

Decoding movie genres with Vision

A Machine Learning Approach to Predict Movie Genres from Posters

CONTRIBUTORS

Ahmeda Cheick
Richard Lumpi
Averine Sanduku
Francesca Scipioni

DATASCI 281 - COMPUTER VISION

April 16, 2025

J.C. Berkeley School of Information

Project Introduction



OVERVIEW

We assembled a dataset of **29,265 IMDb movie posters** and built a **visual-only pipeline** that combines classic HSV & HOG descriptors with deep embeddings (ResNet50, ViT) for **genre classification**.



CHALLANGE

Posters come in wildly varying resolutions, often carry **multiple genre tags**, and exhibit **severe class imbalance**—making single-label, purely visual prediction **especially difficult**.



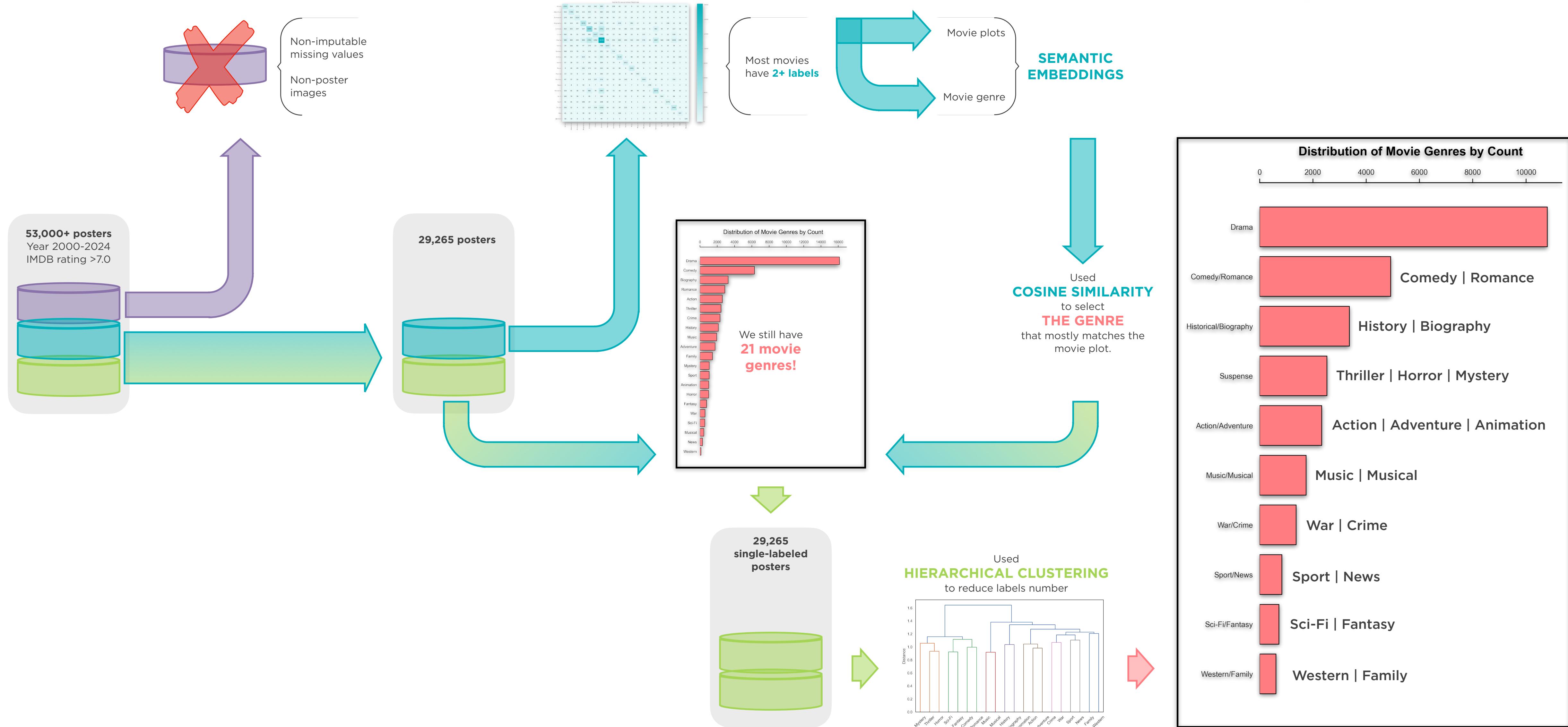
APPLICATION

Develop a **robust classifier** that distills each poster down to its dominant genre—balancing performance, speed, and implementation complexity across simple models (**Logistic/SVM**) and **fine-tuned CNNs**.

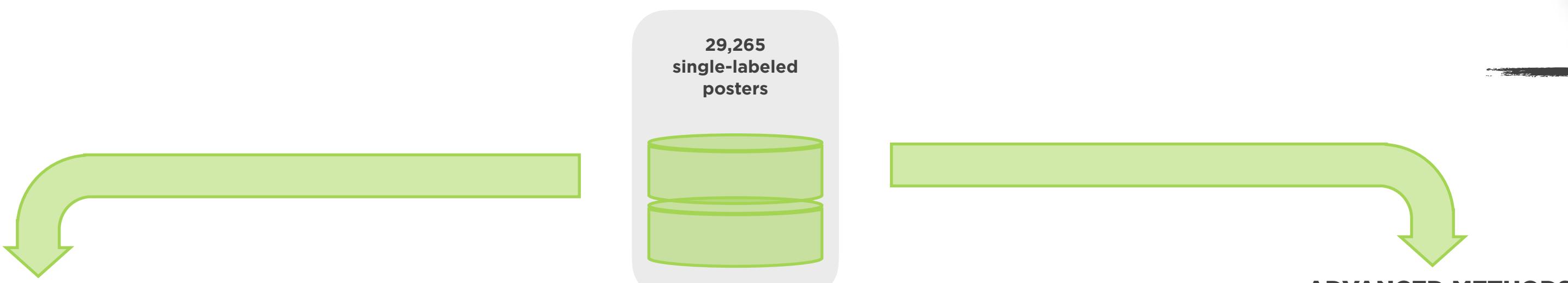


Enable **smarter content recommendation** and **automated cataloging** by inferring genre information directly from poster art—**no text metadata required**.

