

Next Generation Deep Learning Based on Simulators and Synthetic Data

Celso M. de Melo,^{1*} Antonio Torralba,² Leonidas Guibas,³ James DiCarlo,⁴ Rama Chellappa,⁵
Jessica Hodgins⁶

¹ Computational and Information Sciences Directorate, DEVCOM U.S. Army Research
Laboratory, Playa Vista, CA, USA

² Department of Electrical Engineering and Computer Science, Massachusetts Institute of
Technology, Cambridge, MA, USA

³ Computer Science Department, Stanford University, Stanford, CA, USA

⁴ Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology,
Cambridge, MA, USA

⁵ Electrical and Computer Engineering, Johns Hopkins University, Baltimore, MD, USA

⁶ Computer Science Department, Carnegie Mellon University, Pittsburgh, PA, USA

Correspondence concerning this article should be addressed to Celso M. de Melo,
DEVCOM U.S. Army Research Laboratory, Playa Vista, CA 90094, United States. Email:
celso.miguel.de.melo@gmail.com.

Abstract

Deep learning is being successfully applied across multiple domains, yet these models learn in a most artificial way: they require large quantities of labelled data to grasp even simple concepts. Thus, the main bottleneck often is access to supervised data. We review a trend on a potential solution to this challenge: synthetic data. Synthetic data is becoming accessible due to progress in rendering pipelines, generative adversarial models, and fusion models. Moreover, advancement in domain adaptation techniques help close the statistical gap between synthetic and real data. Paradoxically, this artificial solution is also likely to enable more natural learning as we see in biological systems, including continual, multimodal, and embodied learning. Complementary, simulators and deep neural networks will also play a critical role in providing insight on the cognitive and neural functioning of biological systems. We review strengths, opportunities, and novel challenges with synthetic data.

Keywords

Deep neural networks; Synthetic data; Graphics rendering pipelines; Generative adversarial networks; Domain adaptation; Next generation learning.

The Bottleneck is Labeled Data

The last decade has experienced a revolution in interest and investment in deep learning that has enabled successful applications in visual perception, natural language processing, and robotic control, among others [1]. Deep learning's success benefited from converging trends in the development of algorithms to train these models (e.g., backpropagation), the availability of "big data" (e.g., social media), and advances in computational power (e.g., powerful graphical processing units or GPUs). However, despite these initial successes, it is becoming apparent that the current generation of deep neural networks (DNNs) has important practical and theoretical limitations. DNNs are sample-inefficient in that they require large amounts of annotated data (e.g., images of vehicles with bounding boxes) to optimize all its parameters (typically in the order of millions). Therefore, rather than algorithm or computational capability, the availability of annotated data is often the main bottleneck in the development of deep learning models. **Synthetic data** and **simulators** have emerged as a promising solution to this challenge [2]. Synthetic data are, comparatively, easier to generate, inexhaustible, pre-annotated, and less expensive. Synthetic data also have the potential to avoid ethical (e.g., privacy concerns) and practical issues (e.g., security concerns). Synthetic data further introduce unique opportunities in that they enable training data that may be impractical or impossible to collect in the real world.

More fundamentally, DNNs still lack important capabilities seen in biological systems. Humans are able to learn rich representations of the world, including about its (hierarchical) compositional and physical nature [3, 4]. Humans are more efficient learners, often being able to grasp novel concepts from a small sample of examples [5] and in mostly unsupervised fashion [6]. Moreover, human learning is sophisticated often relying on rich interactive experiences, in contrast to static datasets that capture "moments in time" (e.g., ImageNet). Synthetic data and simulators are a new catalyst for these richer representations of the world and more sophisticated forms of learning, including **multimodal learning** [7] (e.g., fusing visual and audio information), **continual learning** [8] (e.g., understanding gradually more complex tasks in sequence), and **embodied learning** [9] (e.g., interactive exploratory play to understand object affordances). Complementary, simulators can be used to gather unique insight on biological systems [10, 11]. By comparing different artificial neural models with respect to how well they simulate cognitive functionality and predict brain activity it becomes possible to test, validate, and extend existent theory [12]. Insofar as simulated data enables training and testing DNNs, thus, it plays an instrumental role in the study of biological systems. Simulated data further present novel opportunities for scientific exploration. By analyzing the properties of DNNs, it is possible to synthesize optimized stimuli to activate specific neural populations with relevant application to the study of brain function [13]. Simulated environments, perhaps even fully immersive (e.g., virtual reality), can further provide a unique opportunity for direct comparison of behavior and neural activation in DNNs vs. humans vs. nonhuman primates in embodied interactive tasks. Synthetic data, therefore, can further the development of artificial neural networks that model critical function we see in biological systems, simultaneously contributing to our understanding of these systems and offering solutions with broad practical relevance.

The article starts by overviewing successful methodologies used to synthesize data for deep learning models, emphasizing the integration of the synthesis and machine learning pipelines. Next, we focus on a central challenge to using simulated data: aligning synthetic data to real data at the pixel and feature level. We then articulate how synthetic data and simulators can enable deep learning solutions that can learn richer representations of the world and learn in sophisticated ways, while simultaneously providing insight on the biological systems they draw inspiration from.

Synthesizing Data and Integrating with the Deep Learning Pipeline

Progress in computer graphics tools, such as game engines (e.g., Unity and Unreal), and the increasing availability of three-dimensional (3D) assets, is making it easier to develop simulators for custom domains (Figure 1A). This approach has been used to synthesize training data for a variety of tasks, including object detection [15-17], object tracking [18, 19], viewpoint estimation [20], semantic segmentation [21-23], robot manipulation and control [24-28], pose estimation [29-31], gaze estimation [32], and activity recognition [33, 34] (for a detailed review of simulators and synthetic datasets see: [35]). Synthetic data, across these diverse domains, often led to improvement in deep neural network performance when tested in real domains, especially when combined with real data. In this approach, synthesis relies on a **computer graphics rendering pipeline**, which takes as input 3D information about the scene (e.g., points in three-dimensional space specifying a vehicle), information about the materials and lighting properties (e.g., vehicle color and light sources), rendering parameters (e.g., rasterization or raytracing algorithm), and produces a 3D visualization of the scene (Figure 1D). Since the pipeline has information about the scene details it can automatically generate error-free ground-truth (e.g., bounding boxes for objects of interest, depth information, and scene segmentation masks). By increasing the amount of 3D information (e.g., the number of 3D vertices specifying the objects of interest) and the sophistication of the algorithms used to render the scene, it is possible to increase the visual realism of the output, i.e., the visual fidelity of the scene when compared to the real world. Similarly, it is possible to increase the motion realism of the output by using 3D motion capture techniques (e.g., for human activity recognition) and sophisticated physics engines (e.g., for robot manipulation). In general, increasing the realism of the synthesized output tends to improve deep learning performance [20, 36-39], though in some cases it is less important [17, 40-43]. Achieving high levels of realism (e.g., as seen in movies), however, can be costly. One alternative approach is to generate synthetic data and then improve realism by using domain adaptation techniques, as discussed in the next section. Another alternative is to use generative adversarial models.

