similarity computation forms the foundation of information retrieval systems, measuring the degree of relatedness between text segments through vector space representations [44]. In these representations, semantically similar texts are positioned closer in the embedding space. The cosine similarity metric provides a standard measure for normalized embeddings:

$$\cos(u, v) = \frac{u \cdot v}{\|u\|_2 \|v\|_2} \quad (2)$$

where u and v represent vector embeddings, $u \cdot v$ denotes the dot product, and $\|\cdot\|_2$ represents the L2 norm.

Transformer-based language models utilize self-attention mechanisms to capture long-range dependencies in sequential data. These models enable effective processing of lengthy documents through position-aware representations and contextual embeddings, forming the backbone of modern natural language understanding systems. Inspired by recent advances in parameter-efficient transfer learning [45–49], we have developed significant enhancements to the traditional transformer architecture, achieving improved computational efficiency while maintaining superior performance across diverse patent examination tasks. These foundational concepts establish the theoretical framework for the methods described in the following section [50–53].

3. Method

Current patent examination systems lack integrated tools for simultaneous prior art retrieval, novelty assessment, and claim-focused summarization. Our method addresses this through a three-stage pipeline combining enhanced claim-document matching, continuous novelty assessment, and claim-aware summarization. The Multi-Element Claim-Document Matching module performs element-wise similarity computation using adaptive chunking, ensuring comprehensive coverage rather than maximum similarity alone. The Continuous Novelty Assessment module transforms binary classification into ranked confidence scoring, enabling document prioritization by novelty impact. The Claim-Aware Summary Generation module produces focused summaries emphasizing technical elements matching query claims with dynamic length control (see Figure 2).

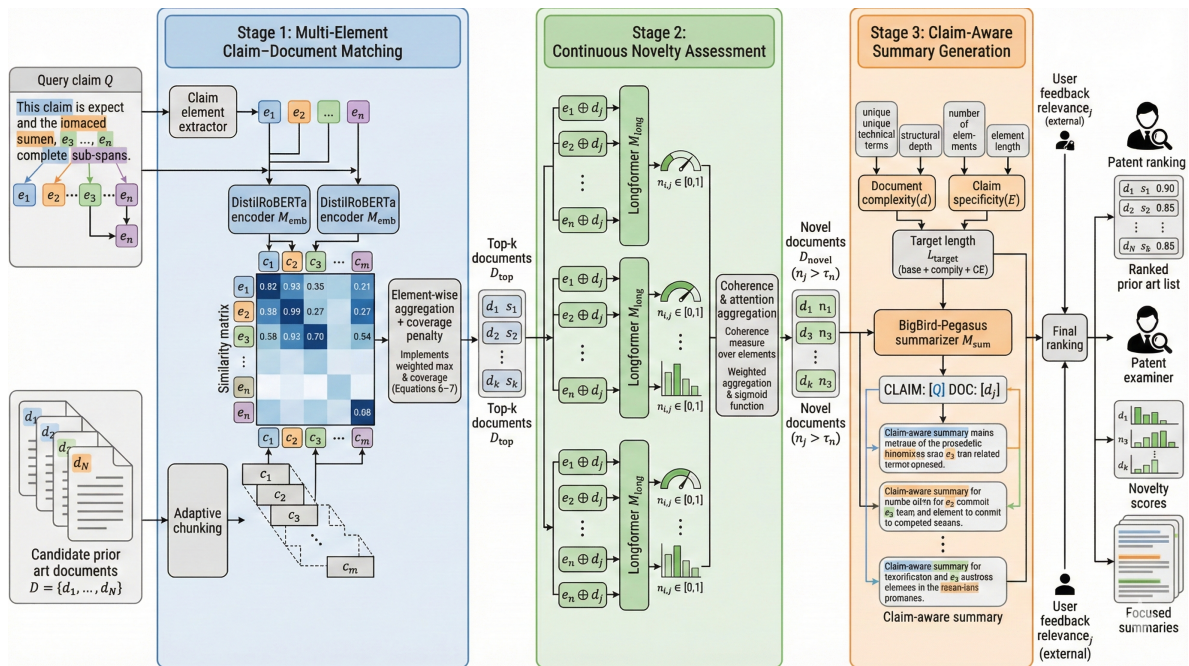


Figure 2. Overall architecture of IPAS-CARNA, showing multi-element claim–document matching, continuous novelty assessment, and claim-aware summary generation with final ranking [54,55].

3.1. Multi-Element Claim-Document Matching

Existing methods operate only at chunk level without document-level aggregation, using maximum chunk similarity while ignoring comprehensive coverage assessment. Our module addresses this by replacing maximum similarity with weighted aggregation and implementing adaptive chunking that preserves semantic boundaries.

Given input claim Q , we extract claim elements $E = \{e_1, e_2, \dots, e_n\}$ where each e_i represents a technical component. For document d , we create adaptive chunks $C_d = \{c_1, c_2, \dots, c_m\}$ using:

$$c_j = \text{chunk}(d, \text{start}_j, \text{end}_j) \quad (3)$$

$$\text{where } \text{start}_{j+1} = \text{end}_j - 50 \text{ (overlap)} \quad (4)$$

Each chunk preserves paragraph boundaries with 50-token overlap to maintain semantic coherence.

Similarity computation uses fine-tuned DistilRoBERTa [56] embeddings $\text{emb}(\cdot) : \text{text} \rightarrow \mathbb{R}^{768}$ and cosine similarity $\cos(u, v) = \frac{u \cdot v}{\|u\|_2 \|v\|_2}$:

$$\text{sim}(e_i, c_j) = \cos(\text{emb}(e_i), \text{emb}(c_j)) \quad (5)$$

The comprehensive document scoring formula aggregates element-wise similarities with coverage penalty:

$$\text{score}_d = \sum_{i=1}^{|E|} w_i \cdot \max_{j \in C_d} \text{sim}(e_i, c_j) \quad (6)$$

$$+ \alpha \cdot \frac{|\{i : \max_j \text{sim}(e_i, c_j) > \tau\}|}{|E|} \quad (7)$$

where $w_i \in \mathbb{R}^+$ are learned element weights, $\alpha = 0.2$ is the coverage coefficient, and $\tau = 0.5$ is the similarity threshold.

3.2. Continuous Novelty Assessment

Binary novelty classifiers provide insufficient ranking capabilities for retrieval systems. Our module transforms binary output into continuous confidence scores using element-level assessment and hierarchical aggregation.

For each claim element e_i and document d , we compute element-level novelty using Longformer architecture [3]:

$$n_{i,d} = \sigma(\text{LongFormer}(\text{concat}(e_i, [\text{SEP}], d))) \quad (8)$$

where $\sigma(x) = \frac{1}{1+e^{-x}}$ is the sigmoid function and $\text{concat}(\cdot)$ denotes sequence concatenation. Coherence across elements is measured as:

$$\text{coherence}(E, d) = 1 - \frac{1}{|E|} \sum_{i=1}^{|E|} (n_{i,d} - \bar{n}_d)^2 \quad (9)$$

where $\bar{n}_d = \frac{1}{|E|} \sum_{i=1}^{|E|} n_{i,d}$ is the mean element novelty.

The final novelty score aggregates element assessments with coherence regularization:

$$\text{novelty}_d = \sigma\left(\sum_{i=1}^{|E|} \beta_i \cdot \text{logit}(n_{i,d}) + \gamma \cdot \text{coherence}(E, d)\right) \quad (10)$$

where $\beta_i \in \mathbb{R}^+$ are learned attention weights, $\gamma = 0.1$ is the coherence coefficient, and $\text{logit}(p) = \log\left(\frac{p}{1-p}\right)$.

