Leveraging GenAl for Authenticity and De-Duplication in CPG Retail Image Data







· Who are we?

Problem Discussion

Some Solutions

Case studies from Industry

Extensions





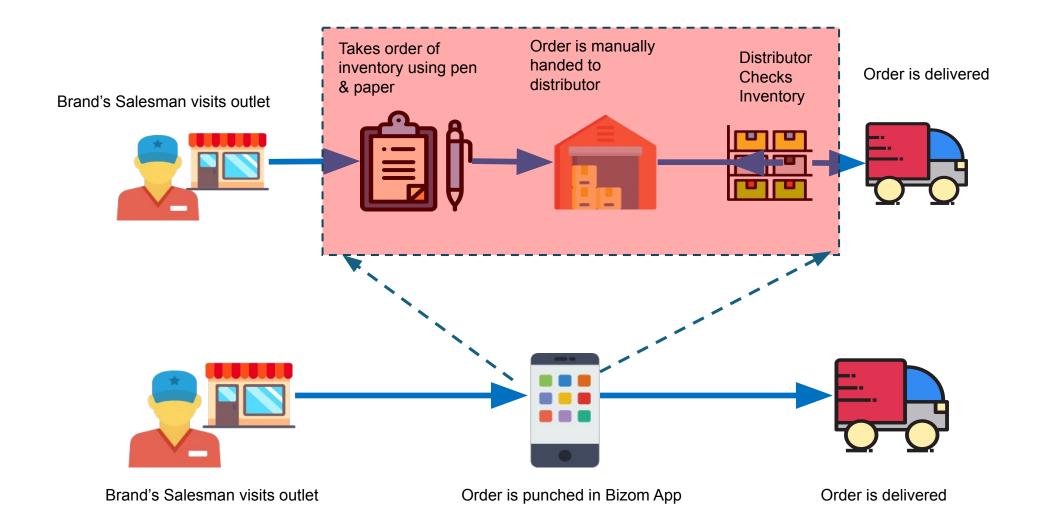
About us @ Bizom



"Bizom is a Sales Force Automation software that helps in digitizing the steps needed for manufacturers (brands) to place their products into various retail shops/outlets (Kiranas)"



