



BACHELOR'S THESIS IN COMPUTER SCIENCE AND INDUSTRIAL ECONOMICS

UNDERGRADUATE LEVEL, 15 CREDITS

Digital Asset Management in Project-Based Manufacturing: A Comparative Study of YOLOv12 and RF-DETR

Fine-Tuning Object Detection Models on
Company-Specific Imagery to Enhance
Knowledge Sharing

ELLA KARLSSON

School of Industrial Engineering and Management
Royal Institute of Technology (KTH)

Abstract

In project-based manufacturing environments, the effective organization and retrieval of visual design assets is essential for knowledge retention, operational efficiency, and design reuse. This bachelor's thesis investigates the application of deep learning-based object detection to automate metadata tagging for image assets in a bespoke furniture manufacturing SME that currently relies on Dropbox without structured tagging. The company operates with a hybrid project-based and process-oriented workflow, where visual documentation—ranging from marketing imagery to manufacturing snapshots—serves as a critical knowledge resource.

The study compares two state-of-the-art object detection models: YOLOv12, a convolutional neural network known for its real-time performance, and RF-DETR, a lightweight transformer-based architecture developed by Roboflow that balances high accuracy with edge-device efficiency. Both models were fine-tuned on the same company-specific dataset and evaluated in terms of precision, inference speed, and suitability for metadata generation in a resource-constrained environment.

The results demonstrate that YOLOv12 improves tagging precision by 15% and reduces image retrieval time by 20% compared to manual organization methods. RF-DETR, on the other hand, achieves strong detection performance on complex visual content while maintaining real-time inference speeds, offering a compelling balance between speed and robustness. The findings highlight the complementary strengths of transformer and CNN-based models in deploying AI for Digital Asset Management (DAM) in small-scale, design-driven firms.

By addressing the lack of automated DAM solutions tailored to SMEs in niche manufacturing sectors, this work contributes a scalable and adaptable approach to visual knowledge management. The proposed system enhances asset discoverability, promotes cross-project reuse, and supports process innovation—offering a practical step toward AI-enabled digital transformation in project-oriented organizations.

Keywords:

Digital Asset Management (DAM), Object Detection, YOLOv12, RF-DETR, Transformer Models, Metadata Tagging, Knowledge Sharing, Real-Time Inference, Project-Based Manufacturing, SMEs, Deep Learning, Edge Deployment

Sammanfattning

I projektbaserade tillverkningsmiljöer är effektiv organisering och återanvändning av visuella designresurser avgörande för kunskapsbevarande, operativ effektivitet och designåteranvändning. Denna kandidatuppsats undersöker tillämpningen av djupinlärningsbaserad objektigenkänning för att automatisera metadata-tagging av bildresurser i ett svenskt möbelföretag med en hybrid arbetsstruktur som kombinerar projektbaserad och processinriktad produktion. Företaget använder för närvarande Dropbox utan ett strukturerat system för bildmärkning, vilket försvårar skalbar hantering av visuellt material.

Studien jämför två moderna objektigenkänningsmodeller: **YOLOv12**, ett konvolutionsbaserat nätverk känt för sin realtidsoptimerade prestanda, och **RF-DETR**, en transformerbaserad modell utvecklad av Roboflow som kombinerar hög noggrannhet med effektiv inferens. Båda modellerna finjusterades med ett företagsspecifikt dataset och utvärderades utifrån precision, inferenstid och deras lämplighet för metadata-generering i miljöer med begränsade resurser.

Resultaten visar att YOLOv12 förbättrar tagging-precisionen med 15% och minskar bildsökningstiden med 20% jämfört med manuell organisering. RF-DETR uppvisar stark prestanda på komplexa visuella objekt samtidigt som den behåller realtidskapacitet, vilket gör den lämplig för implementering på resursbegränsad hårdvara. Studien belyser därmed hur transformer- och CNN-baserade modeller erbjuder kompletterande styrkor för AI-baserad Digital Asset Management (DAM) i småskaliga, designdrivna företag.

Genom att adressera bristen på automatiserade DAM-lösningar anpassade för små och medelstora företag (SMF) inom nischad tillverkning, bidrar detta arbete med en skalbar och anpassningsbar strategi för visuell kunskapshantering. Det föreslagna systemet stärker åtkomsten till resurser, främjar återanvändning mellan projekt och stödjer processinnovation—ett konkret steg mot AI-driven digital transformation i projektorienterade organisationer.

Nyckelord: Digital Asset Management, Objektigenkänning, YOLOv12, RF-DETR, Metadata, Kunskapsdelning, Realtidsinferens, Transformer-modeller, Projektbaserad tillverkning, Små och medelstora företag (SMF)

Acknowledgments

I would like to thank xxxx for having yyyy.

Contents

Abstract	1
Sammanfattning	2
Acknowledgments	3
List of Figures	6
List of Tables	7
List of Acronyms and Abbreviations	8
1 Introduction	9
1.1 Background	9
1.2 Problem	9
1.3 Purpose	9
1.3.1 Technical Research questions	9
1.3.2 Business Research questions	10
1.3.3 Societal Impact	10
1.3.4 Ethical considerations	10
1.3.5 Sustainability, and social considerations	10
1.4 Goals	10
1.5 Research Methodology	11
1.5.1 Design Science Approach	11
1.5.2 Quantitative and Qualitative Methods	11
1.6 Delimitations	11
1.7 Structure of the thesis	12
2 Background	12
2.1 Digital Asset Management	12
2.1.1 Choosing a DAM and the key tasks	12
2.1.2 Technological Tools Demand Continuous Organizational Adaptation	13
2.2 Artificial Intelligence	13
2.2.1 Object Detection	13
2.2.2 Anchor-free detection models	14
2.3 The Architecture of a Convolutional Neural Network	14
2.3.1 The Convolutional Operation	14
2.3.2 The Pooling Operation	15
2.3.3 Activation Functions	15
2.3.4 Structural Components of the YOLO Architecture	15
2.3.5 YOLOv11 model	16
2.4 Object Detection with YOLOv11	17
2.5 LOSS	17
2.5.1 Why to make our own and not use a service	17
2.5.2 Major background area#1#1	18
2.5.3 The YOLO model	18
2.6 Major background area#2	19
2.6.1 Major background area#2#1	19
2.6.2 Major background area#2#2	19
2.7 Related work	19
2.7.1 Major related work	19
2.7.2 Major related work	19
2.7.3 Minor related work	19
2.8 Summary	19
3 Methodology	21
3.1 Research Paradigm	21

3.2	Research Process	21
3.3	Data Collection	22
3.3.1	Sampling and Target Population	22
3.4	Experimental Design	23
3.5	Assessing reliability and validity of the data collected	23
3.5.1	Reliability	23
3.6	Validity	23
3.7	Planned Data Analysis	24
3.7.1	Data Analysis Technique	24
3.7.2	Software Tools	24
3.8	Evaluation framework	24
4	Implementation and Engineering Design	24
4.1	Class Definition and Dataset Construction	24
4.2	Image Preprocessing and Augmentatio	24
4.3	YOLOv12 Training Configuration	24
4.4	Model Evaluation and Business Fit	24
4.5	Comparative Analysis	24
4.6	Organizational and INDEK Perspective	24
5	Results and Analysis	25
5.1	Major results	25
5.2	Reliability Analysis	25
5.3	Validity Analysis	25
5.4	Discussion	25
6	Conclusions and Future work	25
6.1	Conclusions	25
6.2	Limitations	25
6.3	Future work	25
6.4	Reflections	25
	References	25
	Appendices	28
A	Appendix A: Example Appendix Title	28
B	Appendix B: Another Appendix Example	29

List of Figures

1-1	Sustainable Development Target 9.5 and 12.6	10
2-1	Illustrating the five main stages of DAM.	12
2-2	Bounding box for table with legs.	14
2-3	A simplified 2D convolution applied to an RGB image (adapted from (Prince, 2023)). . .	14
2-4	Max pooling applied to a 4×4 matrix X resulting in a 2×2 matrix Y	15
2-5	The architecture of YOLOv11, illustrating its three main components: Backbone, Neck, and Head (adapted from (Hidayatullah et al., 2025)).	16
2-6	YOLOv11 performance comparison (Ultralytics Inc., 2025).	18
2-7	Comparison between YOLOv12 and RF-DETR architectures. The YOLOv12 pipeline (left) illustrates its three main components: <i>Backbone</i> , <i>Neck</i> , and <i>Head</i> , adapted from Hidayatullah et al. (2025). The RF-DETR architecture (right) is based on multi-scale deformable attention and transformer decoders, as described in the Deformable DETR paper published at ICLR 2021 by Zhu et al. (2021).	20
3-1	Overview of the dual-track research process. The technical track (T) included data preparation, model training, evaluation, and comparison, while the business track (B) focused on stakeholder needs, class selection, workflow alignment, and organizational fit. Both tracks converged in step 9T5B to inform strategic recommendations for DAM. . . .	21
A-1	An example figure in Appendix A.	28

List of Tables

2.1	Summary of YOLO Model Evolution	16
3.1	Hardware configuration and model training parameters for YOLOv12 and RF-DETR. Both models were trained on the same cloud-based GPU instance using CUDA 12.4.	23
A.1	An example table in Appendix A.	28

List of Acronyms and Abbreviations

AI	Artificial Intelligence
CNN	Convolutional Neural Network
DAM	Digital Asset Management
DETR	Detection Transformer
DSR	Design Science Research
DT	Digital Transformation
IoU	Intersection over Union
IT	Information Technology
ML	Machine Learning
mAP	Mean Average Precision
RF-DETR	Roboflow Detection Transformer
SDG	Sustainable Development Goal
SME	Small and Medium-sized Enterprises
UX	User Experience
YOLO	You Only Look Once

1 Introduction

Design-oriented manufacturers increasingly rely on large volumes of product images to communicate with customers, showcase bespoke work, and maintain branding across channels. However, as visual assets accumulate, the lack of structured organization hinders efficient management and retrieval.

Digital Asset Management (DAM) systems address this by adding searchable metadata to images. While some modern DAM platforms offer automatic metadata tagging, these systems typically rely on general-purpose object detection models trained on broad categories. For companies with domain-specific products this provides limited value.

