# INFORMATION THEORY IN OPEN-WORLD MACHINE LEARNING FOUNDATIONS, FRAMEWORKS, AND FUTURE DIRECTIONS

**Lin Wang**[*]
Shenzhen Key Laboratory of Neuropsychiatric Modulation,
Shenzhen-Hong Kong Institute of Brain Science,
Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen 518055, China
l.wang@siat.ac.cn

October 20, 2025

## ABSTRACT

Open-World Machine Learning (OWML) aims to develop intelligent systems capable of recognizing known categories, rejecting unknown samples, and continually learning from novel information. Despite significant progress in open-set recognition, novelty detection, and continual learning, the field still lacks a unified theoretical foundation that can quantify uncertainty, characterize information transfer, and explain learning adaptability in dynamic, nonstationary environments. This paper presents a comprehensive review of information-theoretic approaches in open-world machine learning, emphasizing how core concepts such as entropy, mutual information, and Kullback–Leibler divergence provide a mathematical language for describing knowledge acquisition, uncertainty suppression, and risk control under open-world conditions. We synthesize recent studies into three major research axes: information-theoretic open set recognition enabling safe rejection of unknowns, information-driven novelty discovery guiding new concept formation, and information-retentive continuous learning ensuring stable long-term adaptation. Furthermore, we discuss theoretical connections between information theory and provable learning frameworks, including PAC-Bayes bounds, open space risk theory, and causal information flow, to establish a pathway toward provable and trustworthy open-world intelligence. Finally, the review identifies key open problems and future research directions, such as the quantification of information risk, the development of dynamic mutual information bounds, multimodal information fusion, and the integration of information theory with causal reasoning and world-model learning. By bridging fragmented efforts across learning theory, uncertainty quantification, and adaptive AI, this survey provides a unified perspective for building the next generation of self-adaptive, information-aware, and theoretically grounded open-world learning systems.

*Keywords* Machine learning · open-world · information theory · continual learning · open-set recognition · novelty detection

## 1 Introduction

Machine learning has achieved remarkable success in closed-world scenarios(Jordan and Mitchell, 2015; LeCun et al., 2015), where all categories, data distributions(Pan and Yang, 2009), and environmental conditions are assumed to be predefined and stationary(Gama et al., 2014). However, real-world environments are inherently open: novel classes may emerge(Geng et al., 2020), data distributions evolve(Zhang et al., 2022b), and uncertainty persists(Ovadia et al., 2019). These challenges have driven the rapid development of Open-World Machine Learning (OWML), which aims to enable models not only to recognize known categories(Han et al., 2020, 2019) but also to reject and learn from the unknown(Parmar et al., 2023). Although OWML integrates areas such as open set recognition(Scheirer et al., 2014),

---

novelty discovery(Cao et al., 2021), and continual learning(Wang et al., 2024b), the theoretical understanding of how knowledge and uncertainty evolve in open environments remains limited(Zhu et al., 2024).

Most existing approaches to open-world learning rely on heuristic confidence thresholds(Chen et al., 2021a), empirical energy models(Liu et al., 2020b), or incremental memory-based learning(Rebuffi et al., 2017). These methods perform well in practice, but lack a unified mathematical foundation that can explain why and when a system can safely recognize, reject, and adapt. The absence of such a foundation makes it difficult to quantify the limits of adaptability, control the risks of unknown exposure, or provide theoretical guarantees for the stability of continual learning. Therefore, establishing a generalizable and provable framework for OWML has become a fundamental open problem in the field.

Information theory provides a powerful lens for addressing this challenge. Its core quantities, entropy, mutual information(Shannon, 1948), and Kullback-Leibler divergence, naturally measure uncertainty(Kullback and Leibler, 1951), knowledge transfer, and information compression(Varley, 2023). By interpreting the learning process as an information flow from input to representation to decision(Weimar et al., 2025; Tan et al., 2023), information theory allows the dynamics of open-world learning to be described and analyzed mathematically consistent(Dahlke and Pacheco, 2025; Hu et al., 2024b). This perspective unifies multiple tasks, from recognizing the known and rejecting the unknown to continually learning the novel(Mondal et al., 2025; Wen et al., 2025a), under a single information-theoretic formulation.

Recent work has begun to explore this direction. In particular,some research introduced an information-theoretic formulation of open-world learning, interpreting the balance between knowledge retention and unknown suppression as an optimization of mutual information(Dziugaite and Roy, 2018). Subsequent studies extended this idea to continuous learning under open-world distribution shifts(Tan et al., 2024). These advances suggest that information theory may serve as the missing theoretical backbone for understanding and formalizing open-world intelligence.

In this review, we provide the first systematic synthesis of information-theoretic approaches to open-world machine learning. We revisit classical and modern theories of information, analyze their roles across major OWML tasks, and identify the key open problems that remain. We aim to bridge fragmented studies across open-set recognition, novelty discovery, and continual learning into a unified framework that connects empirical progress with theoretical foundations. Finally, we outline open challenges and emerging trends, including information-risk quantification, dynamic mutual information bounds, and the integration of information theory with causality and world models.

## 2    Background: Information Theory and Open-world Machine Learning

Machine learning operates fundamentally as a process of information acquisition, compression, and transmission(Kawaguchi et al., 2023; Chen et al., 2025). While traditional learning theories assume closed and stationary environments, real-world intelligence must continually adapt to uncertainty, novelty, and change(Kim et al., 2025; Wang et al., 2024b). This section provides the conceptual and theoretical background for understanding open-world machine learning through the lens of information theory. We first review the essential principles of information theory, then summarize the foundations and challenges of open-world learning, and finally discuss the conceptual intersection between the two paradigms.

### 2.1    Fundamentals of Information Theory in Learning Systems

Information theory, founded by Claude Shannon(Shannon, 1948), provides a rigorous mathematical framework for quantifying uncertainty and information flow. Its three central quantities—entropy, mutual information, and Kullback–Leibler (KL) divergence—form the backbone of how knowledge, uncertainty, and learning dynamics are measured in modern machine learning systems(Kullback and Leibler, 1951).

Entropy(Bülte et al., 2025; Wood et al., 2024; Zhang et al., 2025) represents the average level of uncertainty in a random variable, making it an essential tool for measuring unpredictability in data, model predictions, or decision boundaries. High entropy indicates uncertainty, while low entropy reflects confident or deterministic predictions. Mutual information (MI)(Dahlke and Pacheco, 2025; Tsur et al., 2023) quantifies the amount of information shared between two variables, capturing the dependency between inputs, representations, and outputs. In representation learning, mutual information measures how much of the input data's relevant structure is preserved in the learned representation. KL divergence(Flynn et al., 2023; Kuzborskij et al., 2024; Roulet et al., 2025), a measure of discrepancy between probability distributions, expresses how much one belief or model deviates from another—such as how much a posterior differs from a prior after learning.

These measures collectively describe the information flow that underlies learning. A learning system can be viewed as a communication channel $X \rightarrow Z \rightarrow Y$, where the input data X is encoded into a latent representation Z that supports the

prediction of task labels Y. The goal of learning is to maximize task-relevant information $I(Z;Y)$ while minimizing redundant or noisy information $I(Z;X)$. This trade-off forms the basis of the Information Bottleneck (IB) principle proposed by(Tishby et al., 2000). The IB principle reframes learning as a compression–relevance optimization problem: the model should compress the input while preserving information useful for predicting the output. Subsequent developments, such as the Deep Variational Information Bottleneck(Alemi et al., 2016) and Information-theoretic Generalization Bounds(Negrea et al., 2019; Neu et al., 2021), have further demonstrated that information quantities can explain generalization, robustness, and adaptation in modern deep learning systems. These foundational ideas provide the theoretical grounding for extending information-theoretic reasoning to open-world environments.

## 2.2 The Foundations of Open-world Machine Learning

Traditional machine learning operates under the closed-world assumption, where all possible classes and data distributions are known and fixed during training and inference. In contrast,OWML relaxes this assumption, confronting models with the challenges of novelty, uncertainty, and non-stationarity(Rios et al., 2024a; Liu et al., 2023; Mundt et al., 2023).

OWML envisions intelligent systems capable of recognizing known categories(Federici et al., 2020), rejecting unfamiliar inputs(Ghassemi and Fazl-Ersi, 2022; Vaze et al., 2022), and continually acquiring new knowledge(De Lange et al., 2021). This paradigm integrates three closely related research areas. The first is open-set recognition(Ge et al., 2017; Wang et al., 2023b), which focuses on identifying and safely rejecting samples that do not belong to any known category. The second is novelty discovery(Han et al., 2022; Zhang et al., 2022a), which involves grouping and interpreting previously unseen samples into coherent new classes. The third is continual learning(De Lange et al., 2021; Mirzadeh et al., 2020), which enables models to assimilate new information while retaining previously acquired knowledge, thereby avoiding catastrophic forgetting. Collectively, these three components constitute the fundamental cycle of open-world learning, in which a model must first reject unknown data, then discover new concepts, subsequently learn them, and finally integrate them into its existing knowledge base.