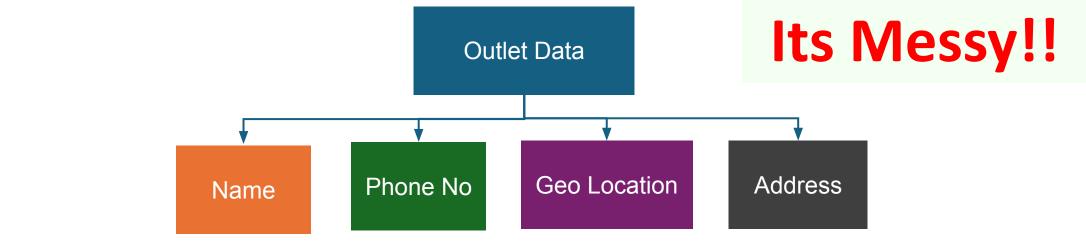
Bizom Simplifies Supply Chain







Retail/Outlet Data

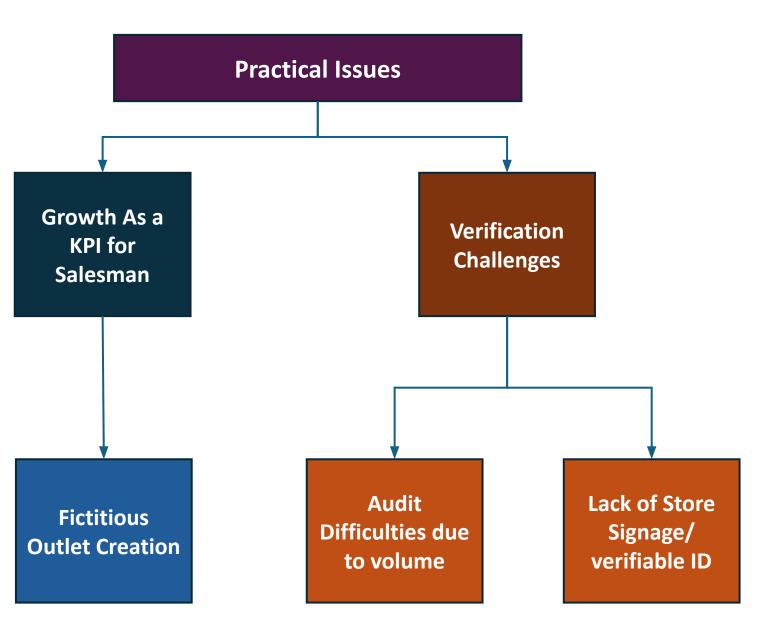


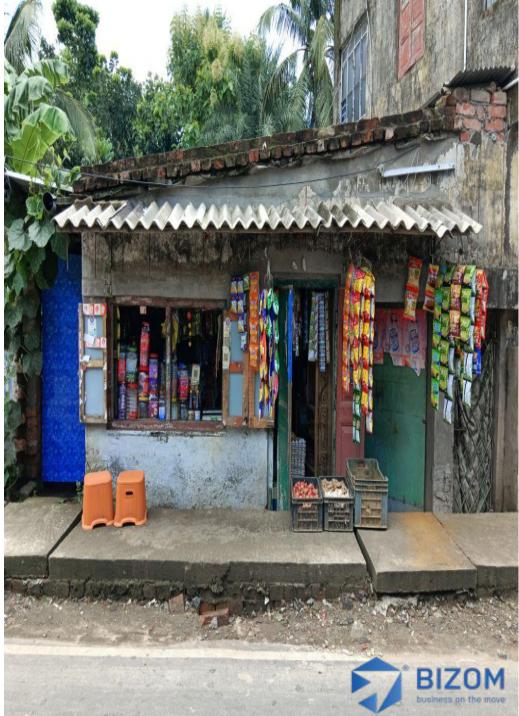


About 10Mn shops look like this !!



Why this mess keeps growing?

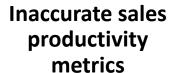




Impact of the inflated Outlet Universe









Distorted sales force deployment



Wasted sales effort & resources



Misaligned targets & incentives

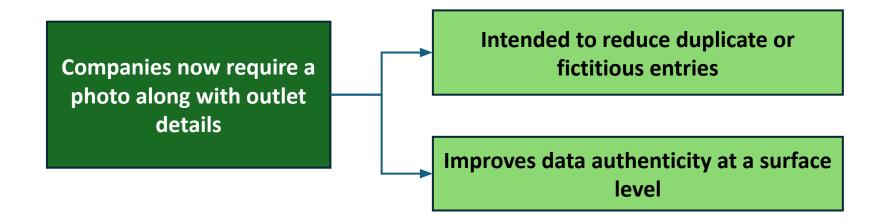


Unreliable reporting



Workaround: Outlet images as part of outlet data





Images can still be reused or manipulated unless further validation measures are in place







Key Challenges in Outlet Image Verification



Rather than resolving the issue, the images have only compounded the problem !!!

Problem #1

How to ensure the uploaded image is genuinely of an outlet, and not unrelated content?

Problem #2

How can we identify if the same image is being reused across different outlets?









Problem #1: Find out if the image represents an outlet

Attributes of a good outlet image:

- Signage visible
- Products Visible
- Store Open
- Day time photo
- Exterior showing the entire outlet
- .







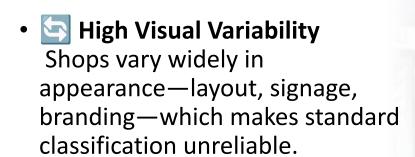








Why Conventional Object Classification Approach fails?



Lack of Training Data
 No consistent dataset exists to train a model that can generalize well across all shop types.



















