



OPEN

A review of some techniques for inclusion of domain-knowledge into deep neural networks

Tirtharaj Dash^{1✉}, Sharad Chitlangia², Aditya Ahuja¹ & Ashwin Srinivasan¹

We present a survey of ways in which existing scientific knowledge are included when constructing models with neural networks. The inclusion of domain-knowledge is of special interest not just to constructing scientific assistants, but also, many other areas that involve understanding data using human-machine collaboration. In many such instances, machine-based model construction may benefit significantly from being provided with human-knowledge of the domain encoded in a sufficiently precise form. This paper examines the inclusion of domain-knowledge by means of changes to: the input, the loss-function, and the architecture of deep networks. The categorisation is for ease of exposition: in practice we expect a combination of such changes will be employed. In each category, we describe techniques that have been shown to yield significant changes in the performance of deep neural networks.

Science is a cumulative enterprise, with generations of scientists discovering, refining, correcting and ultimately increasing our knowledge of how things are. The accelerating pace of development in software and hardware for machine learning—in particular, the area of deep neural networks (DNNs)—inevitably raises the prospect of Artificial Intelligence for Science¹. That is, how can we best use AI methods to accelerate our understanding of the natural world? While ambitious plans exist for completely automated AI-based robot scientists², the immediately useful prospect of using AI for Science remains semi-automated. An example of such a collaborative system is in Fig. 1. For such systems to work effectively, we need at least the following: (1) We have to be able to tell the machine what we know, in a suitably precise form; and (2) The machine has to be able to tell us what it has found, in a suitably understandable form. While the remarkable recent successes of deep neural networks on a wide variety of tasks makes a substantial case for their use in model construction, it is not immediately obvious how either (1) or (2) should be done with deep neural networks. In this paper, we examine ways of achieving (1), that is, the techniques for constructing deep neural networks from data and domain-knowledge concerning the problem. Understanding models constructed by deep neural networks is an area of intense research activity, and good summaries exist elsewhere^{3,4}. To motivate the utility of providing domain-knowledge to a deep network, we reproduce two results from⁵ in Fig. 2, which shows that predictive performance can increase significantly, even with a simplified encoding of domain-knowledge (see Fig. 2a).

It is unsurprising that a recent report on AI for Science¹ identifies the incorporation of domain-knowledge as one of the 3 Grand Challenges in developing AI systems:

“ML and AI are generally domain-agnostic...Off-the-shelf [ML and AI] practice treats [each of these] data-sets in the same way and ignores domain knowledge that extends far beyond the raw data...Improving our ability to systematically incorporate diverse forms of domain knowledge can impact every aspect of AI.”

But it is not just the construction of scientific-assistants that can benefit from this form of man-machine collaboration, and “human-in-the-loop” AI systems are likely to play an increasingly important role in engineering, medicine, healthcare, agriculture, environment and so on⁸. In this survey, we restrict the studies on incorporation of domain-knowledge into neural networks, with 1 or more hidden layers. If the domain-knowledge expressed in a symbolic form (for example, logical relations that are known to hold in the domain), then the broad area of hybrid neural-symbolic systems (see for example,^{9,10}) is clearly relevant to the material in this paper. However, the motivation driving the development of hybrid systems is much broader than this paper, being concerned with general-purpose neural-based architectures for logical representation and inference. Here our goals are more

¹Department of Computer Science and Information Systems, Anuradha and Prashanth Palakurthi Centre for AI Research (APPCAIR), BITS Pilani, K.K. Birla Goa Campus, Goa 403726, India. ²Department of Electrical and Electronics Engineering, Anuradha and Prashanth Palakurthi Centre for AI Research (APPCAIR), BITS Pilani, K.K. Birla Goa Campus, Goa 403726, India. ✉email: tirtharaj@goa.bits-pilani.ac.in

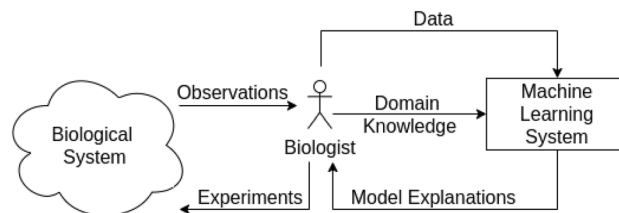


Figure 1. An example of AI for Science. The human-in-the-loop is a biologist. The biologist conducts experiments in a biological system, obtains experimental observations. The biologist then extracts data that can be used to construct machine learning model(s). Additionally, the machine learning system has access to domain knowledge that can be obtained from the biologist. The machine learning system then conveys its explanations to the biologist.

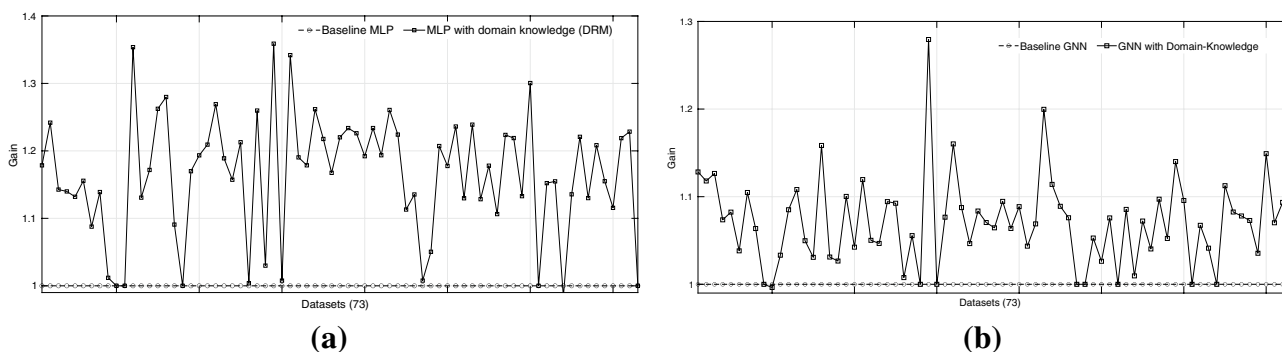


Figure 2. The plots from⁶ showing gains in predictive accuracy of (a) multilayer perceptron (MLP), and (b) graph neural network (GNN) with the inclusion of domain-knowledge. The domain knowledge inclusion method in (a) is a simple technique known as ‘propositionalisation’⁷; and, the method in (b) is a general technique of incorporating domain-knowledge using bottom-graph construction. The results shown are over 70 datasets. No importance to be given to the line joining two points; this is done for visualisation purpose only.

For inhibiting this protein: The presence of a peroxide bridge is relevant. The target site is at most 20Å.	The model should follow that: $p(y = 1 \mathbf{x}) \geq 0.9$ $p(y = 0 \mathbf{x}) \leq 0.1$ Initial weights should be $3n - 2.3$ [13]
(a)	(b)

Figure 3. Informal descriptions of (a) logical; and (b) numerical constraints.

modest: we are looking at the inclusion of problem-specific information into machine-learning models of a kind that will be described shortly. We refer the reader to¹¹ for reviews of work in the broader area of neural-symbolic modelling. More directly related to this paper is the work on “informed machine learning”, reviewed in¹². We share with this work the interest in prior knowledge as an important source of information that can augment existing data. However, the goals of that paper are more ambitious than here. It aims to identify categories of prior knowledge, using as dimensions: the source of the knowledge, its representation, and its point of use in a machine-learning algorithm. In this survey, we are only concerned with some of these categories. Specifically, in terms of the categories in¹², we are interested in implicit or explicit sources of domain-knowledge, represented either as logical or numeric constraints, and used at the model-construction stage by DNNs. Informal examples of what we mean by logical and numeric constraints are shown in Fig. 3. In general, we will assume logical constraints can, in principle, be represented as statements in propositional logic or predicate logic. Numerical constraints will be representable, in principle, as terms in an objective function being minimised (or maximised), or prior distributions on models. We believe this covers a wide range of potential applications, including those concerned with scientific discovery.

Focus of the paper. We adhere to the following informal specification for constructing a deep neural network: given some data D , a structure and parameters of a deep network (denoted by π and θ , respectively), a learner \mathcal{L} attempts to construct a neural network model M that minimises some loss function L . Fig. 4 shows a

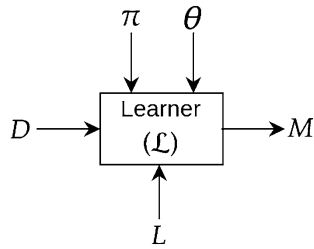


Figure 4. Construction of a deep model M from data (D) using a learner (\mathcal{L}). We use π to denote the structure (organisation of various layers, their interconnections etc.) and θ to denote the parameters (synaptic weights) of the deep network. L denotes the loss function (for example, cross-entropy loss in case of classification).

