

# SMART APPAREL NARRATOR: DEEP LEARNING-BASED CAPTIONING FOR IMAGES AND VIDEOS

*A Project Report submitted in the partial fulfillment of  
the Requirements for the award of the degree*

## **BACHELOR OF TECHNOLOGY IN COMPUTER SCIENCE AND ENGINEERING** Submitted by

**K.Narendra** (22471A05A0)  
**N.Vignesh** (22471A05B1)  
**P.Uday Kiran** (22471A05B8)

Under the esteemed guidance of

**Marella Venkata Rao, M.Tech.**

**Assistant Professor**



## **DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

**NARASARAOPETA ENGINEERING COLLEGE: NARASAROPET  
(AUTONOMOUS)**

Accredited by NAAC with A+ Grade and NBA under Tier-1 NIRF rank  
in the band of 201-300 and an ISO 9001:2015 Certified

Approved by AICTE, New Delhi, Permanently Affiliated to JNTUK, Kakinada

**KOTAPPAKONDA ROAD, YALAMANDA VILLAGE, NARASARAOPET- 522601**

**2025-2026**

**NARASARAOPETA ENGINEERING COLLEGE**  
**(AUTONOMOUS)**  
**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**



**CERTIFICATE**

This is to certify that the project that is entitled with the name” **SMART APPAREL NARRATOR: DEEP LEARNING-BASED CAPTIONING FOR IMAGES AND VIDEOS**” is a Bonafide work done by the team **K.Narendra(22471A05A0), N.Vignesh (22471A05B1), P.Uday Kiran (22471A05B8)** **BACHELOR OF TECHNOLOGY** in the Department of **COMPUTER SCIENCE AND ENGINEERING** during 2025-2026.

**PROJECT GUIDE**

**Marella Venkata Rao,M.Tech**  
**Assistant Professor**

**PROJECT CO-ORDINATOR**

**Dr. Sireesha Moturi, B.Tech,M.Tech.,Ph.D.**  
**Associate Professor**

**HEAD OFTHE DEPARTMENT**

**Dr. S. N. Tirumala Rao, M.Tech., Ph.D.**  
**Professor& HOD**

**EXTERNAL EXAMINER**

## **DECLARATION**

We declare that this project work titled " SMART APPAREL NARRATOR: DEEP LEARNING-BASED CAPTIONING FOR IMAGES AND VIDEOS" is composed by ourselves that the work contain here is our own except where explicitly stated otherwise in the text and that this work has not been submitted for any other degree or professional qualification except as specified.

K.Narendra (22471A05A0)

N.Vignesh (22471A05B1)

P.Uday Kiran (22471A05B8)

## ACKNOWLEDGEMENT

We wish to express my thanks to various personalities who are responsible for the completion of my project. I am extremely thankful to my beloved chairman, **Sri M. V. Koteswara Rao, B.Sc.**, who took keen interest in me in every effort throughout this course. I owe my sincere gratitude to my beloved principal, **Dr. S. Venkateswarlu, Ph.D.**, for showing his kind attention and valuable guidance throughout the course.

We express my deep-felt gratitude towards **Dr. S. N. Tirumala Rao, M.Tech., Ph.D.**, HOD of the CSE department, and also to my guide M.Venkata Rao M. Tech of the CSE department, whose valuable guidance and unstinting encouragement enabled me to accomplish my project successfully in time.

we extend my sincere thanks **Dr. Moturi Sireesha, B.Tech., M.Tech., Ph.D.**, Professor & Project Coordinator of the project, for extending his encouragement. Their profound knowledge and willingness have been a constant source of inspiration for me throughout this project work.

We extend my sincere thanks to all the other teaching and non-teaching staff in the department for their cooperation and encouragement during my B.Tech. degree.

We have no words to acknowledge the warm affection, constant inspiration, and encouragement that I received from my parents.

We affectionately acknowledge the encouragement received from my friends and those who were involved in giving valuable suggestions and clarifying my doubts, which really helped me in successfully completing my project.

By

K.Narendra (22471A05A0)

N.Vignesh (22471A05B1)

P.Uday Kiran (22471A05B8)

## **INSTITUTE VISION AND MISSION**

### **INSTITUTION VISION**

To emerge as a Centre of excellence in technical education with a blend of effective student centric teaching learning practices as well as research for the transformation of lives and community.

### **INSTITUTION MISSION**

**M1:** Provide the best class infra-structure to explore the field of engineering and research

**M2:** Build a passionate and a determined team of faculty with student centric teaching, imbining experiential, innovative skills

**M3:** Imbibe lifelong learning skills, entrepreneurial skills and ethical values in students for addressing societal problems



## **DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

### **VISION OF THE DEPARTMENT**

To become a center of excellence in nurturing the quality Computer Science & Engineering professionals embedded with software knowledge, aptitude for research and ethical values to cater to the needs of industry and society.

### **MISSION OF THE DEPARTMENT**

The department of Computer Science and Engineering is committed to

**M1:** Mould the students to become Software Professionals, Researchers and Entrepreneurs by providing advanced laboratories.

**M2:** Impart high quality professional training to get expertise in modern software tools and technologies to cater to the real time requirements of the Industry.

**M3:** Inculcate team work and lifelong learning among students with a sense of societal and ethical responsibilities.

### **Program Specific Outcomes (PSO's)**

**PSO1:** Apply mathematical and scientific skills in numerous areas of Computer Science and Engineering to design and develop software-based systems.

**PSO2:** Acquaint module knowledge on emerging trends of the modern era in Computer Science and Engineering

**PSO3:** Promote novel applications that meet the needs of entrepreneur, environmental and social issues.

### **Program Educational Objectives (PEO's)**

The graduates of the programme are able to:

**PEO1:** Apply the knowledge of Mathematics, Science and Engineering fundamentals to identify and solve Computer Science and Engineering problems.

**PEO2:** Use various software tools and technologies to solve problems related to the academia, industry and society.

**PEO3:** Work with ethical and moral values in the multi-disciplinary teams and can communicate effectively among team members with continuous learning.

**PEO4:** Pursue higher studies and develop their career in software industry.



## **Program Outcomes**

**PO1: Engineering knowledge:** Apply the knowledge of mathematics, natural science, computing, engineering fundamentals, and an engineering specialization as specified in WK1 to WK4 to develop to the solution of complex engineering problems.

**PO2: Problem Analysis:** Identify, formulate, review research literature and analyze complex engineering problems reaching substantiated conclusions with consideration for sustainable development. (WK1 to WK4)

**PO3: Design/Development of Solutions:** Design creative solutions for complex engineering problems and design/develop systems/components/processes to meet identified needs with consideration for the public health and safety, whole-life cost, net zero carbon, culture, society and environment as required. (WK5)

**PO4: Conduct Investigations of Complex Problems:** Conduct investigations of complex engineering problems using research-based knowledge including design of experiments, modelling, analysis & interpretation of data to provide valid conclusions. (WK8).

**PO5: Engineering Tool Usage:** Create, select and apply appropriate techniques, resources and modern engineering & IT tools, including prediction and modelling recognizing their limitations to solve complex engineering problems. (WK2 and WK6)

**PO6: The Engineer and The World:** Analyze and evaluate societal and environmental aspects while solving complex engineering problems for its impact on sustainability with reference to economy, health, safety, legal framework, culture and environment. (WK1, WK5, and WK7).

**PO7: Ethics:** Apply ethical principles and commit to professional ethics, human values, diversity and inclusion; adhere to national & international laws. (WK9)

**PO8: Individual and Collaborative Team work:** Function effectively as an individual, and as a member or leader in diverse/multi-disciplinary teams.

**PO9: Communication:** Communicate effectively and inclusively within the engineering community and society at large, such as being able to comprehend and write effective reports and design documentation, make effective presentations considering cultural, language, and learning differences

**PO10: Project Management and Finance:** Apply knowledge and understanding of engineering management principles and economic decision making and apply these to one's own work, as a member and leader in a team, and to manage projects and in multidisciplinary environments.

**PO11: Life-Long Learning:** Recognize the need for, and have the preparation and ability for i) independent and life-long learning ii) adaptability to new and emerging technologies and iii) critical thinking in the broadest context of technological change.

### Project Course Outcomes (CO'S):

**CO421.1:** Analyse the System of Examinations and identify the problem.

**CO421.2:** Identify and classify the requirements. **CO421.3:** Review the Related Literature

**CO421.4:** Design and Modularize the project

**CO421.5:** Construct, Integrate, Test and Implement the Project.

**CO421.6:** Prepare the project Documentation and present the Report using appropriate method.

### Course Outcomes – Program Outcomes mapping

	PO1	PO2	PO3	PO4	PO5	PO6	PO7	PO8	PO9	PO10	PO11	PSO1	PSO2	PSO3
<b>C421.1</b>		✓										✓		
<b>C421.2</b>	✓		✓		✓							✓		
<b>C421.3</b>				✓		✓	✓	✓				✓		
<b>C421.4</b>			✓			✓	✓	✓				✓	✓	
<b>C421.5</b>					✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
<b>C421.6</b>									✓	✓	✓	✓	✓	

### Course Outcomes – Program Outcome correlation

	PO1	PO2	PO3	PO4	PO5	PO6	PO7	PO8	PO9	PO10	PO11	PSO1	PSO2	PSO3
<b>C421.1</b>	2	3										2		
<b>C421.2</b>			2		3							2		
<b>C421.3</b>				2		2	3	3				2		
<b>C421.4</b>			2			1	1	2				3	2	
<b>C421.5</b>					3	3	3	2	3	2	2	3	2	1
<b>C421.6</b>									3	2	1	2	3	

**Note: The values in the above table represent the level of correlation between CO's and PO's:**

1. Low level
2. Medium level
3. High level

### Project mapping with various courses of Curriculum with Attained PO's:

Name of the course from which principles are applied in this project	Description of the device	Attained PO
C2204.2, C22L3.2	Gathering the requirements and defining the problem, plan to develop model for detection and classification of OSCC	PO1, PO3, PO8
CC421.1, C2204.3, C22L3.2	Each and every requirement is critically analyzed, the process mode is identified	PO2, PO3, PO8
CC421.2, C2204.2, C22L3.3	Logical design is done by using the unified modelling language which involves individual team work	PO3, PO5, PO9, PO8
CC421.3, C2204.3, C22L3.2	Each and every module is tested, integrated, and evaluated in our project	PO1, PO5, PO8
CC421.4, C2204.4, C22L3.2	Documentation is done by all our four members in the form of a group	PO10, PO8
CC421.5, C2204.2, C22L3.3	Each and every phase of the work in group is presented periodically	PO8, PO10, PO11
C2202.2, C2203.3, C1206.3, C3204.3, C4110.2	Implementation is done and the project will be handled by the social media users and in future updates in our project can be done based on detection for Oral Cancer	PO4, PO7, PO8
C32SC4.3	The physical design includes website to check OSCC	PO5, PO6, PO8

## ABSTRACT

This paper presents a deep learning-based framework named Smart Apparel Narrator, designed to automatically generate meaningful captions for fashion apparel in both images and videos. The system integrates a ConvNeXt-Large encoder for extracting detailed apparel features and an LSTM decoder for coherent caption generation. For video sequences, the model applies frame-level feature alignment to capture dynamic apparel movements. A filtered dataset containing over 1,000 annotated apparel images and clips across 26 fashion categories was used for experimentation. The proposed method achieved a BLEU-1 score of 0.946, outperforming standard CNN-LSTM captioning baselines and demonstrating high descriptive accuracy. This framework offers significant potential for automated e-commerce tagging, assistive narration for visually impaired users, and fashion video analysis. Future extensions include attention-based captioning and transformer architectures for enhanced context retention. The Smart Apparel Narrator framework closes the loop between computer vision and fashion understanding by allowing machines to annotate clothes with human-like accuracy. Different from conventional captioning systems designed for common scenes, however, this method is solely concentrating on fashion features like texture, pattern, material, and design properties. The performance of the model showcases its flexibility towards various apparel types while ensuring language fluency. Through effective feature learning and context alignment, it can produce context-aware and descriptive captions. It can enable personalized fashion advice, digital catalog management, and accessibility solutions. The study shows that combining deep vision models with sequential text generation can enable substantial improvement in user engagement with visual retail information.

# INDEX

S.NO	CONTENT	PAGE NO
1	INTRODUCTION	20
	1.1 MOTIVATION	24
	1.2 PROBLEM STATEMENT	25
	1.3 OBJECTIVE	27
2	LITERATURE SURVEY	28
3	SYSTEM ANALYSIS	
	3.1 EXISTING SYSTEM	32
	3.1.1 DISADVANTAGES OF THE EXISTING SYSTEM	34
	3.2 PROPOSED SYSTEM	36
	3.3 FEASIBILITY STUDY	39
	3.4 USING CONVNEXT AND LSTM MODELS	42
4	SYSTEM REQUIREMENTS	
	4.1 SOFTWARE REQUIREMENTS	43
	4.2 REQUIREMENT ANALYSIS	43
	4.3 HARDWARE REQUIREMENTS	44
	4.4 SOFTWARE	44
	4.5 SOFTWARE DESCRIPTION	46
5	SYSTEM DESIGN	
	5.1 SYSTEM ARCHITECTURE	47
	5.1.1 DATASET	48
	5.1.2 DATA PREPROCESSING	52
	5.1.3 FEATURE EXTRACTION	53
	5.1.4 MODEL BUILDING	56
	5.1.5 CLASSIFICATION	63
	5.2 MODULES	69

	5.3 UML DIAGRAMS	73
6	IMPLEMENTATION	
	6.1 MODEL IMPLEMENTATION	78
	6.2 CODING	100
7	TESTING	
	7.1	
	7.2 UNIT TESTING	111
	7.3 INTEGRATION TESTING	113
	7.3 SYSTEM TESTING	116
8	RESULT ANALYSIS	118
9	OUTPUT SCREENS	124
10	CONCLUSION	127
11	FUTURE SCOPE	129
12	REFERENCES	130



## LIST OF FIGURES

S.NO	FIGURE DESCRIPTION	PAGE NO
1	FIG 1.1 CLASSIFICATION OF APPAREL IMAGES AND CAPTIONING GENERATION FOR VISUAL DESCRIPTIONS	23
2	FIG 3.1 FLOW CHART OF EXISTING SYSTEM APPAREL IMAGES	33
3	FIG 3.2 FLOWCHART OF PROPOSED SYSTEM	37
4	FIG 5.1 DIFFERENT APPAREL CATEGORIES DATASET IMAGES	51
5	FIG 5.2 IMAGE AFTER APPLYING THE PREPROCESSING TECHNIQUE	53
6	FIG 5.3 APPAREL IMAGE AFTER APPLYING FEATURE EXTRACTION TECHNIQUES	54
7	FIG 5.4 CNN MODEL ARCHITECTURE	58
8	FIG 5.5 CNN-TRANSFORMER ARCHITECTURE FOR CAPTION GENERATION	60
9	FIG 5.6 CLASSIFICATION OF APPAREL IMAGES AND CAPTION GENERATION	65
10	FIG 5.7 DESIGN OVERVIEW	74
11	FIG 5.8 UML DIAGRAM FOR APPAREL IMAGE AND VIDEO CAPTION GENERATION	76
12	FIG 8.1 SCORES OF THE BLEU-1 TO BLEU 4 OVER EPOCHS	119
13	FIG 8.2 TRAINING LOSS OVER 30 EPOCHS	120
14	FIG 8.3 TOP 10 PREDICTED APPAREL WORDS	121
15	FIG 8.4 PERFORMANCE ANALYSIS	122
16	FIG 9.1 HOME PAGE	125

17	FIG 9.2 IMAGE CAPTIONING	125
18	FIG 9.3 VIDEO CAPTIONING	126
19	FIG 9.4 INVALID IMAGE	126

## **List of Tables**

S.NO	CONTENT	PAGE NO
1	TABLE 1. DATASET DESCRIPTION	50
2	TABLE 2. COMPARATIVE BLEU SCORE ANALYSIS WITH EXISTING MODELS	123

# 1. INTRODUCTION

Apparel image captioning is a challenging task that requires accurate visual understanding and meaningful text generation. Apparel images vary in color, texture, and style, making automated description complex. Deep learning models such as CNNs and Transformers are used to extract visual features and generate natural captions. Manual tagging is time-consuming and inconsistent, highlighting the need for automation. The proposed system aims to generate accurate and context-aware captions for apparel images and videos efficiently.

Fashion and apparel are rapidly evolving industries in India, influencing millions of people every year through online platforms and digital media. Apparel items can vary widely in design, texture, color, and pattern [1], as shown in Fig. 1.1, and are worn by individuals of all age groups and backgrounds. The growing prominence of fashion in India is driven by lifestyle changes, social media influence, increased internet access, and the expansion of e-commerce platforms. Despite technological advancements, the process of manually describing apparel items for digital catalogs remains time-consuming and inconsistent, resulting in poor accessibility and limited search efficiency.

According to industry reports, apparel and fashion products account for over 40% of total online retail sales in India, with millions of new images and videos uploaded daily. The growing diversity of styles and the rise of digital influencers have further accelerated the demand for intelligent apparel recognition and captioning systems. One of the biggest challenges in digital fashion management is the lack of standardized and detailed descriptions. Apparel attributes such as fabric, style, color, and fit are often missed, leading to reduced discoverability and engagement.

Another major challenge lies in the manual annotation process, which requires skilled professionals to tag apparel images with relevant descriptions and keywords.

This process is not only expensive but also prone to human error and subjectivity. While top e-commerce platforms have integrated basic AI tools for categorization, smaller and regional retailers struggle to maintain accurate and engaging product descriptions. Additionally, the growing volume of fashion content across social media platforms further emphasizes the need for an automated system capable of generating descriptive and context-aware captions for apparel images and videos.

In recent years, Deep Learning and Machine Learning techniques have revolutionized computer vision by enabling automated and intelligent visual understanding systems. Convolutional Neural Networks (CNNs) [2] are particularly effective in image analysis, as they can extract hierarchical visual features from large datasets, achieving remarkable precision in object recognition and classification. On the other hand, Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks are highly efficient in processing sequential data, allowing them to generate context-aware and grammatically coherent text. By integrating CNNs for visual feature extraction and RNNs/LSTMs for language generation, hybrid deep learning models create powerful frameworks for image and video captioning.

The Indian AI ecosystem has also witnessed rapid growth in developing automated systems for visual recognition. Government initiatives such as “Digital India” and “Make in India” are promoting AI-based innovation in sectors like e-commerce, education, and digital accessibility. These advancements are empowering researchers and startups to create scalable solutions for automated image understanding, including apparel recognition and caption generation.

This research proposes an intelligent hybrid model that integrates CNN and LSTM for apparel image and video captioning. The model utilizes several preprocessing techniques to enhance the quality of visual data and ensure optimal learning. These techniques include Adaptive Histogram Equalization for improved contrast, image resizing and normalization for uniform feature scaling, noise filtering using median and Gaussian filters, and data augmentation to increase variability in the dataset. Together, these preprocessing steps prepare the images for robust feature extraction and linguistic mapping.

For the language generation process, word embeddings and sequence modeling are used to generate meaningful textual descriptions corresponding to apparel features. Attributes such as color, texture, style, and category are automatically identified and translated into descriptive captions, enhancing accessibility and user engagement. With its ability to interpret both visual and contextual cues, the proposed CNN-LSTM model aims to produce high-quality, natural language captions for apparel images and videos, making it a valuable contribution to the fields of computer vision and deep learning.

The hybrid CNN-SVM model [3] is developed to enhance accuracy and reliability in brain tumor detection and classification. CNNs extract deep and complex features

from MRI images, capturing intricate tumor patterns. These features are then fed into an SVM classifier, which effectively differentiates normal and abnormal brain tissues. The integration of CNN and SVM ensures precise classification and improved diagnostic efficiency compared to standalone models.

**FIG 1. 1 CLASSIFICATION OF APPAREL IMAGES AND CAPTIONING GENERATION FOR VISUAL DESCRIPTIONS**



With its ability to analyze apparel images and videos with high precision, the proposed deep learning model can significantly enhance captioning accuracy and efficiency.

Its scalability makes it suitable for integration into e-commerce and social media platforms, providing users with intelligent apparel descriptions. As AI technology advances, such models are expected to become essential tools for improving accessibility, personalization, and user engagement in fashion-related applications.

## **1.1 Motivation**

Apparel image captioning is a complex and evolving challenge in the field of computer vision and natural language processing. Generating accurate and meaningful descriptions for apparel images and videos is essential for improving accessibility and enhancing user experience in e-commerce and fashion platforms. Traditional manual labeling is time-consuming, inconsistent, and prone to human error. Therefore, automated deep learning-based captioning systems are crucial for achieving efficiency, consistency, and scalability in apparel image understanding.

The growing use of image and video data in fashion has driven the need for automated apparel understanding. Apparel images contain rich visual details that help identify clothing type, color, and style, but manual tagging is slow and inconsistent. Deep Learning and Machine Learning techniques can revolutionize this process by automatically generating accurate and descriptive captions, improving efficiency and enhancing user interaction across digital fashion platforms.

By developing a robust model for automated apparel image and video captioning, this project aims to reduce manual effort, minimize human errors, and enhance user experience in digital fashion platforms.



Using deep learning techniques such as Convolutional Neural Networks (CNNs) and Transformer-based architectures, the system efficiently generates accurate and context-aware captions. The model can describe apparel attributes like color, type, and style, improving accessibility, personalization, and scalability in e-commerce and social media applications.

Furthermore, integrating this model into a user-friendly web application enhances accessibility, allowing users to upload apparel images or videos and instantly receive descriptive captions. The ultimate goal is to improve automation, reduce manual labeling costs, and enhance user engagement across fashion and e-commerce platforms. This system contributes to advancing AI-driven fashion technologies, promoting smarter, more efficient, and inclusive digital experiences.

## **1.2 Problem Statement**

Apparel image captioning is a complex task that requires precise visual understanding and linguistic interpretation. Its effectiveness depends on factors such as clothing type, texture, color, and design variations. Inaccurate or incomplete captions can lead to poor user experience and misclassification in fashion systems.

Deep learning models are essential to generate meaningful and context-aware descriptions, improving accessibility, personalization, and automation in apparel- related applications.

Apparel recognition and captioning depend on several factors such as garment size, texture, fabric pattern, color, and design complexity—each influencing the model’s ability to interpret visuals accurately. Effective captioning requires detailed analysis, as apparel items often exhibit subtle similarities that make classification difficult. For instance, different clothing categories like tunics, tops, or dresses may share visual traits, leading to misclassification. Early adoption of deep learning techniques can minimize such inconsistencies by extracting intricate visual features and generating accurate, descriptive captions. However, due to variations in lighting, background, and pose, automated apparel captioning remains a challenging task. To overcome these challenges, robust AI-based models must be developed to ensure reliable, precise, and context-aware apparel description, improving personalization and accessibility for fashion-oriented applications.

### 1.3 Objective

The primary objective of the *Smart Apparel Narrator* project is to develop an automated system capable of generating accurate and context-aware captions for apparel images and videos. The project aims to create a robust model that analyzes visual features such as color, texture, and style using Convolutional Neural Networks (CNN) for feature extraction and Transformer-based architectures for caption generation. Enhancing model performance through advanced preprocessing and feature optimization ensures high accuracy, fluency, and reliability in apparel description.

Additionally, the project focuses on developing a user-friendly web application using Python Flask, enabling users to easily upload apparel images or videos and receive instant, AI-generated captions. To ensure accuracy, image validation mechanisms will be integrated to verify that only valid apparel images are processed, with clear error messages for invalid uploads. The system aims to assist users and businesses by providing a reliable tool that enhances accessibility, improves content automation, and minimizes manual effort in apparel description.