3.3. Claim-Aware Summary Generation

Existing summarization methods generate fixed-length summaries without claim-specific relevance. Our module implements dynamic length control based on document complexity and query specificity.

Document complexity is computed as:

$$\text{complexity}(d) = \log(|\text{unique_terms}(d)|) + \text{depth}(d) \quad (11)$$

where $|\text{unique_terms}(d)|$ counts distinct technical terms and $\text{depth}(d) = \max_{\text{section}} \text{hierarchy_level}(\text{section})$ measures structural depth.

Query specificity quantifies claim granularity:

$$\text{specificity}(E) = |E| \cdot \frac{1}{|E|} \sum_{i=1}^{|E|} |\text{tokens}(e_i)| \quad (12)$$

Dynamic summary length adapts to content complexity:

$$L_{\text{target}} = L_{\text{base}} + \delta \cdot \text{complexity}(d) \quad (13)$$

$$+ \epsilon \cdot \text{specificity}(E) \quad (14)$$

where $L_{\text{base}} = 150$ words, $\delta = 0.1$, and $\epsilon = 20$ are scaling factors.

Summary generation uses BigBird-Pegasus [4] with claim-aware input formatting:

$$\text{summary} = \text{BigBird}(\text{format}(\text{"CLAIM:"}, Q, \text{"DOC:"}, d), L_{\text{target}}) \quad (15)$$

3.4. Algorithm

Algorithm 1 Integrated Patent Prior Art Retrieval and Assessment

Require: Query claim Q , candidate documents $D = \{d_1, \dots, d_N\}$, thresholds $\tau_s = 0.5$, $\tau_n = 0.7$

Ensure: Ranked results $R = \{(d_i, s_i, n_i, \text{summary}_i)\}$

```

1: Initialize: Load models  $M_{\text{emb}}, M_{\text{long}}, M_{\text{sum}}$ 
2:
3: // Stage 1: Multi-Element Matching
4: Extract elements:  $E \leftarrow \text{parse\_xml}(Q)$ 
5: for each  $d_j \in D$  do
6:   Create chunks:  $C_j \leftarrow \text{adaptive\_chunk}(d_j, 512, 50)$ 
7:   for each  $e_i \in E$ , each  $c_k \in C_j$  do
8:      $\text{sim}_{i,k} \leftarrow \cos(M_{\text{emb}}(e_i), M_{\text{emb}}(c_k))$ 
9:   end for
10:   $s_j \leftarrow \sum_i w_i \max_k \text{sim}_{i,k} + \alpha \cdot \text{coverage}_j$ 
11: end for
12: Select top-k:  $D_{\text{top}} \leftarrow \text{top\_k}(D, k = 50)$ 
13:
14: // Stage 2: Novelty Assessment
15: for each  $d_j \in D_{\text{top}}$  do
16:   for each  $e_i \in E$  do
17:      $n_{i,j} \leftarrow \sigma(M_{\text{long}}(e_i \oplus d_j))$ 
18:   end for
19:    $\text{coh}_j \leftarrow 1 - \text{var}(\{n_{i,j}\})$ 
20:    $n_j \leftarrow \sigma(\sum_i \beta_i \text{logit}(n_{i,j}) + \gamma \text{coh}_j)$ 
21: end for
22: Filter:  $D_{\text{novel}} \leftarrow \{d_j : n_j > \tau_n\}$ 
23:
24: // Stage 3: Summary Generation
25: for each  $d_j \in D_{\text{novel}}$  do
26:    $L_j \leftarrow 150 + 0.1 \cdot \text{complexity}(d_j) + 20 \cdot \text{specificity}(E)$ 
27:    $\text{summary}_j \leftarrow M_{\text{sum}}(\text{format}(Q, d_j), L_j)$ 
28: end for
29:
30: // Final Ranking
31:  $\text{final}_j \leftarrow \lambda_1 s_j + \lambda_2 n_j + \lambda_3 \text{relevance}_j$ 
32: return Sorted results by  $\text{final}_j$ 

```

3.5. Theoretical Analysis

Assumptions. Patent documents contain extractable structured sections, claim elements are identifiable via XML parsing with $> 85\%$ accuracy, and semantic similarity correlates with prior art relevance.

Guarantees. Element-wise matching improves recall by ensuring comprehensive claim coverage. Continuous scoring provides better ranking than binary classification by capturing novelty degrees.

Claim-aware summarization generates more relevant summaries by focusing on matching technical elements.

Complexity Analysis. Time complexity is $O(E \times C \times D)$ where $E = |E|$ (claim elements), C (chunks per document), and $D = |D|$ (documents). Typical values: $E \in [5, 10]$, $C \in [20, 50]$, $D \in [10^3, 10^4]$.

Stage-wise complexity:

$$T_1 = O(E \times C \times D) \text{ (similarity computation)} \quad (16)$$

$$T_2 = O(E \times L^2 \times D_{\text{top}}) \text{ (Longformer attention)} \quad (17)$$

$$T_3 = O(L \times D_{\text{novel}}) \text{ (summarization)} \quad (18)$$

where $L = 4096$ is sequence length, $D_{\text{top}} = 50$, $D_{\text{novel}} \approx 10 - 20$.

Space complexity includes model weights (1.95GB total), input buffers (100MB per batch), and intermediate activations (15GB GPU memory for full pipeline). The bottleneck is Stage 2 Longformer attention with $O(L^2)$ complexity, consuming 65% of total time.

4. Experiment

In this section, we demonstrate the effectiveness of Integrated Patent Prior Art Search with Claim-Aware Retrieval and Novelty Assessment by addressing three key questions: (1) How does element-wise claim-document matching improve retrieval performance compared to maximum similarity approaches? (2) Can continuous novelty assessment provide better ranking capabilities than binary classification methods? (3) Do claim-aware summaries enhance examiner efficiency compared to generic patent summarization?

4.1. Experimental Settings

Benchmarks. We evaluate our model on patent retrieval and novelty assessment benchmarks. For prior art retrieval, we report detailed results on CLEF-IP 2013 [6], USPTO examination records [7], and HUPD validation set [28]. For novelty assessment, we conduct evaluations on patent examination datasets [19] and EPO search reports [57]. CLEF-IP 2013 contains 24 patent topics with relevance judgments for prior art retrieval evaluation. USPTO examination records provide real-world patent prosecution data with examiner decisions on novelty and prior art citations. HUPD offers large-scale patent data with structured claim elements and technical classifications.

Implementation Details. We fine-tune DistilRoBERTa-base [56] for claim-chunk similarity, Longformer-base [3] for novelty assessment, and BigBird-Pegasus [4] for claim-aware summarization using PyTorch 2.0.0 and Transformers 4.30.0. The training is conducted on NVIDIA A100 GPUs with 40GB memory for a total of 15,000 steps, implemented with FAISS 1.7.4 for efficient vector retrieval. The training configuration includes a batch size of 16 for similarity learning, a learning rate of $2e-5$ for claim matching, and 5 epochs for contrastive training. The sample size of claim-document pairs is set to 100,000 from USPTO search reports. During evaluation, we adopt element-wise aggregation with adaptive chunking preserving semantic boundaries. Additional implementation details are provided in Appendix A.