Arch.	Domain-Knowledge Effect		
	Transform Input Data	Transform Loss Function	Transform Model
MLP	Reformulate feature-representation	(For all architectures) Reformulate regularisation term; Addition of syntactic or semantic constraints with associated penalties; Differential costs of decisions	Changes in layers, hidden units
CNN	Reformulate spatial-representation		As with MLPs, plus changes to connections between units; changes convolution filters
RNN, Transformer	Reformulate sequence-representation		As with MLPs, plus possible changes to attention-mechanism
GNN	Reformulate graph-representation		As with MLPs, plus changes to graph-convolution

Figure 5. Some implications of using domain-knowledge for commonly-used deep network architectures. Although attention-mechanism has also been used recently in many deep network architectures, we mention it only for RNNs and transformers as it is more prominently being used for sequence learning.

diagrammatic representation. Note that: (a) we do not describe how the learner \mathcal{L} constructs a model M given the inputs. But, it would be normal for the learner to optimise the loss L by performing an iterative estimation of the parameters θ , given the model structure π ; and (b) we are not concerned with how the constructed deep model M will be used. However, it suffices to say that when used, the model M would be given one or more data-instances encoded in the same way as was provided for model-construction.

In the literature, domain knowledge—also called background knowledge—does not appear to have an accepted definition, other than that, it refers to information about the problem. This information can be in the form of relevant features, concepts, taxonomies, rules-of-thumb, logical constraints, probability distributions, mathematical distributions, causal connections and so on. In this paper, we use the term “domain-knowledge” to refer to problem-specific information that can directly be translated into alterations to the principal inputs of Fig. 4. That is, by domain-knowledge we will mean problem-specific information that can change: (1) The input data to a deep network; (2) The loss-function used; and (3) The model (that is, the structure or parameters) of the deep network. In a sense, this progression reflects a graded increase in the complexity of changes involved. Figure 5 tabulates the principal implications of this position for commonly-used deep learning architectures.

The rest of the paper is organised as follows: Section 2 describes inclusion of domain-knowledge by augmenting or transformation of inputs; Section 3 describes changes that have been employed to loss-functions; and Section 4 describes biases on parameters and changes to the structure of deep networks. Section 5 outlines some major challenges related to the inclusion of domain-knowledge in the ways we describe. In this section, we also present perspectives on the relevance of the use of domain-knowledge to aspects of Responsible AI, including ethics, fairness, and explainability of DNNs.

Transforming the input data. One of the prominent approaches to incorporate domain-knowledge into a deep network is by changing inputs to the network. Here, the domain-knowledge is primarily in symbolic form. The idea is simple: If a data instance could be described using a set of attributes that not only includes the raw feature-values but also includes more details from the domain, then a standard deep network could then be constructed from these new features. A simple block diagram in Fig. 6 shows how domain knowledge is introduced into the network via changes in inputs. In this survey, we discuss broadly two different ways of doing this: (a) using relational features, mostly constructed by a method called propositionalisation⁷ using another machine learning system (for example, Inductive Logic Programming) that deals with data and background knowledge; (b) without propositionalisation.

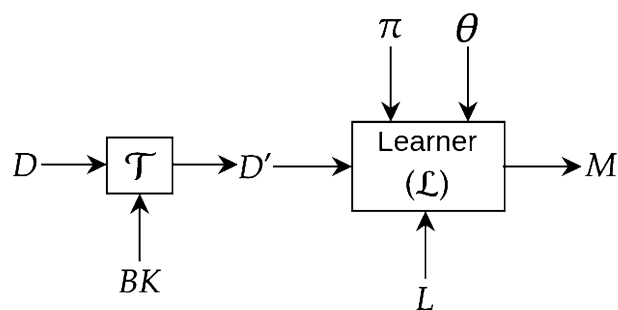


Figure 6. Introducing background knowledge into deep network by transforming data. \mathcal{T} is a transformation block that takes input data D , background knowledge (BK) and outputs transformed data D' that is then used to construct a deep model using a learner \mathcal{L} .

Propositionalisation. The pre-eminent form of symbolic machine learning based on the use of relations in first-order logic is Inductive Logic Programming (ILP)¹⁴, which has an explicit role for domain-knowledge being incorporated into learning. The simplest use of ILP¹⁴ to incorporate n -ary relations in domain knowledge into a neural network relies on techniques that automatically “flatten” the domain-knowledge into a set of domain-specific relational features. Although not all DNNs require data to be a set of feature-vectors, this form of data representation is long-standing and still sufficiently prevalent. In logical terms, we categorise feature-based representations as being encodings in propositional logic. The reader would point out, correctly, that feature-values may not be Boolean. This is correct, but we can represent non-Boolean features by Boolean-valued propositions (for example, a real-valued feature f with value 4.2 would be represented by a corresponding Boolean feature f' that has the value 1 if $f = 4.2$ and 0 otherwise). With the caveat of this rephrasing, it has of course been possible to provide domain-knowledge to neural networks by employing domain-specific features devised by an expert. However, we focus here on ways in which domain-knowledge encoded as rules in propositional and first-order logic has been used to construct the input features for a deep neural network. Techniques for automatic construction of Boolean-valued features from relational domain-knowledge have a long history in the field of ILP^{15–17}, originating from the LINUS⁷. Often called *propositionalisation*, the approach involves the construction of features that identify the conditions under which they take on the value 1 (or 0). For example, given (amongst other things) the definition of benzene rings and of fused rings, an ILP-based propositionalisation may construct the Boolean-valued feature that has the value 1 if a molecule has 3 fused benzene rings, and 0 otherwise. The values of such Boolean-valued features allows us to represent a data instance (like a molecule) as a Boolean-valued feature-vector, which can then be provided to a neural network. There is a long history of propositionalisation: see¹⁸ for a review of some of early use of this technique, and^{19,20} who examine the links between propositionalisation and modern-day use of embeddings in deep neural networks. More clearly, the authors examine that both propositionalisation and embedding approaches aim at transforming data into tabular data format, while they are being used in different problem settings and contexts. One recent example of embedding is demonstrated in²¹ where the authors use different text-embedding approaches such as sentence encoder²² and GPT2²³ to transform textual domain-knowledge into embedding vectors.

A direct application of propositionalisation, demonstrating its utility for deep networks has been its use in Deep Relational Machines (DRMs)²⁴. A DRM is a deep fully-connected neural network with Boolean-valued inputs obtained from propositionalisation by an ILP engine. In²⁵, Boolean-valued features from an ILP engine are sampled from a large space of possible relational features. The sampling technique is refined further in²⁶.

The idea of propositionalisation also forms the foundation for a method known as ‘Bottom Clause Propositionalisation (BCP)’ to propositionalise the literals of a most-specific clause, or “bottom-clause” in ILP. Given a data instance, the bottom-clause is the most-specific first-order clause that entails the data instance, given some domain-knowledge. Loosely speaking, the most-specific clause can be thought of “enriching” the data instance with all domain relations that are true, given the data instance. The construction of such most-specific clauses and their subsequent use in ILP was introduced in²⁷. CILP++²⁸ uses bottom-clauses for data instances to construct feature-vectors for neural networks. This is an extension to CIL²P in which the neural network uses recurrent connections to enforce the background-knowledge during the training²⁹.

Propositionalisation has conceptual and practical limitations. Conceptually, there is no variable-sharing between two or more first-order logic features²⁵. That is, all useful compositions have to be pre-specified. Practically, this makes the space of possible features extremely large: this has meant that the feature-selection has usually been done separately from the construction of the neural network. In this context, another work that does not employ either propositionalisation or network augmentation considers a combination of symbolic knowledge represented in first-order logic with matrix factorization techniques³⁰. This exploits dependencies between textual patterns to generalise to new relations.

Recent work on neural-guided program synthesis also explicitly includes domain-specific relations. Here programs attempt to construct automatically compositions of functional primitives. The primitives are represented as fragments of functional programs that are expected to be relevant. An example of neural-guided program synthesis that uses such domain-primitives is DreamCoder^{31,32}. DreamCoder receives as inputs, the partial specification of a function in the form of some inputs–output pairs, and a set of low-level primitives

represented in a declarative language. Higher-level domain-concepts are then abduced as compositions of these primitives via a neurally-guided search procedure based on a version of the Bayesian “wake-sleep” algorithm³³. The deep networks use a (multi-hot) Boolean-vector encoding to represent functional compositions (a binary digit is associated with each primitive function, and takes the value 1 if and only if the primitive is used in the composite function).

There are methods that do not use an explicit propositionalisation step, but nevertheless amount to re-formulating the input feature-representation. In the area of “domain-adaptation”³⁴, “source” problems act as a proxy for domain-knowledge for a related “target” problem (Superficially, this is also the setting underlying *transfer learning*. However, the principal difference is that source and target problems are closely related in domain-adaptation, but this need not be the case with transfer-learning. Transfer-learning also usually involves changes in both model-parameters and model-structure. Domain-adaptation does not change the model-structure: we consider these points in a later section). There is a form of domain-adaptation in which the target’s input representation is changed based on the source model. In³⁵, for example, a feature-encoder ensures that the feature representation for the target domain that is the same as the one used for the source.