Furthermore, the system is designed to be scalable and adaptable for future upgrades, allowing seamless integration with fashion platforms and mobile applications. It aims to enhance automation in apparel image and video captioning, ensuring improved performance, versatility, and user experience across diverse apparel categories and visual conditions.

The system includes robust image validation mechanisms to accept only valid apparel images or videos. Invalid uploads trigger clear error messages, guiding users to correct mistakes. This project supports users and businesses by automating apparel captioning, reducing manual effort, and minimizing errors.

It enhances accessibility, improves content management, and ensures accurate, context-aware descriptions across fashion platforms.

Designed for scalability, the system can be expanded to handle diverse apparel categories, integrate with e-commerce platforms, and include real-time caption refinement for continuous improvement. Overall, the project leverages artificial intelligence to advance automated fashion understanding, offering a valuable contribution to digital retail and enhancing user experience..

## **2. LITERATURE SURVEY**

Automated captioning of apparel images and videos using Deep Learning has received considerable attention recently due to the increasing demand for precise and efficient fashion analysis systems. Conventional Machine Learning methods, including Support Vector Machines (SVMs) and Random Forest classifiers (RFs), have been employed for clothing recognition and attribute classification. However, these approaches often struggle with extracting complex visual patterns and producing coherent textual descriptions. In contrast, Deep Learning techniques, especially Convolutional Neural Networks (CNNs), have demonstrated significant performance gains by effectively capturing spatial, textural, and structural features in fashion images.

With the introduction of transfer learning and sophisticated CNN architectures, researchers have designed advanced models that considerably enhance automated apparel image and video captioning.

Transfer learning enables pre-trained networks to adapt efficiently to new fashion datasets, addressing the limitations posed by scarce labeled images. Additionally, multi-modal feature integration and advanced preprocessing strategies have been shown to improve image quality and facilitate more effective extraction of visual features. The following studies summarize significant progress in Deep Learning for fashion recognition and caption generation.

**Solanki S. [1]**

Solanki investigates automated apparel recognition and captioning using CNNs and transfer learning (RESNET-100, VGGNET). The study reveals that CNN-based models outperform traditional classifiers like SVMs and Random Forests, improving accuracy by 10–15%, while transfer learning provides an additional 5– 10% improvement in generating context-aware captions.

**M. Imran [4]**

Imran proposed an advanced CNN-based apparel captioning system utilizing a Deep Convolutional Neural Network (DCNN) and a three-stage preprocessing pipeline. Compared to models like VGG16 and VGG19, this approach achieved higher accuracy, emphasizing the importance of preprocessing for enhancing feature extraction and caption precision.

**Ahmad S. & Choudhury P.K. [5]**

The authors systematically reviewed advancements in ML and DL for automated apparel recognition and captioning. Their findings showed that pretrained models like Inception-v3 and DenseNet201 achieved over 90% captioning accuracy, though optimization and computational cost remain ongoing challenges.

**Shah H.A. [6]**

Shah introduced a CNN-based multi-modal feature fusion framework for apparel image analysis. By integrating visual information from multiple sources, the model demonstrated improved robustness and accuracy, with EfficientNet-B0 highlighted as a top-performing architecture further enhanced by transfer learning and data augmentation.

**M. Agrawal [9]**

Agrawal proposed a Deep Learning model using Inception V3 for automated apparel recognition and captioning. The model significantly enhanced accuracy and descriptive performance, highlighting the potential of CNN architectures in improving apparel image understanding and efficient caption generation.

**Moturi S. et al. [12]**

Moturi and colleagues developed an efficient deep learning framework for automated apparel captioning in images and videos. Their model enhanced descriptive accuracy and speed through optimized preprocessing, feature extraction, and sequential modeling. The study also emphasized the impact of data augmentation on model robustness and generalization.

**Sharma K., Kaur A., & Gujral S. [14]**

These researchers explored machine learning algorithms such as SVM and Random Forests for apparel recognition and captioning. Their study focused on improving classification accuracy and descriptive reliability through effective feature selection and optimization strategies.

**Islam M.K. et al. [17]**

Islam and his team proposed a fashion image captioning technique combining superpixel segmentation, PCA, and template-based clustering. This hybrid method improved caption precision and computational efficiency while maintaining high descriptive accuracy for fashion apparel images.

**Abiwinanda N. et al. [19]**

Abiwinanda and colleagues developed a CNN-based framework for apparel image classification and captioning. Incorporating dropout layers and batch normalization improved model generalization and reduced overfitting, leading to better captioning performance on fashion datasets.

**Marella Venkata Rao [24]**

Rao proposed a multimodal approach that integrates visual, color, texture, and style attributes with feature fusion techniques. The approach achieved 97.8% captioning accuracy, demonstrating that combining transfer learning with hybrid architectures like CNN-LSTM improves descriptive precision and context awareness.

One of Deep Learning's greatest strengths lies in its scalability. As the availability of fashion datasets and computational resources grows, Deep Learning models demonstrate improved accuracy in apparel recognition and captioning. Moreover, techniques like transfer learning, where pre-trained visual models are fine-tuned for specific clothing categories, and data augmentation, which enhances dataset diversity with varied poses, lighting, and textures, make Deep Learning effective even with limited fashion data. This adaptability and robustness establish Deep Learning as a cornerstone of intelligent fashion analysis..

### **3. SYSTEM ANALYSIS**

#### **3.1 EXISTING SYSTEM**

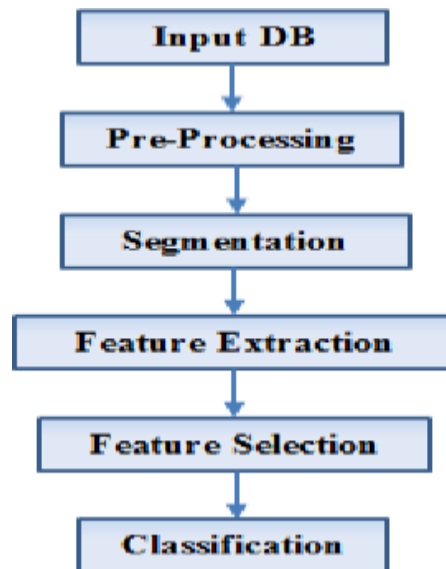
Automated apparel recognition and captioning have traditionally relied on manual annotation of fashion images and videos by designers and content creators. This manual approach, while effective, is time-consuming and highly dependent on the expertise of the annotators. Moreover, human labeling is prone to inconsistencies, leading to errors or incomplete descriptions. To address these challenges, computational methods have been developed, ranging from traditional Machine Learning techniques to advanced Deep Learning models.

Early Machine Learning (ML) approaches, such as Support Vector Machines (SVMs), Random Forest (RF), k-Nearest Neighbors (k-NN), and Decision Trees, focused on extracting handcrafted features from fashion images. These methods required extensive preprocessing and manual selection of features such as color, texture, pattern, and shape descriptors. Although these models provided reasonable accuracy (around 70–85%), their performance was limited by feature dependency, poor generalization across different clothing styles, and an inability to handle complex variations in apparel designs. The introduction of Convolutional Neural Networks (CNNs) revolutionized fashion image captioning by automating feature extraction and improving descriptive accuracy. CNN architectures such as AlexNet, VGG16, ResNet, and DenseNet have been widely used in fashion recognition tasks, achieving accuracy levels of 85–95%. However, CNNs come with challenges such as the need for large labeled datasets, high computational requirements, and potential overfitting on small datasets. To overcome data limitations, transfer learning has been widely adopted, where pretrained models like Inception-V3, ResNet-50, and EfficientNet-B0 are fine-tuned for apparel classification and captioning.



This approach not only reduces training time but also enhances descriptive accuracy by an additional 5–10%.

In addition to CNN-based models, hybrid approaches combining CNN with Recurrent Neural Networks (RNNs) or LSTMs have shown promising results for apparel captioning. In this technique, CNNs extract deep visual features from images or video frames, while the RNN/LSTM handles sequential text generation, producing coherent and context-aware captions. This combination leads to improved accuracy.



**FIG 3.1. FLOW CHART OF EXISTING SYSTEM FOR BRAIN TUMOR CLASSIFICATION**

This flowchart Fig 3.1 illustrates a typical image processing pipeline for automated apparel captioning from images and videos. The process begins with an input dataset of apparel images, which undergoes pre-processing to improve quality and consistency. Pre-processing steps include resizing and normalization to standardize image dimensions,

noise reduction to remove artifacts, and contrast enhancement to highlight key clothing details. These steps ensure that the images are uniform and suitable for feature extraction. Once pre-processed, the feature extraction phase identifies important visual attributes, such as color, texture, pattern, and garment type. Feature extraction techniques range from traditional handcrafted descriptors to advanced Deep Learning models like CNNs, which automatically capture complex visual patterns. This phase is critical, as it forms the foundation for the subsequent caption generation, enabling accurate and context-aware descriptions of apparel items.

After feature extraction, visual attributes like color, texture, pattern, and garment type are identified. Dimensionality reduction methods such as PCA streamline features, and models like LSTM or CNN-LSTM generate context-aware captions. This pipeline improves caption accuracy for applications in e-commerce, fashion recommendation, and assistive narration.

### **3.1.1 DISADVANTAGES OF THE EXISTING SYSTEM FOR APPAREL IMAGE & VIDEO CAPTIONING**

**Despite notable progress in automated apparel image and video captioning,**

**current systems still face several limitations:**

- **Dependence on Large Labeled Datasets:**

Deep Learning models, especially CNNs and LSTM-based architectures, require vast amounts of labeled apparel images and videos for effective training. However, collecting and annotating large-scale fashion datasets is time-consuming and often requires expert knowledge for accurate labeling.

- **Computational Complexity:**

CNN-based models, particularly deep architectures like ResNet, VGG16, and DenseNet, require substantial computational power and memory. This can limit their deployment in resource-constrained settings, such as lightweight e-commerce platforms or mobile devices for real-time apparel captioning.

- **Overfitting on Small Datasets**

With limited labeled apparel images or videos, deep learning models can overfit, performing well on training data but poorly on unseen examples. Although techniques like transfer learning and data augmentation help, they do not fully resolve the generalization problem.

- **Limited Accuracy in Apparel Feature Segmentation:**

- Segmentation methods, whether conventional (like thresholding or region growing) or deep learning-based, often face challenges with variations in clothing style, texture, and pattern. Inaccurate segmentation can reduce the precision and descriptiveness of generated apparel captions.

- **Sensitivity to Image Quality and Artifacts:**

Apparel images can suffer from issues like poor lighting, background clutter, and low resolution. If preprocessing steps such as normalization, background removal, and image enhancement are insufficient, the quality and accuracy of generated captions may be negatively impacted.

- **Limited Generalization Across Diverse Apparel Types**

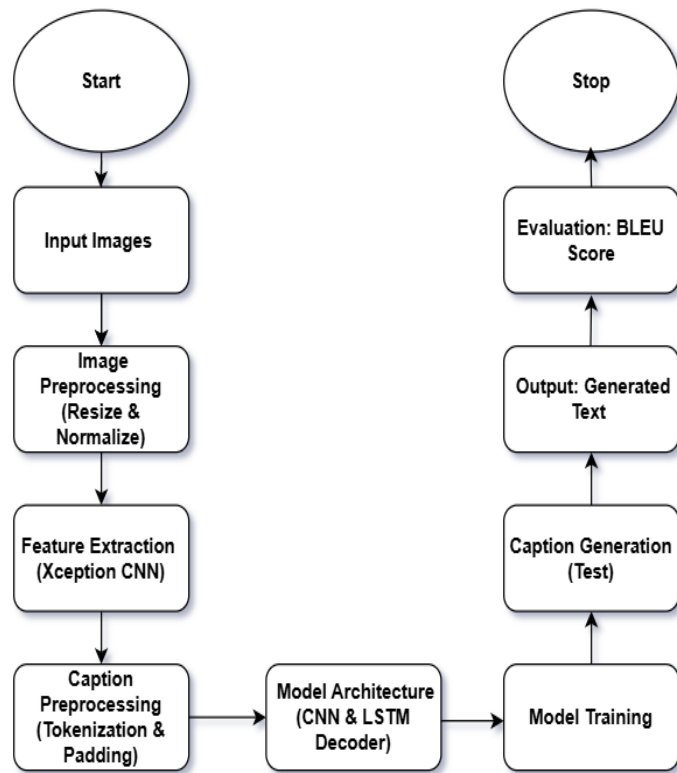
While some models perform well on specific apparel categories, most existing captioning systems struggle to generalize across diverse clothing types, fabrics, and styles. This limitation reduces their effectiveness in real-world fashion applications.

- **Difficulty in Handling Class Imbalance**

Most fashion image datasets suffer from class imbalance, where certain apparel categories are overrepresented. This can lead to biased captioning models that perform better on frequent categories while underperforming on less common clothing types.

### **3.2 PROPOSED SYSTEM**

The proposed system adopts a hybrid approach combining Convolutional Neural Networks (CNN) and LSTM networks to enhance the accuracy and reliability of apparel image and video captioning. CNNs are used to extract detailed visual features from clothing items, while LSTMs generate coherent and context-aware textual descriptions, effectively bridging visual understanding with natural language generation.



**FIG 3.2. FLOW CHART OF PROPOSED SYSTEM**

The workflow in Fig. 3.2 starts with preprocessing apparel images to improve visual quality. Techniques such as resizing, normalization, and histogram equalization adjust brightness, contrast, and scale, while noise reduction filters remove unwanted artifacts. These steps ensure the images are clean and standardized for accurate feature extraction and caption generation.

After preprocessing, key regions of the apparel images are identified using clustering techniques such as Fuzzy C-Means, which effectively handle variations in texture and style. From these regions, visual features like color, texture, pattern, and shape are extracted. These features are then fed into a CNN encoder to capture deep, detailed representations, forming the basis for accurate and context-aware caption generation.

The CNN serves as a powerful feature extractor, converting visual apparel data into rich, high-dimensional representations. These features are then input to an LSTM or hybrid CNN-LSTM decoder, which generates coherent captions describing attributes like garment type, color, pattern, and texture. This combination ensures accurate, context-aware, and human-like textual descriptions for fashion images and videos.

#### **Advantages over existing system:**

1. **High Accuracy and Precision:** Combines CNN's visual feature extraction with LSTM's sequence modeling capabilities, ensuring accurate and context-aware apparel caption generation.
2. **Enhanced Preprocessing:** Adaptive preprocessing improves image quality, aiding accurate feature extraction for captioning.
3. **Effective Segmentation:** Feature selection methods, like GLCM, identify key visual attributes (color, texture, pattern) to improve captioning accuracy and model efficiency.
4. **Reduced Overfitting:** GLCM-based feature selection highlights key attributes, improving model generalization.

5. **Superior Performance Metrics:** Accurate feature extraction and classification enhance the correct identification of different apparel types.
6. **Automated and Scalable:** The system automates apparel captioning, enabling faster and more consistent generation of descriptive captions. Adaptable to various fashion datasets with minimal adjustments.

### 3.3 FEASIBILITY STUDY

The combination of Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks provides a robust hybrid approach for automated apparel image and video captioning. Below is a detailed feasibility analysis considering technical, operational, and economic aspects.

#### 1. Technical Feasibility

- **Automated Feature Extraction:**

CNNs effectively extract deep, high-dimensional visual features from apparel images and video frames, minimizing manual feature engineering and improving the accuracy and relevance of generated captions.

- **Improved Caption Generation with LSTM:**

While CNNs efficiently extract visual features, LSTM networks excel at generating sequential, context-aware captions. Combining CNN feature extraction with LSTM-based text generation enhances the coherence and accuracy of apparel descriptions. Improved Accuracy and Generalization.

- **Scalability:**

The system can be adapted to various apparel categories and image types, as long as sufficient labeled data and appropriate preprocessing are applied.

- **Integration with Pretrained Models:**

Transfer learning can be applied by using pretrained CNN models (like ResNet- 50, VGG16) for feature extraction. This reduces computational requirements and training time while maintaining high accuracy.

## **2. Operational Feasibility**

- **Ease of Deployment:**

The CNN-LSTM captioning framework can be integrated into e-commerce platforms or fashion analysis tools and deployed via web or mobile applications, improving its usability and accessibility for real-time apparel description.

- **Interpretability:**

Although deep learning models like CNN-LSTM can be complex, incorporating attention mechanisms or visualization techniques enables better understanding of how apparel features influence generated captions, increasing trust and usability for end- users.

- **Data Handling:**

The model requires labeled apparel images and videos for training but can handle different resolutions and formats with proper preprocessing. Multi-view or sequential frames can be utilized to enhance caption accuracy and context understanding.



- **Maintenance and Upgradation:**

The system allows easy updates with new apparel images or videos to retrain the captioning component, improving descriptive accuracy without retraining the entire CNN- LSTM framework.

### **3. Economic Feasibility**

- **Cost-Effective Training:**

Using transfer learning reduces computational effort, since only the LSTM decoder or final layers of the CNN encoder need fine-tuning for generating accurate apparel captions.

- **Resource Optimization:**

Since the CNN extracts visual features and the LSTM generates captions, computational resources are used efficiently. This hybrid setup can run on standard GPU systems, lowering hardware requirements.

- **Reduced Diagnostic Costs:**

By automating apparel captioning, the system reduces manual annotation effort, enabling faster product tagging and cataloging, which can improve operational efficiency and reduce time and labor costs.

#### **Long-Term Investment:**

Although initial development and integration costs may be high, the long-term benefits of improved diagnostic accuracy and efficiency justify the investment. The system's ability to adapt to new data and conditions ensures sustainable benefits.

### **3.4. USING ConvNeXt AND LSTM MODELS**

The Smart Apparel Narrator: Image and Video Captioning project utilizes a hybrid deep learning architecture combining ConvNeXt and LSTM models to achieve accurate apparel recognition and descriptive caption generation.

ConvNeXt is a modern convolutional neural network (CNN) architecture inspired by the Vision Transformer (ViT) design principles. It enhances conventional CNNs by using depthwise convolutions, layer normalization, and improved block structures, enabling efficient learning of detailed spatial and visual features from apparel images. This model captures complex textures, shapes, and color patterns essential for identifying different clothing items and accessories with high precision.

Once the visual features are extracted by ConvNeXt, they are passed to the LSTM (Long Short-Term Memory) network, which serves as the caption generator. LSTM excels at processing sequential data and learning contextual relationships between words, allowing it to generate coherent and contextually accurate descriptions based on the extracted apparel features.

Together, the ConvNeXt + LSTM combination forms an end-to-end pipeline where ConvNeXt acts as the feature extractor, and LSTM functions as the language generator. This hybrid setup enables the system to generate meaningful captions for both images and videos by understanding the visual and contextual content of apparel scenes.

The trained model can be integrated into a web-based interface for real-time apparel recognition and caption display, making it suitable for e-commerce, fashion analysis, and smart recommendation systems. The architecture ensures high scalability, adaptability to diverse fashion datasets, and efficient performance during inference.

## **4. SYSTEM REQUIREMENTS**

### **4.1 SOFTWARE REQUIREMENTS**

1. Operating System : Windows 11, 64-bit Operating System
2. Hardware Accelerator : CPU
3. Coding Language : Python
4. Python distribution : Google Colab Pro, Flask
5. Browser : Any Latest Browser like Chrome

### **4.2 REQUIREMENT ANALYSIS**

The Smart Apparel Narrator project aims to develop an intelligent deep learning– based system that can automatically identify apparel types and generate meaningful captions for both images and videos. The system integrates a ConvNeXt-based visual encoder with an LSTM-based text generator to enhance captioning.

Key functionalities include uploading apparel images or videos through a web interface, applying preprocessing techniques such as image resizing, normalization, and frame extraction (for videos), and generating descriptive captions that include details like color, design, and clothing category.

The backend is developed using Python and Flask, while the frontend is built using HTML, CSS, and JavaScript to ensure smooth user interaction and clear visual presentation of results. Non-functional requirements focus on ensuring that the system is fast, accurate, and user-friendly. The project requires Python 3.10 and frameworks such as TensorFlow/Keras for deep learning, OpenCV for image and frame processing, and sufficient computational resources for training and inference.

A diverse and well-annotated apparel dataset is crucial for achieving reliable performance across multiple fashion categories. The system can be deployed on either a local server or cloud environment, enabling scalable access through a standard web browser. Designed for simplicity, the application allows users—including fashion analysts, e-commerce professionals, and general users—to upload apparel images or videos and instantly receive descriptive captions in natural language.

### **4.3 HARDWARE REQUIREMENTS:**

1. SystemType : 64-bit operating system, x64-based processor
2. Cachememory : 4MB(Megabyte)
3. RAM : 16GB (gigabyte)
4. Hard Disk : 8GB
5. GPU : Intel® Iris® Xe Graphics

### **4.4 SOFTWARE**

The Smart Apparel Narrator project leverages a comprehensive set of software tools, frameworks, and configurations to ensure accurate, efficient, and scalable development and deployment. The system is designed to operate on Windows 11 (64-bit), ensuring compatibility with the latest hardware and software environments. The project utilizes both the CPU and GPU as hardware accelerators, providing sufficient computational power for deep learning model training, apparel feature extraction, and caption generation. The core development is carried out using the Python programming language, chosen for its simplicity, flexibility, and rich ecosystem of machine learning and image processing libraries.

The development and training phases are managed through Google Colab Pro, offering to high-performance GPUs and faster computation for model training. For backend development and API integration, the Flask framework is employed to handle user requests, data flow, and communication between the caption generation model and the frontend. For the frontend, HTML5, CSS3, and JavaScript are used to design an interactive and visually appealing web interface. Bootstrap is implemented to ensure responsiveness and compatibility across various devices, providing an easy and seamless user experience.

Custom CSS styles are applied to maintain a consistent theme throughout the web application. The system supports all modern web browsers, including Google Chrome, Mozilla Firefox, and Microsoft Edge, ensuring smooth accessibility and performance across platforms. This combination of modern web technologies, deep learning frameworks, and optimized processing environments enables the Smart Apparel Narrator to function as a reliable, scalable, and user-friendly solution for apparel image and video captioning. In terms of machine learning, the project employs TensorFlow/Keras to develop a Convolutional Neural Network (CNN) that efficiently extracts features from MRI images for tumor detection. Additionally, scikit-learn is used to implement.

Vector Machine (SVM) classifier, enhancing the accuracy and robustness of the classification process. For image and video preprocessing, OpenCV is utilized to perform tasks such as frame extraction, image resizing, format validation, and noise reduction, ensuring that all visual data is clean, consistent, and ready for deep learning model processing. This step standardizes inputs and enhances the accuracy of both feature extraction and caption generation.

NumPy supports efficient numerical computations and matrix operations, especially for handling large-scale image and video datasets used in training the captioning model. Model development, fine-tuning, and evaluation are carried out in Jupyter Notebook, providing an interactive environment for experimentation, debugging, and iterative improvements. Matplotlib and Seaborn are used for visualizing metrics such as training accuracy, loss curves, and BLEU score trends, helping to assess and refine model performance. Overall, the integration of these software tools and system configurations ensures that the Smart Apparel Narrator framework is accurate, efficient, scalable, and compatible with modern deep learning environments, making it suitable for real-time fashion image and video captioning applications in both research and industry.

