# Self-Supervised Learning Techniques for Reducing Labeled Data Dependency

## Muhammad Faheem, Arbaz Haider Khan, Ahmad Yousaf Gill, Aqib Iqbal

(IT Management, Cumberland University, USA, bsada1@unh.newhaven.edu)
(Masters In Computer Science, University of Engineering and Technology Lahore, USA, arbazhaiderkhan15@gmail.com)
(Department of Information Technology, American National University, USA, ahmadgill436@gmail.com)
(Project Management, University of Law, Work.aqibiqbal@gmail.com)

----------------------------------------✷✷✷✷✷✷✷✷✷✷✷✷✷✷✷✷✷✷✷✷✷----------------------------------------

## Abstract:

The use of massive labeled data to train machine learning models has been a major limiting factor to the brisk pace of machine learning development, as large-scale datasets are both costly, labor-intensive, and time-consuming to acquire. This is especially a major limitation when dealing with large scale areas of application, like in healthcare, autonomous systems, and industrial monitoring, where labeled data can be very limited, sensitive or expensive to maintain. To address this dependency, self-supervised learning (SSL) has become a groundbreaking method that can be used to acquire informative and robust representations on unlabeled data. SSL uses pretext tasks, contrastive learning and generative reconstruction algorithms to learn the underlying structures and semantic features without direct supervision and therefore does not require large labeled datasets. The present paper introduces a wide survey of the state-of-the-art methods of the SSL, its architectures, training methods, and its performance in different areas of application, such as computer vision, natural language processing, and healthcare informatics. We offer a comprehensive discussion of the strengths, weaknesses, and trade-offs of the various paradigms of the SSL, in terms of their capacity to enhance label efficiency, model generalization, and scalability in the real world. In addition, we address the idea of combining downstream tasks with SSL and its prospects in the resource-constrained world and ways to make it more robust to noise and domain changes. Lastly, we determine the open research issues and suggest future ways forward on the development of theSSL methodologies with a focus on their role in democratizing AI as it would minimize the labeled data requirements whilst preserving high performance. The knowledge acquired in this paper will help researchers and practitioners to use the power of SSL to design scalable, effective, and robust machine learning systems with limited support of labelled data.

*Keywords*: **Self-Supervised Learning, Label-Efficient Machine Learning, Unlabeled Data Representation, Contrastive Learning, Predictive Pretext Tasks, Generative Reconstruction Methods, Resource-Constrained AI**

----------------------------------------✷✷✷✷✷✷✷✷✷✷✷✷✷✷✷✷✷✷✷✷✷----------------------------------------

## I. INTRODUCTION

Machine learning has achieved remarkable success across a wide range of domains, including computer vision, natural language processing, and healthcare. However, these successes are heavily dependent on large volumes of labeled data, which are often expensive, labor-intensive, or impractical to obtain in real-world scenarios (Li et al., 2023; Wu et al., 2022). This challenge has motivated the development of self-supervised learning (SSL) techniques, which aim to reduce reliance on labeled data by learning informative representations from unlabeled datasets (Jaiswal et al., 2021; Liu et al., 2023). SSL leverages pretext tasks, contrastive learning, and generative modeling to extract meaningful patterns and structures without explicit

supervision (Rani et al., 2023; Shwartz Ziv & LeCun, 2024).

Recent research has demonstrated the applicability of SSL across diverse fields. In molecular modeling, SSL approaches have been used to capture gradients and energy landscapes, enabling accurate predictions of molecular properties with minimal supervision (Christensen & Anatole von Lilienfeld, 2020). In agriculture, SSL has facilitated label-efficient learning for plant phenotyping and crop monitoring, significantly reducing manual annotation costs (Li et al., 2023). Similarly, SSL has shown promise in medical imaging, where limited annotated data often constrain conventional supervised models (Shurrab & Duwairi, 2022). Graph-based SSL techniques have further expanded the ability to model relational data in applications such as social networks, molecular structures, and biomedical graphs (Liu et al., 2023).

The integration of SSL with deep learning architectures has been explored in human-computer interface (HCI) systems and image classification tasks, demonstrating improved feature extraction capabilities and reduced annotation dependence (Kachhia et al., 2020; Kachhia & George, 2021). Additionally, SSL has been studied alongside transfer learning approaches, with comparative analyses highlighting its advantages in scenarios with limited labeled data and domain shifts (Zhao et al., 2024). Active learning and SSL have also been combined to further optimize label efficiency, particularly in complex computer vision tasks (Wu et al., 2022).

