

Bauhaus-Universität Weimar

Fakultät Bauingenieurwesen

IMAGE DEDUPLICATION (IMDEDU)

Olubunmi Emmanuel Ogunleye

Selva Ganapathy Elangovan

Special Project Presentation

Name Of Advisor: Prof. Björn Rüffer

07-03-2025



Introduction/Motivation	03
Aim and Objectives	04
Literature Review	05 - 08
Methodology	09 - 10
Results and Discussion	11 - 12
Conclusion/Findings	13
Limitation of Study	14
Recommendation	15
References	16



Project Presentation Outline

Introduction / Motivation



DATABASE

How Do We Detect Image Duplicates?



Significance #1

Eliminate Multiple Versions of Same Image

Significance #2

Efficient Usage of Data Storage Space

Significance #3

Reduce Search Complexity in Information Retrieval

Aim and Objectives

Aim is to Implement an Algorithm that Detects Near Duplicates.

01

**Implement a way to
gather collection of
images from data
storage**

02

**Define metric to
measure similarities
between images**

03

**Implement a
technique to cluster
similar images**

04

**Determine the best
version of a group
of similar images**

Literature Review

Theory 01

Hash Functions

Mathematical algorithm that transforms data into fixed strings of bit called hash codes.

(Ref: Stallings, 2017)

Theory 02

Cryptographic Hashes

Unique and consistent output for any given input. Minor change result in significantly different hash value.

Collision Resistant. Examples: MD5, SHA-256 (Ref: Schneier, 2017)

Theory 03

Perceptual Hashes

Finding similarities between Images.

Collisions can occur.

Produces similar hash values for similar images.

(Sharma, 2014)

Theory 04

Similarity Distance Measures
Used to cluster similar images. Examples
Euclidean, Hamming, Manhattan etc.

Literature Review Cont'D

Types of Perceptual Hashes (Ref: Viies, 2015)



Average Hash

Averages the pixel values. Computationally Efficient



Median Hash

Uses Median Pixel Values. Robust to Certain Image Modification



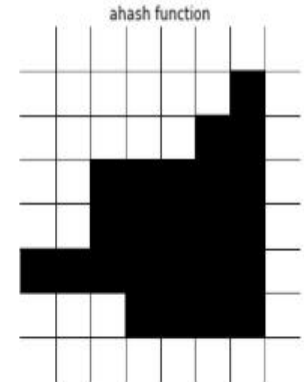
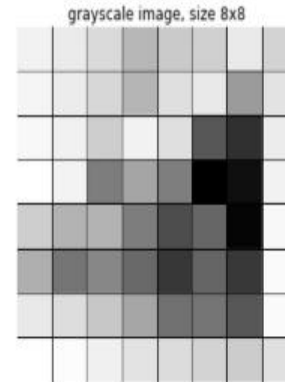
Perceptual Hash

