

Detecting Counterfeit Fashion Products Online using AI and Web Scraping

Tushita Agarwal¹, Devendra Sharma²

¹Pace University, New York, USA

²Xecute Solutions LLC, New Jersey, USA

ARTICLE INFO

ABSTRACT

Received: 20 Sep 2023

Accepted: 25 Nov 2023

An important danger to brand integrity, income, and customer confidence is the widespread presence of fake luxury fashion items in online marketplaces. This study suggests an artificial intelligence (AI) framework that combines picture recognition, natural language processing (NLP), and web scraping methods to detect and flag fake listings on e-commerce platforms. A machine learning (ML) algorithm is used in the present study to identify counterfeit customer evaluations in a dataset. To determine the legitimacy of these reviews and then evaluate their validity, the program makes use of predictive modelling methods. The method employs convolutional neural networks (CNNs) to identify visual differences between genuine and fake products and uses natural language processing (NLP) to examine product evaluations and descriptions for questionable trends after automatically gathering product data and photos from many websites. In the current day, this procedure is especially crucial in domains like business intelligence. In this work, we suggest an AI system that works as a personal assistant for a designer of fashion products. The architecture of the system and all of its parts are shown, with special attention paid to the subsystems for data gathering and data clustering. Our use case scenario involves retrieving datasets of clothing items from two distinct sources and using Natural Language Processes to convert them into a certain format. Comparative results are shown after the two datasets are grouped independently using various mixed-type clustering algorithms, demonstrating the value of the clustering process in the issue of recommending apparel products.

Keywords: Fashion Products, Brand Integrity, Online Marketplaces, AI-Driven, Machine Learning (ML) Algorithm, E-Commerce Platforms, Natural Language Processes, Convolutional Neural Networks (CNNs), Clustering Algorithms, Fashion Product.

I. INTRODUCTION

Global efforts are being made to combat the problem of cybercrime since the risk of fraud exposure has grown dramatically due to the growing usage of the internet and online commerce. In an effort to combat the trafficking of counterfeit products, 21 nations have recently joined forces under the Europol-backed Europe-wide operation Aphrodite.

In 2020, the fashion industry saw a radical transformation. The sector had its worst year ever, with almost three quarters of listed businesses losing money as the coronavirus epidemic sent shockwaves throughout the globe [1]. As the year came to a close, several areas were seeing a second wave of illnesses, which caused disruptions to supply chains and changes in consumer behaviour.

Despite understanding that we will need to take full advantage of the positive aspects of life and business in the next year, a chaotic and frightening year has left us all searching for them. In fact, the McKinsey Global Fashion Index research predicts that fashion firms will report a 90 percent drop in economic profit in 2020, after a 4 percent increase in 2019 [1, 2]. Two possibilities are the focus of our estimates for industry performance in the next year, given the persistent uncertainty.

In recent years, fashion studies have drawn more interest from the computer vision, machine learning, and multimedia fields as it has become an intriguing challenge for computer scientists to handle fashion large data using artificial intelligence (AI) [2]. This study presents the state of fashion research and offers a taxonomy of these

investigations, which range from low-level fashion identification to middle-level fashion comprehension to high-level fashion applications.

Over the last year, there has been a noticeable change in consumer behaviour as a result of individuals staying inside to protect themselves from the virus, travel restrictions, and retail closures worldwide. But as digital consumption continues to dominate and rise in 2021, [1], businesses need to provide more sociable and engaging experiences to entice customers to interact.

This approach has various applications [3, 4], but it is particularly helpful in business intelligence. In order to remain competitive in the twenty-first century, a company must have a strong online presence, and web scraping is a crucial tool for such.

Spammers often utilise unsolicited online communications, or spam, to promote their business's product [5]. Spammers' unfavourable remarks about products from competing companies have grown to be a serious issue for companies that only rely on digital marketing strategies. In today's society, it is usual for spammers to submit phoney evaluations using praise phrases like "awesome," "excellent," or "very good" [5].

In order to get further input on items, those that make unsolicited messages often employ the terms mentioned above [5, 6]. This has become a really serious problem. Utilising effective methods that can identify individuals who engage in this kind of spam is essential. Amazon and Flipkart are enterprises with substantial product inventories, data collecting, and customer feedback [7, 8].

Since data aids in decision-making, it is essential for businesses and organisations. The majority of the data, in particular, are currently accessible online [8]. The first step in any data science research and development process is data collection, which involves gathering data from either private source, such as firm sales records or financial reports, or from public sources, such as periodicals, websites, and open data [9].

The three mains, connected stages of internet scraping are website analysis, website crawling, and information structure [10]. Data mining and web scraping vary in that the latter does not need data analysis, while the former does. Moreover, data mining uses quite complex statistical methods. Because a lot of the necessary functionality can be executed efficiently using a number of readily available tools and frameworks, web scraping is often a very easy procedure [11]. You may create custom HTTP requests with various headers and payloads using most web scraping software.

Additionally, consumers on social media now place a higher emphasis on customer feedback. Before making a purchase, many customers look for human connection to build brand trust since they are leery of one-sided advertising. In this regard, social media provides customers with convenient access to genuine evaluations from other people [11]. In order to browse across brands and influence their purchasing choices, consumers may search posts or tags for a particular brand or product [13]. However, social media's rise has benefited more than just customers and companies. The fight against counterfeit goods has moved to social media. Since counterfeiters are always making new accounts and posts to offer phoney luxury goods for almost nothing, eliminating online counterfeit goods is like playing a game of "whack-a-mole." Owners of trademarks must devote their limited time and resources to trademark monitoring and persistently pursuing the removal of counterfeit listings.

Luxury businesses may now more effectively support the word-of-mouth marketing strategy thanks to the rise of social media [19]. One of the most important sources of market knowledge for customers nowadays is interpersonal conversation about goods and services [59]. According to a 2018 study, 40% of luxury purchases are impacted by what customers see online, highlighting the significance of social media and online platforms for a premium brand's exposure and standing. Conversely, social media offers luxury firms a potent instrument for market research on consumer trends and behaviours [39, 52]. Recent fashion shows, new product releases, or celebrity appearances may influence social media conversations among consumers.

Social media hashtags can assist luxury businesses in navigating and sifting through consumer evaluations and preferences. In order to attract clients that share their social ideals, companies often utilise social media to track their reputation via online influencers [55]. But social media's ease also works in counterfeiters' favour. Social media gives

counterfeiters a shell of anonymity that helps them avoid detection in addition to easy access. Even when its posts are taken down or its account is disabled [25], a counterfeiter may quickly and cheaply create a new account for free to keep selling phoney goods [29]. Additionally, it is very difficult to monitor and trace internet activity in real time due to the large number of fake postings.

For instance, it's predicted that up to 95 million fake accounts may be present on Instagram [20]. Users face chaos as a result of the massive volume of fake posts that many of these bot accounts publish each day. Online trademark enforcement is an unwinnable "whack-a-mole" game because of these difficulties, with enforcers having limited whacking resources for an infinite number of moles[30].

The fashion outfits industry is being driven by fast fashion, which means that retail marketplaces must produce goods more rapidly while still meeting consumer wants and current trends [12, 13]. In order to promote the development of innovative ideas, solve the problem of supply and demand balance, improve customer service, assist designers, and increase overall efficiency, AI-based technologies are applied to a company's whole supply chain. AI techniques have been applied in an increasing number of fashion industry initiatives recently, including those run by Google and Amazon [14].