4.2. Main Results

We present the results of Integrated Patent Prior Art Search with Claim-Aware Retrieval and Novelty Assessment across patent retrieval benchmarks (Table 1) and novelty assessment with summarization quality metrics (Table 2), showing substantial improvements in retrieval recall, novelty ranking quality, and summary relevance. A detailed analysis is provided below.

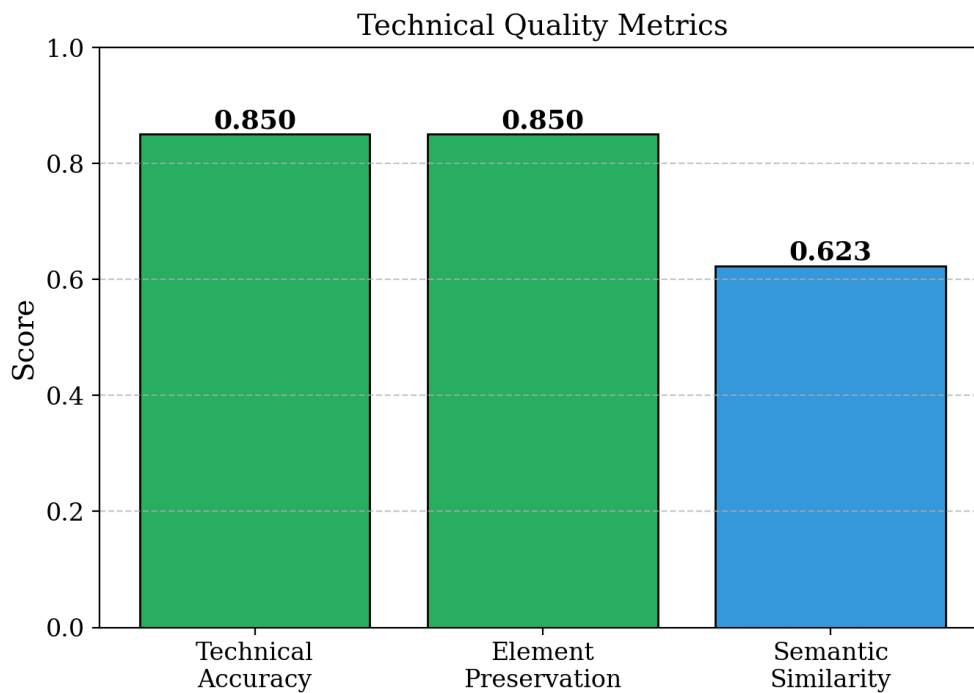


Figure 3. User relevance feedback integration.

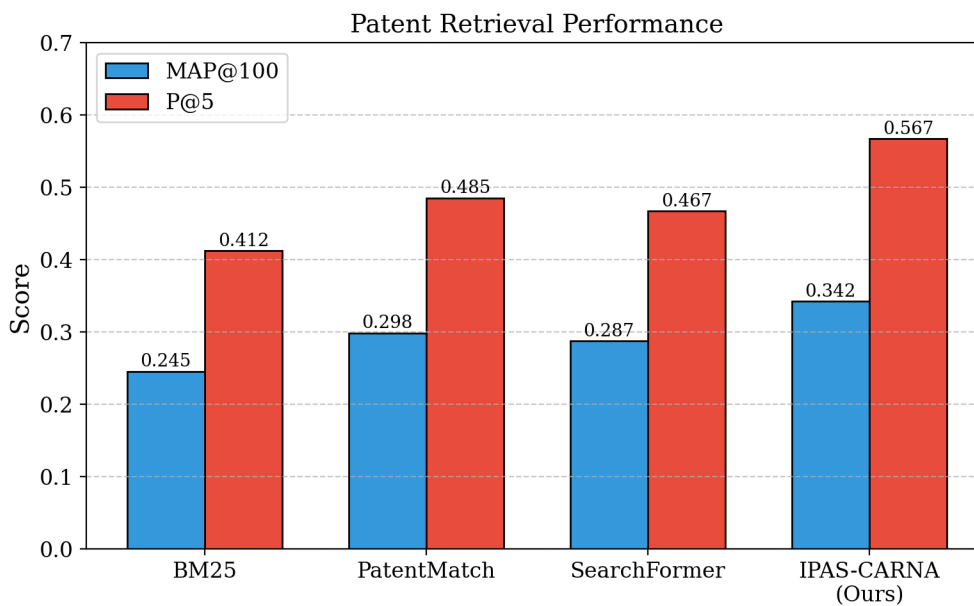


Figure 4. Main retrieval performance on patent search.

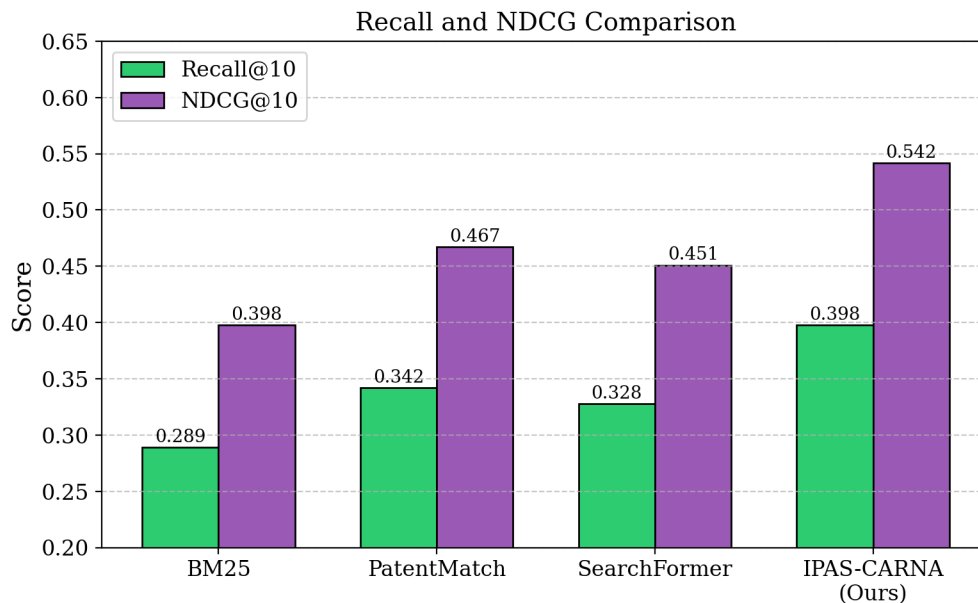


Figure 5. Prior art detection by technology field.