GenAl Characteristics Useful for Retail Domain

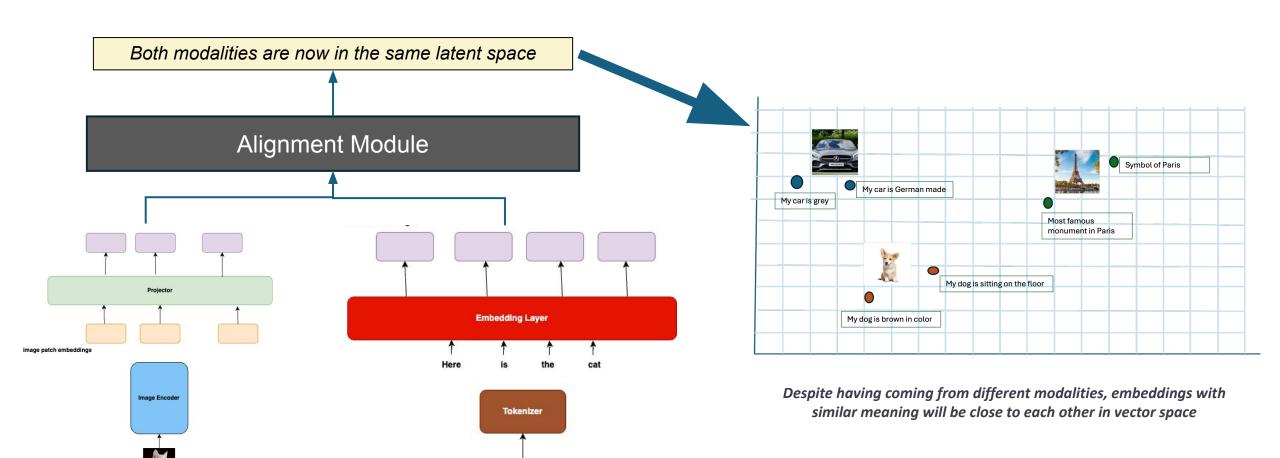
Ability to generate data

GenAl algorithms like GPT are known to generate text in case of Retail it could be Product Descriptions, Report Summaries

Ability to capture context

GenAl algorithms can analyze text, images to create a mental model of various internal relationships in the data for e.g. how the words in a sentence are related to each other, how different pixels in an image are related etc.

Multimodal LLM

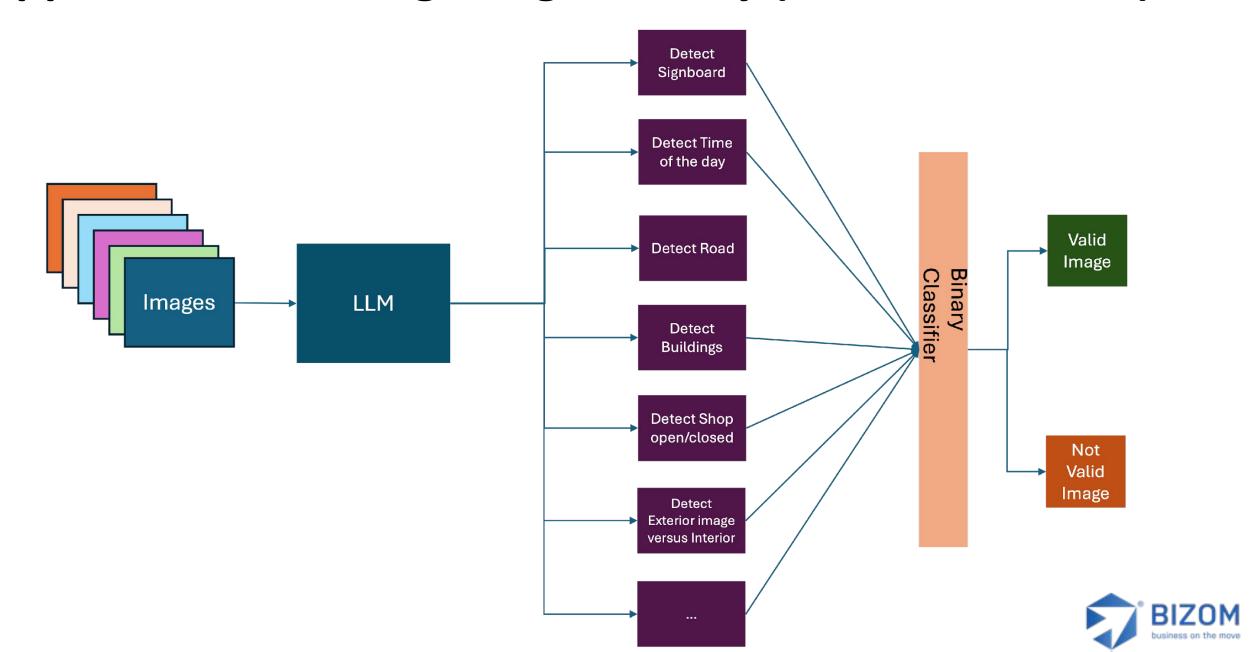


"Here is the cat"

Both the visual and text embeddings can interact now. They can be concatenated together and pass to a LLM to perform some task



Approach for finding Image Quality (Pure LLM based)



MLLM output on Images

fileName	isHangingPro	duct	number_people	isShopOpen	number_buildings	isExteriorVisible	isRoadVisible	timeOfDay	isSignVisible	image
outletData/22726	629.jpg	no '	There are two people near or inside the shop: a man and a woman.	yes	2	exterior	yes	afternoon	yes	
outletData/9431	47.jpg y	yes	There are two people near or inside the shop.	Yes	1	exterior	no	afternoon	no	
outletData/8120	988.jpg	no	There are no humans visible near or inside the shop.	no	1	exterior	no	afternoon	yes	
outletData/2272	2626.jpg	no	There are two people near or inside the shop.	no	5	exterior	yes	afternoon	no	
outletData/17770	097.jpg y	/es	There are two human beings in the shop: a man and a woman.	Yes	1	interior	no	afternoon	no	
outletData/87548	88.jpg		There are no human beings visible in the image.	no	0	exterior	no	night	no	
outletData/22726	622.jpg y	/es	There is only one person, a man, near or inside the shop.	yes	1	interior	no	afternoon	no	