These days, the regular use of e-commerce websites and the massive amount of data collected by fashion businesses provide solutions related to the fashion development processes that apply artificial intelligence techniques [15]. Well-known fashion brands have provided incredible AI-driven solutions, such as the Hugo Boss AI Capsule Collection¹, in which an AI system develops a new collection, and Reimagine Retail, [16], a collaboration between the Fashion Institute of Technology, IBM, and Tommy Hilfiger that aims to improve the process of designing and forecast future market trends.

This paper highlights on the creative side of the fashion industry, the fashion design process. This is accomplished by proposing a sophisticated and semi-autonomous decision-support system for fashion designers. By compiling, organising, and combining data from many sources, this system might act as a personal assistant before recommending fashion items according to the designer's preferences [17]. The information associated with the clothing images is analysed using Natural Language Processing (NLP) techniques, computer vision algorithms are used to extract features from the images and improve their meta-data, and machine learning techniques are used to analyse the raw data and build models that can aid in decision-making [18]. The bulk of the many research that have been published in the field of clothing data analysis have concentrated on product recommendation, dataset construction, clothing classification, and feature extraction from pictures [19].

A selection of characteristics and labels were applied to the 800,000 images that comprise the Deep Fashion collection [19]. The challenge includes the following steps to help you comprehend the features of a garment pictures:

- a) Image retrieval with description,
- b) Learning features for the human body's top and bottom,
- c) Deep learning feature extraction, [20],
- d) Applying pose estimation methods, and
- e) Deep learning for hierarchical feature representation learning [21].

It is hard to overstate how important sophisticated background removal techniques are to meeting these client demands [22]. The correct functioning of customisation algorithms depends on producing high-quality, visually appealing product photographs, which calls for precise and efficient background removal. In the fashion industry, background reduction is a crucial image processing technique that is often used, especially in the e-commerce sector [23]. In order to make educated decisions while making online purchases, customers mostly rely on product images [23]. Their decisions to buy could be heavily impacted by how accurate and high-quality the visual experience is. In addition to eliminating distractions, removing the background from product images helps highlight the key features of the item, making it easier for customers to evaluate and spend.

In the fashion sector, removing the background from product photos may significantly improve their visual appeal. When consumers browse e-commerce websites, they look for high-quality photos that accurately represent the

products they are interested in [11]. The image's clean look and aesthetic appeal are improved by isolating the subject and removing the background. Given that customers make many decisions based on appearance, this is especially important in the fashion industry [13]. Customers may find it simpler to evaluate different products and make informed decisions if product photographs are standardised by removing backgrounds [20]. Backdrop reduction is an essential technique in the fashion industry for creating polished, superior pictures of goods [23].

Because of their ability to adapt to variations in size, orientation, and appearance in fashion images, transformers were especially selected [24]. Their capacity to understand spatial relationships and the semantic meaning of different portions in an image enables accurate and reliable background removal [25]. In practice, fashion datasets may be diverse and big, requiring models with high generalisation abilities. Transformers can effectively employ the vast quantity of training data thanks to their self-attention processes and parallel processing capabilities, which enhance accuracy and performance [25]. By incorporating transformers into the process, professionals in the fashion industry may benefit from automated and efficient background removal techniques that provide superior results [26]. These tactics enhance the visual experience for customers by ensuring that fashion goods are presented in an alluring and contextually appropriate manner [11]. As a result, a significant advancement in the fashion industry is the use of transformers in background-removal techniques, which enable accurate and efficient processing of fashion images while maintaining the highest standards of visual quality.

Customers can easily explore and interact with the platform to locate apparel that fits their own preferences because to its basic user interface, which is accessible via both web and mobile applications. This enables the deep learning and recommendation system to extract and analyse fashion item properties with previously unheard-of accuracy by accurately removing the background. OutfitAI ensures that recommendations are up-to-date and personalised by comparing these features to a comprehensive product database that is regularly gathered from leading e-tailers using transformer-based neural networks [22, 25].

A continual feedback loop with analytics ensures that the system adjusts to shifting user preferences and fashion trends, while a robust cloud architecture provides scalable computing resources and secure data management to support the design. With its advanced technical integration and user-centric design, OutfitAI offers a distinctive approach to outfit-based shopping, making it a model of innovation in retail technology and a vital tool for fashion enthusiasts. Beyond its technical expertise, OutfitAI is committed to data ethics [26], offering robust privacy protections and user control over personal data, unlike less transparent commercial methods. OutfitAI's commitment to ethical buying is shown by the way its recommendation engine gives preference to fashion items made in an environmentally responsible and legal manner [27].

These cloud services provide customers access to pre-built machine learning models that they may use to assess operational effectiveness and other business-related factors, including small and medium-sized enterprises [53]. However, the market isn't solely occupied by the well-known, old companies of big IT. Every year, more businesses join the market as the need for affordable business solutions keeps increasing. As of 2020, machine learning and artificial intelligence are really two of the IT industry's fastest-growing subsectors [49].

Let's return to the data, however. At this point, you may be wondering where the data needed to train the machine learning models comes from. The data that initially trained the model comes from the model's creators for small or mid-sized businesses using machine learning (ML) in a subscription basis, such as a subscriber via Google's Cloud platform [23]. The small firm then enters its own data—such as sales records, location data, inventory data, customer surveys, etc.—that is produced by its own operations into the algorithm as "input data." After processing the data, the model forecasts future encounters with customers [39].

Numerous marketing processes may be streamlined and expedited by AI technology, which can also enhance consumer experiences [55, 58]. AI makes it simple to convert human voice to text. The speech's substance and the appropriate response might then be predicted by further analysing the text. Furthermore, the user information gathered may be used to target a particular user and customise recommendations for goods or services. Voice processing, picture and pattern identification and processing, text processing, and so forth are the primary technologies associated with the use of AI [26]. Examples of AI applications in marketing automation using a variety of technologies are shown in Table 1.

Table 1 AI applications in marketing automation examples.

Technologies	Examples
Voice processing	Requests for purchases are spoken via a device. Using virtual assistants to aid in task execution
Text processing	A virtual assistant-led tour of a retail mall. [29] By answering their questions, a virtual assistant built within a mobile banking app manages customer demands by themselves. A GPS navigation system may indicate the path to the chosen location and recommend nearby or transportation-related attractions.
Image recognition and processing	After determining the skin state of the face, each person chooses the kind of face cream based on an analysis of their data and picture, including the present weather. Using a picture as a starting point for an internet search yields identical results along with related or complimentary goods. [30] In a clothes store, electronic mirrors align the assortment with the customer's look, preferences, and style.
Decision-making	Creating a personal savings strategy after analysing one's account balance, receipts, expenditures, and costs. Based on the traveler's interests, travel destinations (certain neighbourhoods and attractions) are tailored to each user's profile. A chatbot may create recipes utilising the items that are already in the customer's house depending on their preferences. price matching that is dynamic and based on consumers' past purchases and websites visited. Synchronisation of consumer information from every channel via which the brand may be contacted (social media, internet, email, and phone calls).