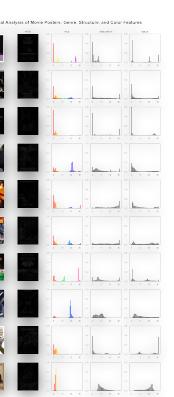
Data & Labels Preparation



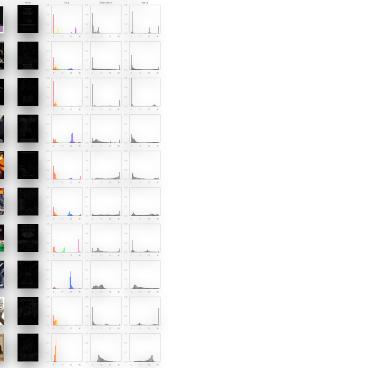
Extracting Features



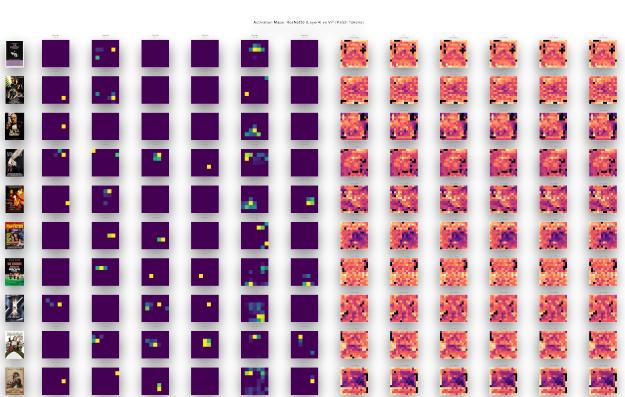
HSV Color Histograms and Moments



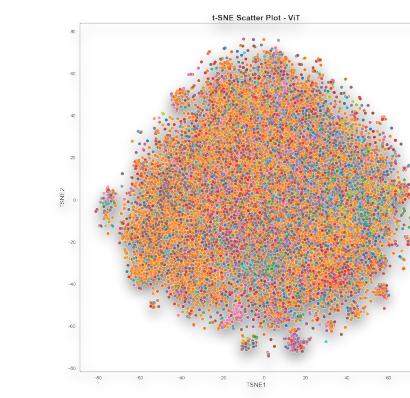
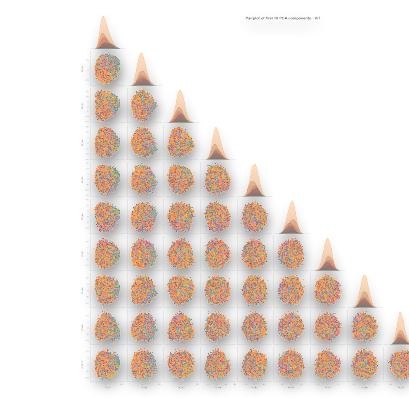
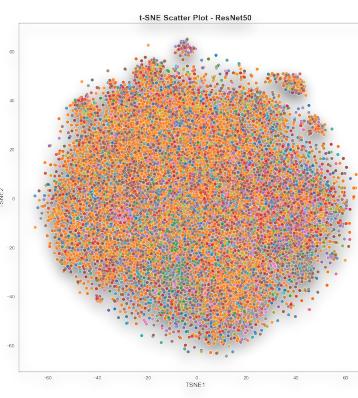
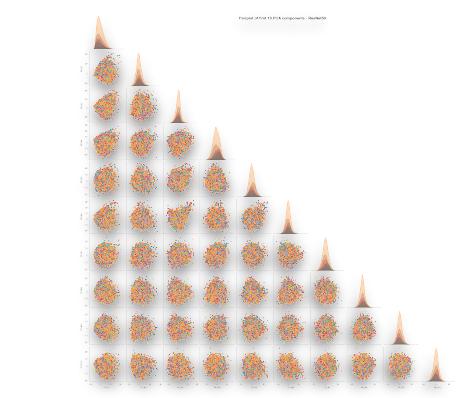
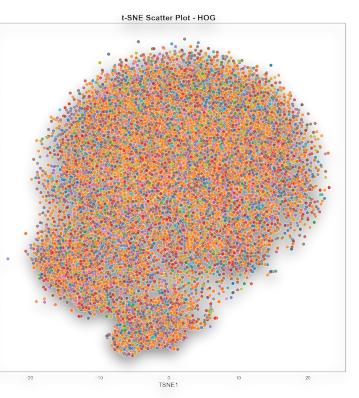
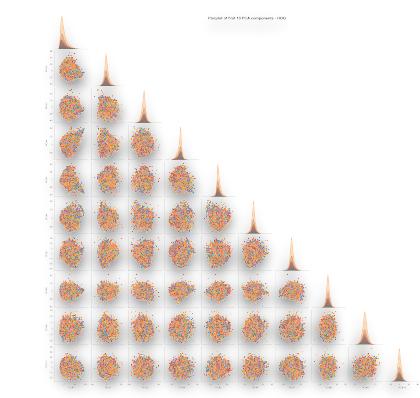
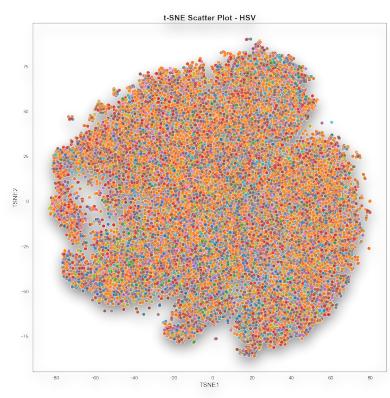
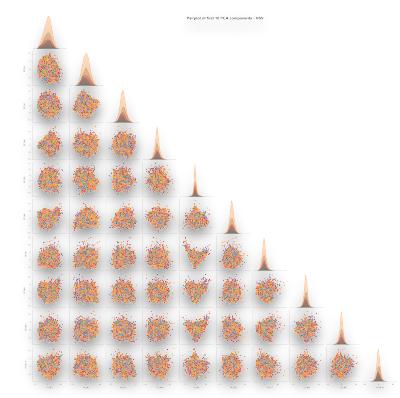
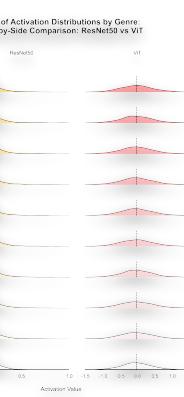
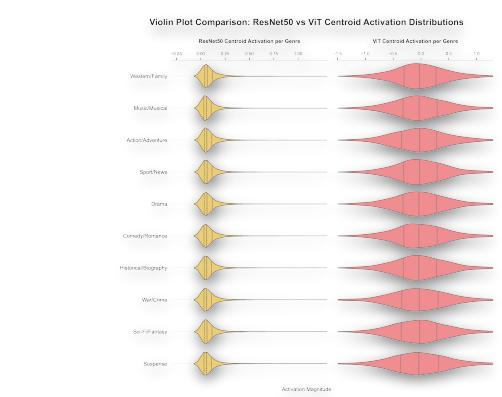
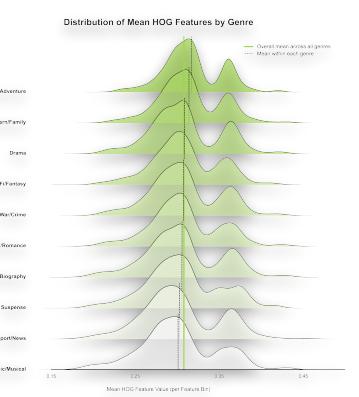
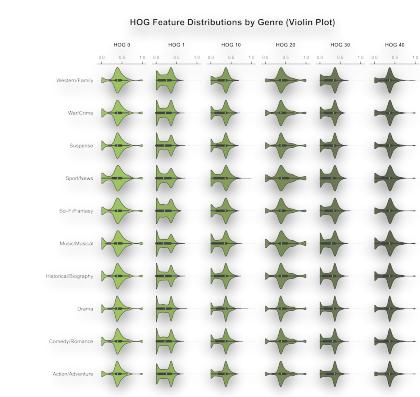
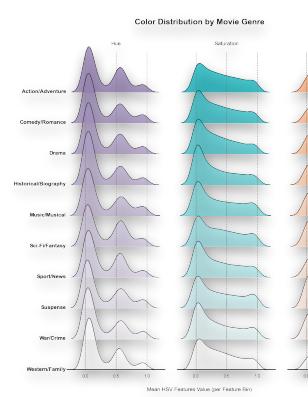
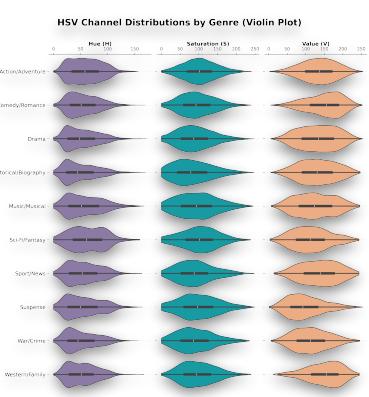
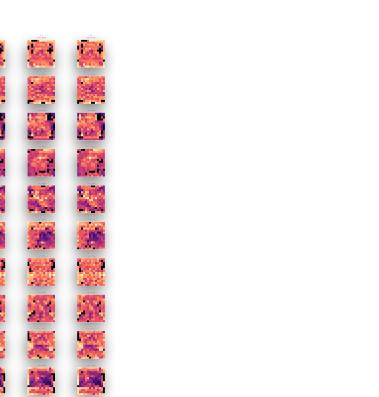
Histogram of Oriented Gradients (HOG)