Generative Adversarial Networks (GANs) are a promising technique to synthesize novel images that match the statistical properties of the training data (Figure 1B) [44] (for a recent survey see: [45]) – for instance, GANs can generate faces of people that don't exist from a training set of existent human faces [46]. GANs consist of two models trained to optimize

opposite objectives (i.e., adversarial): a generator and a discriminator (Figure 1E). The generator learns a lower-dimension latent representation of the training data domain and is able to generate new samples by receiving as input a random vector in the latent space. The discriminator, in turn, learns to distinguish original images from synthesized images. By training the generator and discriminator simultaneously, the generator learns to synthesize better samples, so as to fool the discriminator. GANs are becoming increasingly popular due to the high visual quality of synthesized imagery [46-49], in particular when compared to other generative approaches such as variational autoencoders [45]. However, in its original formulation, it is hard to control the output produced by GANs, though this is still an area of active research. A promising trend consists of conditioning GANs on additional input that characterizes the samples being fed in training (e.g., labels specifying the gender of human faces) [50]. This idea has been extended to allow sophisticated control in the generation of images [46, 51] (e.g., pose and hairstyle of human faces). One challenge with using GANs is that the synthesized imagery is not produced with the associated ground truth data, as for graphics pipelines. However, good progress is being made extending GANs to produce imagery that already comes with detailed annotation, such as images of scenes with automatically generated scene segmentation ground truth [52]. Another recent trend has been to train big generative models (e.g., with billions of parameters and terabytes of data) [53], including language [54] and multimodal models [55], that can subsequently be reused to synthesize novel content and be integrated with other pipelines to solve domain-specific tasks.

A third approach for synthesizing data consists of creating imagery by fusing from multiple data sources (Figure 1C). Often this is accomplished by superimposing virtual objects [56-58] or people [59, 60] on real backgrounds, while ensuring that the virtual entities fit consistently with the background (e.g., by aligning surfaces and lighting). Extending this approach to fuse real entities on real backgrounds brings the extra challenge of cropping the real entities from the original backgrounds. Whereas this could be done manually, GAN-based methods have shown promise in automatically finding the cropping region (i.e., the semantic mask) with minimal annotation (e.g., bounding boxes) [61, 62]. By combining segmentation with domain adaptation techniques, it is further possible to replace in place one type of entity for another (e.g., a bicycle for a motorcycle) while preserving the rest of the image [63, 64].

Finally, we see much promise in integrating the synthesis and learning pipelines. There is a long history of integrating simulators with the learning process in deep reinforcement learning, where it is often impractical or impossible to train in the real world [65, 66]. Reinforcement learning agents learn an action policy (e.g., grasping objects or playing a game) by practicing (millions of times) in simulators [67, 68]. The key distinction is the integration of the learning process with the simulator, rather than relying on a static dataset of simulated data for training. This powerful idea can be extended to support more sophisticated forms of learning, such as continuous (lifelong) learning and embodied (interactive) learning, which we further discuss below. The concept can be applied to supervised learning by using error signals, such as a task classification loss, to optimize data synthesis generation [69] (Figure 1F). When using graphics

rendering pipelines, one challenge is to propagate the error signal through the non-differentiable functions implemented in traditional pipelines. An emerging field, called **neural rendering**, aims to build differentiable rendering pipelines and is showing fast progress in generating controllable visually realistic rendering [70, 71] (for a review see [72]). Integration of deep learning and differentiable rendering pipelines, thus, holds the promise to support the generation of customized curricula for more sample-efficient learning.

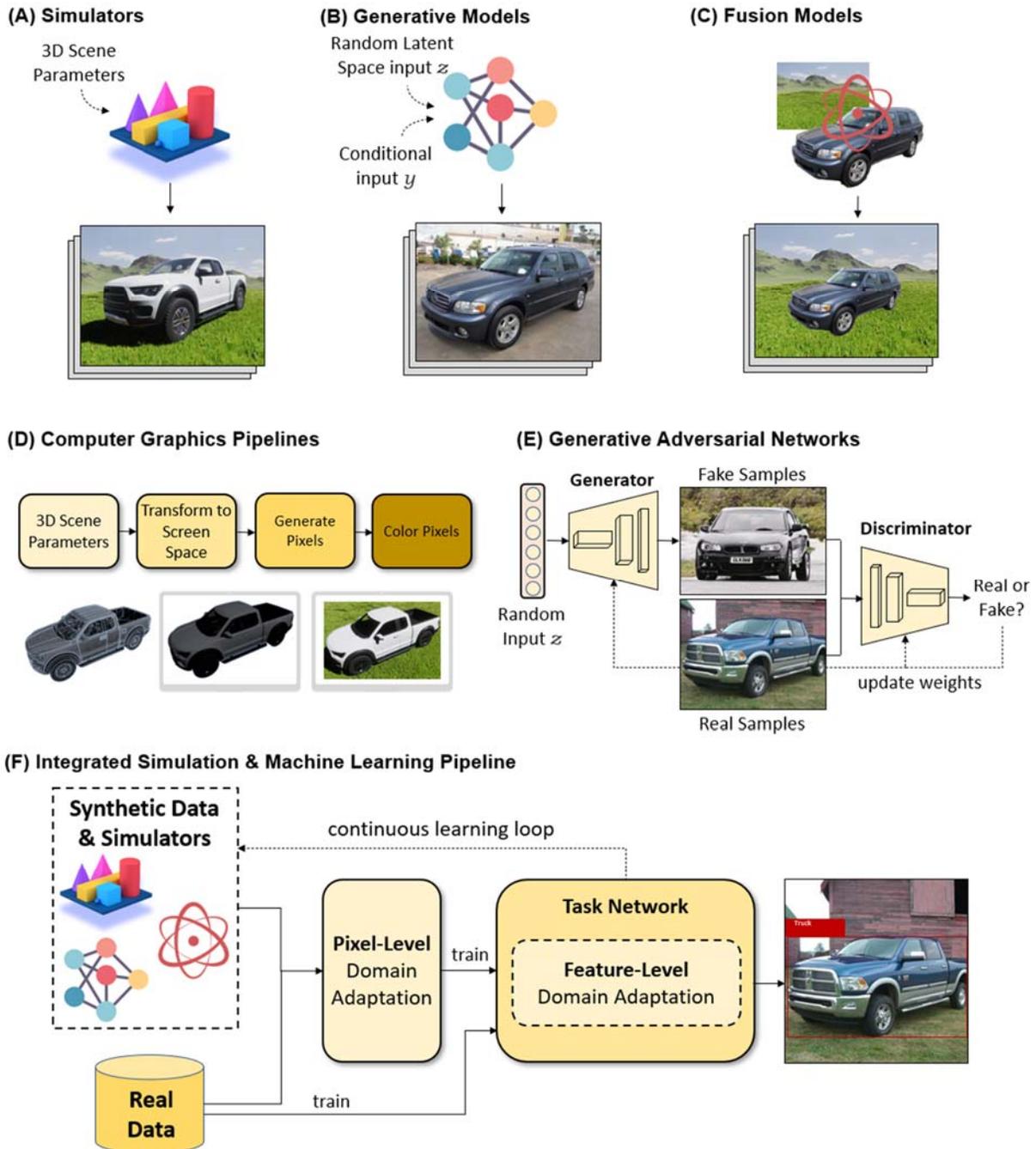


Figure 1. **Data synthesis approaches and integration with the machine learning pipeline.**

Data can be synthesized using computer graphics rendering pipelines (A), generative adversarial models (B), and fusion models (C). The traditional computer graphics pipeline (D) receives as input the 3D information about the scene and renders, through stages, a visualization of the scene in screen space. Generative adversarial networks (E) rely on a generator and discriminator network learning simultaneously on competing objectives, which leads the former to improve the quality of the synthesized imagery. Integrating the synthesis and deep learning pipelines (F) enables more sophisticated learning, such as embodied continuous learning.