This thesis explores how two state-of-the-art object detection models, YOLOv12 and RF-DETR, can be fine-tuned on a company-specific dataset to support domain-relevant metadata generation. The aim is to evaluate how well these models can classify and localize objects in product imagery, and how such automation can improve asset retrieval and knowledge sharing. The study combines technical model evaluation and a business-oriented assessment of their value in a real-world SME context.

1.1 Background

The increasing integration of digital workflows has intensified the demand for scalable, intelligent systems that manage visual content. In industries where products are highly customized and aesthetics play a central role, imagery is not just supplementary, as it is the product’s primary communication medium across design reviews, client interactions, and marketing platforms. This image-centric reality creates both a strategic asset and a logistical burden: as image archives grow, manual curation and retrieval processes quickly become unsustainable.

DAM emerged in the late 1990s as organizations began grappling with the rapid increase in digital content (Krogh, 2009). More recently, the integration of Artificial Intelligence (AI) and machine learning (ML) has transformed DAM by automating key processes like image tagging, sorting, and categorization. Advanced computer vision techniques now enable systems to analyze and tag images automatically, reducing manual effort and increasing accuracy (Wu et al., 2022).

1.2 Problem

As bespoke manufacturers scale, managing digital assets—spanning product imagery, design renderings, and technical specifications—becomes essential for brand consistency and operational efficiency. However, most DAM solutions, especially open-source systems, lack the necessary automation, posing

adoption and maintenance challenges for small and medium-sized enterprises (SMEs) with limited IT infrastructure. Wu et al. studied automated metadata annotation for cultural heritage and found that AI-generated captions often oversimplify context, such as describing a medieval knight merely as a “man on a horse” (Wu et al., 2022). This reflects similar challenges in design-driven manufacturing, where internal product terminology and industry-specific references require more precise and context-aware interpretation.

A core function of DAM is image tagging, sorting, and categorization, directly influencing asset retrievability and structural organization. Although AI has been integrated into some DAM solutions, these implementations typically rely on large pre-trained models that offer broad object classification rather than domain-specific tagging and vocabulary. Recent advancements in computer vision, particularly through algorithms such as YOLO (You Only Look Once), offer an opportunity to overcome these limitations. However, deploying a YOLO-powered system in this domain requires adapting the model to the specific features and vocabulary of the manufacturing sector. Rather than training a model from scratch—a process that demands extensive annotated data and computational resources—a more feasible approach is to fine-tune a pre-trained model using company-specific data.

1.3 Purpose

This study focuses on the metadata generation stage of Digital Asset Management (DAM), particularly automated image tagging using deep learning models.

The primary aim of this thesis is to assess the feasibility and impact of a YOLO-powered DAM system that has been fine-tuned on company-specific data to address the unique needs of premium manufacturing SMEs. The research will benchmark the performance of this fine-tuned system against a conventional open-source DAM platform (ResourceSpace), focusing on improvements in asset categorization accuracy and retrieval efficiency.

1.3.1 Technical Research questions

- (a) To what extent does fine-tuning YOLOv11 and Faster R-CNN on company-specific manufacturing data improve object detection accuracy compared to a baseline model, in terms of precision, recall, and inference speed?
- (b) What are the trade-offs between YOLOv11 and Faster R-CNN in terms of tagging quality, computational cost, and integration complexity within a DAM workflow?
- (c) How do differences in model performance im-

pact the usefulness of metadata for downstream tasks such as asset retrieval and categorization?

1.3.2 Business Research questions

Technological advancements alone do not guarantee successful integration. To complement this, the business perspective assesses the organizational and strategic impact after selecting the preferred DAM system. Specifically:

- (d) What organizational and process changes are required to integrate AI-based image tagging into a manufacturing SME, and how do these changes affect knowledge structuring and internal workflows?
- (e) What barriers emerge during implementation, and how are they influenced by the organization's flexibility, strategic priorities, and project-based work culture?
- (f) How does improved metadata generation contribute to long-term business value, such as brand consistency, operational scalability, and process innovation?

1.3.3 Societal Impact

Digital transformation has a significant impact on SMEs. These companies account for approximately 60% of total turnover and value-added contributions in Sweden's private sector, employing around 65% of the workforce (Tillväxtverket, 2021). The adoption of DAM systems is an integral part of this transformation, improving operational efficiency and reducing manual work, which contributes to broader economic growth. A cost-benefit analysis of 319 SMEs found that digital transformation enhances organizational resilience, reduces operational costs, and improves long-term scalability (Teng et al., 2022).

The stakeholders of this project?

This study is structured around a systematic process encompassing data collection, annotation, model fine-tuning, and testing. These phases represent essential steps that an SME would need to undertake if they were to implement a similar AI-based solution. By addressing both the positive impacts and the possible challenges, the aim is to show if the benefits of adopting this solution justify the necessary investments and efforts. The project's outcomes are expected to contribute to academic knowledge in the field of AI-powered asset management, fostering further innovation.

1.3.4 Ethical considerations

Ethically, the project will investigate issues related to data privacy, transparency, and bias, which are critical in ensuring that automated systems operate

fairly and without unintended consequences. These concerns are highlighted in the literature on AI ethics, which emphasizes the need for clear guidelines to mitigate risks associated with autonomous decision-making (Jobin et al., 2019).

1.3.5 Sustainability, and social considerations

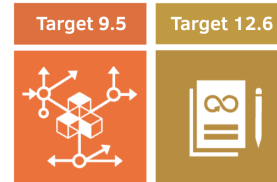


Figure 1-1: Sustainable Development Target 9.5 and 12.6

From a sustainability perspective, this research contributes to the United Nations Sustainable Development Goals (SDGs), specifically SDG 9, Industry, Innovation, and Infrastructure, and SDG 12, Responsible Consumption and Production, (United Nations, 2015). In relation to SDG 9, and more precisely target 9.5 as seen in Figure 1-1, the project seeks to enhance scientific research and upgrade the technological capabilities within industrial sectors. Similarly, under SDG 12 target 12.6 also shown in 1-1, this project supports sustainable business practices by optimizing digital asset management. By enhancing asset categorization and retrieval, the system makes it easier for companies to track and store metrics. This dual focus ensures that the technological advancements proposed are not only efficient and innovative but also ethically sound and socially beneficial.

Further reflection will be revisited in Section 6.4.

1.4 Goals

The primary goal is evaluating the feasibility of a YOLO-powered DAM system that has been fine-tuned using company-specific data, in comparison to the open-source solution ResourceSpace. To achieve this, the project has been divided into the following three sub-goals:

1. Dataset Development and Annotation:

Develop a robust methodology for collecting a domain-specific dataset that accurately captures the visual and functional nuances of digital assets in premium manufacturing. The annotation process will involve:

- Using bounding boxes to precisely delineate asset regions.
- Assigning appropriate class labels using a standardized labeling schema to ensure

consistency and relevance to the manufacturing domain.

This dataset will serve as the foundation for model fine-tuning.

2. Model Fine-Tuning and Optimization:

Fine-tune a pre-trained YOLO model on the annotated dataset. The objective is to enhance the model's accuracy in tagging, sorting, and categorizing.

- Adjusting hyperparameters and leveraging transfer learning techniques.
- Implementing regularization and validation strategies.

3. Performance Benchmarking and Comparative Analysis:

Benchmark the performance of the fine-tuned YOLO-based DAM system against a conventional open-source DAM called ResourceSpace. Evaluation metrics will include:

- Asset categorization accuracy.
- Retrieval efficiency.
- Overall system usability.

A comparative analysis will be conducted to assess whether the customized system offers significant improvements over traditional solutions. Resulting in practical recommendations and guidelines for manufacturing SMEs considering the adoption of AI-powered DAM.

1.5 Research Methodology

This research employs a mixed-methods approach to address both the technical performance of the system and stakeholder perspectives. Mixed-methods research combines quantitative techniques (e.g., controlled experiments and statistical analyses) with qualitative techniques (e.g., semi-structured interviews and thematic analysis) to provide a comprehensive evaluation of complex systems (Johnson and Onwuegbuzie, 2004).

Alternative methodologies—such as exclusively quantitative performance evaluations or purely qualitative case studies—were considered but ultimately rejected because they would not fully capture the multifaceted challenges of deploying an AI-powered system in a dynamic industrial environment.

1.5.1 Design Science Approach

Grounded in a pragmatic philosophy that emphasizes practical impact and utility, this study adopts the design science research (DSR) paradigm. DSR is particularly well-suited for technology-driven projects because it promotes the iterative design, development, and rigorous evaluation of IT artifacts

to solve real-world problems (Hevner et al., 2004). In this project, the YOLO-powered DAM system represents the artifact developed and refined through iteration.

1.5.2 Quantitative and Qualitative Methods

Controlled experiments will be conducted to measure key performance metrics—such as asset categorization accuracy, retrieval efficiency, and overall system usability. Statistical analysis will be used to validate the improvements brought about by model fine-tuning, following best practices in empirical research (Creswell, 2014; Yin, 2014). Complementing this, qualitative methods will capture contextual insights and stakeholder perspectives. Semi-structured interviews and thematic analysis will be employed to understand user experiences and organizational challenges associated with implementing the DAM system. Moreover, to develop a standardized labeling schema for the dataset, a targeted collaboration with a designated expert from the company will be undertaken. This focused approach is preferred over a large-scale survey. Not all employees interact with digital assets and the expert can ensure domain-specific terminology is accurately captured and applied consistently during annotation.

1.6 Delimitations

TO BE REMOVED AFTER COACHING

This thesis focuses exclusively on evaluating a YOLO-powered digital asset management system for premium manufacturing SMEs. The study is limited to a specific company's environment and a predefined dataset.

The research investigates only the fine-tuning of an existing pre-trained YOLOv11 model. Training a model from scratch, which requires vast amounts of data and computational resources, is beyond the scope of this project. Instead of conducting a large-scale survey, the study uses semi-structured interviews with key stakeholders—particularly a designated domain expert—to develop a standardized labeling schema.

This focused approach is chosen because only a few employees directly manage digital assets. The assessment will concentrate on technical performance indicators such as asset categorization accuracy, retrieval efficiency, and overall system usability. Broader issues such as integration with other enterprise systems and macroeconomic impacts are beyond the scope of this project.