However, OWML introduces several fundamental theoretical challenges. First, the label space is no longer fixed(Chen and Liu, 2018), meaning the hypothesis space evolves as the environment changes. Second, data distributions become non-stationary(Farquhar and Gal, 2019; Zhou et al., 2022; Liang et al., 2025; Kurle et al., 2019), making classical risk minimization assumptions invalid. Third, uncertainty becomes multi-faceted—originating from unknown classes(Boult et al., 2019), shifting distributions(Thiagarajan et al., 2022), and partial observability(Jafarzadeh et al., 2022).

Existing approaches often rely on heuristic solutions such as confidence thresholds(Wang et al., 2021; Liu et al., 2020a), distance-based novelty scoring(Cheng and Vasconcelos, 2021), or memory replay strategies(Wang et al., 2022; Raghavan et al., 2019). While effective empirically, these methods lack formal guarantees regarding stability(Sun et al., 2025), adaptability(Bonjour, 2024; Tang et al., 2025), or safety(Mohseni et al., 2022). The absence of a unified theoretical framework—one that can quantify how information about the world evolves—remains the central obstacle to developing robust open-world intelligence(Xue, 2024).

## 2.3 Bridging Information Theory and Open-world Learning

Information theory offers a natural mathematical foundation for formalizing open-world learning(Tishby and Zaslavsky, 2015; Kejriwal et al., 2024). It provides a set of distribution-agnostic tools(Kejriwal et al., 2024; Chen et al., 2021b) to quantify uncertainty(Fakour et al., 2024), learning progress(Li et al., 2024), and adaptation cost(Nguyen et al., 2021b), thus addressing many of the limitations of existing OWML frameworks(Xu et al., 2023; Sun et al., 2021; Wutschitz et al., 2023). From an information-theoretic viewpoint, learning in open environments can be conceptualized as a process of information flow under uncertainty(Rios et al., 2024b; Gawlikowski et al., 2023), where models must decide how much information to extract, retain, and suppress(Achille and Soatto, 2018).

Entropy provides a direct measure of uncertainty in the presence of unknown categories(Wang et al., 2023a). Mutual information quantifies how effectively a model preserves useful knowledge about known tasks(Westphal et al., 2024), while KL divergence captures the adaptation cost when encountering novel distributions(Nguyen et al., 2021a). This framework allows the behaviors of recognition, rejection, and adaptation to be expressed in unified informational terms(Sun et al., 2023).

Recent studies have begun to explore this intersection explicitly. (Li et al., 2025a)formulated open-world learning as an information-optimization problem, balancing knowledge retention and unknown suppression through mutual information.(Li et al., 2025b)extended this to continual learning, modeling open-world adaptation as a dynamic flow of information across time. These approaches demonstrate that information theory is not merely a descriptive tool but a foundational theory for analyzing, quantifying, and even proving properties of open-world learning systems.

Ultimately, the convergence of information theory and OWML reframes the open-world learning challenge from an empirical problem into a quantifiable and provable information process. It enables researchers to reason about learning not only in terms of accuracy or error but also in terms of information gain, uncertainty reduction, and knowledge evolution. This synthesis lays the groundwork for the next section, which introduces formal information-theoretic frameworks for open-world machine learning.

## 3 Information-Theoretic Framework for Open-world Machine Learning

OWML requires learning systems to operate under uncertainty, evolving environments, and continually expanding label spaces(Zhang, 2025; Du et al., 2023a). Traditional empirical risk minimization frameworks assume a fixed distribution and a closed hypothesis space, making them insufficient for reasoning about the dynamics of open environments(Boult et al., 2019; Xie et al., 2024). Information theory offers an elegant and mathematically grounded alternative: it interprets learning as an information flow process(Lakkaraju et al., 2017), where knowledge is acquired, compressed, and transmitted through representations that balance relevance, uncertainty, and adaptability(Cruz et al., 2025). This section introduces the information-theoretic framework that formalizes OWML as a quantifiable process of information exchange.

### 3.1 The Information Flow Perspective of Open-world Learning

In open-world environments, a learning system continuously interacts with uncertain information sources, transforming raw sensory data into structured knowledge. From an information-theoretic perspective, this process can be represented as an information channel connecting three conceptual spaces: the input space X, the representation space Z, and the output space Y. The input data contain both task-relevant and irrelevant information, and the learning objective is to filter and compress this information such that Z preserves only what is useful for predicting Y.

This perspective directly extends the Information Bottleneck (IB) principle proposed by(Tishby et al., 2000), which formulates the learning problem as an optimization of mutual information:

$$\min I(Z; X) - \beta I(Z; Y), \tag{1}$$

where $I(\cdot; \cdot)$ denotes mutual information and $\beta$ controls the balance between input compression and task relevance.

The IB framework interprets learning as an optimal encoding problem, in which a compact representation Z should retain only the information about X that is predictive of Y. This view provides a natural foundation for modeling the learning process as an information flow, bridging the gap between representation learning and decision-making.

In the open-world context, however, this classical formulation must be extended to account for information associated with unknown or novel categories(Bendale and Boult, 2016). Unlike closed-world settings, where all possible outputs are predefined, open-world systems must regulate the information they extract from uncertain or previously unseen inputs(Dhamija et al., 2018; Zhao et al., 2023). This requirement motivates the development of an extended information-theoretic objective for OWML.

### 3.2 The Information-Theoretic Objective for OWML

(Zhou et al., 2024) proposed an information-theoretic formulation that generalizes the Information Bottleneck to open-world learning scenarios. The central insight is that open-world learning involves managing three types of information: (1) compressing redundant information from the input, (2) retaining task-relevant information for known categories, and (3) suppressing misleading information from unknown or uncertain samples. This trade-off is expressed as an optimization problem over mutual information quantities:

$$\min I(Z; X) - \beta I(Z; Y_{\text{known}}) + \gamma I(Z; Y_{\text{unknown}}), \tag{2}$$

where $I(Z; X)$ represents the compression of the input data into a latent representation $Z$; $I(Z; Y_{\text{known}})$ quantifies the task-relevant information preserved for known classes; and $I(Z; Y_{\text{unknown}})$ measures the influence of unknown data on the learned representation. The parameters $\beta$ and $\gamma$ control the trade-off between knowledge retention and novelty suppression.

This formulation unifies the three fundamental aspects of open-world learning—compression, retention, and rejection—into a single mathematical objective. It extends the classical IB framework by incorporating a mechanism for uncertainty regulation: rather than assuming a fixed output distribution, it explicitly models the influence of unknown categories as an information source to be minimized. This objective thus provides a foundation for understanding safe adaptation, where a model can learn effectively from known data while avoiding overconfidence on unknown inputs.

Building upon this formulation, (Li et al., 2025b) extended the objective to continual learning under open-world distribution shifts. They proposed that knowledge retention over time can be represented as the preservation of mutual information between successive representation states:

$$\max_t I(Z_t; Z_{t-1}), \tag{3}$$

where $Z_t$ denotes the latent representation at time step $t$. Maintaining high mutual information between $Z_t$ and $Z_{t-1}$ ensures that new learning does not catastrophically overwrite previously acquired knowledge.

This dynamic regularization term introduces a temporal dimension to the information-theoretic framework, allowing OWML to capture both stability and adaptability across evolving environments.

Together, these formulations transform open-world learning from an empirical problem into a quantifiable process of information regulation. Learning becomes not merely a matter of minimizing error but of managing the balance between information compression (generalization), information preservation (memory), and information suppression (safety).

### 3.3 Extensions and Theoretical Connections

The information-theoretic formulation of OWML is closely connected to several established theoretical frameworks, each of which provides complementary insights.

Relation to the Information Bottleneck (IB): The OWML objective generalizes the IB principle by adding a term for novelty suppression(Xu et al., 2019). While the IB focuses solely on balancing compression and relevance, the OWML formulation introduces an additional constraint to control uncertainty caused by unknown data. This extension allows the model to operate safely beyond the closed-world assumption.

Relation to PAC-Bayes and Generalization Theory: Mutual information plays a central role in bounding generalization error. (Haghifam et al., 2020) showed that the expected generalization gap can be upper-bounded by a function of the mutual information between training data and model parameters. Within OWML, limiting the information a model extracts from unknown inputs ensures tighter, provable bounds on generalization under distribution shifts.

Relation to Open-space Risk Theory: (Scheirer et al., 2012) introduced open-space risk to measure the danger of making confident predictions in unsupported regions of feature space. In information-theoretic terms, these regions correspond to high-entropy, low-mutual-information areas where models should minimize exposure. The OWML framework thus reinterprets open-space risk as an information exposure problem.

Relation to Continual and Lifelong Learning: In continual learning, the mutual information preservation objective $I(Z_t; Z_{t-1})$ serves as a theoretical expression of stability–plasticity balance(Chen et al., 2023). Retaining mutual information across tasks helps maintain previously learned representations, while the addition of controlled new information enables adaptive knowledge expansion.