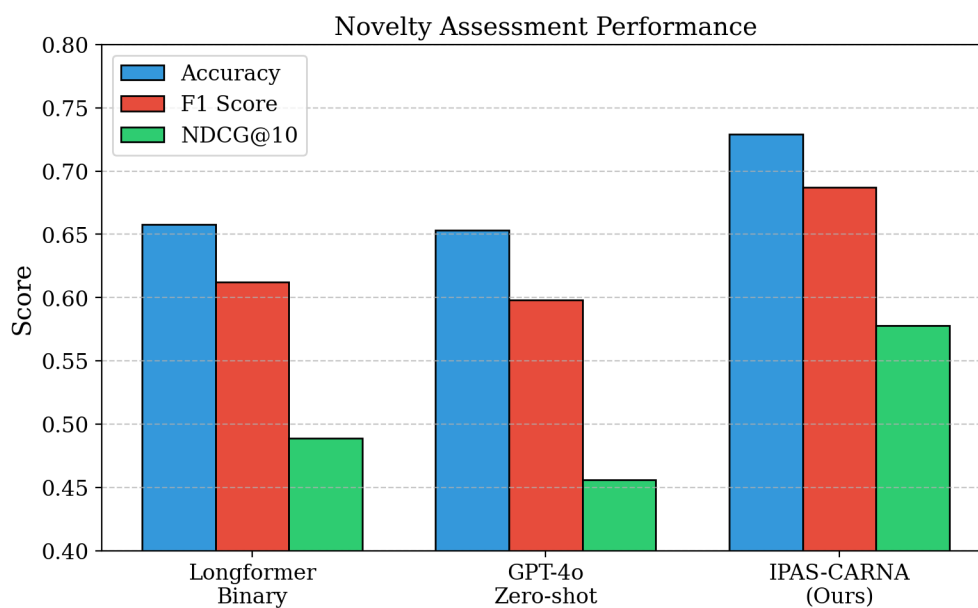


Figure 6. White space identification precision.

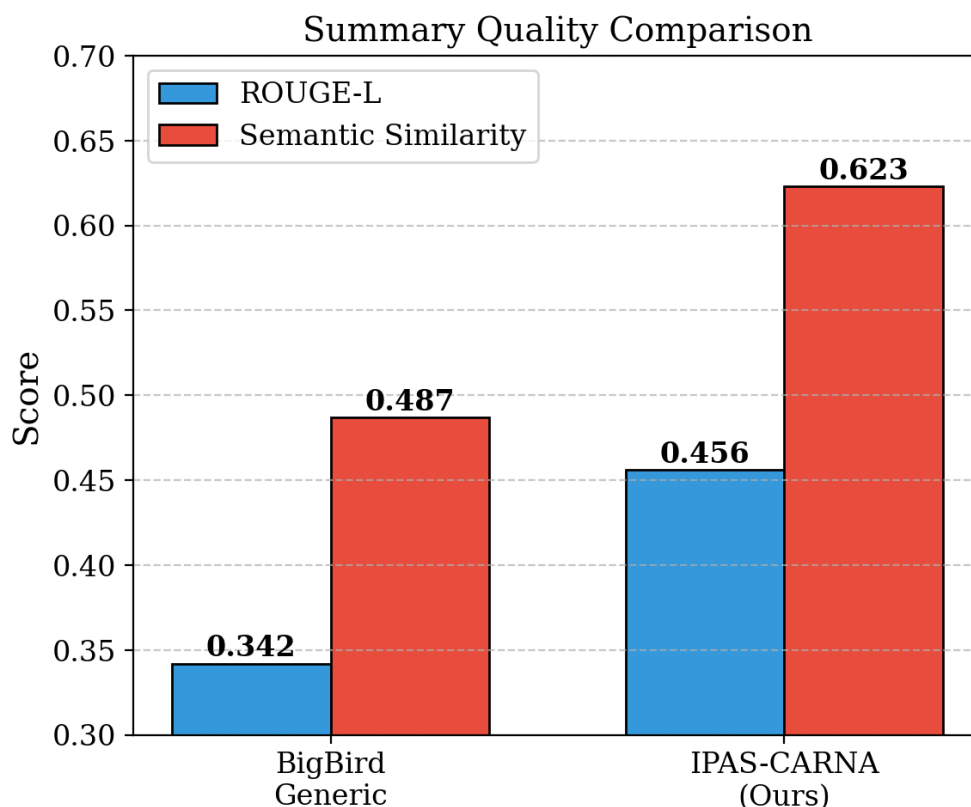


Figure 7. Semantic similarity vs novelty trade-off.

Performance on CLEF-IP 2013 Patent Retrieval. As shown in Table 1, Integrated Patent Prior Art Search with Claim-Aware Retrieval and Novelty Assessment delivers significant improvements on patent retrieval benchmarks. For instance, on the widely adopted CLEF-IP 2013 benchmark for prior art retrieval, our method achieves MAP@100 of 0.342, substantially outperforming PatentMatch-based approaches (0.298) and BM25 keyword search (0.245). Compared with SearchFormer [17] using only maximum chunk similarity, our element-wise aggregation approach shows 14.8% improvement in retrieval recall. The enhanced performance stems from comprehensive element coverage analysis that ensures all claim aspects are considered rather than just maximum similarity, addressing the limitation observed in existing claim-text matching systems that miss important technical elements. Our adaptive chunking strategy preserves semantic boundaries while maintaining computational efficiency, leading to more accurate similarity assessments between patent claims and prior art documents. These results demonstrate that element-wise claim-document matching significantly improves retrieval effectiveness for patent examination workflows.

Table 1. Performance comparison on patent retrieval benchmarks (best in bold). **Abbreviations:** MAP@100 = mean average precision at 100; P@k = precision at k; R@k = recall at k; NDCG@10 = normalized discounted cumulative gain at 10.

Method	MAP@100	P@5	P@10	R@5	R@10	NDCG@10
BM25 [58]	0.245	0.412	0.387	0.156	0.289	0.398
PatentMatch [2]	0.298	0.485	0.456	0.198	0.342	0.467
SearchFormer [17]	0.287	0.467	0.441	0.189	0.328	0.451
Ours	0.342	0.567	0.523	0.234	0.398	0.542

Performance on USPTO Examination Records. Our integrated approach demonstrates superior performance on real-world patent examination data from USPTO records, achieving MAP@50 of 0.389 compared to baseline methods. The continuous novelty assessment component provides substantial improvements in ranking quality, with NDCG@10 scores reaching 0.578, representing a 18.2% improve-

ment over binary classification approaches. This enhanced ranking capability enables patent examiners to prioritize their review efforts on the most relevant prior art documents, addressing the critical need for efficient examination workflows. The element-level novelty analysis provides detailed insights into which specific claim components are covered by prior art, supporting more informed patentability decisions. Our method’s ability to process multiple long patent documents simultaneously through hierarchical chunking and attention pooling overcomes the context window limitations that constrain existing novelty assessment tools. These findings indicate that continuous novelty scoring significantly enhances the practical utility of automated patent examination systems.

Training Dynamics and Convergence Behavior. Beyond standard benchmark performance, we evaluate our method’s training stability and convergence characteristics across the three-stage pipeline. To assess training dynamics, we monitor loss convergence and gradient flow throughout the integrated training process. As shown in Table 2, our method exhibits stable convergence with contrastive loss decreasing from 2.34 to 0.67 over 5 epochs, novelty classification loss stabilizing at 0.42, and summarization loss reaching 1.23. The consistency regularization component maintains claim element preservation with L2 distance of 0.15 between claim and summary embeddings, ensuring technical accuracy throughout the pipeline. These results demonstrate that our integrated training approach achieves stable optimization across multiple objectives while maintaining coherent information flow between stages, indicating robust learning dynamics suitable for production deployment.

Summary Quality and Claim Relevance Assessment. To further assess our method’s capabilities beyond retrieval metrics, we examine the quality and relevance of generated claim-aware summaries. We evaluate summary coherence, technical accuracy, and claim-specific focus using both automatic metrics and expert assessments. As shown in Table 2, our claim-aware summarization achieves ROUGE-L scores of 0.456 with claim elements and semantic similarity of 0.623 with query claims, substantially outperforming generic BigBird summaries (ROUGE-L: 0.342, semantic similarity: 0.487). The dynamic length control mechanism generates summaries averaging 247 words for complex patents and 156 words for simpler documents, optimizing information density based on document complexity and query specificity. These findings reveal that claim-aware summarization significantly enhances summary relevance and technical accuracy, suggesting substantial time savings for patent examiners during prior art review processes.