Issues with pure LLM approach

Inconsistent Accuracy

Highly
Dependent on
Image Quality

Lack of Local Context

Extremely Slow at Scale



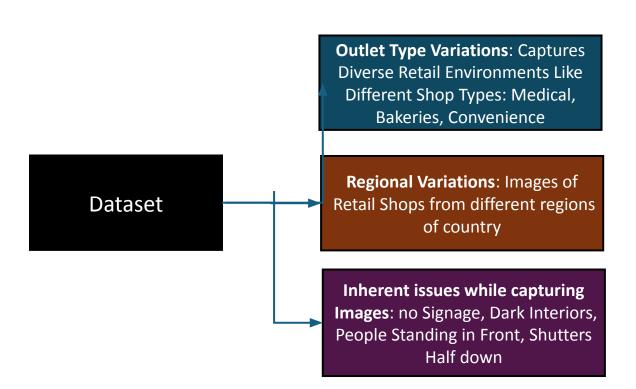
Our Approach: Two stage process

Stage1: **Using MLLMs Generate Pseudo-labels** which act as Training Data for image classifier Stage2: Train a Vision Model like ResNet or ViT to act as an image classifier



Training Setup

Data Set	#of Images
Training Set	8000
Validation Set	1000
Test Set	1000









Dataset Images



Stage1: Pseudo-Label Generation

Domain Specific Prompt Engineering

Create prompts to understand if the image is of a Kirana store.

Prompt Examples

"Is the store board clearly visible and not occluded?"

Multimodal Evaluation

The same input (image + prompt) is fed into three pre-trained MLLMs

Model Ensembling

Each model independently **assesses alignment** between the image and prompt and makes **Binary Decision Output**

Positive – Image aligns with the prompt

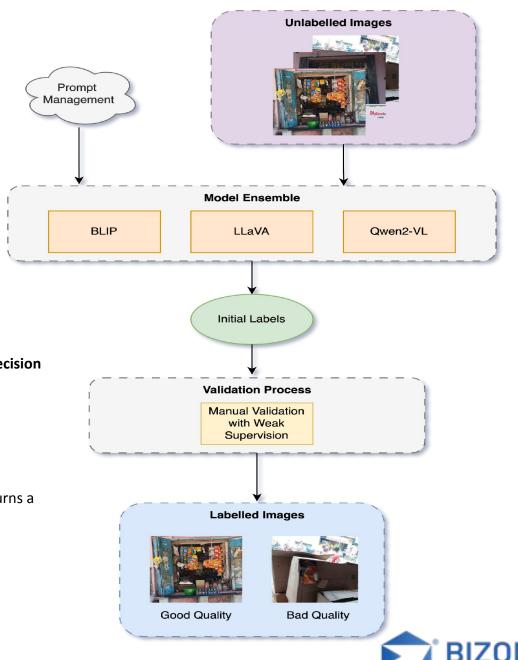
X Negative – Image does not align with the prompt

Unanimous Voting

An image is labelled positive only if all three models agree it is of good quality. If even one model returns a negative decision, the image is labelled negative .

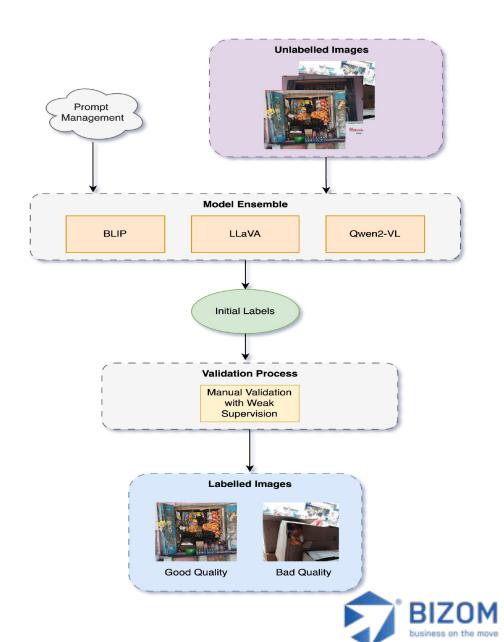
Manual Validation

A quick manual validation is done to ensure that we pass only the correctly labelled images to the classifier training stage