These theoretical connections collectively position information theory as the backbone of open-world learning. It provides a unifying language for describing diverse learning phenomena—recognition, rejection, and adaptation—in terms of quantifiable information dynamics.

## 4 Applications of Information Theory in OWML Subtasks

Information theory not only provides a unifying mathematical foundation for understanding open-world learning but also offers practical tools for analyzing and improving its key subtasks(Liu et al., 2025b). This chapter examines how entropy, mutual information, and divergence-based measures have been applied in three major areas of OWML: open-set recognition, novelty discovery, and continual learning. Each of these subtasks addresses a different phase of the open-world cycle—recognizing the known, discovering the novel, and integrating the new—yet all can be expressed within a shared information-theoretic framework.

### 4.1 Information Theory in Open-set Recognition (OSR)

Open-set recognition (OSR) is the first and most fundamental task in open-world learning(Cao et al., 2025; Nawaz, 2025). Its goal is to enable models to correctly classify known classes while rejecting samples that belong to unseen or unknown categories(Xing et al., 2025; Moazzami et al., 2025). The key challenge is balancing discriminative capability with uncertainty awareness—how to separate known information from informational noise introduced by the unknown.

From an information-theoretic standpoint, OSR can be viewed as minimizing the entropy of known predictions while maximizing the uncertainty for unknown samples(Liu et al., 2025a). Entropy serves as a natural metric of confidence:

low entropy indicates certainty about known classes, whereas high entropy signals uncertainty that can be used to reject unfamiliar inputs(Garg et al., 2022). Formally, OSR can be formulated as an entropy-regularized optimization problem(Hu et al., 2024a):

$$\min \mathbb{E}_{x \in \mathcal{D}_{\text{known}}}[H(Y|X)] \; + \; \lambda \, \mathbb{E}_{x \in \mathcal{D}_{\text{unknown}}}[-H(Y|X)], \qquad (4)$$

where $H(Y|X)$ is the conditional entropy of the output given the input, and $\lambda$ balances classification confidence and rejection safety. Beyond entropy, mutual information has been used to improve rejection decisions. By maximizing $I(Z; Y_{\text{known}})$ while minimizing $I(Z; Y_{\text{unknown}})$, the model learns feature representations that retain discriminative information for known classes but discard information that correlates with unknowns. This information separation principle underpins recent OSR models such as OpenMax(Bendale and Boult, 2016), ARPL(Chen et al., 2021a), and M2IOSR(Sun et al., 2021)). Information-theoretic OSR therefore provides both an intuitive and quantitative measure of model confidence and generalization beyond the closed set.

## 4.2 Information Theory in Novelty Discovery

Once unknown samples are identified, the next step in open-world learning is novelty discovery—the process of grouping, characterizing, and interpreting unknown data to form new classes or concepts(Jin et al., 2024). This task lies at the intersection of unsupervised learning and knowledge expansion(Zhou, 2022), where the model must decide what constitutes "novel" information.

In information-theoretic terms, novelty discovery can be formalized as a maximization of information gain, i.e., the difference between the entropy before and after observing new data(Du et al., 2023b; Lidayan et al., 2025). The objective is to identify clusters or representations that maximize the reduction in uncertainty about the environment(Jafarzadeh et al., 2020):

$$\max_{\theta} \; \mathbb{E}_{x \in \mathcal{D}_{\text{unknown}}} \left[ H(P_{\text{prior}}(Y)) - H(P_{\theta}(Y|X)) \right], \qquad (5)$$

where $H(P_{\text{prior}}(Y))$ denotes prior uncertainty and $H(P_{\theta}(Y|X))$ denotes posterior uncertainty after modeling with parameters $\theta$. The term inside the brackets represents information gain—a measure of how much new information about the world has been discovered.

This information gain formulation aligns with the principles of active learning and Bayesian exploration(Houlsby et al., 2011), where models seek data that maximally reduce uncertainty(Sekar et al., 2020). In practice, novelty discovery methods often combine mutual information with clustering algorithms (e.g., InfoNCE(van den Oord et al., 2019), Deep InfoMax(Hjelm et al., 2019)) to identify meaningful latent structures in the unknown data space. Recent research (Abbasi et al., 2024; Wang et al., 2024a) extends this idea by coupling mutual information maximization with representation disentanglement, ensuring that novel knowledge is structured and separable from known features.

## 4.3 Information Theory in Continual Learning

Continual learning (CL) addresses the problem of learning from a sequence of tasks without catastrophic forgetting(Wang et al., 2024c). In open-world environments, this challenge becomes even more critical as new tasks and distributions emerge continuously(Li et al., 2025a). Information theory provides a principled approach to analyze and mitigate forgetting by quantifying how much information about past knowledge is retained as new learning occurs(Song et al., 2023).

(Li et al., 2025c) formalized this idea by defining mutual information preservation across time. Their objective encourages the model to maintain shared information between consecutive latent representations, ensuring temporal stability in formulation (3), this formulation provides a compact theoretical explanation for methods such as EWC(Kirkpatrick et al., 2017), LwF(Li and Hoiem, 2017), and replay-based continual learning(Wen et al., 2025b), which implicitly preserve information through regularization or memory mechanisms.

Additionally, the trade-off between plasticity (learning new information) and stability (retaining old knowledge) can be expressed as a dual optimization problem over mutual information:

$$\max \; I(Z_t; Y_t) \; - \; \lambda \, D_{\text{KL}}(P(Z_t) \parallel P(Z_{t-1})), \qquad (6)$$

where the first term promotes adaptation to the current task, and the KL-divergence term penalizes excessive deviation from previous representations. This equation directly connects continual learning to the information-theoretic framework of OWML, showing that maintaining controlled information divergence is essential for sustainable knowledge evolution.

Through this perspective, continual learning becomes an information balancing process—preserving relevant past information while incorporating new, task-specific knowledge.

### 4.4 Summary of Information-Theoretic Applications

The integration of information-theoretic principles into open-world learning tasks provides a consistent mathematical foundation for reasoning about uncertainty, discovery, and retention. Information theory thus acts as a bridge across all components of OWML: entropy quantifies uncertainty, mutual information measures knowledge retention, and divergence captures adaptation cost. Together, these measures transform open-world learning into a unified, quantifiable process of information flow and evolution.

## 5 Toward Provable Open-world Learning:Mathematical Foundations

OWML aims to build systems that can operate safely and adaptively under uncertainty. To transition from empirical success to provable intelligence, OWML requires a rigorous mathematical foundation. This chapter explores how information theory, combined with generalization theory, causal inference, and statistical learning principles, can be used to construct provable learning guarantees under open-world conditions. We begin by revisiting classical closed-world theories, extend them to open distributions, formalize information-theoretic risk and bounds, and discuss future directions toward provable open-world intelligence.

### 5.1 Revisiting Closed-world Provability

Traditional machine learning operates under a closed-world assumption, where data are drawn i.i.d. from a fixed distribution $P(X, Y)$, and all classes are known during both training and inference(Zhang, 2025). In this setting, the generalization error—the difference between empirical risk on the training set and true risk on unseen samples—can be bounded using statistical learning theory(Yu et al., 2021).

A classical guarantee(Gross et al., 2025) can be expressed as:

$$\mathcal{R}(f) \leq \hat{\mathcal{R}}(f) + \mathcal{C}(f, n, \delta), \tag{7}$$

where $\mathcal{R}(f)$ denotes the expected risk, $\hat{\mathcal{R}}(f)$ is the empirical risk, and $\mathcal{C}(f, n, \delta)$ is a complexity term depending on model capacity, sample size $n$, and confidence parameter $\delta$.

While this formulation provides solid theoretical grounding for closed-world learning, it collapses under distributional openness—when new classes or data distributions emerge(Chuang et al., 2020). The core limitation is that traditional generalization bounds assume a fixed hypothesis space and a stationary data-generating process, both of which are violated in open-world settings(Feldman and Vondrak, 2019).

### 5.2 Extending Generalization to Open Distributions

In OWML, the learner must adapt to new data domains and evolving label spaces, where the training distribution $P_{\text{train}}(X, Y)$ differs from the test distribution $P_{\text{test}}(X, Y)$(Stojanov et al., 2021). The challenge is thus to establish *transferable* or *adaptive* generalization guarantees.

A common approach is to introduce distributional divergence measures that quantify the difference between source and target distributions(Courty et al., 2016). Let $D_{\text{KL}}(P_{\text{test}} \parallel P_{\text{train}})$ denote the Kullback–Leibler divergence. The expected risk under open-world conditions can be upper-bounded as(Duchi and Namkoong, 2019):

$$\mathcal{R}_{\text{test}}(f) \leq \mathcal{R}_{\text{train}}(f) + \alpha \, D_{\text{KL}}(P_{\text{test}} \parallel P_{\text{train}}), \tag{8}$$

where $\alpha$ is a scaling constant that captures the model's sensitivity to distributional shifts. This inequality reflects a fundamental insight: the cost of openness can be quantified as the information divergence between what has been learned and what is newly encountered.