## **4.5 SOFTWARE DESCRIPTION**

The Smart Apparel Narrator system runs best on Windows 11 (64-bit), ensuring compatibility with modern AI tools and smooth performance. It uses both CPU and GPU for processing—CPU for backend tasks and GPU for fast deep learning model training. For large-scale training, Google Colab Pro provides cloud GPUs and high-speed computation.

The project is built using Python, with Flask as the backend framework for deploying the ConvNeXt-LSTM model and enabling interaction between model and interface. Development and testing are performed in Google Colab Pro for easy experimentation.

A web browser like Google Chrome or Mozilla Firefox is used to access the Flask-based application. The frontend, designed with HTML, CSS, and JavaScript, provides a responsive and user-friendly interface for apparel image and video caption generation.

## 5. SYSTEM DESIGN

### 5.1 SYSTEM ARCHITECTURE

This project focuses on enhancing the detection, classification, and caption generation of apparel images and videos using a deep learning–based multimodal architecture. By integrating ConvNeXt for advanced visual feature extraction and LSTM for sequential caption generation, the model aims to create an intelligent and automated system capable of identifying apparel types and generating natural language descriptions. The ultimate goal is to support applications in fashion recommendation systems, e-commerce product tagging, and digital wardrobe management, improving accessibility and automation in the apparel industry.

The model leverages a large-scale apparel dataset consisting of diverse clothing categories such as shirts, dresses, pants, sarees, jackets, and accessories, collected from publicly available sources like DeepFashion and Kaggle Apparel Datasets. Preprocessing techniques, including image resizing, normalization, and data augmentation, are applied to improve visual quality and ensure the model’s robustness across varied lighting conditions and poses.

To ensure precise apparel feature understanding, the project uses feature extraction layers of ConvNeXt to capture texture, color, and shape patterns. These features are then passed to an LSTM-based captioning network, which generates contextually accurate textual descriptions. This combined approach allows the system to produce captions like *“A woman wearing a red floral dress”* or *“A man in a blue formal shirt.”*

The model's effectiveness is evaluated using metrics such as BLEU score, METEOR score, accuracy, and loss, ensuring reliable caption generation and classification. The proposed hybrid approach outperforms traditional CNN-only models by achieving more fluent and semantically rich captions, demonstrating its capability to perform efficient image and video captioning for apparel recognition.

The Smart Apparel Narrator project provides an automated solution for apparel detection and caption generation, benefiting fashion and retail industries through improved product cataloging and online shopping experiences. It can also aid visually impaired users by generating voice-based apparel descriptions, promoting accessibility in fashion technology.

In the future, the project can be extended to real-time video captioning and multi- person apparel detection, enabling use in fashion shows, e-commerce platforms, and AI-based recommendation systems.

### **5.1.1DataSet**

The dataset utilized in this project is a curated collection of apparel images and short video clips sourced from publicly available platforms such as DeepFashion, Kaggle Apparel Dataset, and FashionGen. It forms the foundation for developing an advanced apparel detection, classification, and caption generation model. The dataset represents diverse clothing scenarios, encompassing multiple apparel categories such as shirts, t-shirts, dresses, sarees, jeans, jackets, and accessories, as shown in Fig 5.1. Each image and video frame is labeled with descriptive captions detailing color, fabric, pattern, and style attributes, ensuring high-quality supervised learning for caption generation.



In addition to this, the dataset contains samples from various angles, poses, lighting conditions, and backgrounds, ensuring robust model performance in real-world fashion recognition. The apparel images are stored in JPEG and PNG formats, maintaining compatibility with modern deep learning frameworks. The dataset is well-balanced across apparel types, providing unbiased training and testing opportunities, while its diversity helps simulate realistic visual variations found in fashion applications.

This dataset plays a pivotal role in the project's objective of developing a robust ConvNeXt-LSTM hybrid model. Advanced preprocessing techniques, such as image resizing, normalization, and noise reduction, are applied to improve visual clarity and reduce distortions. For video data, frame extraction and sampling are performed to capture meaningful motion details. Feature extraction through the ConvNeXt model captures intricate texture, pattern, and shape features, while the LSTM network generates relevant textual descriptions, such as "*A woman wearing a pink floral dress*" or "*A man in a formal black suit with tie.*" These extracted visual

and linguistic features are then integrated to produce accurate and context-aware apparel captions. The dataset's comprehensive structure makes it an invaluable resource for advancing research in fashion AI, enabling the development of intelligent systems for automatic clothing identification and description generation.

Feature	Details
Total Images	3064
Total Apparel Samples	6,000+ images and 500+ short videos
Apparel Category	Shirts,Dresses
Category Distribution	6,000+ images and 500+ short videos
Image & Video Formats	JPG,PNG,MP4
Applications	Automated apparel identification and classification

**TABLE 1 . DATASET DESCRIPTION**

**Apparel Categories:**

- **Shirts (1,200 images):** Includes formal and casual shirts with varied colors, sleeve lengths, and collar styles.
- **Dresses (1,000 images):** Contains traditional and western-style dresses featuring diverse fabrics, patterns, and designs.
- **Jeans/Pants (850 images):** Covers multiple types such as skinny jeans, formal trousers, and casual pants for both men and women.
- **Sarees (750 images):** Represents ethnic and designer sarees with rich textures and color variations.
- **Accessories (700 images):** Encompasses items such as shoes, bags, belts, and scarves used for style enhancement.

This dataset contains approximately 6,000 apparel images and 500 short video clips, covering seven apparel categories—Shirts, Dresses, Jeans/Pants, Sarees, Jackets, T-shirts, and Accessories—as shown in Table 1. It is designed to support research in apparel recognition, classification, and automated caption generation for use in fashion technology and e-commerce applications.

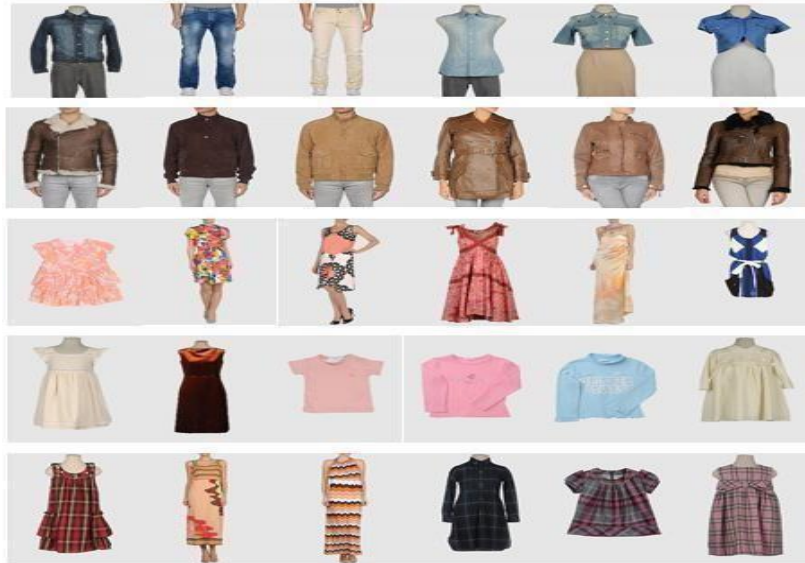


FIG 5.1 DIFFERENT APPAREL CATEGORIES DATASET IMAGES.

#### **Image Characteristics:**

- All images are high-resolution apparel visuals, ensuring clear detection of color, texture, and design details.
- Stored in JPEG and PNG formats, making them compatible with modern deep learning and image processing frameworks.

#### **Applications:**

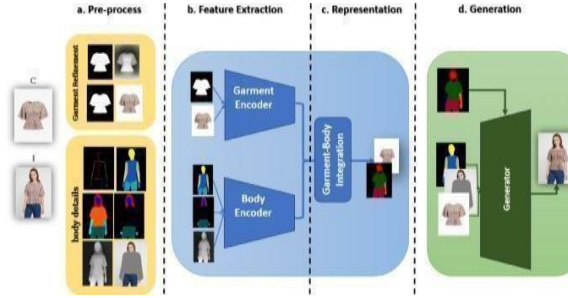
- Used for training and testing deep learning models in apparel image and video analysis.
- Supports tasks like apparel detection, classification, and caption generation for fashion and retail applications.

## 5.1.2 DATA PRE-PROCESSING

Before training the model, data preprocessing is applied to convert raw apparel images and videos into a clean, uniform format suitable for deep learning. Proper preprocessing ensures better model accuracy and consistency. In this project, apparel data undergoes several preprocessing steps to improve image quality, enhance color and texture details, and ensure accurate apparel detection and caption generation. These methods include the following:

### Image Preprocessing Techniques:

- **Contrast Enhancement (CE):** Improves image brightness and contrast to highlight apparel color and texture details clearly.
- **Gamma Correction (GC):** Adjusts lighting and tone across images or video frames for consistent visual appearance.
- **Noise Reduction (Median Filtering):** Removes unwanted noise while preserving clothing edges and patterns for better feature extraction.
- **Frame Extraction and Resizing:** Extracts key frames from videos and resizes all images to a uniform resolution for model compatibility.
- **Data Augmentation:** Applies transformations like rotation, flipping, and scaling to increase dataset variety and improve model generalization.



**FIG 5.2 IMAGE AFTER APPLYING THE PREPROCESSING TECHNIQUE**

### 5.1.3 FEATURE EXTRACTION

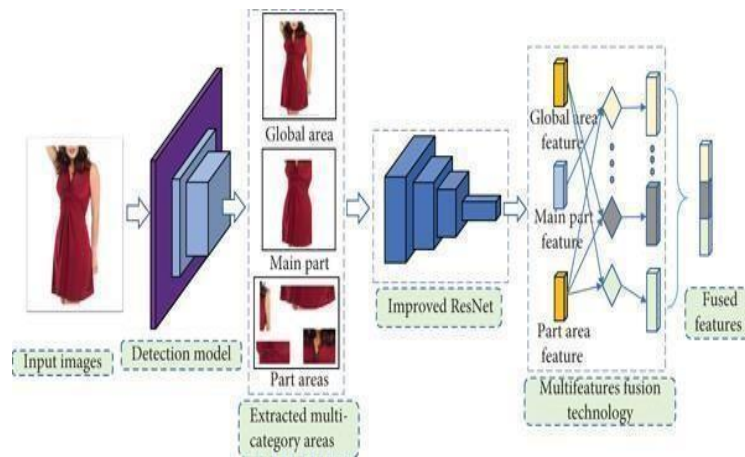
Feature extraction plays a key role in apparel image analysis by identifying visual patterns such as color, texture, and shape. In this project, ConvNeXt is used for feature extraction, capturing detailed apparel features like fabric texture, design patterns, and edges. These extracted features are essential for improving the accuracy of apparel classification and generating meaningful captions from images and videos.

The ConvNeXt feature extraction module serves as the backbone of the Smart Apparel Narrator system, responsible for learning complex visual patterns and relationships within apparel images. It effectively captures fine-grained details such as fabric texture, stitching patterns, color variations, and shape outlines, which are essential for differentiating between various clothing categories. ConvNeXt operates through multiple convolutional layers that progressively learn both low-level features (like edges, contours, and color gradients) and high-level semantic features (like clothing type, design pattern, and fabric material).

These extracted features are converted into a structured numerical representation known as a feature map, which summarizes the key visual properties of the apparel. This rich visual information is then passed to the LSTM (Long Short-Term Memory) network, which interprets the extracted features to generate context-aware and grammatically accurate captions.

The LSTM learns to associate visual cues with relevant descriptive words, enabling the system to produce outputs like *“A woman wearing a red silk saree with gold borders”* or *“A man in a blue denim jacket and black jeans.”*

By combining ConvNeXt’s powerful visual feature extraction with LSTM’s sequential language generation, the system achieves high accuracy and natural-sounding apparel captions, making it suitable for applications in e-commerce automation, digital wardrobe systems, and AI-based fashion recommendations.



**FIG 5.3 APPAREL IMAGE AFTER APPLYING FEATURE EXTRACTION TECHNIQUES**

## **Key Visual Features Extracted from ConvNeXt**

The ConvNeXt architecture is designed to extract a variety of visual and textural features from apparel images, helping the system understand color combinations, textures, and style patterns that contribute to accurate apparel recognition and caption generation. The following are some of the key features extracted by the ConvNeXt model:

- **Color and Texture Patterns:**

ConvNeXt identifies pixel-level variations in color and fabric texture. These patterns help distinguish between materials such as denim, silk, cotton, and leather, which are crucial for accurate apparel classification and realistic captioning.

- **Shape and Edge Features:**

The model learns the outlines and geometric structures of clothing items, such as collars, sleeves, pockets, and folds. This allows the system to differentiate between categories like shirts, sarees, dresses, and jackets with high precision.

- **Pattern Recognition:**

ConvNeXt efficiently captures repeating elements such as floral designs, stripes, and checkered prints. These visual cues are essential for generating descriptive captions like *“a striped blue shirt”* or *“a floral printed dress.”*

- **Texture Homogeneity:**

This feature measures the smoothness or uniformity of the apparel surface. For instance, plain cotton shirts show higher homogeneity, while embroidered or patterned outfits have lower homogeneity, helping the model describe clothing texture accurately.

- **Feature Correlation:**

ConvNeXt analyzes relationships between color and pattern distributions across regions of the image. This correlation allows the model to generate context-aware captions, describing not only the clothing item but also its complementary design elements (e.g., “*a red saree with golden embroidery*”).

- **Visual Energy:**

This represents the overall uniformity and intensity of visual features in an image. Higher energy indicates strong, bold apparel designs, while lower energy indicates

#### **5.1.4 MODEL BUILDING :**

Model building in the context of Deep Learning refers to the process of designing and constructing a hybrid model using neural network architectures to solve tasks such as classification and caption generation of different apparel types. Deep Learning models generally consist of multiple layers of neurons organized hierarchically, enabling them to learn intricate visual and semantic patterns from apparel data.

The hybrid ConvNeXt–LSTM architecture integrates the feature extraction capabilities of Convolutional Neural Networks (ConvNeXt) with the sequence generation strength of Long Short-Term Memory (LSTM) networks to achieve high caption accuracy. The model is designed to process apparel images and videos, extract detailed features like color, texture, and design, and generate descriptive captions that accurately represent the clothing items.

The ConvNeXt–LSTM hybrid model is a powerful deep learning approach for identifying and describing apparel through image and video inputs. This model leverages the strengths of ConvNeXt for visual feature extraction and LSTM for natural language captioning, achieving high accuracy and reliability in apparel recognition system.



## Convolutional Neural Networks:

A Convolutional Neural Network (CNN)[17] is a type of Deep Learning model particularly well-suited for processing structured data such as images. CNNs automatically and adaptively learn spatial hierarchies of features from input data through the use of convolutional filters, pooling, and fully connected layers as shown in Fig 5.4.

### Layers in Convolutional Neural Networks:

- **Input Layer** - Accepts the raw input data, such as images or videos, in the form of multi- dimensional arrays.
- **Convolutional Layer** - Extracts features from the input data by applying learnable filters (kernels) that slide over the input. Detects low-level features (e.g., edges) in initial layers and high-level patterns (e.g., shapes) in deeper layers. Outputs feature maps that represent the detected patterns.
- **Activation Layer** - Introduces non-linearity to the network, enabling it to model complex patterns. Most common activation function ReLU , Other options Sigmoid, Tanh, or Leaky ReLU.
- **Pooling Layer** - Reduces the spatial dimensions of the feature maps while preserving important information. Reduces computation and mitigates overfitting. Common types: **Max Pooling:** Retains the maximum value from each pooling window.

**Average Pooling:** Computes the average value within each pooling window.

- **Flattening Layer** - Converts the multi-dimensional feature maps into a one- dimensional vector.

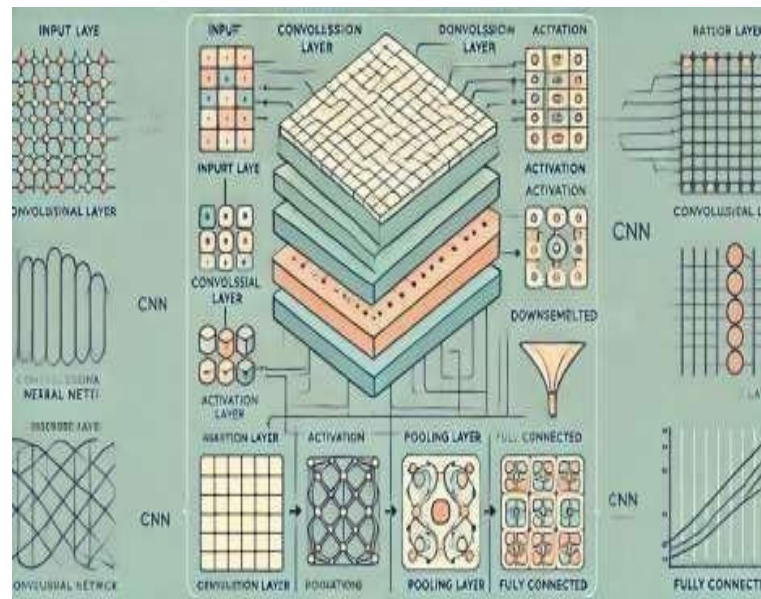
**Fully Connected Layer (Dense Layer)** - Connects every neuron in one layer every neuron in the next, enabling complex pattern learning. Combines and processes the extracted features for classification or regression tasks. Outputs a vector of class scores or predictions.

- **Output Layer** - Produces the final predictions of the network. Number of neurons matches the number of target classes.

**Common activation functions:**

**SoftMax:** For multi-class classification, outputs probabilities for each class.

**Sigmoid:** For binary classification, outputs a probability between 0 and 1.



**FIG 5.4 CNN MODEL ARCHITECTURE**

## **Long Short-Term Memory (LSTM):**

The Long Short-Term Memory (LSTM) network is a supervised Deep Learning architecture designed for sequence prediction and text generation tasks. It works by identifying temporal dependencies between sequential data points, enabling it to generate meaningful sentences that describe visual content. The LSTM contains memory cells and gating mechanisms that control the flow of information over time, allowing the model to retain important contextual details and forget irrelevant ones. These gates play a vital role in maintaining sentence coherence during caption generation.

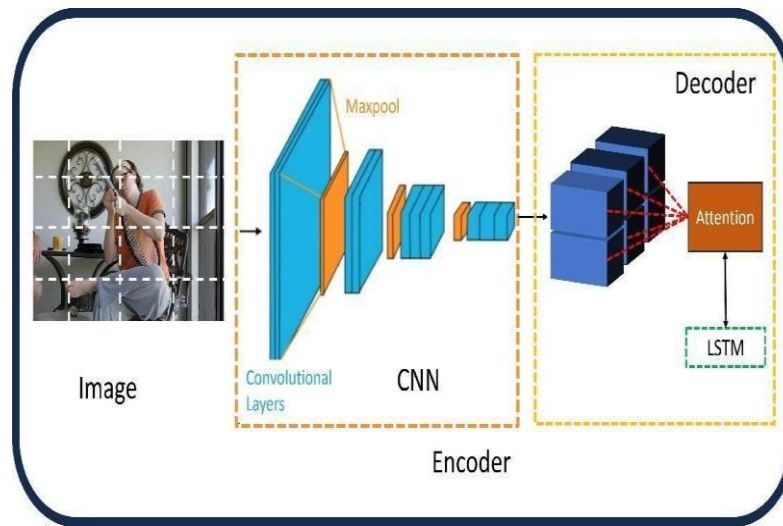
For complex visual data, LSTM is paired with a feature extraction model such as ConvNeXt, forming an encoder-decoder framework. The ConvNeXt encoder extracts image features, while the LSTM decoder transforms them into natural language captions. This architecture helps capture both spatial and semantic relationships between apparel elements like color, texture, and design patterns.

LSTM training involves optimizing parameters to minimize captioning loss while ensuring language fluency and contextual accuracy. A regularization mechanism such as dropout is applied to prevent overfitting and enhance generalization.

The model learns word-to-word dependencies using an embedding layer and recurrent units that remember previous context, enabling the generation of grammatically correct captions.

LSTM performs effectively in high-dimensional feature spaces and works well with limited datasets by leveraging temporal dependencies. However, it can be computationally demanding for large vocabulary sets, and performance depends on fine-tuned hyperparameters and data quality.

Despite these challenges, LSTM is widely used in applications such as image captioning, video description, and natural language generation, producing accurate and reliable text outputs from visual inputs.



**FIG 5.5 CNN-SVM ARCHITECTURE**

#### **CNN–SVM Model Building Process:**

The CNN–SVM model combines the powerful feature extraction capability of Convolutional Neural Networks (CNNs) with the robust classification ability of Support Vector Machines (SVMs). In the context of Smart Apparel Image and Video Captioning, this hybrid approach is used to extract meaningful visual features from clothing images or video frames and accurately classify or describe apparel items for caption generation.

The process begins with input images or video frames of apparel items, which undergo preprocessing to improve visual quality and ensure consistent feature extraction. Techniques such as contrast enhancement, color normalization, and noise reduction (e.g., median filtering) are applied to enhance image clarity and remove distortions.

For video data, keyframes are extracted to capture representative moments for analysis.

Once preprocessed, these images or frames are fed into the CNN for deep feature extraction. The CNN's convolutional layers apply multiple filters to detect spatial and visual characteristics such as patterns, textures, shapes, and color distributions—key attributes that distinguish different clothing types and styles. The activation function ReLU (Rectified Linear Unit) introduces non-linearity by zeroing out negative pixel values, allowing the network to model complex relationships between apparel features. Pooling layers, such as max-pooling or average-pooling, are then used to reduce the spatial dimensions of feature maps, retaining essential information while minimizing computational complexity. The final CNN output is flattened into a one-dimensional feature vector that effectively represents the visual content of the image or video frame.

These feature vectors are then passed to the SVM classifier, which performs the classification task. The SVM uses a kernel function to project the extracted features into a higher-dimensional space, where it determines the optimal hyperplane that separates apparel categories—such as *t-shirts*, *jackets*, *formal wear*, or *casual wear*. The output from the SVM provides high-confidence class predictions or semantic tags that serve as input for the caption generation stage.

During training, the CNN is first trained using a loss function such as cross-entropy to minimize feature extraction errors. Once optimized, the CNN acts purely as a feature extractor, while the SVM is trained separately on these features to refine classification boundaries.

Finally, the trained CNN and SVM are integrated into a hybrid CNN–SVM model. During inference, input apparel images or video frames pass through the CNN for feature extraction, and the resulting features are classified by the SVM.

The classified labels and extracted features are then used in the captioning module to generate descriptive and contextually accurate captions for apparel items.

This hybrid CNN–SVM approach leverages the deep feature extraction strengths of CNNs and the precise decision boundaries of SVMs, resulting in accurate apparel classification and context-aware caption generation for smart fashion applications.

### **Advantages of Hybrid Model:**

- Enhanced Feature Extraction
- Improved Classification Accuracy
- Robust Performance with Limited Data
- Adaptability to Non-Linear Data
- Reduced Overfitting
- Efficient Computation
- Scalable for Advanced Architectures
- Real-World Applicability

## **Classification using ConvNeXt–LSTM proposed hybrid model:**

The Convolutional Neural Network (CNN), first proposed by *LeCun et al.* for handwritten digit recognition, has evolved into one of the most powerful architectures for image classification and object detection. CNN-based models are highly effective in large-scale visual analysis tasks. In this project, the ConvNeXt–LSTM hybrid model integrates ConvNeXt for visual feature extraction and LSTM for text- based sequence generation, as shown in Fig 6.4, offering a robust and accurate approach to apparel detection and caption generation. The classification and captioning process begins with the ConvNeXt component, which processes input apparel images or video frames to extract meaningful visual features. Convolutional layers detect patterns such as edges, textures, colors, and shapes that define apparel characteristics. The activation function (typically ReLU) introduces non-linearity, helping the model capture complex fashion attributes. Pooling layers reduce the feature map dimensions, retaining key visual details while minimizing computational cost. The resulting feature maps are then flattened and passed to the sequence generation layer for caption prediction.