Binary and n -ary relations. An influential form of representing relational domain-knowledge takes the form *knowledge graph*, which are labelled graphs, with vertices representing entities and edges representing binary relations between entities. A knowledge graph provides a structured representation for knowledge that is accessible to both humans and machines³⁶. Knowledge graphs have been used successfully in variety of problems arising in information processing domains such as search, recommendation, summarisation³⁷. Sometimes the formal semantics of knowledge graphs such as domain ontologies are used as sources for external domain-knowledge³⁸. We refer the reader to³⁹ to a comprehensive survey of this form of representation for domain-knowledge.

Incorporation of the information in a knowledge-graph into deep neural models—termed “knowledge-infused learning”—is described in^{40,41}. This aims to incorporate binary relations contained in application-independent sources (like DBPedia, Yago, WikiData) and application-specific sources (like SNOMED-CT, DataMed). The work examines techniques for incorporating relations at various layers of deep-networks (the authors categorise these as “shallow”, “semi-deep” and “deep” infusion). In the case of shallow infusion, both the external knowledge and the method of knowledge infusion are shallow, utilising syntactic and lexical knowledge in word embedding models. In semi-deep infusion, external knowledge is involved through attention mechanisms or learnable knowledge constraints acting as a sentinel to guide model learning. Deep infusion employs a stratified representation of knowledge representing different levels of abstractions in different layers of a deep learning model to transfer the knowledge that aligns with the corresponding layer in the learning process. Fusing the information in a knowledge-graph in this way into various level of hidden representations in a deep network could also allow quantitative and qualitative assessment of its functioning, leading to knowledge-infused interpretability⁴².

There have been some recent advances in introducing external domain-knowledge into deep sequence models. For instance, in³⁸, the authors incorporate domain-specific knowledge into the popular deep learning framework, BERT⁴³ via a declarative knowledge source like drug-abuse ontology. The model constructed here, called Gated-K-BERT, is used for jointly extracting entities and their relationships from tweets by introducing the domain-knowledge using an entity position-aware module into the primary BERT architecture. The experimental results demonstrate that incorporating domain-knowledge in this manner leads to better relation extraction as compared to the state-of-the-art. This work could fall within the category of semi-deep infusion as described in^{40,44}, in their study on learning from electronic health records show that the adjacency information in a medical knowledge graph can be used to model the attention mechanism in an LSTM-based RNN with attention. Whenever the RNN gets an entity (a medical event) as an input, the corresponding sub-graph in the medical knowledge graph (consisting of relations such as *causes* and *is-caused-by*) is then used to compute an attention score. This method of incorporating the medical relations into the RNN falls under the category of semi-deep knowledge infusion. While the above methods use the relational knowledge from a knowledge-graph by altering or adding an attention module within the deep sequence model, a recent method called KRISP⁴⁵ introduces such external knowledge at the output (prediction) layer of BERT. This work could be considered under the category of shallow infusion of domain-knowledge as characterised by⁴⁰.

Knowledge graphs can be used directly by specialised deep network models that can handle graph-based data as input (graph neural networks, or GNNs). The computational machinery available in GNN then aggregates and combines the information available in the knowledge graph (an example of this kind of aggregation and pooling of relational information is in⁴⁶). The final collected information from this computation could be used for further predictions. Some recent works are in^{47,48}, where a GNN is used for estimation of node importance in a knowledge-graph. The intuition is that the nodes (in a problem involving recommender systems, as in⁴⁸, a node represents an entity) in the knowledge-graph can be represented with an aggregated neighbourhood information with bias while adopting the central idea of aggregate-and-combine in GNNs. The idea of encoding a knowledge graph directly for a GNN is also used in Knowledge-Based Recommender Dialog (KBRD) framework developed for recommender systems⁴⁹. In this work, the authors treat an external knowledge graph, such as DBpedia⁵⁰, as a source of domain-knowledge allowing entities to be enriched with this knowledge. The authors found that the introduction of such knowledge in the form of a knowledge-graph can strengthen the recommendation performance significantly and can enhance the consistency and diversity of the generated dialogues. In KRISP⁴⁵, a knowledge-graph is treated as input for a GNN where each node of the graph network corresponds to one specific domain-concept in the knowledge graph. This idea is a consequence of how a GNN operates: it can form more complex domain-concepts by propagating information of the basic domain-concepts along the edges in the knowledge-graph. Further, the authors allow the network parameters to be shared across the domain-concepts

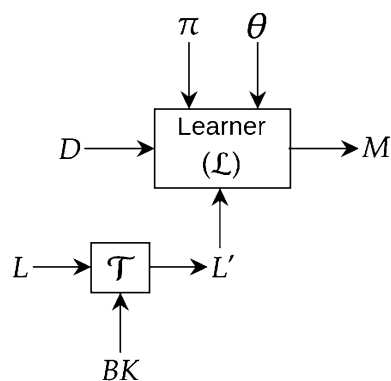


Figure 7. Introducing background knowledge into deep network by transforming the loss function L . \mathcal{T} block takes an input loss L and outputs a new loss function L' by transforming (augmenting or modifying) L based on background knowledge (BK). The learner \mathcal{L} then constructs a deep model using the original data D and the new loss function L' .

with a hope to achieve better generalisation. We note that while knowledge-graph provide an explicit representation of domain-knowledge in the data, some problems contain domain-knowledge implicitly through an inherent topological structure (like a communication network). Clearly, GNNs can accommodate such topological structure just in the same manner as any other form of graph-based relations (see for example:⁵¹).

Going beyond binary relations in knowledge-graphs and treating n -ary relations as hyperedges, a technique called *vertex enrichment* is proposed in⁵. Vertex-enrichment is a simplified approach for the inclusion of domain-knowledge into standard graph neural networks (GNNs). The approach incorporates first-order background relations by augmenting the features associated with the nodes of a graph provided to a GNN. The results reported in⁵ show significant improvements in the predictive accuracy of GNNs across a large number datasets. The simplification used in vertex-enrichment has been made unnecessary in a recent proposal for transforming the most-specific clause constructed by ILP systems employing mode-directed inverse entailment (MDIE²⁷). The transformation converts this clause into a graph can be directly used by any standard GNN⁶. Specifically, the transformation results in a labelled bipartite graph consisting of vertices representing predicates (including domain predicates) and ground terms. This approach reports better predictive performance than those reported in⁵, and includes knowledge-graphs as a special case. Most recently, this method has been combined successfully with deep generative sequence models for generating target-specific molecules, which demonstrates yet another real-world use-case of incorporating domain knowledge into deep networks⁵².

Transforming the loss function

One standard way of incorporating domain-knowledge into a deep network is by introducing “penalty” terms into the loss (or utility) function that reflect constraints imposed by domain-knowledge. The optimiser used for model-construction then minimises the overall loss that includes the penalty terms. Fig. 7 shows a simple block diagram where a new loss term is introduced based on the background knowledge. We distinguish two kinds of domain constraints—syntactic and semantic—and describe how these have been used to introduce penalty terms into the loss function.

Syntactic Loss. The usual mechanism for introducing syntactic constraints is to introduce one or more *regularisation* terms into the loss function. The most common form of regularisation introduces penalties based on model complexity (number of hidden layers, or number of parameters and so on: see, for example,⁵³).

A more elaborate form of syntactic constraints involves the concept of *embeddings*. Embeddings refer to the relatively low-dimensional learned continuous vector representations of discrete variables. Penalty terms based on “regularising embeddings” are used to encode syntactic constraints on the complexity of embeddings.⁵⁴ was an early work in this direction, in which the authors proposed a strategy to establish constraints by designating each node in a Hopfield Net to represent a concept and edges to represent their relationships and learn these nets by finding the solution which maximises the greatest number of these constraints.⁵⁵ was perhaps the first method of regularising embeddings from declarative knowledge encoded in first-order logic. The proposal here is for mapping between logical statements and their embeddings, and logical inferences and matrix operations. That is, the model behaves as if it is following a complex first-order reasoning process, but operates at the level of simple vectors and matrix representations.³⁰ extended this to regularisation by addition of differentiable loss terms to the objective-based on propositionalisation of each first-order predicate. Guo *et al.*⁵⁶ proposed a joint model, called KALE, which embeds facts from knowledge-graphs and logical rules, simultaneously. Here, the facts are represented as ground atoms with a calculated truth value in $[0, 1]$ suggesting how likely that the fact holds. Logical rules (in grounded form) are then interpreted as complex formulae constructed by combining ground atoms with logical connectives, which are then modelled by fuzzy t -norm operators⁵⁷. The truth value that results from this operation is nothing but a composition of the constituent ground atoms, allowing the facts from the knowledge graph to be incorporated into the model.

Li and Srikumar⁵⁸ develop a method to constraint individual neural layers using soft logic based on massively available declarative rules in ConceptNet.⁵⁹ incorporates first-order logic into low dimensional spaces by embedding graphs nodes and represents logical operators as learned geometric relations in the space.⁶⁰ proposed ordering of embedding space based on rules mined from WordNet and found it to better prior knowledge and generalisation capabilities using these relational embeddings.⁶¹ show that domain-based regularisation in loss function can also help in constructing deep networks with less amount of data in prediction problems concerned with cloud computing. In⁶², a knowledge-based distant regularisation framework was proposed that utilises the distance information encoded in a knowledge-graph. It defines prior distributions of model parameters using knowledge-graph embeddings. They show that this results in an optimisation problem for a regularised factor analysis method.

