

A Work Project, presented as part of the requirements for the Award of a Master's degree in Business Analytics from the Nova School of Business and Economics.

Enhancing Product Categorization with LLMs

Dynamic Integration of Taxonomy-Specific Knowledge Graphs
with Large Language Models for Hierarchical Product
Categorization in E-Commerce

Linus Luedecke - 58411

Co-authors of Group Section:

Elias Markwardt - 61794

Kuba Biały - 61678

Leo Ebert Gómez - 59017

Rafael Alejandro Moles Marrero - 61896

Work project carried out under the supervision of:

Prof. Qiwei Han

22-01-2025

Abstract

This research explored techniques to improve Large Language Models performance for Hierarchical Product Classification (HPC), including optimized fine-tuning, optimal prompting techniques, taxonomy-specific Knowledge Graphs, leveraging Retrieval-Augmented Generation, and implementing LLM-based Entity Matching. Tested on benchmark datasets Icecat and WDC-222, these methods significantly enhanced LLMs' ability to solve HPC tasks across various scenarios. Results achieved a hierarchical F1-score (hF) of 0.921, surpassing traditional DL benchmarks (0.85 hF). While not outperforming proprietary models like GPT, the proposed approaches offer a cost-efficient and effective alternative for businesses, demonstrating strong performance without reliance on expensive LLM solutions.

Keywords

Large Language Models, Hierarchical Classification, E-Commerce, In-Context Learning, Fine Tuning, Prompt Engineering, Knowledge Graphs, Retrieval Augmented Generation, Entity Matching

This work used infrastructure and resources funded by Fundação para a Ciência e a Tecnologia (UID/ECO/00124/2013, UID/ECO/00124/2019 and Social Sciences DataLab, Project 22209), POR Lisboa (LISBOA-01-0145-FEDER-007722 and Social Sciences DataLab, Project 22209) and POR Norte (Social Sciences DataLab, Project 22209).

1 Introduction

The e-commerce landscape in 2024 is rapidly evolving, driven by the rise of new players like Temu and the growing impact of social shopping platforms such as TikTok and Instagram, alongside industry leaders like Amazon and Alibaba. Social shopping on social media platforms is reshaping the consumer shopping experience by seamlessly integrating product discovery with purchase. This integration is powered by the accurate categorization of products. With over 1.98 billion products sold on TikTok alone in 2023 (Statista, 2024), there is a growing demand for efficient and automated solutions capable of accurately categorizing products in this rapidly evolving industry.

Beyond e-commerce, product categorization plays a critical role in logistics and international trade, where precise classifications ensure accurate duties and compliance with regulations. Political developments, such as President-elect Donald Trump's proposed tariff reforms targeting specific product categories (Schwartz, 2024), underscore the significant impact that accurate classification can have on businesses. The World Customs Organization posit that misclassifications result in significant financial losses for global businesses each year, including fines, shipping delays, and increased compliance costs (World Customs Organization, 2023). With over 2 billion unique products in the global marketplace, scalable solutions for categorization are essential, driving the need for focused research and innovation.

Techniques to find scalable solutions have evolved over time, beginning with rule-based systems that relied on manually crafted rules. These proved to be both labor-intensive and brittle, often failing to adapt to new categories or datasets. These methods were succeeded by traditional machine learning approaches, which leveraged feature engineering and linguistic patterns to address hierarchical taxonomies. While these approaches improved scalability, they struggled with capturing semantic relationships and hierarchical dependencies. Advancements such as embedding techniques further enhanced semantic representation and classification accuracy, paving the way for more robust solutions. More recently, deep learning models such as Recurrent Neural Networks (RNN) and Long short-term memory (LSTM) have emerged as powerful tools for product categorization, enabling the direct capture of hierarchical dependencies from data while overcoming the limitations of manual feature engineering. Building

on these foundations, pretrained multimodal models (PML) and text-to-text frameworks have pushed the boundaries of scalability and accuracy.

Despite recent advances, product classification (also referred to as "product categorization" in the paper) in a hierarchical context remains challenging for both traditional and deep-learning techniques. Supervised methods rely heavily on labeled data, which can be costly and time-consuming to obtain at scale. Moreover, typical product taxonomies are highly imbalanced, with certain categories dominating others, leading to biased models that struggle to generalize effectively. The diversity of taxonomies across platforms further complicates the task, as models trained on one system may not easily adapt to another. Additionally, variations in product titles, descriptions, brands and images pose significant hurdles for models attempting to standardize and classify products accurately across global markets. These challenges underscore the need for more adaptive and scalable approaches to hierarchical product categorization.

These challenges motivate the exploration of large language models (LLMs) as a promising solution. Leveraging their ability to understand context and semantics to overcome the limitations of traditional and deep-learning methods in hierarchical product categorization. Successfully integrating LLMs into hierarchical product classification represents an exciting opportunity. Offering significant advancements in research by redefining hierarchical classification methodologies, while providing businesses with a cost-efficient alternative to proprietary models.

We hypothesize that employing modern LLM-enabled techniques has the potential to specialize open-sourced LLMs to the task of hierarchical product classification, to achieve performance comparable to state-of-the-art approaches. Their advanced natural language processing capabilities, especially when integrated with domain-specific knowledge, may provide a more accurate and dynamic solution to existing challenges. To investigate, we will employ various approaches to address following research questions:

RQ-A: *Does fine-tuning decoder-only LLMs with hierarchical context enables them to memorize, understand, and accurately predict hierarchical product categories?*

RQ-B: *How can prompting techniques be effectively utilized to improve performance and scalability in hierarchical product classification?*

RQ-C: *How can taxonomy-specific Knowledge Graphs improve the performance of LLMs in hierarchical product classification?*

RQ-D: *How can Retrieval-Augmented Generation enhance product categorization performance, and which approach is most effective?*

RQ-E: *Can LLM-based entity matching be used for hierarchical product classification, and how does it perform compared to current benchmarks?*

Drawing from these research questions, following further hypotheses are made. First, LLM improve generalization compared to non-LLM approaches, leveraging its pre-trained domain-specific knowledge and ability to adapt to diverse contexts. Second, similar results to supervised methods will be achieved at a fraction of the costs and training data usage.

The paper holds the following structure: *Part 2* discusses the theoretical background and related works of hierarchical product classification and the use of LLMs in this domain. In *Part 3* the relevant datasets and benchmarks are introduced followed by a thorough exploratory data analysis in *Part 4*. *Part 5* explains the adopted methodology and outlines experimental setups, forming the basis for the systematic presentation of results in *Part 6*. In *Part 7* a critical analysis of the results is performed, as well as a discussion on limitations and potential future work. Finally, *Part 8* concludes the study with a summary of findings. The subsequent papers tackle the aforementioned research questions individually. *Paper A* investigates various LLM fine-tuning techniques, followed by the evaluation of various prompting methods in *Paper B*. *Paper C* focuses on the integration of domain-specific KGs, while *Paper D* explores the use of various RAG architectures. Finally, *Paper E* introduces an innovative approach to incorporating entity matching into hierarchical product classification.

2 Background and Related Work

This section explores prior studies and related work on traditional product classification, recent advancements in LLMs and their applications in the domains of hierarchical classification and e-commerce.

2.1 Traditional Product Classification

2.1.1 Taxonomies and Hierarchical Classification

The organization of product catalogs into taxonomies or hierarchical structures is a prevalent practice aimed at facilitating efficient search and navigation for users. Examples include Google’s Product Taxonomy and Amazon’s Product Category Tree. While multi-class classification treats taxonomy categories as isolated entities, hierarchical classification leverages the structure to capture inter-category relationships (Krishnan and Amarthaluri, 2019).

Despite its advantages, automated classification within hierarchical taxonomies presents significant challenges. Lower levels of the hierarchy often require distinguishing between highly similar items, while higher levels must accommodate diverse parent categories. These challenges are further compounded by imbalanced datasets, evolving taxonomies, and the demand for scalable, adaptable models (Kozareva et al., 2016).

Recent studies have explored various strategies to address these challenges, with two prominent overarching approaches emerging: flat single classifiers and hierarchical multi-level classifiers. Flat single classifiers aim to classify leaf nodes in a single step, treating all categories independently without considering the hierarchical structure. In contrast, hierarchical multi-level classifiers follow the taxonomy’s structure, categorizing products step by step (Dumais and Chen, 2000). These approaches can further be divided into global methods, which model the entire hierarchy as a single system, and local methods, which focus on smaller segments of the taxonomy, such as specific nodes or levels (Jiang et al., 2022).

2.1.2 Traditional Supervised Methods

Traditional supervised learning methods have been widely explored for product classification on e-commerce platforms, leveraging various algorithms and feature engineering techniques to improve accuracy and efficiency. Several studies have demonstrated the effectiveness of algorithms such as Naïve Bayes, K-Nearest Neighbor (KNN), Decision Tree, Support Vector Machine (SVM), Random Forest, and Logistic Regression when applied to datasets consisting of product titles and descriptions. For instance, Mathivanan et al. (2019) identified KNN as the

best-performing algorithm due to its high accuracy and computational efficiency. In contrast, Patra et al. (2021) emphasized the robustness of Logistic Regression in classifying products.

Domain-specific enhancements to traditional algorithms have further advanced the field. Shen et al. (2011) at eBay Research Lab introduced a feature engineering framework that incorporated rich domain knowledge and linguistic cues, enabling their SVM-based model to address classification challenges in lower levels of hierarchical taxonomies.

Similarly, Kozareva et al. (2016) from Yahoo! Labs developed an automated mechanism for categorizing products based on their titles within a hierarchical taxonomy of 319 categories spanning six levels. By focusing on linguistic patterns inherent in product titles, this system efficiently assigned categories without relying on extensive metadata. Gupta et al. (2016) extended this line of research by incorporating Word2Vec embeddings to create a distributional semantics representation for product descriptions. These embeddings captured semantic nuances in product titles and descriptions, enabling their two-level ensemble approach to better classify products within the taxonomy. By utilizing classifiers corresponding to paths, nodes, and depths within the hierarchy, their method effectively reduced errors, leveraging both semantic relationships and the hierarchical structure for improved classification accuracy.

2.1.3 Deep Learning-Based Methods

The limitations of traditional supervised machine learning methods, particularly their reliance on manual feature engineering and their inability to capture more complex, contextual, or hierarchical relationships in texts, have driven a shift toward deep learning architectures. These models, such as RNNs and LSTMs have revolutionized product classification by learning hierarchical dependencies and semantic structures directly from data. Deep learning approaches are particularly well-suited for handling the complexities of hierarchical classification, where relationships between parent and child categories are crucial for accurate predictions.

Recurrent Neural Networks (RNNs) are widely used in hierarchical product classification due to their ability to model sequential dependencies in text data. Liu et al. (2016) proposed an RNN-based model that leveraged the sequential structure of product titles and descriptions to predict hierarchical categories, effectively capturing dependencies across taxonomy

levels. Building on this, Huang et al. (2019) introduced the Hierarchical Attention-based Recurrent Neural Network (HARNN), which integrates text content with hierarchical category structures. HARNN uses hierarchical attention to capture relationships between texts and taxonomies while modeling level-specific dependencies in a top-down manner. However, RNNs often face challenges such as vanishing gradients, which can hinder their ability to learn long-term dependencies across deep hierarchies.

To address the limitations of standard RNNs, LSTM networks were introduced as a more robust alternative. LSTMs are specifically designed to retain long-term dependencies, making them highly effective for hierarchical classification. Huang et al. (2021) utilized LSTMs to model hierarchical relationships in e-commerce taxonomies, demonstrating improved accuracy in classifying products into deeply nested categories. Their model leveraged the LSTM's gating mechanisms to preserve critical information about parent-child relationships while filtering irrelevant details. Chen et al. (2021) extended this approach by integrating LSTMs with attention mechanisms to prioritize important features in text data, further enhancing the performance of hierarchical product categorization. These studies highlight the strength of LSTMs in handling the depth and complexity of hierarchical structures in product taxonomies.

By learning directly from data and capturing hierarchical dependencies, deep learning architectures offer a powerful solution for large-scale and complex product classification tasks, paving the way for more sophisticated applications in e-commerce.

2.2 LLMs in Product Classification

2.2.1 Emergence of Pre-trained Language Models

The introduction of Transformer-based models (Vaswani et al., 2023) and subsequent pre-trained language models (PLMs) like BERT and GPT (Devlin et al., 2019; Brown et al., 2020) revolutionized NLP, and dethroned previous state-of-the-art LSTM networks for NLP (Melis et al., 2017). PLMs learn general language representations from unlabeled data and can then be adapted to diverse downstream tasks, often outperforming earlier specialised methods (Radford and Narasimhan, 2018).

Domain- and task-adaptive pretraining further enhance performance (Gururangan et al., 2020).

This approach shifts NLP paradigms toward leveraging PLMs for state-of-the-art results across tasks and domains (Min et al., 2021). Moreover, newly introduced in-context learning techniques (Brown et al., 2020) enable rapid adaptation of PLMs to new tasks via zero-shot and few-shot prompts that specialise model solely based on text instructions, reducing the need for expensive training and alterations of models' architectures (Mosbach et al., 2023).

According to Fan et al. (2023), as LLMs are adopted across various industries and domains, researchers continue to explore their expanding capabilities for tasks like NER, QA, and classification. The authors suggest that due to the rapid pace of research in this field, continual investigation into the capabilities and suitability of LLMs for specialized tasks is required.

2.2.2 Recent Advancements in Generative LLMs

Recent advancements in LLMs have produced architectures such as GPT (Brown et al., 2020; Achiam et al., 2024), Claude (AnthropicAI), Google Gemini (Anil et al., 2024), and LLaMa (Touvron et al., 2023; Grattafiori et al., 2024). These models enable rapid adaptation to downstream tasks without designing new architectures or managing complex training.

Sahoo et al. (2024) review classic prompt engineering methods, including zero-shot and few-shot learning, and modern approaches that enhance reasoning abilities. Logic-based techniques (chain-of-thought, self-consistency, tree-of-thoughts) represent a model's generation as a reasoning process. Other approaches integrate external data or third-party outputs (RAG, ReAct) to reduce LLMs hallucinations and improve factual correctness.

Although prompt engineering and RAG can outperform fine-tuned models in some scenarios, fine-tuning excels in specialized tasks (Shin et al., 2023), especially when domain-specific expertise is crucial (MathavRaj et al., 2024). Moreover, fine-tuned models can be combined with techniques such as RAG and advanced prompt engineering for industry-specific applications (de Luis Balaguer et al., 2024; Rangan and Yin, 2024), allowing to create more sophisticated, specialized solutions enabled by LLM-empowered techniques.

2.2.3 Language Models Implementations for Hierarchical Classification

Hierarchical text classification (HTC) has been broadly explored with encoder-only transformers, resulting in complex, specialized approaches. For instance, HiMatch highlights the importance of integrating semantic label information (Chen et al., 2021). HGCLR leverages a Graphformer network to incorporate label hierarchies directly into encoder training (Wang et al., 2022a). HPT enriches model training and inference with hierarchy constraints and hierarchy injection mechanisms (Wang et al., 2022b). Collectively, these strategies demonstrate that encoder-based HTC models has been widely explored by researchers resulting in highly complex and customised model architectures.

Unlike encoder-based approaches that focus on learning vector representations, Sequence-to-Sequence (seq2seq) tasks (e.g. generation, translation, summarization) involve processing sequences of tokens as both input and output. Recursive nature Encoder-decoder make them suitable for text generation tasks (Yang et al., 2023). Risch et al. (2020) pioneered in exploring the use of seq2seq models for hierarchical classification in the context of patent documents. Their findings demonstrated that treating hierarchical labels as sequences significantly improves performance when combined with transformer-based architectures for predicting label encodings as ordered sequences. Numerous frameworks have been proposed to adapt this strategy for hierarchical text classification, highlighting the issues of leveraging semantic meaning of labels, and constraining the model to generate relevant tokens representing existing labels by introducing vocabulary restriction mechanisms, custom Hierarchical Consistency Rate Metric (HCR) or loss functions penalising irrelevant vocabulary (Yu et al., 2022; Jain et al., 2024; Torba et al., 2024).

Recent research has also looked on decoder-only models and in-context learning strategies for HTC in addition to existing encoder-decoder approaches. Gholamian et al. (2024) showed that decoder-only models achieve state-of-the-art outcomes on benchmark datasets when in-context learning techniques were implemented. For instance, Chen et al. (2024b) introduce framework for mixing few-shot in-context learning hierarchical classification with specialized encoder implemented for context retrieval. In their approach they train a very specific method for retrieving in-context learning examples with the use of custom adapted encoder model,

before passing them into LLM prompt. They use hierarchical context by performing single generation per each level of the hierarchy, by using separate prompts at each level. However, even with these encouraging developments, there remains a lack of research on the use of generative decoder-only LLMs for hierarchical categorization, particularly in the e-commerce industry. This points to a promising avenue for further investigation on optimizing decoder-only models for HTC tasks in certain domains.

2.2.4 Examples of LLM Implementations in E-Commerce & Product Categorization

There have been several applications of LLMs in the e-commerce product categorization contexts. A pre-trained language model called E-BERT was developed to include phrase-level and product-level knowledge for overcoming BERT's restrictions to e-commerce tasks (Zhang et al., 2020). Techniques such as Adaptive Hybrid Masking and Neighbor Product Reconstruction were introduced to the model to make it high-performing in several e-commerce-related tasks including product categorization. Also, researchers produced a completely novel approach; from classification to machine translation for product categorization with an encoder-decoder model (Li et al. (2018); Tan et al. (2020)). The latter has the capability to translate product descriptions into a root-to-leaf path within a product taxonomy and outperforms traditional classification methods in terms of predictive accuracy. In addition, the machine translation model can generate new paths between previously unconnected nodes, transforming the taxonomy tree into a directed acyclic graph. This phenomenon improves user navigation and adjusts to new products, which may eventually transform systems of product categorization in e-commerce.

Decoder-only architectures have also been implemented by researchers for the purpose of e-commerce product categorization. Gholamian et al. (2024) approach hierarchical product classification on the WDC-222 dataset by utilizing decoder-only large language models. They compared various LLMs as well as different techniques including Zero-Shot and Few-Shot learning. Additionally, they evaluate flat and hierarchical approaches. Their hierarchical approach comprises using two models, in which the first predicts the top level class, while the second, retrieves the first classifier's output and uses it as input to predict the low level class. In

their results, hierarchical approaches usually underperformed compared to the flat classification methods. Also, their research was focused on various perturbations to the data, and they prove that these decoder LLMs (LLaMa, GPT) seem to be resistant to various data perturbations, which also highlights a potential of LLMs in HTC contexts. Finetuning and prompt engineering were integrated for e-commerce product classification by Cheng et al. (2024), who suggested a unique two-step framework that uses a fine-tuned LLM to predict the top K hierarchical categories for a given product. These are passed to another LLM prompted to reason about the correct final output. However, despite the recent advancements in the field of decoder-only language models and their application in HTC, the topic in the field of e-commerce product categorization remains relatively unexplored.

3 Data and Benchmarks

This section provides an overview of the two datasets used in our study, along with benchmark results for comparison. Both datasets, WDC-222 Gold Standard and Icecat, are widely recognized benchmark datasets. Most pre-processing and cleaning have already been carried out by the Mannheim Data Science Group of the University of Mannheim, which we will briefly summarize here (Nils Richter and Christian Bizer, n.d.).

3.1 Datasets

3.1.1 WDC-222 Gold Standard

The Web Data Commons (WDC)-222 Gold Standard is a manually curated subset of the WDC Training Dataset for Large-scale Product Matching (WDC-LSPM), specifically selecting products from the Icecat dataset. The matching was done using the Global Trade Item Number (GTIN), addressing vendor-specific inconsistencies in identifier terms, such as "sku" or "mpn". Additional columns were created in the WDC-222 dataset, including product titles, descriptions, and category information from Icecat. After cleansing and removing duplicates, the final dataset consists of 2,984 instances with 222 distinct leaf node categories. Similar to Icecat, it is imbalanced, with 100% prevalence for titles and 77% for descriptions. The hierarchy depth

ranges from 2 to 3, with an average of 2.47, mirroring the structure of Icecat. More information in *Part 4: EDA*. As the WDC-222 Gold Standard is derived from the Common Crawl - the largest publicly available web-crawled dataset - it is inherently more heterogeneous than the cleaner Icecat dataset, making it well-suited for testing purposes (Nils Richter and Christian Bizer, n.d.).

3.1.2 Icecat

Icecat is a global provider of multilingual, standardized product data sheets used by e-commerce platforms, ERP systems, and rating portals. It collaborates with 80,000 e-commerce companies and offers two catalogs: the paid "Full Icecat" with 6.5 million product sheets from 24,000 vendors, and the free "Open Icecat," which contains 1 million product sheets from 340 vendors. For our study, we utilized the preprocessed Open Icecat dataset, focusing on the "Computers & Electronics" category, where duplicates and categories with less than 25 products were removed to ensure sufficient training data. The final Icecat dataset contains 765,473 products across 370 categories, with product attributes including title, description, and brand. While every instance has a title and brand, detailed descriptions are available for only 70% of the products (Nils Richter and Christian Bizer, n.d.).

3.2 Benchmarks

In the initial experiments done by Nils Richter and Christian Bizer (n.d.) the performance of various classification methods was assessed on both WDC-222 and Icecat datasets. The lower performance on the noisier WDC-222 dataset, compared to the cleaner Icecat, highlights the challenges posed by noisy data. The table below summarizes the results, with weighted F1 (wF) and hierarchical F1 (hF) scores serving as the primary evaluation metrics. Additional metrics include weighted precision (wP) and weighted recall (wR).