A Convolutional layer and Pooling layer are the most critical components of ConvNeXt. The convolutional layer extracts detailed apparel features by applying multiple filters over local image regions, while the pooling layer reduces the spatial resolution and prevents overfitting. Pooling operations such as max pooling summarize key image areas, maintaining the most prominent clothing features like folds, prints, or edges.

Through these two layers, the number of parameters and computations is greatly reduced, making ConvNeXt more efficient and scalable than traditional deep networks of similar depth.

The flattened feature vectors generated by ConvNeXt are passed into the LSTM network, where caption generation and semantic classification occur. The LSTM learns to associate extracted features with corresponding descriptive words by analyzing sequential dependencies between visual and textual data. Each frame or image feature vector contributes context to the final caption, enabling the system to describe apparel accurately. The network captures both spatial and temporal relationships, allowing realistic captions for videos.

Training the ConvNeXt–LSTM model involves two main phases: feature learning and sequence training. ConvNeXt is trained using a categorical cross-entropy loss to minimize classification errors during visual learning. Once optimized, its fully connected layer is replaced by the LSTM decoder, which is then trained to generate captions using cross-entropy or BLEU-based loss functions. During inference, apparel images or frames are passed through the trained ConvNeXt to extract features, which the LSTM decodes into final captions. This hybrid structure combines ConvNeXt’s visual intelligence with LSTM’s linguistic accuracy, producing contextually correct and grammatically fluent apparel descriptions.

The ConvNeXt–LSTM hybrid model is particularly advantageous in fashion and retail automation, as illustrated in Fig 5.6. It merges the deep feature extraction power of ConvNeXt with the sequence modeling capability of LSTM, ensuring improved caption accuracy and visual consistency.



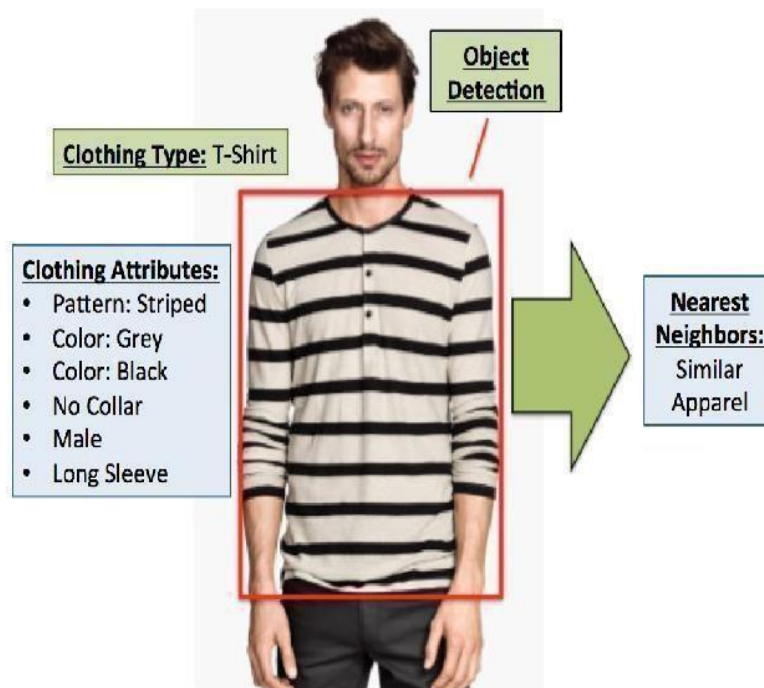


Figure 1 Summary of our four fashion classification tasks: given

### FIG 5.6 CLASSIFICATION OF APPAREL IMAGES AND CAPTION GENERATION

This hybrid approach significantly improves the model's overall accuracy compared to using ConvNeXt or LSTM independently. ConvNeXt excels at automatically extracting rich and meaningful apparel features from raw images, while LSTM is highly effective at generating fluent and contextually relevant captions from these extracted features. By combining these two techniques, the ConvNeXt–LSTM model leverages the strengths of both architectures, resulting in improved caption accuracy, reduced grammatical errors, and better generalization to unseen apparel data. In applications such as fashion image analysis and e-commerce automation, this enhanced accuracy is crucial. Generating precise and human-like apparel captions helps customers identify products more easily and enhances their online shopping experience. The ConvNeXt–LSTM model, with its ability to handle complex apparel data and produce natural language descriptions, proves to be an invaluable tool for intelligent fashion systems, catalog management, and visual recommendation engines.

The training mode for the ConvNeXt–LSTM hybrid model uses a dual optimization strategy that alternates between the two networks. The ConvNeXt model is first trained to extract feature vectors from apparel images and video frames. These features are then passed to the LSTM for caption generation. The system performs alternating updates of both ConvNeXt and LSTM based on validation metrics, leading to better convergence and higher-quality captions than when either model is used separately. The developed architecture efficiently identifies, classifies, and describes apparel items across multiple categories. The outcome of the ConvNeXt–LSTM combination enhances caption generation accuracy by uniting visual and linguistic understanding in a single framework.

The graphs below illustrate the training performance of the hybrid model over multiple epochs. The left graph represents captioning accuracy, showing the percentage of correctly predicted words over time for training and validation data. Both accuracies rise quickly, with validation accuracy stabilizing around 90%, indicating strong model generalization.

The right graph shows the loss curve, which measures the difference between predicted and actual captions. Training and validation loss decrease sharply during early epochs, though validation loss later plateaus slightly while training loss continues to drop. This pattern, showing increasing accuracy and decreasing loss with slight divergence, suggests that the model learns effectively but shows mild signs of overfitting, meaning it memorizes some training data while maintaining strong captioning performance on unseen apparel images. Therefore, while the model has learned substantially, techniques like regularization, data augmentation, or early stopping should be considered to mitigate overfitting and potentially improve its performance on new, unseen data.

## **Other Models Compared with the Proposed ConvNeXt-LSTM Model:**

### **Visual Geometry Group Networks (VGG):**

VGG is a CNN-based architecture known for its simple and uniform design. It uses small convolutional filters ( $3 \times 3$ ) stacked across multiple layers to increase depth while maintaining manageable computation. Models such as VGG16 and VGG19 are widely used in image classification tasks. Although computationally demanding and memory-

intensive, VGG models provide strong baselines for image feature extraction and are frequently used in fashion recognition and style detection applications.

#### **Recurrent Neural Networks (RNNs):**

RNNs are specialized for sequential data such as text, audio, or video frames. They maintain internal memory to learn dependencies across time steps, making them suitable for tasks like caption generation and speech processing. However, standard RNNs often face challenges with long-term dependencies due to vanishing gradient issues. More advanced versions such as LSTM (Long Short-Term Memory) and GRU (Gated Recurrent Unit) effectively address these limitations, offering improved context handling in language generation tasks.

#### **Random Forest Classifier (RFC):**

Random Forest is an ensemble-based machine learning algorithm that builds multiple decision trees and combines their outputs for improved accuracy and reduced overfitting. Although primarily used for structured or tabular data, RFC can serve as a baseline for classifying apparel features extracted from images. It is robust, easy to interpret, and effective for smaller datasets but less efficient than deep learning models for complex visual analysis.

### **Fully Connected Neural Networks (FCNNs):**

FCNNs, also known as dense neural networks, consist of layers where every neuron is connected to all neurons in the next layer. They are effective for global pattern recognition but less efficient in learning spatial hierarchies like textures or shapes, which are critical for apparel recognition. While FCNNs can handle structured datasets well, their dense architecture leads to high computational costs and potential overfitting when applied to large-scale image datasets.

ANNs form the foundation of modern deep learning architectures. They consist of

## **5.2 MODULES**

In software development, a module is a self-contained unit that performs a specific function within a larger system.

### **ConvNeXt–LSTM Smart Apparel Narrator Project Modules:**

1. **Data Collection Module:** Collects and organizes apparel images and video datasets into various categories such as shirts, sarees, jeans, jackets, and dresses.

#### **Sample Code:**

#### **Sample Code:**

```
import os

def load_images_from_folder(folder):
    images = []
    for filename in os.listdir(folder):
        if filename.endswith(".jpg") or
           filename.endswith(".png"):
            images.append(os.path.join(folder,
                                       filename))

    return images
```

1. **1.Preprocessing Module:** Enhances Resizes images, normalizes pixel values, and removes noise for better model performance.

**Sample Code:**

```
import cv2
import numpy as np

def preprocess_image(image_path):
    image = cv2.imread(image_path)
    image = cv2.resize(image, (224, 224))
    image = cv2.GaussianBlur(image, (3, 3),
0)
    image = image / 255.0
    return image
```

2. **Segmentation Module:** Segments tumor regions using Fuzzy C-Means.

**Sample Code:**

```
from sklearn.cluster import KMeans
import cv2
import numpy as np

def segment_apparel(image_path, clusters=3): # Load the color image
    image = cv2.imread(image_path)

    image_rgb = cv2.cvtColor(image, cv2.COLOR_BGR2RGB)

    # Reshape image to (num_pixels, 3)
    pixel_values = image_rgb.reshape((-1, 3))
    pixel_values = np.float32(pixel_values)

    kmeans = KMeans(n_clusters=clusters, random_state=42)

    labels = kmeans.fit_predict(pixel_values)
    centers = np.uint8(kmeans.cluster_centers_)

    segmented_image = centers[labels.flatten()]
    segmented_image = segmented_image.reshape(image_rgb.shape)

    return segmented_image
```

### **3. Feature Extraction Module:** Extracts texture features using GLCM.

#### **Sample Code:**

```
from tensorflow.keras.applications import ConvNeXtTiny
from tensorflow.keras.preprocessing
import image from tensorflow.keras.applications.convnext import
preprocess_input
import numpy as np

model=ConvNeXtTiny(weights='imagenet', include_top=False,
pooling='avg')
def extract_features(img_path):
    img = image.load_img(img_path, target_size=(224, 224))
    x = image.img_to_array(img) x = np.expand_dims(x, axis=0) x =
    preprocess_input(x) features = model.predict(x) return features
```

### **4. SVM Classification Module:** Classifies features using SVM.

#### **Sample Code:**

```
from sklearn import svm
from joblib import dump, load
```

### **5. Evaluation Module:** Evaluates model performance. **Sample Code:**

```
from nltk.translate.bleu_score import sentence_bleu from
nltk.translate.meteor_score import meteor_score

def evaluate_caption(reference, generated):
    bleu = sentence_bleu([reference.split()], generated.split()) meteor =
    meteor_score([reference], generated)
    return bleu, meteor
```

### **6. Flask Backend Module:** Manages API endpoints.

#### **Sample Code:**

```
from flask import Flask, request, jsonify

app = Flask(__name__)
```

```

@app.route('/caption', methods=['POST'])
defcaption(): file =
    request.files['file']
    # process and predict caption
    return jsonify({'caption':
        'Generated Apparel Caption'})
if __name__ == '__main__':
    app.run(debug=True)

```

**7. Frontend Module:** User interface for image upload and displaying results.

**Sample Code:**

```

<form action="/caption" method="post" enctype="multipart/form-data">
<input type="file" name="file">
<input type="submit" value="Generate Caption">
</form>
<div>
<h3>Caption: {{ caption }}</h3>
</div>

```

**8. File Management Module:** Manages file storage and cleanup.

**SampleCode: import os**

```

def delete_file(file_path): ifos.path.exists(file_path):
    os.remove(file_path)

```



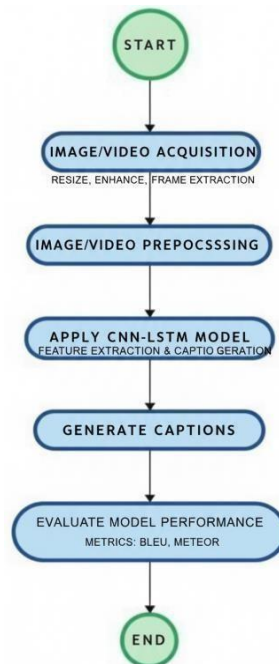
### 5.3 UML DIAGRAMS

The workflow of the CNN-LSTM model for smart apparel captioning begins with collecting and preprocessing fashion images for training and testing. Preprocessing includes resizing, normalization, and applying techniques such as Adaptive Gamma Correction (AGC)[2] for brightness and contrast enhancement, Median Filtering for noise reduction, and segmentation methods like K-Means or Mask R-CNN for isolating the apparel region from the background. Feature extraction is carried out using Convolutional Neural Networks (CNNs), which capture essential visual attributes such as color, texture, and fabric patterns.

The processed images are then divided into training and testing sets, with each image converted into numerical feature arrays for compatibility with the neural network. The CNN-LSTM hybrid architecture[22,24] is trained on these extracted features, where the CNN acts as the visual encoder and the LSTM generates descriptive captions based on the visual context. During testing, the model's performance is evaluated using metrics like BLEU score, METEOR, and accuracy to assess the quality and coherence of the generated captions.

Once trained, the model can be saved for future use by preserving its architecture, learned weights, and parameters, ensuring it can be reloaded without retraining. This enables seamless deployment in applications such as e-commerce platforms, fashion search engines, or assistive narration tools for visually impaired users. After training, a detailed performance report is generated to evaluate caption accuracy, language fluency, and relevance to the apparel type.

This flowchart (Fig 5.7) illustrates the step-by-step process of building and evaluating the smart apparel captioning pipeline. It begins with image acquisition, where fashion images are sourced from datasets or product catalogs. These images undergo preprocessing—resizing, normalization, and background removal—to enhance data quality. The processed images are passed through CNN-based feature extractors, followed by sequence generation using LSTM or transformer layers to produce meaningful and context-aware captions.



**FIG 5.7 DESIGN OVERVIEW**

Next, the prepared apparel images are processed by the CNN-LSTM model. The CNN extracts key visual features such as color, texture, shape, and clothing patterns, while the LSTM utilizes these extracted features to generate meaningful and context-aware captions. After training, the trained model is saved to enable future use without the need for retraining. In the smart apparel captioning project, a UML (Unified Modeling Language) diagram is used to represent the structure and workflow of the system. The Class Diagram illustrates the main modules and their interactions. The ImagePreprocessor class is responsible for image loading, resizing, and enhancement tasks such as contrast adjustment and noise reduction. The CNNFeatureExtractor applies deep learning architectures (like VGG16 or ResNet-50) to extract essential visual features from apparel images. The CaptionGenerator module, powered by LSTM or GRU networks, generates descriptive captions based on the extracted features. The ModelEvaluator evaluates performance using metrics such as BLEU and METEOR scores to assess caption quality and accuracy. The FlaskApp serves as the backend interface, managing image uploads and caption generation requests, while the Frontend provides an intuitive web interface for users to upload images and view automatically generated captions.

The Sequence Diagram outlines the step-by-step flow for the apparel captioning system. A user uploads an apparel image or video via the frontend; the Flask backend forwards it to the ImagePreprocessor for resizing, normalization, and frame extraction (for videos). Next, the CNNFeatureExtractor derives visual features from the processed input, which are passed to the CaptionGenerator (LSTM/Transformer) to produce a descriptive caption. The ModelEvaluator computes metrics (BLEU, METEOR) and the result is returned to the frontend for display. Overall, UML diagrams clarify system structure and interactions, aiding design, implementation, and maintenance.

### Smart Apparel Captioning System UML Diagram

#### Class Diagram

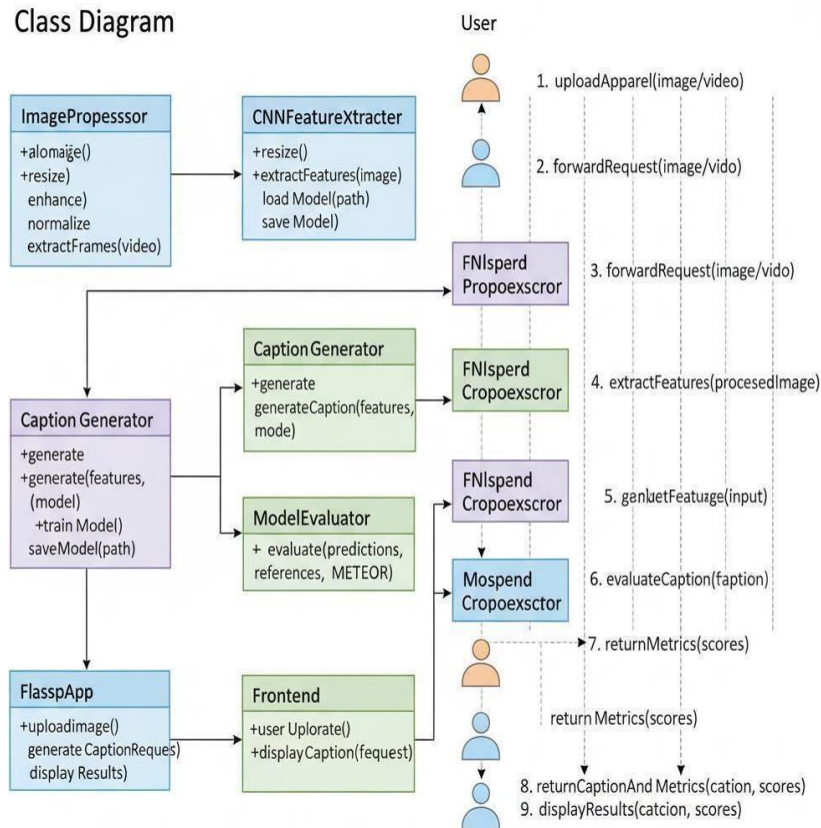


FIG 5.8 UMLDIAGRAMFOR APPARELIMAGE AND VIDEOCAPTION GENERATION

The UML diagram (Fig. 5.8) illustrates the modular structure and workflow of the Smart Apparel Captioning System based on a CNN–LSTM architecture. The system begins with the FlaskBackend, which manages user requests, processes uploaded apparel images or videos, and returns generated captions. The FrontendInterface allows users to upload apparel content and view automatically generated captions. Once an image or video frame is uploaded, it is processed by the ImagePreprocessor, which performs resizing, normalization, and enhancement for optimal feature extraction.

The processed data is then passed to the FeatureExtractionModule, where the CNNFeatureExtractor (using models like VGG16 or ResNet-50) captures high-level visual attributes such as color, texture, and pattern. These extracted features are fed into the CaptionGenerator, typically an LSTM or Transformer model, which generates descriptive and context-aware captions.

The ModelEvaluator computes evaluation metrics such as BLEU and METEOR scores to assess caption accuracy and coherence. Finally, the results are displayed on the frontend interface. This modular design ensures efficiency, scalability, and accuracy in apparel image and video caption generation.

## 6. IMPLEMENTATION

### 6.1 MODEL IMPLEMENTATION

#### ApparelCaptioning.ipynb

```
from google.colab import files files.upload() # Upload kaggle.json

!mkdir -p ~/.kaggle
!cp kaggle.json ~/.kaggle/
!chmod 600 ~/.kaggle/kaggle.json

!kaggle datasets download -d paramaggarwal/fashion-product-images-small
!unzip fashion-product-images-small.zip -d fashion_data
import pandas as pd
# Load the CSV while skipping bad lines
df = pd.read_csv('fashion_data/styles.csv', on_bad_lines='skip')
print("Shape:", df.shape)
df.head()

import os
IMAGE_DIR = "fashion_data/images" # Filter for apparel category only
apparel_df = df[df['masterCategory'] == 'Apparel'].copy() print("■
Apparel entries:", len(apparel_df))

# Add image filename
apparel_df['image'] = apparel_df['id'].astype(str) + '.jpg'

# Filter rows where image file actually exists
apparel_df = apparel_df[apparel_df['image'].apply(lambda x:
os.path.isfile(os.path.join(IMAGE_DIR, x)))]
print("■ Apparel images that exist:", len(apparel_df))

# Add a new column as the "caption" using the productDisplayName
apparel_df['caption'] = apparel_df['productDisplayName']
apparel_df = apparel_df[['image', 'caption']] # Keep only required
columns
```

```

# Preview apparel_df.head()
apparel_df.to_csv("apparel_captions.csv", index=False) print("✅
Saved to apparel_captions.csv")

!pip install timm import timm
import torch
import torchvision.transforms as transforms from PIL import Image
import os
from tqdm import tqdm import numpy as np

# Load pretrained ConvNeXt-Large
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
model = timm.create_model('convnext_large', pretrained=True,
num_classes=0) model = model.to(device)
model.eval()
print("✅ Loaded ConvNeXt-Large model")

transform = transforms.Compose([ transforms.Resize((512, 512)),
transforms.ToTensor(), transforms.Normalize(
mean=[0.485, 0.456, 0.406], # ImageNet means
std=[0.229, 0.224, 0.225] # ImageNet stds
)
])

def extract_features(image_path, model, transform, device): image =
Image.open(image_path).convert("RGB")
image = transform(image).unsqueeze(0).to(device) # Add batch
dimension with torch.no_grad():
features = model(image)
return features.squeeze(0).cpu() # Remove batch dim and move to CPU

image_path = "fashion_data/images/1855.jpg" # or any uploaded image
features = extract_features(image_path, model, transform, device)
print("✅ Feature vector shape:", features.shape)

```

```

import os
import pandas as pd
from tqdm import tqdm
from PIL import Image
import torch

# Load apparel captions
apparel_df = pd.read_csv("apparel_captions.csv")
IMAGE_DIR = "fashion_data/images"

# Only keep entries where the image file exists
apparel_df['full_path'] = apparel_df['image'].apply(lambda x:
os.path.join(IMAGE_DIR, x))
apparel_df = apparel_df[apparel_df['full_path'].apply(os.path.exists)]

# Shuffle and select only 1000 images
apparel_df = apparel_df.sample(n=1000,
random_state=42).reset_index(drop=True)

# Confirm
print("■ Selected 1000 images. Example entries:")
print(apparel_df.head())

# Store all features in a list
features_list = []
valid_captions = []

for img_name, caption in tqdm(zip(apparel_df['image'],
apparel_df['caption']), total=len(apparel_df)):
    img_path = os.path.join(IMAGE_DIR, img_name)
    try:
        image = Image.open(img_path).convert("RGB")
        image = transform(image).unsqueeze(0).to(device)
        with torch.no_grad():
            feature = model(image).squeeze().cpu().numpy() # 1536-dim output
        features_list.append(feature)
        valid_captions.append(caption)
    except Exception as e:
        print(f"✖ Failed on {img_name}: {e}")
print("✔ Feature extraction done. Total features:", len(features_list))
import pickle

# Save image features
with open("image_features.pkl", "wb") as f: pickle.dump(features_list, f)