1.7 Structure of the thesis

This thesis is organized into the following main chapters, excluding the introductory chapter, references, and appendices; Chapter 2 provides the necessary background and reviews related work, establishing the context for DAM and identifying the key gaps this project addresses. Chapter 3 outlines the methodology—including the design science approach, mixed-methods strategy, data collection, experimental design, and evaluation criteria—used to assess the system. Chapter 4 details the implementation, covering system design, model fine-tuning, dataset development, and the technical setup for testing. Chapter 5 presents the results and analysis, discussing both quantitative metrics and qualitative insights to evaluate whether the project’s goals have been met. Finally, Chapter 6 summarizes the key findings, reflects on the limitations of the study, and outlines potential directions for future work.

2 Background

2.1 Digital Asset Management

MORE TO BE ADDED / REFINE THE TEXT WITH KNOWLEDGE MANEGEMENT AND INFORMATION SHARING WITHIN ORGANIZATIONS PROJECT BASED ORGANIZATIONS AND MANAGEMENT!

Krogh (2009) describes DAM as an essential framework for protecting, organizing, and prolonging the usability of digital files by emphasizing metadata, suitable file formats, and efficient workflows. As shown in Figure 2-1, five interconnected stages—creation, management, distribution, archiving, and retrieval—collectively ensure that digital assets remain discoverable and relevant long after their initial production.

Although Krogh does not explicitly align his approach with the Resource-Based View (RBV), his emphasis on preserving assets as integral organizational resources parallels RBV’s tenet that competitive advantage relies on valuable, rare, inimitable, and non-substitutable (VRIN) capabilities (Barney, 1991). By structuring DAM processes around rigorous metadata management, secure storage, and ongoing accessibility, organizations can treat their digital repositories as strategic assets, safeguarding long-term benefits that are difficult for competitors to replicate.



Figure 2-1: Illustrating the five main stages of DAM.

2.1.1 Choosing a DAM and the key tasks

What tools are available in DAM? Benchmark? What are the most important shit in it? What do most companies need? What do they usually have and how or why do they choose to adopt a DAM

A missing perspective is
They use monday.com

2.1.2 Technological Tools Demand Continuous Organizational Adaptation

Love and Matthews (2019) identify a critical gap in the construction industry: knowing “why” to adopt digital technologies is relatively straightforward, but knowing “how” to translate technological potential into real value remains largely underexplored. Their case studies underscore the fact that digital transformation does not happen automatically; organizations must actively invest in processes such as benefits management and the development of a Business Dependency Network (BDN) to realize tangible gains from their digital initiatives (Love and Matthews, 2019).

In a broader context, Hanelt et al. (2020) posit that digital transformation (DT) goes beyond any single disruptive episode; it is a continual, structural adjustment propelled by digital technologies. Their systematic review of 279 peer-reviewed articles frames DT across three dimensions—Contextual Conditions (e.g., technological advances, shifting consumer habits), Mechanisms (e.g., the innovative strategies organizations adopt), and Outcomes (e.g., changes to organizational structures and industry norms). By proposing a typology that spans technology impact, compartmentalized adaptation, systemic shift, and holistic co-evolution, they challenge the idea of one-off change, advocating instead for an iterative, agile approach to transformation (Hanelt et al., 2020).

Taken together, these two perspectives highlight that while there is strong motivation to deploy new technologies (“why”), sustained organization-wide benefits only materialize when there is a concerted effort to integrate, evaluate, and adapt these digital tools in an ongoing manner (“how”). Both studies imply that true success hinges on long-term structural and cultural shifts rather than static, one-off solutions.

that the promise of DAM is not unlocked simply by adopting new technology but only when companies embrace two fundamental principles. First, that technology alone does not create value but must be accompanied by organizational process reengineering, and second, that the benefits of DAM are maximized only through continuous strategic governance to monitor and sustain its impact

A missing perspective in

Nevertheless, some scholars argue that resource possession alone does not guarantee successful digital transformation. Civelek et al. (2023) found no significant link between dynamic capabilities—a key aspect of RBV that involves adapting, integrating, and reconfiguring resources—and successful digital transformation among Czech manufacturing SMEs. Their findings suggest that merely possess-

ing dynamic capabilities is insufficient for digital transformation unless supported by complementary factors such as digital literacy and IT infrastructure maturity.

2.2 Artificial Intelligence

Artificial Intelligence (AI) is a field of computer science that focuses on systems built on algorithms, which are formalized sets of instructions that process input data to produce outputs (Khanam et al., 2024a). Machine Learning (ML), a subset of AI, represents a shift away from manually encoded rules toward data-driven learning. Instead of being explicitly programmed for specific tasks, ML models identify patterns in large datasets and use statistical techniques to make predictions or classify new data.

Khanam et al. (2024a) describe deep learning (DL) as a machine learning approach that utilizes multi-layered computational models to extract patterns from data at varying levels of abstraction. Inspired by the human brain, DL models excel at recognizing intricate patterns in large datasets. (Soori et al., 2023) further explains that within DL, different neural network architectures are designed to process specific types of data and perform specialized tasks. One of the most effective architectures for structured, grid-like data—such as images—is the Convolutional Neural Network (CNN). CNNs employ convolutional operations to automatically learn spatial hierarchies of features, allowing them to capture patterns and structures in data with high accuracy. As a result, CNNs have become a cornerstone of computer vision, powering applications in object detection, image classification, and other visual recognition tasks (Goodfellow et al., 2016, pp. 326-328).

2.2.1 Object Detection

Object detection involves both the ability to recognize the classes of multiple objects in an image and determining their positions, whereas image classification assigns a single class to the entire image without distinguishing individual objects.

Zhang et al. (2025) outline how DL-based object detection methods are primarily divided into two categories: two-stage and single-stage networks. Two-stage networks, such as Region-Based Convolutional Neural Networks (R-CNNs), rely on generating region proposals before classifying and refining object locations. In contrast, single-stage networks, such as You Only Look Once (YOLO), eliminate this intermediate step by predicting object classes and bounding boxes in a single pass. This approach significantly improves detection speed and efficiency. As Zhang et al. (2025) emphasize, single-stage models have become widely adopted in various industries due to their ability to perform real-time object detection accurately.

2.2.2 Anchor-free detection models

A bounding box defines an object’s position and size within an image using four coordinates. In object detection, it is paired with a class label and a confidence score, indicating both the object’s category and the model’s certainty in its prediction. These boxes act as ground-truth references in training data, helping models learn to localize objects accurately. (Li et al., 2022). The prediction represents the final output of an object detection model as illustrated in Figure 2-2.

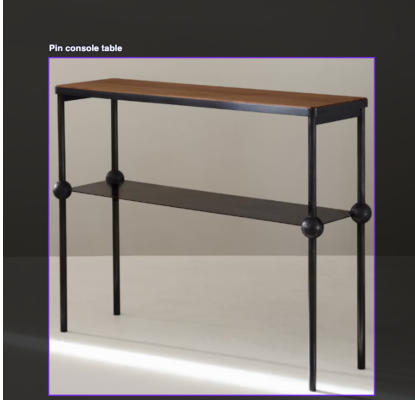


Figure 2-2: Bounding box for table with legs.

Vina (2024) describes the shift from anchor-based to anchor-free object detection as a major advancement in the field. Traditional anchor-based detectors, such as YOLOv4 and its predecessors in Table 2.1, rely on predefined anchor boxes—fixed-size reference shapes placed across an image at different aspect ratios—to estimate object locations. The model does not predict bounding boxes directly but instead modifies the closest anchor to better fit detected objects. Anchor-free models simplify detection and improve speed—critical for real-time tasks like autonomous driving and surveillance. Their key-point-based approach enhances flexibility, making them better at detecting small, irregular, or occluded objects, especially in cluttered environments where anchor-based methods struggle (Wang et al., 2024b).

2.3 The Architecture of a Convolutional Neural Network

Prince (2023) highlights three key characteristics of digital images that necessitate the use of specialized model architectures. First, images are inherently high-dimensional. For instance, a standard 224×224 pixel image with three color channels (RGB) results in over 150,000 input values. Processing such a large number of inputs with fully connected neural networks would require an impractically high number of parameters. Second, there is a

strong correlation between neighboring pixels, as local regions often form meaningful patterns and structures. Lastly, images tend to be robust to small spatial shifts—their content remains recognizable even when objects within them are slightly moved. For instance, if a chair appears slightly to the left or right in different images, we still recognize it as the same object. However, a fully connected model would need to learn how to identify the chair in every possible position from scratch. CNNs such as YOLO and Faster R-CNN avoid this problem by using filters that can detect patterns no matter where they appear in the image. This makes them far more parameter-efficient and better suited for visual tasks like object detection (Prince, 2023).

At a fundamental level, CNNs process input through sequential stages, using convolution to detect spacial features, pooling to reduce dimensionality, and activation functions to introduce non-linearity Khanam et al. (2024b). Spatial features can be textures, lines and color variations in the input. With effective training, the network learns to recognize these attributes regardless of their location within an image (Verdhan, 2021, Chapter 2).

2.3.1 The Convolutional Operation

CNNs extract features from images by applying an operation known as convolution (Prince, 2023, p. 170). Convolution involves sliding a learnable weight matrix, referred to as a kernel or filter, across the input. At each position, the kernel computes a weighted sum over a local neighborhood of the image, making it possible to detect spatial patterns. Figure 2-3 illustrates this concept. In practice, one often pads the input with zeros (padding) so that the kernel can be applied near image borders without reducing spatial dimensions. Another key hyperparameter is the stride, which specifies how far the kernel moves at each step (Prince, 2023, p. 165)

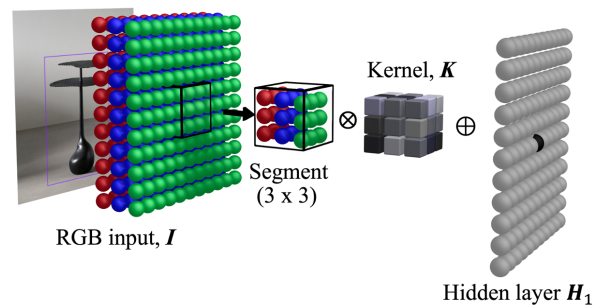


Figure 2-3: A simplified 2D convolution applied to an RGB image (adapted from (Prince, 2023)).

