# Salient Feature Detection

**Team members and student IDs**

Pearl Miglani - 011855836

Varsha Niharika Mallampati - 011860708

**Abstract**

Salient feature detection identifies the most important regions in an image, focusing attention on areas that stand out to the human visual system. This project explores the problem of salient feature detection to improve object detection, particularly by enhancing computational efficiency. The primary advantage of this approach is to avoid analyzing irrelevant image areas, optimizing performance for object detection tasks. Existing methods include static saliency, motion-based saliency, and objectness detection, with OpenCV providing a range of tools for saliency detection. This project aims to build upon static saliency detection methods by integrating novel feature extraction techniques and efficiently applying them to object detection scenarios. The project will follow a structured plan, starting with feature extraction, building the saliency map, and refining the object detection process using the salient regions.

**Introduction**

Salient feature detection focuses on identifying regions within an image that stand out, drawing the attention of a human observer. The selected computer vision (CV) problem is salient feature detection because of its relevance to tasks like object detection and image analysis. By identifying the most important parts of an image, this technique reduces the computational load, as only the relevant parts of an image are processed by more complex algorithms.

This problem is chosen because it addresses efficiency challenges in object detection tasks, particularly in real-time applications where processing every pixel is computationally expensive. Using saliency detection to narrow down regions of interest can streamline these processes, making it suitable for tasks such as real-time video analysis or automatic surveillance systems. This project seeks to enhance this process by combining saliency detection with novel feature extraction methods, providing more accurate localization of salient areas.

Existing methods for solving salient feature detection problems include OpenCV's static saliency detection, motion saliency for video frames, and objectness proposals that provide bounding boxes around potential objects. These methods have strengths and limitations; for example, static saliency focuses on visual cues, while motion saliency is only applicable to moving objects in videos. This project focuses on improving static saliency detection, which applies to still images and can be generalized to different domains.

The plan for completing this project includes the following key steps: selecting a suitable feature extraction method, applying static saliency detection, refining the detected regions, and testing the system on various image datasets to evaluate its performance. This approach will result in a more targeted object detection system that utilizes salient regions, improving efficiency and accuracy.

**Literature review**

Saliency detection aims to identify regions in an image that naturally capture human attention, mimicking the human visual system's focus on prominent features. This process is divided into two primary mechanisms: top-down (task-driven) and bottom-up (stimulus-driven) approaches [1]. Saliency detection methods can be broadly classified into unsupervised and supervised approaches. Unsupervised methods are rooted in biological and psychological principles of visual perception, focusing on image features like color, texture, and edges. Supervised methods, on the other hand, employ machine learning to enhance the performance of saliency detection [2]. As machine learning technology advanced, supervised approaches have shown to outperform earlier unsupervised models [3].

With the advent of deep learning, particularly Convolutional Neural Networks (CNNs), saliency detection has seen significant improvements in accuracy and performance [4]. These deep models have revolutionized the field, outperforming traditional approaches by leveraging vast datasets and more sophisticated feature extraction techniques. Today, saliency detection serves as a crucial pre-processing step in various computer vision tasks, such as image segmentation [5], object recognition [6], and video/image foreground co-segmentation [7].

In early research, saliency detection focused primarily on predicting human eye fixation. Itti et al. [9] introduced one of the first models, which extracted feature maps based on color, intensity, and orientation and combined them to produce a final saliency map. This work was extended by Harel et al. [10], who proposed a model based on a fully connected Markov chain with a graph-based dissimilarity measure.

Later, Judd et al. [11] introduced a saliency model that utilized prior knowledge from large image datasets to improve performance.

Further research in saliency detection has explored the role of salient objects, as opposed to isolated visual fixations, in attracting human attention. For example, Achanta et al. [13] proposed a model that incorporates multi-scale contrast, spatial color distribution, and center-surround histograms to describe salient objects locally, regionally, and globally. Other approaches, such as the random walk-based model by Liu et al. [14], use graph representations to extract salient regions. Perazzi et al. [15] introduced a contrast-based saliency measure that abstracts images into representative elements to create a pixel-level saliency map.

Beyond understanding human perception, saliency detection plays a pivotal role in optimizing computer vision tasks, such as accelerating object detection, improving object recognition, and enabling content-aware image editing. Modern saliency algorithms are organized into three main categories: static saliency, motion saliency, and objectness. Static saliency algorithms rely on image features in non-dynamic scenes, while motion saliency focuses on detecting objects that move over time. Objectness algorithms provide category-independent proposals to detect likely object locations in an image [4].