The Icecat dataset, with its clean and structured nature, consistently achieved high performance across all classification methods. Both flat and hierarchical approaches reached near-perfect weighted and hierarchical F1 scores (0.99).

The WDC-222 dataset, characterized by real-world noise and complexity, exhibited lower

Group Part: Enhancing Product Classification with LLMs

Classification System	Icecat				WDC-222 Gold Standard			
	wP	wR	wF	hF	wP	wR	wF	hF
Initial Results								
Dictionary-based approach	0.79	0.43	0.48	0.55	0.72	0.61	0.62	0.71
Flat traditional classifier	0.98	0.98	0.98	0.99	0.80	0.72	0.73	0.79
LCPN traditional classifier	0.98	0.98	0.98	0.99	0.82	0.73	0.74	0.81
Flat FastText classifier	0.99	0.98	0.99	0.99	0.86	0.73	0.76	0.83
LCPN FastText classifier	0.99	0.98	0.98	0.99	0.85	0.76	0.78	0.84
Flat neural network	0.98	0.98	0.98	0.98	0.84	0.76	0.78	0.84
Neural network + LCPN traditional classifier	0.98	0.98	0.98	0.98	0.84	0.75	0.76	0.84
Neural network + LCPN FastText classifier	0.98	0.98	0.98	0.98	0.84	0.78	0.78	0.85
LLM-Based Results (clean data)								
DeBERTaV3 base Supervised (Flat)	0.98	0.98	0.98	–	0.82	0.71	0.73	–
DeBERTaV3 base Supervised (Hierarchical)	0.97	0.98	0.97	–	0.82	0.69	0.72	–
Llama-2 70b-chat (Flat)	0.50	0.37	0.37	–	0.76	0.51	0.51	–
Llama-2 70b-chat (Hierarchical)	0.65	0.40	0.39	–	0.69	0.42	0.38	–
Llama-2 70b-chat (Few Shot)	0.97	0.96	0.96	–	0.90	0.87	0.86	–
GPT-3.5 ver.: 0613 (Flat)	0.90	0.84	0.84	–	0.92	0.87	0.88	–
GPT-3.5 ver.: 0613 (Hierarchical)	0.88	0.66	0.66	–	0.82	0.65	0.66	–
GPT-3.5 ver.: 0613 (Few Shot)	0.98	0.97	0.97	–	0.94	0.92	0.93	–
GPT-4 ver.: 0613 (Flat)	0.94	0.91	0.91	–	0.95	0.89	0.90	–
GPT-4 ver.: 0613 (Hierarchical)	0.89	0.70	0.70	–	0.85	0.80	0.80	–
GPT-4 ver.: 0613 (Few-shot)	0.99	0.99	0.99	–	0.96	0.94	0.94	–

Table 1: Overview of experimental results: Initial and LLM-based (Gholamian et al., 2024) performance compared to Icecat. Despite this, methods employing hierarchical setups, particularly those integrating neural networks or FastText embeddings, achieved notable improvements. The hierarchical FastText classifier combined with neural networks at the root level achieved the highest hierarchical F1 score (0.85). This indicates the importance of combining deep learning techniques with hierarchical structuring for noisy datasets, establishing a comparable benchmark for this dataset.

The results also highlight a significant challenge in generalization. Models trained on Icecat struggled to maintain performance when tested on WDC-222, suggesting limited robustness to unseen, noisy data. This emphasizes the need for training strategies or model improvements that account for real-world noise to improve generalization across datasets.

Building on these initial results, recent research explored the use of Large Language Models (LLMs), such as GPT-4 and Llama 2, for product classification on the same datasets (Gholamian et al., 2024). These models were evaluated under both clean and noisy data conditions, with perturbations such as abbreviations and truncations simulating real-world challenges. GPT-4 achieved state-of-the-art performance on both datasets in clean-data scenarios, while also demonstrating strong robustness to noise compared to traditional methods. Few-shot

prompting, which provides contextual examples, further enhanced the model's performance, especially under noisy conditions. Although hierarchical configurations with LLMs showed limitations due to error propagation, flat classification using LLMs and few-shot prompting consistently delivered the best results. These findings highlight the potential of LLMs, particularly the closed GPT-4, to handle the variability and complexity of real-world e-commerce data, offering a significant advancement over traditional classification methods.

Key Observations

1. **High Performance on Clean Data:** The Icecat dataset achieved near-perfect scores (0.99) with both flat and hierarchical methods, demonstrating strong performance on clean, structured data.
2. **Challenges with Noisy Data:** The WDC-222 dataset showed lower performance, with hierarchical methods (e.g., FastText + neural networks) reaching 0.85 hF. Models trained on Icecat struggled to generalize to noisy data.
3. **LLMs Show Robustness:** GPT-4, a subscription-based model, achieved state-of-the-art results, excelling in noisy scenarios with few-shot prompting, and outperforms free counterpart LLMs.

4 Exploratory Data Analysis

This exploratory data analysis (EDA) provides a comprehensive comparison of the WDC and Icecat datasets, examining their hierarchical structures, product categories, title lengths, descriptions, and brand distributions to uncover key differences and insights into their respective data compositions.

4.1 Hierarchical Analysis

The hierarchies of the WDC and Icecat datasets share a common foundation in the "Computers & Electronics" category, both featuring the same top-level categories. Notably, Icecat includes a wider range of subcategories, particularly under Computer Cables and Software, highlighting

Group Part: Enhancing Product Classification with LLMs

a more extensive selection of digital solutions and peripherals. While both datasets emphasize consumer electronics, Icecat’s hierarchy reflects a stronger focus on modern trends such as Smart Wearables and advanced Networking components, indicating a broader market scope. Overall, WDC provides a more streamlined structure, while Icecat showcases a richer variety of product offerings within the same foundational categories.

Dataset	Average Depth	Max Depth	Min Depth
WDC	2.34	3	2
Icecat	2.38	3	2

Table 2: Average Category Depth Statistics for WDC and Icecat Datasets

The average depths across both datasets are comparable, suggesting a similar level of granularity in their hierarchical structures. Additionally, the maximum depths are equal, reflecting consistency in the categorization approach for deeper hierarchies. Overall, both datasets exhibit similar structural characteristics, highlighting a shared methodology in their classification systems.

4.2 Categorical Analysis

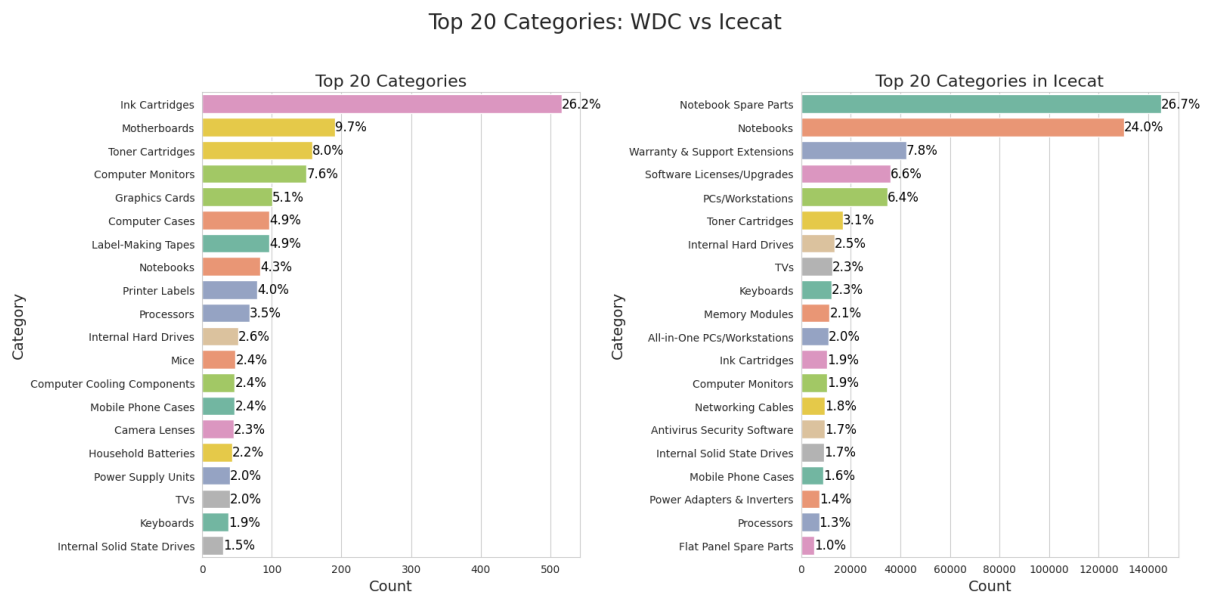


Figure 1: Top 20 Leaf Categories of WDC and Icecat datasets

The comparison between the WDC and Icecat datasets accentuates significant differences in the representation of the products. In WDC, there is higher dependency on consumables, for ex-

Group Part: Enhancing Product Classification with LLMs

ample, "Ink Cartridges" (516 units, 26.18%), "Computer Monitors" (150 units, 7.61%), while other categories, for example, "All-in-One PCs/Workstations" and "Graphics Cards" represent zero counts - proving lack of diversity. Contrarily, Icecat presents a much wider variety of products, particularly in high-demand categories such as "Notebooks" (130,294 units, 23.95%) and "Notebook Spare Parts" (145,020 units, 26.66%).

NaN values at various levels indicate that products are only categorized up to a certain depth. Depth 1 (17 unique categories) and Depth 2 (137 unique categories) contain no NaN-values in the WDC dataset, but it has 924 NaN-values and 103 unique categories in Depth 3. In contrast, the Icecat dataset essentially contains no NaN-values into Depth 1 (17 unique categories) as well as Depth 2 (231 unique categories), but it suffers in Depth 3, where there are 452,976 NaN-values and 161 unique categories. Thus, the two data sources vary with respect to the structure of data, their versatile hierarchy depth across the different categories.

The distribution of Depth 1 categories in WDC and Icecat datasets shows interesting dif-

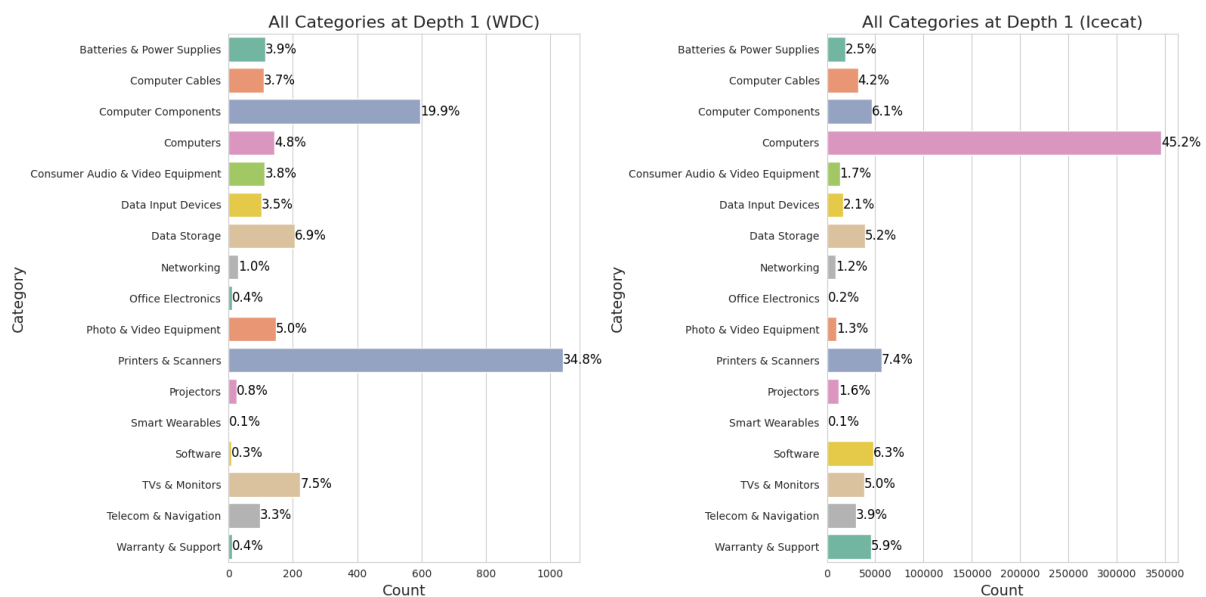


Figure 2: Categories on the first depth in icecat and wdc datasets

ferences regarding the product types they offer (Figure 2). WDC has a very high concentration on "Printers & Scanners," with this category representing 34.8% of its products (1,038 units), followed by "Computer Components" at 19.9% (595 units) and "TVs & Monitors" at 7.5% (223 units). However, there are some categories, such as "Smart Wearables" and "Office Electronics," which do not have much representation in the dataset. On the contrary, Icecat's

Group Part: Enhancing Product Classification with LLMs

distribution is clearly dominated by "Computers," accounting for 45.2% (346,087 units) of its offerings, leaving "Printers & Scanners" with only 7.4% (56,404 units). The difference is indicative of how each dataset is skewed towards its unique categories.

Derived from the difference skewness in Depth 1, the Depth 2 distributions of both datasets

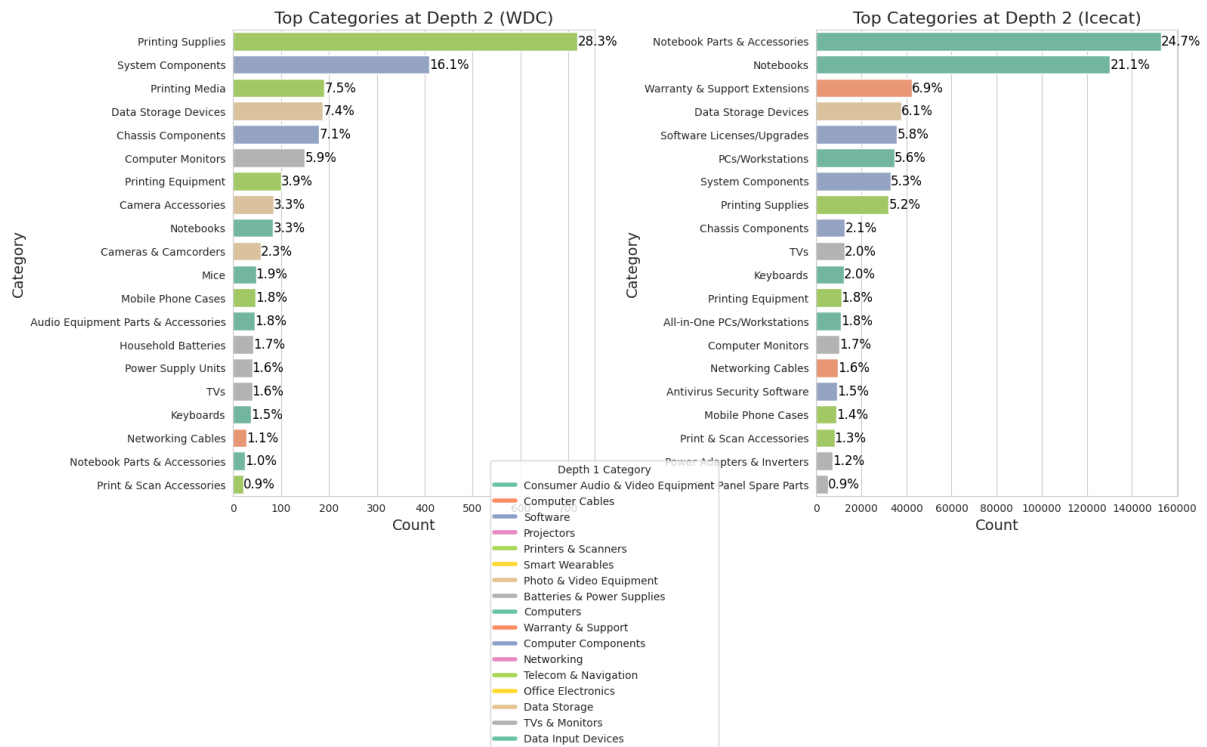


Figure 3: Top Categories on the second depth with the coloring of the first depth

demonstrate quite different structures at depth 2. Of the total categories in WDC, "Printing Supplies" is the most common, with 28.3% (719), followed by "System Components" at 16.1% (409) and "Printing Media" at 7.5% (191). Further categories of in WDC at level included "Data Storage Devices" and "Chassis Components". On the other hand, the prevalent categories at Level 2 in Icecat are "Notebook Parts & Accessories," with 24.7% (152,761) and "Notebooks" with 21.1% (130,294), highlighting the different distributions.

Consequently, also on depth 3 the two datasets contrast greatly. Among all categories in WDC, "Ink Cartridges" alone comprises 25.0% (516), followed by "Motherboards" at 9.3% (191), while "Toner Cartridges" makes up 7.7% (158). The distribution is leaning towards printing and computer supplies. In contrast, Icecat's largest category at level 3 is "Notebook Spare Parts" with 46.4% (145,020), followed by "Toner Cartridges" and "Internal Hard Drives". Once again, standing in stark contrast to WDC. Concluding the categorical analysis, WDC data

Group Part: Enhancing Product Classification with LLMs

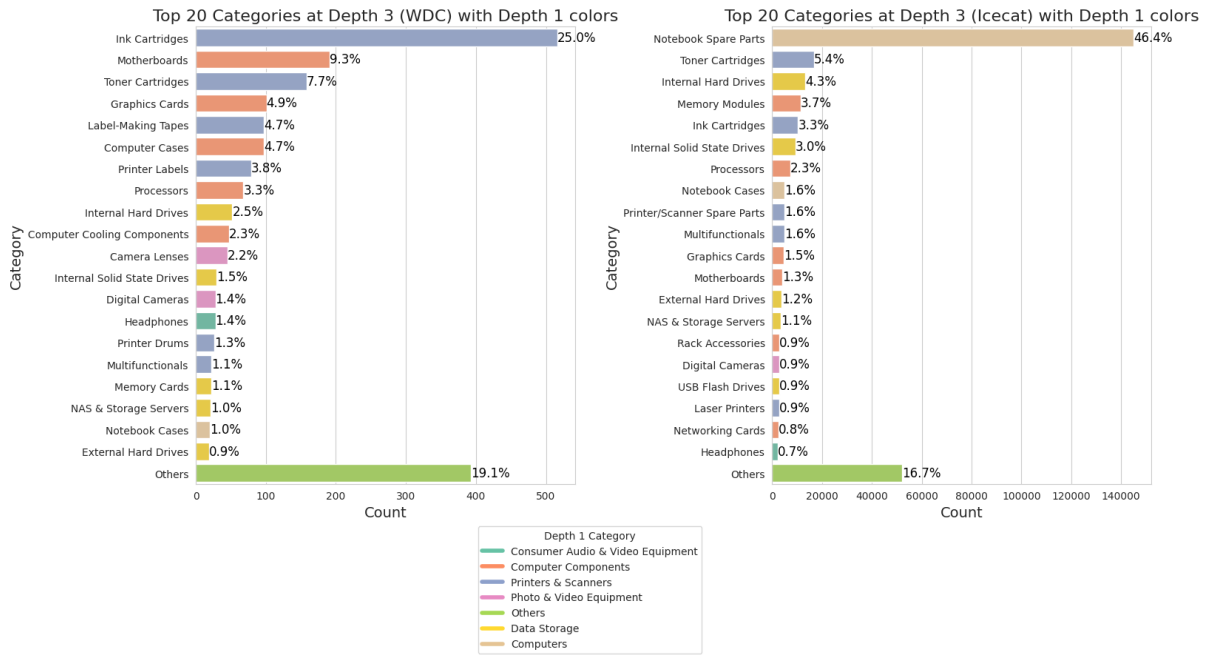


Figure 4: Top Categories on the third depth with the coloring of the first depth. is skewed toward printing products while Icecat leans more on portable computing accessories, indicating how distinct product focuses can arise in different data sources.

4.3 Product Title Analysis

Distribution of Title Lengths:

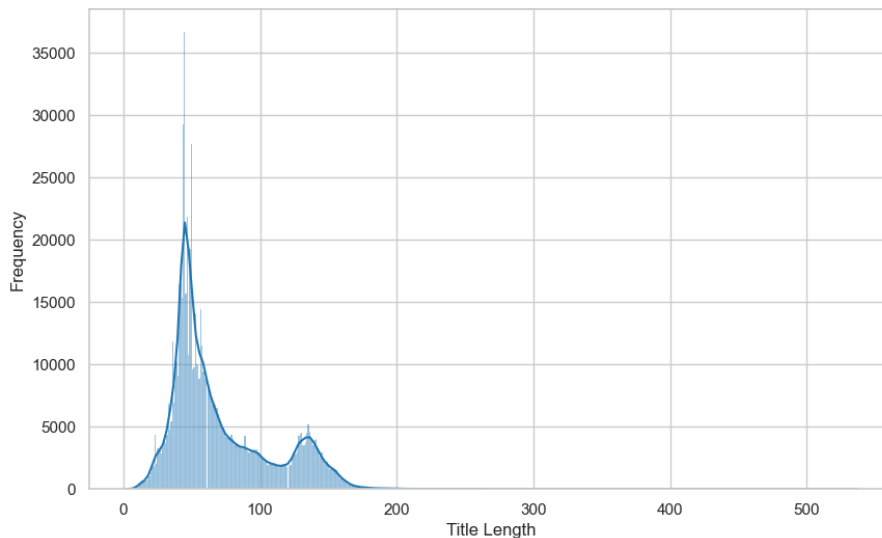


Figure 5: Distribution of Product Title Length Icecat

The analysis of Icecat data reveals that the lengths of product titles are heavily skewed toward shorter titles, with a significant peak observed in the 40-60 character range. With 765,473 data

Group Part: Enhancing Product Classification with LLMs

points, the mean title length is 70.42 characters, with a corresponding standard deviation of 36.26, suggesting some deviation of longer title lengths. The median length is 57 characters, and most titles are found to fall within the quartile range of 45 to 90 characters. Maximum title length interestingly reaches up to 537 characters, a substantial outlier from the average. These findings would suggest that most product titles tend to be short but that there is quite a wide range of lengths, with some few significantly longer outliers.