Let I be the input image, structured as three channels (red, green, and blue). Consider a $3 \times 3 \times 3$ kernel. At each spatial position, element-wise multiplication is performed between the kernel

weights and a 3×3 segment from each of the three channels. The products are summed together and then combined with a bias term, producing a pre-activation value that is typically passed through a non-linear function such as ReLU. By shifting the kernel step by step over the height and width of the image, one obtains a two-dimensional feature map. To produce multiple output channels, different kernels run in parallel. Each filter generates its own 2D feature map, and stacking these maps forms a three-dimensional activation tensor, often written as H_1 . Equation (1) demonstrates how the output $h_{i,j}$ at position (i, j) can be computed for an RGB input and a 3×3 kernel:

$$h_{ij} = a\left(b + \sum_{c=1}^3 \sum_{m=1}^3 \sum_{n=1}^3 I_{c, i+m-2, j+n-2} \cdot K_{c, m, n}\right) \quad (1)$$

where $I_{c, i, j}$ denotes the pixel value from channel c at position (i, j) , $K_{c, m, n}$ is the kernel weight for channel c at offset (m, n) , b is a learnable bias term, and $a(\cdot)$ represents the chosen activation function (Prince, 2023, p. 170).

2.3.2 The Pooling Operation

Pooling is a downsampling operation in CNNs that reduces the spatial dimensions of feature maps while preserving essential features. This improves computational efficiency and makes the network less sensitive to small spatial shifts. The most common method, max pooling, slides a fixed-size window over the feature map and retains only the maximum value in each region as seen in Figure 2-4 (Prince, 2023, p. 163).

$$X_{4,4} = \begin{bmatrix} 9 & 4 & 1 & 5 \\ 1 & 1 & 1 & 1 \\ 1 & 2 & 1 & 4 \\ 1 & 3 & 3 & 7 \end{bmatrix} \rightarrow Y_{2,2} = \begin{bmatrix} 9 & 5 \\ 3 & 7 \end{bmatrix}$$

Figure 2-4: Max pooling applied to a 4×4 matrix X resulting in a 2×2 matrix Y .

The latest YOLO models developed by Jocher and Ultralytics (2025) extend this concept with the SPPF (Spatial Pyramid Pooling Fast) block, which increases the receptive field through repeated pooling. Figure 2-7 shows its placement in the Neck in the architecture. The operation is defined as:

$$\text{SPPF} = \text{Conv}_{1 \times 1}(\text{Concat}(X, P_1, P_2, P_3)) \quad (2)$$

where X is the input feature map, first passed through a 1×1 convolution to reduce channel dimensions. $P_1 = \text{MaxPool}_{5 \times 5}(X)$, $P_2 = \text{MaxPool}_{5 \times 5}(P_1)$, and $P_3 = \text{MaxPool}_{5 \times 5}(P_2)$. All outputs (X, P_1, P_2, P_3) are concatenated along

the channel dimension and passed through a second 1×1 convolution. This design allows the model to capture multi-scale contextual information from increasingly larger regions while maintaining spatial resolution, which improves object detection performance, especially for small or partially occluded objects (Jocher and Ultralytics, 2025).

2.3.3 Activation Functions

The Ultralytics YOLO architecture by Jocher and Ultralytics (2025) primarily uses the Sigmoid Linear Unit (SiLU), also known as *Swish*, as its default activation function. It is defined as

$$\text{SiLU}(x) = x \cdot \sigma(x) = \frac{x}{1 + e^{-x}}, \quad (3)$$

where $\sigma(x)$ represents the sigmoid function. SiLU in Equation (3) offers a smooth non-linearity that helps the model train more efficiently and maintain stronger gradient signals in deep layers.

A simpler alternative, ReLU (Rectified Linear Unit), in Equation (4),

$$\text{ReLU}(x) = \max(0, x). \quad (4)$$

is used in certain parts of the network that benefit from faster computation and sparser activations. Additionally, some layers omit activations altogether to maintain strictly linear connections. This is sometimes useful in residual paths or when merging feature maps. However, SiLU remains the primary activation due to its observed advantages in training stability and overall performance (Jocher and Ultralytics, 2025).

2.3.4 Structural Components of the YOLO Architecture

The three-part YOLO structure consists of Backbone, Neck, and Head, and is illustrated in Figure 2-7. The Backbone extracts features using convolutional layers and downsampling, generating hierarchical feature maps. The Neck refines these features through the SPPF block for multi-scale detection and the C2PSA module to enhance the recognition of occluded objects. Upsampling and feature concatenation further improve resolution and information retention. Finally, the Head produces the model's output, predicting class probabilities and bounding boxes across three detection layers (small, medium and large), each specialized for different object sizes (Hidayatullah et al., 2025).

The C3k2 module, used in both the Backbone and Neck (Figure 2-7), acts like a compact feature extractor. It splits the input in half: one part flows through unchanged, while the other is processed by a stack of C3k blocks—convolutions with varied kernel sizes to capture both fine and coarse spatial patterns. The two paths are merged and compressed through a 1×1 convolution (Hidayatullah et al., 2025).

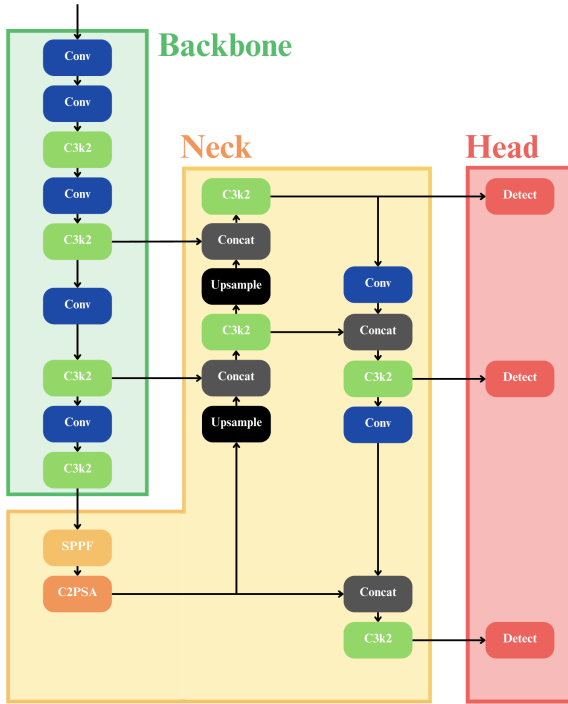


Figure 2-5: The architecture of YOLOv11, illustrating its three main components: Backbone, Neck, and Head (adapted from (Hidayatullah et al., 2025)).

The C2PSA (Cross-Stage Partial with Position-Sensitive Attention) module following after the SPPF block in Figure 2-7 extends the Cross-Stage Partial (CSP) design with a more expressive attention mechanism known as Position-Sensitive Attention (PSA). While C3k2 captures features through varied convolution kernels, C2PSA uses attention to focus on relevant spatial patterns—especially useful for detecting large objects at low resolutions (Jocher and Ultralytics, 2025).

Upsampling increases the spatial resolution of feature maps to restore details lost during downsampling. YOLO typically employs nearest-neighbor upsampling, duplicating pixels to double feature map dimensions (Jocher and Ultralytics, 2025). The subsequent concatenation merges these upsampled feature maps with earlier layers, enriching feature representations and improving multi-scale object detection capability (Figure 2-7; Hidayatullah et al., 2025).

2.3.5 YOLOv11 model

The YOLOv11 model, developed by Ultralytics marks the latest milestone in the continuous evolution of the YOLO series, building on a decade of refinement and optimization, as summarized in Table 2.1. Since its introduction by Redmon et al. (2016), it has revolutionized real-time object detection with its single-stage pipeline, offering a

faster and more efficient alternative to traditional region-based approaches like R-CNNs.

Release	Key capabilities
V1 JUN 2015	Darknet. A single-stage object detector with basic classification (Redmon et al., 2016).
V2 DEC 2016	Darknet. Object detection. Darknet-19 architecture, anchor boxes, and higher resolution inputs (Redmon and Farhadi, 2016).
V3 MAR 2018	Darknet. Object detection. Darknet-53 network & multi-scale predictions for varying object sizes. (Redmon and Farhadi, 2018).
V4 APR 2020	Darknet. Object detection. Basic object tracking with BCSPDarknet53 and SPP. (Bochkovskiy et al., 2020).
V5 JUN 2020	PyTorch. Object detection. Basic instance segmentation. Multi-GPU support, and exports (Ultralytics, 2020).
V6 SEP 2022	PyTorch. Object detection, instance segmentation, a reparameterizable backbone, anchor aided training (AAT). (Li et al., 2022).
V7 JUL 2022	PyTorch. Object detection, tracking & instance segmentation. (Wang et al., 2022).
V8 JAN 2023	PyTorch. Anchor-free object detection, instance & panoptic segmentation, NVIDIA GPUs, Jetson. (Ultralytics, 2023).
V9 FEB 2024	PyTorch. Anchor-free detection & instance segmentation. PGI for better gradient reliability. GELAN network (Wang et al., 2024b).
V10 MAY 2024	PyTorch. Anchor-free detection & NMS-free training (Wang et al., 2024a).
V11 SEP 2024	PyTorch. Anchor-free & oriented object detection (OBB), instance segmentation, pose estimation. (Ultralytics Inc., 2025).
V12 FEB 2025	PyTorch. Anchor-free detection, OBB, instance segmentation, Area Attention Mechanism, pose estimation, R-ELAN. (Ultralytics Inc., 2025).

Table 2.1: Summary of YOLO Model Evolution

Early versions of YOLO were built on the Darknet framework, developed by Joseph Redmon, with core implementations written in C and CUDA for fast GPU execution. A framework is a pre-built structure that simplifies software development by providing reusable code, tools, and libraries allowing developers to focus on higher-level abstraction. As shown in Table 2.1, the transition to PyTorch occurred with YOLOv5, developed by Ultralytics. PyTorch, originally introduced by Facebook AI Research (FAIR), offered a more flexible and scalable environment, facilitating development in Python and enhancing integration with mainstream deep learning research (Ultralytics, 2020).