Closing the Gap Between Synthetic and Real Data

Despite their success in achieving state-of-the-art performance in several visual recognition tasks, neural networks suffer from **domain shift** i.e., the performance of neural networks drops significantly when the test distribution is different from the training distribution, such as when training on synthetic data and testing on real data. To close this gap, several techniques have been developed to enhance the value of synthetic data. **Domain randomization** consists of varying the parameters used to generate the synthetic data, so that the dataset broadly captures the distribution in the target domain [16, 28]. By training on such a diverse dataset, the hope is that the model will be more robust to variation in the target domain and generalize better to novel samples. In some cases, this idea was even pushed to create non-photorealistic versions of the data (e.g., vehicles with random textures) to encourage the model to learn better representations of the target concepts (e.g., features that capture the shape, rather than texture, of vehicles) [17]. Mixing real and synthetic data (**hybrid models**) has also often led, in practice, to a boost in performance, when compared to training only on one type [34]. The idea is that mixing data allows different data types to strengthen training where others may have weaknesses (e.g., synthetic data tends to be more diverse, but real data may capture low level details better).

An increasingly prominent technique is **domain adaptation**, which consists of aligning the synthetic data pixel and feature distribution to the real data (Figure 2). Pixel-level adaptation consists of transferring the style, or visual appearance properties, of the target to the source domain. Approaches based on adversarial generative models are showing increasing success in creating realistic versions of the synthetic data, even without the need for any supervision (i.e., no labels are necessary) [73-77]. Recent promising techniques preserve semantic consistency when translating from source to target through cycle consistency (i.e., the translation needs to learn to go from source to target and back) [75], patch consistency (i.e., image patches in source and target domain should reflect the same content) [76], and leveraging intermediate representations from an integrated computer graphics pipeline (e.g., depth and color masks) [77].

Whereas in pixel alignment the goal is to adjust the visual style of the source domain, in feature-level adaptation the distributional distance between source and target feature spaces is minimized, while simultaneously training a task network (e.g., segmentation model). Visual realism, in this case, is not the main concern, as the focus is on optimization for task performance. This problem is often presented as unsupervised domain adaptation, with labeled

synthetic source data being available, but without labels for the target real data. Several feature alignment approaches have been explored, including through minimization of some distance between source and target distributions [78, 79], weight sharing and discriminators to encourage the network to learn domain invariant representations [80, 81], projecting the distance minimization problem to pixel-space to increase the network capacity and preserve semantic content [82], adapting while accounting for cross-domain label imbalances [83], and learning disentangled internal representations that abstract away irrelevant transformations in the target domain [84]. Often, best results have been achieved by combining pixel and feature alignment approaches [85].

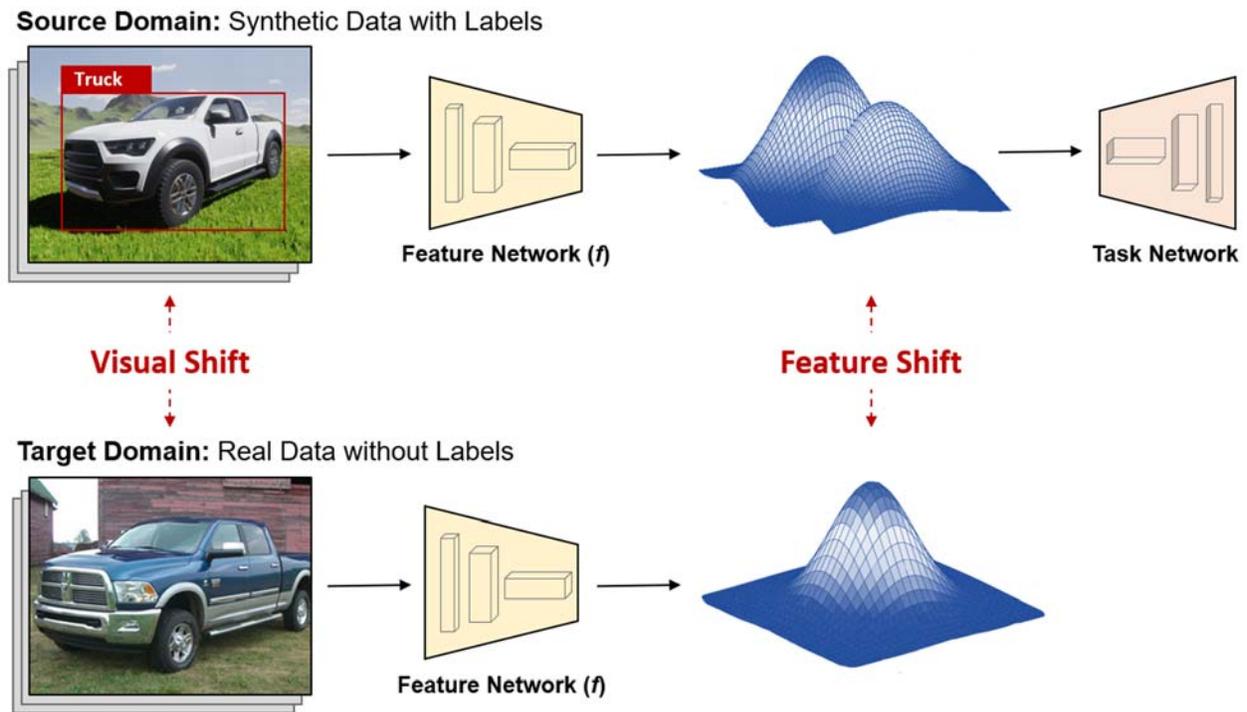


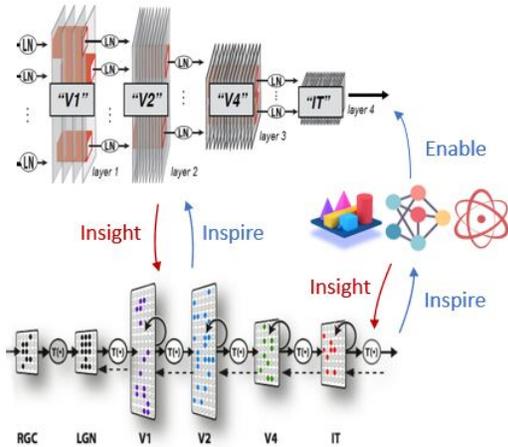
Figure 2. **Domain adaptation at the pixel- and feature-level seeks to close the gap due to visual style and feature distribution shift when moving from synthetic to real data.** In unsupervised domain adaptation, the source domain (synthetic data) is labeled (top row), whereas the target domain (real data) is unlabeled (bottom row). The goal is to close the domain gap by aligning the pixel style of the source to the target domain (i.e., close the visual shift) and learn an embedded representation that is invariant to the domains, while optimizing for a certain downstream task (i.e., close the feature shift).