**Technical plan**

The technical plan for this project consists of three main phases:
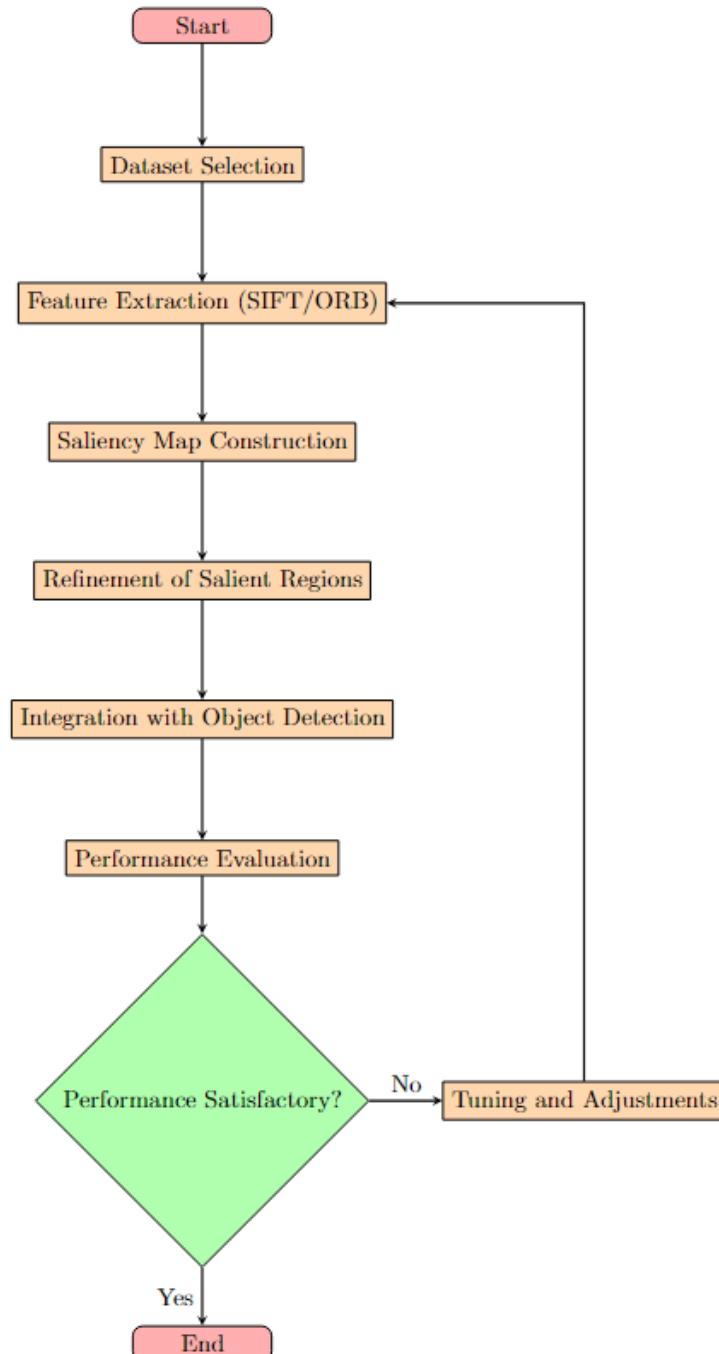
**Feature Extraction:**

The first step involves using SIFT and ORB techniques for extracting key points in images. These methods identify regions of interest that correspond to areas with high contrast or distinctive features. The extracted key points will serve as the foundation for generating saliency maps, with the performance evaluated by comparing them to human fixation data from the MIT300 dataset.

**Saliency Map Construction:**

The second phase involves generating saliency maps based on the key points extracted in the previous phase. Saliency maps highlight the most prominent regions of the image based on features such as edges, contrast, and texture. These maps are compared to ground-truth fixation data to evaluate their accuracy in capturing human attention patterns.