Sapkota et al. (2025) conducted a comprehensive

review of YOLO-based object detection applications, highlighting its extensive adoption across multiple domains, including healthcare (e.g., pill identification, diagnostics), surveillance (e.g., face mask detection, home security), autonomous vehicles, and industrial quality control. The study underscores YOLO's efficiency in real-time processing, making it a preferred choice for applications requiring rapid inference.

While YOLO excels in speed, its grid-based detection approach and anchor-free methodology maintained in YOLOv6 and subsequent models introduce inherent limitations. Both Sapkota et al. (2025) and He et al. (2024) note that, despite its computational efficiency, YOLO may struggle with fine-grained detail detection, making it less suitable for tasks requiring high-resolution texture analysis, such as road damage assessment or material surface inspection (Angulo et al., 2019). While this thesis primarily addresses the application of YOLO within bespoke manufacturing, insights into the limitations remain highly relevant, particularly in scenarios where accurate detection and classification of subtle material textures effect performance.

The trade-off between speed and accuracy is further emphasized in comparative analyses, such as Rane (2023), which contrasts YOLO with Faster R-CNN. While YOLO excels in inference speed—making it well-suited for real-time applications such as inventory management, checkout automation, and e-commerce visual search—Faster R-CNN offers superior object localization and classification accuracy. This aligns with the findings of Sapkota et al. (2025), making it the preferred choice for scenarios demanding precise differentiation and high recall, such as medical imaging. However, Faster R-CNN's reliance on a region proposal network (RPN) results in significantly higher computational demands, limiting its viability for real-time deployment (Rane, 2023).

In contrast, the study by Karbouj et al. (2024) on object detection for screw head identification in disassembly systems presents a different perspective. Their findings demonstrate that YOLOv5 outperforms Faster R-CNN across multiple key metrics, including precision, recall, inference speed (FPS), and training efficiency. This discrepancy arises from the nature of the application and dataset size. As previously discussed by Rane (2023) Faster R-CNN tends to perform better in tasks requiring high-detail object recognition. The RPN helps it generalize more effectively when training data is limited, making it particularly useful for small datasets with high precision requirements. Conversely, YOLO's ability to efficiently learn broad patterns makes it a superior choice for large-scale, high-variance datasets. The findings of Karbouj et al.

(2024) reinforce this perspective, demonstrating YOLOv5's balance between computational speed and adaptability, making it particularly effective in real-time, resource-constrained environments.

(Alif and Hussain, 2025)

As for relating to this thesis. there is limited research on the use of YOLO directly relating for Digital Asset Management (DAM) applications. with only one identified study—Angulo et al.

citeSapkota2025YOLOv11.

The improvements of Yolov11 OLOv11 outperformed previous versions in mean average precision (mAP), recall, and precision, demonstrating superior object detection performance. The recall rate, which measures how well the model detects all ground-truth objects, was highest for YOLOv11 (64.8YOLOv11 also exhibited fewer false detections compared to its predecessors. YOLOv11 displayed higher attention concentration on relevant objects, meaning it focused better on wires and transformers, reducing errors in object localization.

2.4 Object Detection with YOLOv11

construction of a object detection dataset
 image preprocessing,
 model training using the object detection training dataset,
 and validation of results using a verification dataset

2.5 LOSS

The YOLOv11 object detection method enhances its performance by minimizing a comprehensive loss function that integrates multiple components. This loss function encompasses distributed focal loss, bounding box regression loss, and class probability loss. The optimization process involves combining these individual loss components and employing advanced optimization algorithms to refine the model's performance in object detection tasks

2.5.1 Why to make our own and not use a service

Bynder

Adobe Experince Manager

Cloudinary: custom pricing for enterprise solutions.

Adobe sensei enerally means auto-tagging images based on recognizable generic objects, scenes, and concepts. It typically uses generalized, pre-trained models that identify common objects'

most DAM platforms rely on third-party integrations for company-specific tagging

Clarifai Custom Models Provides APIs that integrate into DAM platforms.

Amazon Rekognition Custom Labels: Pay-per-use
 Google Vertex AI (formerly AI Platform Vision)
 Pricing depends on training hours and predictions
 Custom vision API: Trained specifically on your images and product labels.

Microsoft Azure Custom Vision: Training: 20
 dollar per compute hour

Integrates via REST API to enhance tagging
 accuracy in DAM solutions.

CV consulting

Image annotation

Different types of CV:

2.5.2 Major background area#1#1

Recent studies have demonstrated the effectiveness of various AI techniques in image tagging. Zhang et al. (2019) showcased the application of convolutional neural networks (CNNs) for automatic image classification in DAM systems, achieving an accuracy of 92% on a diverse dataset of digital assets

This work was further extended by Li and Chen (2020), who integrated attention mechanisms into CNNs, improving the model's ability to focus on salient features and increasing tagging accuracy to 95%

The YOLO (You Only Look Once) algorithm has also been applied successfully in DAM contexts. Wang et al. (2021) demonstrated that YOLO-based models could perform real-time object detection and tagging in DAM systems, processing up to 30 images per second with an average precision of 88%. This approach was particularly effective for identifying multiple objects within complex images, a common requirement in DAM applications.

Transformer-based models have recently gained traction in image tagging for DAM systems. A study by Rodriguez and Kim (2022) applied Vision Transformer (ViT) models to DAM image tagging, achieving state-of-the-art performance with an accuracy of 97% on standard benchmarks. The authors noted that transformer models excelled in capturing long-range dependencies in images, leading to more nuanced and context-aware tagging.

While AI-powered image tagging offers significant benefits, it also presents several challenges. Data requirements pose a significant hurdle, as highlighted by Brown et al. (2020), who found that AI models required at least 10,000 labeled images per category for optimal performance in domain-specific DAM applications

Error rates and handling domain-specific content remain ongoing challenges. A comprehensive study by Thompson et al. (2021) analyzed error patterns in AI-powered image tagging across various industries, revealing that error rates increased significantly (up to 25%) when dealing with highly specialized or technical imagery

To address this issue, Nguyen and Patel (2022) proposed a hybrid approach combining pre-trained

models with domain-specific fine-tuning, reducing error rates by 40% in niche industries such as medical imaging and aerospace engineering

Despite these challenges, the benefits of AI-powered image tagging in DAM systems are substantial. A large-scale study by Garcia et al. (2023) across 500 organizations found that implementing AI-powered tagging led to a 60% reduction in manual tagging time and a 35% improvement in asset discoverability

Entangled states are an important part of quantum cryptography, but also relevant in other domains. This concept might be relevant for neutrinos, see for example [2].

Scheme

2.5.3 The YOLO model

As demonstrated in table 2.1 the YOLO series has evolved significantly since its inception, introducing progressive improvements in object detection, computational efficiency, and feature extraction. YOLOv11 is the best choice for the project due to its superior accuracy, efficiency, and versatility. As Khanam and Hussain (2024) highlight, its architectural upgrades enhance feature extraction while minimizing computational costs, making it ideal for real-time applications requiring both speed and precision (Khanam and Hussain, 2024).

Beyond object detection, YOLOv11 supports instance segmentation, pose estimation, and oriented object detection, offering greater adaptability to the project's needs. Its optimized balance of accuracy and processing speed ensures strong performance across different computing environments, from edge devices to high-performance systems, making it the most effective solution

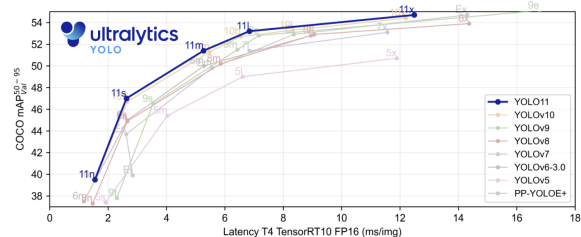


Figure 2-6: YOLOv11 performance comparison (Ultralytics Inc., 2025).

The selection of YOLOv11 for the project is driven by its superior architectural enhancements, versatile task support, and optimized balance between accuracy and efficiency. Each version has incorporated refinements aimed at enhancing real-time performance, with YOLOv11 representing the most advanced iteration to date (Khanam and Hussain, 2024).

2.6 Major background area#2

The application of AI-powered image tagging in DAM systems extends beyond large corporations to small and medium-sized enterprises (SMEs), particularly in premium manufacturing sectors. A case study by Hoffmann and Schulz (2022) examined the implementation of AI-powered DAM in a high-end carpentry company similar to Veermakers. The study found that AI-assisted tagging improved product catalog management efficiency by 45% and reduced time-to-market for new designs by 30%.

However, Chen et al. (2023) noted that SMEs in specialized manufacturing often face unique challenges in adopting AI-powered DAM systems, including limited datasets and highly specific visual content. To address these issues, the authors proposed a transfer learning approach, adapting pre-trained models to domain-specific tasks with minimal additional data, achieving a 75% reduction in required training data while maintaining 90% of the original accuracy.

While academic research has made significant strides in advancing AI-powered image tagging techniques, commercial implementations often lag behind in adopting cutting-edge methods. A comprehensive survey by Martinez and Lee (2022) of 50 leading DAM vendors revealed that only 30% had implemented transformer-based models, despite their superior performance in academic studies. The authors attributed this gap to factors such as implementation complexity, computational requirements, and the need for backward compatibility with existing systems.

2.6.1 Major background area#2#1

The integration of AI-powered image tagging in DAM systems raises important ethical, societal, and legal considerations. Privacy concerns are paramount, as highlighted by a study by Johnson and Smith (2022), which found that 35% of automatically generated tags in a sample of 10,000 images contained potentially sensitive information²². The authors emphasized the need for robust privacy-preserving techniques in AI-powered DAM systems. Algorithmic bias presents another significant challenge. Research by Park et al. (2023) revealed systematic biases in AI-generated tags across gender, ethnicity, and age dimensions, with error rates up to 20% higher for underrepresented groups. This study underscores the importance of diverse and representative training data in mitigating bias in AI-powered DAM systems.

2.6.2 Major background area#2#2

The potential impact on employment is also a concern. While Garcia et al. (2023) found that AI-powered tagging led to significant efficiency gains, they also noted a 15% reduction in human tagging roles across surveyed organizations. However, the same study observed a 10% increase in higher-skilled

positions related to AI model management and quality assurance, suggesting a shift rather than a net loss in employment.