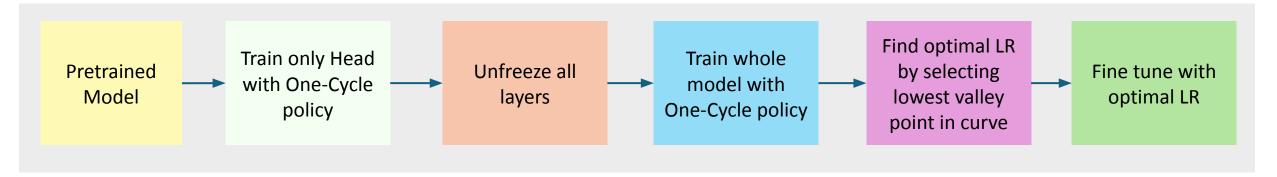


Stage1: Pseudo-Label Generation LLMs contd...

Model	Accuracy	Precision	Recall	F1-score
Blip	0.4869	0.3569	0.9045	0.5118
LLaVA	0.3420	0.3085	0.9771	0.4690
Qwen2-VL	0.7635	0.5587	0.9744	0.7102
Our Approach	0.8066	0.6236	0.8816	0.7305

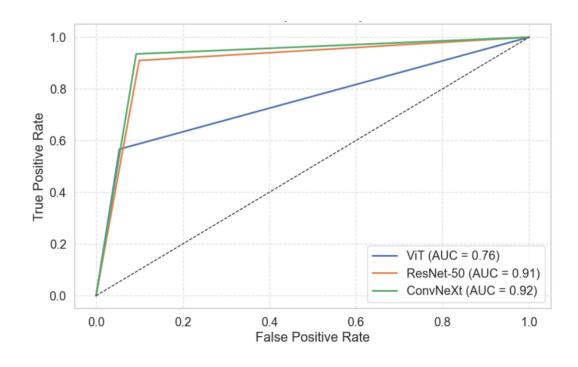


Stage2: Training of Vision Classifier



Model	Accuracy	Precision	Recall	F1-score
ConvNext	0.9150	0.9238	0.9150	0.9170
ResNet-50	0.9030	0.9117	0.9030	0.9052
ViT	0.8410	0.8370	0.8410	0.8316

- ConvNeXt yields the best classification performance but comes with higher computational requirements.
- **ResNet-50** offers a strong balance between performance and efficiency, aligning with real-world deployment constraints.





Results







Not Acceptable Images

Acceptable Images





Advantages of the approach



Cost and computation efficient solution



Achieves >90% accuracy on 1MN outlet images for an enterprise





GPU used only for Pseudo Label creation on a small data sample



Classifier runs entirely on CPU post-labelling



10x faster than GPU-based alternatives



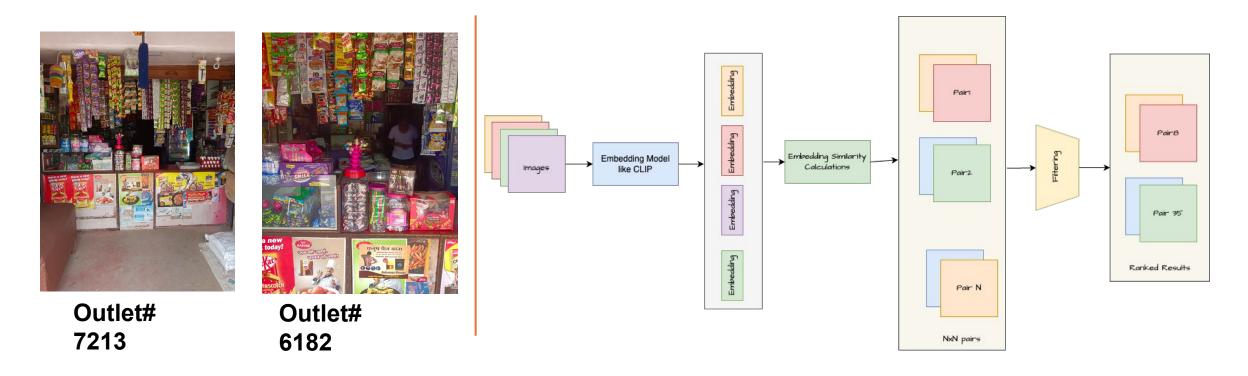
Significantly reduced manual effort needed for training data generation



Automated pipeline ensures speed and reliability



Problem #2: Find out if same image has been used for another outlet



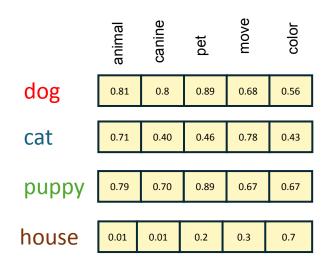
- We employ LLM model to create the contextual vectors from the image data
- The Algorithm is quite fast and finding out duplicates among 10K images took less than 1 minute (NxN comparisons)

