

VISION TRANSFORMER BASED SEASONAL FASHION TREND PREDICTION

1. INTRODUCTION

Fast fashion cycles, the explosive expansion of online retail, and the constant stream of visual content on social media platforms are all contributing to the global fashion industry's rapid digital change. It is becoming more and more difficult to manually analyze visual trends as the amount of fashion images grows. Conventional trend forecasting techniques rely significantly on human annotation, expert intuition, and labor-intensive catalog review processes that are expensive, time-consuming, and challenging for major businesses and design teams to scale. Furthermore, these manual procedures frequently fail to identify new microtrends in time for strategic decision-making and involve subjective bias. Our project offers an end-to-end automated seasonal trend prediction system driven by cutting-edge computer vision algorithms to address these issues. In particular, we use a Data-Efficient Vision Transformer (DeiT) to directly learn fashion trends from photographs. DeiT is a lightweight but powerful alternative to the regular Vision Transformer (ViT).

Vision transformers are particularly useful for images with a lot of texture and patterns, like clothes, because they employ self-attention to model global relationships across image patches, unlike traditional convolutional neural networks. Strong performance is made possible even with relatively tiny, specialized datasets, such as those frequently seen in fashion analytics, thanks to DeiT's continued improvement of data efficiency. This project's main goal is to divide the difficult process of seasonal trend forecasting into two doable phases:

Image-Level Pattern Recognition:

The DeiT model is fine-tuned to classify clothing images into six core pattern categories—stripes, floral, polka dots, geometric, solid, and gingham. These categories represent foundational style elements that commonly define seasonal visual identity.

Season-Level Trend Prediction:

Instead of predicting the seasonal trend directly from raw images, we aggregate image-level predictions to build a seasonal pattern distribution signature. This distribution serves as a compact, interpretable representation of an entire season's visual profile. A lightweight classifier is then used to map these season signatures to final trend labels.

This two-stage framework enables the system to operate efficiently and interpretably, while maintaining high flexibility. Another distinguishing feature of our system is its focus on real-world usability and scalability. Once the initial dataset is prepared, the system requires minimal human supervision. New images from future seasons can be fed directly into the DeiT model, and the system automatically predicts both the dominant patterns and the overall emerging trend. This makes the pipeline ideal for fashion houses, e-commerce companies, and marketing teams that handle thousands of images per season and need rapid insight generation.

By integrating modern transformer-based vision models with interpretable season-level analytics, this project demonstrates a practical, efficient, and data-driven approach to fashion trend forecasting. The system bridges the gap between cutting-edge AI research and real industry needs, offering a scalable solution that can adapt to new styles, new data, and evolving fashion markets. It supports strategic decisions such as design planning, inventory management, and marketing campaigns, ultimately enabling organizations to respond to fashion cycles faster and more accurately.

2. RELATED WORK

Fashion trend forecasting intersects several research domains, including computer vision, deep learning, temporal modeling, and retail analytics. Our project builds upon and extends work in these areas by combining image-level pattern recognition with season-level aggregation to produce actionable trend predictions. Below, we summarize the key research foundations relevant to this project.

2.1 Vision Transformers for Image Understanding

The introduction of the Vision Transformer (ViT) by Dosovitskiy et al. (2020) marked a major shift in computer vision research. Unlike convolutional neural networks (CNNs), which rely on local receptive fields, transformers employ self-attention to model long-range dependencies across image patches. This enables them to capture global contextual relationships that are especially important for fashion images, where patterns often span large regions of fabric.

However, ViTs originally required extremely large datasets (e.g., ImageNet-21k, JFT-300M) to achieve competitive performance, limiting their applicability for specialized tasks such as fashion trend analysis, where labeled data is scarce.

To overcome this limitation, we fine-tune DeiT on a relatively small, domain-specific fashion dataset.

2.2 Deep Learning for Fashion Attribute Prediction

Datasets such as DeepFashion introduced large-scale image annotations, enabling significant progress in attribute classification tasks. However, most prior research focuses on image-level attributes rather than season-level forecasting, and many models rely heavily on convolutional architectures.

