

Title

Computer Vision: Foundations, Algorithms, Applications, and Practical Guide (Thesis-Level)

Author: Generated by Assistant

Purpose: This document provides comprehensive, thesis-level content on computer vision. It is structured to a

Abstract

Computer vision is an interdisciplinary field that develops computational methods to infer information from visual signals—images and videos. This thesis-style document covers physical image formation, feature extraction, classical and deep learning approaches, object detection and segmentation, face recognition fundamentals, metrics and datasets, practical system design, deployment, evaluation, and ethical considerations. It also includes a hands-on experiment using Teachable Machine and a reproducible research checklist.

1. Introduction

1.1 Definition and Goals

Computer vision (CV) studies how to make machines 'see'—i.e., to extract, interpret, and reason about visual

1.2 Roadmap of Topics

This document covers: (a) image formation and representation, (b) classical features and matching, (c) machin

2. Image Formation and Modeling

2.1 The Pinhole Camera Model and Projection

The pinhole camera projects 3D world points $X = [X \ Y \ Z \ 1]^T$ to image points $x = [u \ v \ 1]^T$ via:

$$x \sim K [R \ | \ t] X$$

where K (3×3) is the intrinsic matrix containing focal lengths (f_x, f_y) and principal point (c_x, c_y); R (3×3) is

2.2 Radiometric Image Formation

Image intensity $I(u,v) \approx \int_{\lambda} L(\lambda) S(\lambda) R(\lambda) d\lambda$, where L is illumination spectrum, S is surface spectral reflectance

2.3 Sampling, Aliasing, and Anti-Aliasing

Discrete pixel sampling can cause aliasing. Anti-alias filters (low-pass) are usually applied in the optical or soft

3. Image Representation and Preprocessing

3.1 Color Spaces

Common representations: RGB (device-dependent), grayscale (luminance), HSV (hue, saturation, value), Lab

3.2 Preprocessing Steps

Normalization (zero mean, unit variance), resizing, and histogram equalization. Use CLAHE (Contrast Limited

3.3 Data Augmentation

Essential for generalization: rotations, flips, crops, color jitter, scaling, cutout, mixup, and synthetic domain randomization. Note: augmentation should preserve label semantics.

4. Classical Feature Detection and Description

4.1 Edge and Corner Detectors

Edges are first-order derivative discontinuities; detect via gradient operators (Sobel) and refined using Canny

4.2 Local Invariant Descriptors

SIFT: scale-space extrema detection via DoG (Difference of Gaussians); orientation histograms produce rotated

4.3 Feature Matching and Geometric Verification

Match descriptors using nearest neighbors; use ratio test (Lowe) to filter ambiguous matches. Estimate geometric

5. Machine Learning Foundations for Vision

5.1 Supervised Learning Formulation

Given dataset $D = \{(x_i, y_i)\}_{\{i=1\}^N}$, find parameters θ minimizing $L(\theta) = (1/N) \sum_i (f(x_i; \theta), y_i) + \lambda R(\theta)$, where

5.2 Convolutional Neural Networks (CNNs)

Key components: convolutional layers (shared weights), pooling, non-linearities, skip connections (ResNet). CNNs are highly parallelizable.

5.3 Practical Training Details

Optimizers: SGD with momentum, Adam. Learning rate scheduling, warm restarts, and weight decay are crucial. Use batch normalization to accelerate convergence.

5.4 Transfer Learning

Fine-tuning pretrained networks (ImageNet) reduces data requirements and speeds up convergence. Freeze layers and learn new ones.

6. Image Classification Deep Dive

6.1 Problem Definition

Assign label(s) to an image. In multi-label settings, use sigmoid activation per class with binary cross-entropy.

6.2 Architectures

From AlexNet (2012) to ResNet and EfficientNet. Residual blocks allow training very deep networks by learning

6.3 Regularization and Generalization

Dropout, weight decay, data augmentation, early stopping. Evaluate using confusion matrices to identify class imbalance issues.

