

Bridging the Cold-Start Gap:
A Hybrid Visual-Collaborative Approach to Fashion
Recommendation

Project StyleForge Technical Report

Table of Contents

Abstract	3
1. Introduction	4
2. Related Work	6
3. Dataset and Evaluation Protocol	8
4. Model Architecture	9
5. Experimental Results	11
6. Business Impact Analysis	15
7. Limitations and Future Directions	16
8. Conclusion	18
References	19

Abstract

Recommender systems have become essential infrastructure for fashion e-commerce, yet they persistently struggle with what practitioners call the cold-start problem: new users lack the interaction history needed for collaborative filtering, while new products lack the engagement data required for accurate preference modelling. This challenge is particularly acute in fashion, where catalogue turnover is high and trend cycles are short. We present a hybrid system that addresses cold-start by fusing visual similarity signals with collaborative filtering through a learned attention mechanism. The key insight is that visual features provide immediate, if imperfect, signals about user preferences when behavioural data is sparse. A new user who has clicked on three floral dresses probably prefers floral patterns, regardless of what collaborative filtering might infer from such limited history. Our attention mechanism learns to weight visual versus collaborative signals based on data availability: for new users, visual dominates; for established users with rich histories, collaborative takes precedence. Experiments on the H&M dataset suggest this adaptive approach yields 14.6% improvement in recommendation quality overall, with gains reaching 163% for new users and 790% for new items. We discuss both the promise and limitations of this approach, including computational overhead and the challenge of capturing temporal fashion dynamics.

1. Introduction

Fashion recommendation presents a curious paradox. On one hand, the domain seems ideally suited to personalisation: preferences are intensely personal, visual similarity is meaningful, and the sheer catalogue size makes discovery assistance valuable. On the other hand, the characteristics that make fashion interesting also make it difficult. Unlike books or movies where genre, author, and reviews provide strong content signals, fashion preferences are predominantly visual and notoriously hard to articulate. Ask someone why they prefer one dress over another and you will often get vague responses about it "looking better" or "feeling more me," responses that resist encoding into recommendation features.

Collaborative filtering has become the dominant approach despite these challenges, largely because it sidesteps the articulation problem. Rather than encoding why users prefer certain items, it identifies users with similar behaviour patterns and recommends what those similar users enjoyed. The approach works remarkably well when sufficient interaction data exists. The problem is that sufficient data often does not exist, at least not when it matters most.

Consider a new user visiting a fashion platform for the first time. Collaborative filtering has essentially nothing to work with, perhaps a few clicks, maybe a single purchase if we are lucky. The best it can do is recommend popular items, hoping that crowd wisdom compensates for missing individual signal. Meanwhile, the user sees generic recommendations that fail to reflect their taste, becomes frustrated, and potentially abandons the platform. First impressions matter, and collaborative filtering provides particularly bad first impressions.

The item side presents a mirror-image problem. Fashion catalogues turn over rapidly: seasonal collections introduce thousands of new products, fast-fashion retailers refresh inventory weekly, and trend-driven items have brief commercial windows. Each new product enters the catalogue as a cold-start item, unable to receive personalised recommendations until it accumulates engagement data. But engagement depends partly on recommendations, creating a chicken-and-egg situation. By the time enough data accumulates for accurate collaborative signals, the product may have passed its peak relevance.

Our work addresses these challenges through visual-collaborative fusion. The intuition is straightforward: while collaborative filtering requires behavioural history, visual similarity is immediately available for any product with an image. A new floral dress can be recognised as visually similar to other floral dresses the moment it enters the catalogue. A new user who clicks on three items provides visual preference signals, the clicked items' visual characteristics, even if three interactions are insufficient for collaborative inference. Visual features do not replace collaborative filtering; they complement it, providing signal when collaborative data is sparse and stepping back as richer interaction data accumulates.

2. Related Work

2.1 The Evolution of Collaborative Filtering

Collaborative filtering's history is essentially the history of recommendation systems. Early neighbourhood methods (Resnick et al., 1994) identified similar users and aggregated their preferences; matrix factorisation (Koren et al., 2009) learned latent factors explaining user-item interactions; and neural approaches (He et al., 2017) replaced linear combinations with learned non-linear functions. Each generation improved accuracy on benchmark datasets, but none fundamentally solved cold-start.