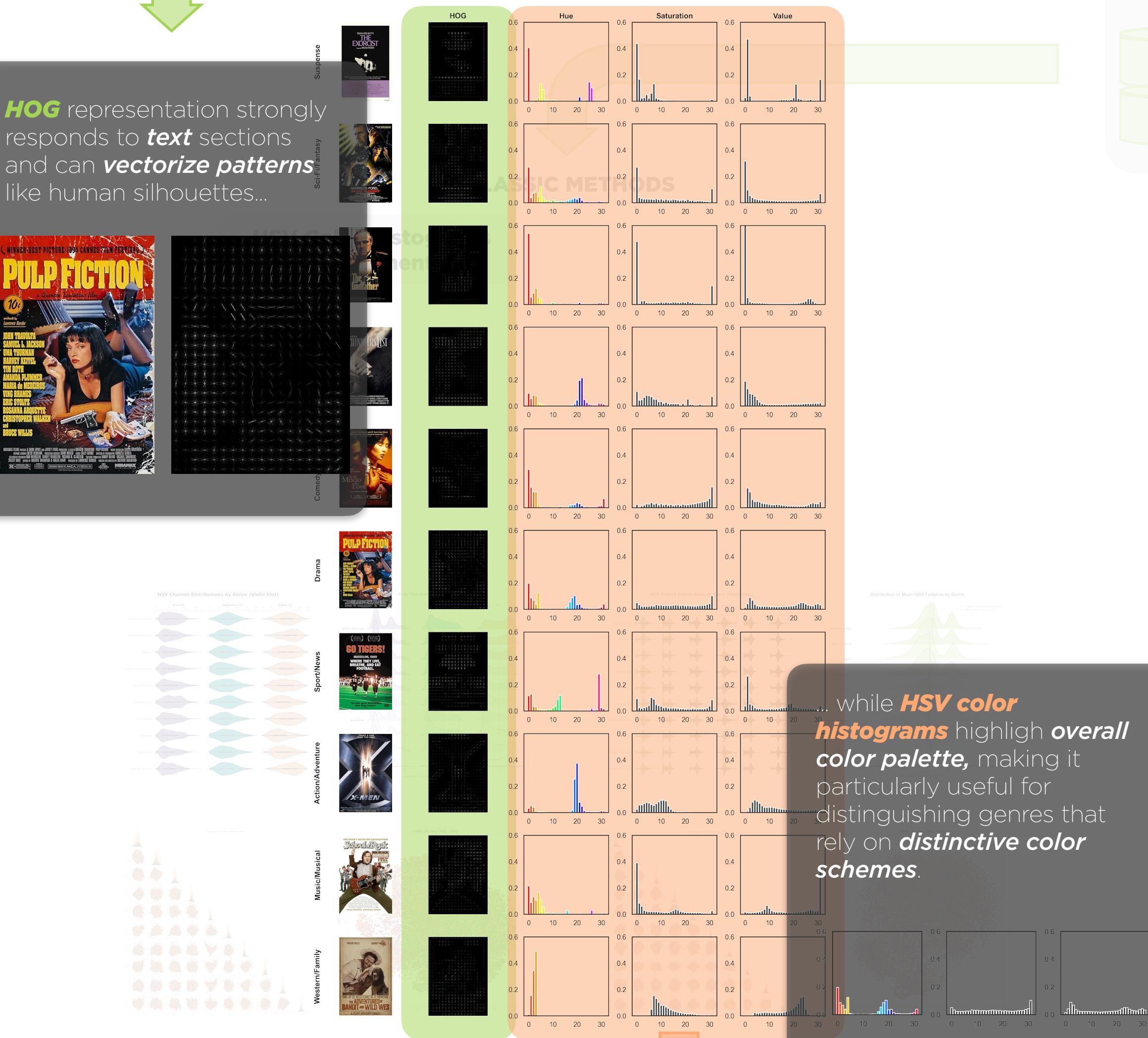
ResNet50 pre-trained CNN



Google's Vision Transformers (ViT)