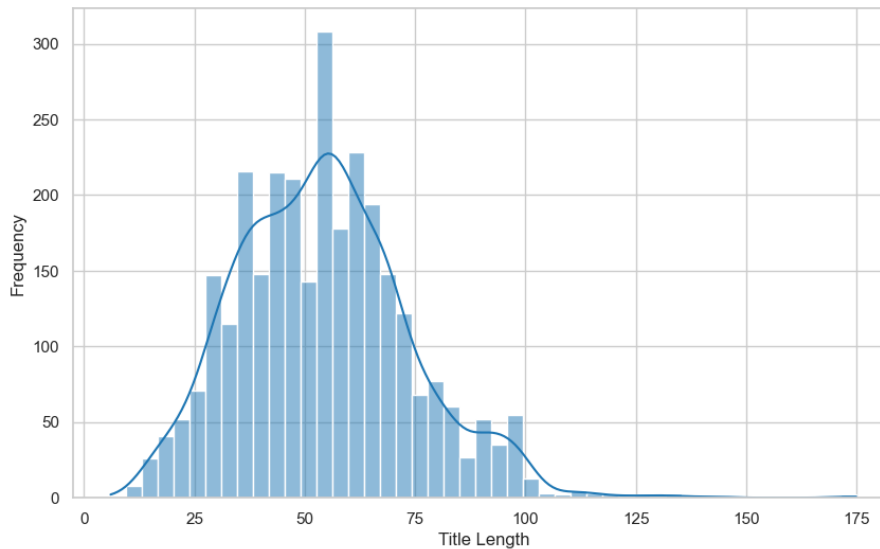


Figure 6: Distribution of Product Title Length WDC

Compared to Icecat, the title lengths of products in the WDC dataset are distributed more symmetrically. There are 2,984 titles. The average title has 54.32 characters with a standard deviation of 19.66, indicating a tighter spread around the mean. The median title length falls at 54 characters, which shows its close affinity to the mean, thus implying that the data is fairly normally distributed, with its mean roughly at the center around 50 to 60 characters. Most title lengths lie in the range between 40 to 66 characters, while the longest title has 175 characters, marking a much smaller variation in the upper range when compared to Icecat data. This further implies that the titles in the WDC dataset are generally shorter and more concise.

Frequent Words

Insights into common phrases or nomenclatures found in product descriptions are the core of this analysis. The study of the frequent words in Icecat product titles. Some of the most used words are "gb," "notebook," "black," and "cm," words that clearly point toward technology

Group Part: Enhancing Product Classification with LLMs

relevant terms, specifications, and descriptive words. Terms like "intel," "sdram," "hdd," and "i7" are related to hardware specifications. Interestingly, when brand names and numbers are eliminated from product titles, the remaining most common words remain related to technical specifications, such as "ssd", "pc", and "display". Domain specific terms and numerical characters are important provide models valuable insights, helping the categorization process.

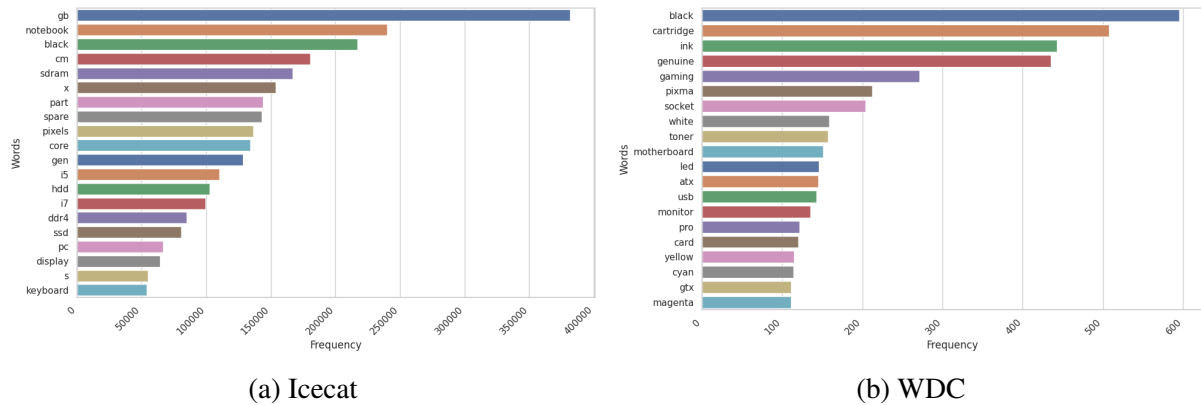


Figure 7: Frequent Words in Product Titles

Results of the analysis of common words in WDC product titles show that they lean heavily on products related to the printing, gaming, and electronics industries. These include common words such as "black", "cartridge", "ink" and "genuine" with a strong dependency on printer products, especially ink and toner cartridges. Other frequent terms such as "gaming", "socket", "motherboard" and "atx" show a focus on computer hardware topics, primarily gaming systems and assembly components. The last terms are more evidence of the visual display technology and connectivity importance, namely "led", "monitor" and "usb." Some of the words also contain "pixma" (from Canon brand) and "pro", meaning that there are also consumer-grade and professional items. This spread of terms is in line with the categories included within WDC. There is a heavy skew toward technology, gaming, and office supplies.

4.4 Description Analysis

WDC: 25,87% of the descriptions are empty. The most frequent words include "summary", "black", and "product". Some numbers like "1" and "2" also appear repeatedly. In some cases, there are multiple observations with the same gtn but with varying descriptions. The most frequent length is 28 characters, the following histogram shows the distribution of descriptions lengths:

Group Part: Enhancing Product Classification with LLMs

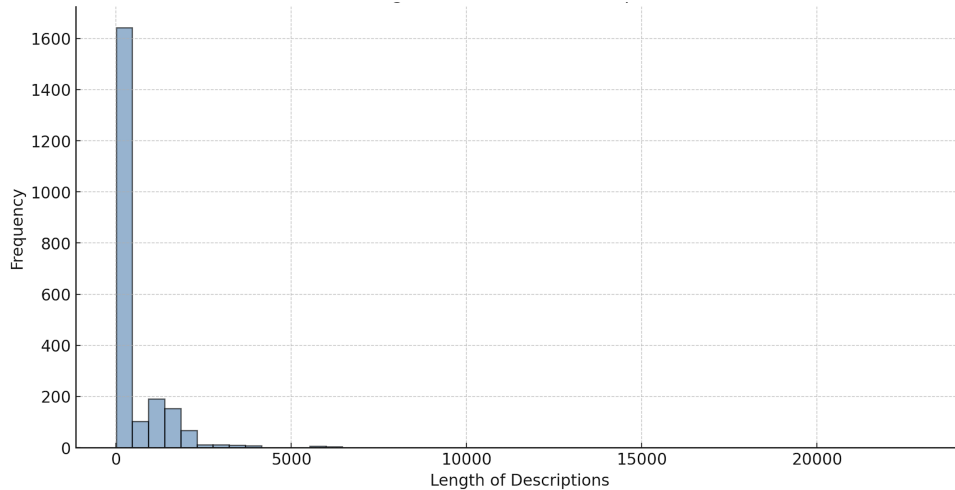


Figure 8: Length Distribution on WDC

The most frequent bigrams in product descriptions highlight technical specifications and features, such as "sf utp" (likely a cable type), "intel xeon", and "xeon processor", highlighting the importance of processor-related terms. Other notable examples include "photo paper" and "533 MHz," indicating specific product categories and technical attributes. Descriptive terms like "summary description" and "auto generated" suggest how some descriptions are structured. In trigrams, "intel xeon processor" is the most frequent. Other common trigrams, such as "cat 5e sf" and "5e sf utp," relate to networking cables, while "processor 06 GHz" and "GHz 512k cache" highlight processor details. Trigrams like "premium a4 70x33" focus on specific paper sizes and formats.

Icecat: The description analysis focus on the columns *SummaryDescription*, *LongSummaryDescription* and *SummaryDescription*. *ShortSummaryDescription* as there are no missing values, and it contains the most representative descriptions for product categorization.

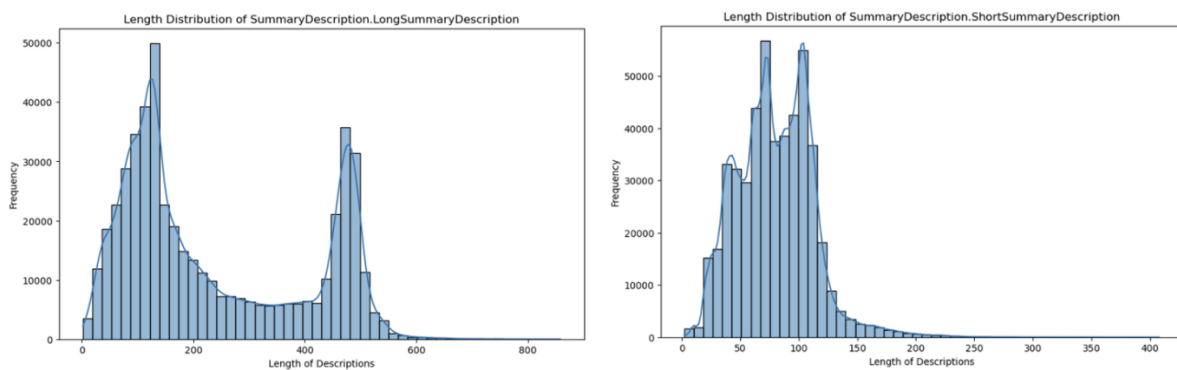


Figure 9: Length Distributions on Icecat Long and Short Descriptions

The length distribution of descriptions shows two distinct patterns for both long and short

Group Part: Enhancing Product Classification with LLMs

summaries. The long summary descriptions exhibit a bimodal distribution with peaks around 150 and 400 characters. The majority of the descriptions cluster in these two ranges, indicating that there are mainly two different products descriptions with varying lengths, possibly depending on the product type or complexity of the description. The short summary descriptions primarily range between 50 and 100 characters, with a much more noticeable bias toward shorter lengths, making them more standardized in format compared to the long summaries.

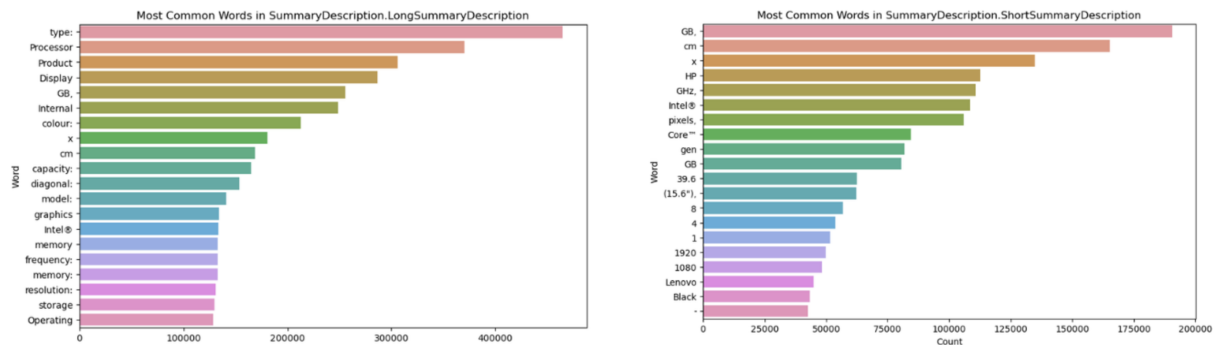


Figure 10: Most common words in Icecat

Most common words: The bar charts illustrate the most common words in both long and short summary descriptions. In the long descriptions, technical specifications such as "Processor", "Product", "Display", "GB", and "Internal" dominate. Similarly, the short descriptions frequently include terms like "GB", "cm", "x", and "HP", which are also focused on product characteristics and technical specs.

4.5 Brand Analysis

The WDC dataset lacked a dedicated brand column, so brands were extracted by taking the set of unique brands present in the Icecat data. WDC titles were then searched for instances of this list, extracted and assigned to the observation's Brand. Within the extracted brands (Fig. 11), the distribution is highly imbalanced, with Canon and Epson dominating Printers & Scanners, while Asus, MSI, and Samsung lead in Computers and Components.

Group Part: Enhancing Product Classification with LLMs

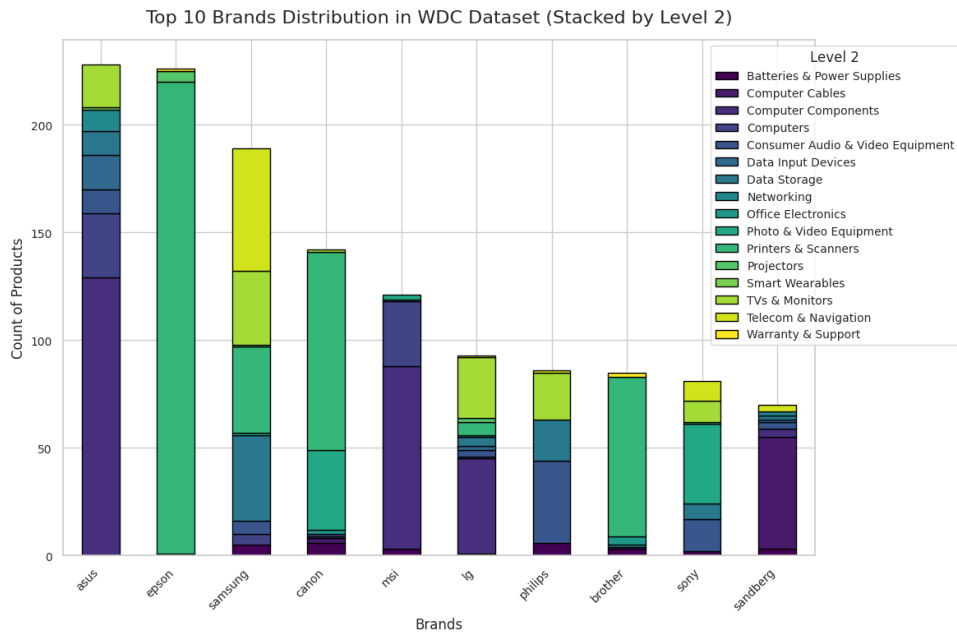


Figure 11: Top 10 Most Frequent Brand in Categories - WDC

Unsurprisingly, with Icecat’s leaning towards ”Computers”, HP overwhelmingly dominates the brand distribution. Lenovo, Acer, and ASUS follow, but with much lower representation. This is highlighted in (Fig. 12) Icecat’s brands consist of key market leaders in computing and warranty-related categories, creating a skewed distribution of brands.

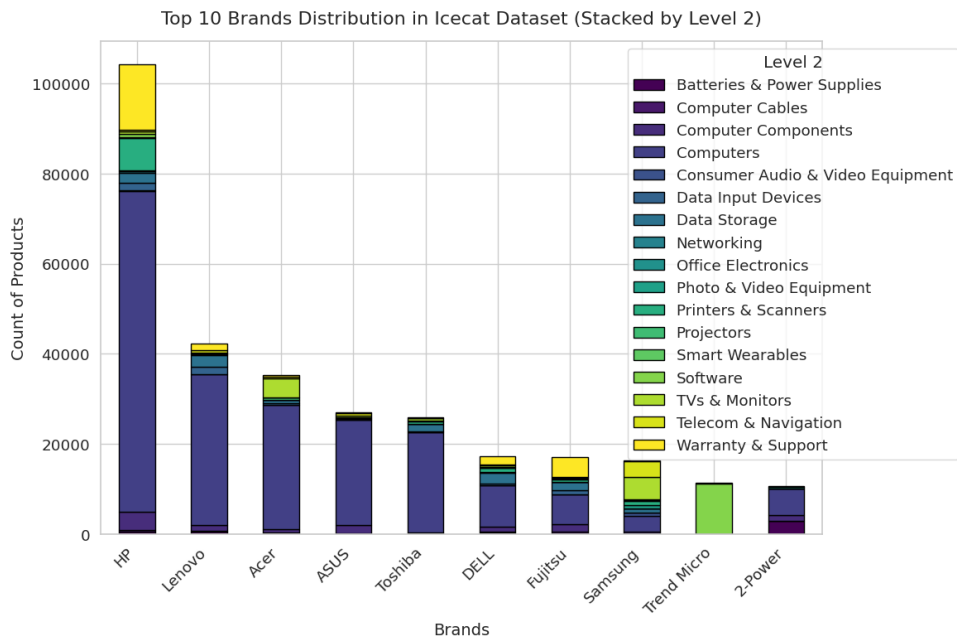


Figure 12: Top 10 Most Frequent Brands in Categories - Icecat

This contrast underscores differing orientations, which are in line with findings from the hierarchical and categorical analysis.

5 Methodology

This paper explores the enhancement of product classification using Large Language Models. To achieve this, five promising approaches are analyzed in detailed deep dives, each focusing on different strategies. All deep dives utilize the same datasets, evaluation metrics, and, in some cases, the same models, which are briefly introduced below. The aim is to conduct a comprehensive comparison, drawing conclusions about which models and approaches perform better in specific scenarios.

5.1 Data Preparation

Different attribute combinations were explored, but for the final models only the title was used. Additionally, the Path and Category columns were utilized for evaluation.

Title	Initial Path	Category Name
ASUS K31CD-IT049T PC 6th gen Intel® Core™ i7 i7-6700 16 GB DDR4-SDRAM 1000 GB HDD Black Tower	Computers & Electronics > Computers > PCs/Workstations	PCs/Workstations
HP 686915-A41 notebook spare part Keyboard	Computers & Electronics > Computers > Notebook Parts & Accessories > Notebook Spare Parts	Notebook Spare Parts
C2G 1m ST/SC Plenum-Rated 9/125 Duplex Single-Mode Fiber Patch Cable - Blue fiber optic cable 39.4" (1 m)	Computers & Electronics > Computer Cables > Fibre Optic Cables	Fibre Optic Cables

Table 3: Three sample rows showcasing the features used in the Icecat dataset.

Further preparation involves exploding the category path to retrieve the actual categories per level. In cases where Level 3 does not exist (i.e., the path has only three entries), the category from Level 2 is propagated to Level 3. This approach is commonly used in hierarchical classification tasks where the depth of categories may vary across different paths (Guo and Zhao, 2021). Lastly, the first-level categories are omitted from the datasets, as they consistently correspond to 'Computers & Electronics' for every observation, making them redundant.

Initial Path	Adjusted Path	Category Level 1	Category Level 2 & Level 3
Computers & Electronics > Computers > PCs/Workstations	Computers > PCs/Workstations > PCs/Workstations	Computers	PCs/Workstations

Table 4: Example row illustrating the Adjusted Path and its derived category levels.

This paper focuses exclusively on the title as input feature, rather than including the description or other attributes, in order to minimize the number of input tokens for the LLM. As highlighted

in Groß et al. (2024), the title contains the most relevant information and is the most important feature for making accurate predictions. Discrepancies in available attributes pose significant challenges in e-commerce Martinek et al. (2024). However, compared to other attributes, product titles generally seem to be consistently available (d’Amato et al., 2018), a pattern supported by the benchmark datasets Icecat and WDC.

Title	Description
acer aspire e1 522 65208g1tmnk specifications	(NaN)
msi gs73vr 7rf 207uk nvidia gtx 1060 6gb gddr5 17 3 uh d ips 4k intel i7 7700hq gaming laptop	gs73vr 7rf stealth pro 4k 207uk 17 3 uh d ips 4k 3840 x 2160 i7 7700hq nvidia gtx 1060 6gb gddr5 256gb ssd 2tb 7200rpm hdd 16gb 2400mhz ddr4 windows 10 home 64 bit 2yr warranty

Table 5: Example Rows from the input WDC-222 dataset, showcasing the inconsistency of Description

5.2 Large Language Models

For the purpose of research of individual LLM enabled techniques (Paper A-E) various LLMs were implemented in experiments. The following provides a non-exhaustive overview of state-of-the-art open-source LLM families referenced in the papers.

5.2.1 Pythia

Pythia is a suite of decoder-only LLMs designed to facilitate analysis across varying training scales (Biderman et al., 2023). It provides eight open-source GPT-NeoX-based models at different sizes, each accompanied by precise training states and hyperparameters. This transparency enables researchers to observe how scaling, optimization, and data affect performance. By offering well-documented internals and standardized transformer-based components, Pythia fosters reproducibility and improves our understanding of language model behavior.

5.2.2 LLaMa

LLaMa (Large Language Model Meta AI) (Touvron et al., 2023) is a family of foundational models (1B–405B parameters) that consistently outperformed competitors upon release. Its latest version, LLaMa 3.2 (November 2024), includes lightweight autoregressive models (1B and 3B parameters) for on-device inference. LLaMa’s architecture incorporates RMSNorm

instead of standard normalization, replaces ReLU with SwiGLU for enhanced performance, and uses rotary positional embeddings instead of absolute positional embeddings.