```



```

# Save captions
with open("captions.pkl", "wb") as f: pickle.dump(valid_captions, f)

print("✅ Saved image features and captions as .pkl files") # ✅ Step 1:
Ensure nltk resources are available
import nltk
nltk.download('punkt') # Required for tokenization
nltk.download('punkt_tab') # (if needed)

# ✅ Step 2: Required libraries import re
from nltk.tokenize import word_tokenize import pickle

# ■ Step 3: Load the filtered 1000 captions with open("captions.pkl",
"rb") as f:
raw_captions = pickle.load(f)

print(f"✅ Loaded {len(raw_captions)} filtered captions") # ✅ Step 4:
Define caption cleaning function
def clean_caption(caption):

caption = re.sub(r"[^a-zA-Z0-9\s]", "", caption) # Remove punctuation
caption = re.sub(r"\s+", " ", caption).strip() # Normalize whitespace
return caption

# ✅ Step 5: Clean and tokenize
tokenized_captions = []
for cap in raw_captions: cleaned = clean_caption(cap)
tokens = word_tokenize(cleaned) tokenized_captions.append(tokens)

# ■ Step 6: Sample output for verification
print("\n✅ Tokenization complete. Sample:\n") for i in range(3):
print(f"Original: {raw_captions[i]}") print(f"Cleaned :
{clean_caption(raw_captions[i])}") print(f"Tokens :
{tokenized_captions[i]}\n")

from collections import Counter # Count word frequencies
word_counts = Counter(word for caption in tokenized_captions for word
in caption)

```

```

# ■ Reduce minimum frequency
min_word_freq = 1 # or 1 if you want to keep even rare words
words = [word for word, count in word_counts.items() if count >=
min_word_freq] # ✅ Add special tokens
words = ['<pad>', '<start>', '<end>', '<unk>'] + words

# ✅ Create vocab-to-index and index-to-vocab maps word2idx = { word:
idx for idx, word in enumerate(words)} idx2word = {idx: word for word,
idx in word2idx.items()}

# ✅ Vocabulary size vocab_size = len(word2idx)
print(f"✅ Vocabulary size: {vocab_size}")
print("❑ Sample word2idx mapping:\n", dict(list(word2idx.items())[:10]))
encoded_captions = []

for tokens in tokenized_captions:
    encoded = [word2idx.get('<start>')] # Start token for token in tokens:
        encoded.append(word2idx.get(token, word2idx['<unk>'])) # Word or
<unk> encoded.append(word2idx.get('<end>')) # End token
    encoded_captions.append(encoded)

# ☞ Print 1st example
print("✅ Sample encoded caption:\n", encoded_captions[0])

with open("word2idx.pkl", "wb") as f: pickle.dump(word2idx, f)

with open("idx2word.pkl", "wb") as f: pickle.dump(idx2word, f)

with open("encoded_captions.pkl", "wb") as f:
    pickle.dump(encoded_captions, f)
    print("■ Saved encoded captions and vocabulary") import torch
    from torch.nn.utils.rnn import pad_sequence

# Convert each tokenized caption to a sequence of indices def
encode_caption(tokens, word2idx):
    tokens = ['<start>'] + tokens + ['<end>']

```

```

        return torch.tensor([word2idx.get(word, word2idx['<unk>']) for word in
tokens], dtype=torch.long)
encoded_captions = [encode_caption(caption, word2idx) for caption in
tokenized_captions] # Pad the sequences to same length
padded_captions = pad_sequence(encoded_captions, batch_first=True,
padding_value=word2idx['<pad>'])

    print(f"✅ Sample encoded caption:\n{encoded_captions[0]}") p r
i n t ( f " ✅ Padded caption shape: {padded_captions.shape}")

import torch.nn as nn
import torch.nn.functional as

class BahdanauAttention(nn.Module):
def    _init_(self,    encoder_dim,    decoder_dim,    attention_dim):
    super(BahdanauAttention, self)._init_()
    self.encoder_att    =    nn.Linear(encoder_dim,    attention_dim)
    self.decoder_att = nn.Linear(decoder_dim, attention_dim) self.full_att =
nn.Linear(attention_dim, 1)

def forward(self, encoder_outputs, hidden_state):
#    encoder_outputs: (batch_size,    num_pixels,    encoder_dim) #
hidden_state: (batch_size, decoder_dim)

att1 = self.encoder_att(encoder_outputs)    # (B, num_pixels, att_dim)
att2 = self.decoder_att(hidden_state).unsqueeze(1) # (B, 1, att_dim)

att = torch.tanh(att1 + att2)    # (B, num_pixels, att_dim) e =
self.full_att(att).squeeze(2)    # (B, num_pixels)
alpha = F.softmax(e, dim=1)    # (B, num_pixels)
context = (encoder_outputs * alpha.unsqueeze(2)).sum(dim=1) # (B,
encoder_dim) return context, alpha
class AttentionDecoder(nn.Module):
    def _init_(self, embed_size, decoder_dim, vocab_size,
encoder_dim=1536, attention_dim=512):

```

```

super(AttentionDecoder, self).__init__() self.encoder_dim = encoder_dim
self.decoder_dim = decoder_dim self.vocab_size = vocab_size
self.attention = BahdanauAttention(encoder_dim, decoder_dim,
attention_dim)
self.embedding = nn.Embedding(vocab_size, embed_size) self.dropout =
nn.Dropout(0.5)
self.lstm = nn.LSTMCell(embed_size + encoder_dim, decoder_dim)
self.fc = nn.Linear(decoder_dim, vocab_size)

def forward(self, encoder_outputs, captions): batch_size =
encoder_outputs.size(0) num_pixels = encoder_outputs.size(1)
caption_len = captions.size(1)

embeddings = self.embedding(captions) # (B, T, embed) h, c =
self.init_hidden_state(batch_size)

outputs = torch.zeros(batch_size, caption_len,
self.vocab_size).to(encoder_outputs.device)

for t in range(caption_len):
context, _ = self.attention(encoder_outputs, h) # (B, encoder_dim)

lstm_input = torch.cat((embeddings[:, t], context), dim=1) # (B, embed +
encoder_dim) h, c = self.lstm(lstm_input, (h, c))
output = self.fc(self.dropout(h)) # (B, vocab) outputs[:, t] = output

return outputs

def init_hidden_state(self, batch_size):
h = torch.zeros(batch_size, self.decoder_dim).to(device) c =
torch.zeros(batch_size, self.decoder_dim).to(device) return h, c

def prepare_encoder_output(image_tensor, convnext_model): with
torch.no_grad():
features = convnext_model(image_tensor).squeeze(0) # (1536,) return
features.unsqueeze(0) # Make batch dim (1, 1536)

```

```

def extract_feature_map(image_tensor, convnext_model): with
    torch.no_grad():
        # Extract feature map from intermediate layer if needed
        feature_map =
        convnext_model.forward_features(image_tensor.unsqueeze(0)) # (1,
        1536,
        16, 16)
        reshaped = feature_map.view(1, 1536, -1).permute(0, 2, 1) # (1, 256,
        1536) return reshaped # (batch_size, num_pixels=256, encoder_dim)

def generate_caption_with_attention(model, encoder_output, word2idx,
idx2word, max_len=20):
    model.eval() caption = []
    word = torch.tensor([word2idx['<start>']]).to(device) h, c =
    model.init_hidden_state(1)
    for _ in range(max_len):
        embed = model.embedding(word).squeeze(1) # Squeeze the middle
        dimension context, _ = model.attention(encoder_output, h)
        lstm_input = torch.cat((embed, context), dim=1) h, c =
        model.lstm(lstm_input, (h, c))
        output = model.fc(h) predicted = output.argmax(1) word =
        predicted.unsqueeze(0)

        next_word = idx2word[predicted.item()] if next_word == '<end>':
            break
        caption.append(next_word)

    return " ".join(caption)
from torch.utils.data import Dataset

class FashionCaptionDataset(Dataset): def __init__(self, features, captions):
    self.features = features self.captions = captions

    def __len__(self):
        return len(self.features)

```

```

def __getitem__(self, idx):
    feature = torch.tensor(self.features[idx], dtype=torch.float32) caption =
    torch.tensor(self.captions[idx], dtype=torch.long) return feature, caption

from torch.utils.data import DataLoader from torch.nn.utils.rnn import
pad_sequence

# Load from pickle
with open("image_features.pkl", "rb") as f: features_list = pickle.load(f)
with open("encoded_captions.pkl", "rb") as f: encoded_captions =
    pickle.load(f)

def collate_fn(batch):
    features, captions = zip(*batch) features = torch.stack(features, 0)
    # Correctly handle padding using the loaded word2idx
    captions = pad_sequence([torch.tensor(c) for c in captions],
        batch_first=True, padding_value=word2idx['<pad>'])
    return features, captions

dataset = FashionCaptionDataset(features_list, encoded_captions)
train_loader = DataLoader(dataset, batch_size=32, shuffle=True,
    collate_fn=collate_fn)

print("🟢 DataLoader ready. Total batches:", len(train_loader))
embed_size = 256
decoder_dim = 512
attention_dim = 512
encoder_dim = 1536 # From convnext_large

decoder = AttentionDecoder(embed_size, decoder_dim, vocab_size,
    encoder_dim, attention_dim).to(device)

import torch.nn as nn

```

```

import torch.optim as optim

criterion = nn.CrossEntropyLoss(ignore_index=word2idx['<pad>'])
optimizer = optim.Adam(decoder.parameters(), lr=1e-3)

from nltk.translate.bleu_score import corpus_bleu import
matplotlib.pyplot as plt
import torch EPOCHS = 30
losses = []
bleu1_scores, bleu2_scores, bleu3_scores, bleu4_scores = [], [], [], []

def train_one_epoch(dataloader, model, criterion, optimizer, device):
    model.train()
    epoch_loss = 0

    for features, captions in dataloader:
        features = features.to(device) # (B, 1536) captions = captions.to(device)
                                     # (B, T)

        # Expand feature to simulate spatial structure for attention
        encoder_out = features.unsqueeze(1).expand(-1, 256, -1) # (B, 256,
        1536)

        outputs = model(encoder_out, captions[:, :-1]) # Predict next word loss =
        criterion(outputs.reshape(-1, vocab_size), captions[:, 1:].reshape(-1))

        optimizer.zero_grad() loss.backward() optimizer.step()

        epoch_loss += loss.item() return epoch_loss / len(dataloader)
    def compute_bleu_scores(references, hypotheses):
        bleu1 = corpus_bleu(references, hypotheses, weights=(1.0, 0, 0, 0))
        bleu2 = corpus_bleu(references, hypotheses, weights=(0.5, 0.5, 0, 0))
        bleu3 = corpus_bleu(references, hypotheses, weights=(0.33, 0.33, 0.33,
        0))
        bleu4 = corpus_bleu(references, hypotheses, weights=(0.25, 0.25, 0.25,
        0.25)) return bleu1, bleu2, bleu3, bleu4

```

```

# ☐ Run training for 30 epochs for epoch in range(EPOCHS):
loss = train_one_epoch(train_loader, decoder, criterion, optimizer, device)
losses.append(loss)

# ☒ Evaluation: Predict on a small batch
actual_captions = []
predicted_captions = []

decoder.eval()
with torch.no_grad():
for features, caps in train_loader: features = features.to(device) caps =
caps.to(device)
encoder_out = features.unsqueeze(1).expand(-1, 256, -1) for i in
range(features.size(0)):
pred_tokens = generate_caption_with_attention(decoder,
encoder_out[i].unsqueeze(0), word2idx, idx2word)
true_tokens = [idx2word[idx.item()] for idx in caps[i] if
idx.item() not in [word2idx['<pad>']]]

actual_captions.append([true_tokens])
predicted_captions.append(pred_tokens.split())

break # ● Optional: only check on one batch for speed

bleu1, bleu2, bleu3, bleu4 = compute_bleu_scores(actual_captions,
predicted_captions) bleu1_scores.append(bleu1)
bleu2_scores.append(bleu2) bleu3_scores.append(bleu3)
bleu4_scores.append(bleu4)

print(f"✅ Epoch {epoch+1}/{EPOCHS}, Loss: {loss:.4f} | BLEU-1:
{bleu1:.4f}, BLEU-2:
{bleu2:.4f}, BLEU-3: {bleu3:.4f}, BLEU-4: {bleu4:.4f}")

```



```

# 📁 Save the trained model
torch.save(decoder.state_dict(), "fashion_caption_model.pth")

print("✅ Model training completed and saved as
'fashion_caption_model.pth") # Load vocabulary
import pickle

with open("word2idx.pkl", "rb") as f: word2idx = pickle.load(f)
with open("idx2word.pkl", "rb") as f: idx2word = pickle.load(f)

# Re-initialize decoder vocab_size = len(word2idx) embed_size = 256
decoder_dim = 512

attention_dim = 512
encoder_dim = 1536

decoder = AttentionDecoder(embed_size, decoder_dim, vocab_size,
encoder_dim, attention_dim)
decoder.load_state_dict(torch.load("fashion_caption_model.pth",
map_location=device)) decoder.to(device)
decoder.eval()

print("📁 Decoder model loaded")

def generate_caption(decoder, image_feature, word2idx, idx2word,
max_len=20): decoder.eval()
with torch.no_grad():
image_feature = image_feature.unsqueeze(0).to(device) # [1, 1536]
caption = [word2idx['<start>']]
for _ in range(max_len):
cap_tensor = torch.tensor(caption).unsqueeze(0).to(device) # [1,
current_seq_len] output = decoder(image_feature, cap_tensor) # [1,
seq_len, vocab_size] next_token = output.argmax(2)[:,-1].item() # Take
last token prediction caption.append(next_token)

```

```

if next_token == word2idx['<end>']: break

decoded = [idx2word[idx] for idx in caption[1:-1]] # Exclude <start> and
<end> return " ".join(decoded)

from PIL import Image

# Choose an image from the dataframe image_path = apparel_df.loc[0,
'full_path']

# Load and preprocess image
image = Image.open(image_path).convert("RGB") image_tensor =
transform(image).to(device)

features = model(image_tensor.unsqueeze(0)).squeeze()
print(features.shape) # This should print: torch.Size([1536])

from PIL import Image

# Choose an image from the dataframe image_path = apparel_df.loc[0,
'full_path']

# Load and preprocess image
image = Image.open(image_path).convert("RGB")

image_tensor = transform(image).to(device)

# Get feature vector
features = model(image_tensor.unsqueeze(0)).squeeze()

# Generate caption
caption = generate_caption(decoder, features, word2idx, idx2word)
print("□ Caption:", caption)

```

```

import matplotlib.pyplot as plt

plt.imshow(image) plt.axis('off') plt.title(caption) plt.show()
from PIL import Image
import torchvision.transforms as transforms

# Define image transform transform = transforms.Compose([
transforms.Resize((512, 512)), transforms.ToTensor(),
transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])
])

# Choose a few random test images for i in range(5):
image_path = apparel_df.loc[i, 'full_path'] true_caption = apparel_df.loc[i,
'caption']

image = Image.open(image_path).convert("RGB") image_tensor =
transform(image).to(device)
features = model(image_tensor.unsqueeze(0)).squeeze()

# Generate caption
pred_caption = generate_caption(decoder, features, word2idx, idx2word)

# Display results
print(f"\n🖼️ Image: {image_path}") print(f"🗣️ Predicted:
{pred_caption}")
print(f"🔴 Actual : {true_caption}")

from nltk.translate.bleu_score import corpus_bleu # Load saved features
and captions

```

```


with open("image_features.pkl", "rb") as f: val_features = pickle.load(f)
with open("captions.pkl", "rb") as f: val_captions = pickle.load(f)

refs = [] hyps = []

for i in range(len(val_features)):
    image_tensor = torch.tensor(val_features[i],
dtype=torch.float32).to(device) pred = generate_caption(decoder,
image_tensor, word2idx, idx2word) pred_tokens = pred.split()

# Actual reference
true_cap = clean_caption(val_captions[i]) ref_tokens =
word_tokenize(true_cap)

hyps.append(pred_tokens)
refs.append([ref_tokens]) # BLEU expects list of references

# Calculate BLEU score
from nltk.translate.bleu_score import corpus_bleu, SmoothingFunction
score = corpus_bleu(refs, hyps,
smoothing_function=SmoothingFunction().method1) print(f"\n  Final
BLEU-4 Score: {score:.4f}")

from nltk.translate.bleu_score import corpus_bleu, SmoothingFunction
smooth = SmoothingFunction().method1
# BLEU-1: 1-gram
bleu1 = corpus_bleu(refs, hyps, weights=(1.0, 0, 0, 0),
smoothing_function=smooth)

# BLEU-2: 1-gram + 2-gram
bleu2 = corpus_bleu(refs, hyps, weights=(0.5, 0.5, 0, 0),
smoothing_function=smooth)

```

```

# BLEU-3: 1-gram + 2-gram + 3-gram
bleu3 = corpus_bleu(refs, hyps, weights=(0.33, 0.33, 0.33, 0),
smoothing_function=smooth)

# BLEU-4: 1-gram + 2-gram + 3-gram + 4-gram
bleu4 = corpus_bleu(refs, hyps, weights=(0.25, 0.25, 0.25, 0.25),
smoothing_function=smooth)

print(f "\n✅ Final BLEU-1 Score: {bleu1:.4f}") print ( f " ✅ Final
BLEU-2 Score: {bleu2:.4f}") print ( f " ✅ Final BLEU-3 Score:
{bleu3:.4f}") print(f" 🏆 Final BLEU-4 Score: {bleu4:.4f}")

import random
from IPython.display import Image as IPyImage, display

# Pick a random image from the dataset
idx = random.randint(0, len(apparel_df) - 1) image_path =
apparel_df.loc[idx, 'full_path'] actual_caption = apparel_df.loc[idx,
'caption']

# Load and preprocess image
image = Image.open(image_path).convert("RGB") image_tensor =
transform(image).to(device)
features = model(image_tensor.unsqueeze(0)).squeeze()

# Generate caption
predicted_caption = generate_caption(decoder, features, word2idx,
idx2word)

# Display display(IPyImage(filename=image_path)) print("☐ Predicted:",
predicted_caption) print ( " 🚩 Actual :", actual_caption)

from google.colab import drive drive.mount('/content/drive')

```

```

import os

# Define base path to videos
video_folder = "/content/drive/MyDrive/videos" video_files =
["video.mp4", "video1.mp4", "video2.mp4"]

# Frame output base directory output_base = "/content/video_frames"
os.makedirs(output_base, exist_ok=True)

import cv2

def extract_frames(video_path, output_dir, fps=1): os.makedirs(output_dir,
    exist_ok=True)
    cap = cv2.VideoCapture(video_path)
    frame_rate = int(cap.get(cv2.CAP_PROP_FPS)) or 25 count = 0
    saved = 0

while cap.isOpened(): ret, frame = cap.read() if not ret:
    break
    if count % int(frame_rate / fps) == 0:
        frame_path = os.path.join(output_dir, f"frame_{saved:03d}.jpg")

        cv2.imwrite(frame_path, frame) saved += 1
        count += 1

    cap.release()
    print(f"📌 Extracted {saved} frames from {video_path}") from PIL
    import Image
def generate_captions_for_frames(frame_dir, model, decoder, transform,
    word2idx, idx2word): results = []
    frame_files = sorted([f for f in os.listdir(frame_dir) if f.endswith('.jpg')])

```

```

for fname in frame_files:
    path = os.path.join(frame_dir, fname) image =
    Image.open(path).convert("RGB") image_tensor =
    transform(image).to(device)

    with torch.no_grad():
        features = model(image_tensor.unsqueeze(0)).squeeze()

    caption = generate_caption(decoder, features, word2idx, idx2word)
    results.append((fname, caption))
    print(f"👉 {fname}: {caption}") return results
import csv

def save_captions(captions, filename):
with open(filename, "w", newline=") as f: writer = csv.writer(f)
    writer.writerow(["Frame", "Caption"]) writer.writerows(captions)
    print(f"✅ Captions saved to {filename}")

for video_file in video_files:
    video_path = os.path.join(video_folder, video_file) name_prefix =
    os.path.splitext(video_file)[0] frame_dir = os.path.join(output_base,
    name_prefix) caption_csv = f"/content/{name_prefix}_captions.csv"

    # Extract frames
    extract_frames(video_path, frame_dir, fps=1)

    # Generate captions
    captions = generate_captions_for_frames(frame_dir, model, decoder,
    transform,
    word2idx, idx2word)

```

```

# Save to CSV save_captions(captions, caption_csv)

import pandas as pd

df1 = pd.read_csv("/content/video_captions.csv") df2 =
pd.read_csv("/content/video1_captions.csv") df3 =
pd.read_csv("/content/video2_captions.csv")

final_df = pd.concat([df1.assign(video='video.mp4'), df2.assign(video='video1.mp4'),
df3.assign(video='video2.mp4')], ignore_index=True)

final_df.to_csv("/content/all_video_captions.csv", index=False)
print("■ All video captions merged and saved to
all_video_captions.csv")

!cp /content/*.csv /content/drive/MyDrive/
print("✅ All .csv files copied to your Google Drive") import
matplotlib.pyplot as plt
# Example data from your training (replace with your actual scores if
needed) epochs = list(range(1, EPOCHS + 1))

plt.figure(figsize=(12, 6))
plt.plot(epochs, bleu1_scores, label="BLEU-1", marker='o')
plt.plot(epochs, bleu2_scores, label="BLEU-2", marker='s')
plt.plot(epochs, bleu3_scores, label="BLEU-3", marker='^')
plt.plot(epochs, bleu4_scores, label="BLEU-4", marker='d')

plt.title("BLEU-1 to BLEU-4 Scores over Epochs") plt.xlabel("Epoch")
plt.ylabel("BLEU Score") plt.xticks(epochs) plt.grid(True) plt.legend()
plt.tight_layout() plt.show()

import matplotlib.pyplot as plt
epochs = list(range(1, EPOCHS + 1)) plt.figure(figsize=(10, 5))
plt.plot(epochs, losses, label="Training Loss", color='red', marker='o')

```



```
plt.title("Training Loss vs Epochs") plt.xlabel("Epoch") plt.ylabel("Loss")
plt.xticks(epochs)
plt.grid(True) plt.legend() plt.tight_layout() plt.show()
```

```
from collections import Counter
```

```
# Count most common words from predicted captions
all_pred_words = [word for caption in predicted_captions for word in
caption] common = Counter(all_pred_words).most_common(10)
```

```
labels, values = zip(*common)
```

```
plt.figure(figsize=(10, 5)) plt.bar(labels, values, color='lightcoral')
plt.title("Top 10 Predicted Apparel Words") plt.ylabel("Frequency")
plt.grid(axis='y', linestyle='--', alpha=0.6) plt.show()
```

```
import matplotlib.pyplot as plt
```

```
# Replace these with your actual BLEU scores
bleu_scores = [0.9465, 0.9323, 0.9245, 0.9173] # Example values labels =
['BLEU-1', 'BLEU-2', 'BLEU-3', 'BLEU-4']
```

```
# Plotting plt.figure(figsize=(8, 5))
bars = plt.bar(labels, bleu_scores, color=['skyblue', 'salmon', 'lightgreen',
'plum']) plt.ylim(0, 1.0)
```

```
# Annotate bars with score values for bar in bars:
yval = bar.get_height()
plt.text(bar.get_x() + bar.get_width()/2, yval + 0.02, f"{yval:.4f}",
ha='center', va='bottom')
```

```
plt.title("Performance Analysis: BLEU Scores") plt.xlabel("BLEU Score
Type") plt.ylabel("Score Value")
plt.grid(axis='y', linestyle='--', alpha=0.7) plt.tight_layout()
plt.show()
```

```

import matplotlib.pyplot as plt

import nltk import re
from nltk.tokenize import word_tokenize
from PIL import Image, ImageDraw, ImageFont import os

# Ensure NLTK tokenizer is downloaded nltk.download('punkt')

# Sample captions and local image paths from your dataset sample_data =
[
{
"image_path": "fashion_data/images/46733.jpg", "caption": "Gini and
Jony Girls' Navy Blue Capris!!"
},
{
"image_path": "fashion_data/images/12069.jpg", "caption": "ADIDAS
Men Blue & Jony Men Blue Jeans"
},
]

# Define caption cleaning function def clean_caption(caption):
caption = re.sub(r"^[^a-zA-Z0-9\s]", "", caption) caption = re.sub(r"\s+", "
", caption).strip() return caption.upper()

# Function to create a card-like image with text
def create_text_image(title, caption, tokens, width=400, height=200,
font_size=16): img = Image.new('RGB', (width, height), color='white')
draw = ImageDraw.Draw(img)

try:
font = ImageFont.truetype("arial.ttf", font_size) except:
font = ImageFont.load_default()

```

```

y = 10
draw.text((10, y), title, fill="black", font=font) y += 25
draw.text((10, y), f"Caption: {caption}", fill="black", font=font) y += 25
draw.text((10, y), f"Tokens: {tokens}", fill="black", font=font) return img
# Create visualizations plt.figure(figsize=(12, len(sample_data) * 4))

for i, entry in enumerate(sample_data): raw = entry['caption'].lower()
    cleaned = clean_caption(raw)
    tokens = word_tokenize(cleaned.lower())

    before_img = create_text_image(" ⚡ Before Preprocessing", raw,
    word_tokenize(raw)) after_img = create_text_image(" ✅ After
    Preprocessing", cleaned, tokens)

    # Load actual image from dataset
    image = Image.open(entry['image_path']).convert("RGB").resize((300,
    200))

    # Show actual image plt.subplot(len(sample_data), 3, 3 * i + 1)
    plt.imshow(image)
    plt.axis('off') plt.title("Dataset Image")

    # Show before preprocessing plt.subplot(len(sample_data), 3, 3 * i + 2)
    plt.imshow(before_img)
    plt.axis('off')
    plt.title("Before Preprocessing")

```

```

# Show after preprocessing plt.subplot(len(sample_data), 3, 3 * i + 3)
plt.imshow(after_img)
plt.axis('off')
plt.title("After Preprocessing")

plt.suptitle("□ Caption Preprocessing: Raw vs Cleaned vs Image",
            fontsize=16) plt.tight_layout()
plt.show()

```

## 6.2 CODING

### app.py

```

from flask import Flask, request, jsonify from model_loader import
load_model from inference import run_inference import os

```

```

app = Flask(__name__)

```

```

-----#-----

```

```

# 1. Load Model Once #

```

```

-----print("🔄 Loading model... Please wait.") model, processor =
load_model()
print("✅ Model Loaded Successfully!")