Uses Discrete Cosine Transform. Robust to scaling and Rotations



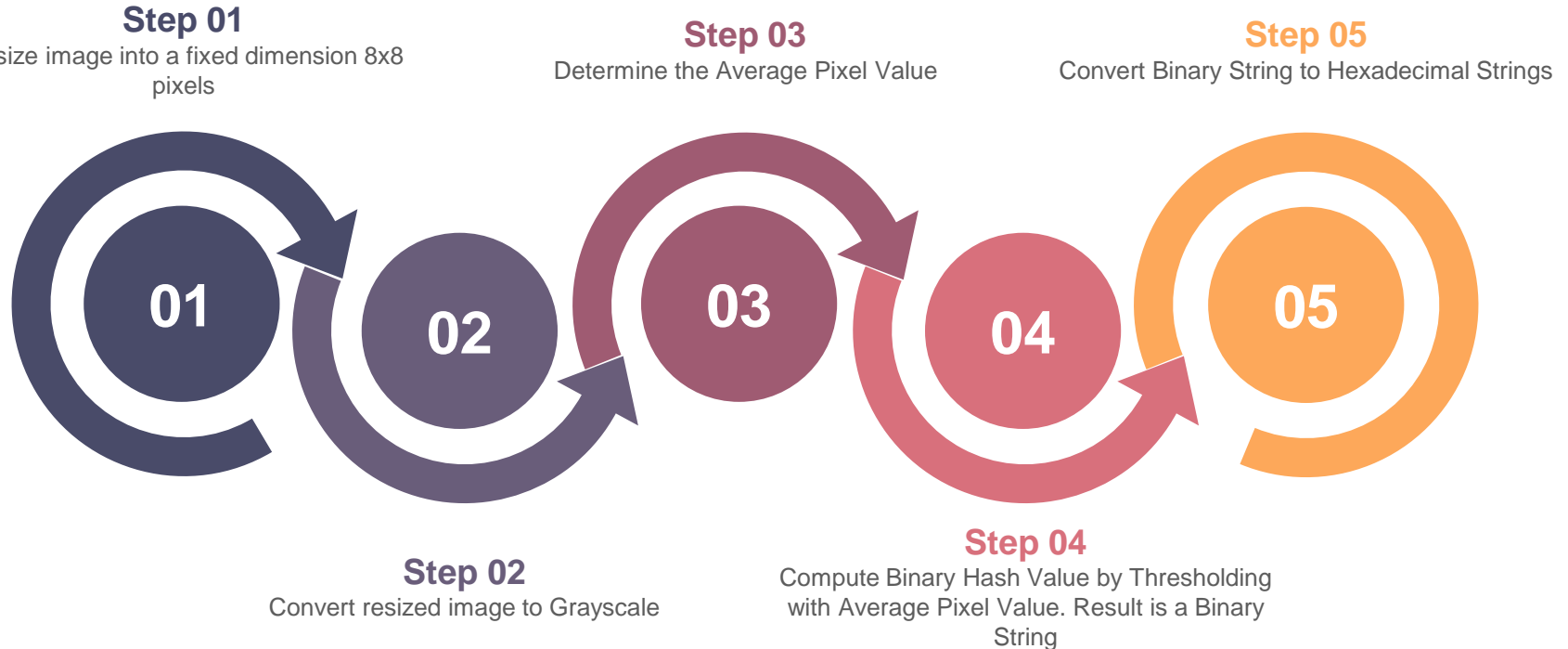
Wavelet Hash

Uses Discrete Wavelet Transformation. Also extract hash values from low frequencies.



Literature Review Cont'D

Average Hashing (Ref: Viies, 2015)



Literature Review Cont'D



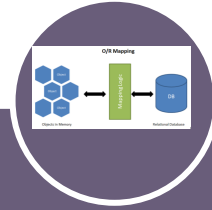
Hamming Distance Metric

Comparing two binary vectors bit by bit.

Differing bits are assigned value of 1

Similar bits are assigned value of 0

(Ref: Hamming, 1950)



SQL Alchemy

Object Relational Mapper (ORM)

Data Storage

Flexibility to Extend

(Ref: Bayer, 2025)



Streamlit

Open Source Python Frameworks for
Interactive Web Applications

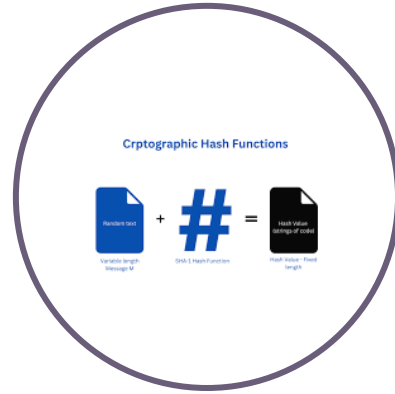
(Ref: Streamlit Community, 2025)

Methodology



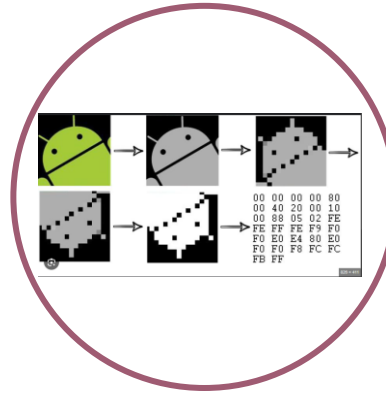
Load Images and Image Properties from Data Storage