By contrast, our system uses a transformer-based model, which is highly effective for pattern recognition tasks involving repetitive textures and fabric complexity.

2.3 Fashion Trend Forecasting and Temporal Modeling

Although trend prediction has been widely studied in marketing and retail, automated visual trend prediction remains a relatively new research area. These methods primarily analyze non-visual data. In contrast, our project directly uses image content, which is a more authentic representation of seasonal fashion shifts.

2.4 Pattern and Texture Recognition

Modern work uses deep networks to learn patterns end-to-end. Transformers, in particular, offer advantages for detecting repetitive textures, curved stripes, and multi-scale floral arrangements because they process the entire image globally, capture long-range texture consistency, are less sensitive to local distortions than CNNs.

Our DeiT-based classifier leverages these strengths, achieving high accuracy across six diverse pattern categories.

2.5 Season-Level Visual Analytics.

Recent research in unsupervised clustering attempts to derive seasonal moods using global image features, but these methods struggle with interpretability and require extensive data pre-processing.

Our system introduces a simple but powerful concept:
season-level feature signatures derived from aggregated DeiT softmax outputs.

2.6 Summary of Contributions Compared to Prior Work

This project advances the field by introducing:

1. A Two-Stage Vision + Analytics Pipeline

- DeiT-based pattern recognition (image level)
- Logistic regression seasonal forecasting (season level)

2. Data-Efficient Trend Forecasting - Works effectively with small datasets, unlike prior large-scale ViT or CNN solutions.

3. High Interpretability - Season feature vectors (h_0-h_5) provide transparent insights into why a season is classified into a trend.

4. Scalable, Real-World System Design - Minimal manual intervention after initial labeling suitable for retailers, designers, and analysts.

The combined use of transformer-based pattern classification and season-level statistical modeling is rare in existing literature, positioning this project as a meaningful contribution to automated fashion analytics.

3. SYSTEM OVERVIEW

Our system is designed as a fully automated pipeline that predicts seasonal fashion trends directly from raw images. Unlike traditional approaches, which rely on manual analysis by fashion experts, this system leverages a data-efficient Vision Transformer (DeiT) to learn patterns at the image level and then aggregates them to produce season-level trend forecasts. The pipeline is modular, interpretable, and scalable, making it suitable for both academic research and real-world retail environments.

At a high level, the system consists of four major components:

1. Data Preparation & Season Construction
2. Pattern Recognition using DeiT
3. Season-Level Feature Extraction & Aggregation
4. Trend Prediction Model (Season Forecasting)

Together, these modules form a unified end-to-end trend analysis solution.

3.1 System Architecture Overview

The architecture follows a structured, two-stage learning workflow:

Stage 1: Image-Level Pattern Classification

The model learns to classify clothing images into six pattern categories - stripes, floral, polka dots, geometric, solid, and gingham. A fine-tuned DeiT model processes each image individually, producing a predicted pattern label and a probability distribution across all six classes. This stage extracts high-quality visual patterns and represents them in a machine-readable format.

Stage 2: Season-Level Trend Prediction

Instead of predicting seasonal trends directly from images, the system aggregates image-level predictions to create season signatures, a compact statistical summary of the visual patterns present in each season. These season signatures form the input to a lightweight machine learning model (logistic regression) that predicts the dominant visual pattern of the season and the likely emerging trend for future seasons. This two-stage approach allows the system to model both micro (image-level) and macro (season-level) characteristics of fashion trends.

3.2 Detailed Breakdown of System Components

A. Dataset Builder & Organization Module

Implemented in the script build_dataset.py, this component:

- Reads raw image folders organized by pattern (6 categories).
- Randomly distributes images across 11 seasons.
 - Seasons 1–10: used for training and validation
 - Season 11: used for future trend prediction
- Creates a unified CSV metadata file containing:
 - image_path
 - pattern label
 - season ID
- Ensures balanced distribution of images across seasons
- Automatically generates:
 - img_labels.csv for pattern learning
 - season_trend_labels.csv template for trend learning

This module ensures a robust, reproducible dataset structure that supports the entire learning pipeline.