AI has enormous potential to enhance existing marketing strategies and provide whole new methods for providing and allocating value to clients. Marketing automation is the use of software to automate marketing tasks including campaign management, customer data integration, and consumer segmentation [29]. With the proper marketing automation technology, sales may be increased at a much-reduced cost and with a significantly reduced human resource need. It is possible to access resources like social network services (SNS), TV shows, retail websites, and more without being limited by the moment or place [59].

The field of computer vision has developed throughout time in various areas, such as information interpretation and picture pattern extraction. Image processing and pattern recognition are combined to form computer vision [36]. Computer vision produces visual comprehension as its output. Thus, computer vision is the field that specialises in obtaining information from pictures, and it is dependent on the computer technology system, whether it is associated with image recognition, quality, or enhancement [19]. Numerous challenging computer vision issues have been resolved using the open, cutting-edge technique known as deep learning. As a result, implementing several sophisticated and novel applications is now feasible. Since digital learning techniques have produced groundbreaking advancements in computer vision and machine learning, deep learning is now the newest trend in machine learning research and development.

Online scraping's definition, stages, and technique are examined in this article, along with its relationships to artificial intelligence, data science, business intelligence, big data, and cyber security. possible future advancements, some of the main benefits of utilising the Python language for web scraping, [28], and a special emphasis on raising awareness of the ethical and legal issues [29].

II. RELATED WORK

Zhao, L. (2019) [32] Techniques for knowledge discovery have long been used in practical areas like corporate intelligence and marketing. Comparably, computer scientists have paid little attention to the fashion industry and other production sectors. Our knowledge of the fashion clothing sector might be greatly improved by using knowledge discovery techniques and large-scale datasets gathered from sites like Twitter and Instagram, as multimedia data from the Web and social media becomes more widely available. Here, the practice of knockoffs is one of the challenges at the core of the modern dynamics and structure of the fashion business. We provide a first description of how brands copy one other's designs by combining Web scraping and network science tools. Research like this could be among the first instances of an emerging subject, which we call and characterise as "fashion informatics."

Lindner, R. (2020) [33] E-commerce fraud has been on the rise at an alarming rate; in Austria, it reached a record high in 2019, growing by 32.3% from the previous year. People's funds are often and seriously harmed by fraudulent retailers, fake subscription services, and fake goods. Manual preventive measures are insufficient due to the high volume of cases submitted for assessment and their increasing frequency. Therefore, it is essential to increase the efficacy of recognising phoney shops and reduce the window of opportunity. Using machine learning approaches, this work provides a way to classify phoney web stores based only on how similar their source code architectures appear.

Rehman, R. U. (2019) [34] Web scraping, sometimes referred to as web harvesting or web data extraction, is essentially the process of extracting data from websites on the World Wide Web. Alternatively said, it may be described as the methodical process of extracting and combining material that has been collected from the internet. One subset of web scraping is price scraping, which is used to extract pricing from e-commerce websites for a variety of products. Bots scrape the website to do this task. Any program or code may be used to create these bots. Pricing scraping may be against the law if a rival lowers his own pricing using the scraped prices in order to increase sales. In order to identify these bots, we may use a number of methods to distinguish between them and persons.

Rahm, E. (2016) [35] With losses believed to be in the billions annually, online goods counterfeiting is becoming a bigger problem. Even if it could be difficult for individuals to identify suspected counterfeits, techniques for doing so must be mostly automated due to the large number of online retailers and products. The authors recommend employing a semi-automatic method to analyse product offers on online platforms and spot possibly fraudulent offers based on a number of factors. Such questionable offers should be manually verified by a domain specialist. The technique compares and groups similar product offers and assesses the suspiciousness of counterfeit goods based on a number of characteristics.

Thankachan, B. (2021) [36] In recent years, the retail market business has expanded to include online product sales as well as the ability for consumers to provide their insightful opinions, recommendations, and ideas. A vast text-based review collection contains a variety of views about various items that are accessible online, which are extracted and identified using opinion summarisation and categorisation methods. SentiWordNet, Random Forest, Naive Bayes, Logistic Regression, and K-Nearest Neighbours algorithms are used in this work to automatically identify the feelings represented in the English text of Amazon and Flipkart items. On the basis of five important criteria, it provides a thorough comparative review of these sentiment analysis methods and approaches. The PCSA system, or Product Comment Summariser and Analyser, was also suggested in the study. It is a general and automated comment analyser that is very successful at determining the polarity of the thoughts and remarks. After summarising the remarks, it divides them into the predetermined positive, negative, and neutral categories.

Wan Aziz, W. A. H. (2020) [37] The growth of the Internet revolution and digital marketing has led to a surge in marketing automation, neuromarketing, and user personalisation. Social engines have produced a wealth of consumer data, which has led to a significant advancement in AI-based marketing applications. With regard to its potential to become a ubiquitous aspect of today's competitive environment, this article aims to discuss the widespread use of artificial intelligence (AI) in marketing. In order to fully comprehend how AI is used in marketing, the fundamentals of the technology are explained. Despite the opportunities presented, it is important to recognise the risks associated with AI in marketing automation when dealing with the dynamic nature of marketing and people's sensitive data. This article calls attention to security concerns pertaining to user privacy and potential harmful activity.

Shukla, S. S. (2020) [38] We address two issues that might be difficult for the e-commerce sector. A challenge that sellers have while submitting product photos to the platform for sale and the resulting manual tagging is one of them. Because of the misclassifications it causes, it is not included in search results. The second issue is the possibility of order placement bottlenecks when a buyer recognises a picture but may not know the correct keywords. In order to fully realise the promise of e-commerce, an image-based search algorithm may allow users to click on an image of an item and look for similar things without typing. In order to address these two issues, we investigate machine learning techniques in this study.

Chinta, S. (2021) [39] The use of machine learning algorithms with big data analytics presents exciting opportunities for improving predictive insights across a range of fields in the age of exponential data expansion. This study offers a thorough framework intended to make it easier to successfully incorporate machine learning methods into large data settings. The suggested approach seeks to enhance model accuracy, optimise data processing, and provide actionable insights by tackling the inherent difficulties of conventional data analysis techniques. The paper identifies gaps in the literature, examines current approaches for integrating machine learning, and investigates the state of big data analytics. This article clarifies the essential elements required for effective implementation, such as data pretreatment, algorithm selection, and performance assessment, by developing an organised methodology. The framework's real-world implementation is shown via case studies, which show significant gains in decision-making and prediction accuracy. The results open the door for further study and business applications by highlighting the revolutionary potential of machine learning in big data analytics.

Johnson, S. (2022) [40] Online contacts between customers and sellers are progressively replacing in-person ones since the creation of the Internet. This chapter reviews the definition of online consumer fraud as well as common fraud tactics. It looks at some of the current tactics that businesses and non-governmental groups use to detect and prevent online consumer fraud. The chapter discusses darknet markets and how they lead to online consumer fraud. Darknet markets are online marketplaces that, among other things, provide anonymity for the selling of both legal and illicit goods. The essay ends with an explanation of how data science methods might be used to help detect and prevent consumer internet fraud. A subtype of the deep web, the dark web anonymises users and servers.