Semantic loss. Penalty terms can also be introduced on the extent to which the model's prediction satisfies semantic domain constraints. For example, the domain may impose specific restrictions on the prediction ("output prediction must be in the range 3 . . . 6"). One way in which such information is provided is in the form of domain-constraints. Penalty terms are then introduced based on the number and importance of such constraints that are violated.

A recent work that is based on loss function is in⁶³. Here the authors propose a semantic loss that signifies how well the outputs of the deep network matches some given constraints encoded as propositional rules. The general intuition behind this idea is that the semantic loss is proportional to a negative logarithm of the probability of generating a state that satisfies the constraint when sampling values according to some probability distribution. This kind of loss function is particularly useful for semi-supervised learning as these losses behave like self-information and are not constructed using explicit labels and can thus utilize unlabelled data.

⁶⁴ proposed a framework to incorporate first-order logic rules with the help of an iterative distillation procedure that transfers the structured information of logic rules into the weights of neural networks. This is done via a modification to the knowledge-distillation loss proposed by Hinton et al.⁶⁵. The authors show that taking this loss-based route of integrating rule-based domain-knowledge allows the flexibility of choosing a deep network architecture suitable for the intended task.

In⁶⁶, authors construct a system for training a neural network with domain-knowledge encoded as logical constraints. Here the available constraints are transferred to a loss function. Specifically, each individual logic operation (such as negation, and, or, equality etc.) is translated to a loss term. The final formulation results in an optimisation problem. The authors extract constraints on inputs that capture certain kinds of convex sets and use them as optimisation constraints to make the optimisation tractable. In the developed system, it is also possible to pose queries on the model to find inputs that satisfy a set of constraints. In a similar line,⁶⁷ proposed domain-adapted neural network (DANN) that works with a balanced loss function at the intersection of models based on purely domain-based loss or purely inductive loss. Specifically, they introduce a domain-loss term that requires a functional form of approximation and monotonicity constraints on the outputs of a deep network. Without detailing much on the underlying equations, it suffices to say that formulating the domain loss using these constraints enforces the model to learn not only from training data but also in accordance with certain accepted domain rules.

Another popular approach that treats domain knowledge as 'domain constraints' is semantic based regularisation^{68,69}. It builds standard multilayered neural networks (e.g. MLP) with kernel machines at the input layer that deal with continuous-valued features. The output of the kernel machines is input to the higher layers implementing a fuzzy generalisation of the domain constraints that are represented in first-order logic. The regularisation term, consisting of a sum of fuzzy generalisation of constraints using t -norm operations, in the cumulative loss then penalises each violation of the constraints during the training of the deep network. In⁷⁰ inject domain knowledge at training time via an approach that combines semantic based regularisation and constraint programming⁷¹. This approach uses the concept of 'propagators', which is inherent in constraint programming to identify infeasible assignments of variables to values in the domain of the variables. The role of semantic-based regularisation is to then penalise these infeasible assignments weighted by a penalty parameter. This is an example of constraints on inputs. In a similar line, In⁷² introduce domain-knowledge into a deep LSTM-based RNN at three different levels: constraining the inputs by designing a filtering module based on the domain-specific rules, constraining the output by enforcing an output range, and also by introducing a penalty term in the loss function.

A library for integrating symbolic domain-knowledge with deep neural networks was introduced recently in⁷³. It provides some effective ways of specifying domain-knowledge, albeit restricted to (binary) hierarchical concepts only, for problems arising in the domain of natural language processing and some subset of computer vision. The principle of integration involves constraint satisfaction using a primal-dual formulation of the optimisation problem. That is: the goal is to satisfy the maximum number of domain constraints while also minimising the training loss, an approach similar to the idea proposed in^{66,67,70}.

While adding a domain-based constraint term to the loss function may seem appealing, there are a few challenges. One challenge that is pointed out in a recent study⁷⁴ is that incorporating domain-knowledge in this manner (that is: adding a domain-based loss to the standard problem-specific loss) may not always be suitable while dealing with safety-critical domains where 100% constraint satisfaction is desirable. One way to guarantee 100% domain-constraint satisfaction is by directly augmenting the output layer with some transformations and then deriving a new loss function due to these transformations. These transformations are such that they guarantee the output of the network to satisfy the domain constraints. In this study, called MultiplexNet⁷⁴, the domain-knowledge is represented as a logical formula in disjunctive normal form (DNF) Here the output (or prediction) layer of a deep network is viewed as a multiplexer in a logical circuit that permits branching in logic.

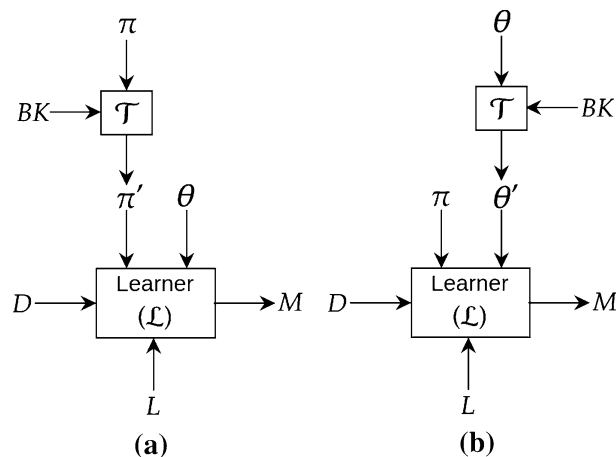


Figure 8. Introducing background knowledge into deep network by transforming the model (structure and parameter). In (a), the transformation block \mathcal{T} takes an input structure of a model π and outputs a transformed structure π' based on background knowledge (BK). In (b), the transformation block \mathcal{T} takes a set of parameters θ of a model and outputs a transformed set of parameters π' based on background knowledge (BK).

That is, the output of the network always satisfies one of the constraints specified in the domain knowledge (disjunctive formula).

The other form of semantic loss could be one that involves a human for post-hoc evaluation of a deep model constructed from a set of first-order rules. In this line, In⁷⁵ proposed an analogical reasoning system intended for discovering rules by training a sequence-to-sequence model using a training set of rules represented in first-order logic. Here the role of domain-knowledge comes post training of the deep sequence model; that is, an evaluator (a human expert) tests each discovered rule from the model by unifying them against the (domain) knowledge base. The domain-knowledge here serves as some kind of a validation set where if the ratio of successful rule unification crosses a certain threshold, then the set of discovered rules are accepted.

Transforming the model

Over the years, many studies have shown that domain knowledge can be incorporated into a deep network by introducing constraints on the model parameters (weights) or by making a design choice of its structure. Figure 8 shows a simple block diagram showing domain knowledge incorporation at the design stage of the deep network.

Constraints on parameters. In a Bayesian formulation, there is an explicit mechanism for the inclusion of domain-knowledge through the use of priors. The regularisation terms in loss-functions, for example, can be seen as an encoding of such prior information, usually on the network's structure. Priors can also be introduced on the parameters (weights) of a network. Explicitly, these would take the form of a prior distribution over the values of the weights in the network. The priors on networks and network weights represent our expectations about networks before receiving any data, and correspond to penalty terms or regularisers. Buntine and Weigend⁷⁶ extensively study how Bayesian theory can be highly relevant to the problem of training feedforward neural networks. This work is explicitly concerned with choosing an appropriate network structure and size based on prior domain-knowledge and with selecting a prior for the weights.

The work by⁷⁷ on Bayesian learning for neural networks also showed how domain-knowledge could help build a prior probability distribution over neural network parameters. In this, the shown priors allow networks to be “self-regularised” to not over-fit even when the complexity of the neural network is increased. In a similar spirit,⁷⁸ showed how prior domain knowledge could be used to define ‘meta-features’ that can aid in defining the prior distribution of weights. These meta-features are additional information about each of the features in the available data. For instance, for an image recognition task, the meta-feature could be the relative position of a pixel (x, y) in the image. This meta information can be used to construct a prior over the weights for the original features.

Transfer learning. Transfer Learning is a mechanism to introduce priors on weights when data is scarce for a problem (usually called the “target” domain). Transfer learning relies on data availability for a problem similar to the target domain (usually called the “source” domain). From the position taken in this paper, domain-knowledge for transfer learning is used to change the structure or the parameter values (or both) for a model for the target problem. The nature of this domain-knowledge can be seen prior distributions on the structure and/or parameter-values (weights) of models for the target problem. The prior distributions for the target model are obtained from the models constructed for the source problem.

In practice, transfer learning from a source domain to a target domain usually involves a transfer of weights from models constructed for the source domain to the network in the target domain. This has been shown to boost performance significantly. From the Bayesian perspective, transfer learning allows the construction of the prior over the weights of a neural network for the target domain based on the posterior constructed in the

source domain. Transfer learning is not limited by the kind of task (such as classification, regression, etc.) but rather by the availability of related problems. Language models are some of the very successful examples of the use of transfer learning, where the models are initially learnt on a huge corpus of data and fine-tuned for other more specialised tasks. In⁷⁹ provides an in-depth review of some of the mechanisms and the strategies of transfer learning. Transfer learning need not be restricted to deep networks only: in a recent study,⁸⁰ proposes a model that transfers knowledge from a neural network to a decision tree using knowledge distillation framework. The symbolic knowledge encoded in the decision tree could further be utilised for a variety of tasks.