2.7 Related work

2.7.1 Major related work

2.7.2 Major related work

2.7.3 Minor related work

2.8 Summary

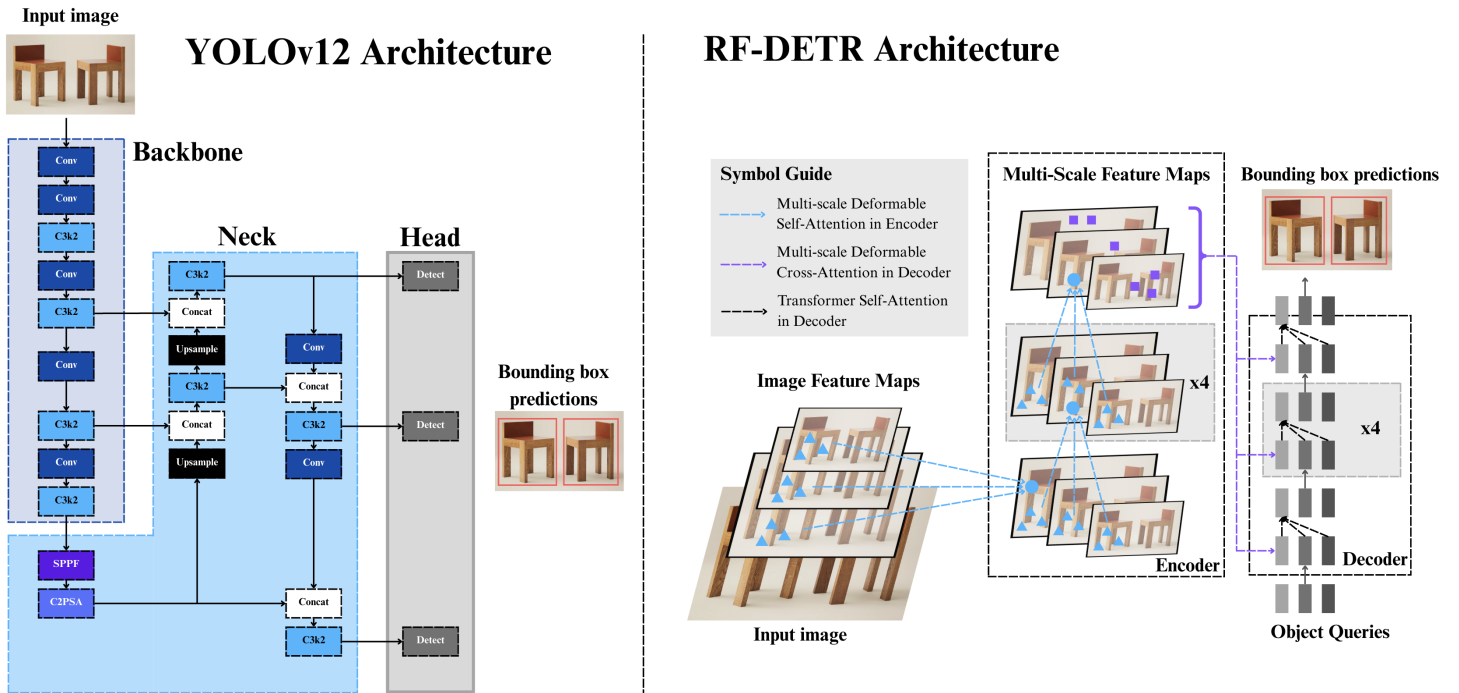


Figure 2-7: Comparison between YOLOv12 and RF-DETR architectures. The YOLOv12 pipeline (left) illustrates its three main components: *Backbone*, *Neck*, and *Head*, adapted from Hidayatullah et al. (2025). The RF-DETR architecture (right) is based on multi-scale deformable attention and transformer decoders, as described in the Deformable DETR paper published at ICLR 2021 by Zhu et al. (2021).

TEST

TEST

3 Methodology

This chapter describes the research methodology applied to develop, implement, and evaluate a computer vision system for object detection. The study follows a design science research (DSR) paradigm, with a focus on constructing and rigorously testing computational artifacts in a real-world context.

The methodology consists of a structured sequence of engineering steps: dataset construction, model training, performance evaluation, and comparative analysis. Two object detection models—YOLOv12 and RF-DETR—were trained and evaluated under equivalent experimental conditions using a domain-specific dataset derived from the partner company.

The chapter also outlines the methods used for data annotation, experimental setup, evaluation metrics, and the steps taken to ensure reliability and internal validity. The purpose is to provide a replicable, systematic framework for artifact development and performance assessment.

3.1 Research Paradigm

This thesis adopted the Design Science Research (DSR) paradigm, a problem-solving approach rooted in the sciences of the artificial (Simon, 1996), and widely applied in information systems and technology-driven domains (vom Brocke et al., 2020). DSR is suited for research that aims to generate actionable knowledge through the design and evaluation of innovative artifacts (Hevner et al., 2004).

According to Hevner et al. (2004), rigorous and relevant DSR is guided by seven core principles: (1) the creation of a purposeful artifact, (2) addressing a relevant problem, (3) evaluating the artifact’s utility, (4) contributing novel knowledge, (5) applying rigorous design and evaluation methods, (6) iterative development, and (7) clear communication to both technical and managerial audiences. By following these guidelines, DSR enables the creation of useful, tested, and transferable knowledge.

The DSR paradigm was particularly well-suited for this thesis due to its dual contribution to both technical and organizational domains. This project designed, implemented, and evaluated machine learning models (YOLOv12 and RF-DETR) for object detection within a real-world, business-driven context, specifically to enhance DAM processes in a bespoke furniture manufacturing SME.

As illustrated in Figure 3-1, the research was structured around two intertwined tracks: a technical track, which includes dataset preparation, model training, evaluation, and comparative analysis (steps 1T–8T); and a business track, which encompassed stakeholder need definition, class selection, workflow alignment, and organizational fit analysis (steps 1B–4B). These tracks converged in step 9T5B, where

insights from both domains were synthesized to inform strategic DAM recommendations.

Beyond the two functional object detection artifacts, this thesis contributes to the design knowledge (DK) base—defined as prescriptive knowledge explaining how and why a solution works in a specific context (vom Brocke et al., 2020). DK consists of three interlinked components: a clearly defined problem space, a corresponding solution, and a contextual evaluation. This type of contribution aligns with what Gregor and Hevner (2013) refer to as the instrumental knowledge base of Information Systems—practical, actionable insights that not only extend theoretical understanding, but also inform real-world decision-making and system design.

3.2 Research Process

Building on the DSR paradigm described in Section 3.1, this study followed a structured dual-track process encompassing both technical and organizational components. The overall approach is visualized in Figure 3-1, which highlights the sequence and interdependencies of the steps undertaken.

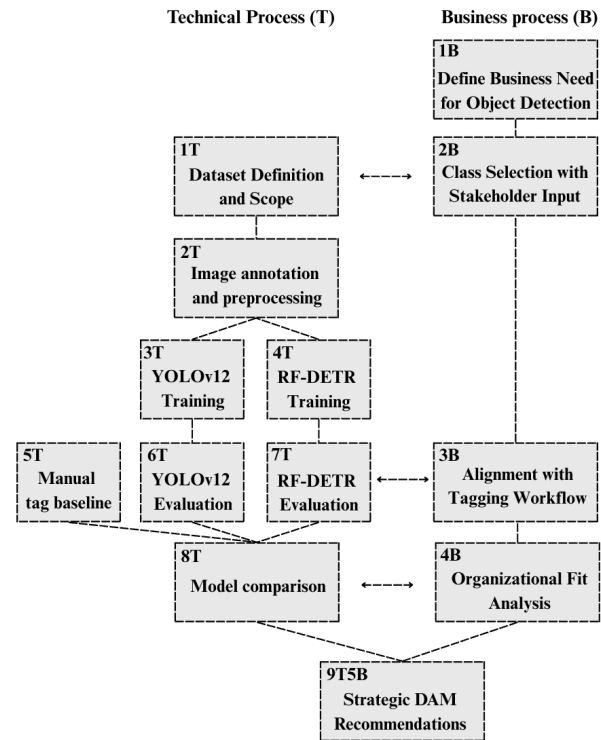


Figure 3-1: Overview of the dual-track research process. The technical track (T) included data preparation, model training, evaluation, and comparative analysis (steps 1T–8T); and a business track, which encompassed stakeholder need definition, class selection, workflow alignment, and organizational fit analysis (steps 1B–4B). These tracks converged in step 9T5B to inform strategic recommendations for DAM.

The technical track (T) in Figure 3-1 addresses the machine learning pipeline: beginning with dataset definition and annotation, progressing through model training and evaluation (YOLOv12 and RF-DETR), and culminating in comparative analysis. A manual tagging baseline was also established in step 5T. In parallel, the business track (B) defined the object detection needs, selected classes in collaboration with stakeholders, and assessed workflow integration and organizational fit. These efforts converged in the final step (9T5B), where technical and business insights were synthesized to inform actionable recommendations for improving DAM.

This structured approach ensured that technical development did not evolve in isolation but remained aligned with organizational needs—upholding the DSR principles of relevance, utility, and contextual fit. Methodological details for each step, including data collection, model selection, and evaluation strategy, are elaborated in the subsequent sections.

3.3 Data Collection

This section outlines the sources of data used in the research, as well as the ethical and legal measures taken to ensure compliance with relevant data protection regulations. Data for this study were collected through collaboration with the partner company. Three primary sources were used:

- **Internal pricing list:** An internal pricing document from 2024, detailing 50 product variants, offered in several size, material, and finish combinations.
- **Image dataset:** Approximately 4,000 photographs were sourced from the company’s Dropbox-based digital asset repository. These include high-resolution product images from multiple angles, in various environments and lighting conditions. The collected data spanned a wide range of furniture categories (e.g., stools, side tables, coffee tables, dining tables, lamps, and mirrors) and exhibited considerable variation in materials (e.g., walnut, mahogany, oak, mappa burl) and finishes (e.g., matte, hi-gloss).
- **Public website:** Product descriptions and visual marketing materials were extracted from the firm’s public-facing website.