Enabling the Next Generation of Deep Learning

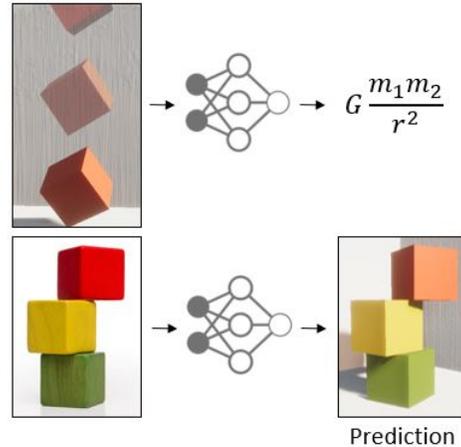
Drawing from cognitive psychology and neuroscience [14], there are several desirable functional and architectural requirements for DNNs. Approaching human-level intelligence

likely requires grasping key concepts related to the physical world and its composition [3, 14, 91], as well as the ability to learn continually, interactively, and multimodally [9, 92]. Here we emphasize the central role synthetic data and simulators play in enabling this next generation of deep learning and, complementary, in providing insight on biological systems (Figure 3).

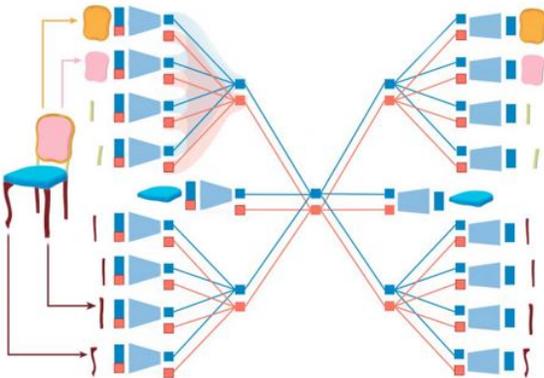
(A) Neural Networks for Scientific Exploration



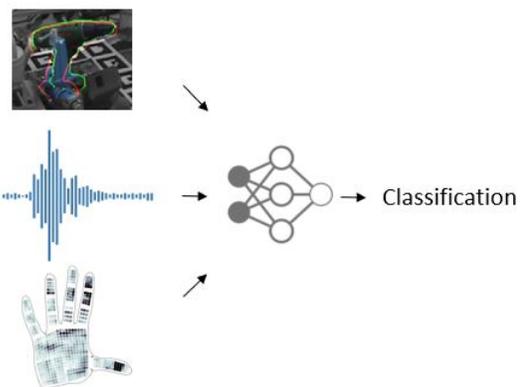
(B) Physics



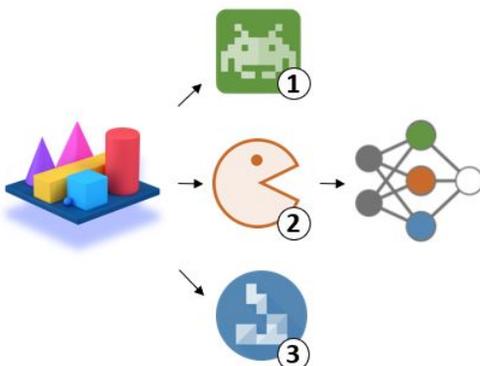
(C) Compositionality



(D) Multimodal Learning



(E) Continual Learning



(F) Embodied Learning

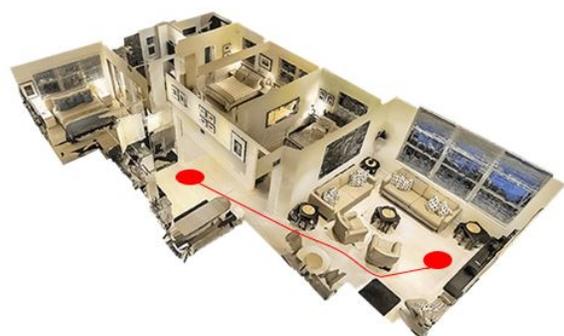


Figure 3. Simulators and synthetic data enable key capability for the next generation of machine learning models. Simulated data and deep neural networks are important tools in the

study of cognitive and neural function in biological systems, enabling exploration of novel (artificial) models for biological function, but also supporting the creation of optimized stimuli to test neural systems (A). Simulated data can be used to learn (known and unknown) physics and make physically realistic predictions (B). Simulated data can be used to teach latent disentangled representations that capture the structure and shape of objects (C). Simulated audio, visual, and haptic data provide redundancy and complementarity that lead to more generalizable internal representations (D). Simulators are ideal for generating a sequence of tasks, which are neither too narrow nor disjoint, to support continual learning without forgetting prior tasks (E). Open-ended interactive exploration in simulators can enable the kind of embodied learning seen in biological systems (F).

Deep Learning for Scientific Exploration

There is a long history of drawing from artificial intelligence to further theory in cognitive psychology and neuroscience [10, 14]. However, DNNs are gaining increasing attention as models of cognitive and neural function due to their ability to learn complex behavior from low-level sensory input, such as image pixels [11]. Some recent successes include predicting behavior and neural activity in perception [86, 87] and memory [88] systems. The benchmarking of different DNNs with respect to how well they predict brain activity allows scientists to test current theory and formulate novel hypotheses about cognitive and neural functioning in biological systems [12]. Given that simulators are able to systematically re-create environmental conditions to test different learning processes and dynamics in DNNs, they are a key enabling technology to the study of biological systems (Figure 3A). Progress in techniques to “open up” DNNs and gather insight on the representations embedded in the hidden layers [89] also introduces the opportunity of retrieving novel post hoc explanations for the functions modeled in DNNs. Here too, synthetic data can be useful to generate stimuli that targets portions of the neural networks to study its function and formulate explanations. For instance, synthetic data has been used to generate optimal stimuli to activate neural subpopulations in primate cortical regions [13]. Synthetic data and simulators also support sophisticated comparison of artificial and biological systems. The use of synthetic materials to study biological systems has a long tradition in cognitive psychology, neuroscience, and artificial intelligence [3, 90]. However, the increasing sophistication and realism of current approaches, as reviewed above, affords novel possibilities. For instance, fully immersive environments (e.g., virtual reality) could support direct comparison of behavior and neural activation in DNNs vs. humans vs. nonhuman primates in embodied interactive tasks. The following subsections review how synthetic data and simulators can help get insight and model key cognitive function.

Physics

At a very young age, humans have a basic understanding of the physical world, such as notions of what constitutes an object and expectations about how they interact with the environment [3, 91]. This knowledge enables mental models about the composition of the world

and predictions about what may happen next [14]. Endowing deep learning models with this type of knowledge, thus, would support more sophisticated learning and inference. Recent work attempts to teach neural networks about physics on the fly (Figure 3B, top). One approach is to train the network with representative examples of the physical domain (e.g., collapsing block towers) and rely on standard learning algorithms to implicitly learn physical knowledge [93]. Another consists in developing specialized architectures that learn physical laws governing the domain (e.g., ball trajectories) [94]. In either case, simulators with an appropriate physics engine are often used to generate the training data [93, 94]. Human intuitive physics have also been argued to rely on a mental simulator to make predictions about the world [4]. In this paradigm, a physics simulator can be explicitly used to make relevant predictions and the challenge, then, is transformed into perceiving the environment (e.g., constructing a scene graph representation through inverse rendering [95]) and feeding that information to the simulator (Figure 3A, bottom). Finally, researchers have also begun embedding physical priors into the learning process (e.g., loss objectives that reflect pertinent physical constraints) to improve transfer from synthetic to real domains [96] and to synthesize more realistic 3D models from 2D imagery [72].