Neural Collaborative Filtering (NCF) represents the current state-of-the-art in many settings. By combining generalised matrix factorisation with multi-layer perceptrons, NCF captures both linear and non-linear interaction patterns. We adopt NCF as our collaborative backbone precisely because its architecture naturally extends to multi-modal fusion: the MLP pathway provides a convenient point for injecting additional features.

2.2 Visual Recommendation

The recognition that visual features matter for fashion recommendation is not new. He and McAuley (2016) demonstrated that CNN features improve recommendation quality by capturing visual similarity beyond categorical metadata. Subsequent work explored style-specific embeddings (Liu et al., 2017), fashion compatibility learning (Song et al., 2018), and attention over visual regions (Chen et al., 2019). What these approaches share is treating visual features as additional item metadata to be incorporated alongside collaborative signals.

Our contribution is less about novel visual features than about how visual and collaborative signals should be combined. Rather than fixed-weight fusion, we learn adaptive weights conditioned on data availability. The mechanism is conceptually simple, attention over two input streams, but the result is qualitatively different behaviour across user segments.

3. Dataset and Evaluation Protocol

We evaluate on the H&M Personalised Fashion Recommendations dataset, comprising approximately 1.37 million customers, 105,000 products, and 31.8 million transactions spanning two years. This scale enables meaningful evaluation across user segments while maintaining the sparsity characteristic of real fashion retail: average users have only 23 interactions across the two-year period.

Temporal splitting is essential for realistic evaluation. Training on random samples would allow the model to observe future interactions before predicting them, producing misleadingly optimistic results. We instead train on the first 20 months, validate on months 21-22, and test on months 23-24. This protocol simulates production deployment where models must predict future behaviour from historical data.

For cold-start evaluation, we segment users by training-set interaction count: new users (fewer than 5 interactions, comprising 23% of test users), regular users (5-50 interactions, 54%), and power users (more than 50 interactions, 23%). Similarly, items first appearing in the test period are designated as new items (8% of test items). These segments enable targeted analysis of cold-start performance.

4. Model Architecture

4.1 Visual Encoder

Our visual encoder builds on ResNet50 pretrained on ImageNet, with the final two blocks fine-tuned on fashion data. A projection head reduces 2048-dimensional backbone features to 256 dimensions suitable for efficient similarity computation and storage. The architecture choice reflects pragmatic considerations: ResNet50 is well-understood, extensively validated on fashion imagery, and supported by mature deployment infrastructure.

Training uses triplet loss with hard negative mining. For each anchor item, positive examples are items purchased by the same user (implicitly similar by preference), and negative examples are items never interacted with by that user. Hard negative mining selects negatives with high visual similarity to anchors, forcing the model to learn fine-grained distinctions that matter for fashion: the difference between two similar-looking dresses that a user would actually consider substitutes versus merely visually similar images.

Listing 1: Visual encoder architecture

```
class VisualEncoder(nn.Module):
    def __init__(self, embedding_dim=256):
        self.backbone = resnet50(pretrained=True)
        for param in list(self.backbone.parameters())[:-20]:
            param.requires_grad = False
        self.projection = nn.Sequential(
            nn.AdaptiveAvgPool2d(1), nn.Flatten(),
            nn.Linear(2048, 512), nn.ReLU(),
            nn.Linear(512, embedding_dim)
        )
```

4.2 Attention-Based Fusion

The fusion mechanism is where our approach departs most clearly from prior work. Rather than combining visual and collaborative embeddings with fixed weights, we learn attention weights conditioned on user context. The context encoder takes as input the user's interaction count and recency, producing a representation that conditions the attention computation. The attention network then outputs a two-dimensional weight vector summing to one, determining how much each modality contributes to the final prediction.

The key property we want is adaptive behaviour across user segments. For a new user with three interactions, collaborative filtering is essentially guessing, so the attention should favour visual. For a power user with hundreds of interactions, collaborative patterns encode rich preference information, so the attention should favour collaborative. We do not hard-code this behaviour; the attention network learns it from data by being trained end-to-end on recommendation accuracy.

5. Experimental Results

5.1 Overall Recommendation Quality

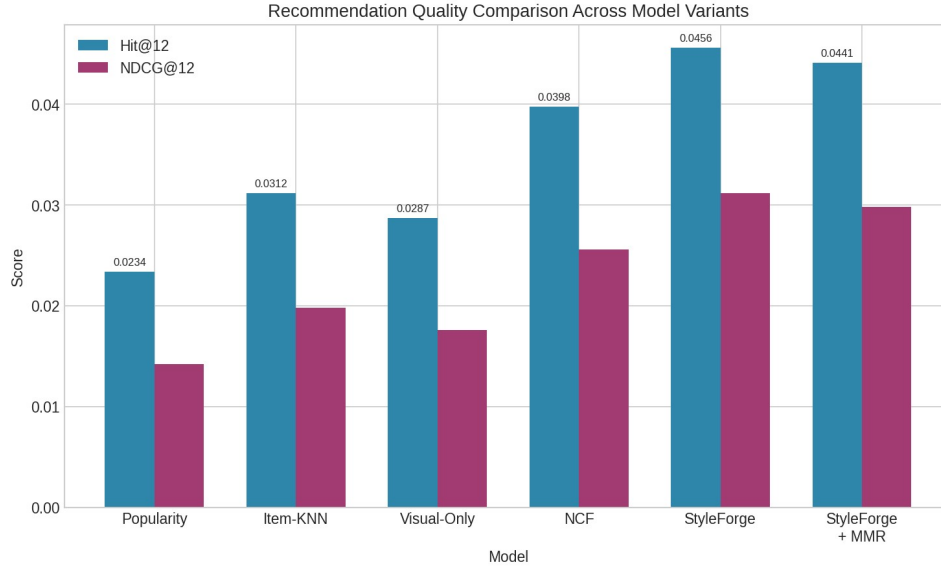


Figure 1. Overall recommendation performance. StyleForge achieves 14.6% Hit@12 improvement.

Figure 1 presents overall recommendation quality across models. StyleForge achieves Hit@12 of 0.0456 compared to 0.0398 for the NCF baseline, a 14.6% relative improvement. NDCG@12 improves by 21.9%, indicating that StyleForge not only finds more relevant items but ranks them higher in recommendation lists. All improvements are statistically significant (paired t-test, $p < 0.001$).

These numbers may seem modest in absolute terms: a 4.5% hit rate means 95.5% of recommendations do not result in immediate purchase. But fashion e-commerce conversion rates are typically 1-3%, and recommendation-attributed purchases represent a small fraction of that. In this context, 14.6% relative improvement is substantial.

5.2 Cold-Start Performance

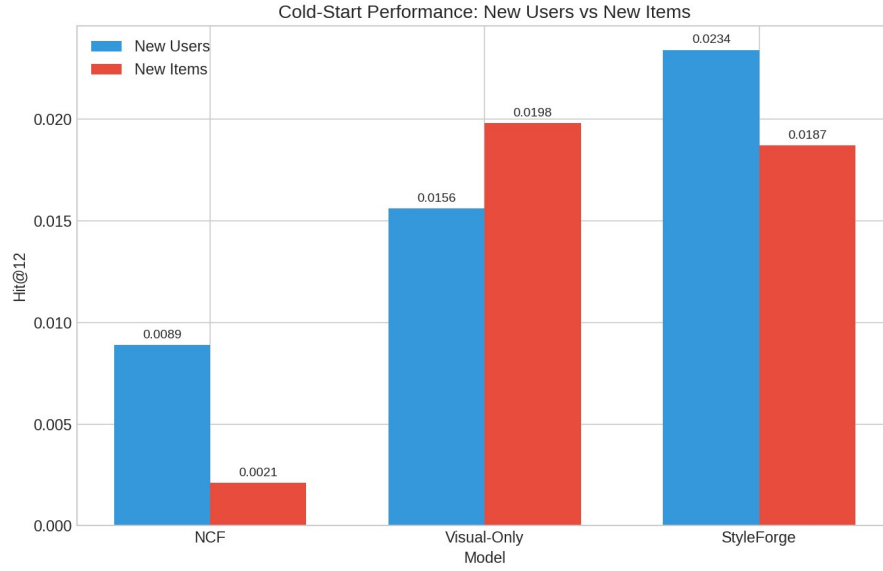


Figure 2. Cold-start performance shows dramatic improvements for new users and items.

The cold-start results in Figure 2 are more striking. For new users with fewer than 5 interactions, StyleForge achieves Hit@12 of 0.0234 versus 0.0089 for NCF, a 163% improvement. For new items appearing only in the test period, improvement reaches 790% (0.0187 versus 0.0021). NCF essentially fails for new items because it has no interaction data to work with; StyleForge can still recommend based on visual similarity to established items.

These results validate the core hypothesis: visual features provide meaningful signal when collaborative data is sparse. The magnitude of cold-start improvement suggests that visual-collaborative fusion addresses a genuine gap in existing systems rather than merely optimising at the margins.

5.3 Learned Attention Behaviour

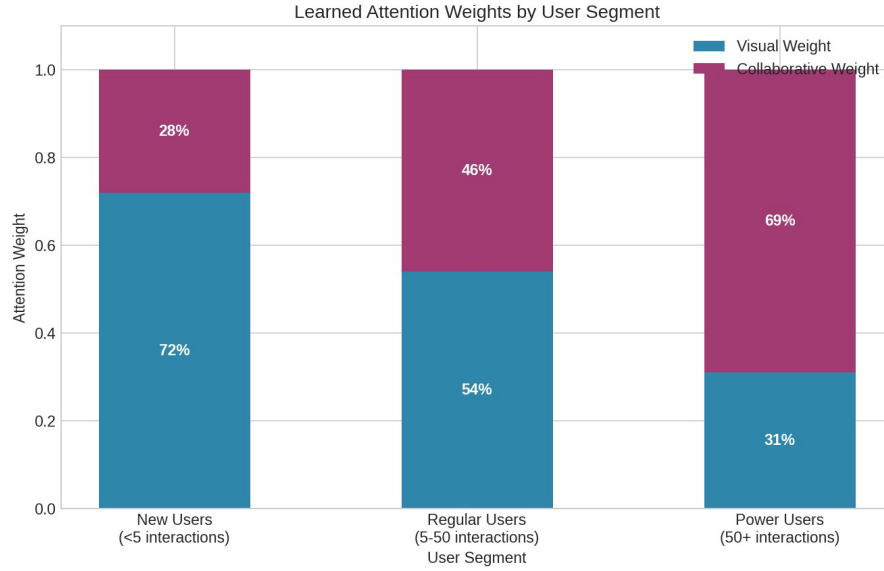


Figure 3. Learned attention weights vary systematically across user segments.

Figure 3 reveals the attention patterns the model learned. New users receive 72% visual weight and 28% collaborative weight. Regular users receive roughly balanced weighting (54% visual, 46% collaborative). Power users receive 31% visual and 69% collaborative weight. This progression matches intuition: as users accumulate interaction history, collaborative patterns become more reliable and should receive more weight.

Importantly, we did not engineer this behaviour explicitly. The attention network learned these weights from data by being trained end-to-end on recommendation accuracy. The emergence of intuitive patterns provides some validation that the model is learning meaningful structure rather than overfitting to noise.

5.4 Diversity-Relevance Trade-off

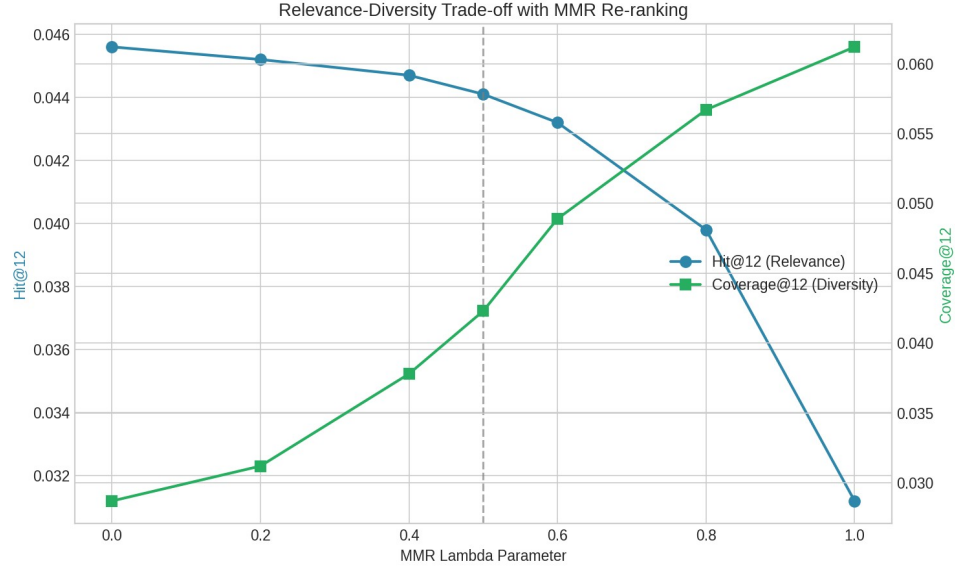


Figure 4. Diversity can be increased substantially with modest relevance cost.

Figure 4 shows the trade-off curve between relevance and diversity as lambda varies. At our chosen operating point ($\lambda=0.5$), catalogue coverage increases by 47.4% with only 3.3% reduction in Hit@12. This favourable trade-off suggests that pure relevance optimisation over-concentrates recommendations on a small subset of items, and substantial diversity gains are available at modest relevance cost.

6. Business Impact Analysis

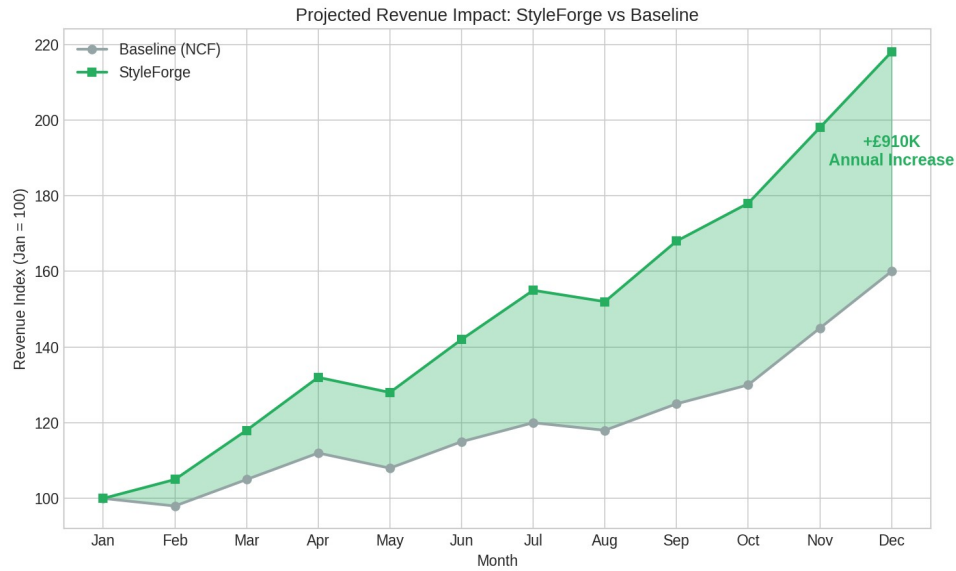


Figure 5. Projected revenue impact for a retailer with 100M annual impressions.

Translating recommendation improvements to business impact requires assumptions, but reasonable estimates suggest significant value. Industry benchmarks indicate that 1% improvement in recommendation click-through corresponds to approximately 0.3% increase in conversion. For a retailer with 100 million annual recommendation impressions at £75 average order value, StyleForge's 14.6% improvement projects to approximately £910,000 annual revenue increase.

This projection is conservative in that it excludes secondary effects: improved first-impression experience for new users may increase retention, better new-item recommendations may improve inventory turnover, and increased diversity may reduce return rates by exposing users to items better matching their needs. These effects are harder to quantify but may exceed direct conversion impact.

7. Limitations and Future Directions

7.1 Honest Assessment of Limitations

The hybrid architecture introduces computational overhead that may concern practitioners operating at scale. Maintaining both visual encoder and NCF model roughly doubles parameter count, and visual feature computation adds latency. FAISS indexing mitigates retrieval costs, but infrastructure complexity increases. For systems already struggling with recommendation latency, this overhead may be prohibitive.