Thankachan, B. (2021) [41] Retail industry sectors have expanded in recent years to include online product sales as well as the ability for consumers to provide insightful comments, ideas, and suggestions. This paper's objective is to provide an automated comment analyser. Additionally, provide a categorisation method and automated comment analyser that can efficiently identify the polarity of consumer comments gathered from the Flipkart and Amazon data domains. The high volume of reviews should be handled by this approach. Five top supervised learning classifiers, including NB, LR, SentiWordNet, RF, and KNN, should be used to classify the comments into positive, negative, and neutral categories. Their experimental findings and difficulties are also covered in the publication. Thus, for a large collection of reviews, this research demonstrates the optimal use of feature extraction, positive-negative sentiment, Amazon site source, and mobile phone in the current algorithms. Preliminary definitions, characteristics of information extraction and retrieval, the function of machine learning, and comment mining were all covered.

Rungta, M. (2017) [42] Web scraping is a major problem in computer science. Popular position-based or structure-based online scraping methods have the drawback of requiring human reconfiguration each time the structure of the web page changes. In this study, we try to solve the problem of information extraction for web pages composed of recurring blocks. We extract these blocks and their constituent properties using a novel classification-based approach. Our approach retrieves product offers with high accuracy from offer-aggregator websites. It is also quite adaptable to changes in a website's structure.

Kayed, M. (2021) [43] It has recently been discovered that e-commerce (EC) websites provide a wealth of helpful information that is beyond the capabilities of human cognitive comprehension. Previous study writers have created opinion summarisation systems based on customer reviews to assist consumers in weighing their options when making a purchase. They disregarded the manufacturers' template material, even though it included the most helpful product details and, unlike reviews, the content was grammatically perfect. Thus, utilising a mix of template data and customer evaluations in two primary stages, this research suggests an approach called SEOpinion (summarisation and

exploration of opinions) to summarise characteristics and identify opinion(s) about them. Initially, a hierarchy of aspects is created from the template using the hierarchical aspect extraction (HAE) step.

Stuenkel, B. (2021) [44] In the hypothetical situation described in this letter, a new technology business scans public websites and webpages to gather information like names, addresses, phone numbers, and even dates of birth. In this hypothetical scenario, the IT firm then aggregates the information into a dataset that contains the PII. After that, the startup tech business uses the dataset to develop its artificial intelligence (AI) skills. A larger software company then acquires the fledgling IT firm. The AI model and training dataset are included in the purchase and will be used in the software product suite, or "code stack," of the acquiring enterprise.

III. METHODS

In this part, the proposed decision-support system for the designers' creative processes is shown. For people who are unfamiliar with the concept of action research preparation, the system's user-friendly design and ability to automatically model the designer's preferences make it simple to use [21]. The system is composed of two interconnected sections:

- **Offline component:** This element does;
 - (a) Gathering information from both external and internal sources,
 - (b) Database administration and data storage, and
 - (c) Data analysis procedures that generate artificial models that provide end users tailored suggestions.
- **Online component:** The user interface (UI) makes up the majority of this component. The graphical user interface (UI) makes it simple for users—who are often fashion designers with little technical expertise—to define their settings, visualise their findings, and comment on the system's output [22]. While the main subsystems and processes are further examined in the next subsections, the overall system architecture is shown in Fig. 1 [23].

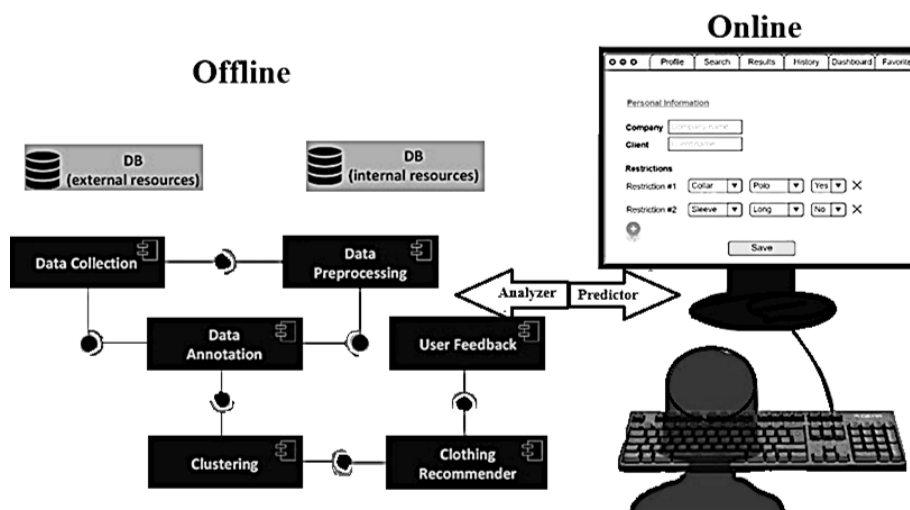


Fig. 1 The system architecture that is suggested. [23]

3.1 Data Collection

Current fashion trends influence the production schedule, rules, and design aesthetics of each company [24]. A creative team's inspiration or starting point is often clothing from the company's previous collections, which they subsequently alter to suit contemporary fashion trends [25]. Designers love to go through the collections of other well-known internet merchants in order to get new ideas. The system uses a web crawler known as the e-shops crawler to do this, which is capable of extracting information on clothes, including photos of items and their meta-data [26].

3.2 Data Processing

For each clothing image, this subsystem has to retrieve the garment attributes from the accompanying meta-data [29]. The attributes listed below were selected from the available meta-data, coupled with a few trustworthy examples:

1. **Product category:** dress, overall, pajamas, shorts, skirts [30].
2. **Product Subcategory:** jacket, coat, T-shirt, leggings.
3. **Length:** short, long, knee.
4. **Sleeve:** short, $\frac{3}{4}$ length, sleeveless.
5. **Collar Design:** shirt collar, peter pan, mao collar.
6. **Neck Design:** V-neck, square neck.
7. **Fit:** regular, slim.

3.3 Meta-Data-Based Clustering

After data collection and annotation, all of the data are available in a common format (row data), which can be analysed using established state-of-the-art techniques. Clustering is a widely used technique for organising data into sets of similar products [30]. By accelerating the look-up subprocess, clustering aids in accelerating the recommendation process when working with huge volumes of data.

The classification of data will determine which clustering algorithm is employed. In general, clothing data may be described using both categorical (like product category) and numerical (like product price) properties [31]. The methods for clustering mixed-type data may be thoroughly examined [32].

3.4 Clothes Suggestions and User Input

Our Clothing Recommender is the most important component of our system, combining all of the previously reported analytical results to create models that provide personalised predictions and product recommendations [22]. User preferences, organisational regulations, and internal and external data are all taken into account. The online component's designer may search for products using keywords thanks to the user interface (UI) [33]. After each product search, the designer may read over the extracted results and save the items they like on their dashboard. Should the user be dissatisfied with the recommendations, they may either refresh their selections or request new ones.

3.5 Using Rich Annotations to Support Sturdy Clothes Recognition and Retrieval

Building clothing databases has been a major factor in recent clothes identification breakthroughs. The number of annotations in existing datasets is restricted, and they struggle to handle the different problems that arise in real-world applications [53]. We provide DeepFashion, a large clothing dataset with thorough annotations, in this paper. It includes more than 800,000 photographs, all of which are extensively annotated with enormous characteristics, landmarks related to apparel, and correspondence of photographs taken in various contexts, such as retail, [27], street photography, and consumer. Such comprehensive annotations facilitate future study and allow for the creation of strong algorithms for clothing identification [52].