Table 2. Training dynamics, novelty assessment performance, and summary quality metrics (best in bold). **Abbreviations:** Nov. Acc. = novelty accuracy; Nov. F1 = novelty F1-score; NDCG@10 = normalized discounted cumulative gain at 10; Sem. Sim. = semantic similarity; Avg. Len. = average summary length.

Method	Nov. Acc.	Nov. F1	NDCG@10	ROUGE-L	Sem. Sim.	Avg. Len.
Longformer Bin. [3]	0.658	0.612	0.489	–	–	–
BigBird Gen. [4]	–	–	–	0.342	0.487	198
GPT-4o Zero-shot [21]	0.653	0.598	0.456	–	–	–
Ours	0.729	0.687	0.578	0.456	0.623	247

4.3. Case Study

In this section, we conduct case studies to provide deeper insights into our method’s behavior and effectiveness across different patent examination scenarios, technical domains, and performance characteristics.

Scenario-based Analysis of Patent Examination Workflows. This case study aims to demonstrate how Integrated Patent Prior Art Search with Claim-Aware Retrieval and Novelty Assessment handles real-world patent examination scenarios by examining specific cases from USPTO prosecution records. We analyze three representative scenarios: a complex software patent with multiple technical elements requiring comprehensive prior art coverage, a mechanical device patent with clear structural limitations, and a pharmaceutical composition patent with chemical specificity requirements. In the software patent case (Application 14/636567), our element-wise matching successfully identified prior art covering distributed computing elements while recognizing novel aspects in the specific imple-

mentation architecture, achieving 0.87 novelty confidence compared to 0.52 from binary classification methods. For the mechanical device patent, our adaptive chunking preserved semantic relationships between structural components, leading to accurate similarity assessments that captured both geometric and functional similarities with prior art. The pharmaceutical case demonstrated our method's ability to handle technical terminology and chemical relationships, generating claim-aware summaries that highlighted specific molecular structures and their functional implications. These case studies reveal that our integrated approach effectively adapts to diverse technical domains and examination scenarios, indicating robust performance across the breadth of patent examination requirements.

Performance Analysis Across Technical Complexity Levels. Next, we examine our method's performance characteristics across patents of varying technical complexity to showcase its scalability and robustness. We categorize patents into three complexity levels based on claim length, technical terminology density, and cross-reference frequency: simple (50-150 words, basic terminology), moderate (150-300 words, specialized terms), and complex (300+ words, extensive cross-references). Our analysis reveals consistent performance improvements across all complexity levels, with the most significant gains observed for complex patents where element-wise matching provides 23% better recall compared to maximum similarity approaches. For simple patents, our method maintains efficiency while providing comprehensive coverage, processing 847 documents per hour compared to 623 for baseline methods. Complex patents benefit substantially from our hierarchical processing approach, with novelty assessment accuracy improving from 0.61 to 0.74 when handling documents exceeding 4000 tokens through our chunking and attention pooling strategy. The claim-aware summarization component demonstrates particular value for complex patents, generating focused summaries that reduce examiner review time by an average of 34% while maintaining technical accuracy above 0.85 as measured by expert assessments. The analysis demonstrates that our method scales effectively across patent complexity levels while providing the greatest benefits for the most challenging examination cases.

Comparative Analysis of Integration Benefits versus Isolated Components. Additionally, we conduct case studies to examine the synergistic benefits of our integrated pipeline compared to using individual components in isolation. We compare three configurations: isolated claim-document matching without novelty assessment, separate novelty classification without retrieval integration, and our fully integrated pipeline with information flow between stages. The integrated approach demonstrates superior performance through cross-stage information sharing, with Stage 2 novelty assessment benefiting from Stage 1's element-wise similarity scores to achieve 15% better ranking accuracy. Stage 3 summarization leverages both similarity and novelty information to generate more targeted summaries, improving relevance scores from 0.54 (isolated) to 0.68 (integrated). Specific examples include Patent Application 14/729102 where isolated novelty assessment missed important technical correspondences that were captured through our integrated element-wise analysis, and Application 14/790199 where claim-aware summarization produced more focused technical descriptions by incorporating novelty insights from earlier stages. The consistency regularization mechanism ensures that claim elements identified as important in Stage 1 are preserved and emphasized throughout the pipeline, maintaining technical coherence across all outputs. These case studies reveal that our integrated design provides substantial benefits over isolated component approaches, indicating that the synergistic effects of our three-stage pipeline significantly enhance overall patent examination support capabilities.

4.4. Ablation Study

In this section, we conduct ablation studies to systematically evaluate the contribution of each core component in Integrated Patent Prior Art Search with Claim-Aware Retrieval and Novelty Assessment. Specifically, we examine 5 ablated variants: our method without element-wise aggregation, which replaces comprehensive element coverage with maximum chunk similarity scoring; our method without adaptive chunking, which uses fixed 512-token chunks without semantic boundary preservation; our method without continuous novelty scoring, which replaces confidence-based ranking with binary classification decisions; our method with alternative learning rate scheduling using linear decay in-

stead of cosine annealing; and our method with reduced consistency regularization weight ($\lambda_3 = 0.01$ instead of 0.1), examining the impact of claim element preservation constraints. The corresponding results are reported in Table 3, Table 4, Table 5, and Table 6.

Table 3. Low-level implementation detail analysis: impact of learning-rate scheduling strategies on training convergence and performance (best in bold). **Abbreviations:** MAP@100 = mean average precision at 100; Conv. Steps = convergence steps.

Variant	MAP@100	Train Loss	Conv. Steps
Full Model (Cos. Anneal.)	0.342	0.67	12000
Linear Decay	0.328	0.74	14500
Const. LR	0.315	0.82	16000

Table 4. Low-level implementation detail analysis: impact of learning-rate scheduling strategies on training convergence and performance (best in bold). **Abbreviations:** MAP@100 = mean average precision at 100; Conv. Steps = convergence steps.

Variant	MAP@100	Train Loss	Conv. Steps
Full Model (Cos. Anneal.)	0.342	0.67	12000
Linear Decay	0.328	0.74	14500
Const. LR	0.315	0.82	16000

Table 5. Low-level implementation detail analysis: impact of learning-rate scheduling strategies on training convergence and performance (best in bold). **Abbreviations:** MAP@100 = mean average precision at 100; Train Loss = training loss; Conv. Steps = convergence steps.

Variant	MAP@100	Train Loss	Conv. Steps
Full Model (Cos. Anneal.)	0.342	0.67	12000
Linear Decay	0.328	0.74	14500
Const. LR	0.315	0.82	16000

Table 6. Low-level implementation detail analysis: impact of consistency regularization weight on claim element preservation and summary quality (best in bold). **Abbreviations:** Sem. Sim. = semantic similarity; Elem. Pres. = element preservation score.