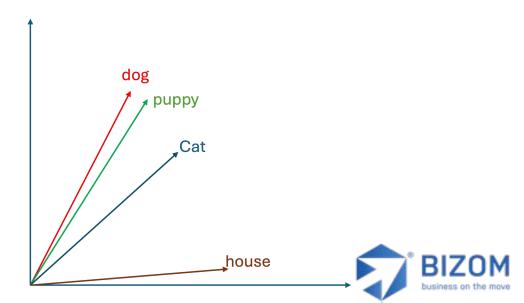


What is an Embedding?

Embedding is the numerical representation of the words (tokens) such that they can be used in the machine learning models. They are also called vectors as they are 1-D arrays of real numbers

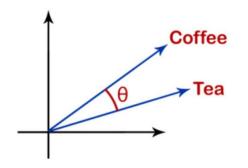
- Embeddings are vectors that allow mathematical operations like computing distances to measure similarity.
- S By default, embeddings are independent and lack relational meaning.
- Neural networks enrich embeddings with contextual understanding, enabling intelligent language-based workflows.
- Dimensionality reduction helps capture **intrinsic relations** between words and grammar.

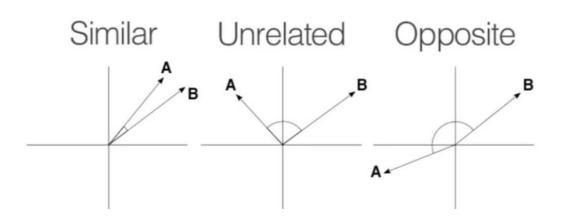




Similarity Measures

• Since Embeddings are vectors we can measure the cosine of the angle between the two vectors to find out the similar and dissimilar vectors





Sentence 1	Sentence 2	Cosine Similarity
The cat sits outside	The dog plays in the garden	0.2838
A man is playing guitar	A woman watches TV	-0.0327
The new movie is awesome	The new movie is so great	0.8939
Jim can run very fast	James is the fastest runner	0.6844
My goldfish is hungry	Pluto is a planet!	0.0454



Our analysis of Different Embedding models

Algorithm	Precision@0.95	Precision@0.90	Precision@0.85	Recall@0.95	Recall@0.90	Recall@0.85
CLIP-VIT-L-14	1	0.8	0.2	0.95	0.93	0.88
Siglip-VIT-B-16	1	0.81	0.12	1	0.98	0.95
ClipSegrd-64-refined	1	0.98	0.82	0.84	0.77	0.6
DINO Base	1	0.99	0.88	0.55	0.5	0.4

DINO Base:

High Precision even at lower thresholds.

Very low recall, risking loss of many valid cases/scenarios

No clear threshold for **TP vs. TN separation**

ClipSegrd-64-refined:

Lower recall compared to CLIP.

High Precision maintained even at a low threshold of 0.85.

No clear threshold for **TP vs. TN separation**

CLIP & SigLip

Both models show high accuracy and recall at a **0.95 threshold**.

Accuracy drops significantly at lower thresholds.

Only high similarity scores are reliable for determining similarity.

Clear threshold for **TP vs. TN separation**



DINO: 0.82, CLIP: 0.93

We are using an ensemble of CLIP and SIGLIP to get

- Optimal Precision
- Separation of Classes (TP vs TN)
- Good Recall



Results from Duplicate Detection Algorithm





Different set of people are visible





The left side photo is the top half of the right side





The right side photo is not showing anything on the shelf while the left side photo shows the proper shop





The camera has been panned to take the right side photo



Success Story: Scaling Image Validation Across Semi Urban India

• 1 Million Images

Captured and processed across **25 states** in semi urban India (mostly Kirana stores).

• 94% Accuracy

Achieved in classifying valid vs. invalid images using a hybrid Al approach (MLLM + trained ResNet).