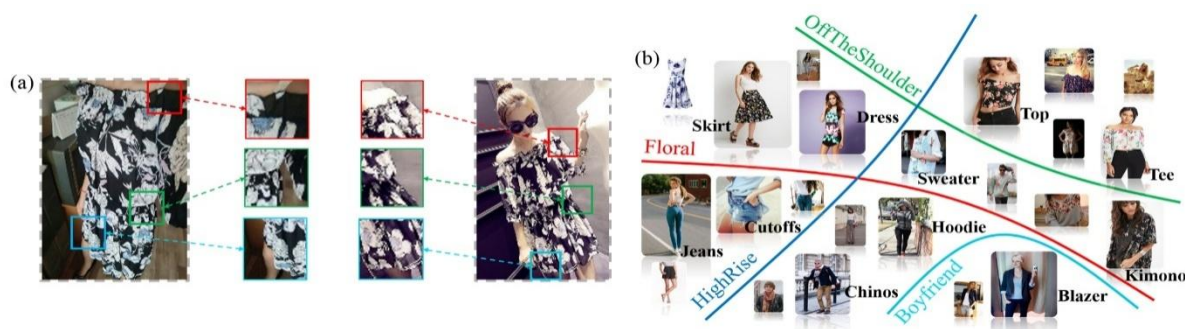


Fig. 2 Deep-Fashion: Using Rich Annotations to Support Sturdy Clothes Recognition and Retrieval. [28]

Classifying clothing is a topic of great interest in scholarly circles. There are typically two groups into which garment categorisation techniques fall. First, the conventional feature extraction techniques for clothing categorisation can also be separated into two categories: one is based on global texture and shape features [56], including Fourier descriptions, geometric invariant distance, local binary patterns (LBP), etc., and the other is based on local feature techniques, such as histogram of orientated gradient (HOG), scale-invariant feature transformation (SIFT), sped up robust features (SURF), etc [8,9,10,11]. Traditional approaches' categorisation accuracy is heavily dependent on the chosen characteristic. The techniques may provide stable and noticeable characteristics with high levels of accuracy in certain situations. The picture of clothes is often diverse, comparable, and intricate.

Robust characteristics cannot be extracted for classification using conventional approaches. Second, convolutional neural networks (CNNs) have been more popular in the categorisation of clothes because to their technical advancements [12,13,14]. Compared to the conventional ways, the method performs better. In order to gain strong clothing attributes, it uses a deep network to extract them without the need for manual settings [15,16]. Deeper networks may perform better overall. The clothing feature is extracted using a convolution network and polling procedures. Nevertheless, CNN does not consider the spatial link between various local characteristics; it merely extracts the picture features [17, 18, 19, 20]. Conventional CNN is unable to overcome the classification performance bottleneck [22].

In order to establish the spatial connection of various characteristics, a new CNN called the capsule network [21,22,23,24] is developed. The dynamic routing method is used to determine the spatial relationship. In vector form, the capsule represents the basic feature unit. The MNIST database, on the other hand, often uses the classic capsule networks [21] for handwritten digit identification, with an input size of 28 x 28 [60]. To extract the multi-scale feature, a unique multi-scale capsule network was developed, using the advantages of the capsule network in spatial characteristics [25, 26]. By using a collection of capsule subspaces rather than just clustering neurones to form capsules, a subspace capsule networking may take use of the concept of capsule networks to describe potential differences in an entity's appearance or implicitly specified features [27]. A multi-lane capsule network is a resource-efficient, separable arrangement of capsule networks that enables parallel processing and achieves excellent accuracy at a lower cost. To put it simply, neither the conventional nor enhanced capsule networks are able to effectively extract strong clothing features [49], making them unsuitable for the precise categorisation of clothing pictures [55].

IV. EXPERIMENTAL RESULTS

A real-world situation is used as a use case to illustrate the importance of the clustering process in the recommendation of clothing items [33]. Our team has partnered with a clothing designer who works for the Greek shop Energiers since he has shown interest in designing the collection for the next season. She finds inspiration in both the current Assos e-shop collections and the clothes designs and productions from the company's previous period [21].

4.1 Datasets

The fashion products from the previous season were taken out of the business database to create the firm dataset, and a web crawler was used to get the relevant E-shop data. Using a variety of labelling for the parameters Product Category, [21], Length, Sleeves, Collar, and Fit, the eShop crawler collected 4674 images for the winter 2020 season. The meta-data of the recovered photographs and a reference to the image's location were stored in a relational database. The meta-data were tokenised and separated into columns by allocating values in the relevant qualities after the plain text was prepared using NLP techniques [22].

4.2 Data Clustering

In this part, the experimental results using the company's and E-shop datasets using the Kmodes, Pam, HAC, FBHC, and VarSel algorithms are shown [23].

Lower Entropy and WSS values and higher Silhouette and Identity values suggest better clustering results. The results of the clustering process differ according on the methods used. Table 1 presents the normalised mutual information of

the methodologies under study in a paired presentation [29]. There is a significant range in the numbers, with most being around 30%.

Table 1 The method evaluated on the Company dataset's normalised mutual information. [30]

	Kmodes	Pam	HAC	FBHC	VarSel
Kmodes	1	0.3295	0.3269	0.3296	0.3296
Pam	0.3985	1	0.3269	0.3189	0.3159
HAC	0.6596	0.3963	1	0.3009	0.3269
FBHC	0.3259	0.3296	0.3296	1	0.2192
VarSel	0.3148	0.3219	0.3189	0.3691	1

A graphical representation of the data transferred throughout the clusters created by the different approaches is provided by the Sankey diagram in Fig. 2 [30]. The Pam approach distributes the data items equally throughout the six groups, as the picture shows, [55], but the Kmodes clustering results have a normal distribution [32].

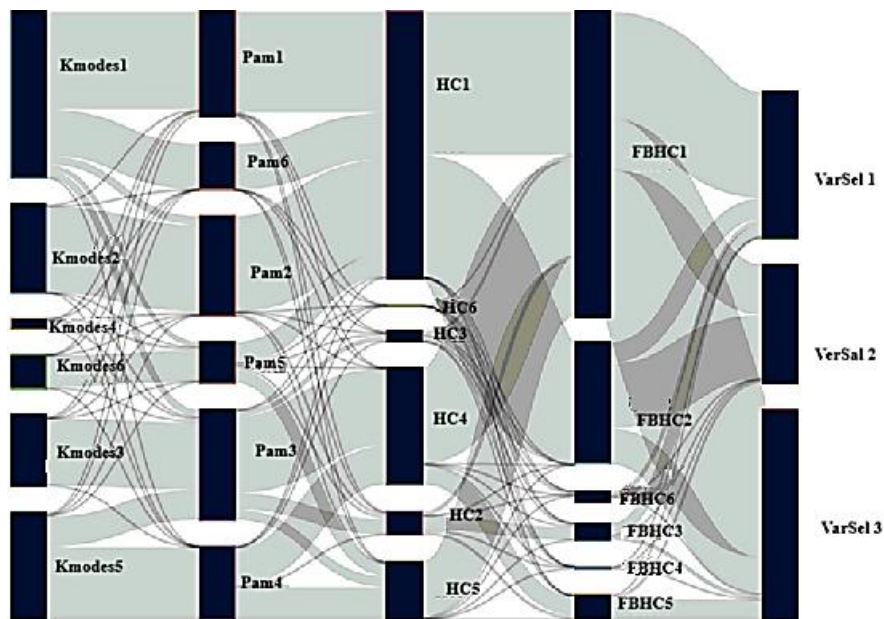


Fig. 1 The clustering results obtained by the methods evaluated on the Company dataset are shown in the Sankey diagram. [33]

These two methods of portioning seem to have comparable distributions [33]. Only three clusters are generated in this case, when normally the VarSel technique would distribute the items uniformly. Hierarchical approaches, on the other hand, result in two large clusters and four considerably smaller clusters, to which the majority of the items are assigned.