5.2.3 Mistral

Mistral 7B (Jiang et al., 2023a) is designed to rival larger models like LLaMA-2-13B with about half the parameters. Achieving efficiency through re-architected transformer blocks and improved attention patterns, Mistral optimizes training pipelines, parameter initialization, and scaling rules. These innovations allow higher-quality performance with fewer resources, making advanced LLMs more accessible. Its latest version, Mistral-7B-v0.3, is available on the HuggingFace Hub.

5.2.4 Instruct Models

Instruct models are fine-tuned to better follow user instructions, often through supervised tuning or Reinforcement Learning from Human Feedback (RLHF). Unlike base models, which are trained broadly and remain open-ended, instruct variants focus on producing responses aligned with specific prompts and guidelines. This makes them more suited for practical tasks employing prompting and in-context learning like information retrieval, and content generation under defined constraints, ensuring outputs are both accurate and useful. Both LLaMa 3.2 and Mistral-7B-v0.3 exemplify this paradigm by offering two open-source variants: a “base” model and an “instruct” model.

5.3 Topic Deep Dives

5.3.1 Fine Tuning LLMs (Paper A)

To investigate the effectiveness of fine-tuning for hierarchical classification, five pre-trained LLMs of varying sizes were fine-tuned on the same dataset using four distinct fine-tuning strategies. Due to computational constraints and evidence suggesting that pre-trained language models perform better on balanced datasets (Weiss and Provost, 2001), the highly imbalanced Icecat dataset was subsampled. A clustering-based subsampling technique (Yen and Lee, 2009)

was employed to preserve the original variance while reducing the training dataset to approximately 30,000 observations. The resampled dataset was subsequently divided into training, validation, and test subsets. For efficiency, all models were trained exclusively on product titles, as shorter inputs accelerate training time (Vaswani et al., 2023) and titles alone have been shown to provide sufficient context for accurate classification (Dai et al., 2020).

The models selected for this analysis, listed in ascending order of size, include Pythia70m, Pythia410m, LLaMa3.2-1B (base and instruct), LLaMa3.2-3B (base and instruct), and Mistral-7B-v0.3. To assess the impact of hierarchical information during training, all models were initially fine-tuned for Flat Classification, where the objective was to generate a single phrase corresponding to the level 3 category. Subsequently, three hierarchical fine-tuning approaches: Hierarchical Completion, Hierarchical Instruct, and Curriculum Learning were applied, as described in detail in Paper A. In these hierarchical setups, the models were tasked with unrestricted causal language modeling to generate sequences of hierarchical labels. Notably, the relatively small Pythia models underwent full fine-tuning, while larger models were fine-tuned using QLoRA, a parameter-efficient fine-tuning technique that significantly reduces the number of trainable parameters while maintaining strong inference performance.

5.3.2 Different Prompting Techniques (Paper B)

This section explores how prompting techniques can be effectively utilized to enhance performance and scalability in hierarchical product classification, a critical challenge in e-commerce. By exploring how the advanced prompting techniques chain-of-thought (CoT) and prompt chaining can excel in HPC, the study uncovers their potential to navigate complex taxonomies and improve classification across multiple hierarchical levels.

Results from this exploration will provide a comparative analysis of developed techniques, offering a roadmap for using optimized Techniques to perform HPC with LLMs. This has significant implications for improving e-commerce workflows, enabling rapid deployment and accurate product categorization while maintaining cost-effectiveness. Future directions include combining methods and refining domain-specific LLM fine-tuning to enhance model capabilities further.

5.3.3 Knowledge Graphs (Paper C)

This paper explores the integration of taxonomy-specific Knowledge Graphs (KGs) into LLMs for hierarchical product classification in e-commerce. KGs represent knowledge in the form of structured triplets (entity – relation – entity), enabling symbolic reasoning. When integrated with LLMs, KGs provide structured, factual external knowledge that not only enhances the precision and interpretability of LLM predictions but also contributes to more grounded reasoning. LLMs, in turn, support the expansion and enrichment of KGs through their advanced natural language processing capabilities (Chen et al., 2024a). This synergy seeks to mitigate the challenges of LLM rigidity and hallucinations by leveraging the interpretability and factual consistency of KGs.

In this study, multiple KG integration techniques are evaluated, ranging from representing only the taxonomy as a KG to extending the taxonomy with new entity nodes using Named Entity Recognition (NER) methods. These approaches are tested across various LLMs using the WDC-222 Gold Standard dataset. All experimental configurations follow the Local Classifier per Parent Node (LCPN) approach and are compared against hierarchical and flat baselines without KG augmentation. The hypothesis posits that taxonomy-specific KGs will enhance both the interpretability and precision of LLMs in hierarchical product classification tasks.

5.3.4 Retrieval Augmented Generation (Paper D)

Retrieval-Augmented Generation (RAG) systems enhance Large Language Models (LLMs) by integrating external knowledge about similar products into the classification process, thereby improving accuracy and relevance. By leveraging RAG, the system can retrieve information which serves as a foundation for more informed and contextualized categorization.

The analysis begins with a comprehensive study of configurations that optimize RAG performance. Various retrieval methods are explored, including standard retrieval techniques, graph-based approaches, and an innovative multi-tiered retrieval framework. Each method is assessed and compared in fetching relevant products, which is important for the classification pipeline.

The approach then, examines the use of retrieval similarity as a standalone method for categorization, independent of LLMs, serving as a baseline for comparison with more advanced

RAG systems. These systems include NaiveRAG, which represents a straightforward implementation, AdvancedRAG, which introduces refined retrieval and classification strategies, GraphRAG, which leverages graph-based retrieval for enhanced contextual understanding, and HybridRAG, a combination of techniques including graph and vector retrieval. To build the prompt different prompting techniques such as Zero-shot, Few-shot, CoT, Prompt Expansion and Prompt chaining are taken into consideration. Overall, the study explores how these RAG systems can be effectively applied to both, flat and hierarchical step-by-step categorization processes.

Overall, this study emphasizes the evolution of RAG systems from basic to advanced configurations, highlighting their potential for product classification.

5.3.5 Entity Matching (Paper E)

Entity matching, often referred to as product matching in the e-commerce context, could play a crucial role in hierarchical product classification. At its core, entity matching aims to identify identical or closely related entities across datasets, even when they are described differently. In e-commerce, this translates to finding the same or similar products listed under varying titles, descriptions, or formats by different vendors.

The process consists of two main steps: blocking and matching. Blocking reduces computational complexity by identifying the k most similar products, narrowing the search space to relevant matches. The matching step identifies the most similar product (based on the title in this study) and uses this "product match" to predict the input's category based on the matched product's category. Similar to methodologies from *Paper 4* and RAG, top- k blocking and similarity scores are used as the baseline, excluding LLMs. LLMs are then integrated on top of this baseline to assess whether LLM-based entity matching can improve performance and achieve comparable results to other approaches tested in this paper.

To reiterate: in contrast to the other four approaches, the LLM is not tasked with directly predicting a category. Instead, it identifies the most similar product in the reference dataset, which then implicitly determines the category.

5.4 Evaluation

The different model performances will be analyzed to draw conclusions about model or dataset characteristics, ultimately enabling the formulation of general insights and recommendations.

5.4.1 Performance Metrics

The WDC-222 publication evaluates several machine learning models using a variety of metrics. Specifically, they employ level-specific weighted Precision (wP), weighted Recall (wR), and weighted F1 (wF) scores at each level of hierarchical prediction for both the Icecat and WDC-222 evaluation datasets. When referring to these metrics, by default, level 3 metrics are meant. Additionally, they implement the hierarchical F1 (hF), precision (hP) and recall (hR) scores (Kiritchenko et al., 2006), which accounts for partial hierarchical predictions, that are easily implementable thanks to HiClass python package (Miranda et al., 2023). In the experiments performed in this paper, we apply the same metrics, as they provide a robust estimation of model performance in real-world scenarios characterized by highly imbalanced datasets, which is typically the case for hierarchical classification tasks. Furthermore, the evaluation of wF and hF scores offers insights into the models' ability to "memorize" the hierarchical structure. The hF score is also calculated in flat scenarios.

Additionally, we have developed a metric to assess the coherence of label name representations in the model outputs. The "Phrase Coverage" (or "Label Coherence") metric quantifies the proportion of classifications at each hierarchical level that correctly identify the label name from the subset of available labels for that level. This metric enables the evaluation of the model's ability to memorize label names at each level, thereby enhancing its understanding of the label hierarchy.

5.4.2 Models' Performance Analysis and Comparison

To evaluate the performance of each approach, we employ a comprehensive set of metrics, fully detailed above. Level-specific metrics reflect how accurately each model identifies correct labels at increasing granularity, while hierarchical metrics measure the severity of misclassifications by considering their placement within the taxonomy. Thus, even when errors occur,

those that remain close to the true class have less impact on the final score. In combination, these metrics provide a balanced evaluation. wF and hF highlight both fine-grained precision and overall taxonomic coherence. Precision and recall offer further insight: models with high precision rarely predict incorrect classes, while those with high recall rarely miss correct ones. By examining all these metrics together, we gain a nuanced understanding of each model's strengths and weaknesses in hierarchical product classification.

5.4.3 Misclassification

A comprehensive analysis of misclassifications is conducted across all the explored approaches, with a focus on high-representation categories such as "Printers & Scanners," "Computer Components," and "TVs & Monitors". Evaluations take place at different hierarchical levels, emphasizing leaf-level (Depth 3) predictions, where overlapping categories like "Ink Cartridges" and "Printer Ink Refills" present significant challenges. Misclassifications are logged and analyzed for trends, highlighting shared challenges across approaches, such as systemic errors in closely related categories. This analysis is crucial for identifying the limitations of classification models and understanding the root causes of errors, enabling targeted improvements in both model performance and taxonomy structure. These insights inform potential taxonomy improvements, including redefining overlapping categories and introducing intermediate subcategories to enhance classification accuracy and reduce ambiguity.

5.4.4 Imbalanced Classes

In this section, the imbalanced class performance across various individual models is analysed. As shown in the EDA section (Categorical Analysis 4.2), the WDC and Icecat datasets are highly imbalanced. The importance of being able to perform in imbalanced scenarios is highly relevant to real-world applications. Product category trees, in particular, are often characterized by such imbalances, making it crucial to develop models that can handle these challenges effectively. For the underrepresented classes, evaluations were performed on a subset of the 100 least-represented categories in the WDC dataset, resulting in 121 observations, with an average of 1.21 observations per class. For the overrepresented classes, the top 5 largest categories were

selected, resulting in 1,213 observations, with an average of 243 observations per class. This further highlights the highly imbalanced nature of the dataset.

5.4.5 Label Coherence Analysis

LLMs, due to their broad training data, often encounter “LLM hallucination,” a phenomenon in which the model produces text closely resembling accurate content yet lacking factual correctness (Huang et al., 2024). This issue is particularly relevant for hierarchical classification tasks reliant on LLM outputs. As noted by Risch et al. (2020), treating hierarchical text classification (HTC) as a sequence-to-sequence (seq2seq) task and incorporating the semantic meaning of labels yields better results than generating encoded label signatures. Consequently, it is crucial that LLM-based solutions produce coherent label names both at a flat level and across hierarchical tiers, since even minor deviations, such as subtle wording differences or typos, can severely affect classification accuracy.

Several studies investigating HTC as a seq2seq task have implemented mechanisms to constrain the model’s vocabulary, ensuring more coherent label generation (Torba et al., 2024; Yu et al., 2022; Jain et al., 2024). The solutions presented in this paper, however, do not rely on such constraints, instead assuming that LLMs, particularly when guided through in-context learning, can inherently produce coherent labels. To assess label coherence, each model’s predictions at all hierarchical levels are compared against established label lists, and the proportion of exact matches is calculated. This coherence metric is critical, as it reflects the model’s ability to produce the correct tokens for each label. Any deviation directly diminishes classification performance, as in real-world scenario products without a coherent label are not assigned to any class.

5.4.6 Generalization

Generalization is crucial for product categorization, helping models perform well across diverse datasets and tackle real-world challenges like data variety and multiple languages. The WDC-222 benchmark dataset is already seen as “unstructured” compared to the cleaner Icecat dataset, making it a good indicator for model generalization. However, we wanted to take this even

further, by testing different languages and title structures. This approach addresses WDC-222's English-only limitation and showcases LLMs' ability to handle unseen input effectively.

To extend the generalization analysis, we created a test dataset consisting of 30 inputs sourced from amazon.de. The sample size is small and does not ensure statistical significance, but it provides an initial indication of model performance on real-world data. We selected following categories from the Icecat taxonomy, representing overrepresented as well as underrepresented classes: Computer Monitors, Ink Cartridges, Digital Photo Frames, Sport Watches, Smartwatches, Activity Trackers and Soundbar Speakers with their respective Paths.

We manually identified Amazon products matching these categories within the Electronics section. To simulate multilingual and structurally diverse datasets, we retrieved titles in multiple languages, by switching the interface language from German to Turkish, English, and Polish. While multilingual product titles on Amazon may come from two sources: vendor-provided translations or "just" machine-generated translations, this approach allowed us to capture a realistic diversity of input data while maintaining consistency within the Electronics category.

Given the differences in model architectures and training approaches, we expect varying levels of generalization across the tested methods. Specifically, we hypothesize that methods with stronger capabilities for handling multilingual and structurally diverse inputs will demonstrate better performance on this new input set. This experiment will highlight the relative strengths and weaknesses of different methods when faced with variations in title structure and language.

5.4.7 Complexity Analysis

The tested LLM approaches differ in their complexity, training data requirements, inference efficiency, and levels of rigidity and adaptability. This section focuses on a theoretical and conceptual discussion of these factors, highlighting the trade-offs and considerations inherent in each method.

6 Results

The solutions chosen for the comparison in this section represent the most suitable methods based on experimental results from Papers A-E. In the following sections, the best-performing approach derived from the individual papers will be referred to as Fine-tuning, Prompting, KG, RAG and Entity-matching.

Method	Description	Model Used
Fine-tuning	Comprises a Mistral-7B-v0.3 LLM fine-tuned with QLoRA for accurate generation of hierarchical text sequences based on product titles as inputs.	Mistral-7B-v0.3
Prompting	A multi-stage top-3 solution that iteratively filters the taxonomy to identify the best-fitting leaf nodes using LLM-based prompting techniques.	Mistral-7B-v0.3
KG	Combines predictions of three weak LLMs. Differences in predictions trigger a re-prompt where taxonomy-specific KGs are injected. KG-augmented predictions are majority voted; if unresolved, an advanced LLM acts as judge.	Mistral-7B-v0.3-Instruct, Llama 3bn, 8bn, 70bn
RAG	A flat Hybrid RAG that combines Knowledge Graphs and vector retrieval. Cohere Reranker v3 reorders and selects the most relevant documents.	Mistral-7B-v0.3-Instruct
Entity Matching	Combines Sentence Embedding + Cosine Similarity as top-k blocking, with Google Flan T5 Large serving as classifier on top using the LLM Select strategy.	Google Flan T5 Large

Table 6: Comparison of Final and Best Performing Models

6.1 Performance Metrics

Method	F1 Level 1	F1 Level 2	wF	hF	wR	wP	hR	hP
Fine-tuning	0.94	0.86	0.80	0.87	0.77	0.88	0.87	0.87
Prompting	0.81	0.74	0.69	0.76	0.65	0.79	0.75	0.77
KG	0.88	0.81	0.76	0.83	0.74	0.83	0.83	0.83
RAG	0.95	0.90	0.86	0.92	0.86	0.88	0.92	0.91
Entity Matching	0.92	0.84	0.77	0.84	0.73	0.88	0.84	0.84

Table 7: Performance Metrics for Different Methods.

The results demonstrate a clear hierarchy in performance among the evaluated approaches. The RAG-based technique emerges as the strongest performer, achieving an wF of 0.860 and a hF of 0.915, thereby indicating superior accuracy at the finest level of product classification as well as exceptional consistency across the entire category structure. Meanwhile, the fine-tuning-based model attains the second-best results (wF = 0.804, hF = 0.874), suggesting that direct

parameter adaptation with domain-specific data remains a reliable strategy to capture complex labeling schemes.

In contrast, the derived prompting framework yields a more modest outcomes ($wF = 0.685$, $hF = 0.7634$), illustrating that relying solely on LLM prompts is less effective in handling the nuanced hierarchy of product categories. Intermediate performances emerge from approaches that leverage knowledge graphs or entity matching techniques (e.g., KG: $wF = 0.759$, $hF = 0.827$; Entity Matching: $wF = 0.771$, $hF = 0.839$), indicating that while such strategies enhance the model's understanding of product categories, they do not fully achieve the performance levels observed with fine-tuning or RAG.

The relationship between the hierarchical and fine-grained metrics also provides meaningful insights. The higher hF scores, relative to wF , underscore that many "mispredictions" are actual "near-misses", often placing items within closely related categories that share the same parent node. This insight suggests that while some misclassifications occur, they often aren't severe in a hierarchical sense, reflecting that the model's understanding of category groupings is reasonably coherent even when it makes mistakes at the finest granularity.

Examining precision and recall patterns further refines these insights. The RAG model shows balanced precision and recall at Level 3 (0.882 and 0.861) and maintains this balance hierarchically (0.912 and 0.918), indicating comprehensive and accurate predictions. By contrast, finetuning favors precision (0.881) over recall (0.773) at Level 3, yet its hierarchical scores ($hPrecision = 0.874$, $hRecall = 0.874$) show that though the model's Level-3 recall lags behind its precision due to strict final-category matching, its balanced hierarchical scores indicate it often places instances on the correct branch of the taxonomy, reflecting partial but meaningful correctness. Prompting, on the other hand, struggles more with recall at Level 3 (0.6501), but the hierarchical perspective ($hRecall = 0.753$) shows its errors still lie within neighboring classes. The KG approach, while achieving moderate, maintains fairly balanced hierarchical precision (0.829) and recall (0.826), indicating it rarely veers far from the correct path. Similarly, Entity Matching achieves a respectable balanced hierarchical recall (0.838) and precision (0.840), reflecting that even when not entirely accurate at the deepest level, predictions remain close to the correct categories. These patterns highlight the value of hierarchical metrics and

confirm that significant portion of misclassifications remain near the correct category node.

6.2 Label Coherence

Model	Level 1	Level 2	Level 3
Fine-tuning	100.00	97.59	95.24
Prompting	96.82	96.82	96.82
KG	99.97	98.83	97.89
RAG	99.87	99.66	99.10
Entity Matching	100.00	100.00	100.00

Table 8: Label Coherence (%) of Models.

The comparison of label coherence across different techniques reveals consistently high performance at all hierarchical levels. Entity Matching achieves perfect coherence because it does not rely on label generation by the model; instead, it selects from existing labels based on LLM-driven entity matching. Other methods also maintain a high degree of coherence, even at the most granular classification levels, ranging from 95.24% (Fine-Tuning) to 99.10% (RAG). The strong coherence of Prompting (96.82%), KG (97.89%), and RAG (99.10%) at level 3 can be attributed to the presence of the target labels directly in the prompt, facilitating accurate label assignments. Notably, the Fine-Tuning model’s 95.24% coherence, despite the prompt consisting solely of the product titles, demonstrates its ability to generate appropriate label names without additional cues.

It is important to note that any shortfall from full coherence directly reduces overall performance. For example, Fine-Tuning’s 4.76 percentage-point gap from perfect coherence translates into a corresponding reduction in its performance metrics, reflecting the model’s diminished capacity to generate entirely correct predictions, ultimately leaving critical predictions unlabeled—a costly inefficiency in any business-driven environment.

6.3 Misclassification

The misclassification analysis (Appendix - 1.1.5) reveals consistent challenges across all approaches, particularly for high-representation categories like "Printers & Scanners," which accounts for the most misclassifications in every method (e.g., 172 for Fine-Tuning, 260 for

Prompting, 220 for Knowledge Graph, and 320 for Entity Matching). Categories with overlapping characteristics, such as "Computer Components" and "TVs & Monitors," also exhibit significant misclassification rates, indicating difficulties in distinguishing between closely related product types. Notably, approximately half of all misclassified data points were predicted incorrectly by at least two approaches, emphasizing the shared limitations and challenges across different methods.

Entity Matching, while effective for smaller categories, produces the highest misclassification rate in "Printers & Scanners," reflecting limitations in handling nuanced distinctions. Prompting and Fine-Tuning distribute errors more broadly but still struggle in high-frequency categories. In contrast, Knowledge Graph and RAG approaches show more evenly distributed misclassification patterns but still fail to handle ambiguous cases effectively.

These results highlight the impact of dataset imbalances and the inherent ambiguity in certain categories. This suggests that combining the strengths of these approaches or addressing class imbalances could further improve performance in hierarchical product classification.

Misclassification analysis at the leaf level (Depth 3) reveals significant challenges in distinguishing closely related categories within hierarchical product classification. For instance, *ink cartridges* were frequently misclassified as *printer ink refills* (110 instances), *ink sticks* (43), *inkjet printers* (41), and *photo printers* (30). These misclassifications suggest that overlap titles and insufficient specificity in category definitions complicate accurate classification. Similarly, *notebooks* were often confused with *pcs/workstations* (59), *graphics cards* (35), and *tablets* (32), highlighting ambiguity in distinguishing between portable computing devices and related components. Other notable misclassifications include *mobile headsets* being assigned to *headphones* (38) and *internal hard drives* to *external hard drives* (30), which reflects overlap in product features that may cause confusion for both machine learning models and human annotators.