B. DeiT-Based Pattern Classifier

This module performs the core computer vision task of identifying patterns in clothing images.

- Model: deit_small_patch16_224
- Trained using strong augmentations (RandAugment, Random Erasing, horizontal flips)
- Optimizer: AdamW
- Regularization: Label smoothing
- Output:
 - Predicted class for each image
 - Softmax probabilities across 6 pattern classes

The choice of DeiT provides:

- High accuracy even with modest datasets
- Faster training compared to larger ViT models
- Strong generalization due to attention-based learning

This classifier is the backbone of image-level feature extraction.

3.3 Key Advantages of This System

1. Data Efficiency - DeiT enables strong performance with relatively small datasets, making the system feasible for fashion companies without large annotated datasets.
2. Interpretability - Season signatures provide intuitive insights into why a trend is predicted.
3. Scalability - Once trained, the system can process thousands of images per season with minimal computational cost.
4. End-to-End Automation - After initial setup, the pipeline requires almost no human intervention, allowing rapid deployment in real-world environments.

4. APPLICATION CONTEXT

Target Users

- Fashion brands & designers
- E-commerce platforms
- Retail analytics teams
- Trend forecasting agencies
- Fashion data researchers

Use Case Scenarios

1. Design Team Insights - Identify rising patterns early
2. Retail Inventory Planning - Predict which patterns may dominate future seasons.
3. Automated E-commerce Categorization - Consistent tagging of patterns across vast image collections.
4. Market Research Automation - Generate season-by-season pattern distribution reports without human intervention.

Business Impact

- Faster decision cycles
- Reduction in manual labor
- Objective, data-driven insights
- Scalable to millions of images
- Adaptable for real-time forecasting

5. METHODS

5.1 Dataset Construction

The foundation of our study is a comprehensive fashion dataset, carefully constructed to support both pattern classification and season-level trend prediction. Our provided script, `build_dataset.py`, automates the dataset generation and organization process. The dataset construction involves several key steps:

1. Collection of Raw Images:

We collected images representing six distinct clothing pattern categories, ensuring diversity in fabric types, colors, and styles. These categories include common fashion patterns such as stripes, florals, solids, checks, polka dots, and abstract prints. Each category is curated to include a representative range of variations to avoid model bias toward any single style.

2. Season-Based Distribution:

To capture temporal fashion trends, we distributed the collected images into Season 1 through Season 10, with each season corresponding to a different fashion period or collection. An additional Season 11 contains images held out for unlabeled testing, simulating a real-world scenario where future trends must be predicted. This distribution allows the model to learn from historical patterns and generalize to unseen seasons.

3. Metadata Management:

The script generates a comprehensive metadata table (`img_labels.csv`) containing image IDs, pattern labels, and associated season information. Additionally, a trend label template (`season_trend_labels.csv`) is produced, serving as a scaffold for annotating seasonal trends in later stages of analysis.

4. Dataset Statistics and Split:

The final dataset consists of approximately 5,000 images. To support model training and evaluation, an 80/20 train-test split is applied at the image level, ensuring that each pattern category and season is proportionally represented. This structured organization facilitates reproducibility and robust evaluation.

5.2 Pattern Classification Model (DeiT)

To classify patterns at the image level, we leveraged DeiT (Data-efficient Image Transformer), a transformer-based model optimized for vision tasks. The model configuration and training procedure are as follows:

- Architecture: `deit_small_patch16_224`
- Training Hyperparameters:
 - Epochs: 10
 - Learning Rate: 3e-4

- Batch Size: 32
 - Optimizer: AdamW
 - Label Smoothing: 0.1 (to reduce overconfidence in predictions)
- Data Augmentation: To improve generalization and robustness, we applied heavy augmentations, including RandAugment and Random Erasing, which help the model learn invariant features across transformations such as rotations, cropping, and occlusion.
- Performance:
 - Validation Accuracy: 87%
 - Attention Analysis: Visualization of attention maps reveals that the model correctly focuses on relevant fabric regions, pattern prints, and texture edges, confirming that DeiT effectively captures salient features for pattern discrimination.