Table 2 presents comparison results for the clustering algorithms based on their performance on the assessment metrics. Shown are the mean values of the assessment metrics [29]. The second-highest results are italicised, while the algorithm's best outcomes are highlighted in boldface.

Table 2 A variety of clustering techniques were used to evaluate the Company dataset. [30]

	Kmodes	Pam	HAC	FBHC	VarSel
#Clusters	6	6	6	6	3
Entropy	0.2196	0.2069	0.5496	0.1489	0.5986
Silhouette	0.0498	0.5949	0.2189	0.5498	0.4189
WSS	0.3652	0.2185	0.6596	0.5412	0.6890
Identity	0.5969	0.9659	0.4896	0.2859	0.6589

There are few labels describing each feature in a cluster, as evidenced by the fact that the hierarchical algorithms outperformed the partition-based algorithms in the metrics related to the distance between the objects of each cluster of machines and the partition-based algorithms in the measures of entropy and identity [20]. Due to its proven ability to distribute data fairly among clusters, we decided to use the Pam technique for the rest of the inquiry in this paper. The data distribution into the six categories that Pam was able to gather is shown in two dimensions in Fig. 3 [29].

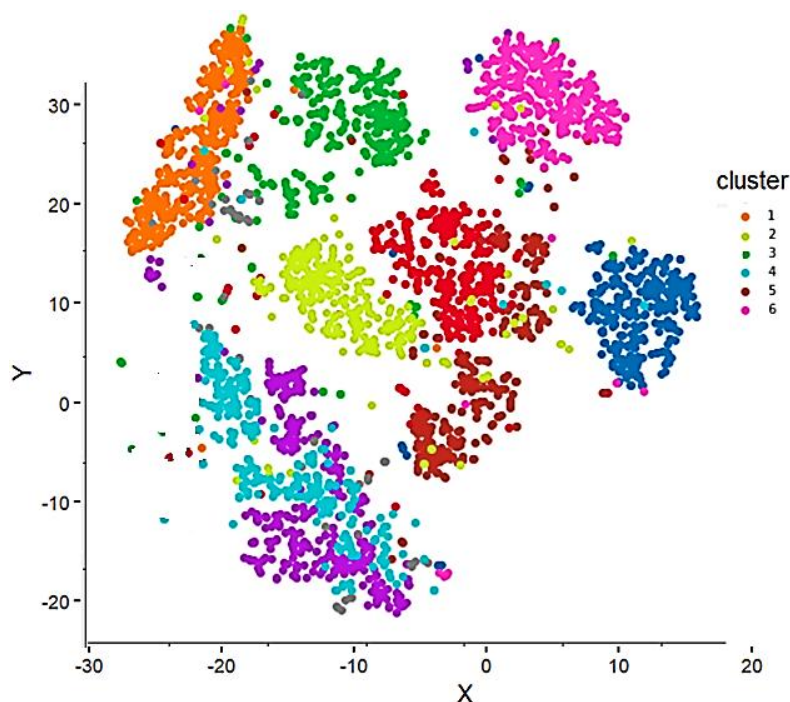


Fig. 3 Pam got a two-dimensional depiction of how the Company dataset was distributed across the six categories. [31]

The centre points of the Company datasets that were obtained using the Pam technique are shown in Tables 3 and 4 [39]. The centroids are computed using the most common attribute values in the row data for each cluster.

Table 3 Using Pam, the clustering centroids were created using the Company dataset. [18]

	Category	Gender	Length	Sleeves	Collar	Neck	Fill
C1	Set	Man	Short	Long	Shirt	Off shoulder	-
C2	Bermuda	Man	Short	Sleeveless	Flat knitted rib	Round	Regular
C3	Blouse	Woman	Medium	Short	Shirt	Round	Regular
C4	Dress	Woman	Short	Sleeveless	Flat knitted rib	Round	-
C5	Blouse	Woman	Short	Short	-	Round	Regular
C6	Leggings	Woman	Capri	-	-	-	Slim

Table 4 The E-shod dataset from the web crawler employing Pam was used for the clustering. [19]

	Category	Gender	Length	Sleeves	Collar	Neck	Fill
C1	Dress	Woman	Medium	Raglan	-	Halter neck	Regular
C2	Shirt	Man	Medium	Long	Stand-up	Collar	Regular
C3	Trousers	Man	Medium	Flared	-	Collar	Cargo

C4	Set	Woman	Knee	Flared	-	Collar	Regular
C5	Romper	Woman	Knee	Flared	-	Off shoulder	Regular
C6	Gardigan	Woman	Knee	Flared	-	V-neck	Regular

The groups' consistency for the criteria Product Category and Gender may be further shown using a heatmap (Fig. 4). Examining the consistency of each group and the distribution of labels across the groups in the two datasets [18] reveals that the Company dataset is separated into six major categories: Set, Bermuda, Blouse for Men and Women, [30], Dress, and Leggings. The following products, on the other hand, set, romper, cardigan, dress, shirt and pants set apart the E-shop dataset [33].

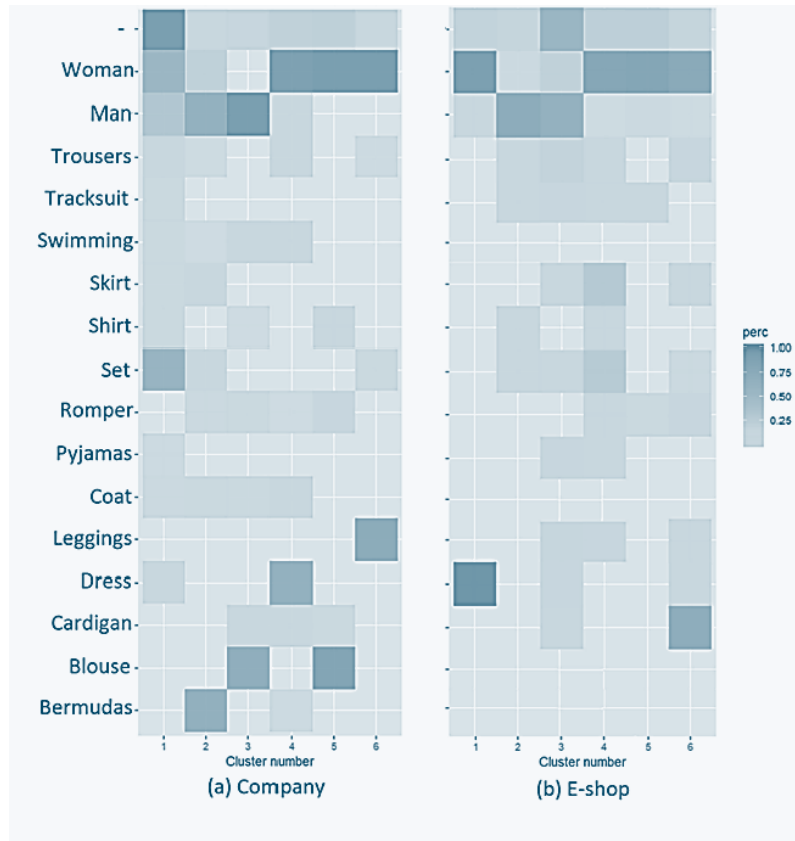


Fig. 4 Label distribution of features for (a) the Company and (b) the E-shop datasets using Pam across the six clusters. [29]

V. DISCUSSION

This research describes a smart technology that automates the routine procedures taken by a fashion designer. Using dictionary mapping and natural language processing (NLP) techniques, the system can transform plain text that goes with the images into clothing features [20], extract information from the designer's company database and online sources, use computer vision to extract novel characteristics from the images, and store all of the data in a relational database in a standard format. Clustering, prediction models, recommender systems, and other advanced machine learning techniques may then be used to handle the processed data [29].