6.4 Example Experiment Design

Dataset split: train/val/test = 70/15/15. Baseline: fine-tune ResNet-50 with initial LR=1e-3 (Adam), batch size=

7. Object Detection: Algorithms and Trade-offs

7.1 Problem Statement

Detect and localize all objects with bounding boxes and class labels. Standard outputs include bounding box coordinates and class confidences.

7.2 Two-Stage Detectors

Faster R-CNN uses an RPN to propose regions; second stage classifies and refines boxes.

Pros: accuracy. Cons: higher latency.

7.3 One-Stage Detectors

YOLO family and SSD predict boxes and classes in a single pass, achieving real-time performance. Modern YOLO v5 is the state-of-the-art.

7.4 Anchor Boxes vs Anchor-Free

Anchor boxes require careful design of scales/aspect ratios to match dataset. Anchor-free methods (centernet) do not require this.

7.5 Losses and Matching

Localization losses: smooth L1, IoU-based losses (GIoU, DIoU) improve bounding box regression. Matching scores: softmax cross-entropy.

8. Segmentation: Semantic and Instance

8.1 Semantic Segmentation

Per-pixel classification. Fully Convolutional Networks (FCNs) replace dense layers with upsampling to preserve spatial information.

8.2 Instance Segmentation

Mask R-CNN extends Faster R-CNN by adding a mask head per ROI. Instance segmentation requires both localization and precise mask prediction.

8.3 Metrics and Challenges

mIoU (mean Intersection over Union) is commonly used. Challenges include class imbalance, small objects, and occlusion.

9. Face Recognition: Pipeline and Methods

9.1 Pipeline Steps

Face detection → alignment (landmarks) → embedding extraction → matching. Modern pipelines use CNNs for all steps.

9.2 Embedding Learning

Triplet loss: $L = \sum [\max(0, \|f(a)-f(p)\|^2 - \|f(a)-f(n)\|^2 + \alpha)]$ encourages anchor-positive distance smaller than anchor-negative distance.

9.3 Practical Considerations

Dataset diversity, data augmentation, and demographic balance are key to fair performance across population groups.

10. Datasets, Annotation, and Benchmarking

10.1 Major Datasets

ImageNet (classification), COCO (detection & segmentation), Pascal VOC (detection), Cityscapes (urban segm)

10.2 Annotation Types and Tools

Bounding boxes, masks, keypoints, attributes. Annotation platforms: CVAT, Labelbox, Supervisely. Quality control: consensus labeling and spot-audits.

10.3 Evaluation Protocols

Use held-out test sets and cross-validation when data is limited. For detection, report mAP at IoU thresholds (0.5, 0.501-0.95)

11. Metrics, Analysis, and Error Diagnosis

11.1 Common Metrics

Classification: accuracy, precision, recall, F1. Detection: mAP, AP@IoU. Segmentation: mIoU, pixel accuracy.

11.2 Error Analysis

Confusion matrices, per-class performance, qualitative inspection of false positives/negatives, and calibration

11.3 Ablation Studies

Systematically vary single components/hyperparameters to attribute performance gains and ensure reproducibility

12. Deployment and System Engineering

12.1 Model Compression and Acceleration

Pruning: remove redundant weights. Quantization: reduce precision (e.g., FP32 → INT8). Knowledge distillation:

12.2 Serving and Latency

Use TensorRT, ONNX Runtime, or TFLite for optimized inference. Consider batching and asynchronous pipelines; measure end-to-end latency including preprocessing and postprocessing.

12.3 Edge vs Cloud

Edge inference reduces latency and preserves privacy; cloud enables larger models and easier updates. Hybrid strategies often balance trade-offs.

13. Advantages and Practical Uses

13.1 Advantages

- Automation: reduces human workload in visual inspection.
- Scalability: processes massive image/video streams.
- Objective measurements: pixel-level accuracy for tasks like medical imaging.

13.2 Use Cases

- Autonomous vehicles: perception stacks for detection, segmentation, lane detection.
- Security & Surveillance: anomaly detection and person tracking.
- Healthcare: radiology image analysis, histopathology.
- Retail: shelf monitoring, demand forecasting.
- Agriculture: crop monitoring, disease detection.