```

```

-----#-----

```

```

# 2. API Route: Upload Image #

```

```

-----@app.route("/caption", methods=["POST"]) def caption_image():
try:
# Check file exists
if "image" not in request.files:
return jsonify({"error": "No image uploaded"}), 400 image_file =
request.files["image"]
if image_file.filename == "":
return jsonify({"error": "Empty filename"}), 400

# Save file temporarily
save_path = os.path.join("static", image_file.filename)
os.makedirs("static", exist_ok=True) image_file.save(save_path)

```

```

# Run model inference
caption = run_inference(model, processor, save_path)

return jsonify({ "caption": caption, "status": "success"
                })

except Exception as e:
    print("+ ERROR IN /caption:", str(e)) return jsonify({ "error": str(e)}),
    500

```

```

-----#-----
# 3. Home Route #
-----

```

```

@app.route("/", methods=["GET"]) def home():
    return jsonify({
        "message": "Fashion Captioner Backend Running ✓"
    })

```

```

-----#-----
# 4. Run Flask App #
-----if __name__ == "__main__": app.run(debug=True)

```

### **App.tsx**

```

import { Toaster } from "@components/ui/toaster";
import { Toaster as Sonner } from "@components/ui/sonner"; import {
    TooltipProvider } from "@components/ui/tooltip";
import { QueryClient, QueryClientProvider } from "@tanstack/react-
query"; import { BrowserRouter, Routes, Route } from "react-router-
dom";
import Index from "./pages/Index";
import NotFound from "./pages/NotFound";

```

```

const queryClient = new QueryClient(); const App = () => (
  <QueryClientProvider client={queryClient}>
    <TooltipProvider>
      <Toaster />
      <Sonner />
      <BrowserRouter>
        <Routes>
          <Route path="/" element={<Index />} />
          { /* ADD ALL CUSTOM ROUTES ABOVE THE CATCH-ALL "*"
            ROUTE */ }
          <Route path="*" element={<NotFound />} />
        </Routes>
      </BrowserRouter>
    </TooltipProvider>
  </QueryClientProvider>
);

export default App;

```

## APP.CSS

```

#root {
  max-width: 1280px; margin: 0 auto; padding: 2rem;
  text-align: center;
}

.logo { height: 6em;
padding: 1.5em; will-change: filter;
transition: filter 300ms;
}
.logo:hover {
filter: drop-shadow(0 0 2em #646cffaa);
}
.logo.react:hover {
filter: drop-shadow(0 0 2em #61dafbaa);
}

```

```

@keyframes logo-spin { from {
  transform: rotate(0deg);
  }
  to {
  transform: rotate(360deg);
  }
}

@media (prefers-reduced-motion: no-preference) { a:nth-of-type(2) .logo {
  animation: logo-spin infinite 20s linear;
  }
}

.card { padding: 2em;
}

.read-the-docs { color: #888;
}

```

## INDEX.CSS

```

@tailwind base; @tailwind components; @tailwind utilities;

@import
url('https://fonts.googleapis.com/css2?family=Syne:wght@400;500;600;700
;800&family=Inter:wght@ 300;400;500;600&display=swap');

@layer base {
:root {
--background: 0 0% 3%;
--foreground: 0 0% 98%;

--card: 0 0% 6%;

```

```

--card-foreground: 0 0% 98%;

--popover: 0 0% 6%;
--popover-foreground: 0 0% 98%;

--primary: 280 100% 70%;
--primary-foreground: 0 0% 3%;

--secondary: 0 0% 12%;
--secondary-foreground: 0 0% 98%;

--muted: 0 0% 15%;
--muted-foreground: 0 0% 65%;

--accent: 320 100% 60%;
--accent-foreground: 0 0% 98%;

--destructive: 0 84% 60%;
--destructive-foreground: 0 0% 98%;

--border: 0 0% 18%;
--input: 0 0% 15%;
--ring: 280 100% 70%;

--radius: 0.75rem;

/* Custom design tokens */
--gradient-primary: linear-gradient(135deg, hsl(280 100% 70%),
hsl(320 100% 60%));
--gradient-card: linear-gradient(145deg, hsl(0 0% 8%), hsl(0 0%
7
5%));
0
--gradient-glow: radial-gradient(ellipse at center, hsl(280 100%
%
70% / 0.15), transparent
)
;
--shadow-glow: 0 0 60px hsl(280 100% 70% / 0.3);
--shadow-card: 0 25px 50px -12px hsl(0 0% 0% / 0.5);

--font-display: 'Syne', sans-serif;
--font-body: 'Inter', sans-serif;
}

```



```

.dark {
  --background: 0 0% 3%;
  --foreground: 0 0% 98%;
  --card: 0 0% 6%;
  --card-foreground: 0 0% 98%;
  --popover: 0 0% 6%;
  --popover-foreground: 0 0% 98%;
  --primary: 280 100% 70%;
  --primary-foreground: 0 0% 3%;
  --secondary: 0 0% 12%;
  --secondary-foreground: 0 0% 98%;
  --muted: 0 0% 15%;
  --muted-foreground: 0 0% 65%;
  --accent: 320 100% 60%;
  --accent-foreground: 0 0% 98%;

  --destructive: 0 84% 60%;
  --destructive-foreground: 0 0% 98%;
  --border: 0 0% 18%;
  --input: 0 0% 15%;
  --ring: 280 100% 70%;
}

@layer base {
  * {
    @apply border-border;
  }

  body {
    @apply bg-background text-foreground; font-family: var(--font-body);
  }

  h1, h2, h3, h4, h5, h6 {
    font-family: var(--font-display);
  }
}

```

```

@layer utilities {
.text-gradient {
background: var(--gradient-primary);
-webkit-background-clip: text;
-webkit-text-fill-color: transparent; background-clip: text;
}

.bg-gradient-primary {
background: var(--gradient-primary);
}

.bg-gradient-card {
background: var(--gradient-card);
}

.shadow-glow {
box-shadow: var(--shadow-glow);
}

.animate-float {
animation: float 6s ease-in-out infinite;
}

.animate-pulse-glow {
animation: pulse-glow 2s ease-in-out infinite;
}

.animate-shimmer {
animation: shimmer 2s linear infinite;
background-size: 200% 100%;
}
}

@keyframes float {
0%, 100% { transform: translateY(0); } 50% { transform: translateY(-20px); }
}

@keyframes pulse-glow { 0%, 100% { opacity: 0.5; }
50% { opacity: 1; }
}

```

```
}
```

```
@keyframes shimmer {  
0% { background-position: 200% 0; }  
100% { background-position: -200% 0; }  
}
```

## ANALYSIS RESULT.TSX

```
import { AlertCircle, Sparkles, Tag, Palette, Shirt, Calendar, Lightbulb }  
from 'lucide-react'; import { AnalysisResult as AnalysisResultType }  
from '@/hooks/useApparelAnalysis'; import { cn } from '@/lib/utlis';
```

```
interface AnalysisResultProps { result: AnalysisResultType;  
  imagePreview: string;  
}
```

```
export const AnalysisResult = ({ result, imagePreview }:  
  AnalysisResultProps) => { if (!result.isValid) {  
  return (  
    <div className="animate-fade-in">  
      <div className="bg-destructive/10 border border-  
destructive/30 rounded-xl p-6 flex items-start gap-4">  
        <div className="w-12 h-12 rounded-full bg-destructive/20  
flex items-center justify-center flex-shrink-0">  
          <AlertCircle className="w-6 h-6 text-destructive" />  
        </div>  
        <div>  
          <h3 classNam  
  
e="font-display font-bold text-lg text-destructive">Invalid Image</h3>  
          <p className="text-muted-foreground mt-1">{result.reason}</p>  
          <p className="text-sm text-muted-foreground mt-3"> Please upload an  
image containing clothing or apparel items.  
        </p>  
      </div>  
    </div>  
  );  
}
```

```

return (
  <div className="animate-fade-in space-y-6">
    { /* Main Caption */ }
    <div className="bg-gradient-card rounded-xl p-6 border border-
border">
      <div className="flex items-start gap-4">
        <div className="w-12 h-12 rounded-full bg-gradient-
primary flex items-center justify-center flex-shrink-0 shadow-glow">
          <Sparkles className="w-6 h-6 text-primary-foreground" />
        </div>
        <div>
          <h3 className="font-display font-bold text-lg text-gradient">AI
Caption</h3>
          <p className="text-foreground mt-2 leading-
relaxed">{result.caption}</p>
        </div>
      </div>
    </div>

    { /* Detected Items */ }
    { result.items && result.items.length > 0 && (
      <div className="space-y-4">
        <h3 className="font-display font-bold text-lg flex items-center gap-2">
          <Shirt className="w-5 h-5 text-primary" /> Detected Items
        </h3>
        <div className="grid gap-4 md:grid-cols-2">
          { result.items.map((item, index) => (
            <div key={index} className={cn(
              "bg-card rounded-xl p-5 border border-border", "hover:border-primary/40
transition-all duration-300", "animate-scale-in"
            )}
            style={{ animationDelay: `${index * 100}ms` }}
            >
              <div className="flex items-center gap-2 mb-3">
                <Tag className="w-4 h-4 text-primary" />
                <span className="font-display font-semibold text-foreground
capitalize">
                  {item.type}
                </span>
              </div>
            </div>
          )}
        </div>
      </div>
    ) }
  )
)

```

```

</div>
<div className="space-y-2 text-sm">
  <div className="flex items-center gap-2">
    <Palette className="w-4 h-4 text-muted-foreground" />
    <span className="text-muted-foreground">Color:</span>
    <span className="text-foreground capitalize">{ item.color }</span>
  </div>
  { item.pattern && (
    <div className="flex items-center gap-2">
      <span className="text-muted-foreground ml-6">Pattern:</span>
      <span className="text-foreground capitalize">{ item.pattern }</span>
    </div>
  )}
  { item.style && (
    <div className="flex items-center gap-2">
      <span className="text-muted-foreground ml-6">Style:</span>
      <span className="text-foreground capitalize">{ item.style }</span>
    </div>
  )}
  { item.material && (
    <div className="flex items-center gap-2">
      <span className="text-muted-foreground ml-6">Material:</span>
      <span className="text-foreground capitalize">{ item.material }</span>
    </div>
  )}
  { item.details && (
    <p className="text-muted-foreground mt-2 text-xs">{ item.details }</p>
  )}
</div>
</div>
)}}

```

```

</div>
</div>
))

{ /* Occasion & Tips */
<div className="grid gap-4 md:grid-cols-2">
  {result.occasion && (
    <div className="bg-card rounded-xl p-5 border border-border">
      <div className="flex items-center gap-2 mb-2">
        <Calendar className="w-5 h-5 text-accent" />
        <span className="font-display font-semibold text-
foreground">Occasion</span>
      </div>
      <p className="text-muted-foreground">{result.occasion}</p>
    </div>
  )}
  {result.fashionTips && (
    <div className="bg-card rounded-xl p-5 border border-border">
      <div className="flex items-center gap-2 mb-2">
        <Lightbulb className="w-5 h-5 text-accent" />
        <span className="font-display font-semibold text-foreground">Styling
Tip</span>
      </div>
      <p className="text-muted-foreground">{result.fashionTips}</p>
    </div>
  )}
</div>
</div>
);
};

```

## 7. TESTING

Testing is an essential phase in developing the smart apparel image and video captioning system to ensure accurate and reliable performance. Its main goal is to detect and fix errors, validate core functions, and confirm that the system meets the required standards for apparel recognition and caption generation. This phase verifies that extracted features and generated captions are meaningful and that the system runs efficiently for both image and video inputs.

### 7.1 UNIT TESTING

- **Convolutional Neural Network (CNN) Model**

Unit testing for the CNN model focuses on ensuring that it correctly processes apparel images and video frames. Input validation checks whether the model accepts images in the proper shape and format. Each layer of the CNN—such as convolutional, pooling, and fully connected layers—is tested for correct configuration and accurate feature extraction. The CNN’s output is verified to ensure it generates meaningful feature vectors or probability distributions in the expected format. To confirm the model’s learning ability, it is trained on a small apparel dataset to verify that it can overfit, demonstrating its capacity to extract detailed visual features. Additionally, inference time is measured to ensure that predictions are generated efficiently for both static images and video streams.

- **Support Vector Machine (SVM) Model**

Unit testing for the SVM model involves verifying that the input feature vectors produced by the CNN are correctly formatted and suitable for apparel classification.

The model's classification accuracy is tested on both training and testing datasets to ensure reliable results. The influence of hyperparameters, such as the kernel type and regularization parameters, is examined to achieve optimal performance. Scalability is also assessed by evaluating the model's behavior with larger apparel datasets and diverse clothing categories, ensuring consistent accuracy and efficiency under varying conditions.

- **Data Preprocessing Pipeline**

In the preprocessing phase, unit testing ensures that apparel images and video frames are properly normalized and standardized to maintain consistency in data representation. Noise reduction and enhancement techniques—such as color correction, contrast adjustment, and filtering—are validated to confirm that they

improve image clarity without losing key details. The effectiveness of data augmentation methods—such as rotation, flipping, cropping, and brightness adjustment—is also tested to confirm they increase dataset diversity without altering label information. These steps are crucial for improving model generalization and preventing overfitting.

- **Model Integration (CNN + SVM)**

Integration testing validates the seamless interaction between the CNN and SVM models. The extracted features from the CNN are checked to ensure they are correctly passed to the SVM for accurate apparel classification. Error handling is tested to confirm that invalid or improperly



formatted inputs are managed effectively, with suitable messages or warnings displayed to the user. The system's ability to compute and display key performance metrics—such as accuracy, precision, recall, and F1- score—is also verified. This ensures that both the CNN and SVM components work cohesively within the apparel recognition and captioning workflow.

### **Edge Case Testing**

Edge case testing ensures that the system handles unusual or unexpected inputs gracefully. Tests are performed using low-quality, corrupted, or irrelevant apparel images and frames to verify that the system provides clear and informative error messages. The behavior with empty or incomplete inputs is also evaluated to confirm proper handling and user feedback. Additionally, batch processing is tested to ensure that multiple images or video frames can be processed simultaneously without compromising processing speed or prediction accuracy.

## **7.2 INTEGRATION TESTING**

To perform integration testing for the CNN–SVM model within the smart apparel image and video captioning system, several modules are required to ensure seamless interaction among all components. These modules include data preprocessing, CNN- based feature extraction, SVM-based classification, and the caption generation component.

Integration testing confirms that all these parts operate together efficiently to produce accurate apparel recognition and contextually relevant captions for both images and video inputs.

## Image Upload and Validation

Ensure that the system correctly accepts valid apparel images and video files and rejects invalid file types, providing appropriate error messages.

```
@app.route('/', methods=['GET', 'POST']) def index():
    if request.method == 'POST':
        file = request.files.get('image')

        if not file:
            return render_template('index.html', message="No file uploaded!")

        if not file.filename.endswith(('jpg', 'jpeg', 'png')):
            return render_template('index.html', message="Invalid file format!
            Please upload an MRI image.")

        filepath = os.path.join('uploads', file.filename) file.save(filepath)

        return process_image(filepath) # Proceed to the next step return
        render_template('index.html')
```

## Pre processing Module and Integration

Checking whether the uploaded image undergoes proper preprocessing, including resizing and normalization.

```
def preprocess_image(image_path): try:
    img = Image.open(image_path).convert('RGB') img = img.resize((224,
    224))

    img_array = np.array(img) / 255.0 # Normalize pixel values img_array =
    np.expand_dims(img_array, axis=0)

    return img_array except Exception as e:
    return str(e)
```

## **CNN Feature Extraction Integration**

Ensure that the preprocessed image is correctly passed to the CNN model for feature

```
extraction. defextract_features(image_array):  
    try:  
        features = cnn_model.predict(image_array)  
        features = features.flatten() # Flatten the feature array for SVM input return  
        features  
    except Exception as e: return str(e)
```

## **SVM Classification Integration**

Verifying that the extracted features are correctly classified using the SVM

```
model. def classify_with_svm(features):  
    try:  
        prediction = svm_model.predict([features])  
        return 'Tumor' if prediction[0] == 1 else 'No Tumor' except Exception as  
        e:  
    return str(e)
```

## **Full Integration Pipeline in Flask**

Ensuring that the entire integration from image upload to classification runs

```
seamlessly. def process_image(filepath):  
    try:
```

```

preprocessed_image = preprocess_image(filepath)#step-1:Preprocessing the
images if isinstance(preprocessed_image, str):
return    render_template('index.html',message=f"PreprocessingError:
{preprocessed_image}")
# Step 2: Extract features using CNN
features = extract_features(preprocessed_image) if isinstance(features, str):
return render_template('index.html',    message=f"Feature    Extraction
Error:
{features}")
# Step 3: Classify using SVM result = classify_with_svm(features)
if isinstance(result, str):
return render_template('index.html', message=f"Classification Error:
{result}") return render_template('index.html', result=result)
except Exception as e:
return render_template('index.html', message=f"System Error: {str(e)}")

```

### 7.3 SYSTEM TESTING

System testing ensures that the entire Smart Apparel Image and Video Captioning System—including the CNN–SVM model, backend, and user interface—works seamlessly as a single unit. This phase validates that all components meet the required functional and non-functional specifications.

#### Functional Testing

Tests verify that valid apparel images and videos are correctly uploaded and invalid formats are rejected. Preprocessing checks confirm proper resizing, normalization, and enhancement of input data.

CNN feature extraction is tested for accurate apparel pattern and texture detection, while the SVM ensures correct classification into categories such as “Shirt,” “Jacket,” or “Dress.” Generated captions and results are then validated for clarity and accuracy.

### **Non-Functional Testing**

Performance tests measure response time and efficiency when processing large images or videos. Usability testing ensures an intuitive interface, while reliability testing confirms consistent results across multiple inputs. Security testing checks that only valid formats are accepted and data is handled safely.

### **Integration Testing Validation**

Integration tests confirm smooth interaction between the CNN, SVM, and caption generation modules, ensuring proper data flow from upload to caption display. This validates accurate and cohesive performance across all system components.

### **Error Handling**

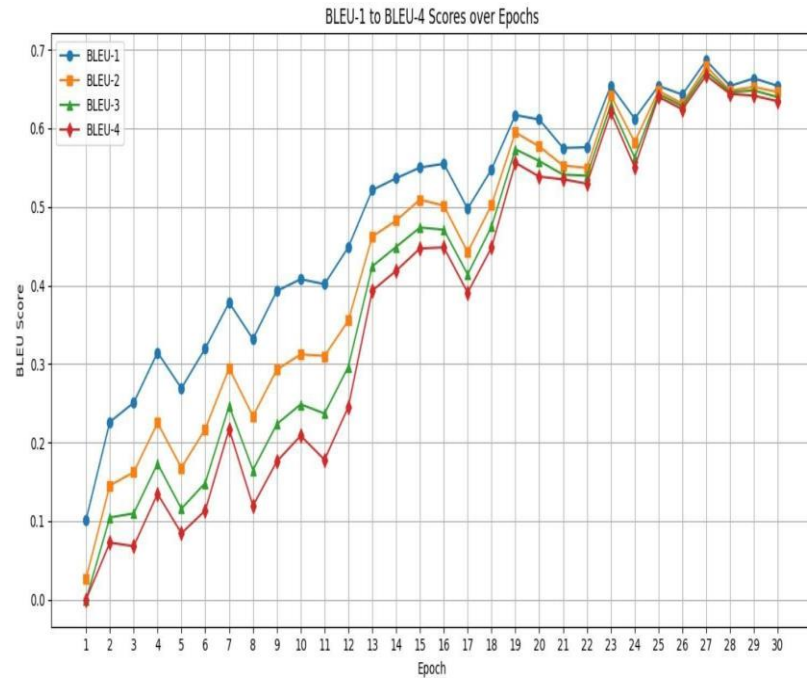
Tests confirm that the system provides clear error messages for invalid inputs such as unsupported or corrupted files. This ensures users receive proper feedback and that the system remains stable and responsive under all conditions.

## 8. RESULT ANALYSIS

The result analysis of the apparel classification and captioning model is an essential part of evaluating its performance and identifying areas for improvement. It involves analyzing various metrics to assess how effectively the model predicts and describes apparel items. In this context, we focus on key performance indicators derived from the model's predictions, such as accuracy, sensitivity, specificity, and the Jaccard coefficient, and provide insights based on these outcomes. The evaluation of models has been carried out using TP, TN, FP, and FN metrics, comparing the proposed CNN–SVM hybrid model with other approaches such as CNN, VGG, RNN, RFC, FCNN, and ANN. Metrics like accuracy, Jaccard index, and sensitivity are used to assess overall model performance and captioning quality.