The approach-specific trends further underscore these challenges. Fine-tuning, while generally accurate, struggled with distinctions such as *mobile phone cases* misclassified as *peripheral device cases* (4.1% of errors), while prompting frequently misclassified *notebooks* as *pcs/workstations* (4.7%) or *tablets* (2.5%). Knowledge graphs improved structural consistency

but showed difficulties with overlapping hardware categories, such as *power supply units* misclassified as *power supplies* (3.2%). RAG demonstrated strong contextual understanding but often confused *mobile headsets* with *headphones* (3.3%). Entity matching faced challenges with rigid mappings, leading to frequent misclassifications of categories like *ink cartridges* that are part of overrepresented categories.

These findings suggest that certain adjustments to the taxonomy could improve classification accuracy. Categories with high misclassification rates, such as *ink cartridges* and *printer ink refills*, could be merged or redefined with clearer boundaries to reduce overlap. Similarly, categories like *notebooks* and *pcs/workstations* might benefit from additional intermediate subcategories that emphasize distinguishing features. For example, adding a subcategory that separates consumer-grade notebooks from professional workstations could aid in resolving ambiguity. Additionally, leveraging hierarchical relationships more effectively, such as requiring more granular context (e.g., intended use or technical specifications), could improve classification at the leaf level. These adjustments would not only improve LLM performance but also align the taxonomy with how humans intuitively group products, thus addressing inherent classification ambiguities.

6.4 Imbalanced Class Performance

6.4.1 Underrepresented Classes

Model	wP	wR	wF	hF	Level 1	Level 2	Level 3
Fine-tuning	0.92	0.47	0.46	0.68	0.85	0.62	0.47
Prompting	0.86	0.32	0.31	0.50	0.69	0.43	0.32
KG	0.91	0.38	0.39	0.56	0.69	0.51	0.38
RAG	0.93	0.39	0.40	0.59	0.75	0.50	0.39
Entity Matching	0.92	0.70	0.70	0.81	0.88	0.80	0.70

Table 9: Performance of methods in underrepresented categories.

Table 9 presents the evaluation metrics for a subset of results focused on underrepresented categories, where Entity Matching consistently outperformed the other methods. With an hF score of 81% and 70% for wR, wF, and level 3 accuracy, Entity Matching significantly outperformed finetuning, which achieved only 49% accuracy for leaf node categories. Other methods performed even worse, with scores below 40%. This difference may be attributed to the fact that

entity matching was the only method where the model was not explicitly prompted to categorize products. Instead, it was tasked with finding an entity match within a subset of similar titles, meaning that if the entity-match belongs to the underrepresented category, the prediction will inherently be correct. In contrast, methods relying on prompts might exhibit hesitancy in assigning products to underrepresented categories due to their niche and specific nature, which could make them appear less likely to the model.

6.4.2 Overrepresented Classes

Model	wP	wR	wF	hF	Level 1	Level 2	Level 3
Fine-tuning	1.00	0.93	0.96	0.97	0.99	0.96	0.93
Prompting	0.99	0.84	0.90	0.92	0.92	0.91	0.84
KG	0.99	0.96	0.97	0.98	0.99	0.98	0.96
RAG	1.00	0.98	0.99	0.99	0.99	0.99	0.98
Entity Matching	0.98	0.74	0.84	0.87	0.97	0.88	0.74

Table 10: Performance of methods in overrepresented categories.

In contrast to the performance observed in underrepresented categories, Table 10 highlights a notable difference in overrepresented categories, where entity matching performs significantly worse than the other models. While the fine-tuning, RAG, and KG-based methods consistently achieve metrics in the high 90s, and prompting scores in the high 80s to low 90s, entity matching lags behind with a wF score of 84 and an hF score of 87, demonstrating a clear underperformance. In overrepresented categories, where a larger variety of products with similar titles exist, the model may struggle to differentiate between subtle variations or incorrectly match to a common but incorrect entity. In contrast, methods like fine-tuning, RAG, and KG-based approaches leverage broader contextual understanding and hierarchical information, enabling them to generalize more effectively and accurately classify products within these categories.

6.5 Generalization

Fine Tuning performs strongly on WDC-222, achieving an remarkable hF of 0.87. On Avg Amazon, however, its hF dips to 0.81, showing a 7% decline. While still a reliable performer overall, this drop reveals its slight struggle with the more varied Amazon datasets. Prompting, on the other hand, delivers moderate results on WDC-222, with an hF of 0.76. Interestingly, it

Method	WDC-222		Avg. Amzn (Δ)		Amzn DE		Amzn EN		Amzn PL		Amzn TR	
	wF	hF	wF	hF	wF	hF	wF	hF	wF	hF	wF	hF
Fine-tuning	0.80	0.87	0.76	0.81 (-7%)	0.81	0.84	0.81	0.86	0.72	0.75	0.76	0.77
Prompting	0.69	0.76	0.73	0.81 (+7%)	0.78	0.87	0.70	0.80	0.72	0.81	0.63	0.77
KG	0.76	0.83	0.79	0.82 (-%1)	0.89	0.89	0.73	0.81	0.78	0.84	0.74	0.74
RAG	0.86	0.92	0.75	0.78 (-15%)	0.81	0.84	0.81	0.84	0.71	0.75	0.67	0.68
Entity Matching	0.77	0.84	0.56	0.68 (-19%)	0.61	0.73	0.70	0.77	0.49	0.63	0.44	0.58

Table 11: Performance metrics across WDC-222 and Amazon datasets. The best performance per dataset is bold, while the overall best hF model performance is greyed out.

improves on Avg Amazon, where the *hF* rises to 0.81, marking a 7% increase. This highlights Prompt’s strong generalization capabilities, as it performs better on the Amazon datasets than on WDC-222. Knowledge Graph stands out for its consistency. It achieves an *hF* of 0.83 on WDC-222 and maintains nearly the same level on Avg Amazon, where the *hF* is 0.82. With only a minor 1% decrease, this method proves its robustness and ability to adapt well across different datasets. RAG emerges as the best-performing method on WDC-222, achieving the highest *hF* of 0.92. However, this success does not carry over to the Avg Amazon datasets, where the *hF* drops sharply to 0.78, a significant 15% decline. While RAG excels on WDC-222, it struggles to generalize across the more diverse Amazon datasets. Entity Matching performs well on WDC-222, recording an *hF* of 0.84. In contrast, it experiences the most substantial drop on Avg Amazon, falling to an *hF* of 0.68, steep 19% decrease. This suggests that Entity Matching has difficulty adapting to the diversity present in the Amazon datasets.

In summary, KG and Prompting demonstrate the best generalization to the Amazon datasets, with minimal drops or even improvements in performance compared to WDC-222. In contrast, RAG and Entity Matching face the largest declines, highlighting challenges in handling the variability of Amazon data. Fine Tuning remains strong but sees a moderate drop. Additionally, DE and EN consistently outperform PL and TR across all methods, likely due to the LLMs’ reliance on English training data, such as Mistral and the fact that the quality of the German titles seem to be better. Within those two categories KG and Prompting achieve their best results for DE, Fine Tuning, EM, and RAG perform highest for EN due to the fact that it employs an english focused re-ranker.

During the analysis, some model predictions seemed like "misclassifications." On closer inspection, they often corrected our initial labels. For instance, a product labeled as "Ink

Cartridges” was accurately refined to “Printer Ink Refill,” highlighting both the models’ ability to capture fine distinctions and the limitations of manual labeling, which led us to re-label parts of the data.

7 Discussion

7.1 Analysis of Key Findings

7.1.1 Strengths and Weaknesses of Approaches

Fine-tuning demonstrates high precision by leveraging sampling techniques and K-means clustering to balance overrepresented classes. This approach ensures a diverse set of centroids, enabling the model to capitalize on these features during training. As a result, it balances hR and hP, providing reliability and depth in its predictions. However, at level 3, the wR is moderate (0.74), indicating some overfitting on more common classes. Despite being trained on specific data, the model exhibits moderately strong generalization abilities, performing reasonably well on the Amazon dataset indicating successful fine-tuning on the task of hierarchical classification in the computers and electronics segment.

The multi-stage top-k implementation, explored as part of different prompting techniques, effectively captures the hierarchical structure without requiring pretraining or additional context. This is evidenced by the difference between hf and level 3 wF scores, highlighting the presence of numerous near misses as the method incrementally refines predictions at each hierarchical level. This technique significantly outperforms both flat and hierarchical baselines, including the highly parameterized LLaMA 70B, demonstrating its ability to enhance the performance of lower-parameterized models. However, its reliance on the pretrained knowledge of the LLM limits its overall performance, as reflected in the lowest hF scores among the approaches. Despite this, the technique exhibits strong generalization on unseen Amazon data by avoiding overfitting, making it highly relevant and applicable to diverse, unseen datasets across multiple languages.

The Knowledge Graph (KG) approach with KG-reasoning and voting mechanisms, provides

contextual knowledge in a hierarchical manner, enabling balanced predictions with minimal drop-offs across hierarchical levels. However, through error propagation and too general contextual information on node attributes, the approach demonstrates lower precision at level 3 (0.834). Additionally, its moderate recall is limited by its reliance on existing relationships reducing its ability to handle edge cases. The KG approach generalizes strongly across the Amazon dataset in various languages by focusing on important entities rather than the entire data, converting raw document information into a structured graph. This reduces overfitting and enhances its applicability to unseen data.

RAG leverages task-specific retrieval to ground the LLM, enhancing its predictions and achieving very strong performance across all hierarchical levels. The combination of retrieval and generation ensures a strong balance between recall and precision at level 3, though its precision is slightly lower compared to finetuning due to occasional irrelevant context being introduced. RAG's generalization, however, is only moderate; the reliance on retrieving specific titles for predictions leads to overfitting, limiting its effectiveness on unseen data across various languages.

Entity matching achieves high precision at level 3 by minimizing false positives, particularly excelling in handling underrepresented classes. The prompt to match entities rather than categorize a title performs better in niche underrepresented categories. Its average recall is dependent on the quality of retrieval, which impacts performance on unseen data. Similar to RAG, entity matching exhibits moderate generalization but struggles with scenarios outside the training data, as the retrieval-based approach does not strongly enhance adaptability to new contexts and languages.

7.1.2 Complexity Analysis

Fine-tuning relies heavily on a subset of 30,000 observations sampled using clustering techniques from the entire dataset, while focusing on finetuning the model to the training data. While the training process is computationally intensive and time-consuming, the resulting model is highly efficient during inference, with fast response times, minimal resource usage, and category-only outputs, limiting token costs. This efficiency makes it suitable for domain-

specific applications. However, the rigidity of this approach is a drawback, as adapting to a new taxonomy or changes in the category structure requires retraining the model, which is both time and resource-intensive.

The multi-stage top-k prompting approach, in contrast, does not require any training data, making it a plug-and-play solution with strong generalization capabilities. By iteratively refining predictions at each level of the hierarchy, this technique leverages the taxonomy structure effectively to enhance classification accuracy without the need for pretraining or additional fine-tuning. It is the least complex of all approaches, requiring a few LLM inferences per observation when employing a local classifier per hierarchical level. Outputs are concise, consisting solely of the predicted category, which ensures token efficiency and eliminates the need for post-processing. Its low input token requirements further enhance its simplicity and reduce costs. Moreover, the setup for adapting to a new taxonomy or updating the hierarchy is instant, requiring no additional training or restructuring.

The Knowledge Graph (KG) approach has minimal training data requirements, needing only a few observations per category to build the graph and load attributes. It shares similarities with prompting in this regard but diverges significantly in complexity. The KG approach involves an inference process that, while offering greater interpretability through reasoning, is less efficient. This inefficiency is amplified by high input token counts, the use of Local Classifier Per Node (LCPN), and the implementation of a voting mechanism, resulting in at least nine LLM inferences per observation. While it provides significant value in scenarios requiring interpretability, its application to tasks such as product categorization may depend on specific use cases and requirements. However, adapting to a new taxonomy is straightforward, as it only requires creating a new KG from the category structure before inference.

Retrieval-Augmented Generation (RAG) is the most reliant on strong training data to ensure effective performance. Its inference outputs are limited to category predictions, but the process is complex and resource-intensive. High input token counts stem from the inclusion of 10 examples per call. While only one LLM inference is made per observation, this process is coupled with proprietary reranking mechanisms. The KG creation, vector database construction, and reranking process add significant computational overhead, resulting in high inference

costs. Adapting to a new taxonomy or making updates to it is similarly complex, requiring re-vectorization of the database, rebuilding of the KG, and reranking before inference.

Entity Matching shares similarities with RAG in its reliance on robust training data but offers a simpler and more streamlined approach. It outputs only the matched title, ensuring 100% label coherence, with average input token counts and a straightforward inference process. The method utilizes a basic vector database and requires just one LLM inference per observation, making it computationally efficient. Adapting to a new taxonomy is also relatively simple, involving only vectorization before proceeding to inference.

In the dynamic world of taxonomies, prompting emerges as the most adaptable approach, followed by the KG and Entity Matching, while finetuning and RAG offer higher accuracy at the cost of increased complexity, reduced generalizability, and greater rigidity.

7.2 Implications

7.2.1 Theoretical Implications

Our results confirm that LLMs can be rapidly adapted to hierarchical classification tasks with satisfactory accuracy. Methods like Prompting and RAG can be easily tested, while fine-tuning enables models to integrate hierarchical knowledge and reuse it during prediction. This supports our hypothesis that generative LLMs, combined with methods fostering effective generation (KG, RAG, Entity Matching), are well-suited for hierarchical classification tasks, often improving label correctness. Moreover, our solutions establish new benchmarks in the E-Commerce and Hierarchical Text Classification domains using WDC-222 with weak LLMs. Although methods suggested mostly do not surpass performance of strong LLM-based (RAG-based) techniques proposed by Gholamian et al. (2024). Implementing these approaches with powerful LLMs, such as GPT-4, has the potential to achieve similar or even higher performance levels.

The conclusions stem from examining five well-structured methodologies involving LLMs. Future improvements may arise from merging these approaches. We outline three promising directions:

Advanced Fine-Tuned Models Implementation: These models inherently maintain proper

hierarchical outputs from minimal inputs, even without in-context learning. Pairing them with prompting and in-context techniques (e.g., KG, RAG) may yield even better results. Since fine-tuned models often lose in-context learning capabilities, in-context-focused, specialized fine-tuning methods like the two-step ProMoT framework (Wang et al., 2024) for maintaining generalization abilities of fine-tuned models. In case of combining fine-tuned model with prompting, Dual Model Approach by Cheng et al. (2024) may lead to improved performance thanks to using fine-tuned model's output as hint for a prompt-based LLM technique.

Class-dependent method selection: Some techniques tend to perform better on certain types of classes. For instance, Entity Matching methods excel at handling underrepresented classes. A custom mechanism could retrieve few-shot examples (based on retrieval methods discussed in Part D) from vectorstores, determine if underrepresented classes dominate these examples, and then apply Entity Matching when needed. Otherwise, other highly efficient methods would apply.

Enhanced Knowledge Graph Retrieval: The Knowledge Graph (Part C) approach can be improved by combining it with retrieval techniques suggested in (Part D). Based on categories returned by the retrieval step, knowledge graph filtering mechanism can be developed (or inherited from the GraphRAG framework described in Paper D), in order to more efficiently provide partial Sub-Graphs as context for the in-context learning approaches, hence, potentially improving the classification performance while reducing prompt size.

7.2.2 Practical Implications

We offer businesses a cost-effective solution for categorizing products across diverse datasets while achieving competitive, state-of-the-art results. Although closed-source LLMs like GPT-4 may deliver better performance, our research demonstrates a viable and affordable alternative through open-source models such as Llama, Mistral, and Google Flan T5. To illustrate the cost difference, let's consider categorizing 100,000 products into a specific taxonomy. Using GPT-4 (8k), each call would process approximately 100 tokens (50 input, 50 output), with an average cost of \$0.0045 per call. This amounts to a total cost of \$450 for 100,000 API calls. In comparison, our approach leverages open-source models, where the only cost involved is

the computational infrastructure, such as a Colab Pro subscription or equivalent IDE services, roughly priced between \$10 and \$50/month. This significant cost disparity highlights the practicality of open-source solutions. Businesses with budget constraints can achieve robust classification results at a fraction of the cost, making these alternatives particularly suitable for small to medium-sized enterprises. While closed LLMs may excel in scenarios requiring maximum accuracy, open-source models offer scalable and cost-efficient performance.

Further, we can provide tailored recommendations for businesses on when to use specific LLM methods, depending on their unique requirements and constraints.

If time is a constraint or attaining training data is costly, we recommend using multi-stage top-k prompting. This approach offers a quick and flexible solution that does not require significant computational resources or training time. This makes it ideal for businesses that need rapid deployment or are working with frequently changing taxonomies.

If the company works with a large dataset of a specific taxonomy, we propose fine-tuning a model tailored to that taxonomy. Fine-tuning achieves high accuracy by adapting the model to domain-specific data, making it the best choice for businesses operating in specialized industries with well-defined category structures, such as fashion or automotive.

If the company deals with many diverse datasets, we suggest using the Knowledge Graph (KG) framework. KG demonstrates excellent generalization capabilities, making it suitable for businesses managing product catalogs with varying structures or those operating across different markets and languages.

If access to a domain-specific unstructured knowledge source is available, we recommend exploring RAG (Retrieval-Augmented Generation). RAG is particularly effective for scenarios involving incomplete or noisy datasets (like WDC-222), as it uses retrieval techniques to enhance accuracy. This is valuable for businesses handling unstructured product descriptions or legacy data systems.

By aligning model selection with specific business needs, companies can maximize the efficiency and accuracy of their product categorization workflows while minimizing costs and resource requirements. Our analysis demonstrates that open-source models, when implemented effectively, can be a cost-effective alternative for various e-commerce and data management

challenges.

7.3 Limitations

This study has several limitations that affect its generalizability and scalability. The focus on the "Computers & Electronics" domain and the use of only one taxonomy, narrows the applicability of the results to other product categories. The work was only generalized on products in the same domain and it was conducted on only 30 manually selected instances. Furthermore, since WDC was sourced from Icecat, the overlap between the WDC and Icecat datasets creates challenges for evaluating methods, as shared information may artificially boost the performance of retrieval-based approaches like RAG, reducing the reliability of comparisons. Dataset imbalances and overlapping subcategories further complicate classification, leading to reduced accuracy across all methods.

Weak learners were used and tested in different frameworks, including Ollama and Huggingface, which may have introduced inconsistencies in performance due to differences in model capabilities. Computational constraints also posed significant challenges. Fine-tuning requires retraining when taxonomies change, making it resource-intensive and impractical for dynamic or large-scale systems. Knowledge Integration methods like RAG and Knowledge Graph are computationally expensive, with high token costs and long inference times, limiting their scalability. While prompting-based techniques are more flexible, they struggle to handle the complexity of deep hierarchies without additional training or fine-tuning.

The hierarchical structure of the datasets, combined with shared information between training and testing sets, raises questions about the robustness of the results. Additionally, the reliance on product titles alone, while practical, may overlook valuable information in attributes like their descriptions. This highlights the need for further research on integrating richer data sources to improve classification accuracy and adaptability.

Improving accuracy is critically important to ensure real-world applicability, as low accuracy can result in significant misclassification, particularly in larger datasets. However, even with improved accuracy, scalability remains a challenge due to issues with label coherence. Failing to assign labels to certain data points necessitates manual labeling, which is not a scalable

solution.

7.4 Future Work

Future work should address the limitations identified in this study while exploring new directions for advancing hierarchical classification with large language models. Integration of the different explored methodologies that combine the strengths of existing approaches, such as fine-tuning, prompting, and knowledge graph-based reasoning can further enhance product categorization.

The current study's focus on computers and electronics domain and uses product titles for categorization. Therefore, future research should test the proposed methods on broader and more diverse domains such as fashion, healthcare, or industrial equipment. There is also potential to explore different taxonomies within the same domain. Furthermore, including richer data like descriptions, specifications, and customer reviews, as this can improve accuracy, especially for underrepresented and rare cases. Incorporating multimodal data, such as images or videos, can significantly enhance classification, not only in areas where text alone falls short but also when used in combination with textual data to improve performance.

Another promising direction is developing class dependent method selection mechanisms. Techniques like entity matching work well for underrepresented classes, while others are better suited for common categories. An adaptive pipeline that chooses methods based on class characteristics could greatly improve performance. A relevant improvement can be done regards, improving label coherence by exploring dynamic vocabulary restriction techniques, which limit the label space contextually during prediction, enhancing precision. Developing more robust training objectives or loss functions for fine-tuning could address the challenge of ambiguous or inconsistent labeling, a key factor for improving scalability.

Cost-effective approaches, such as multi-stage top-k prompting, can be refined for complex hierarchies. Fine-tuning methods could be enhanced for incremental updates, enabling models to adapt to changing taxonomies without full retraining. Future research may focus on reducing token costs and computational demands through strategies like pruning retrieval examples. The interaction between generative and discriminative approaches in hierarchical classification, rep-

resent an other area of research, generative LLMs could complement traditional models through techniques like knowledge distillation or dual-model frameworks. Overall, these advancements aim to create scalable, adaptable, and cost-efficient systems for diverse taxonomies, languages, and data sources, addressing real-world challenges effectively.

8 Conclusion

This study investigated how various techniques can enhance hierarchical product classification using LLMs. Our findings demonstrate that LLMs have significant potential to improve HPC, leveraging their advanced natural language processing capabilities to address challenges such as imbalanced taxonomies, and diverse datasets. While the performance is model-dependent and comes with certain limitations, such as moderate recall in entity matching or computational overhead in RAG, the approaches explored already offer practical implications for real-world applications.