In practice, this divergence term can be estimated or regularized through mutual information constraints, effectively linking distributional generalization with information-theoretic control.

### 5.3 Information-theoretic Risk and Provable Bounds

Information theory provides a natural way to express provable guarantees through mutual information-based risk bounds(Xu and Raginsky, 2017). Building on the PAC-Bayes framework, (Harutyunyan et al., 2021) demonstrated that the generalization error can be bounded in terms of the mutual information between the model parameters $W$ and the training data $D$:

$$\mathbb{E}\left[\mathcal{R}(f_W) - \hat{\mathcal{R}}(f_W)\right] \leq \sqrt{\frac{2I(W; D)}{n}}. \tag{9}$$

This inequality provides a powerful interpretation: models that encode less information about the specific training data generalize better. In open-world learning, this concept extends to controlling the information exchanged between past knowledge, current tasks, and novel environments.

Combining this with the open-world formulation from Chapter 3, we can express the Information-Theoretic Open-world Bound as:

$$\mathcal{R}_{\text{open}}(f) \leq \hat{\mathcal{R}}_{\text{known}}(f) + \sqrt{\frac{2I(Z; X_{\text{known}})}{n}} + \gamma D_{\text{KL}}(P_{\text{unknown}} \parallel P_{\text{known}}), \tag{10}$$

where the first term is the empirical risk on known data, the second term quantifies uncertainty and generalization ability through mutual information, and the third term penalizes deviation caused by unknown distributions. This provides a mathematically grounded definition of open-world risk—the total uncertainty of learning under novelty and distributional drift.

## 5.4 Toward Provable Open-world Intelligence

Developing provable open-world learning requires more than empirical adaptation; it demands formal principles that link uncertainty, adaptability, and stability(Qu et al., 2025). The integration of information theory into statistical learning opens several promising research avenues:

Information Risk Theory: Establishing information risk minimization (IRM) as a generalization of empirical risk minimization (ERM). In IRM, the objective is to minimize the expected information exposure to unknown factors, rather than just empirical error.

Dynamic PAC-Bayes Bounds: Extending PAC-Bayes inequalities to non-stationary and temporally evolving tasks, where the mutual information term becomes time-dependent $I(W_t; D_{t-1})$. This could provide guarantees for continual learning under open-world dynamics.

Causal Information Flow: Incorporating causal reasoning into information-theoretic analysis to distinguish between informative novelty (causally relevant) and spurious novelty (noise). This would enable provable causal generalization in open-world systems.

Information Stability and Safety: Defining safety margins in terms of information bounds—for example, ensuring that $D_{\text{KL}}(P_{\text{unknown}} \parallel P_{\text{known}}) < \epsilon$ for safe model deployment. Such criteria could be used to certify when an open-world model is theoretically safe to operate.

Bridging Theory and Implementation: Translating these bounds into practical algorithms via information regularization, adaptive loss functions, and mutual information estimators. This would allow theory-driven model design where stability, adaptability, and uncertainty are mathematically coupled.

Collectively, these directions move open-world learning toward provable intelligence—systems that not only perform well empirically but also offer verifiable guarantees about their behavior under novel and uncertain conditions. Such a theory would mark the transition from heuristic adaptation to quantified, principled open-world reasoning.

# 6 Comparative Theoretical Analysis

Classical theories such as Statistical Learning Theory(Bartlett et al., 2020), Bayesian Learning(Wilson and Izmailov, 2020), Energy-based Modeling(Song and Kingma, 2021), and Causal Learning(Kaddour et al., 2022) each capture different aspects of intelligence—generalization, uncertainty, representation, and interpretability. However, they are all built upon assumptions of closure, stability, or complete knowledge.This chapter presents a comparative theoretical analysis showing how information theory provides a unified foundation that extends, connects, and generalizes these paradigms under open-world conditions.

## 6.1 Statistical Learning Theory

Statistical Learning Theory (SLT) provides the mathematical foundation for generalization in closed-world settings. Its strength lies in its provable nature: it defines clear relationships between empirical performance(Boucheron et al., 2005), model capacity, and expected risk(Feldman and Vondrak, 2019). This framework assumes that all samples are independent and identically distributed and that the underlying data distribution remains fixed. While SLT has been extraordinarily influential in defining what it means for a model to generalize, it collapses under openness. When new classes appear(Geng et al., 2020) or data distributions shift(Machlanski et al., 2025), the assumption of a stationary environment no longer holds, and classical risk bounds become invalid(Wiles et al., 2021). From an

information-theoretic perspective, SLT can be informationally extended. Information theory transforms risk into an information exchange process, where the uncertainty between model representations and data distributions is explicitly quantified. This shift allows learning guarantees to remain meaningful even when the world is not fixed, providing a bridge between closed-world provability and open-world adaptability. In summary, Statistical Learning Theory remains the formal backbone of provability, but its extension through information measures such as entropy and mutual information enables it to survive in dynamic, open settings.

## 6.2 Bayesian Learning

Bayesian Learning introduces a probabilistic framework for managing uncertainty(Wilson and Izmailov, 2020). Its advantage lies in its strong ability to represent epistemic uncertainty—the uncertainty about model parameters or hypotheses. By maintaining a distribution over model parameters, Bayesian approaches allow reasoning under incomplete information and offer principled probabilistic inference. However, Bayesian learning relies on predefined priors and fixed model structures(Li et al., 2021). In open-world conditions, priors may become outdated, incomplete, or even misleading as the environment evolves(Adel, 2025). This makes traditional Bayesian inference brittle when confronted with novel categories or shifting semantics(Galashov et al., 2024). Information theory complements Bayesian learning by embedding uncertainty management directly into information flow. Instead of assuming fixed priors, information-theoretic learning models uncertainty as a dynamic process of encoding, transmission, and transformation of information. Mutual information becomes the connecting principle—linking the belief-based uncertainty of Bayesian inference to the relevance-based uncertainty of open-world adaptation.

In summary, Bayesian learning interprets uncertainty as belief; information theory interprets it as information relevance. Together, they form a foundation for reasoning under both known and unknown uncertainties.

## 6.3 Energy-based Models

Energy-based Models (EBMs) have achieved remarkable empirical success, particularly in deep learning(Oliva et al., 2025). They define a scalar "energy" function that measures how compatible a configuration of variables is, allowing for flexible modeling of complex dependencies(Xu et al., 2024). EBMs are powerful because they can capture high-dimensional relationships without requiring explicit probability normalization, making them practical and expressive(Song and Kingma, 2021). However, EBMs are largely empirical and lack a comprehensive theoretical framework explaining their learning dynamics, stability, or generalization. Their energy function can be interpreted as an implicit potential landscape, but its theoretical meaning often remains opaque. From the viewpoint of information theory, energy can be reframed as information potential—a representation of how much information is stored or compressed in a given state. Under this interpretation, minimizing energy becomes equivalent to optimizing the information flow within the system, aligning EBMs with information-theoretic principles of efficiency and stability. This provides a theoretical bridge: information theory endows energy-based learning with measurable interpretability, connecting it to the entropy and mutual information structures underlying open-world intelligence.

In summary, EBMs offer practical performance, while information theory offers them a missing theoretical backbone—turning energy minimization into an explicit form of information optimization.

## 6.4 Causal Learning

Causal Learning seeks to uncover and utilize the cause–effect relationships underlying observed data. Its advantage lies in its interpretability and its ability to generalize across interventions and domains(Jiao et al., 2024). Causal inference assumes that the structural relationships between variables remain invariant, allowing predictions to hold even under changes in external conditions. However, causal learning depends on having access to stable causal structures. When new causes emerge or mechanisms evolve—as in open-world environments—this assumption no longer holds. Traditional causal frameworks struggle to adapt to causal novelty or explain newly arising dependencies. Information theory extends causal reasoning by introducing the concept of information–causal flow. This framework treats causation as an information transmission process, where the strength and direction of influence are quantified through changes in information content. By viewing causal links as dynamic information channels, information theory enables models to detect when existing mechanisms break and new ones form. This creates a pathway toward adaptive causal inference, where causality and novelty are jointly represented in the information space. In summary, causal learning ensures interpretability under stable mechanisms; information theory generalizes it to handle causal openness, forming the basis for adaptive, self-updating causal systems.