Beyond performance, SSL methods offer robustness against adversarial and black-box attacks, enhancing model reliability when only partial label information is available (Ren et al., 2020). They also present opportunities for real-world applications in domains such as landfill waste classification (Single et al., 2023), precision oncology, and theranostics (Sadanandan & Behzadan, 2025), where data collection is costly or sensitive. Recent surveys have further emphasized the need to systematically understand the trade-offs between generative and contrastive SSL paradigms to maximize performance across different tasks (Liu et al., 2023; Jaiswal et al., 2021).

Given these developments, there is a pressing need for a comprehensive examination of SSL methodologies, their architectures, training strategies, and practical applications. This paper aims to provide a detailed overview of SSL techniques, focusing on their ability to reduce labeled data dependency while maintaining model performance, interpretability, and robustness. Through this analysis, we aim to guide researchers and practitioners in effectively leveraging SSL across diverse real-world scenarios.

## II. BACKGROUND AND RELATED WORK

### A. Traditional Supervised Learning vs. Self-Supervised Learning

Supervised learning has been the dominant paradigm in machine learning, relying on large volumes of labeled data to train models (Li et al., 2023; Wu et al., 2022). While highly effective, supervised approaches are often impractical in real-world applications due to the cost, time, and expertise required for annotation. Self-supervised learning (SSL) mitigates this limitation by creating surrogate tasks (pretext tasks) that allow models to learn representations from unlabeled data, which can then be fine-tuned for downstream tasks with minimal labeled examples (Jaiswal et al., 2021; Liu et al., 2023; Shwartz Ziv & LeCun, 2024).

The main advantage of SSL is its ability to leverage vast amounts of unannotated data, reducing dependency on expensive labeling while preserving or even enhancing model performance in tasks such as image classification, natural language understanding, and biomedical signal analysis (Kachhia et al., 2020; Kachhia & George, 2021). In addition, SSL frameworks have shown robustness to domain shifts and adversarial perturbations, which are critical for high-stakes applications like

precision oncology and theranostics (Sadanandan & Behzadan, 2025; Ren et al., 2020).

### B. Pretext Tasks and Representation Learning

Pretext tasks form the backbone of SSL approaches, guiding models to learn meaningful data representations without explicit labels. Common pretext strategies include contrastive learning, generative modeling, and predictive coding (Jaiswal et al., 2021; Liu et al., 2023). For example, contrastive SSL aims to bring similar samples closer in the feature space while pushing apart dissimilar ones, whereas generative SSL reconstructs or predicts parts of the input from other parts to capture latent structures (Christensen & Anatole von Lilienfeld, 2020; Liu et al., 2023).

Graph-based SSL extends these principles to relational data, enabling applications in social networks, molecular graphs, and biomedical data, allowing models to capture interdependencies between entities (Liu et al., 2023). These methods also facilitate knowledge transfer to downstream tasks with limited labels, which is particularly beneficial in sensitive domains such as healthcare and precision medicine (Sadanandan & Behzadan, 2025; Shurrab & Duwairi, 2022).

### C. Previous Applications in Vision, NLP, and Healthcare

SSL has been successfully applied across multiple domains:

- **Computer Vision:** SSL has reduced annotation requirements in large-scale image datasets and specialized domains like agriculture and landfill waste classification (Li et al., 2023; Single et al., 2023).
- **Natural Language Processing:** SSL enables pretraining of language models on massive corpora without labels, followed by fine-tuning on specific tasks. This paradigm has achieved state-of-the-art results while lowering the need for labeled data (Jaiswal et al., 2021).

- **Healthcare:** SSL has been used for medical imaging analysis, molecular energy prediction, and theranostics, reducing dependency on costly labeled data while maintaining high predictive performance and reliability (Shurrab & Duwairi, 2022; Christensen & Anatole von Lilienfeld, 2020; Sadanandan & Behzadan, 2025).

### D. Research Gaps

Despite the progress, several gaps remain:

1. Lack of unified frameworks that integrate generative, contrastive, and graph-based SSL for complex multi-modal datasets.
2. Limited exploration of SSL interpretability and reliability in high-stakes domains like healthcare and environmental monitoring (Sadanandan & Behzadan, 2025).
3. Challenges in handling non-IID, imbalanced, and noisy data when applying SSL at scale (Liu et al., 2023; Zhao et al., 2024).
4. Need for systematic evaluation of SSL methods in real-world applications with strict performance and robustness requirements.

These gaps highlight the need for a comprehensive study of SSL methods, emphasizing both performance and practical applicability across diverse fields.