A subcategory of transfer learning is one in which the problem (or task) remains the same, but there is a change in the distribution over the input data from the source and the target. This form of learning is viewed as an instance of domain-adaptation³⁴. Similar to transfer learning, the knowledge is transferred from a source domain to a target domain in the form of a prior distribution over the model parameters. This form of domain-adaptation uses the same model structure as the source, along with an initial set of parameter values obtained from the source model. The parameter values are then fine-tuned using labelled and unlabelled data from the target data⁸¹. An example of this kind of learning is in⁸² where a BERT model is fine-tuned with data from multiple domains. There are some recent surveys along these lines:^{83,84}.

Specialised structures. DNN based methods arguably work best if the domain-knowledge is used to inspire their architecture choices⁸⁵. There are reports on incorporating first-order logic constructs into the structure of the network. This allows neural-networks to operate directly on the logical sentences comprising domain-knowledge.

Domain-knowledge encoded as a set of propositional rules are used to constrain the structure of the neural network. Parameter-learning (updating of the network weights) then proceeds as normal, using the structure. The result could be thought of as learning weighted forms of the antecedents present in the rules. The most popular and oldest work along this line is Knowledge-Based Artificial Neural Network (KBANN)¹³ that incorporates knowledge into neural networks. In KBANN, the domain knowledge is represented as a set of hierarchically structured propositional rules that directly determines a fixed topological structure of a neural network⁸⁶. KBANN was successful in many real-world applications; but, its representational power was bounded by pre-existing set of rules which restricted it to refine these existing rules rather than discovering new rules. A similar study is KBCNN⁸⁷, which first identifies and links domain attributes and concepts consistent with initial domain knowledge. Further, KBCNN introduces additional hidden units into the network and most importantly, it allowed decoding of the learned rules from the network in symbolic form. However, both KBANN and KBCNN were not appropriate for learning new rules because of the way the initial structure was constructed using the initial domain knowledge base.

Some of the limitations described above could be overcome with the proposal of a hybrid system by Fletcher and Obradovic⁸⁸. The system was able to learn a neural network structure that could construct new rules from an initial set of rules. Here, the domain knowledge is transformed into an initial network through an extended version of KBANN's symbolic knowledge encoding. It performed incremental hidden unit generation thereby allowing construction or extension of initial rule-base. In a similar manner, there was a proposal for using Cascade ARTMAP⁸⁹ which could not only construct a neural network structure from rules but also perform explicit cascading of rules and multistep inferencing. It was found that the rules extracted from Cascade ARTMAP are more accurate and much cleaner than the rules extracted from KBANN⁹⁰.

In the late 1990s, Garcez and Zaverucha proposed a massively parallel computational model called CIL²P based on feedforward neural network that integrates inductive learning from examples and domain knowledge, expressed as a propositional logic program⁹¹. A translation algorithm generates a neural network. Unlike KBANN, the approach uses the notion of "bipolar semi-linear" neurons. This allows the proof of a form of correctness, showing the existence of a neural-network structure that can compute the logical consequences of the domain-knowledge. The output of such a network, when combined into subsequent processing naturally incorporates the intended interpretation of the domain predicates. The authors extend this to the use of first-order logic programs: we have already considered this in Sect. 2.

A recent proposal focuses on embedding symbolic knowledge expressed as logical rules⁹¹. It considers two languages of representations: Conjunctive Normal Form (CNF) and decision-Deterministic Decomposable Negation Normal form (d-DNNF), which can naturally be represented as graph structures. The graph structures can be provided to a graph neural network (GNN) to learn an embedding suitable for further task-specific implementations.

Somewhat in a similar vein to the work by²⁹, the work reported in⁶³ considers as a set of propositional statements representing domain constraints. A deep network is then trained to find satisfying assignments for the constraints. Again, once such a network is constructed, it can clearly be used in subsequent processing, capturing the effect of the domain constraints. The network is trained using a semantic loss that we have described in Sect. 3.2.

In⁵⁸ it is proposed to augment a language model that uses a deep net architecture with additional statements in first-order logic. Thus, given domain-knowledge encoded as first-order relations, connections are introduced into the network based on the logical constraints enforced by the domain-relations. The approach is related somewhat to the work in⁹² that does not explicitly consider the incorporation of domain-knowledge but does constrain a deep neural network's structure by first grounding a set of weighted first-order definite clauses and then turning them into propositional programs.

We note that newer areas are emerging that use representations for domain-knowledge that go beyond first-order logic relations. This includes probabilistic first-order logic, as a way of including uncertain domain-knowledge⁹³. One interesting way this is being used is to constrain the training of "neural predicates", which represent probabilistic relations that are implemented by neural networks, and the framework can be trained in

an end-to-end fashion^{93,94}. In DeepProbLog⁹³, for example, high-level logical reasoning can be combined with the sub-symbolic discriminative power of deep networks. For instance, a logic program for adding two digits and producing the output sum is straightforward. However, what if the inputs are images of the corresponding digits? Here, a deep network is used to map an image to a digit, while a (weighted) logic program, written in ProbLog⁹⁵, for the addition operation is treated as the symbolic domain knowledge. The ProbLog program is extended with a set of ground neural predicates for which the weights correspond to the probability distribution of classes of digits (0 ...9). The parameters (weights of predicates and weights of neural network) are learned in an end-to-end fashion. A recent approach called DeepStochLog⁹⁴ is a framework that extends the idea of neural predicates in DeepProbLog to definite clause grammars⁹⁶. The reader may note that although DeepProbLog and DeepStochLog do not really transform the structure of the deep network, we are still considering these methods under the heading of specialised structures. This is because of the fact that the hybrid architecture is a tightly coupled approach combining probabilistic logic and deep neural networks.

One of the approaches involves transformation of a probabilistic logic program to graph-structured representation. For instance, in kLog⁹⁷ the transformed representation is an undirected bipartite graph in the form of 'Probabilistic Entity-Relationship model'⁹⁸ which allows the use of a graph-kernel⁹⁹ for data classification purpose, where each data instance is represented as a logic program constructed from data and background-knowledge. Another approach uses weighted logic programs or *templates* with GNNs¹⁰⁰ demonstrating how simple relational logic programs can capture advanced graph convolution operations in a tightly integrated manner. However, it requires the use of a language of Lifted Relational Neural Networks (LRNNs)¹⁰¹.

An interesting proposal is to transform facts and rules, all represented in (weighted) first-order logic into matrix (or tensor) representations. Learning and inference can then be conducted on these matrices (or tensors)^{102,103} allowing faster computation. NeuralLog¹⁰⁴, for example, extends this idea and constructs a multilayered neural network, to some extent, similar to the ones in LRNN consisting of fact layer, rule layer and literal layer etc. The learning here refers to the updates of the weights of the rules. Another work that translates domain-knowledge in first-order logic into a deep neural network architecture consisting of the input layer (grounded atoms), propositional layer, quantifier layer and output layer is⁶⁸. Similar to LRNN, it uses the fuzzy *t*-norm operators for translating logical OR and AND operations.

Further emerging areas look forward to providing domain-knowledge as higher-order logic templates (or "meta-rules": see¹⁷ for pointers to this area). To the best of our knowledge, there are, as yet, no reports in the literature on how such higher-order statements can be incorporated into deep networks.

Challenges and concluding remarks

We summarise our discussion on domain-knowledge as constraints in Table 1. We now outline some challenges in incorporating domain-knowledge encoded as logical or numerical constraints into a deep network. We first outline some immediate practical challenges concerning the logical constraints:

- There is no standard framework for translating logical constraints to neural networks. While there are simplification methods which first construct a representation of the logical constraint that a standard deep network can consume, this process has its limitations as described in the relevant section above.
- Logic is not differentiable. This does not allow using standard training of deep network using gradient-based methods in an end-to-end fashion. Propagating gradients via logic has now been looked at in¹⁰⁵, but the solution is intractable and does not allow day-to-day use.
- Many neural network structures are directed acyclic graphs (DAGs). However, transforming logical formulae directly into neural network structures in the manner described in some of the discussed works can introduce cyclic dependencies, which may need a separate form of translations.

There are also practical challenges concerning the numerical constraints:

- We have seen that the numerical constraints are often provided with the help of modification to a loss function. Given some domain-knowledge in a logical representation, constructing a term in loss function is not straightforward.
- Incorporating domain-knowledge via domain-based loss may not be suitable for some safety-critical applications.
- The process of introducing a loss term often results in a difficult optimisation problem (sometimes constrained) to be solved. This may require additional mathematical tools for a solution that can be implemented practically.
- Deep network structures constrained via logical domain-knowledge may not always be scalable to large datasets.