14. Limitations and Ethical Concerns

14.1 Technical Limitations

- Data dependency and domain shift.
- Failure modes under occlusion, extreme lighting, or adversarial conditions.

14.2 Ethical and Societal Concerns

- Privacy: face recognition and surveillance implications.
 - Bias & Fairness: datasets reflecting societal biases lead to unequal outcomes.
- Mitigation strategies: balanced datasets, fairness-aware training, audits, and human-in-the-loop validation.

15. Advanced Topics and Research Directions

15.1 Self-Supervised Learning

Contrastive learning (SimCLR, MoCo) and masked image modeling reduce labeled data needs by learning rich representations.

15.2 Multimodal Models

CLIP and similar approaches learn joint image-text embeddings enabling zero-shot recognition.

15.3 3D and Neural Rendering

NeRFs (Neural Radiance Fields) model scenes for novel view synthesis, opening avenues in 3D reconstruction.

16. Reproducible Research and Experimentation Checklist

16.1 Data and Code Management

Version data and code. Provide train/val/test splits, random seeds, and environment specifications.

16.2 Reporting

Report hyperparameters, compute budget, and dataset statistics. Include failure cases and qualitative results.

16.3 Example Hyperparameter Table

Model: ResNet-50 fine-tune. Optimizer: Adam. LR: 1e-4. Batch size: 64. Weight decay: 1e-4.

Epochs: 30. Augmentation: flip, crop, color jitter.

17. Practical Activity (Detailed): Image Classifier with Teachable Machine

17.1 Objective

Build and evaluate a simple image classification model without code using Teachable Machine. Demonstrates end-to-end pipeline: data collection, training, evaluation, and export.

17.2 Step-by-step Guide

1. Open teachablemachine.withgoogle.com and select 'Image Project'.
2. Create classes. Aim for at least 50 diverse examples per class for meaningful results.
3. Provide data via webcam or upload; ensure variation in lighting and background.
4. Train: select transfer learning settings or default. Observe training/validation accuracy.
5. Evaluate: use test images and view confusion. Inspect failure modes and retrain with more data or augmentation.
6. Export: TensorFlow.js (web), TensorFlow SavedModel (research), or TFLite (mobile). Deploy to a static site.

17.3 Technical Lessons

Teachable Machine uses transfer learning on a mobile-friendly backbone. For production problems, Teachable Machine is a prototyping tool but lacks rigorous dataset versioning and advanced augmentation controls.

18. Sample Project: From Data to Deployment (Worked Example)

18.1 Problem Statement

Detect and classify retail products on shelf images for stock monitoring.

18.2 Data Strategy

Capture shelf images across stores, annotate bounding boxes and product IDs. Use a balanced sampling strategy across store lighting conditions.

18.3 Model and Pipeline

Use Faster R-CNN or YOLOX as baseline. Train with mosaic augmentation for small objects; use class-balanced sampling. After training, quantize model to INT8 and deploy on edge devices with a small server for updates.

18.4 Evaluation

Report per-class AP, overall mAP, precision-recall curves, and runtime throughput (images/sec). Conduct A/B

19. Ethics & Governance Checklist

19.1 Data Collection Consent

Record consent where required. Anonymize personally identifiable information where possible.

19.2 Bias Audits

Run demographic performance breakdowns. If disparities exist, collect more data or consider algorithmic fairness techniques.

19.3 Human Oversight

Design systems with human-in-the-loop for critical decisions; log decisions for audits.

20. Conclusion

Computer vision has matured from geometric approaches to data-driven deep learning systems. It powers numerous applications but requires careful data engineering, evaluation, and ethical oversight. This document provides a foundation for both teaching and starting applied projects. Continuous validation and responsible deployment are essential.