Compositionality

One way to address the complexity of modeling complex synthetic environments is through compositionality. Scenes in the world are decomposable into stable entities, animate (like people or animals) or inanimate (like vehicles or furniture), that we can generically call objects. Such a decomposition is useful, because objects recur in scenes in many different arrangements, but maintain their appearance, properties, and functionality. Similarly, objects themselves consist of parts that have a simpler structure and are often shared across related semantic categories (e.g., chairs, tables, beds, etc. all can have legs). There is a long history in computer vision [97] and computer graphics (e.g., scene graphs) demonstrating the utility of exploiting compositionality in both analysis and synthesis tasks, static and dynamic. Compositional representations are typically hierarchical groupings, based on many possible criteria, including spatial proximity, symmetry, causality, functionality, and others [98, 99], related to principles studied in Gestalt psychology. In synthesis settings, such hierarchies provide natural scaffolds for editing operations, allowing convenient manipulation respecting the semantics of the object or scene – and facilitating the generation of multiple variations.

The machine learning era has created the need for annotated compositional data. In the 3D object domain, datasets that provide fine-grained part decompositions have begun emerging, where objects are mapped into manually curated hierarchies [100]. Hierarchical neural nets, as well as hierarchical convolutional graph networks, have been used in the synthetic generation of objects and scenes, incorporating joint structure and geometry synthesis [101-103]. Scenes naturally exhibit more compositional variability than objects, as many of their constituent entities are mobile or movable. In that setting probabilistic formulations make sense, supporting rule statistics for a generative probabilistic scene grammar to be learned from data [69] and grammar productions themselves to be inferred [104]. Many variations are possible and a recent survey of generative 3D models for objects and scenes is available [105]. As with all generative

approaches, defining appropriate losses for assessing the quality of the generated compositions remains a challenge. Generative compositional models can be conditioned on partial scans, images, or even language. Ideally, one looks for low-dimensional parametrizations of compositional variability with disentangled parameters. Composition often reflects function and structured models can be useful in simulators, with either real or qualitative physics. It has also been suggested that compositionality will be a key attribute in building machines that think more like people [3].

Multimodal Learning

Humans experience and learn about the world through multiple senses, including vision, hearing, and touch [9]. The ability of multiple sensory neural structures to participate in the same function [106] enables redundancy and self-supervision in learning. Redundancy pertains to the ability to learn to perform a task using different modalities (e.g., vision or touch to grasp an object). Self-supervision pertains to the ability of different sensory systems to educate each other about performing a task (e.g., visual-haptic feedback to reach for an object inside a transparent container). In deep learning systems, it is also possible to use redundant and complementary information from multiple data sources to learn more robust and generalizable concept representations [107]. This is perhaps best exemplified by models that integrate audio and visual information, which often co-occur in nature, and learn correspondences that enable predictions on visual tasks from audio information [7] and vice-versa [108]. Recently, haptic information was further shown to be useful for learning features that are pertinent to visual recognition tasks [109]. Given its relevance to building robust robotics and autonomous systems, there has also been considerable interest in merging RGB camera information with complementary sensors, such as depth, LiDAR, and infrared [110]. However, multi-modal sensor data often requires alignment or registration, a non-trivial task. Synthetic data generation can mostly alleviate the need for data alignment as data generation is under our control. This motivated the simulation of various sensor modalities, often by enhancing rendering pipelines with specialized physics engines [111, 112]. A recent promising trend is to develop open-ended simulators that support multimodal training (e.g., physically realistic audio-visual data), as well as explorative incremental learning [113] (more on this in the ‘Embodied Learning’ subsection).

Continual Learning

Humans and animals are remarkably apt at adapting to a changing environment and learning continuously [8, 114]. Replicating this capability in deep learning models would support learning of a potentially infinite series of tasks (e.g., detection of a growing number of categories). Therefore, researchers have explored several mechanisms to support this type of learning, often taking inspiration from biological systems. One approach prevents older tasks from being forgotten by protecting the weights relevant to those tasks [115], similarly to synaptic plasticity mechanisms in biological brains. Another approach integrates memory systems to support replay and episodic memory of relevant prior information [116]. Yet another approach replicates modularity in the brain, often achieved through interactive expansion of the network

parameters while simultaneously trying to meet sparsity constraints [117]. A common challenge to these methods is the need to have a representative set of tasks, which is neither too narrow nor disjoint, to train continual learning algorithms. Due to the difficulty of collecting these training sets from the real world, researchers have often resorted to simulated tasks, such as the Atari game suite [67] and robot manipulation tasks [118], to train these algorithms. A natural extension is the development of simulated open-ended environments [92, 113] that would not only enable lifelong, but structured [119], continual learning. A practical consideration in this setting is designing computationally efficient data generation given the extended training timelines [120].

Embodied Learning

Exploration is essential for human learning. Babies acquire foundational knowledge about the compositionality and the affordances of the physical world through free play with objects in their environments [3, 91]. This interactive engagement leads to rich time-locked correlated visual, haptic, and auditory feedback that contributes to the formation of general internal representations of concepts. The idea that aspects of human intelligence are grounded and emerge from embodied interaction with the world has been associated with learning of basic concepts (e.g., intuitive physics [93]) but also sophisticated symbolic systems (e.g., language [121]). Consequently, researchers noted that, in contrast to training from static datasets that capture moments in time, interactive explorative learning could lead deep learning systems to acquire more robust and generalizable representations of objects, actions, and functions [122, 123]. This paradigm shift calls for datasets that, rather than capturing the world from a third-person perspective, represent first-person experiences. Whereas datasets have started emerging to support embodied learning [124, 125], collecting this type of data is particularly labor intensive [126]. Accordingly, researchers have started developing open-ended physically realistic simulated environments that provide multimodal feedback [92, 113].

The notion of embodied learning implies, at a fundamental level, knowledge about the three-dimensional properties of the world. To understand, for instance, how to interact with a novel object it is necessary to understand its 3D affordances [122]. Whereas this information is readily available in simulators, there is also research in **inverse rendering** that tries to retrieve this information directly from 2D imagery [95, 127, 128]. Reconstructing 3D shapes, however, requires training data with multiple views of the target object or scene, which is seldom available in practice. In promising recent work, though, researchers attempt to automatically retrieve, or disentangle, implicit 3D information from the latent space in GANs [129]. Simulators, inverse rendering, and **latent space disentanglement** techniques, therefore, establish a comprehensive foundation to enabling embodied learning in deep learning models.

Concluding Remarks and Future Challenges

The next generation of deep neural networks will be able to learn rich models of the world in continual, multimodal, and embodied fashion, matching cognitive capability only seen in biological systems. Simulators and synthetic data will play a central role in this transformation. The current generation of deep learning models is limited by access to high

quality training data. This challenge will only be exacerbated due to increased scrutiny of data privacy and security practices. Current research shows that synthetic data can be successfully leveraged to train deep learning models, especially when used in conjunction with domain adaptation techniques that align, statistically at the pixel and feature levels, synthetic and real data. This trend is bound to become more pervasive as it is becoming easier to synthesize realistic data due to impressive advancement in computer graphics rendering pipelines, generative adversarial models, and fusion models. However, beyond meeting current demands for data, synthetic data will meet novel demands. Open-ended interactive multimodal simulation will shift the training paradigm from static datasets usually from a third-person perspective to first-person embodied experiences datasets, which are difficult to collect in the real world. Integration of the synthesis and learning pipelines will support continuous life-long structured learning more similarly to how humans learn and, thus, likely to produce richer, robust, and generalizable knowledge about the world. The paradox of using synthetic data to model natural forms of learning may thus, in practice, be no paradox at all.