It is possible to consider representing domain-knowledge not as logical or numeric constraints, but through statements in natural language. Recent rapid progress in the area of language models, for example, the models based on attention^{106,107} raises the possibility of incorporating domain-knowledge through conversations. While the precision of these formal representations may continue to be needed for the construction of scientific assistants, the flexibility of natural language may be especially useful in communicating commonsense knowledge to day-to-day machine assistants that need to an informal knowledge of the world^{108,109}. Progress in this is being made (see, for example, <https://allenai.org/aristo>), but there is much more that needs to be done to make the language models required accessible to everyday machinery.

Principal approach	Work (reference)	Type of learner
Transforming Data	DRM ^{24,25}	MLP
	CILP++ ²⁸	MLP
	R-GCN ⁴⁶	GNN
	KGCN ⁶¹	GNN
	KBRD ⁴⁹	GNN
	DG-RNN ⁴⁴	RNN
	DreamCoder ³²	DNN*
	Gated-K-BERT ³⁸	Transformer
	VEGNN ¹⁵	GNN
	BotGNN ⁶	GNN
	KRISP ⁴⁵	GNN, Transformer
Transforming Loss	IPKFL ⁷⁸	CNN
	ILBKRME ³⁰	MLP
	HDNNLR ⁶⁴	CNN, RNN
	SBR ⁶⁸	MLP
	SBR ⁶⁹	CNN
	DL2 ⁶⁶	CNN
	Semantic Loss ⁶³	CNN
	LENSR ⁹¹	GNN
	DANN ⁶⁷	MLP
	PC-LSTM ⁷²	RNN
	DomiKnowS ⁷³	DNN*
	MultiplexNet ⁷⁴	MLP, CNN
	Analogy Model ⁷⁵	RNN
Transforming Model	KBANN ⁸⁶	MLP
	Cascade-ARTMAP ⁸⁹	ARTMAP
	CIL ² p ²⁹	RNN
	DeepProbLog ⁹³	CNN
	LRNN ¹⁰¹	MLP
	TensorLog ¹⁰³	MLP
	Domain-Aware BERT ⁸²	Transformer
	NeuralLog ¹⁰⁴	MLP
DeepStochLog ⁹⁴	DNN*	

Table 1. Some selected works, in no particular order, showing the principal approach of domain knowledge inclusion into deep neural networks. For each work referred here, we show the type of learner with following acronyms: Multilayer Perceptron (MLP), Convolutional Neural Network (CNN), Recurrent Neural Network (RNN), Graph Neural Network (GNN), Adaptive Resonance Theory-based Network Map (ARTMAP), DNN* refers to a DNN structure dependent on intended task. We use ‘MLP’ here to represent any neural network, that conforms to a layered-structure that may or maynot be fully-connected. RNN also refers to sequence models constructed using Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU) cells.

More broadly, incorporating domain-knowledge into learning is highlighted in¹ as one of the Grand Challenges facing the foundations of AI and ML. Addressing this challenge effectively is seen as being relevant to issues arising in automated model-construction like data-efficiency and constraint-satisfaction. Additionally, it is suggested that developing a mapping of internal representations of the model to domain-concepts maybe necessary for acceptable explanations for the model’s predictions and for developing trust in the model.

It is now accepted that trust comes through understanding of how decisions are made by the machine-constructed models¹¹⁰, and what are the determining factors in these decisions. One important requirement of machine-constructed models in workflows with humans-in-the-loop is that the models are human-understandable. Domain-knowledge can be used in two different ways to assist this. First, it can constrain the kinds of models that are deemed understandable. Secondly, it can provide concepts that are meaningful for use in a model. Most of the work in this review has been focused on improving predictive performance. However, the role of domain-knowledge in constructing explanations for deep network models is also being explored (see, for example,¹¹¹). However, that work only generates *post hoc* explanations that are locally consistent. Explanatory deep network models that identify true causal connections based on concepts provided as domain-knowledge remain elusive.

Domain-knowledge can also be used to correct biases¹¹² built into a deep network either declaratively, through the use of constraints, or through the use of loss functions that include “ethical penalty” terms. Demonstrations

of the use of domain-knowledge driven, ethics-sensitive machine learning have been available in the literature for some time¹¹³. Can these carry over to the construction of deep network models? This remains to be investigated.

The issues raised above all go beyond just the “how” questions related to the incorporation of domain-knowledge into deep neural networks. They provide pointers to why the use of domain-knowledge may extend beyond its utility for prediction.

Received: 5 August 2021; Accepted: 20 December 2021

Published online: 20 January 2022

References

1. Stevens, R. *et al.* Ai for science. Tech. Rep., Argonne National Lab.(ANL), Argonne, IL (United States) (2020).
2. Kitano, H. Artificial intelligence to win the nobel prize and beyond: Creating the engine for scientific discovery. *AI Mag.* **37**, 39–49 (2016).
3. Lipton, Z. C. The mythos of model interpretability. arXiv preprint [arXiv:1606.03490](https://arxiv.org/abs/1606.03490) (2016).
4. Arrieta, A. B. *et al.* Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. arXiv preprint [arXiv:1910.10045](https://arxiv.org/abs/1910.10045) (2019).
5. Dash, T., Srinivasan, A. & Vig, L. Incorporating symbolic domain knowledge into graph neural networks. *Mach. Learn.* 1–28 (2021).
6. Dash, T., Srinivasan, A. & Baskar, A. Inclusion of domain-knowledge into gnn using mode-directed inverse entailment. arXiv preprint [arXiv:2105.10709](https://arxiv.org/abs/2105.10709) (2021).
7. Lavrač, N., Džeroski, S. & Grobelnik, M. Learning nonrecursive definitions of relations with linus. In *European Working Session on Learning*, 265–281 (Springer, 1991).
8. Tomašev, N. *et al.* Ai for social good: Unlocking the opportunity for positive impact. *Nat. Commun.* **11**, 1–6 (2020).
9. Garcez, A. S. d., Broda, K. B. & Gabbay, D. M. *Neural-symbolic learning systems: foundations and applications* (Springer, 2012).
10. Raedt, L. d., Dumančić, S., Manhaeve, R. & Marra, G. From statistical relational to neuro-symbolic artificial intelligence. In Bessiere, C. (ed.) *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, 4943–4950 (International Joint Conferences on Artificial Intelligence Organization, 2020). <https://doi.org/10.24963/ijcai.2020/688>. Survey track.
11. Besold, T. R. *et al.* Neural-symbolic learning and reasoning: A survey and interpretation. arXiv preprint [arXiv:1711.03902](https://arxiv.org/abs/1711.03902) (2017).
12. von Rueden, L. *et al.* Informed machine learning—a taxonomy and survey of integrating prior knowledge into learning systems. *IEEE Trans. Knowl. Data Eng.* (2021).
13. Towell, G. G., Shavlik, J. W. & Noordewier, M. O. Refinement of approximate domain theories by knowledge-based neural networks. In *Proceedings of the eighth National conference on Artificial intelligence*, vol. 861866 (Boston, MA, 1990).
14. Muggleton, S. Inductive logic programming. *New Gen. Comput.* **8**, 295–318 (1991).
15. Muggleton, S. & de Raedt, L. Inductive logic programming: Theory and methods. *J. Logic Program.* **19–20**, 629–679 (1994). <https://www.sciencedirect.com/science/article/pii/0743106694900353>. Special Issue: Ten Years of Logic Programming.
16. Muggleton, S. *et al.* Ilp turns 20. *Mach. Learn.* **86**, 3–23 (2012).
17. Cropper, A., Dumančić, S. & Muggleton, S. H. Turning 30: New ideas in inductive logic programming. In Bessiere, C. (ed.) *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, 4833–4839 (2020). <https://doi.org/10.24963/ijcai.2020/673>.
18. Kramer, S., Lavrač, N. & Flach, P. Propositionalization Approaches to Relational Data Mining, 262–291 (Springer). *Heidelberg* https://doi.org/10.1007/978-3-662-04599-2_11 (2001).
19. Lavrac, N., Skrlj, B. & Robnik-Sikonja, M. Propositionalization and embeddings: Two sides of the same coin. *Mach. Learn.* **109**, 1465–1507. <https://doi.org/10.1007/s10994-020-05890-8> (2020).
20. Vig, L., Srinivasan, A., Bain, M. & Verma, A. An investigation into the role of domain-knowledge on the use of embeddings. In Lachiche, N., & Vrain, C. (eds.) *Inductive Logic Programming—27th International Conference, ILP 2017, Orleans, France, September 4-6, 2017, Revised Selected Papers*, vol. 10759 of *Lecture Notes in Computer Science*, 169–183 (Springer, 2017). https://doi.org/10.1007/978-3-319-78090-0_12.
21. Gaur, M., Roy, K., Sharma, A., Srivastava, B. & Sheth, A. “who can help me?”: Knowledge infused matching of support seekers and support providers during covid-19 on reddit. arXiv preprint [arXiv:2105.06398](https://arxiv.org/abs/2105.06398) (2021).
22. Cer, D. *et al.* Universal sentence encoder for english. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 169–174 (2018).
23. Radford, A. *et al.* Language models are unsupervised multitask learners. *OpenAI Blog.* **1**, 9 (2019).
24. Lodhi, H. Deep relational machines. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* (2013).
25. Dash, T., Srinivasan, A., Vig, L., Orhobor, O. I. & King, R. D. Large-scale assessment of deep relational machines. In *International Conference on Inductive Logic Programming*, 22–37 (Springer, 2018).
26. Dash, T., Srinivasan, A., Joshi, R. S. & Baskar, A. Discrete stochastic search and its application to feature-selection for deep relational machines. In *International Conference on Artificial Neural Networks*, 29–45 (Springer, 2019).
27. Muggleton, S. Inverse entailment and progol. *New Gen. Comput.* **13**, 245–286 (1995).
28. França, M. V., Zaverucha, G. & Garcez, A. S. Fast relational learning using bottom clause propositionalization with artificial neural networks. *Mach. Learn.* **94**, 81–104 (2014).
29. Avila Garcez, A. S. & Zaverucha, G. Connectionist inductive learning and logic programming system. *Appl. Intell.* (1999).
30. Rocktäschel, T., Singh, S. & Riedel, S. Injecting logical background knowledge into embeddings for relation extraction. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 1119–1129 (Association for Computational Linguistics, Denver, Colorado, 2015). <https://www.aclweb.org/anthology/N15-1118>.
31. Ellis, K., Morales, L., Meyer, M. S., Solar-Lezama, A. & Tenenbaum, J. B. Dreamcoder: Bootstrapping domain-specific languages for neurally-guided bayesian program learning. In *Proceedings of the 2nd Workshop on Neural Abstract Machines and Program Induction* (2018).
32. Ellis, K. *et al.* Dreamcoder: Growing generalizable, interpretable knowledge with wake-sleep bayesian program learning. arXiv preprint [arXiv:2006.08381](https://arxiv.org/abs/2006.08381) (2020).
33. Hinton, G. E., Dayan, P., Frey, B. J. & Neal, R. M. The “wake-sleep” algorithm for unsupervised neural networks. *Science* **268**, 1158–1161 (1995).
34. Ben-David, S. *et al.* A theory of learning from different domains. *Mach. Learn.* **79**, 151–175 (2010).
35. Dong, T. *et al.* Generative adversarial network-based transfer reinforcement learning for routing with prior knowledge. *IEEE Trans. Network Service Manag.* (2021).
36. Purohit, H., Shalin, V. L. & Sheth, A. P. Knowledge graphs to empower humanity-inspired ai systems. *IEEE Internet Comput.* **24**, 48–54 (2020).