Images of clothing are widely accessible online, particularly via e-commerce platforms. Retrieving such pictures has lately drawn a lot of interest from fields like computer vision and multimedia processing, and it is important for both commercial and social applications [49]. However, such issues are difficult to resolve due to the wide range of clothing styles and appearances, as well as the sheer number of different categories and characteristics [41]. Additionally, the labelling that online vendors give for real-world photos are sometimes inaccurate or lacking [42]. Furthermore, effective learning is hindered by the imbalance among those picture categories [49].

The system's architecture is presented in the study [29], with particular attention paid to the data collecting and processing operations and the clustering techniques that may be used to the grouping of the row data [49]. A real-world use case example was provided to further demonstrate how well the clustering technique solved the product suggestion issue [19].

To extract additional relevant elements from unannotated photos, the Data Annotations approach will be improved in the future [55]. Moreover, [30], the prolonged use of the product pricing and sales history may greatly enhance the model-creation process, resulting in more logical and personalised recommendations for the designers [33].

Additional steps might be taken to improve the usability and functionality of the user interface, which designers will use to create dashboards, search products, save system-recommended products, and enter their own preferences [29]. Last but not least, even though fashion designers will test and assess the recommended goods in more authentic use case scenarios, a more extensive set of trials using new datasets and methodology is necessary.

REFERENCES

- [1] Amed, I., Berg, A., Balchandani, A., Hedrich, S., Rolkens, F., Young, R., & Ekelof, J. (2021). The state of fashion 2021. Business Fashion and McKinsey and Company.
- [2] Gu, X., Gao, F., Tan, M., & Peng, P. (2020). Fashion analysis and understanding with artificial intelligence. *Information Processing & Management*, 57(5), 102276.
- [3] Bouwmans T, Javed S, Sultana M, Jung SK (2019) Deep neural network concepts for background subtraction: a systematic review and comparative evaluation. *Neural Netw* 117:8–66.
- [4] Bouwmans T, Silva C, Marghes C, Zitouni MS, Bhaskar H, Frelicot C (2018) On the role and the importance of features for background modeling and foreground detection. *Computer Science Review* 28:26–91.
- [5] Cordts M, Omran M, Ramos S, Rehfeld T, Enzweiler M, Benenson R, Franke U, Roth S, Schiele B (2016) The cityscapes dataset for semantic urban scene understanding. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp 3213–3223.
- [6] Deng J, Dong W, Socher R, Li LJ, Li K, Fei-Fei L (2009) Imagenet: a large-scale hierarchical image database. In: *2009 IEEE Conference on computer vision and pattern recognition*, pp 248–255.
- [7] Benatia, M.A., Baudry, D., Louis, A., 2022. Detecting counterfeit products by means of frequent pattern mining. *J. Ambient Intell. Humaniz. Comput.* 13, 3683–3692.
- [8] Bruce, M., 2009. Combating counterfeiting and piracy: The commonwealth initiative. *Commonw. Law Bull.* 35, 703–711.
- [9] Calderoni, F., Favarin, S., Garofalo, L., Sarno, F., 2014. Counterfeiting, illegal firearms, gambling and waste management: an exploratory estimation of four criminal markets. *Glob. Crime* 15, 108–137.
- [10] Coulter, K.S., Petty, R.D., 2012. Using the law to protect the brand on social media sites: A three “M”s framework for marketing managers. *Manag. Res. Rev.* 35, 758–769.
- [11] eCommerce report 2021 [WWW Document], n.d. . Statista.
- [12] Gai, K., Qiu, M., Zhao, H., Dai, W., 2016. Anti-Counterfeit Scheme Using Monte Carlo Simulation for Ecommerce in Cloud Systems. *Proc. - 2nd IEEE Int. Conf. Cyber Secur. Cloud Comput. CSCloud 2015 - IEEE Int. Symp. Smart Cloud IEEE SSC 2015* 74–79.
- [13] Edward F. Sherman, Amending Complaints to Sue Previously Misnamed or Unidentified Defendants After the Statute of Limitations Has Run: Questions Remaining from the Krupski Decision, 15 NEV. L.J. 1329, 1345 (2015)
- [14] Sophia Guan, 93% of Consumer Engagement with Luxury Brands Happens on Instagram, DIGIMIND (Nov. 30, 2018).
- [15] James Ray, Trademark Enforcement: A More Nuanced Game Than Whack-a-Mole, IPWATCHDOG (Oct. 23, 2018),
- [16] Frederick Mostert, Study on Approaches to Online Trademark Infringement, WORLD INTELL. PROP. ORG. 7 (Sept. 1, 2017).
- [17] Christine Lee, Wise Up: The Big Mistakes Luxury Brands Are Making with China’s Gen Z, JING DAILY (May 28, 2018).