Several open issues, nevertheless, remain with respect to synthesizing data that is optimal for deep learning models (see Outstanding Questions). From a modeling perspective, it is essential to assess how similar is the learning and decision process in DNNs when compared to biological systems. Performance on existent datasets may provide insight on the model's predictive ability, but the explanation for the prediction can be obscure. Progress in techniques to dissect and visualize the internal representations of DNNs [89] will likely play an essential role in retrieving these explanations. Furthermore, synthetic data can be systematically created (e.g., with increasing levels of complexity) precisely to study how internal representations are built. Synthetic data is also ideal for exploration [10], not only allowing creation of stimuli to study brain behavior [86, 87], but also to create stimuli (e.g., virtual environments) for sophisticated interactive comparison of behavioral outcome and neural activation of artificial vs. biological systems. From a practical perspective, prior to deploying DNNs in the real world, one needs to provide some assurances that systems built using synthetic data will perform close to systems that were built using data collected by real sensors. Such assurances will require theoretically sound metrics for synthetic data quality that go beyond subjective impressions (e.g., "looks good") and performance on benchmark datasets. Considerable investment has been made developing simulators for mainstream domains (e.g., driving), yet another practical difficulty is that there are still no sophisticated simulators for other, perhaps more complex, domains (e.g., social interaction). Nevertheless, good progress is being made in furthering these types of simulations (e.g., cognitive models of emotion and social expression [130]), as well as using simulated environments to facilitate the collection of data for these domains (e.g., virtual environments to study social interaction [90]). A more fundamental challenge, however, may be whether people will trust and adopt systems trained exclusively, or mostly, with synthetic data – e.g., would people trust a self-driving car that was trained on simulators? It is, therefore, important to understand the differences, not only in terms of performance, but in terms of representation in feature space between models trained with synthetic vs. real data. Here too,

visualization and dissection techniques [89] will likely play a crucial role in explaining how synthetic networks work and help build trust. Nevertheless, despite these challenges, synthetic data introduces a unique opportunity that is worth exploring to enable a new generation of deep learning models that are not limited by available data, but by our imagination alone.

Glossary

Computer graphics rendering pipeline: a sequence of algorithms that produces a three-dimensional visualization in screen space from parameters that describe the scene, such as object 3D and material information, lighting and camera properties, and rendering parameters.

Continual learning: the process of learning a sequence of tasks without forgetting about how to perform earlier tasks in the sequence.

Domain randomization: techniques to create a dataset that is diverse and broadly representative of a target domain, for the purpose of increasing the robustness and generalizability of a deep learning model.

Domain adaptation: techniques that seek to align the statistical properties across domains (e.g., synthetic and real), so that deep learning models training in one domain can be deployed in another.

Domain shift: change in the domain distribution that occurs when a deep learning model is trained with data from one domain (e.g., synthetic) and tested on another (e.g., real).

Embodied learning: the process of learning from multimodal information obtained through interactive exploration of the environment.

Hybrid models: deep learning models trained with a mix of real and synthetic data.

Inverse rendering: The process of automatically retrieving, from 2D imagery, scene attributes such as 3D object information, lighting properties, and camera parameters.

Latent space disentanglement: techniques that seeks to learn a lower dimension representation (e.g., letter category, rotation, and color) of a high dimension space (e.g., images of letters) to support classification and generation in the lower dimension space.

Multimodal learning: the process of learning knowledge from time-locked synchronized information from multiple sensors, such as audio, visual, and haptic input.

Neural rendering: controlled rendering of 3D realistic imagery using deep learning models. In contrast to traditional rendering pipelines, neural rendering pipelines are differentiable and can acquire 3D and physics knowledge from 2D training data.

Simulators: software that is able to generate, often in real time, data and ground-truth annotation from metadata for deep learning models, such as three-dimensional imagery of a scene and segmentation masks.

Synthetic data: data used to train and test deep learning models that is created by artificial means, such as by rendering pipelines, GANs, and fusion models.

Acknowledgments

This research was supported by the US Army. The content does not necessarily reflect the position or the policy of any Government, and no official endorsement should be inferred. We would like to thank Raghuvveer Rao, Achuta Kadambi, and Judy Hoffman for their insightful comments in the preparation of this manuscript.

References

1. LeCun, Y. *et al.* (2015) Deep learning. *Nature* 521, 436-444.
2. Rao, R. *et al.* (2021) Synthetic environments for artificial intelligence (AI) and machine learning (ML) in multi-domain operations. DEVCOM Army Research Laboratory; 2021 May. Report No.: ARL-TR-9198.
3. Lake, B. *et al.* (2016). Building machines that learn and think like people. *Behavioral and Brain Sciences* 40, E253.
4. Battaglia, P. *et al.* (2013) Simulation as an engine of physical scene understanding. *Proceedings of the National Academy of Sciences U.S.A.* 5, 18327-18332.
5. Xu, F. and Tenenbaum, J. Word learning as Bayesian inference. *Psychological Review* 114, 245-272.
6. Barlow, H. (1989) Unsupervised learning. *Neural Computation* 1, 295-311.
7. Owens, A. *et al.* (2016). Ambient sound provides supervision for visual learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
8. Hadsell, R. *et al.* (2020) Embracing change: Continual learning in deep neural networks. *Trends in Cognitive Sciences* 24, 1028-1040.
9. Smith, L. and Gasser, M. (2005) The development of embodied cognitions: Six lessons from babies. *Artificial Life* 11, 13-29.
10. Cichy, M. and Kaiser, D. (2019) Deep neural networks as scientific models. *Trends in Cognitive Sciences* 25, 305-317.
11. Saxe, A. *et al.* (2021) If deep learning is the answer, what is the question? *Nature Reviews Neuroscience*, 22, 55-67.
12. Schrimpf, M. *et al.* (2020). Integrative benchmarking to advance neurally mechanistic models of human intelligence. *Neuron*, 108.
13. Bashivan, P. *et al.* (2019). Neural population control via deep image synthesis. *Science* 364: eaav9436.
14. Hassabis, D. *et al.* (2017) Neuroscience-inspired artificial intelligence. *Neuron* 95, 245-258.

15. Alhaija, H. *et al.* (2018) Augmented reality meets computer vision: Efficient data generation for urban driving scenes. *International Journal of Computer Vision* 126, 961-972.
16. Prakash, A. *et al.* (2019) Structured domain randomization: Bridging the reality gap by context-aware synthetic data. In *Proceedings of the International Conference on Robotics and Automation (ICRA)*.
17. Tremblay, J. *et al.* (2018) Training deep networks with synthetic data: Bridging the reality gap by domain randomization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*.
18. Gaidon, A. *et al.* (2016) Virtual worlds as proxy for multi-object tracking analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
19. Muller, M. *et al.* (2018) Sim4cv: A photo-realistic simulator for computer vision applications? *International Journal of Computer Vision* 126, 902-919.
20. Movshovitz-Attias, Y. *et al.* (2016) How useful is photorealistic rendering for visual learning? In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*.
21. Ros, G. *et al.* (2016) The SYNTHIA dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
22. Richter, S. *et al.* (2016) Playing for data: Ground truth from computer games. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
23. Shafaei, A. *et al.* (2016) Play and learn: Using video games to train computer vision models. In *Proceedings of the British Machine Vision Conference (BMVC)*.
24. Sizikoval, E. *et al.* (2016) Enhancing place recognition using joint intensity - depth analysis and synthetic data. In *Proceedings of the European Conference on Computer Vision (ECCV), VARVAI Workshop*.
25. Wijmans, E. *et al.* (2019) Embodied question answering in photorealistic environments with point cloud perception. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
26. Bousmalis, K. *et al.* (2018) Using simulation and domain adaptation to improve efficiency of deep robotic grasping. In *Proceedings of the International Conference on Robotics and Automation (ICRA)*.
27. Saxena, A. *et al.* (2008) Robotic grasping of novel objects using vision. *International Journal of Robotics Research* 27, 157-173.
28. Tobin, R. *et al.* (2017) Domain randomization for transferring deep neural networks from simulation to the real world. In *Proceedings of the International Conference on Intelligent Robots and Systems (IROS)*.