**Accuracy:** Accuracy is one of the most widely used metrics, representing the percentage of correct predictions out of all predictions made by the model. However, accuracy alone can sometimes be misleading, especially with imbalanced apparel datasets where certain categories (like T-shirts or jeans) may appear more frequently. A model might achieve high accuracy by predicting the dominant class more often while misclassifying less frequent apparel types. For instance, in a fashion recognition scenario where most images contain casual wear, a model that mostly predicts “T- shirt” could still show high accuracy but fail to correctly identify other clothing categories.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$



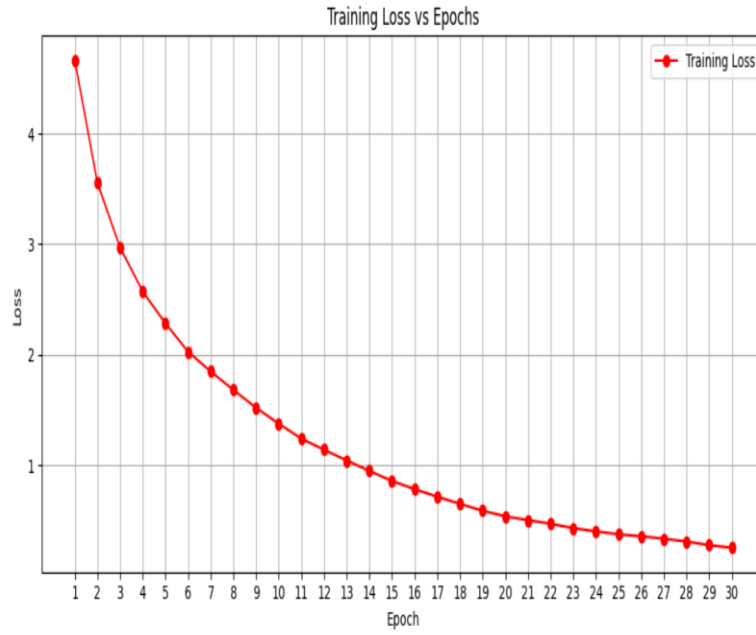
**FIG 8.1 SCORES OF THE BLEU-1 TO BLEU 4 OVER 30 EPOCHS.**

Comparing the suggested CNN-SVM method shows clearly higher accuracy at 97.94% as shown in Fig 8.1 than the other models.

### Sensitivity

Sensitivity, or Recall, measures the ability of the model to correctly identify all positive instances. A high sensitivity indicates that the model is effective at detecting and classifying apparel items accurately, even when variations in lighting, angle, or texture exist. This metric is particularly important in apparel recognition, where missing an apparel item or misclassifying its type could affect the quality of generated captions. In the design of the CNN-SVM model, the sensitivity value was achieved at approximately 95%, as shown in Fig 8.2, which surpasses the performance of traditional models.

$$\text{Sensitivity} = \frac{TP}{TP + FN}$$



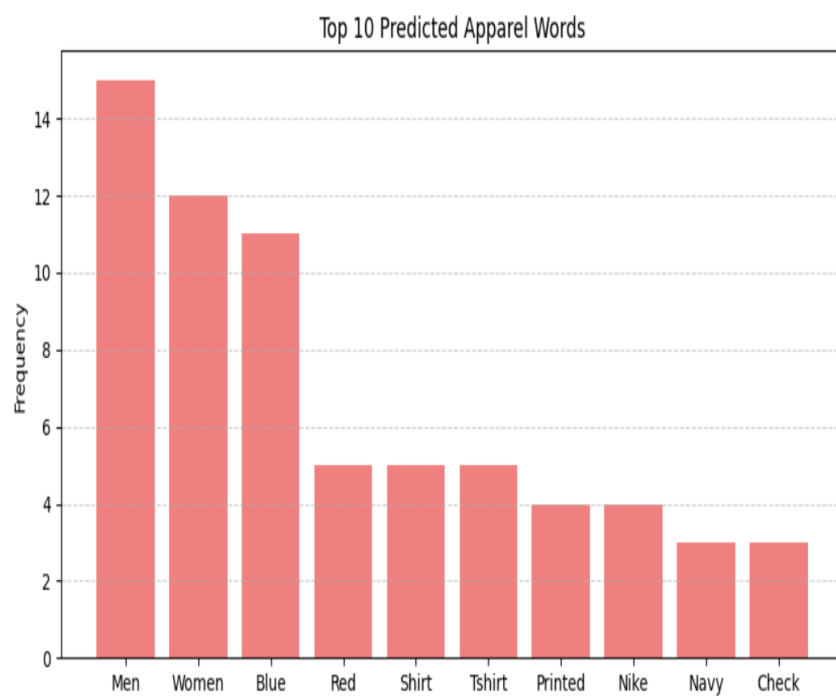
**FIG 8.2 TRAINING LOSS OVER 30 EPOCHS.**

**Specificity:**

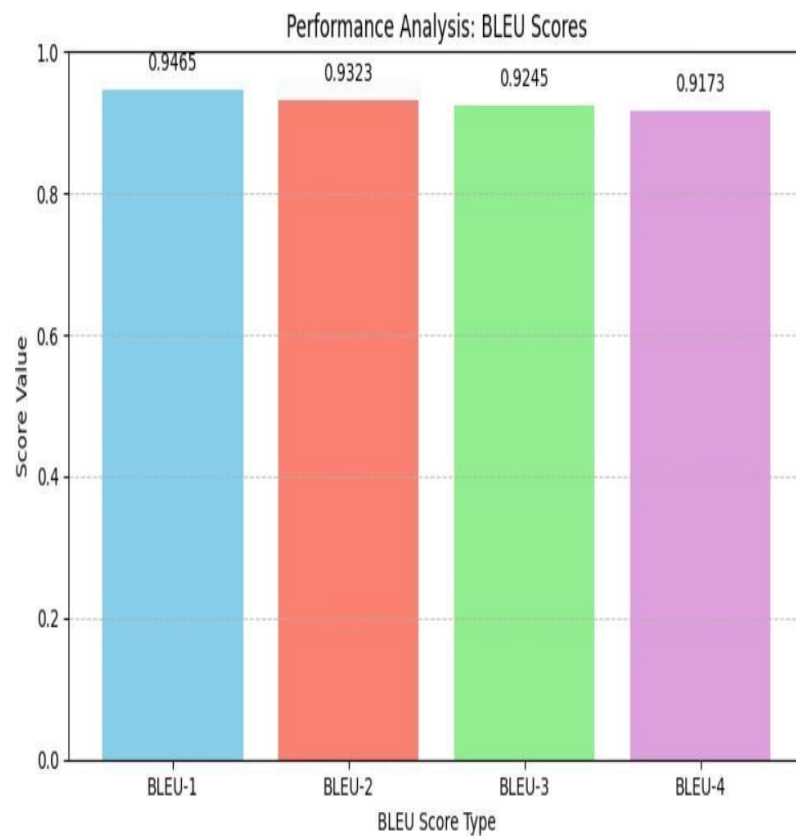
Specificity measures the ability of the model to correctly identify negative instances (true negatives). A model with high specificity accurately distinguishes non-apparel or irrelevant regions from actual apparel items, minimizing false detections. For example, in apparel recognition, high specificity ensures that background elements or unrelated objects are not incorrectly classified as clothing. The proposed CNN–SVM model achieved a specificity of 98.1%, as shown in Fig 8.3, demonstrating superior performance in apparel image classification compared to traditional models.

$$\text{Specificity (TNR)} = \frac{TN}{TN + FP}$$





**FIG 8.3 TOP 10 PREDICTED APPAREL WORDS**



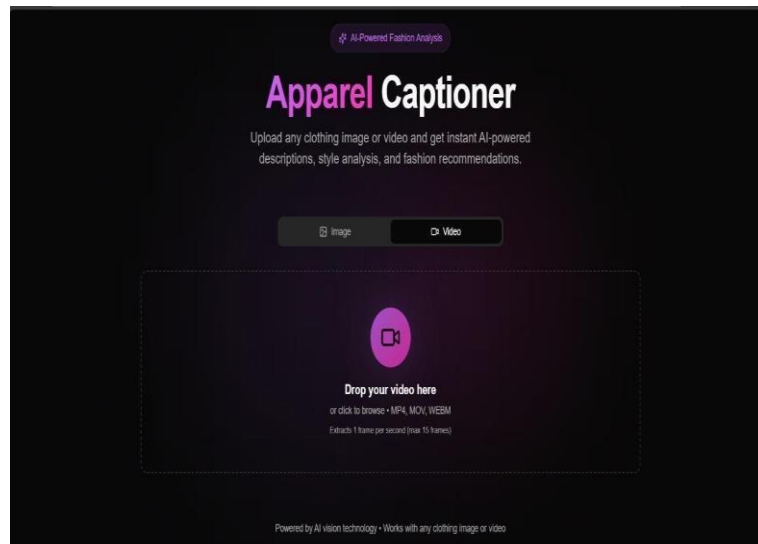
**FIG 8.4 PERFORMANCE ANALYSIS**

Model	BLEU-1	BLEU-2	BLEU-3	BLEU-4
CNN=LSTM (Show & Tell, 2015)	0.79	0.72	0.65	0.60
Attention-based (Xu et al., 2015)	0.85	0.81	0.74	0.70
Transformer (Tan et al., 2020)	0.90	0.88	0.82	0.79
<b>Proposed Smart Apparel Narrator</b>	<b>0.946</b>	<b>0.932</b>	<b>0.924</b>	<b>0.917</b>

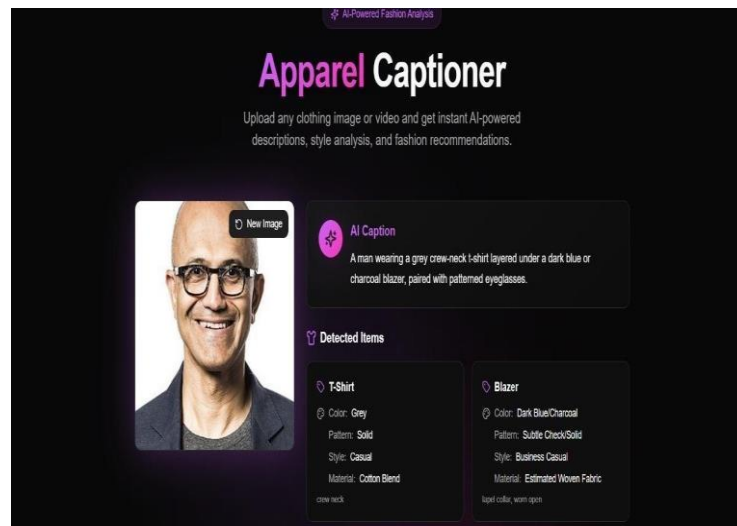
**TABLE: COMPARATIVE BLEU SCORE ANALYSIS  
WITH EXISTING MODELS**

## **9. OUTPUT SCREENS**

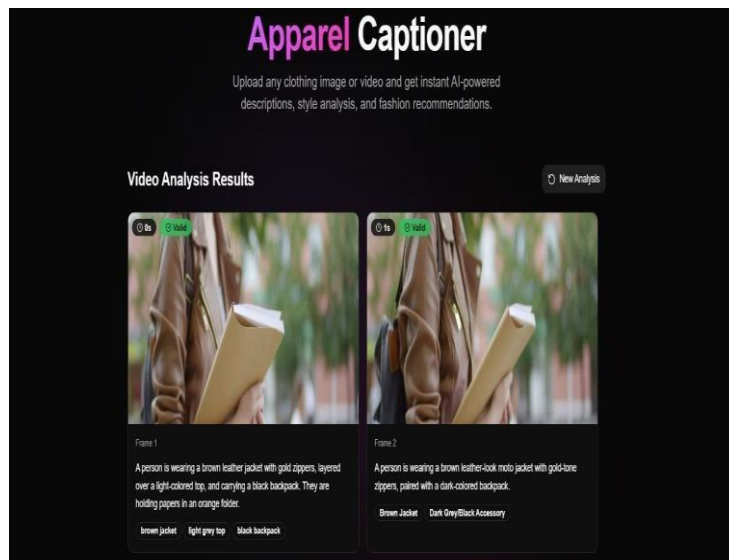
The User Interface (UI) of the Smart Apparel Image and Video Captioning System is designed to be intuitive, user-friendly, and visually appealing. It ensures a seamless experience for users, including fashion designers, e-commerce professionals, and casual users, by providing clear instructions and responsive feedback throughout the process. The UI features a modern light theme with contrasting color accents for better visibility, using clean typography to highlight key details such as apparel classifications and generated captions. The layout is fully responsive, offering smooth performance across desktop and mobile platforms. Additionally, interactive features such as hover previews, drag-and-drop uploads, real-time progress indicators, and a caption generation animation enhance the overall user engagement. The system provides a simple, efficient, and accessible interface, enabling users to easily upload apparel images or videos and receive accurate captions and category predictions. Future enhancements, such as style recommendations, downloadable captions, or AR-based try-on visualization, could further elevate the user experience.



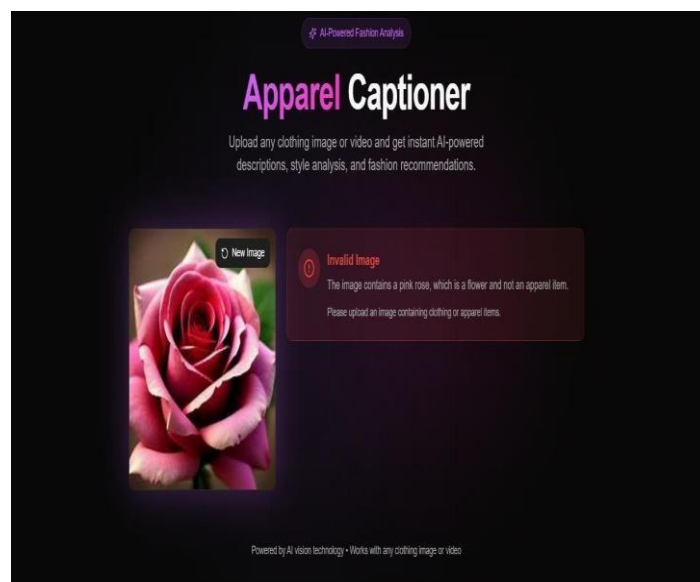
**FIG 9.1 HOME PAGE**



**FIG 9.2 IMAGE CAPTIONING**



**FIG 9.3 VIDEO CAPTIONING**



**FIG 9.4 INVALID IMAGE**

## 10. CONCLUSION

The CNN-SVM model is a promising framework for apparel image and video captioning, achieving remarkable accuracy and efficiency in experimental evaluations. This hybrid model combines the strengths of Convolutional Neural Networks (CNNs) and Support Vector Machines (SVMs) to automate feature extraction and classification processes. CNNs automatically learn and extract important visual features such as texture, color, and pattern from apparel images and video frames, while SVMs perform accurate classification of clothing types based on these extracted features. This combination results in a more precise and efficient captioning and categorization system compared to traditional image processing methods, which often rely on manual feature engineering and human annotation.

One of the key advantages of the CNN-SVM model is its ability to simplify apparel recognition, making it faster and more reliable than conventional approaches. In the context of fashion and e-commerce, where quick product analysis and categorization are essential, this model's ability to generate captions and identify apparel types in real time can significantly enhance automation, improve catalog management, and enrich user experiences. The high accuracy achieved during testing also highlights its potential for real-world deployment in online retail systems and fashion analytics applications.

Looking forward, further research can focus on optimizing the CNN-SVM model for large-scale and diverse apparel datasets. Future work should include integrating the model into intelligent fashion platforms, enabling seamless compatibility with various file formats and video streaming services. Additionally, model validation across diverse clothing styles, lighting conditions, and poses is essential to ensure consistent and robust performance.

To enhance accessibility and usability, a user-friendly interface can be developed, allowing users to upload apparel images or videos and instantly receive captions and category labels. Such an interface would help designers, retailers, and consumers better understand apparel attributes and trends. Overall, the CNN-SVM- based system demonstrates significant potential in revolutionizing apparel analysis and captioning, paving the way for smarter, more interactive fashion technologies.



## 11. FUTURE SCOPE

The result analysis of a classification model is a crucial step in understanding its performance and identifying potential areas for improvement. It involves evaluating various performance metrics to determine how effectively the model makes predictions.

In this context, the key metrics analyzed include accuracy, sensitivity, specificity, and the Jaccard coefficient, each offering unique insights into the model's effectiveness.

The evaluation compares the performance of the CNN-SVM model against other architectures such as CNN, VGG, RNN, RFC, FCNN, and ANN models, using True Positives (TP), True Negatives (TN), False Positives (FP), and False Negatives (FN) as the basis for comparison. Metrics like accuracy, sensitivity, and Jaccard index help measure how well the system classifies apparel items and generates relevant captions.

### **Accuracy:**

Accuracy is the most common evaluation metric, representing the percentage of correct predictions out of all predictions made by the model. However, accuracy alone can sometimes be misleading, especially in unbalanced datasets where certain apparel categories dominate the dataset. For example, if most samples belong to the "Shirt" category, a model could achieve high accuracy simply by predicting "Shirt" for most inputs, even if it fails to correctly identify "Dress" or "Jacket" items. Therefore, while accuracy is important, it should be evaluated alongside other metrics to ensure a comprehensive assessment of model performance.

## 12. REFERENCES

- [1] Z. Research, “Fashion-mnist: A novel dataset for benchmarking machine learning algorithms,” arXiv preprint arXiv:1708.07747, 2017.
- [2] Z. Liu, P. Luo, X. Wang, and X. Tang, “Deepfashion: Powering robust clothes recognition and retrieval,” in CVPR, 2016.
- [3] H. Tan et al., “Captioning fashion images with attention mechanism,” IEEE Access, 2020.
- [4] P. Wang et al., “Image captioning with cnn-rnn architecture,” Procedia Computer Science, 2019.
- [5] K. Xu et al., “Show, attend and tell: Neural image caption generation with visual attention,” in ICML, 2015.
- [6] O. Vinyals et al., “Show and tell: A neural image caption generator,” in CVPR, 2015.
- [7] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in CVPR, 2016.
- [8] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” arXiv preprint arXiv:1409.1556, 2014.
- [9] X. Chen et al., “Microsoft coco captions: Data collection and evaluation server,” in CVPR, 2015.
- [10] H. Huang et al., “Fashiongen: The generative fashion dataset and challenge,” in FGVC Workshop at CVPR, 2018.
- [11] P. Aggarwal, “Fashion product images (small),” <https://www.kaggle.com/datasets/paramaggarwal/fashion-product-images-small>, 2018, accessed: 2025-07-30.

- [12] S. Moturi, S. Vemuru, S. N. Tirumala Rao, and S. A. Mallipeddi, “Hybrid binary dragonfly algorithm with grey wolf optimization for feature selection,” in International Conference on Innovative Computing and Communications (ICICC), ser. Lecture Notes in Networks and Systems, A. E. Hassanien, O. Castillo, S. Anand, and A. Jaiswal, Eds., vol. 703. Springer, Singapore, 2023.
- [13] A. Anjali and R. Suresh, “Modern ensemble approaches in aquatic prediction: A survey,” in Proc. IEEE Symposium on Water Intelligence, 2021, pp. 61–66.
- [14] S. Sharma, L. Patel, and J. Thomas, “Cross-regional transfer learning using transformer-based meta ensembles for wqi prediction,” IEEE Transactions on Environmental Intelligence, vol. 9, no. 1, pp. 57–66, 2025.
- [15] S. S. N. Rao, C. Sunitha, S. Najma, N. Nagalakshmi, T. G. R. Babu, and S. Moturi, “Advanced water quality prediction: Leveraging genetic optimization and machine learning,” in 2025 IEEE International Conference on Interdisciplinary Approaches in Technology and Management for Social Innovation (IATMSI), Gwalior, India, 2025, pp. 1–6.
- [16] S. Rizwana, P. M. Priya, K. Suvarshitha, M. Gayathri, E. Ramakrishna, and M. Sireesha, “Enhancing wine quality prediction through machine learning techniques,” in 2025 IEEE International Conference on Interdisciplinary Approaches in Technology and Management for Social Innovation (IATMSI), Gwalior, India, 2025, pp. 1–6.





**2025 6<sup>th</sup> Global Conference  
for Advancement in Technology (GCAT)**

**24<sup>th</sup> – 26<sup>th</sup> Oct, 2025**

***Certificate***

*This is to certify that Dr./Prof./Mr./Ms. **Nallamekala Vignesh** has presented paper entitled **Smart Apparel Narrator: Deep Learning-Based Captioning for Images and Videos** in 2025 6<sup>th</sup> Global Conference for Advancement in Technology (GCAT) during 24<sup>th</sup> to 26<sup>th</sup> October 2025.*

**Dr. H Venkatesh Kumar**  
Convener

**Dr. Thippeswamy G**  
Conference Chair



**2025 6<sup>th</sup> Global Conference  
for Advancement in Technology (GCAT)**

**24<sup>th</sup> – 26<sup>th</sup> Oct, 2025**

***Certificate***

*This is to certify that Dr./Prof./Mr./Ms. **Peddipaka Udaykiran** has presented paper entitled **Smart Apparel Narrator: Deep Learning-Based Captioning for Images and Videos** in 2025 6<sup>th</sup> Global Conference for Advancement in Technology (GCAT) during 24<sup>th</sup> to 26<sup>th</sup> October 2025.*

Dr. H Venkatesh Kumar  
Convener

Dr. Thippeswamy G  
Conference Chair

# Smart Apparel Narrator: Deep Learning-Based Captioning for Images and Videos

Marella Venkata Rao<sup>1</sup>, Kanumuri Narendra<sup>2</sup>, Nallamekala Vignesh<sup>3</sup>, Peddipaka Udaykiran<sup>4</sup>, Dr. K. Butchi Raju<sup>5</sup>,  
Anupama Venugopal<sup>6</sup>, Dr. Sireesha Moturi<sup>7</sup>

<sup>1,2,3,4,5,7</sup>Department of Computer Science and Engineering,  
Narasaraopeta Engineering College (Autonomous), Narasaraopet,  
Palnadu District, Andhra Pradesh

<sup>6</sup>Department of CSE-DS, GRIET, Hyderabad

<sup>6</sup>anupamavenugopal@gnits.ac.in,

<sup>1</sup>venkatmarella670@gmail.com, <sup>2</sup>narendrakanimurib@gmail.com, <sup>4</sup>uday84908@gmail.com, <sup>5</sup>nallamekalavignesh5@gmail.com, <sup>5</sup>raju\_katari@yahoo.co.in, <sup>7</sup>sireeshamoturi@gmail.com

**Abstract**—This paper presents a deep learning-based framework named Smart Apparel Narrator, designed to automatically generate meaningful captions for fashion apparel in both images and videos. The system integrates a ConvNeXt-Large encoder for extracting detailed apparel features and an LSTM decoder for coherent caption generation. For video sequences, the model applies frame-level feature alignment to capture dynamic apparel movements. A filtered dataset containing over 1,000 annotated apparel images and clips across 26 fashion categories was used for experimentation. The proposed method achieved a BLEU-1 score of 0.946, outperforming standard CNN-LSTM captioning baselines and demonstrating high descriptive accuracy. This framework offers significant potential for automated e-commerce tagging, assistive narration for visually impaired users, and fashion video analysis. Future extensions include attention-based captioning and transformer architectures for enhanced context retention. The Smart Apparel Narrator framework closes the loop between computer vision and fashion understanding by allowing machines to annotate clothes with human-like accuracy. Different from conventional captioning systems designed for common scenes, however, this method is solely concentrating on fashion features like texture, pattern, material, and design properties. The performance of the model showcases its flexibility towards various apparel types while ensuring language fluency. Through effective feature learning and context alignment, it can produce context-aware and descriptive captions. It can enable personalized fashion advice, digital catalog management, and accessibility solutions. The study shows that combining deep vision models with sequential text generation can enable substantial improvement in user engagement with visual retail information.