- [18] F. Yang, A. Kale, Y. Bubnov, L. Stein, Q. Wang, H. Kiapour, and R. Piramuthu, "Visual search at ebay," in Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ser. KDD '17. New York, NY, USA: ACM, 2017, pp. 2101-2110.
- [19] H. Hu, Y. Wang, L. Yang, P. Komlev, L. Huang, X. S. Chen, J. Huang, Y. Wu, M. Merchant, and A. Sacheti, "Web-scale responsive visual search at bing," in Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ser. KDD '18. New York, NY, USA: ACM, 2018, pp. 359-367.
- [20] Y. Jing, D. Liu, D. Kislyuk, A. Zhai, J. Xu, J. Donahue, and S. Tavel, "Visual search at pinterest," in Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ser. KDD '15. New York, NY, USA: ACM, 2015, pp. 1889-1898.
- [21] Wang, W.; Xu, Y.; Shen, J.; Zhu, S.C. Attentive fashion grammar network for fashion landmark detection and clothing category classification. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, 18–22 June 2018; pp. 4271–4280.
- [22] Zhang, S.; Song, Z.; Cao, X.; Zhang, H.; Zhou, J. Task-aware attention model for clothing attribute prediction. IEEE Trans. Circuits Syst. Video Technol. 2019, 30, 1051–1064.
- [23] Zhang, Y.; Zhang, P.; Yuan, C.; Wang, Z. Texture and shape biased two-stream networks for clothing classification and attribute recognition. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, Seattle, WA, USA, 13–19 June 2020; pp. 13538–13547.
- [24] An, L.; Li, W. An integrated approach to fashion flat sketches classification. Int. J. Cloth. Sci. Technol. 2014, 26, 346–366.
- [25] Berg, T.L.; Ortiz, L.E.; Kiapour, M.H.; Yamaguchi, K. Parsing clothing in fashion photographs. In Proceedings of the IEEE Conference on Computer Vision & Pattern Recognition, Providence, RI, USA, 16–21 June 2012.
- [26] Yang, X.; Yuan, S.; Tian, Y. Assistive clothing pattern recognition for visually impaired people. IEEE Trans. Hum. Mach. Syst. 2014, 44, 234–243.
- [27] Chen, H.; Gallagher, A.; Girod, B. Describing clothing by semantic attributes. In European Conference on Computer Vision; Springer: Berlin/Heidelberg, Germany, 2012; pp. 609–623.
- [28] Green, R.T., Smith, T., 2002. Executive Insights: Countering Brand Counterfeiters. J. Int. Mark. 10, 89– 106.
- [29] Alrashed, T., Almahmoud, J., Zhang, A. X., Karger, D. R. (2020). ScrAPIr: Making Web Data APIs Accessible to End Users. Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems (pp. 1–12). New York, NY, USA: Association for Computing Machinery.
- [30] Asma. (2021, September 29). The Future of Web Scraping Services.
- [31] Banerjee, R. (2014). Website Scraping. Happiest Minds Technologies Pvt. Ltd.
- [32] Copeland, L., Ciampaglia, G. L., & Zhao, L. (2019). Fashion informatics and the network of fashion knockoffs. First Monday.
- [33] Beltzung, L., Lindley, A., Dinica, O., Hermann, N., & Lindner, R. (2020, December). Real-time detection of fake-shops through machine learning. In 2020 IEEE International Conference on Big Data (Big Data) (pp. 2254-2263). IEEE.
- [34] Bali, A., & Rehman, R. U. (2019). Detection of Price Scraping using Behavioral Analysis.
- [35] Arnold, P., Wartner, C., & Rahm, E. (2016). Semi-automatic identification of counterfeit offers in online shopping platforms. Journal of Internet commerce, 15(1), 59-75.
- [36] Dadhich, A., & Thankachan, B. (2021). Social & juristic challenges of AI for opinion mining approaches on Amazon & Flipkart product reviews using machine learning algorithms. SN Computer Science, 2(3), 180.
- [37] Wan Abdul Rahman, W. F., Che Fauzi, A. A., Wan Husain, W. S., Che Hassan, S. H., Nik Kamaruzaman, N. N., & Wan Aziz, W. A. H. (2020). The Usage of artificial intelligence in marketing automation: potentials and pitfalls. Journal of Mathematics and Computing Science (JMCS), 6(2), 1-8.
- [38] Li, F., Kant, S., Araki, S., Banger, S., & Shukla, S. S. (2020). Neural networks for fashion image classification and visual search. arXiv preprint arXiv:2005.08170.
- [39] Chinta, S. (2021). Integrating Machine Learning Algorithms in Big Data Analytics: A Framework for Enhancing Predictive Insights.
- [40] Soldner, F., Kleinberg, B., & Johnson, S. (2022). Trends in online consumer fraud: A data science perspective. In A Fresh Look at Fraud (pp. 167-191). Routledge.

- [41] Dadhich, A., & Thankachan, B. (2021). Sentiment analysis of amazon product reviews using hybrid rule-based approach. In Smart Systems: Innovations in Computing: Proceedings of SSIC 2021 (pp. 173-193). Singapore: Springer Singapore.
- [42] Ujwal, B. V. S., Gaiind, B., Kundu, A., Holla, A., & Rungta, M. (2017, December). Classification-based adaptive web scraper. In 2017 16th IEEE International Conference on Machine Learning and Applications (ICMLA) (pp. 125-132). IEEE.
- [43] Mabrouk, A., Redondo, R. P. D., & Kayed, M. (2021). Seopinion: summarization and exploration of opinion from e-commerce websites. *Sensors*, 21(2), 636.
- [44] Stuenkel, B. (2021). Personal information and artificial intelligence: Website scraping and the California Consumer Privacy Act. *Colo. Tech. LJ*, 19, 429.
- [45] M. Brenner, "Artificial Intelligence Top Benefits & Uses of AI in Digital Marketing", Marketing Insider Group, 2020.
- [46] K. Jarek, A. L. Kozminskiego, and G. Mazurek, "Marketing and Artificial Intelligence," *Cent. Eur. Bus. Rev.*, vol. 8, no. June, pp. 46–55, 2019.
- [47] M. Al-Rifaie, M. Bishop, "Weak vs. Strong Computational Creativity", ISB/IACAP World Congress 2012 - 5th AISB Symposium on Computing and Philosophy: Computing, Philosophy and the Question of Bio-Machine Hybrids, Part of Alan Turing, pp.61 – 67, 2012.
- [48] Bhatia, M. A. (2016). Artificial Intelligence–Making an Intelligent personal assistant. *Indian J. Comput. Sci. Eng.*, 6, 208-214.
- [49] Broucke, S. V., Baesens, B. (2018). Practical Web Scraping for Data Science: Best Practices and Examples with Python. (1st, Ed.) Apress.
- [50] Chamoso, P., Bartolomé, Á., García-Retuerta, D., Prieto, J., De La Prieta, F. (2020). Profile generation system using artificial intelligence for information recovery and analysis. *Journal of Ambient Intelligence and Humanized Computing*, 11(11), 4583-4592.
- [51] Strudel R, Garcia R, Laptev I, Schmid C (2021) Segmenter: transformer for semantic segmentation. In: Proceedings of the IEEE/CVF international conference on computer vision, pp 7262–7272 33.
- [52] Szegedy C, Liu W, Jia Y, Sermanet P, Reed S, Anguelov D, Erhan D, Vanhoucke V, Rabinovich A (2015) Going deeper with convolutions. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1–9 34.
- [53] Van Gansbeke W, Vandenhende S, Georgoulis S, Van Gool L (2021) Unsupervised semantic segmentation by contrasting object mask proposals. In: Proceedings of the IEEE/CVF international conference on computer vision, pp. 10052–10062 35.
- [54] Wang W, Zhou T, Yu F, Dai J, Konukoglu E, Van Gool L (2021) Exploring cross-image pixel contrast for semantic segmentation. In: Proceedings of the IEEE/CVF international conference on computer vision, pp 7303–7313.
- [55] J. Flowers, "Strong and Weak AI: Deweyan Considerations". In: AAAI Spring Symposium: Towards Conscious AI Systems, 2018.
- [56] S. Schneider, "Superintelligent AI and the Post-biological Cosmos Approach", pp. 178 – 198, 2017.
- [57] K.Nurm, "The Possibilities and Potential Risks of Using Artificial Intelligence In Marketing – A Literature Review", pp. 1 – 69, 2020.
- [58] C. Kumar, GN, "Artificial Intelligence: Definition, Types, Examples, Technologies", 2019. [11] A. Mathematics, "Artificial Intelligence - The Marketing Game Changer," vol. 119, no. 17, pp. 1881–1890, 2018.
- [59] Hidayati, S.C.; You, C.W.; Cheng, W.H.; Hua, K.L. Learning and recognition of clothing genres from full-body images. *IEEE Trans. Cybern.* 2017, 48, 1647–1659.