### 6.5 Comparative Theoretical Summary

The comparison across the four paradigms reveals that information theory acts as a meta-theoretical bridge—linking provability, uncertainty, representation, and causality under a unified formalism of information flow. Where classical theories assume closure and stability, information theory reinterprets them as special cases of information transformation within evolving environments. Statistical learning offers provability, Bayesian learning offers uncertainty reasoning,

Table 1: Comparative analysis of major theoretical frameworks and their relationship with Information Theory in OWML

| Theoretical Framework | Advantages | Limitations | Relation to Information Theory |
|---|---|---|---|
| **Statistical Learning Theory** | Theoretically complete and provable | Relies on closed-world distribution assumptions, unable to handle open-world dynamics | Can be extended through information measures to provide provable bounds under open distributions |
| **Bayesian Learning** | Strong ability to handle uncertainty | Sensitive to prior assumptions and vulnerable to novel classes | Can be integrated with mutual information to model dynamic uncertainty |
| **Energy-based Models** | Excellent empirical performance | Lack rigorous theoretical interpretation | Can be viewed as Information Potential Models, where energy reflects informational capacity |
| **Causal Learning** | Strong interpretability and transferability | Depends on stable structures, difficult to adapt to new causal mechanisms | Can be extended to Info-Causal Flow, integrating causality with information transmission |

energy-based models offer practical expressiveness, and causal learning offers interpretability. Yet all rely on the assumption of a closed and stationary world. Information theory unifies these paradigms under a single conceptual lens: learning as a process of dynamic information regulation. It redefines risk as information transfer, uncertainty as information entropy, energy as information potential, and causality as information flow. Through this perspective, OWML transforms machine learning from a static optimization paradigm into a living system of information exchange, capable of reasoning, adapting, and proving its behavior in open and uncertain worlds.

## 7 Open Problems and Future Directions

OWML remains far from a mature theoretical discipline. Information theory has provided a powerful foundation for reasoning about uncertainty, adaptability, and provability, but many scientific questions remain unanswered. This chapter outlines key open problems and proposes several promising research directions that could define the next stage of development for information-theoretic open-world learning. Each direction highlights a core challenge, its underlying scientific question, and a potential conceptual goal for future exploration.

### 7.1 Open Problems Overview

The following table summarizes the principal future directions and open theoretical challenges in information-theoretic OWML.

### 7.2 Directional Analysis

#### 7.2.1 Information Risk Theory

One of the most fundamental open problems is how to define and measure information risk in open-world learning. Traditional risk is expressed in terms of error or loss under fixed distributions, but in open environments, models face not only predictive errors but also information exposure—how much of the unknown world they fail to represent or misinterpret. A future Open-space Information Risk framework would quantify this exposure as a balance between known, unknown, and novel informational components, providing a new axis of theoretical provability for open systems.

Table 2: Future research directions and open problems in information-theoretic Open-world Machine Learning

| Direction | Scientific Question | Potential Research Objective |
|---|---|---|
| Information Risk Theory | How to define open information risk under non-stationary and uncertain conditions | Establish an Open-space Information Risk framework for quantifying information exposure in open environments |
| Dynamic Information Flow | How do information generalization bounds evolve over time in continual and incremental learning scenarios | Develop Temporal Mutual Information formulations for time-dependent information retention and adaptation |
| Information and Causality | How can information flow explain or predict causal flow in evolving systems | Build an Info-Causal Fusion model that unifies information theory with dynamic causal inference |
| Multimodal Information Theory | How can mutual information across different modalities be measured and combined effectively | Define Cross-modal Information Bounds for aligning heterogeneous sensory and symbolic representations |
| Self-adaptive Agents | Can an intelligent agent quantify and regulate its own information boundary | Design Self-Information-Aware Agents capable of evaluating and optimizing their internal informational states |

### 7.2.2 Dynamic Information Flow

Most existing information-theoretic bounds, such as mutual information constraints, are static and assume a single learning stage. In open-world scenarios, however, learning unfolds over time: information is accumulated, transformed, and sometimes forgotten. Understanding how generalization and retention evolve dynamically requires formulating Temporal Mutual Information—a time-dependent measure capturing how information propagates and stabilizes across sequential tasks. This direction could lay the groundwork for a dynamic, provable continual learning theory.

### 7.2.3 Information and Causality

A long-term scientific challenge lies in merging causal reasoning with information-theoretic flow. While causality explains why events occur, information theory explains how signals about those events are transmitted and processed. The question of whether information flow can serve as a sufficient descriptor for causal flow remains open. A unified Info-Causal Fusion framework could describe both mechanisms and communication processes in adaptive systems, allowing models to infer not only correlations but evolving causal structures in open environments.

### 7.2.4 Multimodal Information Theory

As real-world systems increasingly rely on multimodal data—text, vision, speech, sensors—the challenge is to define mutual information consistently across modalities that differ in structure and scale. Information theory provides a natural bridge, yet cross-modal dependencies are often nontrivial and asymmetric. Developing Cross-modal Information Bounds could help quantify how much information from one modality contributes to another, enabling robust multimodal integration, alignment, and transfer in open settings.

### 7.2.5 Self-adaptive Information-aware Agents

Perhaps the most ambitious direction is creating agents capable of self-quantifying their information state. Such an agent would be aware of its informational limits—how much it knows, how much it ignores, and when it should seek new knowledge. This leads to the concept of a Self-Information-Aware Agent, which can measure and regulate its own informational entropy, dynamically adjust its learning capacity, and maintain safety margins based on information divergence. This self-referential information control may form the foundation for autonomous, self-regulating open-world intelligence.

### 7.3 Theoretical Outlook

These research directions collectively point toward a unifying theory of Provable Information Dynamics—a mathematical and conceptual framework capable of describing how information evolves, interacts, and stabilizes in open systems. By bridging the gap between uncertainty, adaptation, and causality, information theory could provide the first truly general foundation for intelligent behavior under non-stationary conditions.

Future progress will depend on three key developments:

1. Formalization: Establishing precise definitions of open information quantities (entropy, divergence, risk).

2. Estimation: Designing scalable estimators for temporal and multimodal mutual information.

3. Integration: Embedding these measures into learning systems that can reason, explain, and adapt beyond static environments.

Together, these advances could transform open-world learning from an empirical engineering paradigm into a scientifically grounded theory of information-driven intelligence.

## 8 Conclusion

Open-world Machine Learning (OWML) represents a new frontier in artificial intelligence—one where systems must not only learn from limited, incomplete, or shifting data but must also recognize when the world itself changes. Traditional learning paradigms—rooted in the closed-world assumption—struggle to handle novelty, uncertainty, and non-stationarity. Information theory provides the missing bridge: a unified theoretical language for quantifying, explaining, and ultimately mastering learning under openness.

Throughout this review, we have shown that information theory connects the fundamental components of open-world intelligence: risk, generalization, rejection, learning, and cognition. It redefines risk as information exposure, generalization as information preservation, rejection as uncertainty maximization, and continual learning as information flow regulation. Within this view, cognition itself can be understood as an emergent property of adaptive information dynamics—an ongoing process of compressing, retaining, and transforming information in response to environmental change.

By embedding open-world learning within an information-theoretic framework, we move from heuristic adaptation to mathematically grounded reasoning. This shift enables the development of provable learning bounds, dynamic information control, and safe adaptive systems capable of operating in non-stationary, uncertain environments. Information theory thus acts not only as a descriptive tool but as a generative principle—a foundation for designing learning systems that are both accountable and adaptive.

Looking forward, the evolution of open-world intelligence will depend on advancing this information-driven paradigm. Future learning systems must understand their own informational limits, reason about unknowns, and dynamically regulate their internal knowledge flow. They must move beyond static optimization and embrace continuous adaptation—balancing what is known, what is learnable, and what remains unknowable. In this way, information theory will guide the transformation from learning systems that merely react to their data toward intelligent agents that reason about information itself.

In summary, information theory offers not just mathematical precision but philosophical coherence: it unifies the diverse mechanisms of open-world learning into a single, principled narrative—a narrative in which intelligence is nothing more, and nothing less, than the art of managing information in an ever-changing world.

## References

Abbasi, R., Rohban, M. H., and Baghshah, M. S. (2024). Deciphering the role of representation disentanglement: Investigating compositional generalization in clip models. In *European Conference on Computer Vision*, pages 35–50. Springer.

Achille, A. and Soatto, S. (2018). Information dropout: Learning optimal representations through noisy computation. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2897–2905.

Adel, T. (2025). The bayesian approach to continual learning: An overview. *arXiv preprint arXiv:2507.08922*.

Alemi, A. A., Fischer, I., Dillon, J. V., and Murphy, K. (2016). Deep variational information bottleneck. *arXiv preprint arXiv:1612.00410*.

Bartlett, P. L., Long, P. M., Lugosi, G., and Tsigler, A. (2020). Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117:30063–30070.

Bendale, A. and Boult, T. E. (2016). Towards open set deep networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1563–1572.

Bonjour, T. (2024). *Towards Novelty-Resilient AI: Learning in the Open World*. PhD thesis, Purdue University.

Boucheron, S., Bousquet, O., and Lugosi, G. (2005). Theory of classification: A survey of some recent advances. *ESAIM: probability and statistics*, 9:323–375.

Boult, T. E., Cruz, S., Dhamija, A. R., Gunther, M., Henrydoss, J., and Scheirer, W. J. (2019). Learning and the unknown: Surveying steps toward open world recognition. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 9801–9807.