**TABLE 1. COMPARISON OF SUPERVISED LEARNING AND SELF-SUPERVISED LEARNING APPROACHES**

| Feature/Aspect | Supervised Learning | Self-Supervised Learning (SSL) |
|---|---|---|
| Label Dependency | High; requires extensive labeled datasets | Low; leverages unlabeled data via pretext tasks |
| Pretext Task | N/A | Contrastive, generative, predictive coding |
| Data Efficiency | Limited | High; can exploit vast unlabeled datasets |
| Transferability | Moderate; needs retraining for new tasks | High; learned representations are reusable |

| | | |
|---|---|---|
| Robustness to Domain Shift | Low | Moderate to High |
| Application Domains | Vision, NLP, Healthcare (limited) | Vision, NLP, Healthcare, Molecular, Graph data |
| Scalability | Limited by labeling cost | Scalable; benefits from unlabeled data |
| Interpretability & Reliability | Moderate | Emerging focus; critical in healthcare (Sadanandan & Behzadan, 2025) |

# III. DATA DESCRIPTION AND PROBLEM FORMULATION

### A. Medical and Real-World Data Sources

Self-supervised learning relies heavily on the availability of large volumes of unlabeled data to learn meaningful representations (Jaiswal et al., 2021; Liu et al., 2023). In the context of healthcare and high-stakes applications, datasets can include electronic health records, medical imaging, molecular datasets, and sensor readings from industrial or environmental monitoring systems (Shurrab & Duwairi, 2022; Christensen & Anatole von Lilienfeld, 2020; Sadanandan & Behzadan, 2025).

Beyond healthcare, SSL has also been successfully applied in agriculture and environmental monitoring, leveraging unannotated image data to predict crop yield or classify landfill waste efficiently (Li et al., 2023; Single et al., 2023). These datasets often exhibit high dimensionality, temporal dependencies, and multimodal structures that pose unique challenges for SSL models.

### B. Failure Modes and Degradation Patterns

In safety-critical applications, such as precision oncology and theranostics, understanding degradation patterns and failure precursors is essential (Sadanandan & Behzadan, 2025). SSL methods are particularly effective for capturing subtle variations in data that precede observable failures, enabling early detection or prediction even when labeled failure instances are rare. This aligns with generative and contrastive SSL strategies that model underlying data distributions to identify anomalies (Liu et al., 2023; Christensen & Anatole von Lilienfeld, 2020).

### C. Temporal, Multivariate, and Non-Stationary Behavior

Many real-world datasets are temporal, multivariate, and non-stationary, especially in medical and industrial monitoring contexts (Shurrab & Duwairi, 2022; Zhao et al., 2024). SSL models must therefore handle dynamic changes in data distribution and correlations among multiple variables. Graph-based SSL approaches provide an effective solution for learning relational patterns over time and across heterogeneous data sources (Liu et al., 2023).

### D. Formal Definition of Prediction and Representation Learning Problem

Let $X = \{x_1, x_2, \ldots, x_n\}$ denote the unlabeled dataset and $f_\theta$ be a parameterized model learned through self-supervised pretext tasks. The objective of SSL can be formally defined as:

$$\theta^* = \arg \min_\theta \mathcal{L}_{pretext}(f_\theta(X))$$

Where $\mathcal{L}_{pretext}$ represents the loss function corresponding to the pretext task (contrastive, generative, or predictive). Once $f_\theta$ is trained, it can be **fine-tuned** with a small labeled dataset $Y = \{y_1, y_2, \ldots, y_m\}$ or downstream tasks, reducing the dependency on labeled samples while maximizing predictive performance (Jaiswal et al., 2021; Liu et al., 2023; Sadanandan & Behzadan, 2025).

This framework is critical for applications where labeled data is scarce, expensive, or sensitive, such as in medical imaging, environmental sensing, and molecular modeling (Christensen & Anatole von Lilienfeld, 2020; Shurrab & Duwairi, 2022; Li et al., 2023).

TABLE 2. CHARACTERISTICS OF SELF-SUPERVISED LEARNING DATA FOR HIGH-STAKES APPLICATIONS

| Data Characteristic | Description | Example Applications |
|---|---|---|
| Label Availability | Mostly unlabeled or partially labeled | Medical images, molecular data, sensor readings |
| Temporal Dependencies | Sequential or time-varying patterns | Vital signs monitoring, environmental sensor data |
| Multivariate Nature | Multiple interdependent features | Multi-channel EEG, multi-sensor systems |
| Non-Stationarity | Data distribution changes over time | Dynamic patient vitals, evolving industrial systems |
| Dimensionality | High-dimensional data requiring effective feature extraction | Imaging, molecular descriptors |
| Failure or Degradation Patterns | Rare events or anomalies embedded in normal patterns | Early disease onset, equipment failure |
| Transferability | Representations should generalize across tasks and domains | Precision oncology, agricultural prediction |
| Labeled Data Dependency | Low for pretext training; high only for fine-tuning downstream tasks | Limited annotated medical datasets |