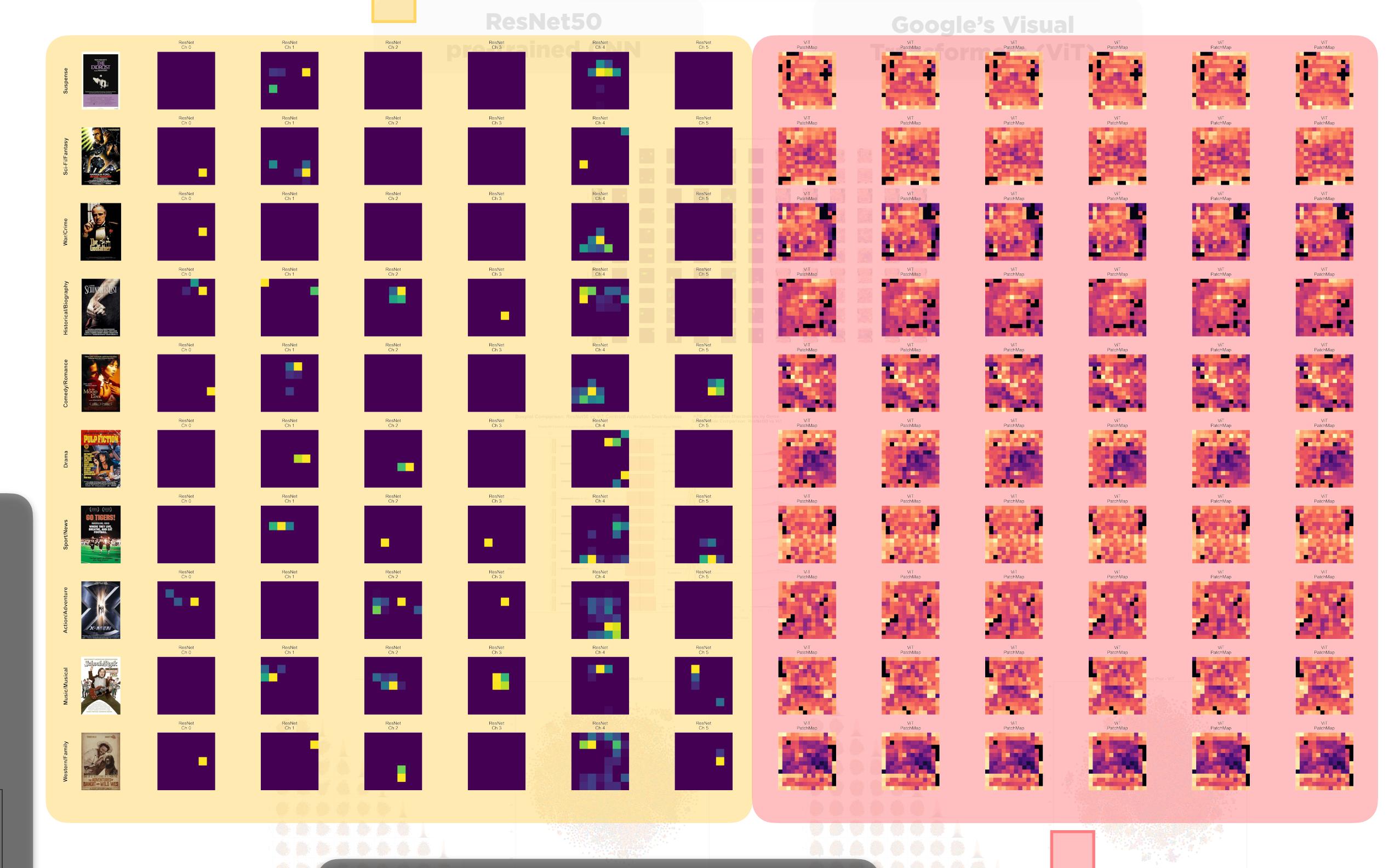
Extracting Features



29,265 single-labeled posters

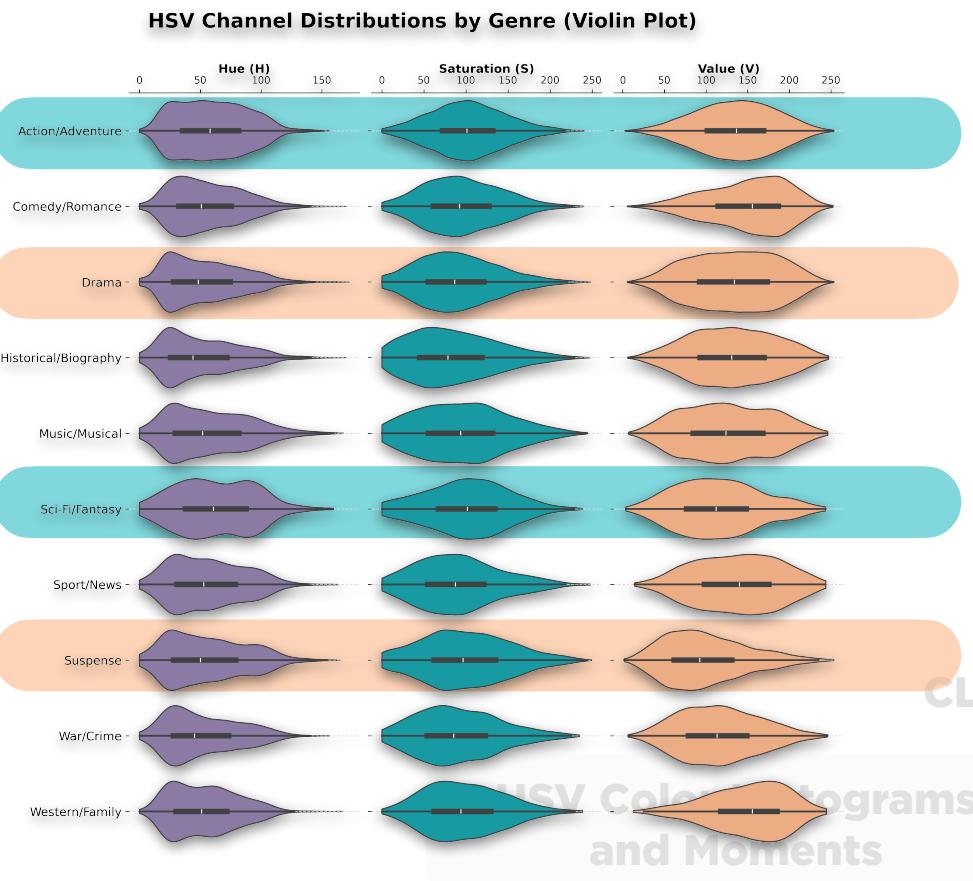


The **ResNet50's convolutional layers** focus on **local textures and shapes**—bright spots in each patch reflect strong feature responses to edges, text, or distinctive silhouettes



While the **ViT** patch tokens capture **broader contextual relationships** earlier in the network, using **self-attention** to integrate information across the entire poster

Extracting Features

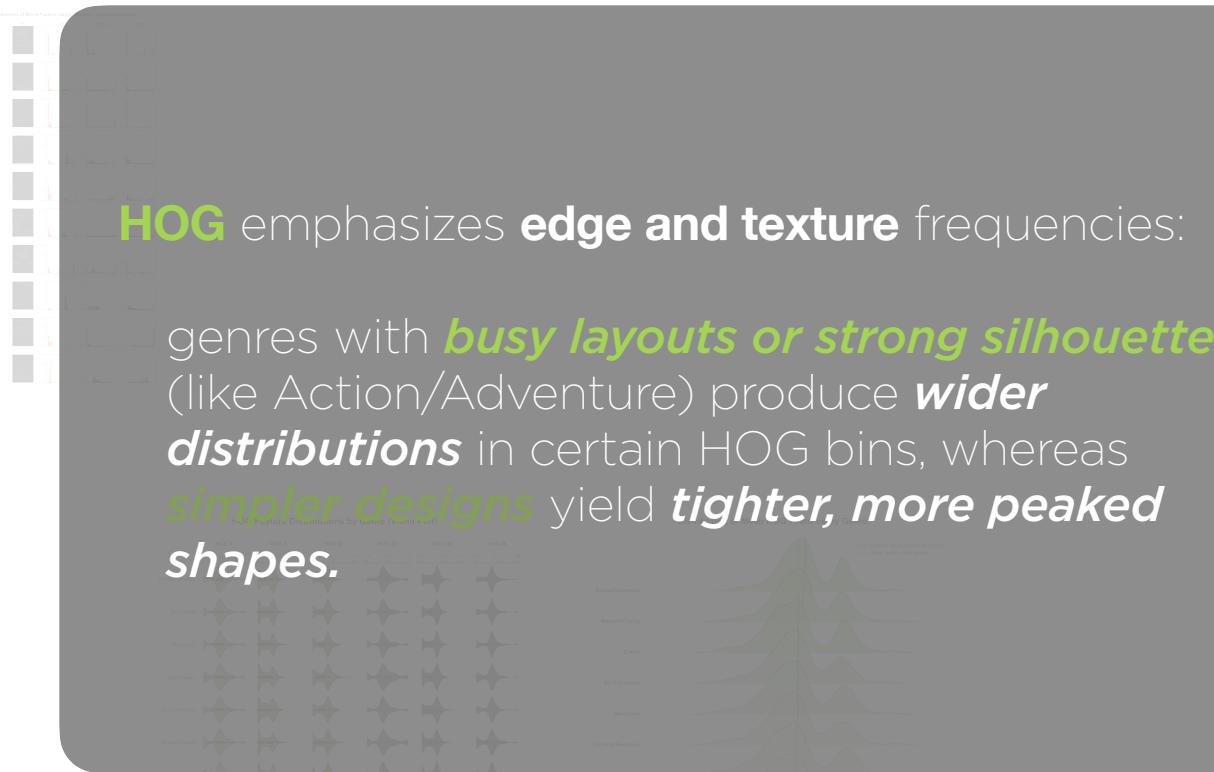


HSV captures color distributions:
genres with **bright, saturated palettes** (e.g., Sci-Fi/Fantasy) show distinct **peaks**, while more **muted categories** cluster around **narrower hue or value ranges**.