(Local Drive or Hard Drive)



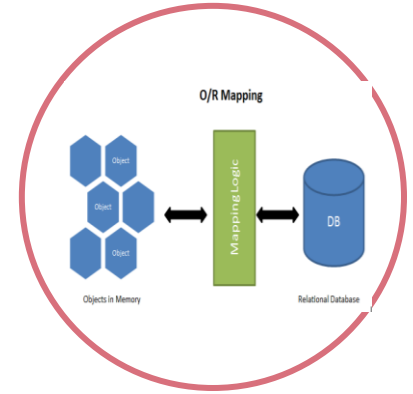
Compute Cryptographic Hash Values as Key to Store Image Metadata in SQL Database

SHA-256 was used



Compute Perceptual Hashes. Involves creating 3 instances along side original.

Average hash was used.

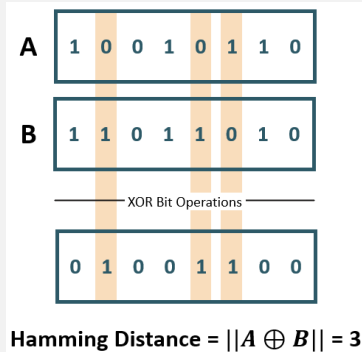


SQLAlchemy to store Image MetaData and Perceptual Hashes into DB. Cryptographic Hash used as Key

Methodology Cont'D

Distance Metrics

- Using Hamming Distance.



Similarity Threshold

- Implementing a slider to allow user defined similarity threshold.
- Limits of the threshold are based on max. and min. value of the minimum pairwise hamming distance.

Image Ranking

- From Similar Image Cluster, a ranking criteria is implemented to choose best fit.
- This ranking criteria is based on image with best resolution.

Visualization

- Streamlit is used to create web based UI.
- This helps users interact with the built algorithm.
- Takes in user input and returns output.

Result and Discussion



Cryptographic Hash Result

64-character long hexadecimal string. Proved effective for storing Images MetaData.



Perceptual Hash Result

16-character long hexadecimal string. Proved effective for detecting similarities between images.



Similarity Measure

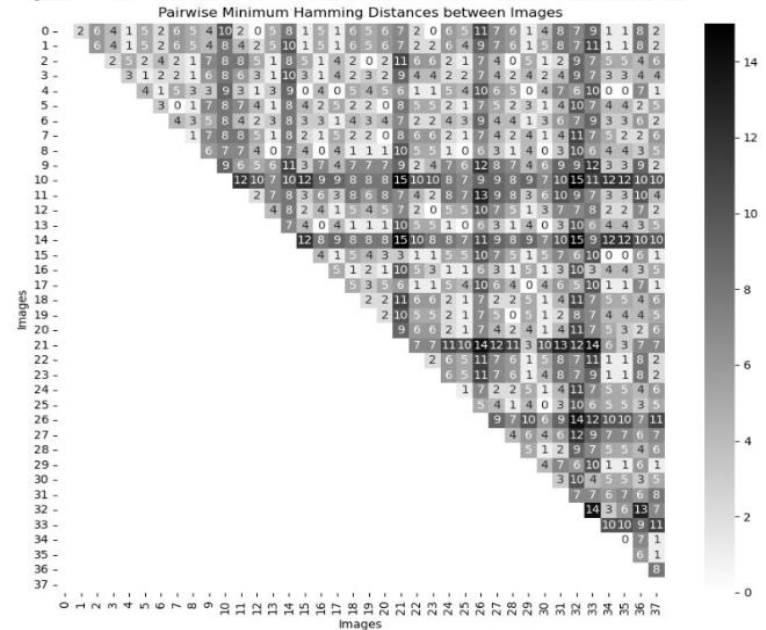
Hamming Distance was effective in measuring similarity between images. Similarity Matrix was plotted.