29. Hattori, H. *et al.* (2018) Augmented reality meets computer vision: Efficient data generation for urban driving scenes. *International Journal of Computer Vision* 126, 1027-1044.
30. Ionescu, C. *et al.* (2014) Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36, 1325-1339.
31. Shotton, J. *et al.* (2013) Efficient human pose estimation from single depth images. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35, 2821-2840.
32. Shrivastava, A. *et al.* (2017) Learning from simulated and unsupervised images through adversarial training. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
33. Dosovitskiy, A. *et al.* (2015) FlowNet: Learning optical flow with convolutional networks. In *Proceedings of the International Conference on Computer Vision (ICCV)*.
34. de Melo, C. *et al.* (2020) Vision-based gesture recognition in human-robot teams using synthetic data. In *Proceedings of the International Conference on Intelligent Robots and Systems (IROS)*.
35. Nikolenko, S. (2019) Synthetic data for deep learning. *ArXiv*, 1909.11512.
36. Planche, B. *et al.* (2017) DepthSynth: Real-time realistic synthetic data generation from CAD models for 2.5d recognition. In *Proceedings of the International Conference on 3D Vision (3DV)*.
37. Tsirikoglou, A. *et al.* (2017) Procedural modeling and physically based rendering for synthetic data generation in automotive applications. *CoRR*, abs/1710.06270.
38. Wrenninge, M. and Unger, J. (2018) Synscapes: A photorealistic synthetic dataset for street scene parsing. *ArXiv*, 1810.08705.
39. Zhang, Y. *et al.* (2016) Physically-based rendering for indoor scene understanding using convolutional neural networks. *ArXiv*, 1612.07429.
40. Howard, A. *et al.* (2017) Mobilenets: Efficient convolutional neural networks for mobile vision applications. *ArXiv*, 1704.04861.
41. Hu, G. *et al.* (2018) Frankenstein: Learning deep face representations using small data. *IEEE Transactions on Image Processing* 27, 293-303.
42. Lopez, A. *et al.* (2017) From virtual to real world visual perception using domain adaptation-The DPM as example, In: Csurka, G. (ed.) *Domain Adaptation in Computer Vision Applications*, pp. 243-258. Springer International Publishing, Cham.
43. Mayer, N. *et al.* (2018) What makes good synthetic training data for learning disparity and optical flow estimation? *International Journal of Computer Vision* 126, 942-960.
44. Goodfellow, I. *et al.* (2014) Generative adversarial networks. In *Proceedings of the Neural Information Processing (NIPS)*.
45. Liu, M.Y. *et al.* (2020) Generative adversarial networks for image and video synthesis: Algorithms and applications. *ArXiv*, 2008.02793.

46. Karras, T. *et al.* (2019) A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
47. Miyato, T. *et al.* (2018) Spectral normalization for generative adversarial networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
48. Brock, A. *et al.* (2019) Large scale GAN training for high fidelity natural image synthesis. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
49. Azadi, S. *et al.* (2019) Semantic bottleneck scene generation. *ArXiv*, 1911.11357.
50. Mirza, M. and Osindero, S. (2014) Conditional generative adversarial nets. *ArXiv*, 1411.1784.
51. Niemeyer, M. and Geiger, A. (2021) GIRAFFE: Representing scenes as compositional generative neural feature fields. In *Proceedings of Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
52. Zhang *et al.* (2021) DatasetGAN: Efficient labeled data factory with minimal human effort. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
53. Bommasani, R. *et al.* (2021) On the opportunities and risks of foundational models. *ArXiv*, 2108.07258v2.
54. Brown, T. *et al.* (2020) Language models are few-shot learners. *ArXiv*, 2005.14165v4.
55. Radford, A. *et al.* (2021) Learning transferable visual models from natural language supervision. *ArXiv*, 2103.00020v1.
56. Dwibedi, D. *et al.* (2017) Cut, paste and learn: Surprisingly easy synthesis for instance detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
57. Alhaija, H. *et al.* (2018) Augmented reality meets computer vision: Efficient data generation for urban driving scenes. *International Journal of Computer Vision* 126, 961-972.
58. Georgakis, G. *et al.* (2017) Synthesizing training data for object detection in indoor scenes. In *Proceedings of the Robotics: Science and Systems (RSS)*.
59. Hattori, H. *et al.* (2015) Learning scene-specific pedestrian detectors without real data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
60. Fernández, C. *et al.* (2011) Augmenting video surveillance footage with virtual agents for incremental event evaluation. In *Pattern Recognition Letters* 32, 878-889.
61. Remez, T. *et al.* (2018) Learning to segment via cut-and-paste. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
62. Ostyakov, P. *et al.* (2018) SEIGAN: towards compositional image generation by simultaneously learning to segment, enhance, and inpaint. *ArXiv*, 1811.07630.

63. Liang, X. *et al.* (2018) Generative semantic manipulation with mask-contrasting GAN. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
64. Chen, X. *et al.* (2018) Attention-GAN for object transfiguration in wild images. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
65. Mnih, V. *et al.* (2015) Human-level control through deep reinforcement learning. *Nature* 518, 529-533.
66. Botvinick, M. *et al.* (2019) Reinforcement learning, fast and slow. *Trends in Cognitive Sciences* 23, 408-422.
67. Mnih, V. *et al.* (2013) Playing Atari with deep reinforcement learning. In *Proceedings of the NIPS Deep Learning Workshop*.
68. Akkaya, I. *et al.* (2019) Solving Rubik's cube with a robot hand. *ArXiv*, 1910.07113.
69. Kar, A. *et al.* (2019) Meta-sim: Learning to generate synthetic datasets. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
70. Eslami, S., *et al.* (2018) Neural scene representation and rendering. *Science* 360, 1204-1210.
71. Kato, H. *et al.* (2018) Neural 3d mesh renderer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
72. Tewari, A. *et al.* (2020) State of the art on neural rendering. *Computer Graphics Forum* 39, 701-727.
73. Bousmalis, K. *et al.* (2017) Unsupervised pixel-level domain adaptation with generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
74. Shrivastava, A., *et al.* (2017) Learning from simulated and unsupervised images through adversarial training. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
75. Zhu, J.-Y. *et al.* (2017) Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the International Conference on Computer Vision (ICCV)*.
76. Park, T. *et al.* (2020) Contrastive learning for unpaired image-to-image translation. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
77. Richter, S. *et al.* (2021) Enhancing photorealism enhancement. *ArXiv*, 2105.04619.
78. Long, M. and Wang, J. (2015) Learning transferable features with deep adaptation networks. In *Proceedings of the International Conference on Machine Learning (ICML)*.
79. Sun, B. and Saenko, K. (2016) Deep CORAL: correlation alignment for deep domain adaptation. In *Proceedings of the ICCV workshop on Transferring and Adapting Source Knowledge in Computer Vision (TASK-CV)*.
80. Liu, M.Y. and Tuzel, O. (2016) Coupled generative adversarial networks. In *Proceedings of Advances in Neural Processing Systems (NeurIPS)*.
81. Tzeng, E. *et al.* (2017) Adversarial discriminative domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