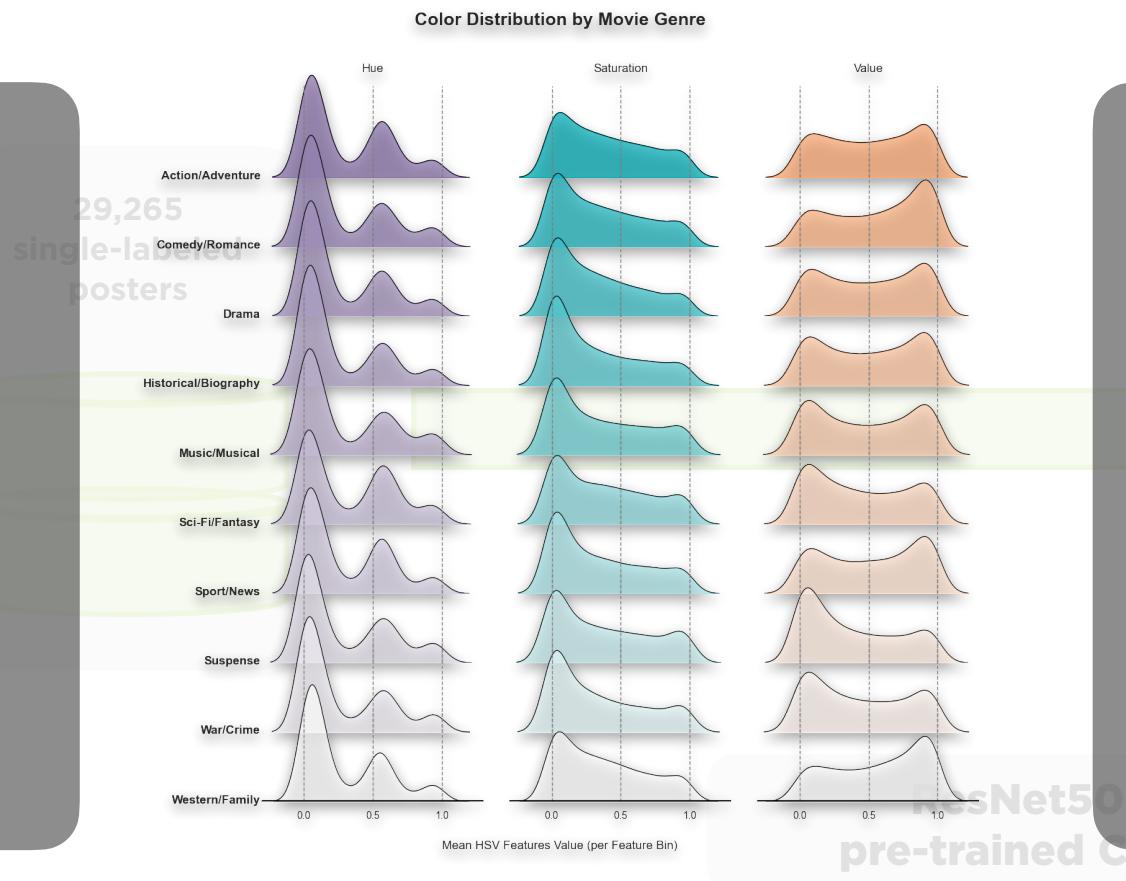
CLASSIC METHODS

Histogram of Oriented Gradients (HOG)

Violin Plot Comparison: ResNet50 vs ViT Centroid Activation Distributions



HOG emphasizes edge and texture frequencies:
genres with **busy layouts or strong silhouettes** (like Action/Adventure) produce **wider distributions** in certain HOG bins, whereas **simpler designs** yield **tighter, more peaked shapes**.

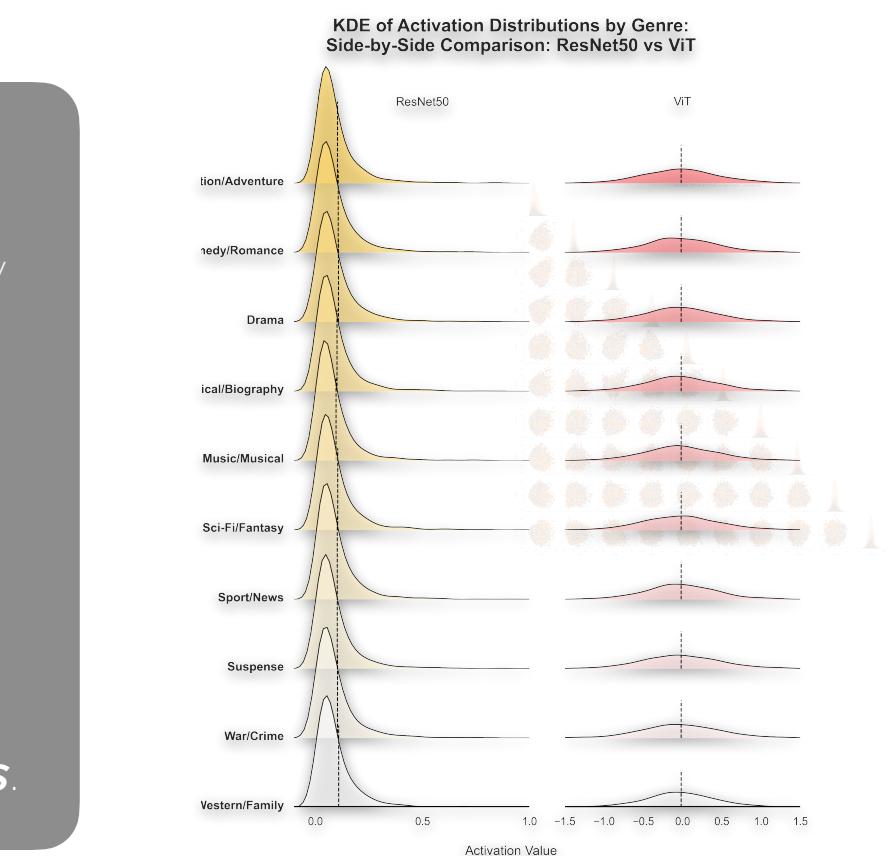
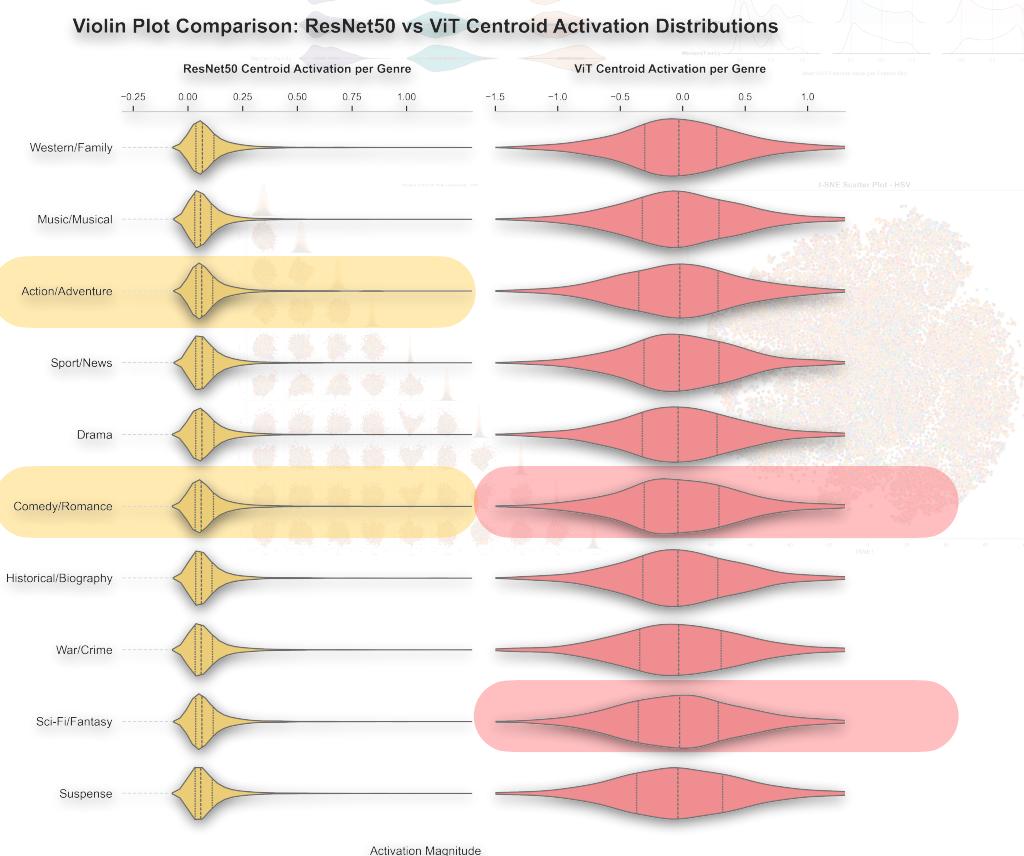
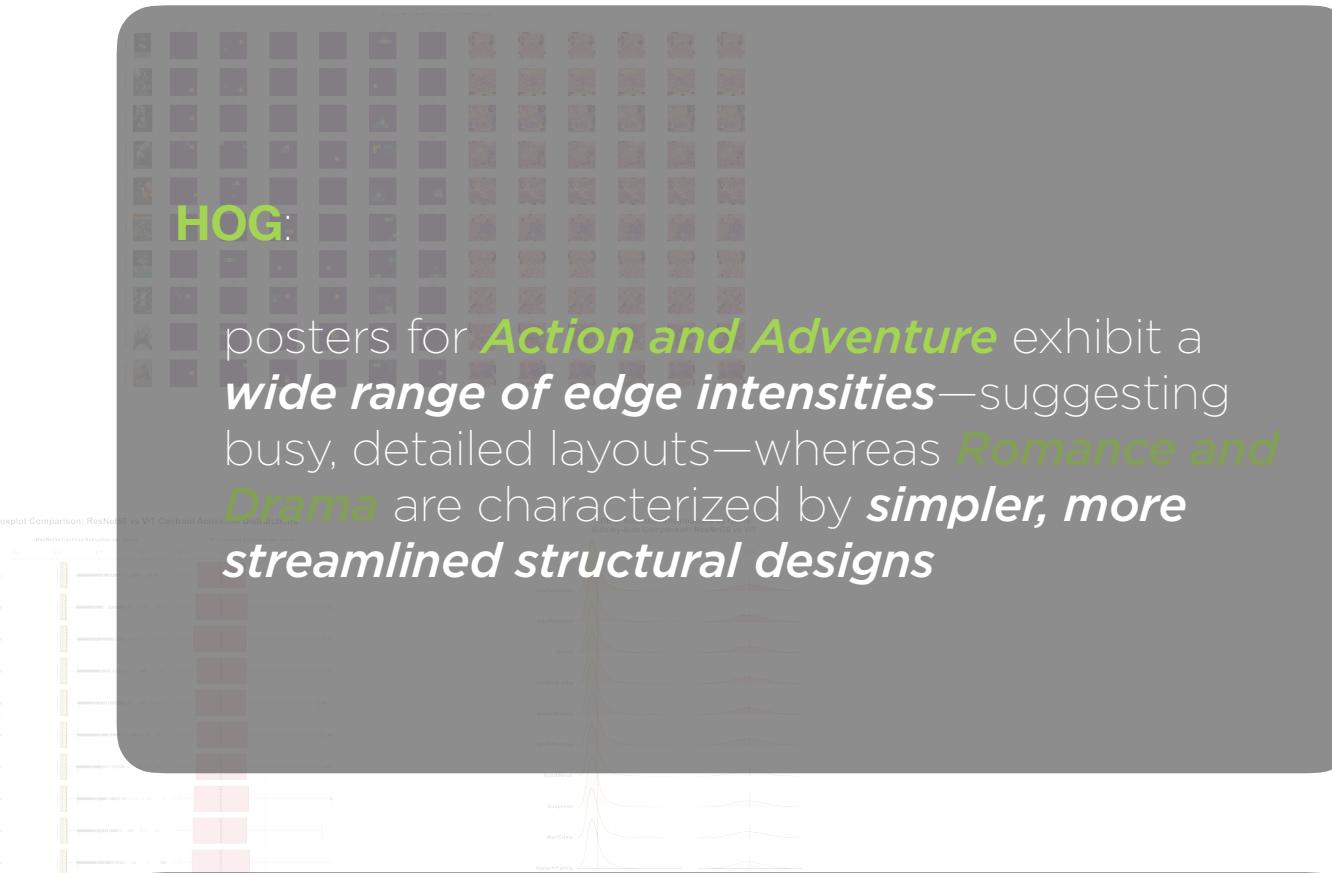
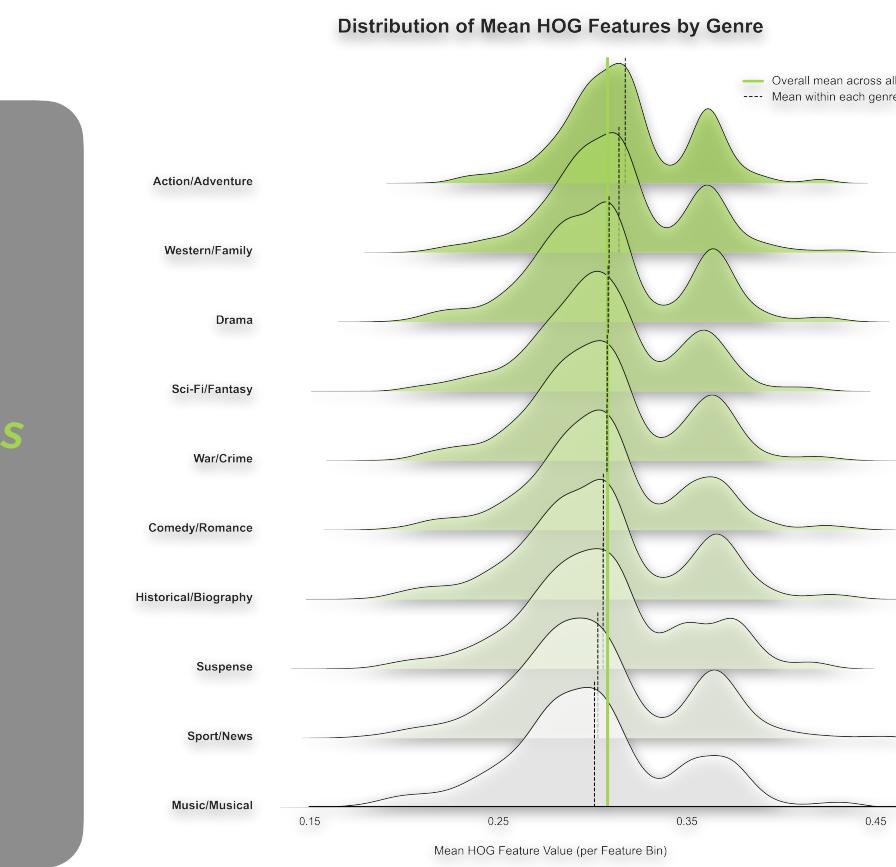