Appendix A: Mathematical Details

A.1 Camera Intrinsics

$K = [[f_x, s, c_x], [0, f_y, c_y], [0, 0, 1]]$ where s is skew (often 0). Use homogeneous coordinates for projective

A.2 Epipolar Geometry

For calibrated cameras, $x^T E x = 0$ where $E = [t]_x R$ is the essential matrix; recover R and t up to scale via SVD

A.3 Convolution as Linear Operation

Convolution with kernel K over input I : $(I * K)(i,j) = \sum_m \sum_n I(i+m, j+n) K(m,n)$. In practice, implement via fast convolution

Appendix B: Recommended Reading and Papers

- Hartley & Zisserman, Multiple View Geometry (camera geometry).
- Goodfellow, Bengio, Courville, Deep Learning (foundations).
- Papers: AlexNet (2012), Faster R-CNN (2015), ResNet (2016), YOLO series, Mask R-CNN (2017), SimCLR

References

Representative citations: Krizhevsky et al., 2012; He et al., 2016; Ren et al., 2015; Redmon et al., YOLO; He et al., Mask R-CNN. For formal citations, consult original papers and canonical textbooks.

Extended Discussion

Troubleshooting checklist: inspect dataset balance, check for label leakage, visualize predictions, evaluate on out-of-sample data, calibrate confidence thresholds, and perform sanity checks like training on a small subset to ensure learnability. Hyperparameter tuning: use grid or Bayesian search; log experiments with tools (Weights & Biases, TensorBoard). For object detection, start with pretrained COCO weights, set initial learning rate 1e-3 for Adam or 0.01 for SGD, and run 12–24 epochs with step LR decay.

Extended Discussion

Troubleshooting checklist: inspect dataset balance, check for label leakage, visualize predictions, evaluate on out-of-sample data, calibrate confidence thresholds, and perform sanity checks like training on a small subset to ensure learnability. Hyperparameter tuning: use grid or Bayesian search; log experiments with tools (Weights & Biases, TensorBoard). For object detection, start with pretrained COCO weights, set initial learning rate 1e-3 for Adam or 0.01 for SGD, and run 12–24 epochs with step LR decay.

Extended Discussion

Troubleshooting checklist: inspect dataset balance, check for label leakage, visualize predictions, evaluate on out-of-sample data, calibrate confidence thresholds, and perform sanity checks like training on a small subset to ensure learnability. Hyperparameter tuning: use grid or Bayesian search; log experiments with tools (Weights & Biases, TensorBoard). For object detection, start with pretrained COCO weights, set initial learning rate 1e-3 for Adam or 0.01 for SGD, and run 12–24 epochs with step LR decay.

Extended Discussion

Troubleshooting checklist: inspect dataset balance, check for label leakage, visualize predictions, evaluate on out-of-sample data, calibrate confidence thresholds, and perform sanity checks like training on a small subset to ensure learnability. Hyperparameter tuning: use grid or Bayesian search; log experiments with tools (Weights & Biases, TensorBoard). For object detection, start with pretrained COCO weights, set initial learning rate 1e-3 for Adam or 0.01 for SGD, and run 12–24 epochs with step LR decay.

Extended Discussion

Troubleshooting checklist: inspect dataset balance, check for label leakage, visualize predictions, evaluate on out-of-sample data, calibrate confidence thresholds, and perform sanity checks like training on a small subset to ensure learnability. Hyperparameter tuning: use grid or Bayesian search; log experiments with tools (Weights & Biases, TensorBoard). For object detection, start with pretrained COCO weights, set initial learning rate 1e-3 for Adam or 0.01 for SGD, and run 12–24 epochs with step LR decay.

Extended Discussion

Troubleshooting checklist: inspect dataset balance, check for label leakage, visualize predictions, evaluate on out-of-sample data, calibrate confidence thresholds, and perform sanity checks like training on a small subset to ensure learnability. Hyperparameter tuning: use grid or Bayesian search; log experiments with tools (Weights & Biases, TensorBoard). For object detection, start with pretrained COCO weights, set initial learning rate 1e-3 for Adam or 0.01 for SGD, and run 12–24 epochs with step LR decay.