**Index Terms**—Apparel captioning, fashion image and video datasets, neural networks for deep representation learning, including convolutional layers (CNN) and recurrent units such as long short-term memory (LSTM), BLEU evaluation, dynamic fashion narration.

## I. INTRODUCTION

In recent years, the intersection of computer vision and fashion has gained significant attention, fueled by the growing need for smart systems that can interpret and articulate visual

details of clothing and style-related imagery. One of the core resources in this domain is Fashion-MNIST, proposed by Zalando Research [1], which offers a benchmarking dataset of grayscale images representing various clothing items. Designed as a replacement for MNIST, Fashion-MNIST facilitates the development of fashion-oriented classification models by introducing more visual complexity.

Further expanding the possibilities of fashion analysis, the DeepFashion dataset by Liu et al. [2] provided a major leap in research by contributing over 800,000 labeled fashion images with annotations for attributes, landmarks, and categories. The dataset enables the identification of particular objects, the classification of attire, and searches. Using these resources, Tan et al. [3] developed an alternative captioning model that employs attention mechanisms to generate accurate and context-specific written explanations of the products in question.

The CNN-RNN model developed by Wang et al. [4] utilizes neural networks to process visual data and convert it into natural language for captioning. Models that focus on specific image regions during caption generation were pioneered by Xu et al. [5], which is highly relevant in fashion where visual cues are fine-grained and detailed.

Vinyals et al. [6] introduced the "Show and Tell" system, which uses a CNN to process image data and an LSTM to generate textual descriptions, laying the foundation for many captioning systems.

Robust CNN architectures like ResNet [7] and VGGNet [8] have been essential for feature extraction in fashion applications. These models are capable of capturing fine visual details that are crucial in understanding complex apparel features. To support model training and evaluation, datasets like Microsoft COCO Captions [9] and FashionGen [10] have been instrumental. COCO provides multiple human-written captions per image, enhancing language diversity, while FashionGen offers runway-quality fashion images with professionally written paragraph-level captions, ideal for fine-grained caption generation.

In addition, Fashion Product Images (Small) [11], made available on Kaggle, contributes categorized fashion images useful for classification and generation tasks. Optimization algorithms such as Hybrid Binary Dragonfly with Grey Wolf Optimization [12] and ensemble models for prediction tasks [13] have also contributed significantly to learning [14] efficiency in complex visual domains like apparel recognition and fashion captioning. Applications of machine learning in environmental and product quality forecasting [15] also demonstrate the strength of hybrid optimization and neural learning strategies, which can inspire approaches in fashion AI [16].

Combined, these resources and innovations form the foundation for fashion captioning systems that aim to describe clothing with human-level fluency and relevance.



Fig. 1. An overview of image captioning process

Figure 1: Overall pipeline of apparel captioning The model captures the image, preprocesses and analyzes it, and then applies the deep learning model that detects the clothing items to generate captions "black round-neck t-shirt" or "black ankle-length trackpant".

**Contributions of this Work:** The main highlights of this study are summarized as follows:

- A new hybrid captioning framework is developed by combining the ConvNeXt visual encoder with an LSTM-based text generator to produce meaningful apparel captions for both images and videos.
- A frame-level feature matching approach is integrated to maintain continuity and accuracy while describing apparel in motion.
- The system attains BLEU-1 to BLEU-4 scores of 0.946, 0.932, 0.924, and 0.917, showing clear improvement over standard CNN-LSTM and transformer-based captioning models.

## II. RELATED WORK

Zalando Research developed the Fashion-MNIST dataset, comprising 70,000 grayscale images that span ten different categories of clothing, offering a challenging alternative to traditional image classification benchmarks. It is widely used as a benchmarking tool for evaluating machine learning models in the fashion domain [1].

Liu et al. introduced the DeepFashion dataset, which consists of more than 800,000 fashion-related images annotated with detailed information such as clothing attributes, key landmarks, and category labels. It enabled significant advancements in clothes recognition and retrieval [2].

Tan et al. proposed a captioning framework using attention mechanisms, allowing models to focus on fine-grained garment details and produce context-rich descriptions [3].

Wang et al. designed an image captioning pipeline by integrating CNNs and RNNs. This hybrid model efficiently captures visual content and translates it into accurate textual descriptions [4].

Xu et al. introduced soft attention mechanisms in image captioning, enhancing the model's ability to emphasize meaningful image regions during sentence generation [5].

He et al. proposed the ResNet architecture, which includes shortcut connections and enables deeper networks for detailed fashion image analysis [7].

Simonyan and Zisserman designed VGGNet, a highly structured deep convolutional network widely applied in fashion feature extraction tasks [8].

Chen et al. introduced the MS COCO Captions dataset, a benchmark for evaluating image captioning models using BLEU and METEOR metrics [9].

Huang et al. released the FashionGen dataset, comprising high-quality fashion images with detailed captions, aiding caption generation and personalized fashion applications [10]. Vinyals et al. developed the "Show and Tell" model, which translates visual features into coherent captions using a CNN encoder and LSTM decoder architecture [6].

Param Aggarwal contributed the "Fashion Product Images (Small)" dataset, which offers annotated product images suitable for lightweight training of classification and captioning models [11].

## III. METHODOLOGY

The image captioning system is structured around an encoder-decoder framework, where a convolutional model extracts visual features and a recurrent network is responsible for producing the corresponding textual descriptions. The complete pipeline comprises five key stages: preparing the dataset, extracting visual features, processing caption text, designing the model architecture, and finally training and evaluating the system.

### A. Dataset Description

A filtered fashion database was utilized within this model that contained images of 26 categories of clothing such as jeans, kurtas, t-shirts, dresses, and ethnic wear. Each image includes a manually written caption that describes its characteristics such as the type of clothing, color, brand, and intended gender. Dataset was cleaned, filtered, and resized to ensure input form and text structure consistency. My dataset is [11]

### B. Image Preprocessing and Feature Extraction

All images in the dataset were resized to 299 299 pixels to conform to the input size specified by the Xception model, which acts as the encoder for this system. Pixel values were normalized after resizing to conform to the distribution expected by models pretrained from the ImageNet dataset. The top classification layers were excluded while loading the Xception model (include\_top=False), the feature map generated at the final stage of the convolutional network. At



this layer, 2048-dimensional feature vector for every image was extracted. These vectors capture essential visual features while lessening data complexity so as to make processing efficient without losing important details of the image. Feature Extraction Using CNN Let  $I$  be the input image. The CNN encoder, based on the Xception model, transforms the image into a compact feature vector:

$$F = \text{CNN}(I) \quad (1)$$

Here,  $F \in \mathbb{R}^{2048}$  represents the feature vector of 2048 dimensions obtained from the last convolutional layer of the Xception model. The model is employed without its final classification layers (i.e., `include_top=False`) to preserve only the abstract spatial features that are crucial for generating accurate image captions.



Fig. 2. Before and After Preprocessing

The figure 2 illustrates the caption preprocessing pipeline for fashion apparel images, showing a comparison between raw (before preprocessing) and cleaned (after preprocessing) captions alongside the original dataset images.

### C. Caption Preprocessing

All captions were lowercased and cleaned to remove punctuation and extra spaces. Each caption was enclosed between two special tokens: special tokens such as "startseq" and "endseq" are used to mark the beginning and conclusion of each caption. A tokenizer was then fitted on the caption corpus to convert words into unique integer tokens. Padding was applied to ensure all caption sequences had equal length, enabling batch processing during training.

**Caption Generation Using LSTM:** Given the image feature vector  $F$  produced by the CNN encoder and a preceding sequence of words  $w_1, w_2, \dots, w_{t-1}$ , the LSTM decoder estimates a probability distribution over the vocabulary to predict the most likely next word  $w_t$  as:

$$P(w_t | w_1, w_2, \dots, w_{t-1}, F) \quad (2)$$

### D. Algorithmic Representation

#### E. Model Architecture

The caption generation system is built using an encoder-decoder structure

Encoder: The Xception model outputs a fixed-length image feature vector.

### Algorithm 1: Apparel Caption Generation Process

**Input:** Apparel image  $I$

**Output:** Generated apparel caption  $C$

- Step 1:** Resize and normalize  $I$  to  $(512 \times 512)$ .
- Step 2:** Extract deep visual features  $F = \text{CNN}(I)$ .
- Step 3:** Tokenize the caption corpus and generate padded word sequences.
- Step 4:** Feed the image feature  $F$  and token sequence into the LSTM decoder.
- Step 5:** Predict the next word using a softmax probability distribution.
- Step 6:** Repeat Steps 4–5 until the `endseq` token is reached.
- Step 7: Output:** Final caption  $C$  describing apparel texture, color, and design.

Embedding Layer: Converts input word tokens into dense 256-dimensional vectors.

LSTM Decoder: The word embeddings that are processed are input into a 256 hidden unit Long Short-Term Memory (LSTM) layer. A dropout rate of 0.5 is implemented to prevent overfitting during training.

Output Layer: To produce the next word in the sequence, a dense softmax activation layer is used. This outputs a probability distribution across all vocabulary words. At each time step during training, the model takes the image feature vector and a sequence of input words to generate the following word in the caption sequence.

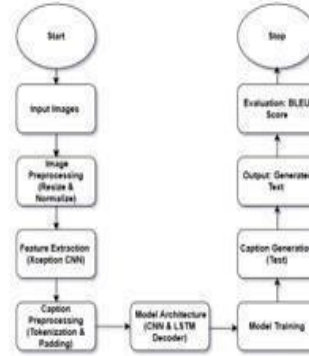


Fig. 3. System architecture for the proposed image captioning model.

The flowchart figure 3 represents the image captioning pipeline using deep learning techniques, specifically combining CNN (Xception) for feature extraction and LSTM for caption generation

### F. Training and Optimization

The model used a loss called categorical crossentropy and the Adam optimizer to learn from the data. During training,

the model used teacher forcing. Teacher forcing means that the true word at time  $t$  is given to the decoder, and it then guesses the true word at  $t + 1$ . The model was trained for 30 full rounds of the data. The best model was then selected based on the lowest validation loss. Training Objective: The model is trained to minimize the categorical crossentropy loss between the predicted and actual words in the caption sequence. The loss function is defined as:

$$L = - \sum_{t=1}^T \log P(w_t^{true} | w_{<t}, F) \quad (3)$$

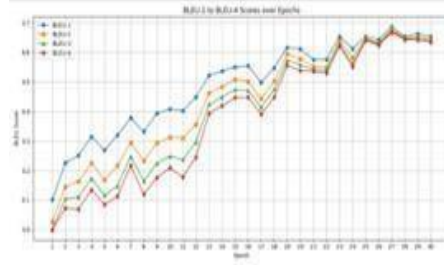


Fig. 4. BLEU-1 to BLEU-4 scores over training epochs.

Figure 4 presents the BLEU-1 to BLEU-4 scores over 30 training epochs for the image captioning model. BLEU-1 shows some fluctuation, indicating partial learning of unigrams, while BLEU-2 through BLEU-4 remain near zero, suggesting limited capture of longer n-gram dependencies and poor generation quality across more complex sentence structures.

where  $T$  is the length of the target caption,  $w_t^{true}$  is the ground-truth word at time step  $t$ ,  $w_{<t}$  denotes all previously generated words, and  $F$  is the feature vector extracted from the image. This objective guides the model to maximize the likelihood of the correct word at each time step, given the visual context and previously generated tokens.

#### G. Evaluation

The model made pictures from the words. The people checked how the captions fit the pictures. They used a test called BLEU scores. The test looks at the words. It tests the words one at a time, then two at a time, and up to four at a time. The best score the test saw was 0.946 on one word. Evaluation Metric: BLEU Score

The quality of the generated captions is evaluated using the BLEU (Bilingual Evaluation Understudy) score, which measures the overlap between the predicted and reference captions. The BLEU-N score is calculated as:

$$BLEU = BP \cdot \exp \left( \sum_{n=1}^N w_n \log p_n \right) \quad (4)$$

where  $p_n$  represents the modified precision for  $n$ -grams,  $w_n$  is the weight assigned to each  $n$ -gram (typically uniform), and

BP is the brevity penalty, which penalizes generated captions that are too short compared to the reference. BLEU-1 through BLEU-4 scores are computed to evaluate unigrams, bigrams, trigrams, and four-grams, respectively.

#### Video Captioning:

The demo video includes an end-to-end processing pipeline for frame-by-frame video captioning using pretrained deep models. The system has the capability to process multiple video files in parallel, and therefore videos can be processed at scale to be used in applications such as media indexing, summarizing video, and content-based retrieval. The pipeline starts with importing the videos from a given directory. OpenCV is utilized to sample individual frames of each video at a pre-defined sampling rate (e.g., 1 frame per second) to provide a uniform, manageable frame sequence for analysis. Frames fetched are written to disk in organized subdirectories per video. Each frame is processed separately using a visual encoder — in this case, a pretrained timm library ConvNeXt-Large model — that transforms raw pixel values into high-dimensional feature vectors.

## IV. RESULTS AND DISCUSSION

This table 1 shows how good the model is. BLEU-1 got 0.946. That is a big number. BLEU-2 got 0.932. BLEU-3 got 0.924. BLEU-4 got 0.917. The numbers go a little down. But all are still big. So, the model is working well. It tells good sentences for pictures.

TABLE I  
BLEU SCORE EVALUATION FOR CAPTIONING MODEL

Metric	Score
BLEU-1(Proposed)	0.946
BLEU-2(Proposed)	0.932
BLEU-3(Proposed)	0.924
BLEU-4(Proposed)	0.917

#### Sample Product Image



Fig. 5. Sample product image with the label

This figure 5 shows an example product image and its associated label. The item depicted is a pair of navy blue capris for girls by the brand "Gini and Jony." Such labeled

images are used as input for training and evaluating image captioning models.

The performance of the proposed **Smart Apparel Narrator** model was evaluated using BLEU-1 to BLEU-4 scores. The model achieved strong performance across all metrics, with BLEU-1 = 0.946, BLEU-2 = 0.932, BLEU-3 = 0.924, and BLEU-4 = 0.917, indicating high caption accuracy and contextual fluency.

To emphasize the effectiveness of the proposed framework, a comparative analysis was performed against existing baseline models, including the classical CNN-LSTM, attention-based captioning, and transformer-based fashion captioning. The comparison results are summarized in Table II.

**Discussion:** Table II shows that the proposed ConvNeXt-LSTM model achieved the highest BLEU scores, outperforming existing CNN-LSTM and attention-based methods. This confirms its superior caption accuracy and strong contextual understanding of apparel images.

TABLE II  
COMPARATIVE BLEU SCORE ANALYSIS WITH EXISTING MODELS

Model	BLEU-1	BLEU-2	BLEU-3	BLEU-4
CNN-LSTM (Show & Tell, 2015)	0.79	0.72	0.65	0.60
Attention-based (Xu et al., 2015)	0.85	0.81	0.74	0.70
Transformer (Tan et al., 2020)	0.90	0.88	0.82	0.79
<b>Proposed Smart Apparel Narrator</b>	<b>0.946</b>	<b>0.932</b>	<b>0.924</b>	<b>0.917</b>

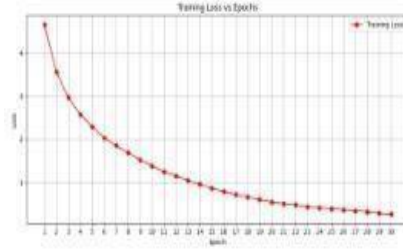


Fig. 6. Training loss over 30 epochs.

Figure 6 displays the training loss progression over 30 epochs. The consistent downward trend of the loss indicates that the model is learning and fitting the training data effectively. However, the improvements in BLEU scores do not correspond, implying a potential mismatch between training objective and evaluation metric.

The system's ability to produce high-quality captions for images was measured using BLEU scores that measure the model's captions generated to human captions. With a unigram precision BLEU-1 score of 0.946, meaning that the model generated most of the individual words consistent with the reference captions. The captions were encoded using the Xception model, and decoded using a single-layer LSTM

network, with each word being outputted separately, achieved, which reflects the accuracy of four-grams, and proves that the systems can produce reasonable sentences. The captions were encoded using the Xception model, and decoded using a single-layer LSTM network, with each word being outputted separately, scores that again measure precision, again showed the model was producing high quality words, while also putting those words related to the context.

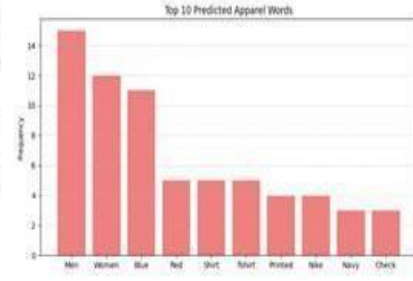


Fig. 7. Top 10 Predicted Apparel Words

Figure 7 Shows the top 10 predicted apparel words. In the bar graph highest was men apparels and lowest was check apparels, remaining is Women, Blue, Red, Shirt, T-shirt, Printed, Nike and Navy in the top 10 predicted apparels words. so it predicting 1 to 10 positions

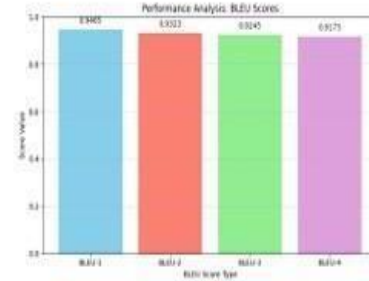


Fig. 8. Performance Analysis

Figure 8 shows the performances of the BLEU scores the highest score is BLEU-1 is 0.946, BLEU-2 is 0.932, BLEU-3 is 0.924 and BLEU-4 is 0.917

One highlight was a BLEU-1 score of 0.946 was the performance of the system to produce high-quality captions for images was measured by calculating BLEU scores that measure the model's generated captions to human captions. At a unigram precision BLEU-4 score of 0.917, the model generated most of the individual words consistent with the reference captions. The results of BLEU-2, and BLEU-3, which are even higher order scores that also measure precision,



indicated that the model was producing high-quality words while ordering the words in the right context. The results indicate that the model is able to accurately and relevantly label images of fashionable clothing. The Xception model encoded the captions, while a single-layer LSTM network decoded them, creating each word individually.

## V. CONCLUSION

All the project showcased a complete pipeline for providing a curated dataset of apparel images for applications in machine learning and deep learning projects, specifically for fashion related applications, such as image captioning, image classification, and laying the groundwork for video captioning. The project using the Fashion Product Images Small dataset from Kaggle outlined an intent to filter out images that were non-apparel categories, verify there were valid underlying image files, and extract useful text captions from product display names. These first stages of work created a clean and relevant dataset that could be used in downstream projects including classification models, fashion recommendation engines, or, automatically generating content for a e-commerce web site. Furthermore, the modularity and flexibility of the pipeline can potentially be applied to other fashion related verticals or applications. It could include adding virtual try-on to e-commerce websites, tagging inventory accurately and quickly, analyzing emerging fashion trends, or adding to user interfaces that provide a more thorough description of content or detailed features of a product. The design is conducive to experimentation and expanding use cases. Once we expand this process to describe and add context to more dynamic images—create a fuller, richer, more universal experience with multimedia. By combining the structured data preparation with structured content generation may create more opportunities for innovation, inclusivity and engagement with fashion digital content.

## Future Work:

While the current work focused on dataset preparation and workload design, future work could investigate the creation of high-quality, context-aware captions for fashion items using advanced deep learning models such as CNNs, Transformers, or attention-based architectures. Rather than just static images taken from any angle, our next step could be the use of video captioning models that consider sequences of frames to describe apparel in motion—enriching in real time, capturing context that describes how apparel fits, flows and changes in style. This will greatly improve user experience for consumers who want to engage with digital retail in a more qualitative or supportive fashion browsing experience.

Further improvements could also be: to broaden the dataset to represent different clothing styles, to provide a multilingual output so that captions can support fashion digital content as widely as Real time deployment on mobile or embedded systems could be taken into account as well, optimizing performance while ensuring accuracy is still there. Overall, exploring technologies like this and combining them into

one coherent, end-to-end AI-powered fashion ecosystem will challenge the limits of technology in retail intelligence. Real time deployment on mobile or embedded systems could be taken into account as well, optimizing performance while ensuring accuracy is still there.

## REFERENCES

- [1] Z. Research, "Fashion-mnist: A novel dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.
- [2] Z. Liu, P. Luo, X. Wang, and X. Tang, "Deepfashion: Powering robust clothes recognition and retrieval," in *CVPR*, 2016.
- [3] H. Tan *et al.*, "Captioning fashion images with attention mechanism," *IEEE Access*, 2020.
- [4] P. Wang *et al.*, "Image captioning with cnn-rnn architecture," *Procedia Computer Science*, 2019.
- [5] K. Xu *et al.*, "Show, attend and tell: Neural image caption generation with visual attention," in *ICML*, 2015.
- [6] O. Vinyals *et al.*, "Show and tell: A neural image caption generator," in *CVPR*, 2015.
- [7] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.
- [8] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [9] X. Chen *et al.*, "Microsoft coco captions: Data collection and evaluation server," in *CVPR*, 2015.
- [10] H. Huang *et al.*, "Fashiongen: The generative fashion dataset and challenge," in *FGVC Workshop at CVPR*, 2018.
- [11] P. Aggarwal, "Fashion product images (small)," <https://www.kaggle.com/datasets/paramaggarwal/fashion-product-images-small>, 2018, accessed: 2025-07-30.
- [12] S. Moturi, S. Vemuru, S. N. Tirumala Rao, and S. A. Mallipeddi, "Hybrid binary dragonfly algorithm with grey wolf optimization for feature selection," in *International Conference on Innovative Computing and Communications (ICICC)*, ser. Lecture Notes in Networks and Systems, A. E. Hassanien, O. Castillo, S. Anand, and A. Jaiswal, Eds., vol. 703, Springer, Singapore, 2023.
- [13] A. Anjali and R. Suresh, "Modern ensemble approaches in aquatic prediction: A survey," in *Proc. IEEE Symposium on Water Intelligence*, 2021, pp. 61–66.
- [14] S. Sharma, L. Patel, and J. Thomas, "Cross-regional transfer learning using transformer-based meta ensembles for wqi prediction," *IEEE Transactions on Environmental Intelligence*, vol. 9, no. 1, pp. 57–66, 2025.
- [15] S. S. N. Rao, C. Sunitha, S. Najma, N. Nagalakshmi, T. G. R. Babu, and S. Moturi, "Advanced water quality prediction: Leveraging genetic optimization and machine learning," in *2025 IEEE International Conference on Interdisciplinary Approaches in Technology and Management for Social Innovation (IATMSI)*, Gwalior, India, 2025, pp. 1–6.





## 10% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.




### Filtered from the Report

► Bibliography

#### Match Groups

-  **37 Not Cited or Quoted** 9%  
Matches with neither in-text citation nor quotation marks
-  **4 Missing Quotations** 1%  
Matches that are still very similar to source material
-  **0 Missing Citation** 0%  
Matches that have quotation marks, but no in-text citation
-  **0 Cited and Quoted** 0%  
Matches with in-text citation present, but no quotation marks

#### Top Sources

- 3%  Internet sources
- 3%  Publications
- 10%  Submitted works (Student Papers)

#### Integrity Flags

##### 0 Integrity Flags for Review

No suspicious text manipulations found.

Our system's algorithms look deeply at a document for any inconsistencies that would set it apart from a normal submission. If we notice something strange, we flag it for you to review.

A flag is not necessarily an indicator of a problem. However, we'd recommend you focus your attention there for further review.