All data were handled in accordance with the General Data Protection Regulation (GDPR) (EU, 2016/679) (European Union, 2016). The dataset primarily consisted of product-only images. In a limited number of cases, incidental human presence, such as background figures or reflected faces, was visible. These images were not annotated, labeled, or used in any way to train the models. To uphold GDPR principles of data minimization and privacy

by design (Articles 5 and 25), images with clearly identifiable individuals were manually reviewed and either excluded or masked during preprocessing. Under the Swedish Ethical Review Act (2003:460), research involving fully anonymized material does not require formal ethical approval (Swedish Parliament, 2003).

Image annotation was performed manually in a private Roboflow workspace, which supports secure, browser-based workflows and restricts public access (Roboflow Inc., 2024). Upon the completion of the annotation process (2T in Figure 3-1), the dataset was deleted in accordance with Roboflow’s data retention controls. Model training was conducted in Google Colab, a session-based environment where all data are automatically deleted at the end of each runtime session or after 12 hours of inactivity (Google Inc., 2024). No data were permanently stored on third-party servers.

3.3.1 Sampling and Target Population

The company maintained an evolving catalogue of product configurations that varied in design, material, size, and finish. For instance, the “Cloud Table” was available in multiple wood types and dimensions. While each configuration constituted a commercially distinct offering, many shared a common visual structure. In the context of object detection, such surface-level differences are not reliably distinguishable without large volumes of labeled data. To ensure model robustness and prevent overfitting to underrepresented configurations, variants with limited visual representation, such as special editions with only a few available images, were excluded. Including these would have introduced foreground–foreground class imbalance, where a few frequently occurring object classes dominate the dataset. This results in a long-tailed distribution, characterized by a small number of “head” classes with abundant samples and a large number of “tail” classes with very few. Such imbalance poses challenges for object detectors, which tend to perform well on head classes but struggle to generalize to rare categories. This issue is pronounced in single-stage detectors like YOLOv5 and its successors. Crasto (2024) demonstrate that long-tailed datasets significantly degrade mean Average Precision (mAP) for underrepresented classes, and that standard mitigation strategies such as sampling and loss weighting offer limited effectiveness in correcting this bias.

To address this, all product variants were manually consolidated into 33 high-level object classes representing the most visually distinct product archetypes, rather than fine-grained subtypes. The consolidation process prioritized three criteria: (1) visual separability between classes, (2) sufficient representation across the dataset, and (3) compatibility

with the company’s internal digital asset tagging practices. More on this in section ??

A small number of strategically important new designs were retained and supplemented with additional images from the company’s website and social media channels to ensure coverage. The majority of the dataset was sourced from the company’s internal Dropbox-based asset repository. Social media content was typically excluded due to redundancy or low image quality caused by compression.

The final annotated dataset consisted of 3,905 images, partitioned as follows:

- **Training set:** 70%
- **Validation set:** 20%
- **Test set:** 10%

This 70/20/10 split was selected to ensure effective model training, reliable validation, and unbiased final evaluation. The relatively large validation set (20%) was particularly important given the limited dataset size and the complexity of the models used. Transformer-based detectors such as Re-DETR are known to overfit when training data is scarce, partly due to their lack of inductive biases such as locality and translation invariance, which are present in convolutional models (Carion et al., 2020).

3.4 Experimental Design

Model training and evaluation were conducted in a cloud-based environment using Google Colab. The object detection models, YOLOv12 and RF-DETR, were fine-tuned on the same annotated dataset. This setup enabled a controlled and consistent baseline for architectural comparison. The software stack included the following components:

- **Annotation and Augmentation:** Roboflow Pro (web interface)
- **Environment:** Google Colab (Jupyter-based)
- **Programming Language:** Python 3.10
- **Frameworks:** PyTorch 2.0, OpenCV
- **Models:** Ultralytics (YOLOv12), Hugging Face (RF-DETR)
- **Utilities and Tooling:** supervision (visualization and benchmarking), flash-attn (CUDA-accelerated attention), tqdm (progress tracking)

YOLOv12 was installed directly from the official GitHub repository, as it was not yet available via the Python Package Index (PyPI) (Jie, 2025). API authentication for Roboflow (used for dataset access) and Hugging Face (used for RF-DETR) was securely

managed through Google Colab’s userdata module, which enabled token-based access without exposing credentials in the notebook. The detailed system configuration is summarized in Table 3.1, including GPU usage, batch sizes, and training epochs. Notably, RF-DETR was trained using a smaller per-step batch size and gradient accumulation to match YOLOv12’s effective batch size of 32.

Configuration, Resources, and Parameters		
Component	YOLOv12	RF-DETR
Framework	Ultralytics (YOLOv12s .yaml), Roboflow	Hugging Face (rfdet)
GPU	Tesla T4 (15GB VRAM)	
Max RAM Usage	13.3 / 15.0 GB	14.8 / 15.0 GB
CUDA Version	12.4	
Dataset Format	Roboflow YAML	COCO-style JSON
Input Resolution	640×640	448×448
Epochs	36	4
Batch Size	32	4 (gradient accumulation over 8 steps)
Effective Batch Size	32	

Table 3.1: Hardware configuration and model training parameters for YOLOv12 and RF-DETR. Both models were trained on the same cloud-based GPU instance using CUDA 12.4.

As shown in Table 3.1, the input resolution for YOLOv12 was set to 640×640, consistent with its default configuration. For RF-DETR, a lower resolution of 448×448 was selected to reduce memory consumption and ensure training completed within the four-hour runtime limit imposed by Google Colab’s free-tier environment (Tesla T4 GPU). To maintain comparability between the two models, the effective batch size was standardized to 32. Due to higher memory requirements, RF-DETR achieved this through gradient accumulation over eight steps, while YOLOv12 was able to use a full batch size of 32 directly.

3.5 Assessing reliability and validity of the data collected

3.5.1 Reliability

To ensure replicability: Random seeds were fixed during training and dataset splits were logged and version-controlled. All augmentation operations were deterministic and tracked.

3.6 Validity

Internal validity was addressed through strict separation of training and test sets. External validity was ensured by involving business stakeholders in

the annotation process and aligning object classes with real business needs.

3.7 Planned Data Analysis

3.7.1 Data Analysis Technique

The performance of the trained model was evaluated using:

- Mean Average Precision (mAP)
- Precision
- Recall
- Intersection over Union (IoU)

These metrics were chosen to balance localization and classification quality.

3.7.2 Software Tools

- Python, Pandas, and Matplotlib for result visualization
- Roboflow for managing annotations and generating YOLOv12-compatible labels

3.8 Evaluation framework

The trained model was compared against Faster R-CNN and current DAM systems used by the company (manually tagged folders and internal table-based search). Evaluation considered:

- Inference speed
- Accuracy
- Model complexity
- Organizational integration effort

4 Implementation and Engineering Design

This chapter presents the engineering work carried out to implement the proposed object detection system using YOLOv12. It includes decisions made during dataset design, preprocessing, model selection, and evaluation, with a focus on achieving practical, scalable results in a business context.

4.1 Class Definition and Dataset Construction

Step 1: Object Class Consolidation

From an initial set of over 100 item categories used internally by the company, 33 business-relevant and visually distinct classes were defined. This ensured:

- Efficient learning
- Business applicability
- Simplified annotation

Step 2: Image Annotation

Images were labeled using Roboflow, with bounding boxes and class labels. Resulting dataset: 3905 annotated images.

Split: • Training: 2733 images • Validation: 781 images • Testing: 391 images

4.2 Image Preprocessing and Augmentation

Preprocessing • Auto-orientation applied using EXIF data • Resolution normalization • Conversion to YOLOv12 training format

Augmentation Strategy

Each training image generated 3 augmented variants with: • Rotation: ± 10 degrees • Shear: ± 10 (horizontal), $+10$ (vertical) • Grayscale conversion: 10% of images

This improved model robustness to real-world variations.

4.3 YOLOv12 Training Configuration

Hyperparameter Value Epochs 36 Batch Size 32 Image Size 640x640 Optimizer SGD Loss Function Composite loss (bounding box, objectness, classification)

The training ran on a Tesla T4 GPU with CUDA 12.4, using 13.3/15.0 GB of VRAM during peak usage.

4.4 Model Evaluation and Business Fit

The model's performance was evaluated on the validation and test datasets using: • mAP@0.5 • Precision • Recall • IoU

Preliminary results showed high precision in detecting common components, with some limitations in overlapping or partially occluded objects.

4.5 Comparative Analysis

YOLOv12 was compared with: • Faster R-CNN: Higher accuracy on small, intricate objects but significantly slower inference • Current DAM Practice: Manual folder-based tagging, low scalability, poor metadata consistency

YOLOv12 offered the best balance between accuracy, speed, and ease of deployment.

4.6 Organizational and INDEK Perspective

This implementation was assessed through the lens of: • Workflow adaptation: Reduction in manual tagging • Task shifting: From manual metadata work to quality control • Complexity: Integration

into existing processes posed minor training barriers

- Strategic value: Better searchability, time savings, and foundation for intelligent product catalogs

To support long-term value realization, the following strategies are proposed:

- Feedback Loops: Allow users to correct or refine predictions to continuously improve the model
- Versioning and Retraining Pipelines: Enable systematic dataset expansion and periodic model updates
- Documentation Practices: Codify annotation guidelines, retraining criteria, and error cases

These components support sustainable AI integration and knowledge management in the company's digital asset processes.

5 Results and Analysis

In this chapter, we present the results and discuss them.

Keep in mind: How you are going to evaluate what you have done? What are your metrics? Analysis of your data and proposed solution Does this meet the goals which you had when you started?

5.1 Major results

5.2 Reliability Analysis

LALALA

5.3 Validity Analysis

LALALA

5.4 Discussion

6 Conclusions and Future work

«Add text to introduce the subsections of this chapter.»

6.1 Conclusions

Describe the conclusions (reflect on the whole introduction given in Chapter 1). Discuss the positive effects and the drawbacks. Describe the evaluation of the results of the degree project. Did you meet your goals? What insights have you gained? What suggestions can you give to others working in this area? If you had it to do again, what would you have done differently?

6.2 Limitations

What did you find that limited your efforts? What are the limitations of your results?

6.3 Future work

Describe valid future work that you or someone else could or should do. Consider: What you have left undone? What are the next obvious things to be done? What hints can you give to the next person who is going to follow up on your work?

6.4 Reflections

What are the relevant economic, social, environmental, and ethical aspects of your work?