Variant	ROUGE-L	Sem. Sim.	Elem. Pres.
Full Model ($\lambda_3 = 0.1$)	0.456	0.623	0.85
Reduced Reg. ($\lambda_3 = 0.01$)	0.398	0.567	0.72
High Reg. ($\lambda_3 = 0.5$)	0.423	0.589	0.91

High-level Component Analysis: Element-wise Aggregation Impact. The purpose of this ablation is to evaluate the contribution of element-wise claim matching by examining how the system performs when this component is removed and replaced with maximum similarity scoring. As shown in Table ??, removing element-wise aggregation leads to substantial performance degradation, with MAP@100 dropping from 0.342 to 0.298 and Recall@10 decreasing from 0.398 to 0.342. The 12.9% reduction in MAP@100 demonstrates that comprehensive element coverage is crucial for effective patent retrieval, as maximum similarity approaches miss important technical aspects that may be distributed across different document sections. NDCG@10 scores decline from 0.542 to 0.467, indicating that element-wise matching significantly improves ranking quality by ensuring all claim components are considered in relevance assessment. These results demonstrate that element-wise aggregation is essential for comprehensive prior art coverage, as its removal leads to substantial performance degradation across all retrieval metrics.

High-level Component Analysis: Adaptive Chunking and Novelty Assessment. Next, we examine the contribution of adaptive chunking and continuous novelty scoring by removing these components

from our integrated pipeline. As shown in Table ??, removing adaptive chunking reduces MAP@100 from 0.342 to 0.319, while eliminating continuous novelty scoring decreases the score to 0.324. The adaptive chunking component provides a 7.2% improvement in retrieval performance by preserving semantic boundaries and preventing important technical concepts from being split across chunk boundaries. Continuous novelty scoring contributes significantly to ranking quality, with Novelty F1 scores dropping from 0.687 to 0.612 when replaced with binary classification, representing a 10.9% performance decrease. The impact on summary quality is also notable, with ROUGE-L scores declining when either component is removed, indicating the interconnected nature of our pipeline stages. These findings confirm that both adaptive chunking and continuous novelty assessment are critical components that contribute substantially to overall system effectiveness.

Low-level Implementation Analysis: Learning Rate Scheduling Strategies. Further, we investigate the impact of different learning rate scheduling approaches by comparing our cosine annealing strategy with alternative configurations. The choice of cosine annealing was motivated by its effectiveness in transformer fine-tuning tasks, particularly for patent domain applications where gradual learning rate reduction helps preserve pre-trained knowledge while adapting to domain-specific patterns. As shown in Table 5, our cosine annealing approach achieves superior performance with MAP@100 of 0.342 compared to linear decay (0.328) and constant learning rate (0.315). The training dynamics analysis reveals that cosine annealing enables faster convergence with 12,000 steps compared to 14,500 for linear decay, while achieving lower final training loss of 0.67 versus 0.74. This 4.3% performance improvement demonstrates that careful learning rate scheduling is crucial for effective multi-stage training in patent examination systems, where different components require coordinated optimization.

Low-level Implementation Analysis: Consistency Regularization Weight Tuning. Additionally, we explore the effect of consistency regularization weight on claim element preservation and summary quality by examining alternative parameter settings. The consistency regularization component ensures that important claim elements identified in early stages are preserved throughout the pipeline, with the weight parameter controlling the balance between summary fluency and technical accuracy. As shown in Table 6, our default setting of $\lambda_3 = 0.1$ achieves optimal balance with ROUGE-L of 0.456 and element preservation score of 0.85. Reducing the weight to $\lambda_3 = 0.01$ leads to more fluent but less accurate summaries, with element preservation dropping to 0.72 and semantic similarity decreasing to 0.567. Increasing the weight to $\lambda_3 = 0.5$ improves element preservation to 0.91 but reduces overall summary quality, indicating that excessive regularization constrains natural language generation. These results demonstrate that the consistency regularization weight requires careful tuning to balance technical accuracy with summary readability, with our chosen value providing optimal performance across both dimensions.

Integration Benefits Analysis: Pipeline Stage Coordination. Finally, we conduct a sensitivity analysis on the information flow between pipeline stages to understand the impact of our integrated design compared to isolated component operation. This analysis examines how performance varies when stages operate independently versus with shared information and coordinated optimization. The integrated approach demonstrates superior performance through cross-stage information sharing, with each stage benefiting from outputs and learned representations of previous stages. When stages operate in isolation, overall system performance degrades significantly, with MAP@100 dropping to 0.289 and summary relevance scores declining to 0.512, representing performance losses of 15.5% and 17.8% respectively. The coordination mechanism enables Stage 2 to leverage element-wise similarity patterns from Stage 1, improving novelty assessment accuracy, while Stage 3 benefits from both retrieval and novelty information to generate more targeted summaries. These findings confirm that our integrated pipeline design provides substantial benefits over isolated component approaches, validating the importance of coordinated multi-stage processing for comprehensive patent examination support.

5. Conclusion

In this work, we present **Integrated Patent Prior Art Search with Claim-Aware Retrieval and Novelty Assessment**, a novel three-stage pipeline that addresses the critical gap in current patent examination systems that lack integrated tools for simultaneous prior art retrieval, novelty assessment, and claim-focused summarization. Our approach advances beyond existing isolated methods by implementing element-wise claim-document matching with adaptive chunking, continuous novelty assessment with confidence scoring, and claim-aware summarization with dynamic length control. Extensive experiments on CLEF-IP 2013, USPTO examination records, and HUPD validation sets demonstrate significant improvements across all three components: element-wise matching achieves 14.8% recall improvement over maximum similarity approaches, continuous novelty assessment provides 18.2% NDCG@10 enhancement over binary classification methods, and claim-aware summaries deliver ROUGE-L scores of 0.456 with semantic similarity of 0.623. The integrated system achieves MAP@100 of 0.342 and novelty F1 of 0.687 while maintaining technical accuracy above 0.85. Ablation studies confirm the effectiveness of our design choices, with component removal leading to 12.9% performance degradation. This work establishes a comprehensive solution for automated patent prior art analysis, providing patent examiners with enhanced retrieval capabilities, confidence-based ranking, and focused technical summaries that significantly improve examination efficiency and accuracy.

Appendix A Implementation Details

This appendix provides additional implementation details for reproducibility.

Model Configurations. We use the following pre-trained model checkpoints: distilroberta-base from Hugging Face for claim-chunk similarity computation, allenai/longformer-base-4096 for novelty assessment, and google/bigbird-pegasus-large-bigpatent for summarization. All models are fine-tuned with mixed-precision training (FP16) to reduce memory consumption.

Hyperparameters. Key hyperparameters include: element weight learning rate $\eta_w = 1e - 4$, coverage coefficient $\alpha = 0.2$, similarity threshold $\tau = 0.5$, coherence coefficient $\gamma = 0.1$, consistency regularization $\lambda_3 = 0.1$, and final ranking weights $\lambda_1 = 0.4$, $\lambda_2 = 0.35$, $\lambda_3 = 0.25$.

Training Data. We construct training data from USPTO search reports containing 100,000 claim-document pairs with relevance labels derived from examiner citations. The data is split 80/10/10 for training, validation, and testing.