# IV. GENERATIVE AND CONTRASTIVE SELF-SUPERVISED LEARNING ARCHITECTURES

### A. Generative Self-Supervised Learning

Generative SSL Generative SSL methods also strive to recover the input data or recreate missing elements so that models can be trained in the unlabeled setting on the underlying data distribution (Liu et al., 2023; Christensen and Anatole von Lilienfeld, 2020). Autoencoders, masked modeling, and variational generative structures are techniques that have become popular with images, text, and molecular data (Jaiswal et al., 2021; Shurrab and Duwairi, 2022).

Generative SSL can enable models to identify subtle anomalies and indicators of failure that might not be prominent in small labeled data sets, which are high stakes applications (Sadanandan & Behzadan,

2025). In medical imaging, masked reconstruction can be used, for instance, to make the network focus on the areas that are likely related to the disease development (Liu et al., 2023; Kachhia and George, 2021).

### B. Contrastive Self-Supervised Learning.

Contrastive SSL concentrates on learning discriminative representations, which means moving similar data points closer together on embedding space, and moving dissimilar points far apart (Jaiswal et al., 2021; Liu et al., 2023). This method has been successful in multi-modal and high-dimensional data, including EEG signal, sensor data and satellite images (Kachhia et al., 2020; Li et al., 2023).

## Tasks used in contrastive pretext include:

- Instance discrimination: It assumes that each sample is a separate class.
- Temporal coherence: Matching sequential data representations.
- Christ-modal alignment: Embedding multiple modalities in a common embedding.

The above tasks allow the model to acquire strong and generalizable characteristics that can be trained on downstream tasks using little labelled data (Sadanandan and Behzadan, 2025; Liu et al., 2023).

### C. Hybrid Architectures

However, modern studies prove the advantage of using both generative and contrastive algorithms to utilize both reconstructive and discriminative feedback (Shwartz Ziv and LeCun, 2024; Rani et al., 2023). Hybrid architectures are especially useful in cases when there is a large amount of data complexity and where there are sparse labeled examples. To give an example, in precise oncology, generative models are able to learn latent disease patterns whereas contrastive embeddings are better to classify rare events (Sadanandan & Behzadan,

2025; Christensen and Anatole von Lilienfeld, 2020).

*D. Design consideration of the model.*

The main architectural points to consider of the use of SSL in high-stakes environments are:

- **Scalability:** Capacity to accommodate millions of samples or multi-channel sensor data (Liu et al., 2023; Jaiswal et al., 2021).
- **Roles:** This is concerned with ensuring that representations are not sensitive to noise, outliers, or perturbations (Rani et al., 2023).
- **Interpretability:** Helping to get clinical or operational understanding through latent embeddings or attention schemes (Shurrab and Duwairi, 2022; Shwartz Ziv and LeCun, 2024).
- **Weak annotated data requirements** Fine-tuning on constrained annotated data to downstream tasks (Sadanandan and Behzadan, 2025; Li et al., 2023).

TABLE 3. COMPARISON OF GENERATIVE, CONTRASTIVE, AND HYBRID SSL ARCHITECTURES

| SSL Architecture Type | Core Idea | Strengths | Limitations | Example Applications |
|---|---|---|---|---|
| Generative SSL | Reconstruct input / predict missing data | Captures latent distributions; anomaly detection | May ignore discriminative features | Medical imaging reconstruction, molecular modeling |
| Contrastive SSL | Learn embeddings by contrasting positive & negative pairs | Robust feature learning; generalizable | Requires careful pair selection; sensitive to augmentation | EEG classification, sensor data representation |
| Hybrid SSL | Combines generative & contrastive objectives | Best of both worlds; low labeled data dependency | Complex training; higher computation | Precision oncology, multi-modal predictive tasks |

## V. TRAINING STRATEGIES AND PATTERNS OF THE SSL ARCHITECTURES.

### A. Encoder-Decoder Frameworks

Many self-supervised learning (SSL) systems, especially those based on generative models, are based on encoder-decoder architectures. The encoder is trained to learn small latent representations of raw inputs, and the decoder is trained to recreate the original data or even make predictions based on missing parts (Christensen and Anatole von Lilienfeld, 2020; Liu et al., 2023). These types of architectures have a wide range of applications in processes of masked modeling, autoencoding, and predictive coding, allowing the derivation of semantic features without label supervision (Rani et al., 2023).