HSV reveals clear difference in color distributions:
genres such as **Comedy and Musical** burst with **vibrant, saturated hues**, while **Suspense and War/Crime** favor **muted, darker palettes**

ADVANCED METHODS

ResNet50 pre-trained CNN

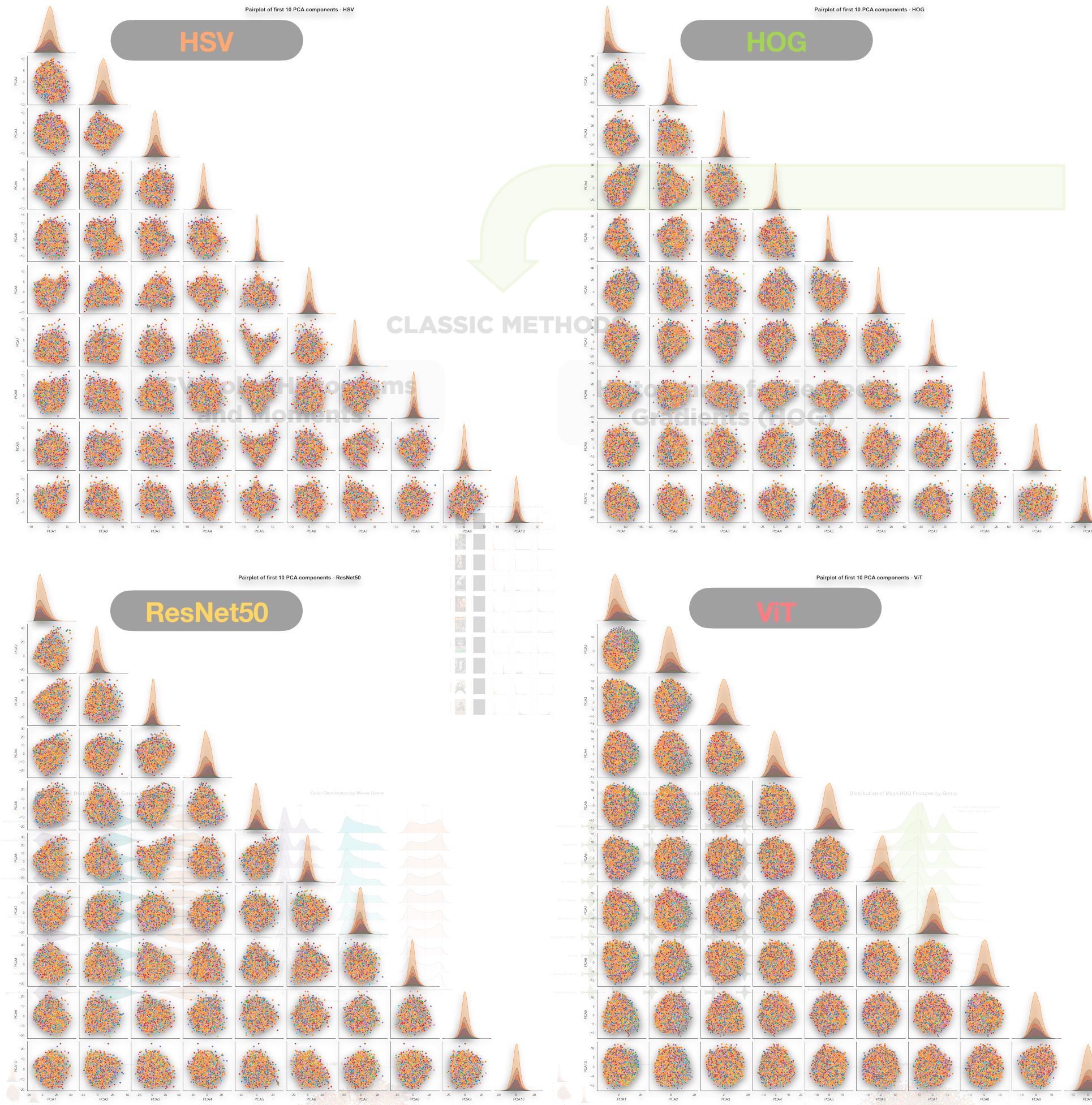
Google's Visual Transformers (ViT)