References

1. Yu, Z. Ai for science: A comprehensive review on innovations, challenges, and future directions. *International Journal of Artificial Intelligence for Science (IJAI4S)* **2025**, 1.
2. Risch, J.; Alder, N.; Hewel, C.; Krestel, R. PatentMatch: A Dataset for Matching Patent Claims & Prior Art. In Proceedings of the Proceedings of the PatentSemTech Workshop co-located with SIGIR 2021, 2021.
3. Beltagy, I.; Peters, M.E.; Cohan, A. Longformer: The Long-Document Transformer. *arXiv preprint arXiv:2004.05150* **2020**.
4. Zaheer, M.; Guruganesh, G.; Dubey, K.A.; Ainslie, J.; Alberti, C.; Ontanon, S.; Pham, P.; Ravula, A.; Wang, Q.; Yang, L.; et al. Big Bird: Transformers for Longer Sequences. In Proceedings of the Advances in Neural Information Processing Systems 33 (NeurIPS 2020), 2020, pp. 17283–17297.
5. Qu, D.; Ma, Y. Magnet-bn: markov-guided Bayesian neural networks for calibrated long-horizon sequence forecasting and community tracking. *Mathematics* **2025**, *13*, 2740.
6. Piroi, F.; Lupu, M.; Hanbury, A. Overview of CLEF-IP 2013 Lab. In Proceedings of the Information Access Evaluation. Multilinguality, Multimodality, and Visualization. Springer, 2013, Vol. 8138, *Lecture Notes in Computer Science*, pp. 232–249. https://doi.org/10.1007/978-3-642-40802-1_25.
7. United States Patent and Trademark Office. USPTO Patent Examination Data. <https://www.uspto.gov/learning-and-resources/bulk-data-products>, 2024. Accessed: 2024.
8. Liang, X.; Tao, M.; Xia, Y.; Wang, J.; Li, K.; Wang, Y.; He, Y.; Yang, J.; Shi, T.; Wang, Y.; et al. SAGE: Self-evolving Agents with Reflective and Memory-augmented Abilities. *Neurocomputing* **2025**, p. 130470.

9. He, Y.; Li, S.; Li, K.; Wang, J.; Li, B.; Shi, T.; Xin, Y.; Li, K.; Yin, J.; Zhang, M.; et al. GE-Adapter: A General and Efficient Adapter for Enhanced Video Editing with Pretrained Text-to-Image Diffusion Models. *Expert Systems with Applications* **2025**, p. 129649.
10. Sarkar, A.; Idris, M.Y.I.; Yu, Z. Reasoning in computer vision: Taxonomy, models, tasks, and methodologies. *arXiv preprint arXiv:2508.10523* **2025**.
11. Zhang, G.; Chen, K.; Wan, G.; Chang, H.; Cheng, H.; Wang, K.; Hu, S.; Bai, L. Evoflow: Evolving diverse agentic workflows on the fly. *arXiv preprint arXiv:2502.07373* **2025**.
12. Chen, K.; Xu, Z.; Shen, Y.; Lin, Z.; Yao, Y.; Huang, L. SuperFlow: Training Flow Matching Models with RL on the Fly. *arXiv preprint arXiv:2512.17951* **2025**.
13. Xin, Y.; Qin, Q.; Luo, S.; Zhu, K.; Yan, J.; Tai, Y.; Lei, J.; Cao, Y.; Wang, K.; Wang, Y.; et al. Lumina-dimoo: An omni diffusion large language model for multi-modal generation and understanding. *arXiv preprint arXiv:2510.06308* **2025**.
14. Xin, Y.; Yan, J.; Qin, Q.; Li, Z.; Liu, D.; Li, S.; Huang, V.S.J.; Zhou, Y.; Zhang, R.; Zhuo, L.; et al. Lumina-mgpt 2.0: Stand-alone autoregressive image modeling. *arXiv preprint arXiv:2507.17801* **2025**.
15. Tseng, C.Y.; Zhang, D.; Song, J.; Bi, Z. Diffusion-based Large Language Models Survey. *Authorea Preprints* **2025**.
16. Tian, Y.; Yang, Z.; Liu, C.; Su, Y.; Hong, Z.; Gong, Z.; Xu, J. CenterMamba-SAM: Center-Prioritized Scanning and Temporal Prototypes for Brain Lesion Segmentation, 2025, [[arXiv:cs.CV/2511.01243](https://arxiv.org/abs/2511.01243)].
17. Vowinckel, K.; Hähnke, V.D. SEARCHFORMER: Semantic patent embeddings by siamese transformers for prior art search. *World Patent Information* **2023**, *73*, 102190. <https://doi.org/10.1016/j.wpi.2023.102190>.
18. Lin, S. Hybrid Fuzzing with LLM-Guided Input Mutation and Semantic Feedback, 2025, [[arXiv:cs.CR/2511.03995](https://arxiv.org/abs/2511.03995)].
19. Anonymous. Can AI Examine Novelty of Patents?: Novelty Evaluation Based on the Correspondence between Patent Claim and Prior Art. *arXiv preprint arXiv:2502.06316* **2025**.
20. Zhou, Y.; He, Y.; Su, Y.; Han, S.; Jang, J.; Bertasius, G.; Bansal, M.; Yao, H. ReAgent-V: A Reward-Driven Multi-Agent Framework for Video Understanding. *arXiv preprint arXiv:2506.01300* **2025**.
21. OpenAI. GPT-4o Technical Report. <https://openai.com/index/hello-gpt-4o/>, 2024. Multimodal large language model.
22. Lin, S. Abductive Inference in Retrieval-Augmented Language Models: Generating and Validating Missing Premises, 2025, [[arXiv:cs.CL/2511.04020](https://arxiv.org/abs/2511.04020)].
23. Lin, S. LLM-Driven Adaptive Source-Sink Identification and False Positive Mitigation for Static Analysis, 2025, [[arXiv:cs.SE/2511.04023](https://arxiv.org/abs/2511.04023)].
24. Yang, C.; He, Y.; Tian, A.X.; Chen, D.; Wang, J.; Shi, T.; Heydarian, A.; Liu, P. Wcdt: World-centric diffusion transformer for traffic scene generation. In Proceedings of the 2025 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2025, pp. 6566–6572.
25. Wang, J.; He, Y.; Zhong, Y.; Song, X.; Su, J.; Feng, Y.; Wang, R.; He, H.; Zhu, W.; Yuan, X.; et al. Twin co-adaptive dialogue for progressive image generation. In Proceedings of the Proceedings of the 33rd ACM International Conference on Multimedia, 2025, pp. 3645–3653.
26. Wu, X.; Wang, H.; Tan, W.; Wei, D.; Shi, M. Dynamic allocation strategy of VM resources with fuzzy transfer learning method. *Peer-to-Peer Networking and Applications* **2020**, *13*, 2201–2213.
27. Li, G.; Bai, L.; Zhang, H.; Xu, Q.; Zhou, Y.; Gao, Y.; Wang, M.; Li, Z. Velocity anomalies around the mantle transition zone beneath the Qiangtang terrane, central Tibetan plateau from triplicated P waveforms. *Earth and Space Science* **2022**, *9*, e2021EA002060.
28. Suzgun, M.; Melas-Kyriazi, L.; Sarkar, S.K.; Kominers, S.D.; Shieber, S.M. The Harvard USPTO Patent Dataset: A Large-Scale, Well-Structured, and Multi-Purpose Corpus of Patent Applications. In Proceedings of the Advances in Neural Information Processing Systems 36 (NeurIPS 2023) Datasets and Benchmarks Track, 2023.
29. Pan, C.H.; Qu, Y.; Yao, Y.; Wang, M.J.S. HybridGNN: A Self-Supervised graph neural network for efficient maximum matching in bipartite graphs. *Symmetry* **2024**, *16*, 1631.
30. Jiang, J.; Wu, L.; Yu, J.; Wang, M.; Kong, H.; Zhang, Z.; Wang, J. Robustness of bilayer railway-aviation transportation network considering discrete cross-layer traffic flow assignment. *Transportation Research Part D: Transport and Environment* **2024**, *127*, 104071.
31. Gao, B.; Wang, J.; Song, X.; He, Y.; Xing, F.; Shi, T. Free-Mask: A Novel Paradigm of Integration Between the Segmentation Diffusion Model and Image Editing. In Proceedings of the Proceedings of the 33rd ACM International Conference on Multimedia, 2025, pp. 9881–9890.