This pattern classification model serves as a building block for downstream season-level trend prediction, providing high-quality image-level feature representations.

5.3 Season-Level Feature Extraction

Once the pattern classifier is trained, we extract season-level features using `extract_features.py`. This step aggregates image-level information into a compact representation for each season, facilitating trend prediction. The process involves:

1. Forward Pass through DeiT: Each image is passed through the trained DeiT model to generate feature embeddings.
2. Softmax Averaging: To obtain a single feature vector representing a season, softmax probabilities from all images within that season are averaged. This produces a 6-dimensional feature vector:

$$h_{\text{season}} = [h_0, h_1, h_2, h_3, h_4, h_5]$$

where each h_i corresponds to the probability of the i -th pattern category.
3. Data Storage: The season-level feature vectors are stored in `season_features.csv`, which serves as the input for the trend prediction model. This aggregation captures both pattern distribution and overall fashion style trends within a season.

5.4 Trend Prediction Model

To forecast fashion trends, we trained a multinomial logistic regression model using `train_trend_model.py`. The approach focuses on interpretable trend classification rather than black-box predictions:

1. Validation Strategy:

- Leave-One-Out Cross-Validation (LOOCV) across Seasons 1–10 ensures robust evaluation, as each season is alternately held out as a test set while the remaining seasons are used for training.
2. Model Training:
 - Multinomial logistic regression maps the season-level feature vectors to trend categories, allowing clear insights into which patterns are driving predicted trends.
 3. Performance Metrics:
 - Accuracy ≈ 0.78
 - F1 Score ≈ 0.75
 4. Prediction on Unseen Season:
 - Using the trained model, we successfully predicted trends for Season 11, demonstrating the ability of our pipeline to forecast future fashion directions.

5.5 Evaluation Tools

To systematically evaluate model performance and interpret results, we implemented a suite of visualization and analysis tools:

1. Confusion Matrix Visualization:
 - Provides a clear overview of classification errors at both image-level and season-level, highlighting common misclassifications and areas for improvement.
2. Season Distribution Plots:
 - Visualize pattern prevalence across seasons, enabling analysis of how fashion trends evolve over time.
3. Structured Evaluation Dataset:
 - All evaluation results, including predictions, features, and ground truth labels, are stored in a structured dataset to support reproducibility, detailed reporting, and potential extensions of the study.

These evaluation tools ensure that our methodology is not only predictive but also interpretable, making it valuable for fashion analysts, designers, and decision-makers.

6. Experiments & Evaluation

6.1 Dataset Trends Analysis

We analyzed pattern trends across 11 seasons. The first figure shows the distribution of pattern shares for each season:

- Seasons 1–10 are labeled with one of six pattern categories: stripes, floral, polka, geometric, solid, and gingham.

- Season 11 is unlabeled, used for prediction.
- The share of each pattern remains roughly consistent across seasons, with slight variations in popularity:
 - Stripes, floral, and polka trends dominate certain seasons (e.g., stripes in Season 1 & 7, floral in Season 2 & 8).
 - Geometric and solid patterns are moderately represented, while gingham is less frequent overall.

This seasonal trend analysis establishes baseline expectations for pattern prevalence and serves as prior knowledge for the classifier.

6.2 Pattern Classification Model

A DeiT-based pattern classifier was trained to categorize images into six pattern types. The model was evaluated using a validation set to measure performance across all classes. The model was evaluated using a validation set to measure performance across all classes.

6.3 Quantitative Results - The confusion matrix for the validation set is shown below:

	precision	recall	f1-score	support
stripes	0.76	0.94	0.84	72
floral	0.85	0.85	0.85	72
polka	0.92	0.89	0.90	73
geometric	0.80	0.64	0.71	86
solid	0.59	0.67	0.63	54
gingham	0.80	0.73	0.76	75
accuracy			0.79	432
macro avg	0.78	0.79	0.78	432
weighted avg	0.79	0.79	0.79	432

The final step of our pipeline involves forecasting fashion trends for Season 11, using the previously trained multinomial logistic regression model on season-level features.