Bülte, C., Sale, Y., Löhr, T., Hofman, P., Kutyniok, G., and Hüllermeier, E. (2025). An axiomatic assessment of entropy-and variance-based uncertainty quantification in regression. *arXiv preprint arXiv:2504.18433*.

Cao, K., Brbic, M., and Leskovec, J. (2021). Open-world semi-supervised learning. *arXiv preprint arXiv:2102.03526*.

Cao, W., Yao, X., Xu, Z., Liu, Y., Pan, Y., and Ming, Z. (2025). A survey of zero-shot object detection. *Big Data Mining and Analytics*, 8:726–750.

Chen, G., Peng, P., Wang, X., and Tian, Y. (2021a). Adversarial reciprocal points learning for open set recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):8065–8081.

Chen, J., Fang, Y., Khisti, A., Özgür, A., and Shlezinger, N. (2025). Information compression in the ai era: Recent advances and future challenges. *IEEE Journal on Selected Areas in Communications*.

Chen, Q., Shui, C., Han, L., and Marchand, M. (2023). On the stability-plasticity dilemma in continual meta-learning: Theory and algorithm. *Advances in Neural Information Processing Systems*, 36:27414–27468.

Chen, Q., Shui, C., and Marchand, M. (2021b). Generalization bounds for meta-learning: An information-theoretic analysis. *Advances in Neural Information Processing Systems*, 34:25878–25890.

Chen, Z. and Liu, B. (2018). *Lifelong machine learning*. Morgan & Claypool Publishers.

Cheng, J. and Vasconcelos, N. (2021). Learning deep classifiers consistent with fine-grained novelty detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1664–1673.

Chuang, C.-Y., Torralba, A., and Jegelka, S. (2020). Estimating generalization under distribution shifts via domain-invariant representations.

Courty, N., Flamary, R., Tuia, D., and Rakotomamonjy, A. (2016). Optimal transport for domain adaptation.

Cruz, S., Doctor, K., Funk, C., and Scheirer, W. (2025). Open issues in open world learning. *AI Magazine*, 46(2):e70001.

Dahlke, C. and Pacheco, J. (2025). Flow-based variational mutual information: Fast and flexible approximations. In *The Thirteenth International Conference on Learning Representations*.

De Lange, M., Aljundi, R., Masana, M., Parisot, S., Jia, X., Leonardis, A., Slabaugh, G., and Tuytelaars, T. (2021). A continual learning survey: Defying forgetting in classification tasks. *IEEE transactions on pattern analysis and machine intelligence*, 44(7):3366–3385.

Dhamija, A. R., Günther, M., and Boult, T. (2018). Reducing network agnostophobia. *Advances in Neural Information Processing Systems*, 31.

Du, S., Fang, Z., Lan, S., Tan, Y., Günther, M., Wang, S., and Guo, W. (2023a). Bridging trustworthiness and open-world learning: An exploratory neural approach for enhancing interpretability, generalization, and robustness. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 8719–8729.

Du, Y., Kosoy, E., Dayan, A., Rufova, M., Abbeel, P., and Gopnik, A. (2023b). What can ai learn from human exploration? intrinsically-motivated humans and agents in open-world exploration. In *Neurips 2023 workshop: Information-theoretic principles in cognitive systems*.

Duchi, J. and Namkoong, H. (2019). Variance-based regularization with convex objectives. *Journal of Machine Learning Research*, 20(68):1–55.

Dziugaite, G. K. and Roy, D. M. (2018). Data-dependent pac-bayes priors via differential privacy. *Advances in neural information processing systems*, 31.

Fakour, F., Mosleh, A., and Ramezani, R. (2024). A structured review of literature on uncertainty in machine learning & deep learning. *arXiv preprint arXiv:2406.00332*.

Farquhar, S. and Gal, Y. (2019). A unifying bayesian view of continual learning. *arXiv preprint arXiv:1902.06494*.

Federici, M., Dutta, A., Forré, P., Kushman, N., and Akata, Z. (2020). Learning robust representations via multi-view information bottleneck. *arXiv preprint arXiv:2002.07017*.

Feldman, V. and Vondrak, J. (2019). Generalization bounds for uniformly stable algorithms.

Flynn, H., Reeb, D., Kandemir, M., and Peters, J. (2023). Pac-bayes bounds for bandit problems: A survey and experimental comparison. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(12):15308–15327.

Galashov, A., Titsias, M., György, A., Lyle, C., Pascanu, R., Teh, Y. W., and Sahani, M. (2024). Non-stationary learning of neural networks with automatic soft parameter reset. *Advances in Neural Information Processing Systems*, 37:83197–83234.

Gama, J., Žliobaitė, I., Bifet, A., Pechenizkiy, M., and Bouchachia, A. (2014). A survey on concept drift adaptation. *ACM computing surveys (CSUR)*, 46(4):1–37.

Garg, S., Balakrishnan, S., and Lipton, Z. (2022). Domain adaptation under open set label shift. *Advances in Neural Information Processing Systems*, 35:22531–22546.

Gawlikowski, J., Tassi, C. R. N., Ali, M., Lee, J., Humt, M., Feng, J., Kruspe, A., Triebel, R., Jung, P., Roscher, R., et al. (2023). A survey of uncertainty in deep neural networks. *Artificial Intelligence Review*, 56(Suppl 1):1513–1589.

Ge, Z., Demyanov, S., Chen, Z., and Garnavi, R. (2017). Generative openmax for multi-class open set classification. *arXiv preprint arXiv:1707.07418*.

Geng, C., Huang, S.-j., and Chen, S. (2020). Recent advances in open set recognition: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3614–3631.

Ghassemi, N. and Fazl-Ersi, E. (2022). A comprehensive review of trends, applications and challenges in out-of-distribution detection. *arXiv preprint arXiv:2209.12935*.

Gross, R., Con, R., and Yaakobi, E. (2025). Improved constructions of linear codes for insertions and deletions.

Haghifam, M., Negrea, J., Khisti, A., Roy, D. M., and Dziugaite, G. K. (2020). Sharpened generalization bounds based on conditional mutual information and an application to noisy, iterative algorithms. *Advances in Neural Information Processing Systems*, 33:9925–9935.

Han, K., Rebuffi, S.-A., Ehrhardt, S., Vedaldi, A., and Zisserman, A. (2020). Automatically discovering and learning new visual categories with ranking statistics. *arXiv preprint arXiv:2002.05714*.

Han, K., Rebuffi, S.-A., Ehrhardt, S., Vedaldi, A., and Zisserman, A. (2022). Autonovel: Automatically discovering and learning novel visual categories. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(10):6767–6781.

Han, K., Vedaldi, A., and Zisserman, A. (2019). Learning to discover novel visual categories via deep transfer clustering. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8401–8409.

Harutyunyan, H., Raginsky, M., Ver Steeg, G., and Galstyan, A. (2021). Information-theoretic generalization bounds for black-box learning algorithms. *Advances in Neural Information Processing Systems*, 34:24670–24682.

Hjelm, R. D., Fedorov, A., Lavoie-Marchildon, S., Grewal, K., Bachman, P., Trischler, A., and Bengio, Y. (2019). Learning deep representations by mutual information estimation and maximization.

Houlsby, N., Huszár, F., Ghahramani, Z., and Lengyel, M. (2011). Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*.

Hu, J., Liu, W., Chang, H., Ma, B., Shan, S., and Chen, X. (2024a). An information theoretical view for out-of-distribution detection. In *European Conference on Computer Vision*, pages 418–435. Springer.

Hu, S., Lou, Z., Yan, X., and Ye, Y. (2024b). A survey on information bottleneck. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(8):5325–5344.

Jafarzadeh, M., Dhamija, A. R., Cruz, S., Li, C., Ahmad, T., and Boult, T. E. (2020). A review of open-world learning and steps toward open-world learning without labels. *arXiv preprint arXiv:2011.12906*.

Jafarzadeh, M., Dhamija, A. R., Cruz, S., Li, C., Ahmad, T., and Boult, T. E. (2022). A review of open-world learning and steps toward open-world learning without labels.

Jiao, L., Wang, Y., Liu, X., Li, L., Liu, F., Ma, W., Guo, Y., Chen, P., Yang, S., and Hou, B. (2024). Causal inference meets deep learning: A comprehensive survey. *Research*, 7:0467.

Jin, Y., Xiong, Y., Fang, J., Wu, X., He, D., Jia, X., Zhao, B., and Yu, P. (2024). Beyond the known: Novel class discovery for open-world graph learning.

Jordan, M. I. and Mitchell, T. M. (2015). Machine learning: Trends, perspectives, and prospects. *Science*, 349(6245):255–260.

Kaddour, J., Lynch, A., Liu, Q., Kusner, M. J., and Silva, R. (2022). Causal machine learning: A survey and open problems. *arXiv preprint arXiv:2206.15475*.

Kawaguchi, K., Deng, Z., Ji, X., and Huang, J. (2023). How does information bottleneck help deep learning? In *International conference on machine learning*, pages 16049–16096. PMLR.