- Image Deduplication Accuracy
- •96% at 0.95 similarity threshold
- •85% at 0.90 similarity threshold

Advantage of the Deduplication Approach

No training of Images required. Only the dump of the images is needed

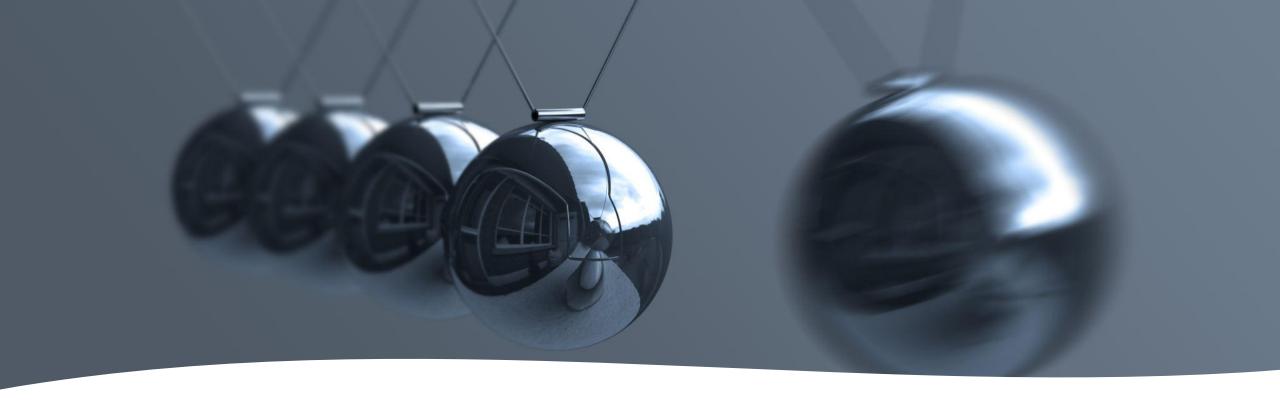
Fast execution

Works for any domain

Able to find out matches beyond the pixel comparison

Can be applied to different geographies of the world





Some Extensions of the two approaches



Outlet Image Search: Finding needle in the Hay

"I'm planning to launch a new product priced at \$20. Which stores are likely to be ready for it?"

Approach 1 (Intuition Based)

I have the sales data I know 5\$ pack sells over here let me place 20\$ pack also

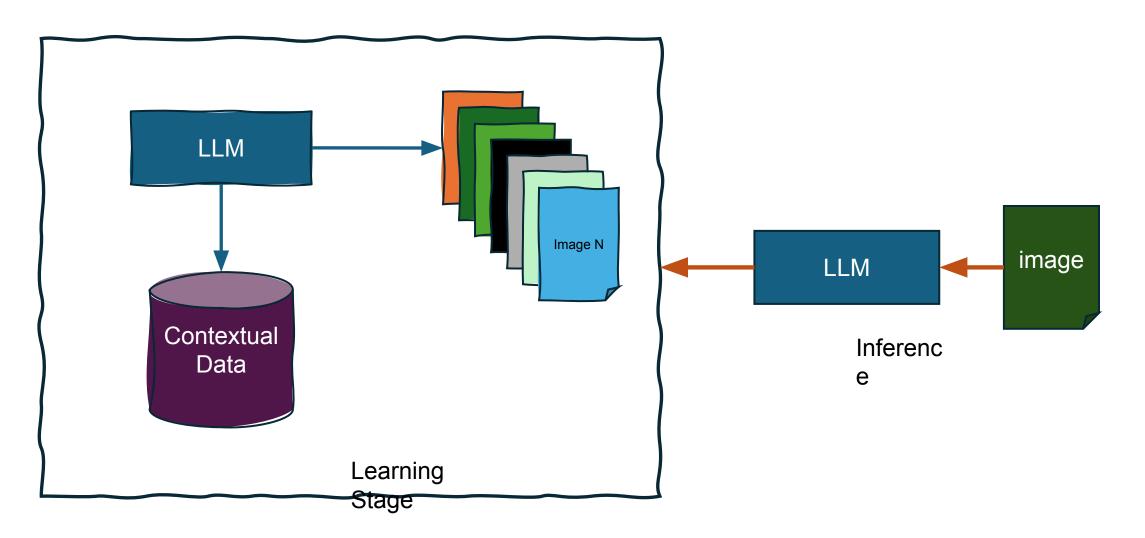
Approach 2 (Visual Inspection)