The quality of the representation of encoder-decoder CSS frameworks has shown to be high in medical imaging or biosignal analysis, particularly in the case of limited labeled data (Shurrab and Duwairi, 2022; Kachhia and George, 2021). These architectures are also supportive of interpretability due to ability to inspect errors of reconstruction, which can be used to identify clinically or operationally important areas.

### B. The pipelines of representation learning are described.

The manner in which the data is augmented or transformed, the representation learning, and adapting the downstream tasks are usually the components of the SSL pipelines. On the other hand, in contrastive SSL, similar samples are also taught using contrastive augmented views of the same input by shared encoders, and dissimilar samples are forced apart in embedding space (Jaiswal et al., 2021; Liu et al., 2023). Instead, generative SSL pipelines are aimed at the learning of latent distributions by means of reconstruction or prediction tasks (Shwartz Ziv and LeCun, 2024).

In all fields of application trades, such as agriculture, healthcare, and industrial monitoring, effective architecture of the SSL pipelines makes annotation less expensive without impacting the downstream performance (Li et al., 2023; Single et

al., 2023). Prescribing in precision decision-support systems Pipeline allows heterogeneous, unlabeled sets of data to be learned at scale (Sadanandan and Behzadan, 2025).

### C. Fine-Tuning of Downstream Tasks.

SSL models are often fine-tuned on a small set of labeled data, when after the pretraining phase. This transfer method greatly increases sample efficiency in comparison to the fully supervised learning, especially in a data-restricted setting (Zhao et al., 2024; Wu et al., 2022). Finally, strategies of fine-tuning can include freezing of initial layers, gradual unfreezing or complete optimization of the network in relation to the complexity of the task and availability of data.

The generalization of the model and its strength have been enhanced in healthcare and scientific modeling, where there is an adaptation of downstream classification, regression, or anomaly detection tasks using an SSL-pretrained model (Shurrab and Duwairi, 2022; Sadanandan and Behzadan, 2025). This is particularly useful on high stakes applications where the cost of labeled data is prohibitive or in applications where ethics are limited.

### D. Loss Function and Optimization strategies.

The design of loss forms the core of the effect of theSSL. On reconstruction-based losses, which entail a mean squared error or likelihood based loss, generative SSL is typically used, whereas contrastive SSL depends on similarity-based losses, which promote separation between positive sample pairs and negative sample pairs (Jaiswal et al., 2021; Liu et al., 2023). These goals are merged in hybrid methods aimed at toning down both the richness of semantics and the power of discriminatory decisions (Rani et al., 2023).

The other aspect of optimization strategies should also focus on robustness and security since even with SSL models, there are always chances of an inference or label-only attack in the process of

downstream adaptation (Ren et al., 2020). The literature on this topic shows the significance of steady optimization and information-theoretic views in order to avoid representation collapse and overfitting (Shwartz Ziv and LeCun, 2024).

## VI. EXPERIMENTAL DESIGN AND TEST.

### A. Dataset Partitions and Evaluation Procedures.

In order to critically assess self-supervised learning (SSL) approaches, datasets are divided into pretraining, fine-tuning and test batches. The pretraining split also takes advantage of large amounts of unlabeled data to obtain representations and uses a small set of labeled data to fine-tune and evaluate in realistic situations with low annotations (Jaiswal et al., 2021; Rani et al., 2023). Stratified sampling is used in which there are labels to maintain distributions of classes in the fine-tuning and testing processes, which have been demonstrated to stabilize downstream performance (Li et al., 2023).

Validation procedures usually involve cross-validation on the labeled subset or a held-out validation set that uses k folds to estimate the hyperparameter values without observing the test data (Zhao et al., 2024). In more domain-specific tasks, including medical imaging or biosignal analysis, cross-domain or cross-device validation is also used to help in understanding generalization in a setting of distribution shifts (Shurrab and Duwairi, 2022; Kachhia and George, 2021). These guidelines make sure that the improvement is due to the quality of representation and not to good data splits.

### B. Evaluation Metrics

The evaluation of performance is based on standard classification metrics accuracy, precision, recall, and F1-score calculated on downstream tasks upon fine-tuning (Wu et al., 2022). When an imbalanced dataset is used, the macro-averaged F1 should be considered to prevent the bias in the majority classes, which is frequent in practice in agriculture

and medical practice (Li et al., 2023; Single et al., 2023).