Streamlit for Visualization

Performed effectively to provide efficient user interaction with algorithm.

```
{
  '1f1a20e54273cf28d5342de29e7591b0829cede977040303bc8cd73aac871212d': {
    'Perceptual_Hash': '3f3b380d1e3c6400',
    'Perceptual_Hash_rotation_15degrees_Clockwise': '1e7efefefff7f7e70',
    'Perceptual_Hash_rotation_15degrees_CounterClockwise': '787f7fffffe7e0e',
    'Perceptual_Hash_Shear_Low_Difference': '3f3bba3e1efefee00'
  },
  '0046be088fdafdb8d921d937ce11973058b21d8f26a974d4009d6d26705f5b44': {
    'Perceptual_Hash': '0040f8fcfeffffff',
    'Perceptual_Hash_rotation_15degrees_Clockwise': '0e7efefefff7f7e78',
    'Perceptual_Hash_rotation_15degrees_CounterClockwise': '787e7fffffe7e1e',
    'Perceptual_Hash_Shear_Low_Difference': '00fefefefefefee00'
  }
}
```



Result and Discussion Cont'D

IMDEDU - Image Deduplication

Compare images using perceptual hashes and Hamming distances.

File Path

Enter folder path:

e.g., /path/to/images

Process Folder

Upload Image

Upload an image



Drag and drop file here

Limit: 200MB per file • PNG, JPG, JPEG

Browse files



Best Matching Image

Image Similarity Finder

Deploy

IMDEDU - Image Deduplication

Compare images using perceptual hashes and Hamming distances.

File Path

Enter folder path:

C:\APPL\SELVA\Hogal 2024

Process Folder

Loaded 35 images successfully

Database updated with image metadata

Download All Image Metadata as CSV

Upload Image

Upload an image



Drag and drop file here

Limit: 200MB per file • PNG, JPG, JPEG

Browse files



Best Matching Image

Image Similarity Finder

Deploy

Result and Discussion Cont'D

IMDEDU - Image Deduplication

Compare images using perceptual hashes and Hamming distances.

File Path

Enter folder path:

C:\APPU SEIWA\Jungai 2024

Process Folder

Upload Image

Upload an image

Drag and drop file here

13.JPG 12.6MB


Uploaded image

Hamming Distance Threshold

Adjust similarity threshold

0 10 20


Best Matching Image



Best Match: 13.JPG (Dist: 0)

File Name: 13.JPG
Hamming Distance: 0
Dimensions: 600x400
File Size: 12896.0 KB
File Creation Date: 2024-12-21 12:20:46.042108
File Extension: .jpg