Fine-tuning with hierarchical information proved effective for domain-specific classifications, achieving a strong balance between recall and precision, albeit with some limitations at deeper hierarchical levels. Multi-stage prompting emerged as a scalable and cost-effective solution, excelling in adaptability and generalization without requiring pretraining or large datasets. However, its performance remains inferior to advanced methods like fine-tuning or RAG when high precision is essential. Knowledge Graphs added value through structured contextual reasoning, balancing predictions and generalizing well to unseen data, although precision at deeper levels was sometimes affected by error propagation. RAG excelled in performing HPC on WDC but encountered limitations in generalization due to reliance on specific examples. Finally, entity matching excelled in precision for underrepresented classes, making it an ideal choice for stable and monolingual taxonomies, though its generalization was limited by retrieval quality.

The results confirm that the different explored techniques improve HPC using LLMs. The developed frameworks provide viable alternatives to traditional supervised methods, achieving comparable results with reduced training data requirements. RAG and Fine Tuning managed

Group Part: Enhancing Product Classification with LLMs

to beat the supervised benchmark of 0.85 hF set by Nils Richter and Christian Bizer (n.d.). By strategically employing the developed methods, businesses can implement adaptable and cost-efficient classification systems tailored to their needs. Future advancements, such as integrating these methodologies and further developing them hold promise for continued improvements in HPC with LLMs. These insights pave the way for scalable and robust solutions to meet the demands of an increasingly dynamic and competitive marketplace.

Paper C: Dynamic Integration of Taxonomy-Specific Knowledge Graphs with Large Language Models for Hierarchical Product Categorization in E-Commerce

C - 1 Introduction

The rapid scaling of large language models (LLMs) to billions of parameters has significantly improved performances across a diverse range of tasks (Minaee et al., 2024). Gholamian et al. (2024) demonstrate that large language models (LLMs) can outperform traditional supervised approaches in product classification across a range of configurations and scenarios.

Despite their achievements, LLMs have come under the scrutiny for occasionally lacking factual knowledge and generating hallucinations. Research indicates that LLMs often struggle to retrieve memorized facts and knowledge from their training data, resulting in the production of inaccurate or misleading statements (Ji et al., 2023).

The possession of outdated knowledge, combined with the inflexibility of pre-trained parameters and the high difficulty and cost of fine-tuning, can lead to a less adaptable and constrained LLM (Dong et al., 2024).

Similarly, LLMs have faced criticism for their lack of interpretability and transparency. Due to their black-box nature, with knowledge embedded in their parameters, the inference process of LLMs is challenging to validate and interpret (Zhao et al., 2024). Although some LLMs can explain their predictions through chain-of-thought reasoning, these explanations are often affected by the hallucination issue, limiting their reliability and applicability in high-stakes scenarios.

To address the above LLM limitations, researchers have recently turned to Knowledge Graphs (KGs) as a complementary approach. KGs store facts as structured triplets (entity – relation – entity), allowing them to encode knowledge with symbolic reasoning, generating interpretable results. This makes them particularly valuable for classification tasks, as they can store domain-specific knowledge in a structured format, enabling more accurate and explain-

able outcomes (Xu et al., 2024).

Moreover, the synergized integration of LLMs and KGs has emerged as a promising area of research. KGs are used to enhance LLMs by providing structured factual external knowledge, but also enhancing interpretability. LLMs contribute to the construction and expansion of KGs through their natural language processing capabilities (Chen et al., 2024a). This synergy aims to address the rigidity and hallucination issues of LLMs while leveraging the interpretability and precision of KGs.

Despite growing interest in the unification of LLMs and KGs, there remains limited research of their integration within hierarchical classification tasks. Representing taxonomies as KGs offers a structured framework, enabling the encoding of domain-specific hierarchical knowledge that can complement predictive capabilities of LLMs.

This paper addresses a critical gap in hierarchical product classification (HPC) by integrating taxonomy-specific Knowledge Graphs (KGs) into Large Language Model (LLM)-based classification. The key contributions of this work are as follows:

- **Dynamic KG-LLM Integration Framework:** Propose a novel, adaptive framework that dynamically integrates KGs into LLM prompts based on task-specific requirements.
- **Error Mitigation and Recovery Mechanism:** Introduce a robust hierarchical safety-net mechanism to detect, mitigate, and recover from error propagation through the taxonomy.
- **Empirical Validation of the Framework:** Validate the proposed framework through extensive experimentation on benchmark datasets.

C - 2 Background and Related Work

This section focuses on the background and existing work on KGs, focusing specifically on their integration with LLMs.

C - 2.1 Knowledge Graphs

Knowledge Graphs (KGs) are a widely used method for representing structured knowledge. They organize information as a collection of triplets. In a KG, nodes represent entities, which can refer to real-world objects or abstract concepts. Edges define relationships between these entities, capturing their semantic or logical connections. Each fact in a KG is represented as a triple in the format $(subject\ entity, relation, object\ entity)$, commonly denoted as (e_s, r, e_o) (Jiang et al., 2023b).

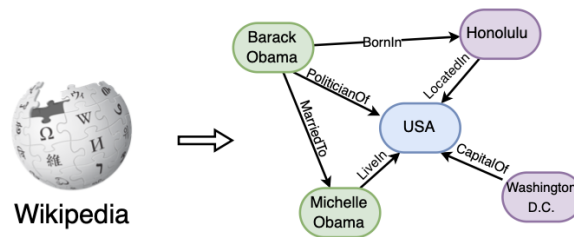


Figure 13: Example of a Knowledge Graph

Figure 13 illustrates a KG, highlighting how entities and their relationships are represented as nodes and edges. In a Knowledge Graph, each node can possess attributes that provide additional information about the entity it represents. For instance, in Figure 13, the node for "Barack Obama" may include attributes such as birth date, occupation, and nationality. Similarly, the "USA" node could have attributes like population, GDP, and capital city. These attributes enrich the graph's data, enabling more detailed queries and analyses (Bayram et al., 2020). There are several types of KGs based on the type of information they represent.

Encyclopedic KGs: These are the most common type of KGs, representing general real-world knowledge. They are typically constructed using a variety of extensive data sources, such as expert knowledge, databases, and encyclopedias. The largest knowledge graph, known as *Knowledge Ocean (KO)*, currently contains over 4 million entities and 1.7 billion relations (Pan et al., 2024).

Commonsense KGs: These KGs capture knowledge about everyday concepts, such as objects, events, and their relationships. Unlike encyclopedic KGs, commonsense KGs focus on modeling implicit knowledge, often extracted from textual sources. For example, a commonsense KG might include a fact like $(Bed, UsedFor, Sleep)$. A well-known example is *Concept-*

Net, which contains a wide range of commonsense concepts and relationships to help machines understand the meaning of human language (Pan et al., 2024).

Domain-Specific KGs: These KGs are tailored to represent knowledge within a specific domain, such as medicine, the sciences, or finance. While smaller in size compared to encyclopedic KGs, domain-specific KGs are often more precise and reliable due to their specialized focus. This study will focus on the integration Domain-Specific KGs (Pan et al., 2024). The construction of KGs relies on knowledge extraction techniques to generate structured triples. Two key techniques are Named Entity Recognition (NER) and Relation Extraction (RE). NER involves identifying and extracting entities from text. Common types of entities include categories such as *CARDINAL*, *ORDINAL*, and *DATE*, as well as *PERSON*, *ORG* (organization), and *LOC* (location). RE focuses on identifying and classifying the relationships between entities within textual content, enabling the formation of meaningful connections between the extracted entities. (Jiang et al., 2023b)

C - 2.2 Knowledge Graphs: Application in LLMs

The integration of KGs into LLMs and the construction of KGs using LLMs has become a focal point of recent research, aiming to enhance the models' comprehension and generation of knowledge-intensive content.

C - 2.2.1 KG-enhanced LLMs

Recent studies have investigated enhancing LLMs with KGs by integrating them into the inference process. Early approaches focused on incorporating KG triplets into LLM inputs using attention mechanisms (Liu et al., 2020) or by attaching graph encoders to LLM encoders (Wang et al., 2018). With the growing success of pre-trained LLMs, the focus has shifted towards prompting fixed pre-trained LLMs with graphical inputs (Choudhary and Reddy, 2024).

By leveraging knowledge retrieved from KGs, these methods can substantially enhance the performance of LLMs, particularly in accessing domain-specific knowledge (Lewis et al., 2021). Additionally, KGs have been employed to improve the interpretability of LLMs, enabling a clearer understanding of their reasoning processes and the facts they rely on (Petroni

et al., 2019).

C - 2.2.2 LLMs in KG Construction

Recent advancements have leveraged LLMs' ability in automating the construction of KGs by exploiting their advanced capabilities in natural language understanding. Researchers are now exploring unified approaches that integrate NER and RE to create structured triplets, combining entities and their relationships. These triplets serve as the building blocks for constructing meaningful KGs, enabling the transformation of unstructured raw text into structured knowledge representations (Kumar et al., 2020). This integrated approach demonstrates the potential of LLMs to streamline KG construction and enhance the utility of extracted knowledge.

Despite these advancements, there remains limited research on representing taxonomies as KGs, specifically for hierarchical classification tasks. Taxonomies, with their inherently structured and hierarchical nature, are well-suited to be modeled as KGs, yet their potential in this context has been underexplored.

Building on this gap, the next section outlines the methodology for constructing taxonomy-specific KGs and dynamically integrating them with LLMs to address challenges in hierarchical classification.

C - 3 Methodology

This section provides a systematic approach to integrating KGs into the product categorization pipeline, the approach consists of three main elements. **Extracting entities and relationships** from the dataset to construct the KG. **Integrating KGs into LLM inputs** by representing constructed KGs in formats that are compatible with LLMs. Finally, leverage the LLM's ability to **reason over the integrated KG** by processing subgraph information along with the given product data.

C - 3.1 Entity Extraction and KG Creation

In this paper, three broad types of domain-specific KGs are evaluated. **Taxonomy-Only KGs** represent the taxonomy as a graph with categories as nodes and "contains" as the only relationship, without any additional attributes or nodes. **Taxonomy-Enriched KGs** build on top of the taxonomy-only KG, by incorporating corpus-based information as attributes stored within the category nodes, while maintaining "contains" as the sole relationship. Finally, **NER-Extended KGs** expand the taxonomy by extracting additional entities from the corpus and linking them to the relevant taxonomy nodes.

To create the **NER-Extended KGs**, two distinct NER techniques are applied. The first approach utilizes spaCy's 'gliner' model, a pre-trained entity recognition tool that specializes in extracting entities from text. The second technique leverages the LLM itself for NER, extending the KG by identifying relevant entities within the corpus. These entities are then linked to their corresponding nodes in the taxonomy based on identified relationships. The prompt and model setup used in the NER process are detailed in Appendix - 2.

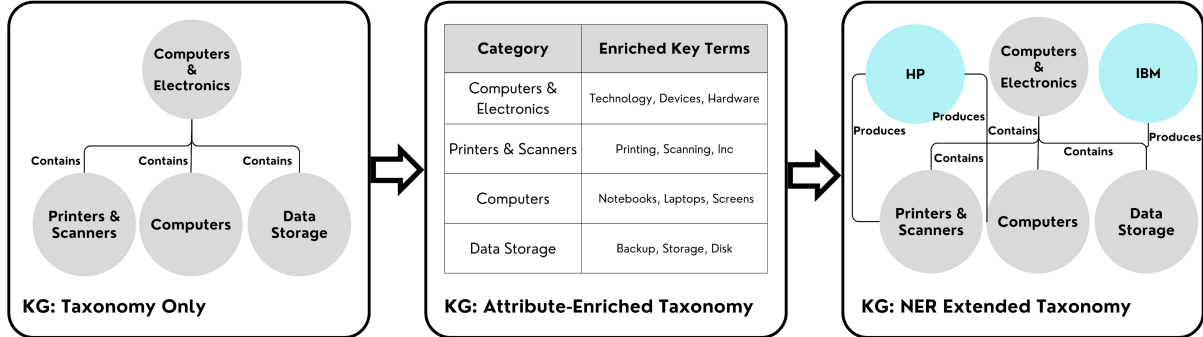


Figure 14: Three Types of Knowledge Graph Development and Enrichment Tested in this Paper

C - 3.2 Integrating KGs into LLM Inputs

C - 3.2.1 Verbalization

There has been extensive research on integrating KGs into LLM prompts. Building on the approach presented in Baek et al. (2023), this paper incorporates subgraph information into LLM prompts in the form of factual triplets. Since LLMs process textual inputs, while factual triplets are represented symbolically in KGs, the triplets must be converted into a textual format before

injection. This transformation, referred to as verbalization, converts (s, r, o) into a linear textual string by concatenating the subject, relation, and object. For example, the triplet *Computers, contains, Printers & Scanners* is verbalized as: "*Computers contains Printers & Scanners,*" which is then used as input for the LLM. This linear verbalization approach has been observed to perform effectively in LLM prompting.

Integrating the transformed KG information into the LLM prompt comes with a further key challenge: the number of triplets in the entire KG is too large and most of them unrelated, misleading and possibly confusing the LLM. To address this, a method of retrieving only the relevant triples for each hierarchical level is proposed, reducing the amount of unnecessary information fed into the model.

C - 3.2.2 Task-Relevant Subgraph Retrieval: Local Classifier per Parent Node

Instead of processing all KG triplets simultaneously, this paper employs a LCPN approach, where each hierarchical level is handled independently (Alagoz, 2023). This significantly reduces the number of triplets required for each prediction. By restricting the subgraph to the child nodes of the previous level's LLM output, the model focuses only on the relevant triplets for the current level. For instance, in Level 1 predictions, the previous LLM output ("*Computers & Electronics*") serves as the parent node, and the model considers only the triplets associated with this node. The output prediction of the LLM is then used as the input parent node for the next level's prediction. If the model predicts *Printers & Scanners* the subgraph included in the subsequent classifier's prompt will contain all child nodes and relevant attributes of the *Printers & Scanners* node. This process continues until a leaf node is reached. Figure 15 illustrates this process.

This approach was chosen for its simplicity, as it avoids the complexities of vectorization and other advanced retrieval techniques, making it easier and more cost-effective to implement and scale (Zhu et al., 2024). Unlike RAG architectures, which focus on optimizing retrieval and generation, this method primarily tests how to effectively incorporate KGs into LLMs for hierarchical classification tasks. Additionally, it offers better generalization, as it does not rely on training data and avoids the need for semantic similarity matching and ensuring robustness

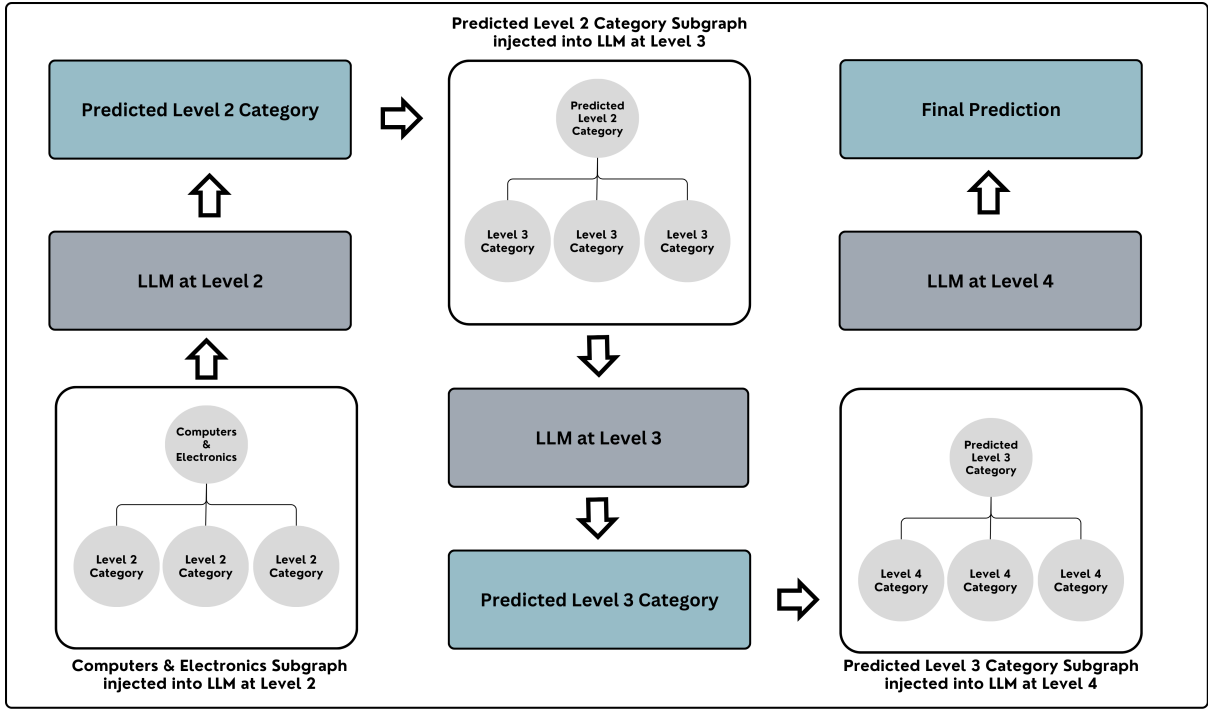


Figure 15: Subgraph Retrieval Process for Hierarchical Classification with LCPN even when new data is introduced.

C - 3.3 Prompting for Knowledge Graph Reasoning

To generate the final result and prediction, the LLM is provided with a prompt consisting of four components: a system instruction, the product title, the verbalized subgraph, and the graph reasoning instructions. A detailed prompt template along with the final LLM output are provided in Table 12. The final output includes a predicted category, enclosed in < and >, along with a step-by-step inference process and the graph reasoning pathway that led to the categorization decision.

The prompting and reasoning instruction techniques used in this approach are inspired by the successful methods employed in the MindMap framework (Wen et al., 2024), which integrates KGs into LLMs to facilitate structured reasoning. Similar to how MindMap enhances model reasoning by providing a clear logical structure and step-by-step reasoning, the current approach combines graph-based evidence with hierarchical category structures to guide LLMs in generating accurate predictions. MindMap’s use of reasoning chains has been shown to improve model transparency and accuracy, which is directly applied in this method to ensure that the model not only makes a prediction but also justifies it with a well-defined reasoning pro-

Task Description: Act as an AI product categorization assistant. Your task is to categorize the product based on the provided title and evidence.

Product Title: *Exemplar Title*

Retrieved Knowledge Graph: *Verbalized Subgraph*

Explicit Task: Based on the provided evidence, select the most appropriate child category for the product.

Reasoning Instructions:

Step 1: Analyze the Graph-based Evidence to identify any hierarchical relationships (parent-child category connections) that align with the product title.

Step 2: Combine insights from the hierarchical relationships and contextual information to infer the most accurate category.

Outputs:

Output 1: The predicted category enclosed in < and >. For example, if the category is 'printers & scanners', the output should be <printers & scanners>.

Output 2: Explain the inference process as a step-by-step reasoning chain. Include:

Which graph-based evidence was most relevant and why.

How the hierarchical relationships and contextual information from the graph-based evidence contributed to the final decision.

Exemplar Title: *Canon Fax L120 Laser 33 6kbit/s Fax Machine*

Output 1 (Category): <office electronics>

Output 2 (Inference): The product title "Canon Fax L120 Laser 33 6kbit/s Fax Machine" contains the keyword "fax machine," which clearly indicates that the product belongs to the "office electronics" category. The graph-based evidence further supports this by showing that the "computers & electronics" category has a child category called "office electronics," which also shares common words like "fax" and "machine." Therefore, the most relevant graph-based evidence is the "office electronics" category within "computers & electronics."

Table 12: Prompt Template: KG Reasoning for Categorization and Example Output process. Examining the output in Table 12 portrays how KG reasoning improves interpretability of LLM decisions.

C - 4 Experiments

This section evaluates various combinations of KG construction methods, node attribute generation approaches, and different stages of integrating the KG into the prompt to optimize the hierarchical product classification process. The methods are compared against basic LLM flat and hierarchical baselines without KGs.

C - 4.1 Datasets

The experiments are conducted on two publicly available datasets, Icecat (Open Icecat) and WDC-222 (Nils Richter and Christian Bizer, n.d.). For the creation of the KG and attributes, a balanced Icecat sample of approximately 15,000 observations is used. The sampling ensures 50 observations per category, except for categories where the number of observations is less than 50, in which case the sample size is denoted as $n < 50$. This sampling technique is employed for computational efficiency and generalization ability. A balanced approach is used as it ensures that each node has sufficient data to accurately generate attributes and node extensions (Cattaneo et al., 2022). The entire WDC-222 dataset is used to evaluate the various configurations.

C - 4.2 LLM Configuration & Baseline Configuration

The experiments were conducted using open-source LLMs, including Llama 3B-Instruct, Llama 8B-Instruct, and Mistral 7B-Instruct, with baseline evaluations additionally performed on Llama 70B-Instruct. Re-runs for consistency were not conducted due to the high inference costs; however, following the methodology outlined in Gholamian et al. (2024), the temperature parameter was set to 0 to reduce potential variations in LLM outputs across different runs. The Hugging Face API was utilized to access these models, and the maximum number of new tokens was capped at 200 to strike a balance between sufficient tokens for meaningful reasoning and the need to control computational costs. The decision to use open-source, moderately-sized models allows for the evaluation of accessible and cost-effective models in KG-enhanced classification tasks.

Baselines use the same LLM configurations as the other testing setups and serve as benchmarks in addition to those discussed in Section 3.2 of the group part. No knowledge augmentation from KGs is applied to the baselines. The **flat baseline** is provided with the same system instructions and prompt as shown in Table 12. The only difference being, instead of the retrieved and KG, the leaf nodes are presented as a list from which the LLM must select one. Additionally, the reasoning instructions are omitted in the flat baseline prediction. The **hierarchical baseline** is given the same instructions and prompt, but follows the LCPN approach

outlined in 15. However, unlike the KG-augmented model, the hierarchical baseline does not use KG injection or reasoning. Instead, it is provided with a list of available categories at each level, based on the prediction of the previous level.