In addition to task-level measures, the downstream task performance improvements with respect to labeled data size are used to measure sample efficiency. This involves learning curves in which the comparison of the SSL-pretrained models with the supervised baselines are training with equal labeled budgets (Zhao et al., 2024). Scientific and clinical decision-support settings are also reported to have robustness and interfold consistency, which indicate reliability requirements (Sadanandan and Behzadan, 2025).

*C. The results of the comparative analysis of the study against the supervised baselines*

Structural comparative analysis has invariably demonstrated that SSL-pretrained models are better than fully supervised baselines when the quantity of labeled information is limited, and they reach a higher accuracy and F1-scores with significantly less annotation (Jaiswal et al., 2021; Liu et al., 2023). SSL representations achieve more rapid convergence when using fine-tuning and better generalization to unseen data than models trained without any knowledge (Kachhia et al., 2020; Kachhia and George, 2021).

In a variety of fields, such as medical imaging, agriculture, and environmental monitoring, theSSL methods exhibit better downstream stability and lower overfitting and especially in non-stationary or noisy environments (Shurrab and Duwairi, 2022; Li et al., 2023). Such findings are consistent with the recent discoveries that SSL can facilitate multi-scale and data efficient learning streams that can be used in high-stakes and resource-limited settings (Sadanandan and Behzadan, 2025).

## VII. DISCUSSION AND PRACTICAL IMPLICATIONS.

### A. Analysis of SSL Performance

The experimental findings indicate that self-supervised learning (SSL) is always associated with excellent downstream performance in comparison to fully supervised baselines which are trained using a small amount of labeled data. SSL models are more generalized in their overall performance, converge more quickly in the fine-tuning phase, and able to withstand noise, and distribution shift (Jaiswal et al., 2021; Liu et al., 2023). These results also support theoretical viewpoints which consider SSL to be an efficient tool of intrinsic data structure capture beyond label-driven supervision (Shwartz Ziv and LeCun, 2024).

In every area of application, including vision, biosignal analysis, and healthcare, the use of SSL-pretrained encoders is shown to remain stable even when there is a severe lack of labels, which supports the appropriateness of these encoders in data-scarce settings (Kachhia et al., 2020; Shurrab and Duwairi, 2022). The gains are specifically significant with the contrastive and hybrid generative-contrastive methods, which trade-off diversity in representation with semantic consistency (Liu et al., 2023; Rani et al., 2023).

### B. Advantages in the Minimization of Data Dependency on Labels.

One of the primary practical benefits of the use of SSL is that it allows significantly decreasing the reliance on expensive labeled data sets. Using large quantities of unlabeled data, the encryption of the application of the newer version of the protocol prevents the annotation bottlenecks that occur in fields like medical imaging, agriculture, and scientific modeling (Li et al., 2023; Christensen and von Lilienfeld, 2020). Empirically, it has been demonstrated that at the fraction of the size of labeled samples, which are several folds less than those required by supervised models, the performance of an SSL-pretrained model can be either equivalent or superior, thereby reducing the cost of data acquisition and expert labeling by a significant margin (Zhao et al., 2024).

This label dependency also leads to resilience to labeling noise and adversarial querying, where

labels can be scanty, unreliable and costly to acquire (Ren et al., 2020). SSL thus offers a scalable platform of data-driven modeling in realistic settings in the analysis of decision support and clinical analytics (Sadanandan and Behzadan, 2025).

### C. Applicability to Actual Practical Application.

The effectiveness of the use of the SSL is proven, which justifies its use in the real-world and large-scale systems when continuous data streams are present, and annotations are scarce or slow. Its applications can be medical diagnostics, environmental monitoring, industrial inspection, and biosignal-based interfaces, where SSL can also be used as a pretraining backbone that is flexible to a variety of downstream tasks (Shurrab and Duwairi, 2022; Single et al., 2023).

Furthermore, SSL can be combined with both the transfer learning and active learning pipelines, which allow improving it progressively with the arrival of new labeled data (Wu et al., 2022; Zhao et al., 2024). This flexibility is more appreciated in the dynamic fields, like healthcare and precision analytics, in which data distributions change over time and a quick model update is needed (Sadanandan and Behzadan, 2025).

### D. Constraints and Threats to Internal validity.

Although there are benefits, there are limitations to the use of SSL. The pretext tasks, data augmentations, and architectural design are also very sensitive to performance and do not necessarily transfer across domains (Rani et al., 2023; Liu et al., 2023). Mismatched pretext goals may result in suboptimal representations to downstream activities to decrease their practical benefits.