37. Sheth, A., Padhee, S. & Gyrard, A. Knowledge graphs and knowledge networks: The story in brief. *IEEE Internet Comput.* **23**, 67–75 (2019).
38. Yadav, S. *et al.* “when they say weed causes depression, but it’s your fav antidepressant”: Knowledge-aware attention framework for relationship extraction. *PLoS one* **16**, e0248299 (2021).
39. Hogan, A. *et al.* Knowledge graphs. arXiv preprint [arXiv:2003.02320](https://arxiv.org/abs/2003.02320) (2020).
40. Kursuncu, U., Gaur, M. & Sheth, A. Knowledge infused learning (k-il): Towards deep incorporation of knowledge in deep learning. arXiv preprint [arXiv:1912.00512](https://arxiv.org/abs/1912.00512) (2019).
41. Sheth, A., Gaur, M., Kursuncu, U. & Wickramarachchi, R. Shades of knowledge-infused learning for enhancing deep learning. *IEEE Internet Comput.* **23**, 54–63 (2019).
42. Gaur, M., Faldu, K. & Sheth, A. Semantics of the black-box: Can knowledge graphs help make deep learning systems more interpretable and explainable?. *IEEE Internet Comput.* **25**, 51–59 (2021).
43. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT (1)*, 4171–4186 (2019). <https://aclweb.org/anthology/papers/N/N19/N19-1423>.
44. Yin, C., Zhao, R., Qian, B., Lv, X. & Zhang, P. Domain knowledge guided deep learning with electronic health records. In *2019 IEEE International Conference on Data Mining (ICDM)*, 738–747 (IEEE, 2019).
45. Marino, K., Chen, X., Parikh, D., Gupta, A. & Rohrbach, M. Krisp: Integrating implicit and symbolic knowledge for open-domain knowledge-based vqa. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14111–14121 (2021).
46. Schlichtkrull, M. *et al.* Modeling relational data with graph convolutional networks. In *European Semantic Web Conference*, 593–607 (Springer, 2018).
47. Park, N., Kan, A., Dong, X. L., Zhao, T. & Faloutsos, C. Estimating Node Importance in Knowledge Graphs Using Graph Neural Networks, 596–606 (Association for Computing Machinery, New York, NY). USA <https://doi.org/10.1145/3292500.3330855> (2019).
48. Wang, H., Zhao, M., Xie, X., Li, W. & Guo, M. Knowledge graph convolutional networks for recommender systems. In *The World Wide Web Conference, WWW ’19*, 3307–3313 (Association for Computing Machinery, New York, NY, USA, 2019). <https://doi.org/10.1145/3308558.3313417>.
49. Chen, Q. *et al.* Towards knowledge-based recommender dialog system. arXiv preprint [arXiv:1908.05391](https://arxiv.org/abs/1908.05391) (2019).
50. Lehmann, J. *et al.* Dbpedia-a large-scale, multilingual knowledge base extracted from wikipedia. *Semant. Web* **6**, 167–195 (2015).
51. Zhuang, Z., Wang, J., Qi, Q., Sun, H. & Liao, J. Toward greater intelligence in route planning: A graph-aware deep learning approach. *IEEE Syst. J.* **14**, 1658–1669 (2019).
52. Dash, T., Srinivasan, A., Vig, L. & Roy, A. Using domain-knowledge to assist lead discovery in early-stage drug design. *bioRxiv* (2021). <https://www.biorxiv.org/content/early/2021/07/09/2021.07.09.451519>.
53. Kukačka, J., Golkov, V. & Cremers, D. Regularization for deep learning: A taxonomy. arXiv preprint [arXiv:1710.10686](https://arxiv.org/abs/1710.10686) (2017).
54. Fu, L. M. Introduction to knowledge-based neural networks. *Knowl.-Based Syst.* (1995).
55. Rocktäschel, T., Bosnjak, M., Singh, S. & Riedel, S. Low-dimensional embeddings of logic. In *Proceedings of the ACL 2014 Workshop on Semantic Parsing*, 45–49 (Association for Computational Linguistics, Baltimore, MD, 2014). <https://www.aclweb.org/anthology/W14-2409>.
56. Guo, S., Wang, Q., Wang, L., Wang, B. & Guo, L. Jointly embedding knowledge graphs and logical rules. In *Proceedings of the 2016 conference on empirical methods in natural language processing*, 192–202 (2016).
57. Hájek, P. *Metamathematics of fuzzy logic*, vol. 4 (Springer Science & Business Media, 2013).
58. Li, T. & Srikumar, V. Augmenting neural networks with first-order logic. In *ACL 2019 - 57th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference* (2020). [arXiv:1906.06298](https://arxiv.org/abs/1906.06298).
59. Hamilton, W. L., Bajaj, P., Zitnik, M., Jurafsky, D. & Leskovec, J. Embedding logical queries on knowledge graphs. In *NeurIPS* (2018).
60. Demeester, T., Rocktäschel, T. & Riedel, S. Lifted rule injection for relation embeddings. [arXiv:1606.08359](https://arxiv.org/abs/1606.08359) (2016).
61. Li, L. *et al.* Domain knowledge embedding regularization neural networks for workload prediction and analysis in cloud computing. *J. Inf. Technol. Res.* **11**, 137–154. <https://doi.org/10.4018/JITR.2018100109> (2018).
62. Takeishi, N. & Akimoto, K. Knowledge-based distant regularization in learning probabilistic models. arXiv preprint [arXiv:1806.11332](https://arxiv.org/abs/1806.11332) (2018).
63. Xu, J., Zhang, Z., Friedman, T., Liang, Y. & Van Den Broeck, G. A semantic loss function for deep learning with symbolic knowledge. In *35th International Conference on Machine Learning, ICML 2018* (2018). [arXiv:1711.11157](https://arxiv.org/abs/1711.11157).
64. Hu, Z., Ma, X., Liu, Z., Hovy, E. & Xing, E. Harnessing deep neural networks with logic rules. [arXiv:1603.06318](https://arxiv.org/abs/1603.06318) (2016).
65. Hinton, G., Vinyals, O. & Dean, J. Distilling the knowledge in a neural network. arXiv preprint [arXiv:1503.02531](https://arxiv.org/abs/1503.02531) (2015).
66. Fischer, M. *et al.* DL2: Training and querying neural networks with logic. In *ICML* (2019).
67. Muralidhar, N., Islam, M. R., Marwah, M., Karpatne, A. & Ramakrishnan, N. Incorporating Prior Domain Knowledge into Deep Neural Networks. In *Proceedings - 2018 IEEE International Conference on Big Data, Big Data 2018* (2019).
68. Diligenti, M., Gori, M. & Sacca, C. Semantic-based regularization for learning and inference. *Artif. Intell.* **244**, 143–165 (2017).
69. Diligenti, M., Roychowdhury, S. & Gori, M. Integrating prior knowledge into deep learning. In *2017 16th IEEE International Conference on Machine Learning and Applications (ICMLA)*, 920–923 (IEEE, 2017).
70. Silvestri, M., Lombardi, M. & Milano, M. Injecting domain knowledge in neural networks: a controlled experiment on a constrained problem. arXiv preprint [arXiv:2002.10742](https://arxiv.org/abs/2002.10742) (2020).
71. Rossi, F., Van Beek, P. & Walsh, T. *Handbook of constraint programming* (Elsevier, 2006).
72. Luo, X., Zhang, D. & Zhu, X. Deep learning based forecasting of photovoltaic power generation by incorporating domain knowledge. *Energy* **225**, 120240 (2021).
73. Faghihi, H. R. *et al.* Domiknows: A library for integration of symbolic domain knowledge in deep learning. arXiv preprint [arXiv:2108.12370](https://arxiv.org/abs/2108.12370) (2021).
74. Hoernle, N., Karampatsis, R. M., Belle, V. & Gal, K. Multiplexnet: Towards fully satisfied logical constraints in neural networks. arXiv preprint [arXiv:2111.01564](https://arxiv.org/abs/2111.01564) (2021).
75. Honda, H. & Hagiwara, M. Analogical reasoning with deep learning-based symbolic processing. *IEEE Access* **9**, 121859–121870 (2021).
76. Buntine, W. L. & Weigend, A. Bayesian back-propagation. *Complex Syst.* **5** (1991).
77. Neal, R. M. *Bayesian Learning for Neural Networks*. Ph.D. thesis, University of Toronto, CAN (1995). AAINN02676.
78. Krupka, E. & Tishby, N. Incorporating prior knowledge on features into learning. In *AISTATS* (2007).
79. Zhuang, F. *et al.* A comprehensive survey on transfer learning. *Proc. IEEE* **109**, 43–76 (2020).
80. Liu, X., Wang, X. & Matwin, S. Improving the interpretability of deep neural networks with knowledge distillation. arXiv preprint [arXiv:1812.10924](https://arxiv.org/abs/1812.10924) (2018).
81. Tzeng, E., Hoffman, J., Darrell, T. & Saenko, K. Simultaneous deep transfer across domains and tasks. In *Proceedings of the IEEE international conference on computer vision*, 4068–4076 (2015).
82. Du, C., Sun, H., Wang, J., Qi, Q. & Liao, J. Adversarial and domain-aware bert for cross-domain sentiment analysis. In *Proceedings of the 58th annual meeting of the Association for Computational Linguistics*, 4019–4028 (2020).
83. Wang, M. & Deng, W. Deep visual domain adaptation: A survey. *Neurocomputing* **312**, 135–153 (2018).