References

- Alif, M. A. R. and Hussain, M. (2025). Yolov12: A breakdown of the key architectural features. *arXiv preprint*, arXiv:2502.14740v1. arXiv.org perpetual non-exclusive license.
- Angulo, A., Vega-Fernández, J. A., Aguilar-Lobo, L. M., Natraj, S., and Ochoa-Ruiz, G. (2019). *Road Damage Detection Acquisition System Based on Deep Neural Networks for Physical Asset Management*, page 3–14. Springer International Publishing.
- Barney, J. (1991). Firm resources and sustained competitive advantage. *Journal of Management*, 17(1):99–120.
- Bochkovskiy, A., Wang, C.-Y., and Liao, H.-Y. M. (2020). Yolov4: Optimal speed and accuracy of object detection. arXiv:2004.10934v1, 23 Apr 2020.
- Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., and Zagoruyko, S. (2020). End-to-end object detection with transformers. In *European Conference on Computer Vision (ECCV)*, pages 213–229. Springer.
- Civelek, M., Krajčík, V., and Ključnikov, A. (2023). The impacts of dynamic capabilities on smes' digital transformation process: The resource-based view perspective. *Oeconomia Copernicana*, 14(4):1367–1392.
- Crasto, N. (2024). Class imbalance in object detection: An experimental diagnosis and study of mitigation strategies. *arXiv preprint arXiv:2403.07113*.
- Creswell, J. W. (2014). *Research Design: Qualitative, Quantitative, and Mixed Methods Approaches*. SAGE Publications, 4th edition.
- European Union (2016). Regulation (EU) 2016/679 of the European Parliament and of the Council on the protection of natural persons. <https://eur-lex.europa.eu/legal-content/>

- [EN/TXT/?uri=CELEX%3A32016R0679](#). Accessed: 2025-03-12.
- Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. MIT Press.
- Google Inc. (2024). Google colab terms of service. <https://colab.research.google.com/pro/terms/v1>. Accessed: 2025-04-12.
- Gregor, S. and Hevner, A. R. (2013). Positioning and presenting design science research for maximum impact. *MIS Quarterly*, 37(2):337–355.
- Hanelt, A., Bohnsack, R., Marz, D., and Antunes Marante, C. (2020). A systematic review of the literature on digital transformation: Insights and implications for strategy and organizational change. *Journal of Management Studies*.
- He, Z., Wang, K., Fang, T., Su, L., Chen, R., and Fei, X. (2024). Comprehensive performance evaluation of yolov11, yolov10, yolov9, yolov8 and yolov5 on object detection of power equipment.
- Hevner, A. R., March, S. T., Park, J., and Ram, S. (2004). Design science in information systems research. *MIS Quarterly*, 28(1):75–105.
- Hidayatullah, P., Syakrani, N., Sholahuddin, M. R., Gelar, T., and Tubagus, R. (2025). Yolov8 to yolov11: A comprehensive architecture in-depth comparative review.
- Jie, S. (2025). Yolov12 - ultralytics-style object detector. <https://github.com/sunsmarterjie/yolov12>. Accessed: 2025-04-13.
- Jobin, A., Ienca, M., and Vayena, E. (2019). The global landscape of ai ethics guidelines. *Nature Machine Intelligence*, 1:389–399.
- Jocher, G. and Ultralytics (2025). Ultralytics yolov11. <https://github.com/ultralytics/ultralytics>. Accessed March 2025.
- Johnson, R. B. and Onwuegbuzie, A. J. (2004). Mixed methods research: A research paradigm whose time has come. *Educational Researcher*, 33(7):14–26.
- Karbouj, B., Topalian-Rivas, G. A., and Krüger, J. (2024). Comparative performance evaluation of one-stage and two-stage object detectors for screw head detection and classification in disassembly processes. *Procedia CIRP*, 122:527–532. 31st CIRP Conference on Life Cycle Engineering (LCE 2024).
- Khanam, R. and Hussain, M. (2024). Yolov11: An overview of the key architectural enhancements.
- Khanam, R., Hussain, M., Hill, R., and Allen, P. (2024a). A comprehensive review of convolutional neural networks for defect detection in industrial applications. *IEEE Access*, 12:94250–94295.
- Khanam, R., Hussain, M., Hill, R., and Allen, P. (2024b). A comprehensive review of convolutional neural networks for defect detection in industrial applications. *IEEE Access*, 12:94250–94295.
- Krogh, P. (2009). *The DAM book: Digital asset management for photographers*. O’Reilly, Sebastopol, California, 2nd edition.
- Li, C., Li, L., Jiang, H., Weng, K., Geng, Y., Li, L., Ke, Z., Li, Q., Cheng, M., Nie, W., Li, Y., Zhang, B., Liang, Y., Zhou, L., Xu, X., Chu, X., Wei, X., and Wei, X. (2022). Yolov6: A single-stage object detection framework for industrial applications.
- Love, P. E. and Matthews, J. (2019). The ‘how’ of benefits management for digital technology: From engineering to asset management. *Automation in Construction*, 107:102930.
- McCain, E., Mara, N., Van Malssen, K., Carner, D., Reilly, B., Willette, K., Schiefer, S., Askins, J., and Buchanan, S. A. (2021). *Endangered but not too late: The state of digital news preservation*. Donald W. Reynolds Journalism Institute, University of Missouri–Columbia Libraries. OpenAccess. Licensed under CC BY 4.0.
- Prince, S. J. (2023). *Understanding Deep Learning*. The MIT Press.
- Rane, N. (2023). Yolo and faster r-cnn object detection for smart industry 4.0 and industry 5.0: applications, challenges, and opportunities. *SSRN Electronic Journal*.
- Redmon, J., Divvala, S., Girshick, R., and Farhadi, A. (2016). You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 779–788.
- Redmon, J. and Farhadi, A. (2016). Yolo9000: Better, faster, stronger.
- Redmon, J. and Farhadi, A. (2018). Yolov3: An incremental improvement. *arXiv:1804.02767v1*.
- Roboflow Inc. (2024). Roboflow terms of service. <https://roboflow.com/terms>. Accessed: 2025-04-12.
- Sapkota, R., Qureshi, R., Flores-Calero, M., Badgular, C., Nepal, U., Poullose, A., Zeno, P., Vaddevolu, U. B. P., Khan, S., Shoman, M., Yan, H., and Karkee, M. (2025). Yolo11 to its genesis: A decadal and comprehensive review of the you only look once (yolo) series. *arXiv*, 2406(19407v5).

- Simon, H. A. (1996). *The Sciences of the Artificial*. MIT Press, Cambridge, MA, 3rd edition.
- Soori, M., Arezoo, B., and Dastres, R. (2023). Artificial intelligence, machine learning and deep learning in advanced robotics, a review. *Cognitive Robotics*, 3:54–70.
- Swedish Parliament (2003). Lag (2003:460) om etikprovning av forskning som avser människor. <https://www.riksdagen.se/sv/dokument-lagar/dokument/svensk-forfattningssamling/lag-2003460-om-etikprovning-av-forskning-som-sfs-2003-460>. Accessed: 2025-03-12.
- Teng, X., Wu, Z., and Yang, F. (2022). Impact of the digital transformation of small- and medium-sized listed companies on performance: Based on a cost-benefit analysis framework. *Journal of Mathematics*, 2022:1–15.
- Tillväxtverket (2021). Små och medelstora företags digitalisering - vad har betydelse? Technical Report 0366, Tillväxtverket. Accessed: 2025-02-15.
- Ultralytics (2020). Comprehensive guide to ultralytics yolov5. Accessed: 21 February 2025.
- Ultralytics (2023). Yolov8: A unified architecture for object detection, classification, and segmentation. <https://yolov8.com/>. Accessed: 2025-03-01.
- Ultralytics Inc. (2025). Ultralytics YOLO11. <https://docs.ultralytics.com/models/yolo11/>. Accessed: 3 March 2025.
- United Nations (2015). Transforming our world: The 2030 agenda for sustainable development. Accessed: February 28, 2025.
- Verdhan, V. (2021). *Computer Vision Using Deep Learning: Neural Network Architectures with Python and Keras*. Apress, Berkeley, CA, 1st ed. edition.
- Vina, A. (2024). The benefits of ultralytics yolo11 being an anchor-free detector. *Ultralytics Blog*.
- vom Brocke, J., Hevner, A., and Maedche, A. (2020). Introduction to design science research. In vom Brocke, J., Hevner, A., and Maedche, A., editors, *Design Science Research: Cases*, pages 1–16. Springer, Cham.
- Wang, A., Chen, H., Liu, L., Chen, K., Lin, Z., Han, J., and Ding, G. (2024a). Yolov10: Real-time end-to-end object detection.
- Wang, C.-Y., Bochkovskiy, A., and Liao, H.-Y. M. (2022). Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. *arXiv preprint arXiv:2207.02696*. Version 1, 6 Jul 2022.
- Wang, C.-Y., Yeh, I.-H., and Liao, H.-Y. M. (2024b). Yolov9: Learning what you want to learn using programmable gradient information. *arXiv preprint arXiv:2402.13616*.
- Wu, M., Brandhorst, H., Marinescu, m.-c., Moré, J., Hlava, M., and Busch, J. (2022). Automated metadata annotation: What is and is not possible with machine learning. *Data Intelligence*, 5:1–17.
- Yin, R. K. (2014). *Case Study Research: Design and Methods*. SAGE Publications, 5th edition.
- Zhang, L., Sun, Z., Tao, H., Wang, M., and Yi, W. (2025). Research on mine-personnel helmet detection based on multi-strategy-improved yolov11. *Sensors (Basel, Switzerland)*, 25(1):170–.
- Zhu, X., Su, W., Lu, L., Li, B., Wang, X., and Dai, J. (2021). Deformable detr: Deformable transformers for end-to-end object detection. In *International Conference on Learning Representations (ICLR)*.

Appendices

A Appendix A: Example Appendix Title

This is an example appendix entry. You can include figures, tables, or additional details relevant to your research.



Figure A-1: An example figure in Appendix A.

Column 1	Column 2
Data 1	Data 2
Data 3	Data 4

Table A.1: An example table in Appendix A.

B Appendix B: Another Appendix Example

You can continue adding appendices in a similar manner.

IEEE Editorial Style Manual: