# A Comprehensive Comparative Study of Intuitive Physics Modeling in Machine Learning Trained with Cartoon and Realistic Data

**Songyu Yang**

College of Computer Science and Technology, Zhejiang University, HangZhou, China

3220104690@zju.edu.cn

**Abstract.** This study delves into the influence of training data types—specifically cartoon versus realistic visual datasets—on the development of intuitive physics modeling in machine learning. Intuitive physics, the inherent human ability to understand and predict the physical properties and dynamics of objects, presents a significant challenge for current AI systems to replicate accurately. Leveraging YOLOv5, a cutting-edge object detection model, this research systematically evaluates the cognitive understanding and performance of AI models trained on distinct types of visual data. The findings reveal that the visual complexity inherent in the training datasets plays a crucial role in shaping the model's ability to generalize and accurately perform intuitive physics tasks. Models trained on cartoon datasets exhibited different learning patterns and generalization capabilities compared to those trained on realistic data, providing valuable insights into the role of data representation in AI training. This research offers both theoretical advancements in understanding AI's cognitive limitations and practical guidance for designing AI systems that can interact with the physical world more effectively. Ultimately, the study contributes to bridging the gap between human cognition and machine learning, pushing the boundaries of what AI can achieve in modeling complex, real-world phenomena.

**Keywords:** Intuitive physics, cartoon data, realistic data, model generalization, AI systems.

## 1. Introduction

**Research Background:** The intersection of artificial intelligence (AI) and cognitive science has generated considerable interest in replicating human-like intuitive physics within machine learning models. Intuitive physics refers to the innate human ability to understand and predict the physical properties and dynamics of objects—a skill that is evident even in early childhood. Despite advancements in AI, current systems still struggle to achieve the level of physical understanding that humans naturally possess, revealing a significant gap between human cognition and machine learning capabilities. This disparity underscores the need for further research to enhance AI's ability to model and understand the physical world as intuitively as humans do.

**Research Problem:** Recent efforts to bridge this gap have focused on developing AI models inspired by developmental psychology, particularly those that leverage visual cognition to teach machines about intuitive physics. For example, the use of the violation-of-expectation (VoE) paradigm has demonstrated that AI can learn physical concepts from visual data, though the depth of understanding still lags behind that of humans [1]. Additionally, metrics like Multi-Scale Structural Complexity (MSSC) have been

proposed to evaluate visual complexity, influencing how AI systems learn from varying data types [2]. Moreover, integrating shared representations across visual and language models has shown potential in improving AI's generalization capabilities across multiple tasks [3]. Despite these advancements, a crucial question remains: How does the nature of the training data—whether cartoon-like or realistic—affect the development of intuitive physics in AI models?

**This Work:** This study seeks to explore the impact of different types of visual datasets—cartoon versus realistic—on the training of machine learning models in the context of intuitive physics. By comparing the cognitive understanding and performance of models trained on these diverse datasets, this research aims to shed light on how the nature of training data influences the development of intuitive physics in AI. The findings from this study will not only contribute to the theoretical understanding of AI's cognitive capabilities but also provide practical insights for designing more effective AI systems capable of interacting with the physical world in a human-like manner.

## 2. Relevant Theories

### 2.1. Overview of intuitive physics

Intuitive physics is a critical aspect of human cognition, enabling individuals to understand and predict the behavior of physical objects in the environment. According to research, intuitive physics is fundamental to embodied intelligence, as it is essential for practical actions and provides a foundation for conceptual knowledge and compositional representation. Even infants demonstrate an early understanding of fundamental physical principles, such as object permanence, continuity, and solidity, which are essential for interacting with the physical world. These foundational concepts are critical for developing AI models that attempt to replicate human-like intuitive physics [4]. Despite significant advancements in AI, current systems still struggle to match the intuitive physics understanding of young children, revealing a gap between human cognitive abilities and machine learning models.

In efforts to bridge the gap between human and machine understanding, researchers have drawn inspiration from developmental psychology. For example, models using the VoE paradigm can probe AI systems' understanding by measuring their responses to physically possible and impossible scenarios, providing insights into how these systems perceive physical interactions [5]. Studies show that AI models inspired by these developmental principles can acquire some understanding of intuitive physics, but they still fall short of fully replicating the depth and adaptability of human cognition.

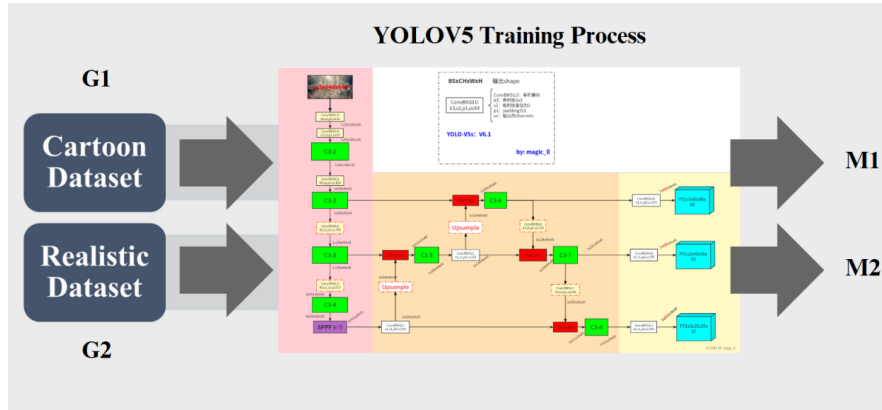### 2.2. Characteristics of cartoon and realistic data

The type of data used to train AI models significantly impacts their performance in modeling intuitive physics. Cartoon data, with its simplified and exaggerated features, allows models to focus on fundamental physical principles without the distractions of real-world noise and complexity. This can lead to more accurate modeling of basic physical interactions, such as object stability and motion prediction [6].

In contrast, realistic data, characterized by detailed textures, lighting, and complex object interactions, provides a closer approximation to real-world scenarios but introduces challenges related to visual complexity. Measures such as Multi-Scale Structural Complexity (MSSC) have been proposed to quantify this visual complexity, offering a way to assess how different datasets might impact a model's learning process. Studies have shown that MSSC correlates strongly with human perceptions of complexity, making it a valuable tool for evaluating training datasets [7].

### 2.3. Applications of machine learning and deep learning in intuitive physics

Machine learning, particularly deep learning, has seen significant applications in modeling intuitive physics. One approach is to utilize models like YOLOv5, which excels in real-time object detection, allowing for accurate predictions of physical interactions in dynamic environments [8]. By leveraging shared representations across tasks, such models can generalize better across different scenarios, enhancing their applicability in various domains [9].

However, the effectiveness of these models often depends on the complexity of the training data. Research indicates that models trained on data with balanced levels of visual complexity—neither too simple nor overly complex—tend to perform better in generalizing to new, unseen scenarios. This finding underscores the importance of carefully curating training datasets to optimize model performance in intuitive physics tasks [10]. As shown in Figure 1.



**Figure 1.** YOLOv5 Training Process: Comparative Analysis of Cartoon and Realistic Datasets (Photo credit: Original).

## 3. System Analysis and Comparative Study

### 3.1. Description of datasets and models

In this study, we conducted a comprehensive analysis using two distinct types of visual datasets: cartoon data and realistic data. The cartoon dataset was characterized by simplified, exaggerated features with minimal visual noise, providing a clear and controlled environment for training machine learning models. This dataset was designed to focus on fundamental physical interactions, such as collisions, object stability, and basic motion, which are easier to model due to the absence of complex textures and lighting variations.

In contrast, the realistic dataset was composed of images and videos that closely mimic real-world conditions. This dataset included detailed textures, varied lighting conditions, and more complex object interactions, making it a challenging environment for machine learning models. The realistic dataset was intended to test the models' ability to generalize intuitive physics understanding to scenarios that closely resemble real-life physical events.

For model implementation, we employed YOLOv5, a state-of-the-art object detection model known for its speed and accuracy in identifying objects in images and video sequences. YOLOv5 was chosen due to its effectiveness in real-time object detection, which is crucial for analyzing physical interactions in dynamic environments. The model was trained separately on both the cartoon and realistic datasets to evaluate how the nature of the visual data impacts its ability to model intuitive physics.

The primary focus of the analysis was to assess the performance of YOLOv5 in predicting physical interactions, such as object permanence, continuity, and collisions, in both datasets. The model's predictions were evaluated against a set of predefined physical principles to determine its accuracy and ability to generalize from training data to unseen scenarios. This approach allowed us to compare how different types of visual complexity influence the model's understanding and prediction of physical events, providing insights into the relationship between data complexity and intuitive physics performance.

### 3.2. Model training and evaluation

The models were trained using both cartoon and realistic datasets, with training parameters optimized for each dataset to ensure the best possible performance. During evaluation, the models were tested on

unseen data from both cartoon and realistic environments to assess their generalization capabilities. Metrics such as accuracy in object detection, collision prediction, and motion tracking were used to evaluate performance. The results were analyzed to understand how the visual complexity of the training data influenced the models' intuitive physics reasoning.

### 3.3. The Relationship between visual complexity and intuitive physics performance

The relationship between visual complexity and a model's performance in intuitive physics tasks is crucial in understanding how different data types impact machine learning outcomes. Visual complexity, quantified by measures such as Multi-Scale Structural Complexity (MSSC), plays a significant role in how models process and learn from data. Our study revealed that the complexity of visual data directly affects the performance of YOLOv5 in modeling intuitive physics.

When trained on cartoon data, YOLOv5 demonstrated high accuracy in predicting physical interactions. The simplified nature of the cartoon dataset allowed the model to focus on the fundamental aspects of physical events without being hindered by the noise and intricacies present in more complex visual environments. This resulted in better performance in tasks such as object tracking, collision detection, and motion prediction.

However, when the model was trained on the realistic dataset, its performance varied significantly. The increased visual complexity, including detailed textures and varied lighting, posed challenges for the model, leading to a slight decrease in accuracy. This suggests that while realistic data provides a richer and more diverse set of training scenarios, it also requires the model to handle a higher level of complexity, which can affect its ability to generalize intuitive physics principles.

The comparative analysis highlighted that while both datasets have their merits, the choice between cartoon and realistic data should depend on the specific goals of the machine learning application. For applications requiring precise and clear understanding of basic physical principles, cartoon data may be more effective. In contrast, realistic data is essential for developing models that need to perform in real-world conditions, despite the additional challenges posed by higher visual complexity.

In conclusion, the study underscores the importance of considering visual complexity when designing training datasets for machine learning models aimed at replicating human-like intuitive physics. The balance between simplicity and realism in visual data is key to optimizing model performance across different types of physical reasoning tasks.

## 4. Challenges and Future Directions

### 4.1. Limitations of the current study

In this study, titled A Comprehensive Comparative Study of Intuitive Physics Modeling in Machine Learning Trained with Cartoon and Realistic Data: Analyzing the Impacts on Cognitive Understanding and Model Performance, several limitations were identified that may impact the generalizability and applicability of the findings.

Firstly, the study was limited by the scope and type of datasets used. The cartoon dataset, while useful for focusing on fundamental physical principles, may not fully capture the complexity and variety of physical interactions found in more realistic scenarios. Similarly, the realistic dataset, despite its complexity, may introduce noise and distractions that could obscure the model's ability to learn and generalize from the data effectively. This dichotomy suggests that the datasets, while complementary, may not represent the full spectrum of visual environments in which AI systems need to operate.

Secondly, the use of YOLOv5, although effective for object detection, may not be the most suitable model for all aspects of intuitive physics. YOLOv5 excels in real-time object detection, but its ability to model more intricate physical interactions and causal relationships between objects could be enhanced by integrating it with other deep learning models or using more specialized approaches like physics-informed neural networks.

Thirdly, the study focused primarily on visual complexity and its impact on intuitive physics modeling, potentially overlooking other factors that could influence model performance, such as

temporal dynamics, object material properties, and the role of prior knowledge in physical reasoning. Additionally, the study's reliance on pre-defined physical principles as benchmarks may have limited the exploration of more complex or emergent physical behaviors that could arise in diverse environments.

Lastly, the generalizability ofof the findings is constrained by the specific training and evaluation conditions used in this study. The results may vary significantly in different contexts, such as with different types of models, alternative datasets, or in real-world applications where environmental variables are less controlled.

### 4.2. Future research directions

Building on the findings of this study, several future research directions are proposed to address the identified limitations and advance the field of intuitive physics modeling in machine learning.

Expansion of Datasets: Future research should consider expanding the variety and complexity of datasets used for training and evaluating intuitive physics models. Incorporating datasets that blend elements of both cartoon and realistic visuals, as well as introducing more diverse physical interactions, could provide a more comprehensive training environment. This would allow models to learn from a broader range of scenarios, potentially improving their ability to generalize to real-world conditions.

Integration of Advanced Models: While YOLOv5 is effective for object detection, future studies should explore integrating it with other models that are better suited for capturing complex physical interactions. Combining YOLOv5 with physics-informed neural networks, graph-based models, or reinforcement learning frameworks could enhance the model's ability to reason about physical causality and predict the outcomes of interactions with greater accuracy.

Incorporating Temporal and Material Dynamics: Future research should also focus on incorporating temporal dynamics and material properties into intuitive physics modeling. Understanding how objects move and interact over time, and how their physical properties (such as mass, friction, and elasticity) influence these interactions, is crucial for developing more robust models. This could involve the use of simulations or real-world data capturing dynamic interactions to train models on more complex physical reasoning tasks.

Exploring Emergent Behaviors: Another promising direction is the exploration of emergent physical behaviors that arise in complex systems. This could involve creating environments where the physical principles are not pre-defined but instead emerge from the interactions of multiple objects. Models could then be evaluated on their ability to discover and understand these emergent behaviors, pushing the boundaries of intuitive physics modeling.

Application to Real-World Problems: Finally, future research should focus on applying intuitive physics models to real-world problems, such as robotics, autonomous vehicles, and augmented reality. These applications require models that can handle the unpredictability and complexity of real environments. Research should explore how intuitive physics models can be integrated into these systems to improve their safety, efficiency, and overall performance.

### 5. Conclusion

The findings indicate that while cartoon data allows models to concentrate on fundamental physical principles without the distraction of visual noise, realistic data introduces a level of complexity that better approximates real-world scenarios. This complexity, however, can challenge the model's ability to generalize intuitive physics concepts effectively. The study highlights the trade-offs between employing simplified versus realistic data in training machine learning models and underscores the importance of selecting appropriate datasets based on the specific goals of the application. Moreover, this research demonstrates the potential of integrating insights from cognitive science with advanced machine learning techniques to create models that more closely replicate human-like intuitive physics. It also points out the limitations of current approaches and suggests several avenues for future research, including the expansion of datasets, integration of advanced models, and application of these models to real-world scenarios.

In conclusion, this study contributes to the ongoing effort to develop AI systems capable of understanding and interacting with the physical world more naturally and effectively. By advancing the understanding of the relationship between visual complexity and model performance, this research hopes to pave the way for more robust and generalizable AI systems capable of intuitive physical reasoning in diverse environments.

## References

[1]  Piloto L S, et al. Intuitive physics learning in a deep-learning model inspired by developmental psychology. Nature Human Behaviour, 2022. DOI: 10.1038/s41562-022-01394-8.

[2]  Kravchenko A., Bagrov A. A., Katsnelson M. I., Dudarev V. Multi-scale structural complexity as a quantitative measure of visual complexity. arXiv preprint, 2024. arXiv:2408.04076. DOI: 10.48550/arXiv.2408.04076.

[3]  Anonymous. Cross-modal Learning with Shared Representations in Vision and Language Models. arXiv preprint, 2024. arXiv:2202.06481. DOI: 10.48550/arXiv.2202.06481.

[4]  Riochet R., Castro M. Y., Bernard M., et al. Intphys 2019: A benchmark for visual intuitive physics understanding. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2021, 44(9): 5016-5025.

[5]  Zhu X., Guo C., Feng H., Huang Y., Feng Y., Wang X., Wang R. A Review of Key Technologies for Emotion Analysis Using Multimodal Information. Cogn. Comput., 2024, 1(1): 1-27.

[6]  Riochet R., Castro M. Y., Bernard M., et al. Intphys: A framework and benchmark for visual intuitive physics reasoning. arXiv preprint, 2018. arXiv:1803.07616

[7]  Vicovaro M. Grounding Intuitive Physics in Perceptual Experience. Journal of Intelligence, 2023, 11(10): 187.

[8]  Suomala J., Kauttonen J. Human's intuitive mental models as a source of realistic artificial intelligence and engineering. Frontiers in psychology, 2022, 13: 873289.

[9]  Bates C. J., Yildirim I., Tenenbaum J. B., et al. Modeling human intuitions about liquid flow with particle-based simulation. PLoS computational biology, 2019, 15(7): e1007210.

[10]  Mitko A., Fischer J. When it all falls down: The relationship between intuitive physics and spatial cognition. Cognitive research: principles and implications, 2020, 5: 1-13.