Kejriwal, M., Kildebeck, E., Steininger, R., and Shrivastava, A. (2024). Challenges, evaluation and opportunities for open-world learning. *Nature Machine Intelligence*, 6(6):580–588.

Kim, G., Xiao, C., Konishi, T., Ke, Z., and Liu, B. (2025). Open-world continual learning: Unifying novelty detection and continual learning. *Artificial Intelligence*, 338:104237.

Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., et al. (2017). Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526.

Kullback, S. and Leibler, R. A. (1951). On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86.

Kurle, R., Cseke, B., Klushyn, A., Van Der Smagt, P., and Günnemann, S. (2019). Continual learning with bayesian neural networks for non-stationary data. In *International Conference on Learning Representations*.

Kuzborskij, I., Jun, K.-S., Wu, Y., Jang, K., and Orabona, F. (2024). Better-than-kl pac-bayes bounds. In *The Thirty Seventh Annual Conference on Learning Theory*, pages 3325–3352. PMLR.

Lakkaraju, H., Kamar, E., Caruana, R., and Horvitz, E. (2017). Identifying unknown unknowns in the open world: Representations and policies for guided exploration. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31.

LeCun, Y., Bengio, Y., and Hinton, G. (2015). Deep learning. *nature*, 521(7553):436–444.

Li, A., Boyd, A., Smyth, P., and Mandt, S. (2021). Detecting and adapting to irregular distribution shifts in bayesian online learning. *Advances in neural information processing systems*, 34:6816–6828.

Li, S., Xu, R., Xiu, J., Zheng, Y., Feng, P., Yang, Y., and Liu, X. (2024). Robust multi-agent reinforcement learning by mutual information regularization.

Li, Y., Lai, G., Yang, X., Li, Y., Bonsangue, M., and Li, T. (2025a). Exploring open-world continual learning with knowns-unknowns knowledge transfer. *arXiv preprint arXiv:2502.20124*.

Li, Y., Wang, X., Yang, X., Bonsangue, M., Zhang, J., and Li, T. (2025b). Improving open-world continual learning under the constraints of scarce labeled data. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pages 1647–1658.

Li, Z. and Hoiem, D. (2017). Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947.

Li, Z., Xing, Y., Wang, X., Cai, Y., Zhou, X., and Zhang, X. (2025c). Estimating global phase synchronization by quantifying multivariate mutual information and detecting network structure. *Neural Networks*, 183:106984.

Liang, J., He, R., and Tan, T. (2025). A comprehensive survey on test-time adaptation under distribution shifts. *International Journal of Computer Vision*, 133(1):31–64.

Lidayan, A., Du, Y., Kosoy, E., Rufova, M., Abbeel, P., and Gopnik, A. (2025). Intrinsically-motivated humans and agents in open-world exploration. *arXiv preprint arXiv:2503.23631*.

Liu, B., Mazumder, S., Robertson, E., and Grigsby, S. (2023). Ai autonomy: Self-initiated open-world continual learning and adaptation. *AI Magazine*, 44(2):185–199.

Liu, J., Lin, Z., Padhy, S., Tran, D., Bedrax Weiss, T., and Lakshminarayanan, B. (2020a). Simple and principled uncertainty estimation with deterministic deep learning via distance awareness. *Advances in neural information processing systems*, 33:7498–7512.

Liu, W., Wang, X., Owens, J., and Li, Y. (2020b). Energy-based out-of-distribution detection. *Advances in neural information processing systems*, 33:21464–21475.

Liu, W., Yu, G., Wang, L., and Liao, R. (2025a). An information-theoretic framework for out-of-distribution generalization with applications to stochastic gradient langevin dynamics. *IEEE Transactions on Information Theory*.

Liu, Z., Lu, J., Xuan, J., and Zhang, G. (2025b). Learning latent and changing dynamics in real non-stationary environments. *IEEE Transactions on Knowledge and Data Engineering*.

Machlanski, D., Riley, S., Moroshko, E., Butler, K., Dimitrakopoulos, P., Melistas, T., Chanchal, A., McDonagh, S., Silva, R., and Tsaftaris, S. A. (2025). A shift in perspective on causality in domain generalization.

Mirzadeh, S. I., Farajtabar, M., Pascanu, R., and Ghasemzadeh, H. (2020). Understanding the role of training regimes in continual learning. *Advances in Neural Information Processing Systems*, 33:7308–7320.

Moazzami, K., Son, S., Lin, J., Lee, S. M., Son, D., Lee, H., Lee, J., and Lee, S. (2025). Open set recognition for endoscopic image classification: A deep learning approach on the kvasir dataset. *arXiv preprint arXiv:2506.18284*.

Mohseni, S., Wang, H., Xiao, C., Yu, Z., Wang, Z., and Yadawa, J. (2022). Taxonomy of machine learning safety: A survey and primer. *ACM Computing Surveys*, 55(8):1–38.

Mondal, S., Jiang, Z., and Sundaramoorthi, G. (2025). A variational information theoretic approach to out-of-distribution detection. *arXiv preprint arXiv:2506.14194*.

Mundt, M., Hong, Y., Pliushch, I., and Ramesh, V. (2023). A wholistic view of continual learning with deep neural networks: Forgotten lessons and the bridge to active and open world learning. *Neural Networks*, 160:306–336.

Nawaz, M. (2025). Beyond closed-set assumptions: Advancing open-set problem with adaptive learning. Doctoral Consortium, PAKDD 2025. Accessed October 13, 2025.

Negrea, J., Haghifam, M., Dziugaite, G. K., Khisti, A., and Roy, D. M. (2019). Information-theoretic generalization bounds for sgld via data-dependent estimates. *Advances in Neural Information Processing Systems*, 32.

Neu, G., Dziugaite, G. K., Haghifam, M., and Roy, D. M. (2021). Information-theoretic generalization bounds for stochastic gradient descent. In *Conference on Learning Theory*, pages 3526–3545. PMLR.

Nguyen, A. T., Tran, T., Gal, Y., Torr, P. H., and Baydin, A. G. (2021a). Kl guided domain adaptation. *arXiv preprint arXiv:2106.07780*.

Nguyen, Q. P., Low, B. K. H., and Jaillet, P. (2021b). An information-theoretic framework for unifying active learning problems. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 9126–9134.

Oliva, P. F. V., Akyildiz, O. D., and Duncan, A. (2025). Uniform-in-time convergence bounds for persistent contrastive divergence algorithms. *arXiv preprint arXiv:2510.01944*.

Ovadia, Y., Fertig, E., Ren, J., Nado, Z., Sculley, D., Nowozin, S., Dillon, J., Lakshminarayanan, B., and Snoek, J. (2019). Can you trust your model's uncertainty? evaluating predictive uncertainty under dataset shift. *Advances in neural information processing systems*, 32.

Pan, S. J. and Yang, Q. (2009). A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359.

Parmar, J., Chouhan, S., Raychoudhury, V., and Rathore, S. (2023). Open-world machine learning: applications, challenges, and opportunities. *ACM Computing Surveys*, 55(10):1–37.

Qu, Y., Tang, Y., Zhang, C., Cai, X., Yuan, X., and Zhang, W. (2025). Dual-space contrastive learning for open-world semi-supervised classification. *IEEE Transactions on Neural Networks and Learning Systems*.

Raghavan, A., Hostetler, J., and Chai, S. (2019). Generative memory for lifelong reinforcement learning.

Rebuffi, S.-A., Kolesnikov, A., Sperl, G., and Lampert, C. H. (2017). icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2001–2010.

Rios, A. S., Ndiour, I. J., Datta, P., Sydir, J., Tickoo, O., and Ahuja, N. (2024a). Uncertainty quantification in continual open-world learning. *arXiv preprint arXiv:2412.16409*.

Rios, A. S., Ndiour, I. J., Datta, P., Sydir, J., Tickoo, O., and Ahuja, N. (2024b). Uncertainty quantification in continual open-world learning.

Roulet, V., Liu, T., Vieillard, N., Sander, M. E., and Blondel, M. (2025). Loss functions and operators generated by f-divergences. *arXiv preprint arXiv:2501.18537*.

Scheirer, W. J., de Rezende Rocha, A., Sapkota, A., and Boult, T. E. (2012). Toward open set recognition. *IEEE transactions on pattern analysis and machine intelligence*, 35:1757–1772.

Scheirer, W. J., Jain, L. P., and Boult, T. E. (2014). Probability models for open set recognition. *IEEE transactions on pattern analysis and machine intelligence*, 36(11):2317–2324.

Sekar, R., Rybkin, O., Daniilidis, K., Abbeel, P., Hafner, D., and Pathak, D. (2020). Planning to explore via self-supervised world models. In *International conference on machine learning*, pages 8583–8592. PMLR.

Shannon, C. E. (1948). A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423.

Song, Y. and Kingma, D. P. (2021). How to train your energy-based models. *arXiv preprint arXiv:2101.03288*.