C - 4.3 Experimental Configurations

In Figure 14, the various types of KGs tested in this paper are outlined. These include: **Taxonomy Only**, which captures parent-child relationships in the form of verbalized triplets; **Attribute-enriched Taxonomy**, where attributes are added to the taxonomy nodes. In this experiment, the top-10 most frequently used words for each category are retrieved and loaded into the corresponding nodes for prompt retrieval. Finally, **NER-extended KGs** are evaluated, incorporating two types of Named Entity Recognition (NER). First, spaCy’s ‘gliner’ model and second, the LLM itself is leveraged for named entity extraction to further enrich the KG. Initial results indicated that the additional information introduced at Level 1 caused confusion for the LLM. To test this hypothesis, a modified configuration was introduced in which the NER techniques are applied only at Level 2 for comparison. For all configurations, other than the flat baseline, a zero-shot LCPN approach is employed, using the same prompt template as shown in Figure 12.

C - 4.4 Results

In Table 13, the various configurations are evaluated based on wP, wR, wF and hF scores. Additionally, the table displays accuracy metrics per level. Hierarchical approaches outperform the flat baseline for all models except LLaMA 70B, aligning with the findings of Gholamian et al. (2024).

For LLaMA 3B, the hierarchical baseline outperforms the flat baseline and any KG configuration. Evaluation metrics move opposite to input tokens, with performance dropping as more tokens are provided. The lower-parameter LLM appears to struggle with increased inputs and reasoning tasks, leading to lower accuracies. In contrast, LLaMA 8B and Mistral 7B show different patterns.

For LLaMA 7B, the hierarchical baseline, the attribute-enriched taxonomy, and the LLM-

Model	Architecture	wP	wR	wF	hF	lv1	lv2	lv3
Llama 3B	Baseline - Flat	0.04	0.02	0.02	0.07	0.08	0.04	0.02
	Baseline - Hierarchical	0.60	0.56	0.54	0.67	0.77	0.62	0.56
	Taxonomy Only	0.44	0.38	0.39	0.48	0.55	0.42	0.38
	Taxonomy Top 10	0.35	0.12	0.16	0.25	0.42	0.16	0.13
	Taxonomy NER (Spacy - Gliner)	0.23	0.14	0.16	0.18	0.48	0.28	0.26
	Taxonomy NER (Level 2)	0.42	0.41	0.41	0.58	0.70	0.49	0.41
	Taxonomy NER (LLM - Level 3)	0.49	0.39	0.39	0.56	0.70	0.48	0.39
Llama 8B	Baseline - Flat	0.67	0.33	0.36	0.50	0.42	0.39	0.33
	Baseline - Hierarchical	0.74	0.66	0.67	0.73	0.80	0.68	0.66
	Taxonomy Only	0.70	0.62	0.62	0.71	0.79	0.67	0.62
	Taxonomy Top 10	0.71	0.64	0.65	0.73	0.79	0.70	0.67
	Taxonomy NER (Spacy - Gliner)	0.46	0.45	0.42	0.57	0.63	0.47	0.45
	Taxonomy NER (Level 2)	0.61	0.54	0.54	0.68	0.79	0.54	0.54
	Taxonomy NER (LLM - Level 3)	0.70	0.63	0.63	0.73	0.79	0.68	0.63
Mistral 7B	Baseline - Flat	0.18	0.12	0.13	0.28	0.28	0.18	0.12
	Baseline - Hierarchical	0.65	0.50	0.54	0.66	0.75	0.56	0.50
	Taxonomy Only	0.70	0.61	0.63	0.72	0.79	0.68	0.61
	Taxonomy Top 10	0.68	0.59	0.60	0.70	0.76	0.65	0.59
	Taxonomy NER (Spacy - Gliner)	0.31	0.17	0.19	0.30	0.29	0.20	0.17
	Taxonomy NER (Level 2)	0.64	0.56	0.57	0.69	0.80	0.62	0.56
	Taxonomy NER (LLM - Level 3)	0.68	0.61	0.62	0.73	0.80	0.68	0.61
Llama 80B	Baseline - Flat	0.81	0.56	0.62	0.72	0.62	0.59	0.56
	Baseline - Hierarchical	0.70	0.62	0.64	0.71	0.79	0.64	0.62
Framework	Framework - With Safety Net	0.83	0.74	0.76	0.83	0.87	0.81	0.74

Table 13: Results

NER extended taxonomy all achieve an hF score of 0.73. Interestingly, while the hierarchical baseline achieves the highest accuracy at Level 1, the KG-configurations exhibit smaller accuracy drop-offs between Levels 1 and 2, resulting in the same hF score, despite error propagation from Level 1. It seems that the KG boosts accuracy at deeper levels but may hinder LLM predictions at earlier ones, as further evidenced by the strong performance of Level 2 SpaCy-NER compared to the regular SpaCy-NER configuration.

For Mistral, LLM-NER outperforms all other configurations on an hF basis and achieves the highest wR. All KG configurations, except for the spaCy-NER at full hierarchy, outperform the hierarchical baseline. Analyzing the accuracy fall-offs, the KG-injected models outperform the hierarchical baseline at Level 2 predictions, highlighting the effectiveness of KG injection

at deeper levels. Even the taxonomy-only KG improves over the hierarchical baseline, showing that triplet representation and KG-reasoning is effective with Mistral-7B, as the contextual knowledge between the baseline and taxonomy-only KG is the same.

The error propagation issue in the LCPN approach is evident in the results, as high accuracies at Level 1 typically lead to higher final accuracies. Examining the accuracies at Level 1, the fall-offs from Level 1 to 2 and from Level 2 to 3 and the eventual hF score, we can draw the following conclusions: the injection of KGs and reasoning over the KG has mixed effects. In some cases, it confuses the LLM, as seen with LLaMA 3B and Level 1 predictions for LLaMA 7B. However, in other scenarios the KG augmented knowledge and performed reasoning boosts accuracies, as seen with the strong performing KG configurations in Mistral and the deeper levels for LLaMA 7B predictions.

It is hypothesized that in observations where the LLM can classify accurately without KG injection, adding additional knowledge confuses the LLM. However, when the task is more complex, the injection of a KG benefits the LLM. Based on these findings, a final framework is proposed that models LLM uncertainty to determine when to inject the KG. This framework combines the two best-performing configurations, namely the hierarchical-baseline and LLM NER—while adding a safety-net flat classifier with an advanced LLM to address the LCPN error propagation dilemma. The framework is illustrated in Figure 16.

C - 4.5 Framework

The proposed framework, presented in Figure 16, operates as follows: three weak LLMs are initially prompted, without KG-augmentation. Using the hierarchical-baseline, to predict the first level of the taxonomy for a specific observation. If all predictions match, the consensus will be used as the final prediction for this level, which then serves as input for the next one. In case of a mismatch, indicating uncertainty among the LLMs, they are re-prompted with KG injection. This step leverages the finding from the results that KG injection improves LLM performance when the models are uncertain. If a majority prediction emerges from the re-prompting, it is selected as the final prediction for the next level. In cases of three different predictions, the advanced LLM steps in as a judge to make the final decision. This iterative

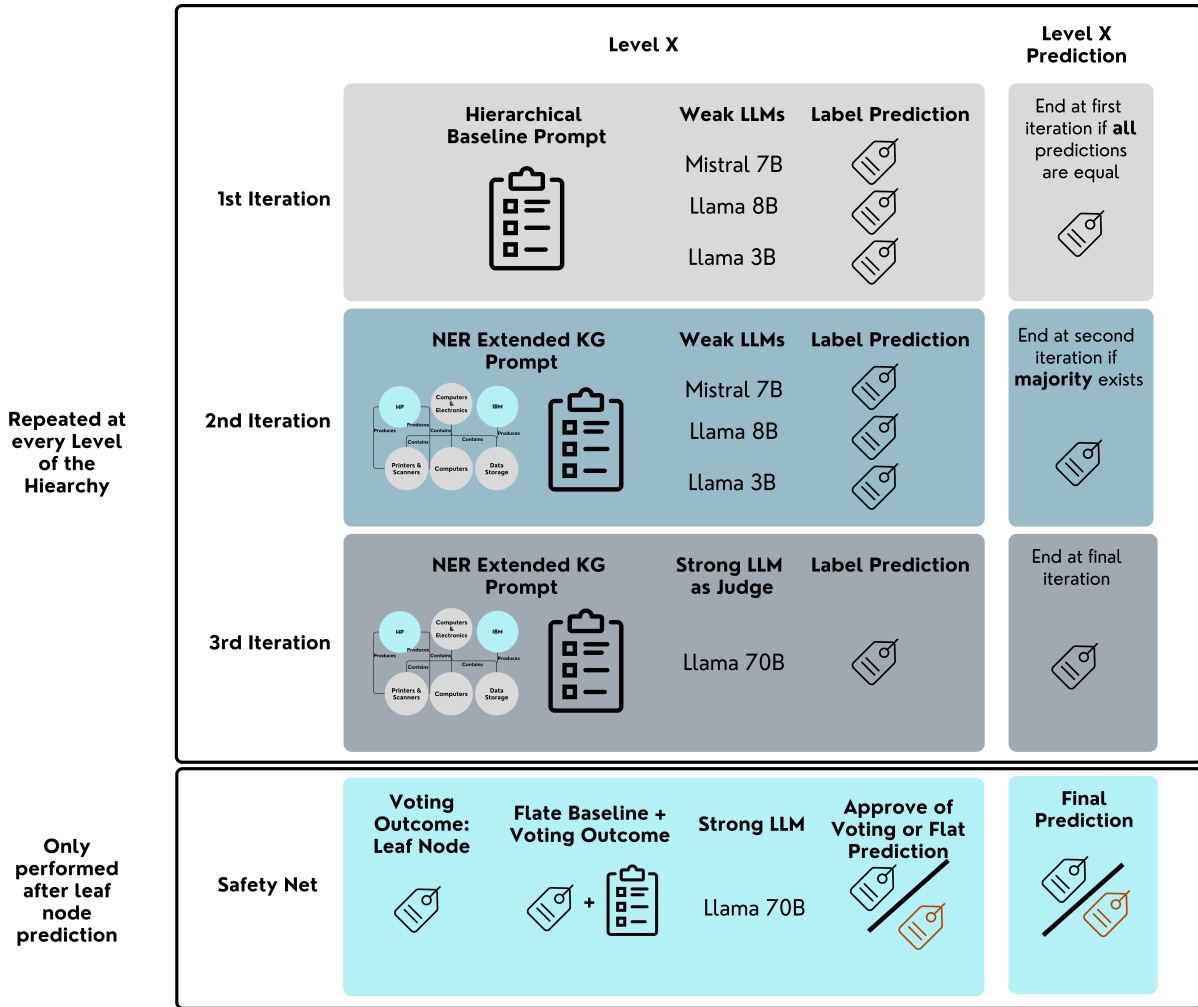


Figure 16: Framework for KG-Augmented Product Classification

process continues until the leaf node is reached. Upon reaching the leaf node, the final prediction is presented to an advanced LLM, acting as a safety-net to mitigate error propagation. Without this safety net, errors in earlier levels would not be able to be corrected. The safety net is an advanced LLM flat classifier for the leaf nodes, applied only if it disagrees with the voting mechanism’s output. This framework integrates both LCPN and flat classification, akin to the method proposed by Alagoz (2023).

As shown in Table 13, the framework outperforms all other configurations across every metric. With an hF score of 0.83, it exceeds the performance of the advanced LLM by 12 percentage points. To evaluate the actual impact of KG injection, the following analysis is conducted. Of the 8,952 predictions (2,984 observations x 3 levels), 5,470 were resolved by the weak LLMs in the first iteration. 2,712 were solved by the KG-injected LLMs, and only 768 required a third iteration with the strong LLM. Among the 2,712 predictions resolved by

KG-injected prompts, accuracy improved by a significant 21.2% compared to the predictions made in the first iteration by the non-KG augmented prompt. This clearly demonstrates that KG injection and reasoning significantly enhances the predictive power of LLMs in product categorization, particularly when the model is uncertain. Further, more detailed results per level and prediction source can be seen in Figure 23 in Appendix - 2

C - 5 Conclusion, Limitations and Future Work

This research evaluates various configurations of KG-injection into LLM prompts for HPC and introduces a novel adaptive framework that dynamically integrates KGs into LLM prompts based on task-specific requirements. Following an LCPN + flat scheme with voting mechanisms, LLM uncertainty and error propagation are successfully addressed. It outperforms all other configurations and performs on par with state-of-the-art supervised baselines while using significantly less data. This supports the hypothesis that KGs enhance the precision of LLMs in HPC tasks.

However, the analysis has limitations. Specifically, the performance of KGs in the LLaMA 3B model is unstable, and no definitive conclusions can be drawn. The framework's reliance on multiple LLMs per prediction and its LCPN structure incurs significant computational costs. While most predictions are made by low-cost weak LLMs, the process requires substantial time and computational power. Another limitation is the inability to incorporate full category paths into prompts due to token limits. The current LCPN approach may not fully leverage the power of KGs, and improving the subgraph retrieval process could enable the injection of full paths for specific KG subgraphs, which could be explored in future research.

Future studies could build on this work by testing KG integration with few-shot learning approaches and optimizing attribute generation. Additionally, research could explore enhancing the safety net mechanisms or finding more efficient methods for handling error propagation in LCPN structures. Furthermore, integrating external KGs, rather than relying solely on taxonomy-specific KGs, could eliminate the KG creation process and facilitate predictions with zero training data.

References

- J. Achiam, S. Adler, S. Agarwal, et al. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- C. Alagoz. Performance improvement in multi-class classification via automated hierarchy generation and exploitation through extended lcpn schemes, 2023. URL <https://arxiv.org/abs/2310.20641>.
- R. Anil, S. Borgeaud, J.-B. Alayrac, et al. Gemini: A family of highly capable multimodal models, 2024. URL <https://arxiv.org/abs/2312.11805>.
- AnthropicAI. The claude 3 model family: Opus, sonnet, haiku. URL <https://api.semanticscholar.org/CorpusID:268232499>.
- J. Baek, A. F. Aji, and A. Saffari. Knowledge-augmented language model prompting for zero-shot knowledge graph question answering, 2023. URL <https://arxiv.org/abs/2306.04136>.
- E. Bayram, A. Garcia-Duran, and R. West. Node attribute completion in knowledge graphs with multi-relational propagation, 2020. URL <https://arxiv.org/abs/2011.05301>.
- S. Biderman, H. Schoelkopf, Q. Anthony, H. Bradley, K. O’Brien, E. Hallahan, M. A. Khan, S. Purohit, U. S. Prashanth, E. Raff, A. Skowron, L. Sutawika, and O. van der Wal. Pythia: A suite for analyzing large language models across training and scaling, 2023.
- T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>.
- A. Cattaneo, D. Justus, H. Mellor, D. Orr, J. Maloberti, Z. Liu, T. Farnsworth, A. Fitzgibbon, B. Banaszewski, and C. Luschi. Bess: Balanced entity sampling and sharing for large-scale knowledge graph completion, 2022. URL <https://arxiv.org/abs/2211.12281>.
- H. Chen, Q. Ma, Z. Lin, and J. Yan. Hierarchy-aware label semantics matching network for hierarchical text classification. In C. Zong, F. Xia, W. Li, and R. Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4370–4379, Online, Aug. 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.337. URL <https://aclanthology.org/2021.acl-long.337>.
- H. Chen, X. Shen, Q. Lv, J. Wang, X. Ni, and J. Ye. Sac-kg: Exploiting large language models as skilled automatic constructors for domain knowledge graphs, 2024a. URL <https://arxiv.org/abs/2410.02811>.
- H. Chen, Y. Zhao, Z. Chen, M. Wang, L. Li, M. Zhang, and M. Zhang. Retrieval-style in-context learning for few-shot hierarchical text classification, 2024b. URL <https://arxiv.org/abs/2406.17534>.
- Z. Cheng, W. Zhang, C.-C. Chou, Y.-Y. Jau, A. Pathak, P. Gao, and U. Batur. E-commerce product categorization with LLM-based dual-expert classification paradigm. In S. Kumar, V. Balachandran, C. Y. Park, W. Shi, S. A. Hayati, Y. Tsvetkov, N. Smith, H. Hajishirzi, D. Kang, and D. Jurgens, editors, *Proceedings of the 1st Workshop on Customizable NLP: Progress and Challenges in Customizing NLP for a Domain, Application, Group, or Individual (CustomNLP4U)*, pages 294–304, Miami, Florida, USA, Nov. 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.customnlp4u-1.22. URL <https://aclanthology.org/2024.customnlp4u-1.22>.

Bibliography

- N. Choudhary and C. K. Reddy. Complex logical reasoning over knowledge graphs using large language models, 2024. URL <https://arxiv.org/abs/2305.01157>.
- J. Dai, T. Wang, and S. Wang. A deep forest method for classifying e-commerce products by using title information. In *2020 International Conference on Computing, Networking and Communications (ICNC)*, page 1–5. IEEE, Feb. 2020. doi: 10.1109/icnc47757.2020.9049751. URL <http://dx.doi.org/10.1109/ICNC47757.2020.9049751>.
- M. A. de Luis Balaguer, V. Benara, R. L. de Freitas Cunha, R. de M. Estevao Filho, T. Hendry, D. Holstein, J. Marsman, N. Mecklenburg, S. Malvar, L. Nunes, R. Padilha, M. Sharp, B. L. B. Silva, S. Sharma, V. Aski, and R. Chandra. Rag vs fine-tuning: Pipelines, tradeoffs, and a case study on agriculture. *ArXiv*, abs/2401.08406, 2024. URL <https://api.semanticscholar.org/CorpusID:267027552>.
- J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019. URL <https://arxiv.org/abs/1810.04805>.
- Y. Dong, H. Zhang, C. Li, S. Guo, V. C. M. Leung, and X. Hu. Fine-tuning and deploying large language models over edges: Issues and approaches, 2024. URL <https://arxiv.org/abs/2408.10691>.
- S. T. Dumais and H. Chen. Hierarchical classification of web content, 2000. URL <https://api.semanticscholar.org/CorpusID:1194024>.
- C. d’Amato, P. Ristoski, P. Petrovski, P. Mika, and H. Paulheim. A machine learning approach for product matching and categorization. *Semant. Web*, 9(5):707–728, Jan. 2018. ISSN 1570-0844. doi: 10.3233/SW-180300. URL <https://doi.org/10.3233/SW-180300>.
- L. Fan, L. Li, Z. Ma, S. Lee, H. Yu, and L. Hemphill. A bibliometric review of large language models research from 2017 to 2023, 2023.
- S. Gholamian, G. Romani, B. Rudnikowicz, and S. Skylaki. Llm-based robust product classification in commerce and compliance, 2024. URL <https://arxiv.org/abs/2408.05874>.
- A. Grattafiori, A. Dubey, A. Jauhri, et al. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- L. Groß, R. Walter, N. Zoppi, and A. R. Justus. Rebuild product classification. Masters thesis, Nova School of Business and Economics, Lisbon, Portugal, 2024.
- S. Guo and H. Zhao. Hierarchical classification with multi-path selection based on granular computing. *Artif. Intell. Rev.*, 54(3):2067–2089, Mar. 2021. ISSN 0269-2821. doi: 10.1007/s10462-020-09899-2. URL <https://doi.org/10.1007/s10462-020-09899-2>.
- V. Gupta, H. Karnick, A. Bansal, and P. Jhala. Product classification in e-commerce using distributional semantics, 2016. URL <https://arxiv.org/abs/1606.06083>.
- S. Gururangan, A. Marasović, S. Swayamdipta, K. Lo, I. Beltagy, D. Downey, and N. A. Smith. Don’t stop pretraining: Adapt language models to domains and tasks. In D. Jurafsky, J. Chai, N. Schluter, and J. Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.740. URL <https://aclanthology.org/2020.acl-main.740>.
- L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin, and T. Liu. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, Nov. 2024. ISSN 1558-2868. doi: 10.1145/3703155. URL <http://dx.doi.org/10.1145/3703155>.