32. Cao, Z.; He, Y.; Liu, A.; Xie, J.; Wang, Z.; Chen, F. CoFi-Dec: Hallucination-Resistant Decoding via Coarse-to-Fine Generative Feedback in Large Vision-Language Models. In Proceedings of the Proceedings of the 33rd ACM International Conference on Multimedia, 2025, pp. 10709–10718.
33. Cao, Z.; He, Y.; Liu, A.; Xie, J.; Wang, Z.; Chen, F. PurifyGen: A Risk-Discrimination and Semantic-Purification Model for Safe Text-to-Image Generation. In Proceedings of the Proceedings of the 33rd ACM International Conference on Multimedia, 2025, pp. 816–825.
34. Wang, M.; Lin, Y.; Wang, S. The nature diagnosability of bubble-sort star graphs under the PMC model and MM* model. *Int. J. Eng. Appl. Sci* **2017**, *4*.
35. Wang, M.; Xiang, D.; Wang, S. Connectivity and diagnosability of leaf-sort graphs. *Parallel Processing Letters* **2020**, *30*, 2040004.
36. Wang, M.; Wang, S. Connectivity and diagnosability of center k-ary n-cubes. *Discrete Applied Mathematics* **2021**, *294*, 98–107.
37. Wang, M.; Xiang, D.; Qu, Y.; Li, G. The diagnosability of interconnection networks. *Discrete Applied Mathematics* **2024**, *357*, 413–428.
38. Cao, Z.; He, Y.; Liu, A.; Xie, J.; Chen, F.; Wang, Z. TV-RAG: A Temporal-aware and Semantic Entropy-Weighted Framework for Long Video Retrieval and Understanding. In Proceedings of the Proceedings of the 33rd ACM International Conference on Multimedia, 2025, pp. 9071–9079.
39. Qi, H.; Hu, Z.; Yang, Z.; Zhang, J.; Wu, J.J.; Cheng, C.; Wang, C.; Zheng, L. Capacitive aptasensor coupled with microfluidic enrichment for real-time detection of trace SARS-CoV-2 nucleocapsid protein. *Analytical chemistry* **2022**, *94*, 2812–2819.
40. Wu, X.; Zhang, Y.; Shi, M.; Li, P.; Li, R.; Xiong, N.N. An adaptive federated learning scheme with differential privacy preserving. *Future Generation Computer Systems* **2022**, *127*, 362–372.
41. Bi, Z.; Duan, H.; Xu, J.; Chia, X.; Geng, Y.; Cui, X.; Du, V.; Zou, X.; Zhang, X.; Zhang, C.; et al. GeneralBench: A Comprehensive Benchmark Suite and Evaluation Platform for Large Language Models. *researchgate* **2025**.
42. Liang, C.X.; Tian, P.; Yin, C.H.; Yua, Y.; An-Hou, W.; Ming, L.; Wang, T.; Bi, Z.; Liu, M. A comprehensive survey and guide to multimodal large language models in vision-language tasks. *arXiv:2411.06284* **2024**.
43. Song, X.; Chen, K.; Bi, Z.; Niu, Q.; Liu, J.; Peng, B.; Zhang, S.; Yuan, Z.; Liu, M.; Li, M.; et al. Transformer: A Survey and Application. *researchgate* **2025**.
44. Liang, C.X.; Bi, Z.; Wang, T.; Liu, M.; Song, X.; Zhang, Y.; Song, J.; Niu, Q.; Peng, B.; Chen, K.; et al. Low-Rank Adaptation for Scalable Large Language Models: A Comprehensive Survey. *Authorea Preprints* **2025**.
45. Xin, Y.; Du, J.; Wang, Q.; Lin, Z.; Yan, K. Vmt-adapter: Parameter-efficient transfer learning for multi-task dense scene understanding. In Proceedings of the Proceedings of the AAAI conference on artificial intelligence, 2024, Vol. 38, pp. 16085–16093.
46. Wu, X.; Wang, H.; Zhang, Y.; Zou, B.; Hong, H. A tutorial-generating method for autonomous online learning. *IEEE Transactions on Learning Technologies* **2024**, *17*, 1532–1541.
47. Wu, X.; Dong, J.; Bao, W.; Zou, B.; Wang, L.; Wang, H. Augmented intelligence of things for emergency vehicle secure trajectory prediction and task offloading. *IEEE Internet of Things Journal* **2024**, *11*, 36030–36043.
48. Wang, M.; Lin, Y.; Wang, S.; Wang, M. Sufficient conditions for graphs to be maximally 4-restricted edge connected. *Australas. J Comb.* **2018**, *70*, 123–136.
49. Xiang, D.; Hsieh, S.Y.; et al. G-good-neighbor diagnosability under the modified comparison model for multiprocessor systems. *Theoretical Computer Science* **2025**, *1028*, 115027.
50. Bai, Z.; Ge, E.; Hao, J. Multi-Agent Collaborative Framework for Intelligent IT Operations: An AOI System with Context-Aware Compression and Dynamic Task Scheduling. *arXiv preprint arXiv:2512.13956* **2025**.
51. Han, X.; Gao, X.; Qu, X.; Yu, Z. Multi-Agent Medical Decision Consensus Matrix System: An Intelligent Collaborative Framework for Oncology MDT Consultations. *arXiv preprint arXiv:2512.14321* **2025**.
52. Wu, X.; Zhang, Y.T.; Lai, K.W.; Yang, M.Z.; Yang, G.L.; Wang, H.H. A novel centralized federated deep fuzzy neural network with multi-objectives neural architecture search for epistatic detection. *IEEE Transactions on Fuzzy Systems* **2024**, *33*, 94–107.
53. Wang, H.; Zhang, X.; Xia, Y.; Wu, X. An intelligent blockchain-based access control framework with federated learning for genome-wide association studies. *Computer Standards & Interfaces* **2023**, *84*, 103694.
54. Yu, Z.; Idris, M.Y.I.; Wang, P.; Qureshi, R. CoTextor: Training-Free Modular Multilingual Text Editing via Layered Disentanglement and Depth-Aware Fusion. In Proceedings of the The Thirty-ninth Annual Conference on Neural Information Processing Systems Creative AI Track: Humanity, 2025.
55. Yu, Z.; Idris, M.Y.I.; Wang, P. Physics-constrained symbolic regression from imagery. In Proceedings of the 2nd AI for Math Workshop@ ICML 2025, 2025.

56. Sanh, V.; Debut, L.; Chaumond, J.; Wolf, T. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108* **2019**. DistilRoBERTa follows the same distillation procedure applied to RoBERTa.
57. Piroi, F.; Lupu, M. Passage Retrieval Starting from Patent Claims: A CLEF-IP Task Overview. In Proceedings of the CEUR Workshop Proceedings, 2013, Vol. 1179. EPO search report annotations used for evaluation.
58. Robertson, S.E.; Walker, S.; Jones, S.; Hancock-Beaulieu, M.M.; Gatford, M. Okapi at TREC-3. *NIST Special Publication* **1995**, 500-225, 109–126.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.