Similar Images



Download Results as CSV

Conclusions / Findings

Cryptographic Hashes
efficient for storing and
retrieval of image
information

Perceptual Hash effective
for identifying
Similarities

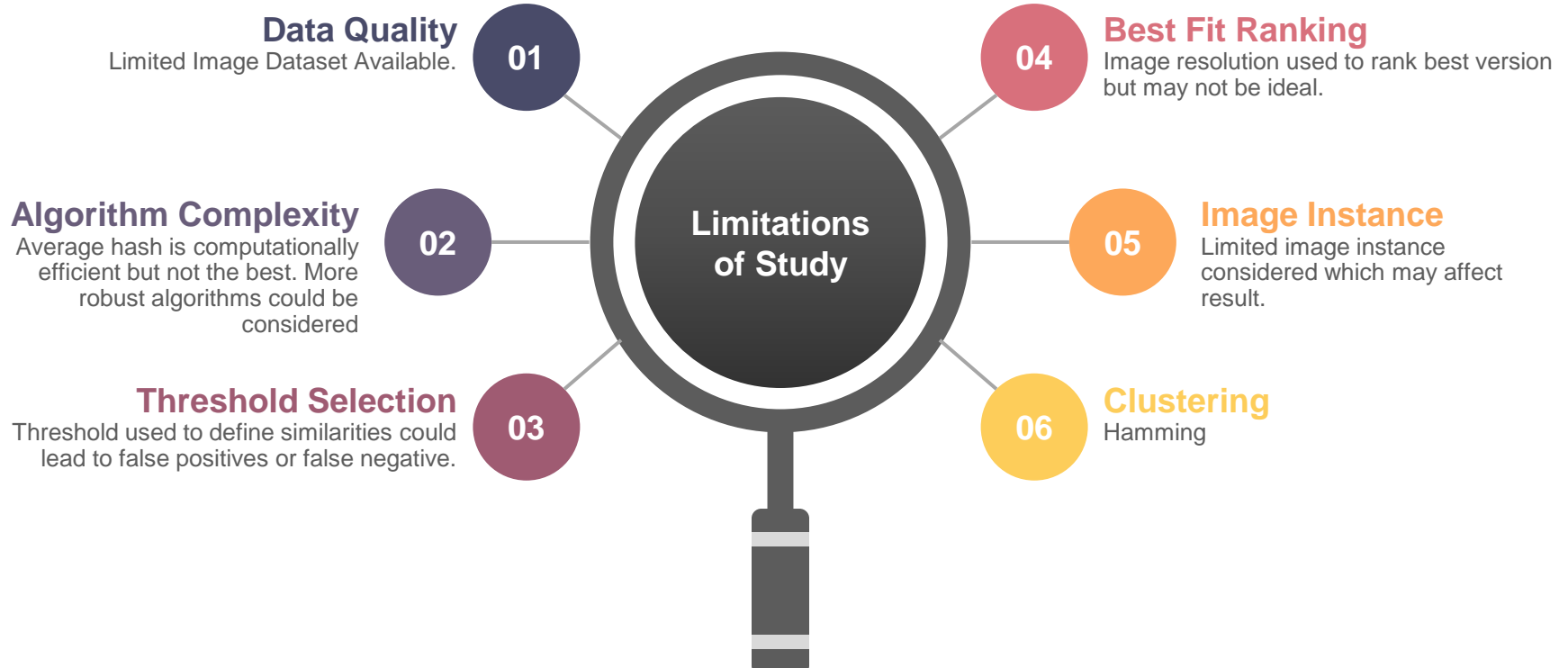
Hamming Distance
proved effective for
measuring similarities

SQL Alchemy effective
for ORM and Image
MetaData Storage

Streamlit provided user
friendly GUI

Proposed Hash
technique effective for
Image Deduplication

Limitations of Study



Recommendations



Recommendation 1

Research should be expanded to include more sophisticated algorithms such as pHash to improve accuracy.

Recommendation 2

Exploring modern methods like machine learning-based hashing techniques could offer greater flexibility and adaptability to identify similarities.

Recommendation 3

Future work could involve testing larger and more diverse datasets including images from sources such as social media, scientific data and websites.

References

1

Aradhana and D. S. M. Ghosh. Review Paper on Secure Hash Algorithm With Its Variants. *International Journal of Technical Innovation in Modern Engineering & Science*, 3(5):01–07, Nov. 2021.

2

M. Bayer. SQLAlchemy - The Database Toolkit for Python, 2025. URL <https://www.sqlalchemy.org/>. Version 2.0, Accessed: 2025-03-04.

3

B. Schneier. *Applied Cryptography: Protocols, Algorithms and Source Code in C*. Wiley, 2017.

4

R. W. Hamming. Error Detecting and Error Correcting Codes. *Bell System Technical Journal*, 29(2):147–160, 1950.

5

S. Sharma. Distance distributions and runtime analysis of perceptual hashing algorithms. *Journal of Visual Communication and Image Representation*, 104:104310, 2024. doi: 10.1016/j.jvcir.2024.104310.

6

W. Stallings. *Cryptography and Network Security: Principles and Practice*. Pearson, 7th Global Edition edition, 2017.

7

V. Viies. Possible Application of Perceptual Image Hashing. Tallinn University of Technology, Faculty of Information Technology Department of Computer Engineering, Master Thesis, 2015.

8

Streamlit Community. Streamlit: Turn Python Scripts into Beautiful WebApps, 2025. URL <https://docs.streamlit.io/get-started>. Version 1.42.0, Accessed: 2025-02-05.



QUESTIONS

**THANK YOU
FOR
LISTENING!**