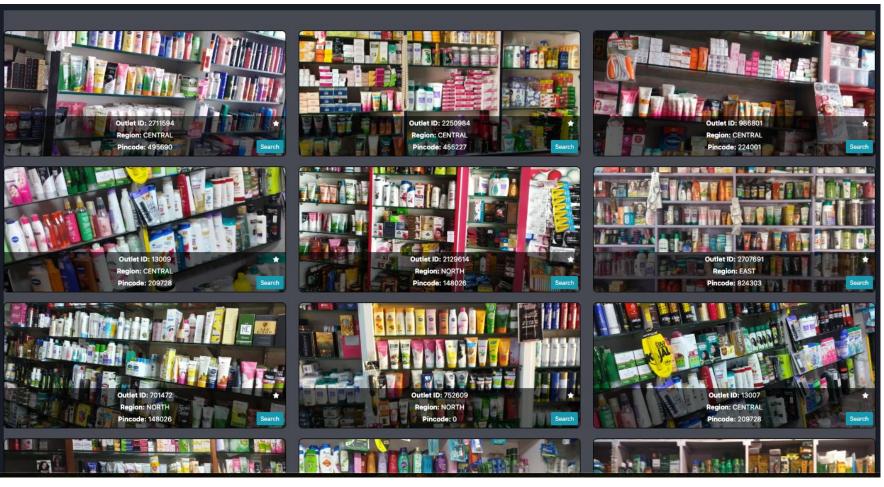
Can we help in finding the right outlets?





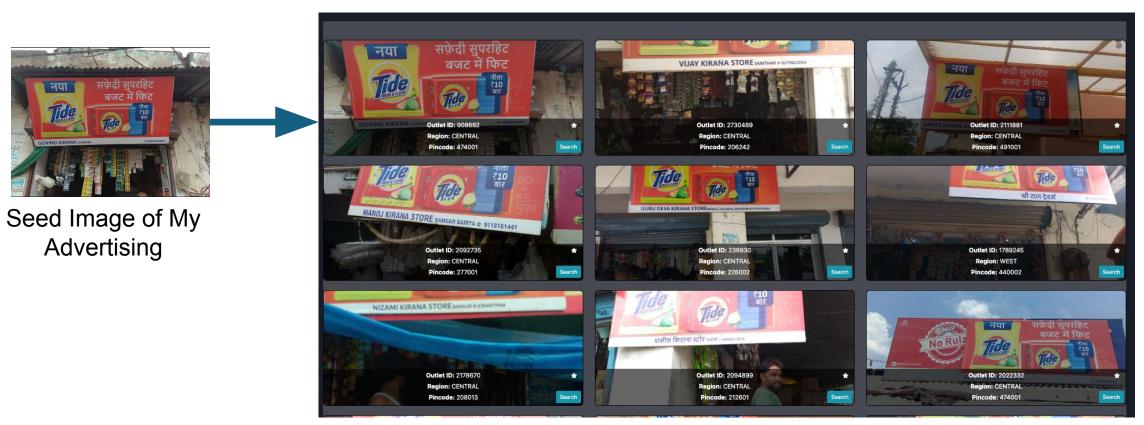
Usecase#1: In which outlets can I place my 20\$ shampoo







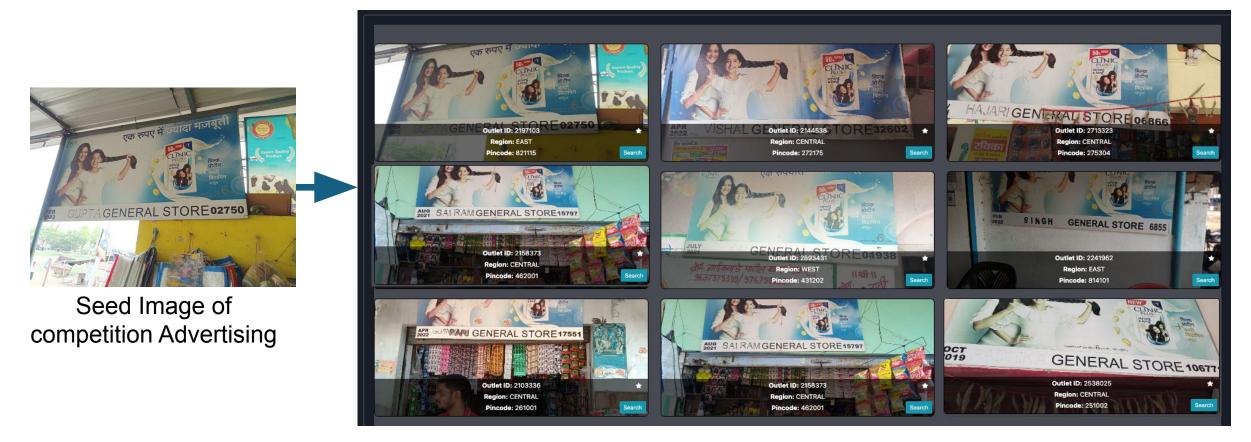
Usecase#2: Is my Trade spend on Advertising giving good visibility against competition



Outlets in my universe which show my advertising



Usecase#1: Is my Trade spend on Advertising giving good visibility against competition



Outlets in my universe which show my competition

Data will suggest if my advertising has more visibility than the competition. I might be adding more budget or reducing the budget for Trade promotion