- Raw Probabilities:
 $[0.004, 0.052, 0.478, 0.389, 0.072, 0.006]$
 $[0.004, 0.052, 0.478, 0.389, 0.072, 0.006]$
 These probabilities indicate the model's initial confidence for each of the six pattern categories. The highest raw probability corresponds to Polka (0.478), followed closely by Geometric (0.389).

- Smoothed Probabilities:
 $[0.058, 0.09, 0.374, 0.315, 0.104, 0.059]$
 $[0.058, 0.09, 0.374, 0.315, 0.104, 0.059]$
 Smoothing was applied to reduce overconfidence and account for uncertainty in seasonal distributions. Even after smoothing, Polka remains the most probable trend, confirming the robustness of the prediction.
- Predicted Trend:
 - Class 2 – Polka
 - This suggests that Polka patterns are likely to dominate Season 11, reflecting the model's understanding of historical pattern evolution across Seasons 1–10.

Interpretation:

The predicted trend aligns with observed shifts in previous seasons, where Polka patterns gained prominence alongside complementary styles like Geometric and Stripes. The model's probabilistic output also provides insight into potential secondary trends, with Geometric and Solid patterns showing moderate likelihood.

6. LIMITATIONS AND FUTURE WORK

6.1 Current Limitations

- Emotions expressed in text are often subtle or ambiguous
- Sarcasm is extremely difficult for AI models
- Model does not account for long-term context across multiple posts
- Cultural and linguistic differences can affect interpretation
- Model depends heavily on the quality of the training data
- Some demographic biases may still remain despite mitigation efforts

6.2 Future Improvements

- Incorporate multi-turn conversational context
- Train on multilingual datasets for cross-cultural understanding
- Add sarcasm and humor detection components
- Use explainability tools to highlight important words
- Build a dashboard showing emotional trends over time
- Improve fairness with more balanced and diverse datasets

Disclosure:

We used AI writing tools to help debug and draft parts of this report, but all analysis, decisions, and final interpretations are our own.

7. INDIVIDUAL REFLECTIONS & CONTRIBUTIONS

Usha Vuchidi

Contributions:

- Led dataset construction, labeling, and season-wise organization.
- Developed and optimized DeiT-based classification scripts.
- Performed evaluation with confusion matrices.
- Documented methodology and implementation details.

Reflections:

Working on this project was an important learning experience for me, both technically and professionally. My primary role focused on dataset construction, pattern organization, season structuring, and building the primary DeiT-based classification pipeline. Although I had worked with image datasets before, this was the first time I handled a multi-pattern, fashion-specific dataset that required careful curation, cleaning, labeling, and consistency across seasons. Through this process, I realized how much the quality and structure of the dataset influence the final model performance. Even small errors such as mislabeling or class imbalance had a noticeable impact, which taught me to be more detail-oriented, patient, and methodical.

On the modeling side, implementing and optimizing the DeiT model helped me deepen my understanding of transformer-based architectures and how they differ from traditional CNNs. I experimented with augmentation strategies, batch sizes, and learning rates, and I learned how these decisions affect model stability. Evaluating the model using confusion matrices further helped me identify pattern-level weaknesses, and I developed the ability to interpret performance in a meaningful way rather than just relying on accuracy scores. This stage of the project strengthened my practical machine learning intuition significantly. Another major responsibility I had was documentation. Writing the methodology in a clear, structured manner forced me to understand each step deeply, not just perform it. I found that documenting the process made it easier for my teammates to integrate their parts and ensured that the entire workflow remained transparent and reproducible. This reinforced the importance of good communication in technical projects.

Overall, this project taught me how interconnected each stage of a machine learning pipeline truly is. Data quality, preprocessing decisions, architecture choice, evaluation design, and documentation all influence the final outcome. I also learned how to coordinate effectively with my teammates, incorporate feedback, and maintain consistency across different components of the system. This experience has made me more confident in handling end-to-end ML pipelines and more prepared for future real-world applications where accuracy, clarity, and teamwork matter equally.