84. Ramponi, A. & Plank, B. Neural unsupervised domain adaptation in nlp—a survey. In *Proceedings of the 28th International Conference on Computational Linguistics*, 6838–6855 (2020).
85. Berner, J., Grohs, P., Kutyniok, G. & Petersen, P. The modern mathematics of deep learning. arXiv preprint [arXiv:2105.04026](https://arxiv.org/abs/2105.04026) (2021).
86. Towell, G. G. & Shavlik, J. W. Knowledge-based artificial neural networks. *Artif. Intell.* **70**, 119–165 (1994).
87. Fu, L. M. Knowledge-based connectionism for revising domain theories. *IEEE Trans. Syst. Man. Cybern.* **23**, 173–182 (1993).
88. Fletcher, J. & Obradovic, Z. Combining prior symbolic knowledge and constructive neural network learning. *Connect. Sci.* **5**, 365–375 (1993).
89. Tan, A. H. Cascade ARTMAP: Integrating neural computation and symbolic knowledge processing. *IEEE Trans. Neural Netw.* (1997).
90. Towell, G. G. & Shavlik, J. W. Extracting refined rules from knowledge-based neural networks. *Mach. Learn.* **13**, 71–101 (1993).
91. Xie, Y., Xu, Z., Kankanhalli, M. S., Meel, K. S. & Soh, H. Embedding symbolic knowledge into deep networks. In *Advances in Neural Information Processing Systems* (2019).
92. Sourek, G., Aschenbrenner, V., Zelezny, F., Schockaert, S. & Kuzelka, O. Lifted relational neural networks: Efficient learning of latent relational structures. *J. Artif. Intell. Res.* **62**, 69–100. <https://doi.org/10.1613/jair.1.11203> (2018).
93. Manhaeve, R., Dumancic, S., Kimmig, A., Demeester, T. & De Raedt, L. Deepprolog: Neural probabilistic logic programming. *Adv. Neural Inf. Process. Syst.* **31**, 3749–3759 (2018).
94. Winters, T., Marra, G., Manhaeve, R. & Raedt, L. D. Deepstochlog: Neural stochastic logic programming. [arXiv:2106.12574](https://arxiv.org/abs/2106.12574) (2021).
95. De Raedt, L., Kimmig, A. & Toivonen, H. Prolog: A probabilistic prolog and its application in link discovery. In *IJCAI*, vol. 7, 2462–2467 (Hyderabad, 2007).
96. Pereira, F. C. & Warren, D. H. Definite clause grammars for language analysis—a survey of the formalism and a comparison with augmented transition networks. *Artif. Intell.* **13**, 231–278 (1980).
97. Frasconi, P., Costa, F., De Raedt, L. & De Grave, K. klog: A language for logical and relational learning with kernels. *Artif. Intell.* **217**, 117–143 (2014).
98. Heckerman, D., Meek, C. & Koller, D. Probabilistic entity-relationship models, prms, and plate models. *Introduction to statistical relational learning* 201–238 (2007).
99. Vishwanathan, S. V. N., Schraudolph, N. N., Kondor, R. & Borgwardt, K. M. Graph kernels. *J. Mach. Learn. Res.* **11**, 1201–1242 (2010).
100. Sourek, G., Zelezny, F. & Kuzelka, O. Beyond graph neural networks with lifted relational neural networks. arXiv preprint [arXiv:2007.06286](https://arxiv.org/abs/2007.06286) (2020).
101. Sourek, G., Aschenbrenner, V., Zelezny, F., Schockaert, S. & Kuzelka, O. Lifted relational neural networks: Efficient learning of latent relational structures. *J. Artif. Intell. Res.* **62**, 69–100 (2018).
102. Serafini, L. & Garcez, A. d. Logic tensor networks: Deep learning and logical reasoning from data and knowledge. arXiv preprint [arXiv:1606.04422](https://arxiv.org/abs/1606.04422) (2016).
103. Cohen, W., Yang, F. & Mazaitis, K. R. Tensorlog: A probabilistic database implemented using deep-learning infrastructure. *J. Artif. Intell. Res.* **67**, 285–325 (2020).
104. Guimarães, V. & Costa, V. S. Neurallog: a neural logic language. arXiv preprint [arXiv:2105.01442](https://arxiv.org/abs/2105.01442) (2021).
105. Evans, R. & Grefenstette, E. Learning explanatory rules from noisy data. *J. Artif. Intell. Res.* **61**, 1–64. <https://doi.org/10.1613/jair.5714> (2018).
106. Vaswani, A. *et al.* Attention is all you need. In *Advances in neural information processing systems*, 5998–6008 (2017).
107. Brown, T. B. *et al.* Language models are few-shot learners. arXiv preprint [arXiv:2005.14165](https://arxiv.org/abs/2005.14165) (2020).
108. Tandon, N., Varde, A. S. & de Melo, G. Commonsense knowledge in machine intelligence. *SIGMOD Rec.* **46**, 49–52 (2018).
109. Zhao, J., Khashabi, D., Khot, T., Sabharwal, A. & Chang, K.-W. Ethical-advice taker: Do language models understand natural language interventions? In *Findings of the Association for Computational Linguistics: ACL-IJCNLP*, 4158–4164 (2021).
110. Jobin, A., Ienca, M. & Vayena, E. The global landscape of AI ethics guidelines. *Nat. Mach. Intell.* **1**, 389–399 (2019).
111. Srinivasan, A., Vig, L. & Bain, M. Logical explanations for deep relational machines using relevance information. *J. Mach. Learn. Res.* **20**, 1–47 (2019).
112. Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K. & Galstyan, A. A survey on bias and fairness in machine learning. *ACM Comput. Surv. (CSUR)* **54**, 1–35 (2021).
113. Anderson, M., Anderson, S. & Armen, C. *Medethex: Toward a medical ethics advisor* (Caring Machines, In AAAI Fall Symposium, 2005).

Acknowledgements

AS is the Class of 1981 Chair Professor at BITS Pilani; a Visiting Professorial Fellow at UNSW Sydney; and a TCS Affiliate Professor.

Author contributions

This paper was conceived and written by T.D. and A.S. S.C. and A.A. contributed to some sections of an earlier version of the paper (available at: <https://arxiv.org/abs/2103.00180>).

Competing interests

The authors declare no competing interests.

Additional information

Correspondence and requests for materials should be addressed to T.D.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022