Song, Y., Wang, P., Xiong, W., Zhu, D., Liu, T., Sui, Z., and Li, S. (2023). Infocl: Alleviating catastrophic forgetting in continual text classification from an information theoretic perspective. *arXiv preprint arXiv:2310.06362*.

Stojanov, P., Li, Z., Gong, M., Cai, R., Carbonell, J., and Zhang, K. (2021). Domain adaptation with invariant representation learning: What transformations to learn? *Advances in Neural Information Processing Systems*, 34:24791–24803.

Sun, X., Ding, H., Zhang, C., Lin, G., and Ling, K.-V. (2021). M2iosr: Maximal mutual information open set recognition.

Sun, Y., Shi, Z., and Li, Y. (2023). A graph-theoretic framework for understanding open-world semi-supervised learning. *Advances in Neural Information Processing Systems*, 36:23934–23967.

Sun, Y., Wang, X., Fu, J., Lu, C., and Zhou, B. (2025). R$^2$AI: Towards resistant and resilient ai in an evolving world.

Tan, X., Xie, T., Zuo, Z., and Zhang, X. (2024). Exploiting open-world data for adaptive continual learning.

Tan, Z., Yang, J., Huang, W., Yuan, Y., and Zhang, Y. (2023). Information flow in self-supervised learning. *arXiv preprint arXiv:2309.17281*.

Tang, P., Xu, Y., Jiao, Y., Zhang, M., Song, Y., and Ding, G. (2025). Similarity-adaptive framework for semi-supervised open-world specific emitter identification. *IEEE Transactions on Information Forensics and Security*.

Thiagarajan, J. J., Anirudh, R., Narayanaswamy, V., and Bremer, P.-T. (2022). Single model uncertainty estimation via stochastic data centering.

Tishby, N., Pereira, F. C., and Bialek, W. (2000). The information bottleneck method. *arXiv preprint physics/0004057*.

Tishby, N. and Zaslavsky, N. (2015). Deep learning and the information bottleneck principle. In *2015 ieee information theory workshop (itw)*, pages 1–5. Ieee.

Tsur, D., Goldfeld, Z., and Greenewald, K. (2023). Max-sliced mutual information. *Advances in neural information processing systems*, 36:80338–80351.

van den Oord, A., Li, Y., and Vinyals, O. (2019). Representation learning with contrastive predictive coding.

Varley, T. F. (2023). Information theory for complex systems scientists. *arXiv preprint arXiv:2304.12482*.

Vaze, S., Han, K., Vedaldi, A., and Zisserman, A. (2022). Open-set recognition: a good closed-set classifier is all you need?

Wang, C., Luo, S., Pei, J., Liu, X., Huang, Y., Zhang, Y., and Yang, J. (2023a). An entropy-awareness meta-learning method for sar open-set atr. *IEEE Geoscience and Remote Sensing Letters*, 20:1–5.

Wang, H., Pang, G., Wang, P., Zhang, L., Wei, W., and Zhang, Y. (2023b). Glocal energy-based learning for few-shot open-set recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7507–7516.

Wang, H., Vaze, S., and Han, K. (2024a). Hilo: A learning framework for generalized category discovery robust to domain shifts. *arXiv preprint arXiv:2408.04591*.

Wang, L., Zhang, X., Su, H., and Zhu, J. (2024b). A comprehensive survey of continual learning: Theory, method and application. *IEEE transactions on pattern analysis and machine intelligence*, 46(8):5362–5383.

Wang, L., Zhang, X., Su, H., and Zhu, J. (2024c). A comprehensive survey of continual learning: Theory, method and application. *IEEE transactions on pattern analysis and machine intelligence*, 46(8):5362–5383.

Wang, Q., Fink, O., Van Gool, L., and Dai, D. (2022). Continual test-time domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7201–7211.

Wang, Y., Li, B., Che, T., Zhou, K., Liu, Z., and Li, D. (2021). Energy-based open-world uncertainty modeling for confidence calibration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9302–9311.

Weimar, M., Rachbauer, L. M., Starshynov, I., Faccio, D., Adilova, L., Bouchet, D., and Rotter, S. (2025). Fisher information flow in artificial neural networks. *Physical Review X*, 15(3):031072.

Wen, W., Gong, T., Zhang, Y., Gao, Z., Zhang, W., and Liu, Y.-J. (2025a). Information-theoretic generalization bounds of replay-based continual learning. *arXiv preprint arXiv:2507.12043*.

Wen, W., Gong, T., Zhang, Y., Gao, Z., Zhang, W., and Liu, Y.-J. (2025b). Information-theoretic generalization bounds of replay-based continual learning. *arXiv preprint arXiv:2507.12043*.

Westphal, C., Hailes, S., and Musolesi, M. (2024). Information-theoretic state variable selection for reinforcement learning. *arXiv preprint arXiv:2401.11512*.

Wiles, O., Gowal, S., Stimberg, F., Alvise-Rebuffi, S., Ktena, I., Dvijotham, K., and Cemgil, T. (2021). A fine-grained analysis on distribution shift. *arXiv preprint arXiv:2110.11328*.

Wilson, A. G. and Izmailov, P. (2020). Bayesian deep learning and a probabilistic perspective of generalization. *Advances in neural information processing systems*, 33:4697–4708.

Wood, D., Papamarkou, T., Benatan, M., and Allmendinger, R. (2024). Model-agnostic variable importance for predictive uncertainty: an entropy-based approach. *Data Mining and Knowledge Discovery*, 38(6):4184–4216.

Wutschitz, L., Köpf, B., Paverd, A., Rajmohan, S., Salem, A., Tople, S., Zanella-Béguelin, S., Xia, M., and Rühle, V. (2023). Rethinking privacy in machine learning pipelines from an information flow control perspective.

Xie, T., Zhang, J., Bai, H., and Nowak, R. (2024). Deep active learning in the open world. *arXiv preprint arXiv:2411.06353*.

Xing, H., Boey, K. Z., and Cheng, G. (2025). Towards open-world human action segmentation using graph convolutional networks. *arXiv preprint arXiv:2507.00756*.

Xu, A. and Raginsky, M. (2017). Information-theoretic analysis of generalization capability of learning algorithms. *Advances in neural information processing systems*, 30.

Xu, H., Liu, B., Shu, L., and Yu, P. (2019). Open-world learning and application to product classification. In *The World Wide Web Conference*, pages 3413–3419.

Xu, X., Perin, G., and Picek, S. (2023). Ib-rar: Information bottleneck as regularizer for adversarial robustness.

Xu, X., Qin, Y., Mi, L., Wang, H., and Li, X. (2024). Energy-based concept bottleneck models: Unifying prediction, concept intervention, and probabilistic interpretations. *arXiv preprint arXiv:2401.14142*.

Xue, C. (2024). *Advancing AI Capabilities for Dynamic Physical Environments: Transitioning from Closed-World Problem Solving to Open-World Challenges*. PhD thesis, The Australian National University (Australia).

Yu, Y., Yang, Z., Dobriban, E., Steinhardt, J., and Ma, Y. (2021). Understanding generalization in adversarial training via the bias-variance decomposition.

Zhang, J. (2025). Ai for the open-world: the learning principles. *arXiv preprint arXiv:2504.14751*.

Zhang, L., Qi, L., Yang, X., Qiao, H., Yang, M.-H., and Liu, Z. (2022a). Automatically discovering novel visual categories with self-supervised prototype learning. *arXiv preprint arXiv:2208.00979*.

Zhang, M., Fernández-Torres, M.-Á., Cohrs, K.-H., and Camps-Valls, G. (2025). Calibration and uncertainty quantification for deep learning-based drought detection. *International Journal of Applied Earth Observation and Geoinformation*, 140:104563.

Zhang, Z., Wang, X., Zhang, Z., Li, H., Qin, Z., and Zhu, W. (2022b). Dynamic graph neural networks under spatio-temporal distribution shift. *Advances in neural information processing systems*, 35:6074–6089.

Zhao, X., Ma, Y., Wang, D., Shen, Y., Qiao, Y., and Liu, X. (2023). Revisiting open world object detection. *IEEE Transactions on Circuits and Systems for Video Technology*, 34:3496–3509.

Zhou, K., Liu, Z., Qiao, Y., Xiang, T., and Loy, C. C. (2022). Domain generalization: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 45(4):4396–4415.

Zhou, Z.-H. (2022). Open-environment machine learning. *National Science Review*, 9(8):nwac123.

Zhou, Z.-H., Fang, S., Zhou, Z.-J., Wei, T., Wan, Y., and Zhang, M.-L. (2024). Continuous contrastive learning for long-tailed semi-supervised recognition. *Advances in Neural Information Processing Systems*, 37:51411–51435.

Zhu, F., Ma, S., Cheng, Z., Zhang, X.-Y., Zhang, Z., Tao, D., and Liu, C.-L. (2024). Open-world machine learning: A systematic review and future directions. *arXiv preprint arXiv:2403.01759*.