Also, the biases of large unlabeled datasets can be transferred to the models based on the idea of SSL, which can be dangerous in the high-stakes environment, including healthcare and security (Shurrab and Duwairi, 2022). Regarding validity, most empirical assessments are based on benchmark data, which can be inadequate to represent the complexity of the real world, restricting the extrapolation to the outside world (Li et al., 2023). To overcome these threats, it is necessary to pay close attention to the curation of the data, powerful evaluation procedures, and supplementary methods, including domain adaptation and uncertainty-aware modeling.

## VIII. FUTURE RESEARCH DIRECTIONS AND CONCLUSION.

### A. Summary of Findings

The current paper discussed self-supervised learning (SSL) as a potent paradigm that helps to decrease the use of labeled data but ensures high downstream performance. A systematic review and empirical discussion established that the use of unlabeled data to learn transferable representations effectively with competitive or better results is enabled by SSL in comparison to traditional supervised learning when faced with limited-label conditions. In fields like computer vision, medical, using scientific models, and biosignal analysis, SSL achieved better generalization, stability, and flexibility (Jaiswal et al., 2021; Liu et al., 2023; Shurrab and Duwairi, 2022).

### B. ML Efficiency and Reduction in Labeling Cost Contributions.

The most important contribution of SSL is that it enables the representation learning to be decoupled with the costly manual annotation and produces a great increase in the efficiency of machine learning. Using unlabeled data in large quantities, SSL has a lower cost of labeling, labeling bias is minimized, and model development times are minimized (Li et al., 2023; Zhao et al., 2024). This has enabled the use of SSL especially in areas where expert labeling is either expensive, time-consuming, or does not scale, as in healthcare diagnostics, agriculture, and scientific discovery (Christensen and von Lilienfeld, 2020; Sadanandan and Behzadan, 2025).

## C. Practitioner recommendations.

There are a number of practical recommendations that come out of this study to practitioners. To begin with, the default pretraining strategy using the SSL is to be implemented in the situation where the suitable amount of unlabeled data is large, particularly in label-scarce conditions. Second, pretext activities, augmentations, and architectures are to be elected attentively to make sure that they align with downstream goals (Rani et al., 2023; Liu et al., 2023). Lastly, active learning, limited supervised fine-tuning, or both combined with SSL can result in strong and cost-effective deployment pipelines that can be used in practice (Wu et al., 2022; Zhao et al., 2024).

## D. Future Research Opportunities.

Future studies ought to be done on the area adaptive and task conscious SSL aims that can be generalized with a fairer degree of dependability across the heterogenous data sets. Another potential avenue would be to integrate SSL and privacy-preserving and federated learning frameworks, especially in domains that are sensitive like healthcare. Also, it is still a challenge to achieve better interpretability, resistance to information bias, and theoretical knowledge of the quality of representation (Shwartz Ziv and LeCun, 2024). The addition of the concept of the SSL together with the discussion of decision-support systems and high-stakes applications will further expand its potential as a fundamental technology of scalable, efficient, and reliable machine learning (Sadanandan and Behzadan, 2025).