ResNet50:
focuses on local texture details, producing **sharp activation peaks in visually dense genres**, while **calmer genres yield lower, more concentrated responses**

ViT activations:
patch-based, global attention strategy generates **smoother, more uniform activation patterns**, yet still distinguishes genres.

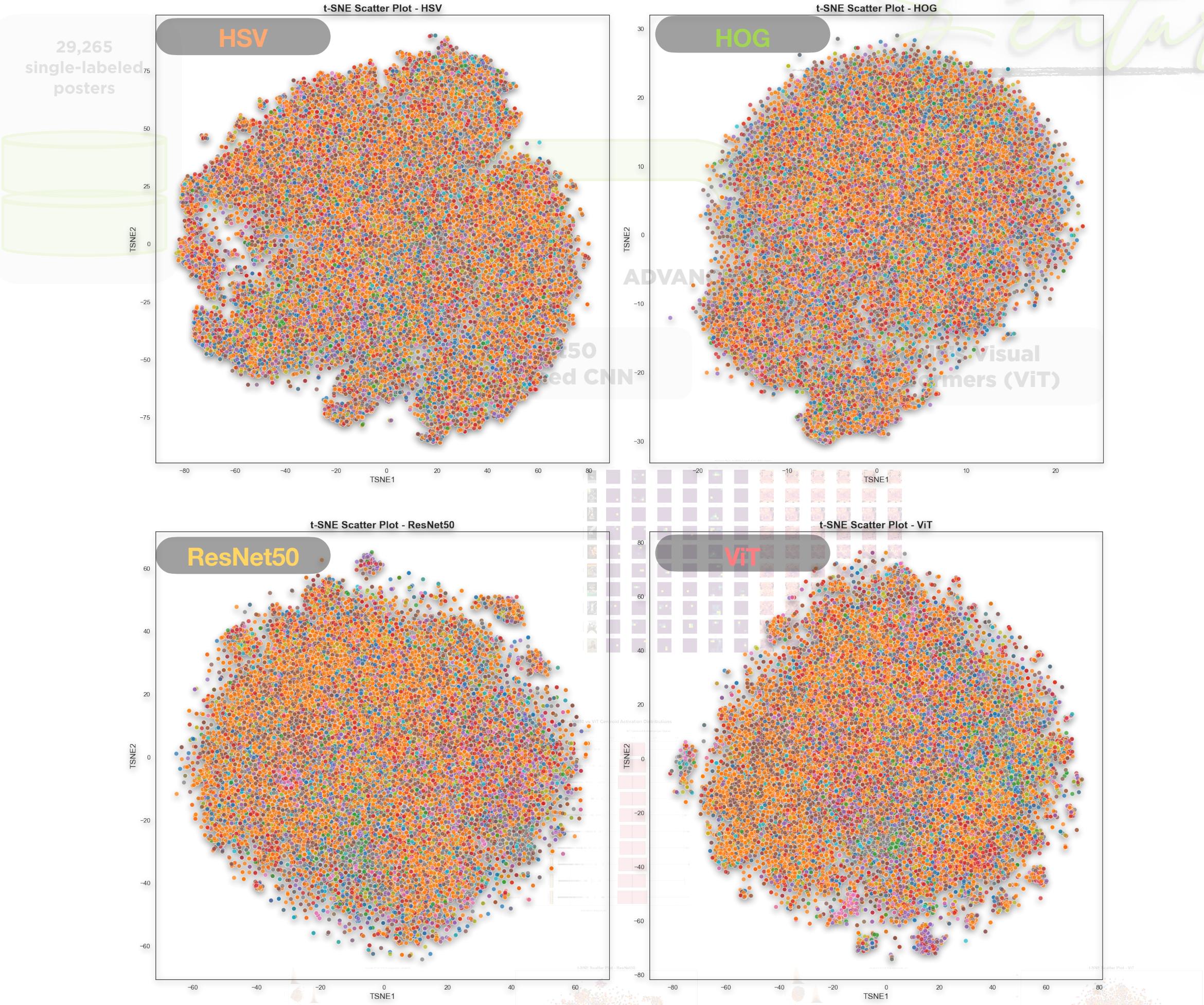
Pair plots of first 10 PCA components



HSV and **HOG** tend to produce **overlapping clusters**, indicating that their lower-level representations capture **some common visual characteristics across genres**.

In contrast, the deep features from **ResNet50** and **ViT** exhibit more **distinct clustering**, suggesting that their learned, high-level representations offer **better separability** among movie poster genres.

Scatter plot of t-SNE 2 components

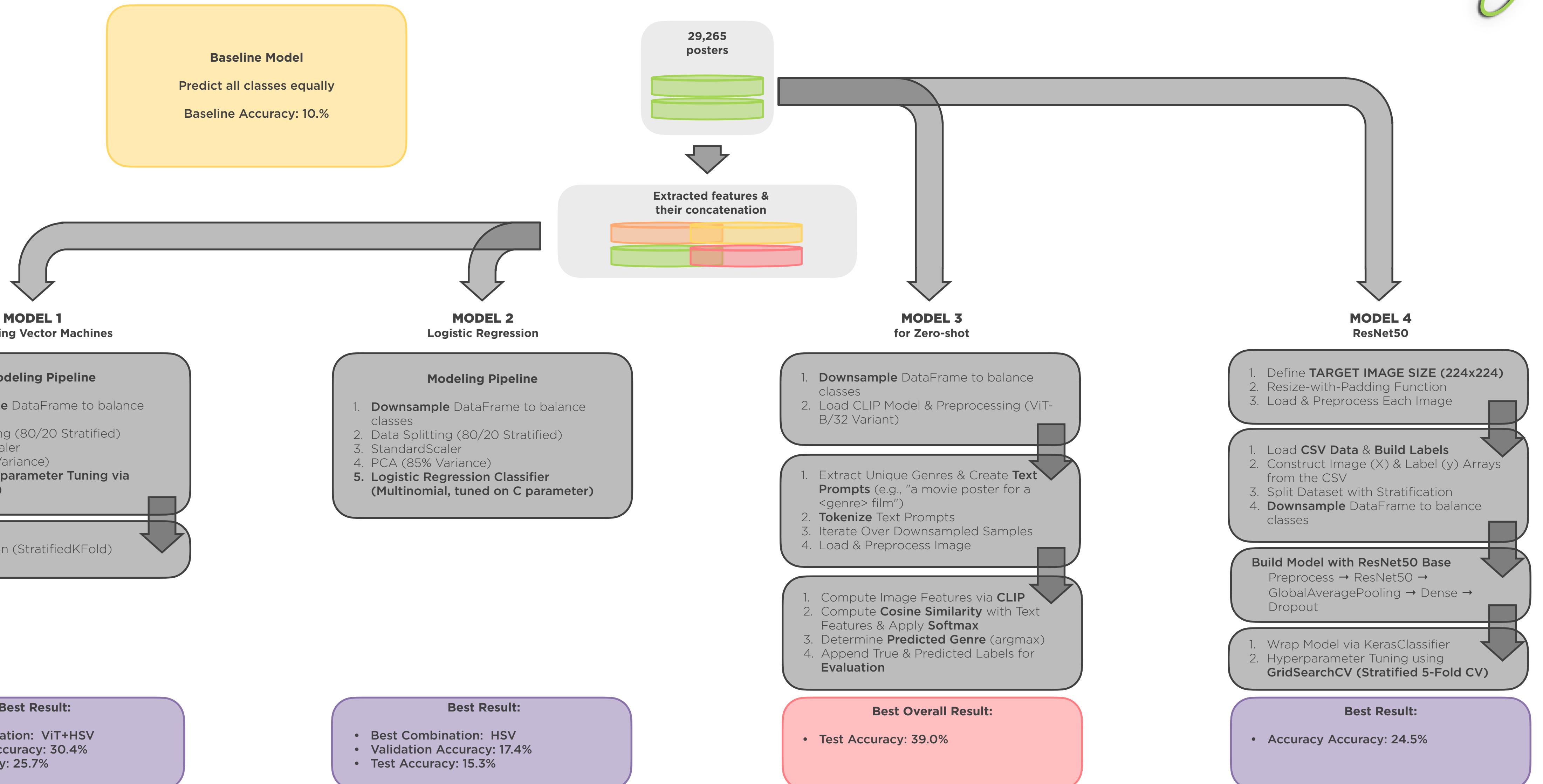


HSV and **HOG** yield relatively **diffuse clusters** with considerable overlap, suggesting that these low-level features **capture shared visual cues** across genres.

In contrast, the deep features from **ResNet50** and **ViT** form more **distinct, coherent clusters**, indicating their superior ability to **capture high-level semantic information** that differentiates movie poster genres.

Extracting
Features

Modeling



Modeling Results

	Action/Adventure	Comedy/Romance	Drama	Historical/Biography	Music/Musical	Sci-Fi/Fantasy	Sport/News	Suspense	War/Crime	Western/Family	Predicted Label
Action/Adventure	163	37	22	28	13	42	39	32	48	43	
Comedy/Romance	76	302	105	85	54	45	68	39	57	153	
Drama	132	333	367	220	134	99	119	201	226	330	
Historical/Biography	37	46	59	169	67	36	78	46	80	56	
Music/Musical	29	22	10	41	117	34	30	19	26	21	
Sci-Fi/Fantasy	24	8	9	14	9	41	7	17	10	7	
Sport/News	12	5	6	15	11	8	73	6	12	19	
Suspense	30	32	31	37	28	60	18	151	97	23	
War/Crime	26	17	22	35	16	19	18	32	78	11	
Western/Family	15	23	11	10	5	3	5	3	6	43	

 Drama is the most confused class, and the best predicted path the same time

 CLIP for Zero-shot is the best performing model

 Suspense and Comedy/Romance have good predictions across all the models

	Action/Adventure	Comedy/Romance	Drama	Historical/Biography	Music/Musical	Sci-Fi/Fantasy	Sport/News	Suspense	War/Crime	Western/Family	Predicted Label
Action/Adventure	68	65	29	34	17	41	28	67	32	86	
Comedy/Romance	82	250	70	67	34	54	62	109	55	201	
Drama	174	360	173	183	71	144	117	384	178	377	
Historical/Biography	57	83	41	87	37	40	57	108	68	96	
Music/Musical	30	37	23	32	15	45	28	82	26	31	
Sci-Fi/Fantasy	24	16	5	3	4	25	3	41	10	15	
Sport/News	18	27	13	13	9	8	12	30	13	24	
Suspense	34	43	29	27	26	48	21	191	60	28	
War/Crime	20	18	19	30	16	23	12	59	40	37	
Western/Family	19	25	4	8	3	6	12	5	6	36	

	Action/Adventure	Comedy/Romance	Drama	Historical/Biography	Music/Musical	Sci-Fi/Fantasy	Sport/News	Suspense	War/Crime	Western/Family	Predicted Label
Action/Adventure	101	9	111	44	160	13	75	57	34	16	
Comedy/Romance	9	88	169	45	179	6	36	66	11	11	
Drama	8	29	248	64	135	5	22	71	30	8	
Historical/Biography	11	10	101	171	162	8	54	47	48	8	
Music/Musical	10	10	52	72	410	10	19	18	14	5	
Sci-Fi/Fantasy	62	13	96	43	163	89	51	73	13	17	
Sport/News	25	4	34	31	81	10	366	21	38	10	
Suspense	36	10	97	39	110	17	22	242	35	12	
War/Crime	35	4	116	65	91	6	29	129	135	10	
Western/Family	10	19	195	51	172	4	35	57	11	66	

MODEL 3
CLIP for Zero-shot

	Action/Adventure	Comedy/Romance	Drama	Historical/Biography	Music/Musical	Sci-Fi/Fantasy	Sport/News	Suspense	War/Crime	Western/Family	Predicted Label
Action/Adventure	150	34	23	36	15	29	44	45	60	31	
Comedy/Romance	77	326	115	80	67	20	78	55	81	85	
Drama	175	374	401	206	84	34	129	232	341	185	
Historical/Biography	47	48	60	155	63	5	93	56	124	23	
Music/Musical	33	33	21	42	87	9	29	38	41	16	
Sci-Fi/Fantasy	34	11	12	15	7	15	11	25	11	5	
Sport/News	16	10	11	18	14	3	60	7	20	8	
Suspense	38	28	52	34	17	20	21	156	124	17	
War/Crime	26	21	23	32	17	6	13	45	82	9	
Western/Family	21	32	15	7	3	1	9	9	9	18	

MODEL 4
ResNet50

Conclusions



KEY QUANTITATIVE INSIGHTS

- **HSV / HOG baselines** extract in milliseconds but plateau at ~15% test accuracy.
- **Best “classical” pipeline** (ViT embeddings + HSV → linear SVM) reaches 26% test accuracy.
- **CLIP for Zero-shot** tops out at **~39% test accuracy** but requires ~20× more compute than the ViT+HSV SVM for equivalent performance.



WHY ACCURACY STALLS

- Purely single-label classification **discards multi-genre cues**—many posters straddle several categories—so recall on minority genres suffers and overall accuracy stagnates.



COMPUTER-VISION NEXT STEPS

- 1. Multi-Label Training**
Align the objective with real poster semantics by predicting all applicable genres.
- 2. End-to-End Transformer Fine-Tuning**
Unfreeze and train larger ViT variants (e.g. ViT-L/16) on poster data for richer, domain-specific representations.
- 3. Genre-Aware Augmentation**
Use adaptive color jitter, random erasing, or GAN-based style transfers to simulate varied poster styles.
- 4. Domain Adaptation & Metric Learning**
Emphasize genre-specific visual cues via adversarial domain alignment or triplet-loss training.
- 5. Multimodal Fusion**
Combine poster imagery with plot embeddings, keywords, or trailer frames to push accuracy beyond 30-40%.



BOTTOM LINE

- Deep embeddings unlock gains, but breaking the 30% barrier on single-genre poster classification will demand richer labels, smarter augmentations, and seamless integration of multiple modalities.