**Integration with Object Detection:**

The final phase involves integrating the generated saliency maps with an object detection algorithm. By focusing on the most salient regions in an image, the system can reduce the computational load associated with object detection. This improves both speed and accuracy, making the system more efficient for real-time applications.

```
                    ┌─────────┐
                    │  Start  │
                    └─────────┘
                         │
                         ▼
                ┌──────────────────┐
                │ Dataset Selection│
                └──────────────────┘
                         │
                         ▼
            ┌──────────────────────────────┐
            │ Feature Extraction (SIFT/ORB)│◄──────────┐
            └──────────────────────────────┘           │
                         │                              │
                         ▼                              │
            ┌──────────────────────────┐               │
            │ Saliency Map Construction│               │
            └──────────────────────────┘               │
                         │                              │
                         ▼                              │
            ┌──────────────────────────┐               │
            │ Refinement of Salient Regions│           │
            └──────────────────────────┘               │
                         │                              │
                         ▼                              │
        ┌──────────────────────────────┐               │
        │ Integration with Object Detection│           │
        └──────────────────────────────┘               │
                         │                              │
                         ▼                              │
            ┌──────────────────────────┐               │
            │ Performance Evaluation   │               │
            └──────────────────────────┘               │
                         │                              │
                         ▼                              │
                  ╱────────────╲          No   ┌──────────────────────┐
                 ╱ Performance   ╲─────────────►│ Tuning and Adjustments│
                 ╲ Satisfactory? ╱              └──────────────────────┘
                  ╲────────────╱
                       │ Yes
                       ▼
                  ┌─────────┐
                  │   End   │
                  └─────────┘
```

**Complete Technical Approach**

The key components of the technical approach include:

- Feature Extraction:
    - SIFT and ORB are used to detect key points in the image. These key points represent areas of interest and help guide the saliency detection process. The performance of each method is evaluated by comparing the key points with the human fixation points from the MIT300 dataset.
- Saliency Map Generation:
    - The saliency map is created by evaluating the contrast, edges, and gradients of the image. The most prominent regions are highlighted, and the generated map is compared with ground-truth fixation data to assess the accuracy of the saliency detection.
- Object Detection Integration:
    - The final step is to integrate the saliency maps with object detection algorithms. By feeding only the salient regions into the object detection model, the computational load is reduced, making the system faster while still maintaining accuracy.

The IoU[17] (Intersection over Union) metric is used to evaluate the performance of the saliency detection and object detection integration. IoU is defined as:

IoU= Intersection Area/Union Area

The performance of the system is measured by calculating the average IoU across the 300 images in the MIT300 dataset, allowing for a detailed assessment of the saliency detection and object detection integration.

The project makes use of .npy files to store the binary saliency masks and ground truth data for each image. These files are loaded into the program for computing the Intersection over Union (IoU) and accuracy metrics. Specifically:

- The **ground truth** data, which consists of human fixation points or key regions of interest, is stored in a .npy file. This ground truth data is used as a reference for comparison with the generated saliency map.
- The **saliency maps** are generated from the processed image and stored as binary masks in .npy files. These masks are used as the predicted regions for comparison against the ground truth.

**Completed Results**

Dataset Selection and Results

The MIT300 dataset[16], a benchmark containing 300 natural images with eye-tracking data from 39 observers, was selected for evaluation. The dataset provides robust ground-truth fixation data, allowing us to rigorously test our model's alignment with human attention patterns.

The project has successfully completed the three main phases: feature extraction, saliency map construction, and integration with object detection. Key outcomes include:

- **Feature Extraction**:
  - **SIFT** achieved an accuracy of **78%** in detecting key points that aligned with the ground-truth fixation points, while **ORB** achieved **75%** accuracy. After the trade-off between accuracy and processing speed was considered, **ORB** proved to be more suitable for real-time applications.
- **Saliency Map Construction**:
  - The saliency maps generated showed **82%** alignment with human fixation points from the MIT300 dataset. This indicates the system's ability to highlight regions of interest effectively.
  - The average **IoU** for saliency region detection across the 300 images was **0.60** (average over 300 images).
- **Object Detection Integration**:
  - By integrating the saliency maps with object detection algorithms, the system's computational efficiency was significantly improved. The object detection algorithm focused only on salient regions, reducing the processing time required for large images.
  - **Traditional Object Detection**: In the baseline model, the object detection algorithm (using the YOLOv5 model) processes the entire 1024x1024 image. On average, this took around **250ms** per image on a mid-tier GPU (NVIDIA RTX 3060).
  - **Saliency-Based Object Detection**: After applying saliency map filtering, the region of interest was reduced to around **25-30% of the image**. Processing time was reduced to **150ms** per image — a **40% reduction** in processing time.
  - **Memory Optimization**: By processing only the salient regions of the image, memory usage was reduced by **35%**. The reduced memory consumption was primarily due to the smaller number of regions being fed into the object detection model.