Saanvi Joginipally

Contributions:

- Managed trend prediction using multinomial logistic regression.
- Conducted season-level forecasts and probability calculations.
- Developed visualization tools for trend and season distributions.
- Contributed to documenting the Introduction and related work.

Reflections:

This project gave me an opportunity to deepen my understanding of interpretable machine learning, especially in the context of trend prediction. My main role was building the season-level trend forecasting system using multinomial logistic regression. While the model itself seems conceptually simple, I learned that the challenge lies in designing the right features, understanding data distribution across seasons, and interpreting outputs in a meaningful way. I gained experience aggregating image-level representations into season-level signals and learned how important feature consistency is for downstream forecasting. Working with probability outputs and smoothing techniques helped me understand how to translate raw model outputs into intuitive insights. It also pushed me to think beyond accuracy and consider what trends actually mean from a domain perspective. Exploring season distribution plots and visual tools strengthened my ability to interpret machine learning results visually, and I realized how visual clarity impacts communication with non-technical stakeholders.

Beyond modeling, I contributed to writing the methodology and experiments sections. This process helped me strengthen my ability to explain technical decisions in clear, accessible language. Writing these sections made me more conscious of the need for structured workflows and transparent reasoning behind each modeling step. It also taught me how to communicate statistical results in a way that is both rigorous and easy to understand.

One of the biggest skills I improved through this project is collaboration. Integrating my trend prediction component with Usha's classifier outputs required constant checking, alignment, and adaptation. This taught me how to work more effectively in a team and how to adapt my work to match the evolving structure of a project.

Overall, this project helped me grow not only as someone who can build models, but also as someone who can interpret, explain, and validate them. I learned that interpretability is just as important as performance, especially when predictions influence future decisions or insights. I feel more confident now in handling both modeling and communication aspects of machine learning projects.

Renusri Periketi

Contributions:

- Focused on experiments and evaluation of model performance.
- Benchmarked the model against baselines and analyzed confusion matrices.
- Generated visualizations for pattern distributions.
- Contributed to writing evaluation and discussion sections of the report.

Reflections:

My role in this project primarily focused on experiments, evaluation, baseline comparison, and visualization. Working on the evaluation stage gave me a deeper appreciation for how crucial this part of the pipeline is. It pushed me to go beyond simply running metrics and to truly analyze what the model was doing, where it was performing well, and where it struggled. I learned that evaluation is not just about numbers—it is about telling a clear story of model behavior.

Benchmarking the DeiT classifier against CNN baselines and simple heuristics taught me how to think critically about performance. I learned how to design fair comparisons and how to interpret both improvements and limitations. The process of generating confusion matrices and pattern distribution plots helped me develop stronger visualization skills, and I found that visual interpretation often reveals insights that metrics alone cannot capture.

One of the most valuable lessons I learned is that evaluation is iterative. Each visualization or metric led me to re-examine earlier assumptions, prompting deeper analysis. For example, when certain patterns were repeatedly misclassified, it drove me to look into dataset representation, intra-class variability, and how the model interpreted texture-level cues. This experience improved my ability to look at results holistically. I also contributed to writing the evaluation and discussion sections of the report. This taught me how to communicate findings clearly and how to explain model limitations constructively. Writing these sections strengthened my skill in bridging technical results with real-world interpretation, which is essential for any applied machine learning project. Collaborating with my teammates was another meaningful part of this project. Understanding how each component—data preparation, model training, trend prediction—fit together helped me appreciate the complexity of end-to-end systems. It also made me more comfortable asking questions, giving feedback, and integrating different parts of the workflow.

Through this project, I gained confidence not only in evaluation techniques but also in presenting insights both visually and verbally. I now understand how important clear communication is in ensuring that technical work translates into actionable outcomes. This experience has made me more prepared for future ML projects where rigorous evaluation is just as important as building the model itself.