82. Sankaranarayanan, S. *et al.* (2018). Generate to adapt: Aligning domains using generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
83. Prabhu, V. *et al.* (2021) SENTRY: Selective entropy optimization via committee consistency for unsupervised domain adaptation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
84. Litany, O. *et al.* (2021) Representation learning through latent canonicalizations. In *Proceedings of the Winter Conference on Applications of Computer Vision (WACV)*.
85. Hoffman, J. *et al.* (2017) CyCADA: Cycle-consistent adversarial domain adaptation. In *Proceedings of the International Conference on Machine Learning (ICML)*.
86. Khaligh-Razavi, S.-M. and Kriegeskorte, N. (2014) Deep supervised, but not unsupervised, models may explain IT cortical representation. *PLoS Computational Biology* 10, e1003915.
87. Yamins, D.L.K. *et al.* (2014) Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the National Academy of Sciences U.S.A.* 111, 8619-8624.
88. Lloyd, K. *et al.* (2012) Learning to use working memory: A reinforcement learning gating model of rule acquisition in rats. *Frontiers in Computational Neuroscience*, <https://doi.org/10.3389/fncom.2012.00087>.
89. Bau, D. *et al.* (2020) Understanding the role of individual units in a deep network. In *Proceedings of the National Academy of Sciences U.S.A.* 117, 30071-30078.
90. Blascovich, J. *et al.* (2002) Immersive virtual environment technology as a methodological tool for social psychology. *Psychological Inquiry* 13, 103-124.
91. Spelke, E. and Kinzler, E. (2005) Core knowledge. *Developmental Science* 10, 89-96.
92. Savva, M. *et al.* (2019). Habitat: A platform for embodied AI research. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
93. Lerner *et al.* (2016) Learning physical intuition of block towers by example. In *Proceedings of the International Conference on Machine Learning (ICML)*.
94. Chari, P. *et al.* (2019) Visual physics: Discovering physical laws from videos. *ArXiv*, 1911.11893.
95. Ost, J. *et al.* (2021) Neural scene graphs for dynamic scenes. In *Proceedings of Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
96. Chen, Z. *et al.* (2021) PSD: Principled synthetic-to-real dehazing guided by physical priors. In *Proceedings of Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
97. Felzenszwalb, P. and Huttenlocher, D. (2005) Pictorial structures for object recognition. *International Journal of Computer Vision* 61, 55-79.
98. Lun, Z. *et al.* (2017) Learning to group discrete graphical patterns. In *Proceedings of SIGGRAPH ASIA*.

99. Wang, Y. *et al.* (2011) Symmetry hierarchy of man-made objects. *Computer Graphics Forum* 30, 287-296.
100. Mo, K. *et al.* (2019). PartNet: A large-scale benchmark for fine-grained and hierarchical part-level 3D object understanding. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
101. Li, J. *et al.* (2017) GRASS: Generative recursive autoencoders for shape structures. In *Proceedings of SIGGRAPH*.
102. Mo, K. *et al.* (2020) StructureNet: Hierarchical graph networks for 3D shape generation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
103. Li, M. *et al.* (2019) *GRAINS*: Generative Recursive Autoencoders for INdoor scenes. *ACM Transactions on Graphics* 38, 1-16.
104. J. Devaranjan *et al.* (2020) *Meta-Sim2: Unsupervised learning of scene structure for synthetic data generation*. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
105. Chaudhuri, S. *et al.* (2020) Learning generative models of 3D structures. *Eurographics, STAR State-of-the-Art-Report*.
106. Edelman, G. and Gally, J. (2001) Degeneracy and complexity in biological systems. *Proceedings of the National Academy of Sciences U.S.A* 98, 13763-13768.
107. Baltrušaitis, T. *et al.* (2017) Multimodal machine learning: A survey and taxonomy. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41, 423-443.
108. Arandjelović, R. and Zisserman, A. (2017) Look, listen and learn. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
109. Sundaram, S. *et al.* (2019) Learning the signatures of the human grasp using a scalable tactile glove. *Nature* 569, 698-702.
110. Grigorescu, S. *et al.* (2020) A survey of deep learning techniques for autonomous driving. *Journal of Field Robotics* 37, 362-386.
111. Wu, B. *et al.* (2019) SqueezeSegV2: Improved model structure and unsupervised domain adaptation for road-object segmentation from a LiDAR point cloud. In *Proceedings of the International Conference on Robotics and Automation (ICRA)*.
112. Ahmadibeni, A. *et al.* (2020) Automatic target recognition of aerial vehicles based on synthetic SAR imagery using hybrid stacked denoising auto-encoders. In *Proceedings of the SPIE 11393, Algorithms for Synthetic Aperture Radar Imagery XXVII, 113930J*.
113. Gan *et al.* (2020) ThreeDWorld: A platform for interactive multi-modal physical simulation. *ArXiv*, 2007.04954.
114. Cichon, J. and Gan, W.-B. (2015) Branch-specific dendritic Ca²⁺ spikes cause persistent synaptic plasticity. *Nature* 520, 180-185.
115. Kirkpatrick, J. *et al.* (2017) Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences U.S.A.* 114, 3521-3526.

116. Aljundi, R. *et al.* (2019) Online continual learning with maximally interfered retrieval. In *Proceedings of Advances in Neural Processing Systems (NeurIPS)*.
117. Bellec, G. *et al.* (2018) Deep rewiring: Training very sparse deep networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
118. Yu, T. *et al.* (2021) Meta-World: A benchmark and evaluation for multi-task and meta learning. *ArXiv*, 1910.10897.
119. Bengio, Y. *et al.* (2009) Curriculum learning. In *Proceedings of the International Conference on Machine Learning (ICML)*.
120. Mason, K. *et al.* (2019) An “on the fly” framework for efficiently generating synthetic big data sets. In *Proceedings of the IEEE International Conference on Big Data*.
121. McNeill, D. (1992) *Hand and Mind: What gestures reveal about thought*. The University of Chicago Press.
122. Mandikal, P. and Grauman, K. (2021) Learning dexterous grasping with object-centric visual affordances. In *Proceedings of the International Conference on Robotics and Automation (ICRA)*.
123. Nagarajan, T. *et al.* (2021) Ego-Topo: Environment affordances from egocentric video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
124. Sigurdsson, G. *et al.* (2018) Charades-Ego: A large-scale dataset of paired third and first person videos. *ArXiv*, 1804.09626.
125. Damen, D. *et al.* (2018) Scaling egocentric vision: The EPIC-KITCHENS dataset. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
126. Grauman, K. *et al.* (2021) Ego4D: Around the world in 3,000 hours of egocentric video. *ArXiv*, 2110.07058.
127. Wang, N. *et al.* (2018) Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
128. Mescheder, L. *et al.* (2019) Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
129. Zhang, Y. *et al.* (2021) Image GANs meet differentiable rendering for inverse graphics and interpretable neural rendering. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
130. Marsella, S. *et al.* (2010) Computational models of emotion. In K. Scherer., T. Bänziger, and E. Roesch (Eds.). *A blueprint for an affectively competent agent: Cross-fertilization between Emotion Psychology, Affective Neuroscience, and Affective Computing*. Oxford University Press, New York, NY.