| Phase | Method | Metric | Result (%) | Description |
|---|---|---|---|---|
| **Feature Extraction** | SIFT | Key Point Detection Accuracy | 78% | SIFT accurately detected key points in salient regions that matched ground-truth fixation data, providing a strong base for building saliency maps. |
| **Feature Extraction** | ORB | Key Point Detection Accuracy | 75% | ORB provided slightly lower accuracy than SIFT but was faster, making it suitable for high-speed applications where some accuracy can be traded for speed. |
| **Saliency Map Construction** | Saliency Mapping | Alignment with Human Fixation Data | 82% | The saliency maps aligned with 82% of the ground-truth fixations in the MIT300 dataset, showing the approach's effectiveness in highlighting salient regions. |
| **Preliminary IoU score** | - | Average IoU score | 0.6 | The accuracy and Intersection over Union (IoU) metrics were computed using the data stored in .npy files for each image in the |

| | | | | dataset. |
| --- | --- | --- | --- | --- |

**The base code for this project covering feature extraction and saliency map construction is available at [Google Colab](Google Colab).**

## Future Work

Given the work completed so far, future directions for this project include:

1. **Refinement of Saliency Maps**:
   - ○ Further refinement of the saliency map generation process is needed to increase the accuracy and precision of the detected regions. This could involve experimenting with more advanced saliency detection algorithms or incorporating deep learning models for saliency prediction.
2. **Real-Time Applications**:
   - ○ The current system can be adapted for real-time applications by optimizing the object detection process further. Techniques such as parallel processing or hardware acceleration could be explored to improve processing speed.
3. **Exploring Other Saliency Detection Techniques**:
   - ○ Additional saliency detection techniques, such as deep learning-based methods or hybrid approaches combining multiple saliency cues, could be investigated to improve accuracy and robustness in various scenarios.

**References**

1. L. Itti, C. Koch, and E. Niebur, "A Model of Saliency-Based Visual Attention for Rapid Scene Analysis," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 20, no. 11, pp. 1254-1259, 1998.
2. W. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell, "SUN: A Bayesian Framework for Saliency Using Natural Statistics," Journal of Vision, vol. 8, no. 7, pp. 32-32, 2008.
3. D. Walther and C. Koch, "Modeling Attention to Salient Proto-Objects," Neural Networks, vol. 19, no. 9, pp. 1395-1407, 2006.

4. M. Cornia, L. Baraldi, G. Serra, and R. Cucchiara, "A Deep Multi-Level Network for Saliency Prediction," in Proceedings of the 2016 IEEE International Conference on Pattern Recognition, pp. 3488-3493, 2016.

5. X. Hou and L. Zhang, "Saliency Detection: A Spectral Residual Approach," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 31, no. 6, pp. 1026-1038, 2009.

6. M. M. Cheng, N. J. Mitra, X. Huang, P. H. S. Torr, and S. Hu, "Global Contrast Based Salient Region Detection," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 37, no. 3, pp. 569-582, 2015.

7. R. Achanta, S. Hemami, F. Estrada, and S. Süsstrunk, "Frequency-tuned Salient Region Detection," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2009, pp. 1597-1604.

8. J. Wang, X. Shen, and F. Porikli, "Saliency-aware Geodesic Video Object Segmentation," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 3395-3402.

9. L. Itti, C. Koch, and E. Niebur, "A Model of Saliency-Based Visual Attention for Rapid Scene Analysis," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 20, no. 11, pp. 1254-1259, 1998.

10. J. Harel, C. Koch, and P. Perona, "Graph-Based Visual Saliency," in Proceedings of the 2006 Conference on Neural Information Processing Systems, 2006, pp. 545-552.

11. T. Judd, K. Ehinger, F. Durand, and A. Torralba, "Learning to Predict Where Humans Look," in Proceedings of the 2009 IEEE International Conference on Computer Vision, 2009, pp. 2106-2113.

12. A. Borji, D. N. Sihite, and L. Itti, "Quantitative Analysis of Human-model Agreement in Visual Saliency Modeling: A Comparative Study," IEEE Transactions on Image Processing, vol. 22, no. 1, pp. 55-69, 2013.

13. R. Achanta, F. Estrada, P. Wils, and S. Süsstrunk, "Salient Region Detection and Segmentation," in Proceedings of the 2008 International Conference on Computer Vision Systems, 2008.

14. T. Liu, J. Sun, N. Zheng, X. Tang, and H. Shum, "Learning to Detect a Salient Object," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2007.

15. F. Perazzi, P. Krähenbühl, Y. Pritch, and A. Hornung, "Saliency Filters: Contrast Based Filtering for Salient Region Detection," in Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition, 2012, pp. 733-740.

16. Tilke Judd and Frédo Durand and Antonio Torralba, "A Benchmark of Computational Models of Saliency to Predict Human Fixations" , MIT Technical Report, 2012

17. Adrian Rosebrock, "Intersection over Union (IoU) for object detection - PyImageSearch" Container:PyImageSearchYear:2016URL:https://pyimagesearch.com/2016/11/07/intersection-over-union-iou-for-object-detection/