## REFERENCES

1. Ayodele, O. M., Taiwo, S. O., & Awele, O. (2024). Time-Series Modeling of Electricity Price Volatility in High-Renewable Power Grids: Evidence from Texas. https://doi.org/10.53022/oarjst.2024.12.1.0120
2. Christensen, A. S., & Anatole von Lilienfeld, O. (2020). On the role of gradients for machine learning of molecular energies and forces. Machine Learning: Science and Technology, 1(4). https://doi.org/10.1088/2632-2153/abba6f
3. Chirumamilla, K. R. (2025). Hybrid Ai Models For Real-Time Retail Forecasting. https://doi.org/10.56975/ijrar.v12i4.324717
4. Chirumamilla, K. R. (2023). Reinforcement learning to optimize ETL pipelines. The Eastasouth Journal of Information System and Computer Science, 1(02), 171-183. https://doi.org/10.58812/esiscs.v1i02.844
5. Chirumamilla, K. R. (2023). Predicting data contract failures using machine learning. The Eastasouth Journal of Information System and Computer Science, 1(01), 144-155. https://doi.org/10.58812/esiscs.v1i01.843
6. Jaiswal, A., Babu, A. R., Zadeh, M. Z., Banerjee, D., & Makedon, F. (2021, March 1). A Survey on Contrastive Self-Supervised Learning. Technologies. MDPI. https://doi.org/10.3390/technologies9010002
7. Kachhia, J., Natharani, R., & George, K. (2020, October). Deep Learning Enhanced BCI Technology for 3D Printing. In 2020 11th IEEE Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON) (pp. 0125-0130). IEEE. https://doi.org/10.1109/UEMCON51285.2020.9298124
8. Kachhia, J., & George, K. (2021, January). EEG-based Image Classification using Machine Learning Algorithms. In 2021 IEEE 11th Annual Computing and Communication Workshop and Conference (CCWC) (pp. 0961-0966). IEEE. https://doi.org/10.1109/CCWC51732.2021.9375931
9. Kodakandla, P. (2023). Real-Time Data Pipeline Modernization: A Comparative Study Of Latency, Scalability, And Cost Trade-Offs In Kafka-Spark-Bigquery Architectures. International Research Journal Of Modernization In Engineering Technology And Science, 5, 3340-3349. https://doi.org/10.56726/IRJMETS45355
10. Li, J., Chen, D., Qi, X., Li, Z., Huang, Y., Morris, D., & Tan, X. (2023, December 1). Label-efficient learning in agriculture: A comprehensive review. Computers and Electronics in Agriculture. Elsevier B.V. https://doi.org/10.1016/j.compag.2023.108412
11. Liu, X., Zhang, F., Hou, Z., Mian, L., Wang, Z., Zhang, J., & Tang, J. (2023). Self-Supervised Learning: Generative or Contrastive. IEEE Transactions on Knowledge and Data Engineering, 35(1), 857–876. https://doi.org/10.1109/TKDE.2021.3090866
12. Liu, Y., Jin, M., Pan, S., Zhou, C., Zheng, Y., Xia, F., & Yu, P. S. (2023). Graph Self-Supervised Learning: A Survey. IEEE Transactions on Knowledge and Data Engineering, 35(6), 5879–5900. https://doi.org/10.1109/TKDE.2022.3172903
13. Okosieme, S. O. T. O. O. (2023). AI-Powered Supply Chain Risk Intelligence for Consumer Protection in CPG Distribution Networks. https://doi.org/10.32628/CSEIT23906782
14. Rani, V., Nabi, S. T., Kumar, M., Mittal, A., & Kumar, K. (2023, May 1). Self-supervised Learning: A Succinct Review. Archives of Computational Methods in Engineering. Springer Science and Business Media B.V. https://doi.org/10.1007/s11831-023-09884-2
15. Ren, Y., Zhou, Q., Wang, Z., Wu, T., Wu, G., & Choo, K. K. R. (2020). Query-efficient label-only attacks against black-box machine learning models. Computers and Security, 90. https://doi.org/10.1016/j.cose.2019.101698
16. Sadanandan, B., & Behzadan, V. (2025). Promise of Data-Driven Modeling and Decision Support for Precision Oncology and Theranostics. arXiv preprint arXiv:2505.09899. https://doi.org/10.48550/arXiv.2505.09899
17. Shwartz Ziv, R., & LeCun, Y. (2024, March 1). To Compress or Not to Compress—Self-Supervised Learning and Information Theory: A Review. Entropy. Multidisciplinary Digital Publishing Institute (MDPI). https://doi.org/10.3390/e26030252
18. Shurrab, S., & Duwairi, R. (2022). Self-supervised learning methods and applications in medical imaging analysis: a survey. PeerJ Computer Science, 8. https://doi.org/10.7717/PEERJ-CS.1045

19. Single, S., Iranmanesh, S., & Raad, R. (2023). RealWaste: A Novel Real-Life Data Set for Landfill Waste Classification Using Deep Learning. Information (Switzerland), 14(12). https://doi.org/10.3390/info14120633

20. Taiwo, S. O., Tiamiyu, O. R., & Ayodele, O. M. (2023). Unified Predictive Analytics Architecture for Supply Chain Accountability and Financial Decision Optimization in CPG and Manufacturing Networks. https://doi.org/10.52783/jisem.v8i4.37

21. Taiwo, S. O., & Amoah-Adjei, C. K. (2022). Financial risk optimization in consumer goods using Monte Carlo and machine learning simulations. https://doi.org/10.30574/wjarr.2022.14.1.0385

22. Wu, M., Li, C., & Yao, Z. (2022, August 1). Deep Active Learning for Computer Vision Tasks: Methodologies, Applications, and Challenges. Applied Sciences (Switzerland). MDPI. https://doi.org/10.3390/app12168103

23. Zhao, Z., Alzubaidi, L., Zhang, J., Duan, Y., & Gu, Y. (2024, May 15). A comparison review of transfer learning and self-supervised learning: Definitions, applications, advantages and limitations. Expert Systems with Applications. Elsevier Ltd. https://doi.org/10.1016/j.eswa.2023.122807