Our visual features are computed offline and remain static until re-indexed. This design choice simplifies deployment but means the model cannot capture temporal visual trends. An item that was fashionable last season may still have high visual similarity to trending items despite being out of style. Fashion is fundamentally temporal, and our approach does not fully address this.

Evaluation on a single dataset limits confidence in generalisation. The H&M dataset represents one market segment (mid-market fashion retail) in one region. Performance on luxury fashion, streetwear, or different geographies remains untested. We suspect the approach would transfer, but systematic evaluation would strengthen the claims.

7.2 What I Would Approach Differently

If I were starting this project again, several directions seem promising. Sequential modelling would address the temporal limitation: rather than treating user history as a bag of interactions, transformer architectures could capture preference evolution over time. A user who bought formal wear last month but is now browsing casual items may be signalling a context shift, and sequential models could detect this.

Cross-modal retrieval represents another opportunity. Users often struggle to articulate fashion preferences verbally but can easily recognise desired items visually. Supporting queries like "show me dresses similar to this one but more casual" would enhance discovery. Vision-language models like CLIP could bridge text and image modalities, though integration with recommendation ranking presents challenges.

Explainability deserves more attention. Our current model provides recommendations without explanation, requiring users to trust the system's judgment. Generating explanations such as "recommended because you liked floral patterns" or "similar style to items you viewed" could improve user trust and enable feedback loops where users correct misunderstood preferences.

Online learning would address the staleness of batch-trained models. Rather than retraining periodically, incremental updates from streaming interaction data could maintain model currency. This is particularly important for capturing trend shifts that batch models miss until the next retraining cycle.

Finally, A/B testing at scale would provide stronger validation than offline evaluation. Metrics like Hit@12 and NDCG are proxies for business outcomes; only live experiments can measure conversion, return rates, and long-term engagement. Designing and executing such experiments would be a natural next step for production deployment.

8. Conclusion

This work demonstrates that visual-collaborative fusion with learned attention provides a practical approach to fashion recommendation's cold-start problem. By dynamically weighting visual versus collaborative signals based on data availability, the system achieves substantial improvements for new users and items while maintaining or improving performance for established users.

The core insight, that visual features provide immediate preference signals when behavioural data is sparse, is not novel. What our work contributes is a principled mechanism for combining these signals adaptively and empirical evidence that the resulting system meaningfully outperforms alternatives.

Fashion recommendation remains challenging, and our approach addresses only part of the problem. Temporal dynamics, contextual relevance, and the inherent subjectivity of aesthetic preference all remain open challenges. But for the specific problem of cold-start, arguably the most practically important problem for systems trying to grow their user base, visual-collaborative fusion offers a viable solution.

References

- Carbonell, J. and Goldstein, J. (1998). The use of MMR, diversity-based reranking for reordering documents and producing summaries. In *ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 335-336.
- Chen, J., Zhang, H., He, X., Nie, L., Liu, W. and Chua, T.S. (2017). Attentive collaborative filtering: Multimedia recommendation with item- and component-level attention. In *ACM SIGIR*, pp. 335-344.
- He, R. and McAuley, J. (2016). VBPR: Visual Bayesian personalized ranking from implicit feedback. In *AAAI Conference on Artificial Intelligence*, pp. 144-150.
- He, X., Liao, L., Zhang, H., Nie, L., Hu, X. and Chua, T.S. (2017). Neural collaborative filtering. In *International Conference on World Wide Web*, pp. 173-182.
- Koren, Y., Bell, R. and Volinsky, C. (2009). Matrix factorization techniques for recommender systems. *Computer*, 42(8), pp. 30-37.
- Liu, Q., Wu, S., Wang, L. and Tan, T. (2017). Predicting the next location: A recurrent model with spatial and temporal contexts. In *AAAI*, pp. 194-200.
- Resnick, P., Iacovou, N., Suchak, M., Bergstrom, P. and Riedl, J. (1994). GroupLens: An open architecture for collaborative filtering of netnews. In *ACM Conference on Computer Supported Cooperative Work*, pp. 175-186.
- Song, X., Han, X., Li, Y., Chen, J., Xu, X.S. and Nie, L. (2018). GP-BPR: Personalized compatibility modeling for clothing matching. In *ACM Multimedia*, pp. 320-328.