Bibliography

- R. Huang, C. Wei, B. Wang, J. Yang, X. Xu, S. Wu, and S. Huang. Well performance prediction based on long short-term memory (lstm) neural network. *Journal of Petroleum Science and Engineering*, 208:109686, 10 2021. doi: 10.1016/j.petrol.2021.109686.
- W. Huang, E. Chen, Q. Liu, Y. Chen, Z. Huang, Y. Liu, Z. Zhao, D. Zhang, and S. Wang. Hierarchical multi-label text classification: An attention-based recurrent network approach, 2019. URL <https://api.semanticscholar.org/CorpusID:207758666>.
- V. Jain, M. Rungta, Y. Zhuang, Y. Yu, Z. Wang, M. Gao, J. Skolnick, and C. Zhang. Higen: Hierarchy-aware sequence generation for hierarchical text classification, 2024. URL <https://arxiv.org/abs/2402.01696>.
- Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, and P. Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, Mar. 2023. ISSN 1557-7341. doi: 10.1145/3571730. URL <http://dx.doi.org/10.1145/3571730>.
- A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. de las Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, L. R. Lavaud, M.-A. Lachaux, P. Stock, T. L. Scao, T. Lavril, T. Wang, T. Lacroix, and W. E. Sayed. Mistral 7b, 2023a.
- T. Jiang, D. Wang, L. Sun, Z. Chen, F. Zhuang, and Q. Yang. Exploiting global and local hierarchies for hierarchical text classification, 2022. URL <https://arxiv.org/abs/2205.02613>.
- X. Jiang, C. Xu, Y. Shen, X. Sun, L. Tang, S. Wang, Z. Chen, Y. Wang, and J. Guo. On the evolution of knowledge graphs: A survey and perspective, 2023b. URL <https://arxiv.org/abs/2310.04835>.
- S. Kiritchenko, S. Matwin, R. Nock, and A. F. Famili. *Learning and Evaluation in the Presence of Class Hierarchies: Application to Text Categorization*, page 395–406. Springer Berlin Heidelberg, 2006. ISBN 9783540248408. doi: 10.1007/11766247_34. URL http://dx.doi.org/10.1007/11766247_34.
- Z. Kozareva, Q. Li, K. Zhai, and W. Guo. Recognizing salient entities in shopping queries, 01 2016.
- A. Krishnan and A. Amarthaluri. Large scale product categorization using structured and unstructured attributes, 2019. URL <https://arxiv.org/abs/1903.04254>.
- A. Kumar, A. Pandey, R. Gadia, and M. Mishra. Building knowledge graph using pre-trained language model for learning entity-aware relationships. *2020 IEEE International Conference on Computing, Power and Communication Technologies (GUCON)*, pages 310–315, 2020. URL <https://api.semanticscholar.org/CorpusID:226266162>.
- P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W. tau Yih, T. Rocktäschel, S. Riedel, and D. Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks, 2021. URL <https://arxiv.org/abs/2005.11401>.
- M. Y. Li, S. Kok, and L. Tan. Don’t classify, translate: Multi-level e-commerce product categorization via machine translation. *ArXiv*, abs/1812.05774, 2018. URL <https://api.semanticscholar.org/CorpusID:55702283>.
- P. Liu, X. Qiu, and X. Huang. Recurrent neural network for text classification with multi-task learning, 2016. URL <https://arxiv.org/abs/1605.05101>.
- W. Liu, P. Zhou, Z. Zhao, Z. Wang, Q. Ju, H. Deng, and P. Wang. K-bert: Enabling language representation with knowledge graph. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34: 2901–2908, 04 2020. doi: 10.1609/aaai.v34i03.5681.

Bibliography

- A. Martinek, S. Łukasik, and A. H. Gandomi. Text-based product matching – semi-supervised clustering approach, 2024. URL <https://arxiv.org/abs/2402.10091>.
- J. MathavRaj, V. Kushala, H. Warriar, and Y. Gupta. Fine tuning llm for enterprise: Practical guidelines and recommendations. *ArXiv*, abs/2404.10779, 2024. URL <https://api.semanticscholar.org/CorpusID:269188062>.
- N. Mathivanan, N. Ghani, and R. Mohd Janor. Analysis of k-means clustering algorithm: A case study using large scale e-commerce products, 11 2019.
- G. Melis, C. Dyer, and P. Blunsom. On the state of the art of evaluation in neural language models, 2017. URL <https://arxiv.org/abs/1707.05589>.
- B. Min, H. Ross, E. Sulem, A. P. B. Veyseh, T. H. Nguyen, O. Sainz, E. Agirre, I. Heinz, and D. Roth. Recent advances in natural language processing via large pre-trained language models: A survey, 2021. URL <https://arxiv.org/abs/2111.01243>.
- S. Minaee, T. Mikolov, N. Nikzad, M. Chenaghlu, R. Socher, X. Amatriain, and J. Gao. Large language models: A survey, 2024. URL <https://arxiv.org/abs/2402.06196>.
- F. M. Miranda, N. Köhnecke, and B. Y. Renard. Hiclass: a python library for local hierarchical classification compatible with scikit-learn. *Journal of Machine Learning Research*, 24(29):1–17, 2023. URL <http://jmlr.org/papers/v24/21-1518.html>.
- M. Mosbach, T. Pimentel, S. Ravfogel, D. Klakow, and Y. Elazar. Few-shot fine-tuning vs. in-context learning: A fair comparison and evaluation. In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Findings of the Association for Computational Linguistics: ACL 2023*, pages 12284–12314, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.779. URL <https://aclanthology.org/2023.findings-acl.779>.
- Nils Richter and Christian Bizer. Wdc-222 gold standard for hierarchical product categorization, n.d. URL <https://data.dws.informatik.uni-mannheim.de/largescaleproductcorpus/categorization/>. Accessed: 2024-11-11.
- Open Icecat. Open icecat catalog. URL <https://icecat.de/>. Accessed: 2024-11-11.
- S. Pan, L. Luo, Y. Wang, C. Chen, J. Wang, and X. Wu. Unifying large language models and knowledge graphs: A roadmap. *IEEE Transactions on Knowledge and Data Engineering*, 36(7):3580–3599, July 2024. ISSN 2326-3865. doi: 10.1109/tkde.2024.3352100. URL <http://dx.doi.org/10.1109/TKDE.2024.3352100>.
- A. Patra, V. Vivek, B. R. Shambhavi, K. S. Sindhu, and S. A. Balaji. Product classification in e-commerce sites, 2021. URL <https://api.semanticscholar.org/CorpusID:235082768>.
- F. Petroni, T. Rocktäschel, P. Lewis, A. Bakhtin, Y. Wu, A. H. Miller, and S. Riedel. Language models as knowledge bases?, 2019. URL <https://arxiv.org/abs/1909.01066>.
- A. Radford and K. Narasimhan. Improving language understanding by generative pre-training, 2018. URL <https://api.semanticscholar.org/CorpusID:49313245>.
- K. K. Rangan and Y. Yin. A fine-tuning enhanced rag system with quantized influence measure as ai judge. *Scientific Reports*, 14, 2024. URL <https://api.semanticscholar.org/CorpusID:268033066>.

Bibliography

- J. Risch, S. Garda, and R. Krestel. Hierarchical document classification as a sequence generation task. In *Proceedings of the ACM/IEEE Joint Conference on Digital Libraries in 2020, JCDL '20*, page 147–155, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450375856. doi: 10.1145/3383583.3398538. URL <https://doi.org/10.1145/3383583.3398538>.
- P. Sahoo, A. K. Singh, S. Saha, V. Jain, S. S. Mondal, and A. Chadha. A systematic survey of prompt engineering in large language models: Techniques and applications. *ArXiv*, abs/2402.07927, 2024. URL <https://api.semanticscholar.org/CorpusID:267636769>.
- B. Schwartz. Ceos want trump to change course on tariffs. he isn't budging. *The Wall Street Journal*, 2024. URL <https://www.wsj.com/economy/trade/trump-tariff-plan-business-lobbying-8f02ccea>. Accessed: 2024-12-15.
- D. Shen, J.-D. Ruvini, M. Somaiya, and N. Sundaresan. Item categorization in the e-commerce domain, 10 2011.
- J. Shin, C. Tang, T. Mohati, M. Nayebi, S. Wang, and H. Hemmati. Prompt engineering or fine tuning: An empirical assessment of large language models in automated software engineering tasks, 2023. URL <https://arxiv.org/abs/2310.10508>.
- Statista. Gross merchandise value (gmv) of tiktok shops worldwide as of 2024, by product category, 2024. URL <https://www.statista.com/statistics/1461508/gross-merchandise-value-tiktok-shops-worldwide>. Accessed: 2024-12-15.
- L. Tan, M. Y. Li, and S. Kok. E-commerce product categorization via machine translation. *ACM Trans. Manage. Inf. Syst.*, 11(3), July 2020. ISSN 2158-656X. doi: 10.1145/3382189. URL <https://doi.org/10.1145/3382189>.
- F. Torba, C. Gravier, C. Laclau, A. Kammoun, and J. Subercaze. A Study on Hierarchical Text Classification as a Seq2seq Task. In *46th European Conference on Information Retrieval (ECIR 2024)*, GLASGOW, United Kingdom, Mar. 2024. URL <https://hal.science/hal-04423348>.
- H. Touvron, L. Martin, K. Stone, et al. Llama 2: Open foundation and fine-tuned chat models, 2023. URL <https://arxiv.org/abs/2307.09288>.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need, 2023. URL <https://arxiv.org/abs/1706.03762>.
- X. Wang, P. Kapanipathi, R. Musa, M. Yu, K. Talamadupula, I. Abdelaziz, M. Chang, A. Fokoue, B. Makni, N. Mattei, and M. Witbrock. Improving natural language inference using external knowledge in the science questions domain, 09 2018.
- Y. Wang, S. Si, D. Li, M. Lukasik, F. Yu, C.-J. Hsieh, I. S. Dhillon, and S. Kumar. Two-stage llm fine-tuning with less specialization and more generalization, 2024. URL <https://arxiv.org/abs/2211.00635>.
- Z. Wang, P. Wang, L. Huang, X. Sun, and H. Wang. Incorporating hierarchy into text encoder: a contrastive learning approach for hierarchical text classification. In S. Muresan, P. Nakov, and A. Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7109–7119, Dublin, Ireland, May 2022a. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.491. URL <https://aclanthology.org/2022.acl-long.491>.

Bibliography

- Z. Wang, P. Wang, T. Liu, B. Lin, Y. Cao, Z. Sui, and H. Wang. Hpt: Hierarchy-aware prompt tuning for hierarchical text classification, 2022b. URL <https://arxiv.org/abs/2204.13413>.
- G. Weiss and F. Provost. The effect of class distribution on classifier learning: An empirical study. *Tech Rep*, 09 2001.
- Y. Wen, Z. Wang, and J. Sun. Mindmap: Knowledge graph prompting sparks graph of thoughts in large language models, 2024. URL <https://arxiv.org/abs/2308.09729>.
- World Customs Organization. Annual report 2022-2023, 2023. Accessed: 2024-11-13.
- Z. Xu, M. J. Cruz, M. Guevara, T. Wang, M. Deshpande, X. Wang, and Z. Li. Retrieval-augmented generation with knowledge graphs for customer service question answering. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR 2024, page 2905–2909. ACM, July 2024. doi: 10.1145/3626772.3661370. URL <http://dx.doi.org/10.1145/3626772.3661370>.
- J. Yang, H. Jin, R. Tang, X. Han, Q. Feng, H. Jiang, B. Yin, and X. Hu. Harnessing the power of llms in practice: A survey on chatgpt and beyond, 2023. URL <https://arxiv.org/abs/2304.13712>.
- S.-J. Yen and Y.-S. Lee. Cluster-based under-sampling approaches for imbalanced data distributions. *Expert Systems with Applications*, 36(3):5718–5727, Apr. 2009. ISSN 0957-4174. doi: 10.1016/j.eswa.2008.06.108. URL <http://dx.doi.org/10.1016/j.eswa.2008.06.108>.
- C. Yu, Y. Shen, Y. Mao, and L. Cai. Constrained sequence-to-tree generation for hierarchical text classification, 2022. URL <https://arxiv.org/abs/2204.00811>.
- D. Zhang, Z. Yuan, Y. Liu, F. Zhuang, and H. Xiong. E-bert: A phrase and product knowledge enhanced language model for e-commerce. 2020. URL <https://api.semanticscholar.org/CorpusID:221516807>.
- H. Zhao, F. Yang, B. Shen, H. Lakkaraju, and M. Du. Towards uncovering how large language model works: An explainability perspective, 2024. URL <https://arxiv.org/abs/2402.10688>.
- Y. Zhu, J.-C. Gu, C. Sikora, H. Ko, Y. Liu, C.-C. Lin, L. Shu, L. Luo, L. Meng, B. Liu, and J. Chen. Accelerating inference of retrieval-augmented generation via sparse context selection, 2024. URL <https://arxiv.org/abs/2405.16178>.

Appendix

Appendix - 1 Group Work

Appendix - 1.1 Exploratory Data Analysis

Appendix - 1.1.1 Hierarchical Analysis

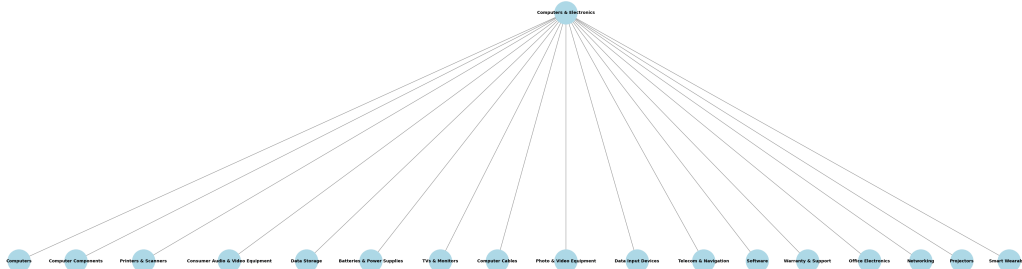


Figure 17: First and Second Levels of Product Categories in WDC and Icecat

Appendix - 1.1.2 Title Analysis

Table 14: Descriptive Statistics Summary for Icecat and WDC

Statistic	Icecat Value	WDC Value
Count	765,473	2,984
Mean	70.42	54.32
Standard Deviation	36.26	19.66
Minimum	2.00	6.00
25th Percentile	45.00	40.00
Median (50th Percentile)	57.00	54.00
75th Percentile	90.00	66.00
Maximum	537.00	175.00

Appendix - 1.1.3 Description Analysis

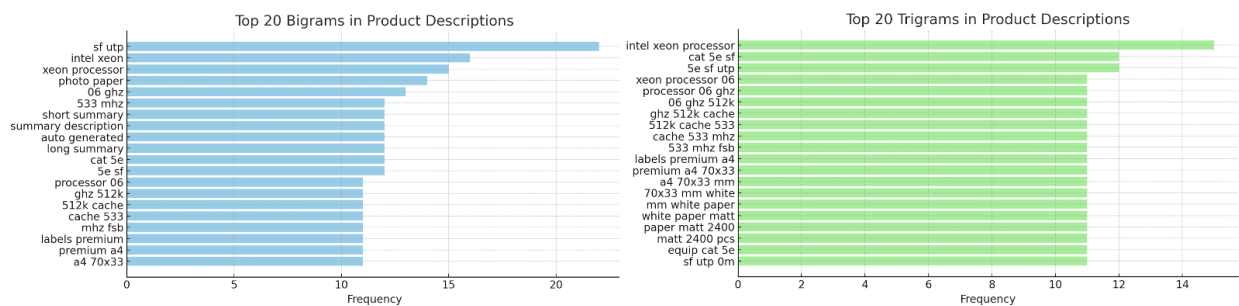


Figure 18: Top 20 Bigrams and Trigrams in Product Descriptions - WDC

Appendix

Appendix - 1.1.4 Brand Analysis

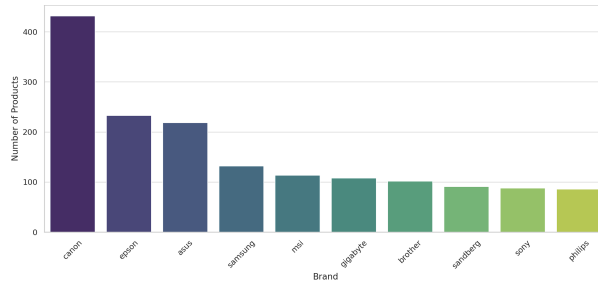


Figure 19: Top 10 Most Frequent Brands - WDC

Appendix

Appendix - 1.1.5 Misclassifications Analysis

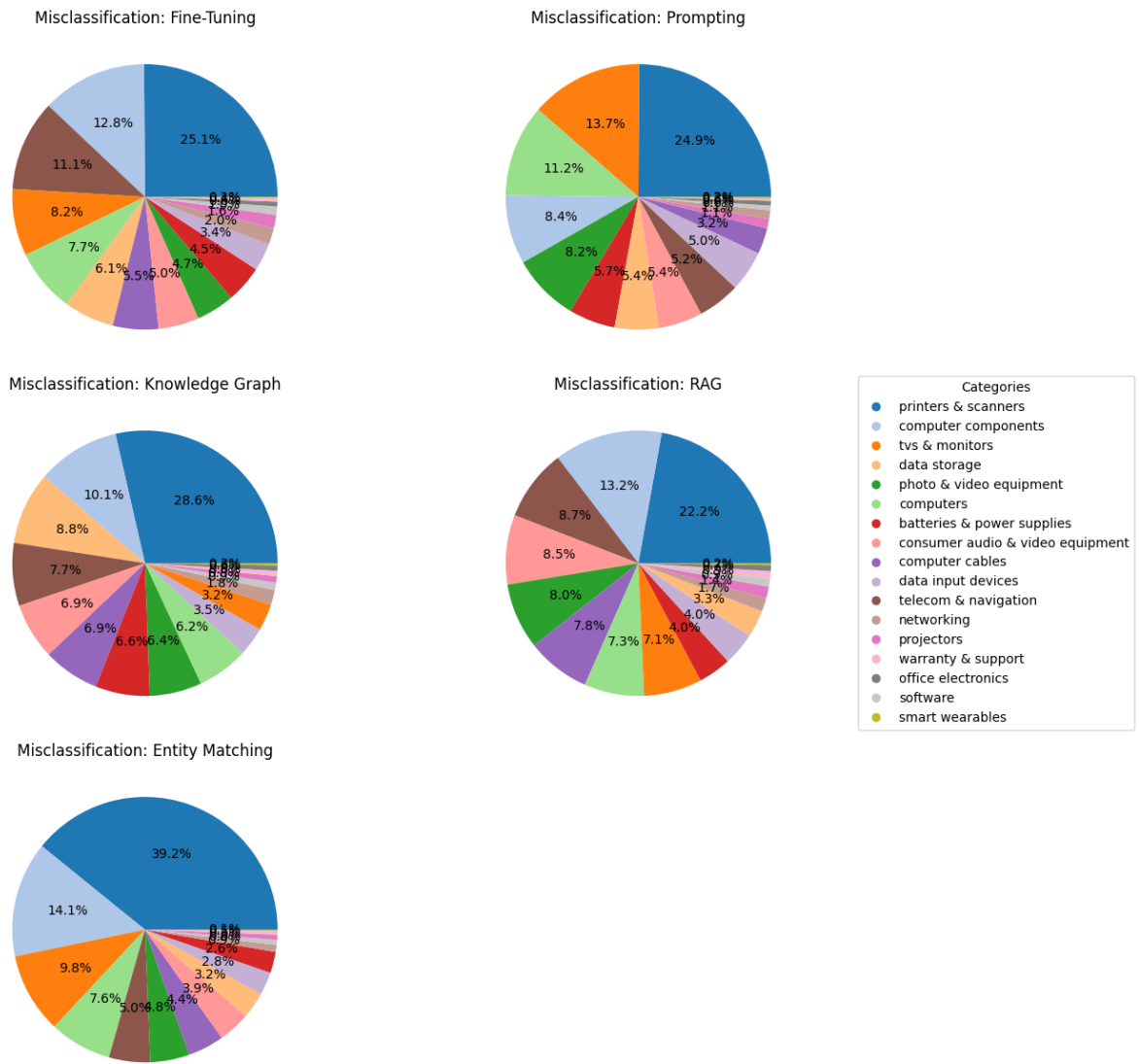


Figure 20: Normalized Misclassification

Appendix - 2 Knowledge Graphs

```

@Language.factory("gliner_spacy")
def create_gliner_spacy_component(nlp, name, labels):
    def gliner_spacy_component(doc):
        # Add your custom processing logic here
        return doc
    return gliner_spacy_component

nlp = spacy.load("en_core_web_sm")
nlp.add_pipe("gliner_spacy", config={'labels': ['electronics producer']})

def extract_org_entities(text, chunk_size=50000):
    """
    Extract organization entities from text by processing in chunks.
    Args:
        text (str): The input text to process.
        chunk_size (int): The maximum size of each chunk to process.
    Returns:
        List[str]: List of unique organization names extracted from the text.
    """
    org_entities = set() # Use a set to avoid duplicates

    for i in range(0, len(text), chunk_size):
        chunk = text[i:i + chunk_size]
        doc = nlp(chunk)

        # Add found organizations to the set
        org_entities.update(ent.text for ent in doc.ents if ent.label_ == "ORG")

    return list(org_entities) # Convert back to a list

```

Figure 21: Spacy Gliner

```

def create_prompt_for_orgs(title):
    return f"""
    Identify organizations or brands mentioned in the following product title.

    Title: {title}

    Output the names of organizations or brands as a comma-separated list (e.g., Sony, Microsoft, Samsung).
    """

```

Figure 22: LLM NER Prompt

Appendix

	Level	Source	Accuracy
0	2	Overall	0.87
1	2	KG-Enhanced LLM (Majority Vote)	0.72
2	2	Strong LLM (Llama 70B)	0.66
3	2	Weak LLM (All Equal)	0.94
4	3	Overall	0.80
5	3	KG-Enhanced LLM (Majority Vote)	0.77
6	3	No Child Categories	1.00
7	3	Strong LLM (Llama 70B)	0.58
8	3	Weak LLM (All Equal)	0.85
9	4	Overall	0.74
10	4	KG-Enhanced LLM (Majority Vote)	0.42
11	4	No Child Categories	1.00
12	4	Strong LLM (Llama 70B)	0.46
13	4	Weak LLM (All Equal)	0.90

Figure 23: Framework Results per Source and Level

	Level	Weak LLM	KG-enhanced LLM	Strong LLM (Llama 70B)
0	Level 2	69.235925	23.022788	7.741287
1	Level 3	47.687668	46.045576	6.233244
2	Level 4	66.387399	21.816354	11.762735

Figure 24: Framework LLM Usage per Level