# Tensor Logic: The Language of AI

**Pedro Domingos**                                        PEDROD@CS.WASHINGTON.EDU

*Paul G. Allen School of Computer Science & Engineering*
*University of Washington*
*Seattle, WA 98195-2350, USA*

## Abstract

Progress in AI is hindered by the lack of a programming language with all the requisite features. Libraries like PyTorch and TensorFlow provide automatic differentiation and efficient GPU implementation, but are additions to Python, which was never intended for AI. Their lack of support for automated reasoning and knowledge acquisition has led to a long and costly series of hacky attempts to tack them on. On the other hand, AI languages like LISP an Prolog lack scalability and support for learning. This paper proposes tensor logic, a language that solves these problems by unifying neural and symbolic AI at a fundamental level. The sole construct in tensor logic is the tensor equation, based on the observation that logical rules and Einstein summation are essentially the same operation, and all else can be reduced to them. I show how to elegantly implement key forms of neural, symbolic and statistical AI in tensor logic, including transformers, formal reasoning, kernel machines and graphical models. Most importantly, tensor logic makes new directions possible, such as sound reasoning in embedding space. This combines the scalability and learnability of neural networks with the reliability and transparency of symbolic reasoning, and is potentially a basis for the wider adoption of AI.

**Keywords:** deep learning, automated reasoning, knowledge representation, logic programming, Einstein summation, embeddings, kernel machines, probabilistic graphical models

## 1. Introduction

Fields take off when they find their language. Physics took off when Newton invented calculus, and couldn't have done so without it. Maxwell's equations would be unusable without Heaviside's vector calculus notation. As mathematicians and physicists like to say, a good notation is half the battle. Much of electrical engineering would be impossible without complex numbers, and digital circuits without Boolean logic. Modern chip design is made possible by harware description languages, databases by relational algebra, the Internet by the Internet Protocol, and the Web by HTML. More generally, computer science would not have gotten far without high-level programming languages. Qualitative fields also depend critically on their terminology. Even artists rely on the idioms and stylistic conventions of their genre for their work.

A field's language saves its practitioners time, focuses their attention, and changes how they think. It unites the field around common directions and decreases entropy. It makes key things obvious and avoids repeatedly hacking solutions from scratch.

Has AI found its language? LISP, one of the first high-level programming languages, made symbolic AI possible. In the 80s Prolog also became popular. Both, however, suffered from poor scalability and lack of support for learning, and were ultimately displaced, even within AI, by general-purpose languages like Java and C++. Graphical models provide a

*lingua franca* for probabilistic AI, but their applicability is limited by the cost of inference. Formalisms like Markov logic seamlessly combine symbolic and probabilistic AI, but are also hindered by the cost of inference.

Python is currently the *de facto* language of AI, but was never designed for it, and it shows. Libraries like PyTorch and TensorFlow provide important features like automatic differentiation and GPU implementation, but are of no help for key tasks like automated reasoning and knowledge acquisition. Neurosymbolic AI seeks to ameliorate this by combining deep learning modules with symbolic AI ones, but often winds up having the shortcomings of both. In sum, AI has clearly not found its language yet.

There are clear desiderata for such a language. Unlike Python, it should hide everything that is not AI, allowing AI programmers to focus on what matters. It should facilitate incorporating prior knowledge into AI systems and reasoning automatically over it. It should also facilitate learning automatically, and the resulting models should be transparent and reliable. It should scale effortlessly. Symbolic AI has some of these properties and deep learning has others, but neither has all. We therefore need to merge them.

Tensor logic does this by unifying their mathematical foundations. It is based on the observation that essentially all neural networks can be constructed using tensor algebra, all symbolic AI using logic programming, and the two are fundamentally equivalent, differing only in the atomic data types used.

I begin with a brief review of logic programming and tensor algebra. The core of the paper defines tensor logic and describes its inference and learning engines. I then show how to elegantly implement neural networks, symbolic AI, kernel machines and graphical models in it. I show how tensor logic enables reliable and transparent reasoning in embedding space. I propose two approaches to scaling it up. The paper concludes with a discussion of other potential uses of tensor logic, prospects for its wide adoption, and next steps toward it.

## 2. Background

### 2.1 Logic Programming

The most widely used formalism in symbolic AI is logic programming (Lloyd, 1987). The simplest logic programming language, which suffices for our purposes, is Datalog (Greco and Molinaro, 2016). A Datalog program is a set of *rules* and *facts*. A fact is a statement of the form $r(o_1, \ldots, o_n)$, where $r$ is a relation name and the $o$'s are object names. For example, $\texttt{Parent}(\texttt{Bob}, \texttt{Charlie})$ states that Bob is a parent of Charlie, and $\texttt{Ancestor}(\texttt{Alice}, \texttt{Bob})$ that Alice is an ancestor of Bob. A rule is a statement of the form $A_0 \leftarrow A_1, \ldots, A_m$, where the arrow means "if", commas denote conjunction, and each of the $A$'s has the form $r(x_1, \ldots, x_n)$, with $r$ being a relation name and the $x$'s being variables or object names. For example, the rule

$$\texttt{Ancestor}(x, y) \leftarrow \texttt{Parent}(x, y)$$

says that parents are ancestors, and the rule

$$\texttt{Ancestor}(x, z) \leftarrow \texttt{Ancestor}(x, y), \texttt{Parent}(y, z)$$

says that $x$ is $z$'s ancestor if $x$ is $y$'s ancestor and $y$ is $z$'s parent. Informally, a rule says that its left-hand side or *head* is true if there are known facts that make all the relations

on its right-hand side or *body* simultaneously true. For example, the rules and facts above imply that `Ancestor(Alice, Charlie)` is true.

In database terminology, a Datalog rule is a series of *joins* followed by a *projection*. The (natural) join of two relations $R$ and $S$ is the set of all tuples that can be formed from tuples in $R$ and $S$ having the same values for the same arguments. When two relations have no arguments in common, their join reduces to their Cartesian product. The projection of a relation $R$ onto a subset $G$ of its arguments is the relation obtained by discarding from the tuples in $R$ all arguments not in $G$. For example, the rule

$$\text{Ancestor}(x, z) \leftarrow \text{Ancestor}(x, y), \text{Parent}(y, z)$$

joins the relations `Ancestor(x, y)` and `Parent(y, z)` on `y` and projects the result onto $\{x, z\}$; the tuples `Ancestor(Alice, Bob)` and `Parent(Bob, Charlie)` yield the tuple `Ancestor(Alice, Charlie)`.

Two common inference algorithms in logic programming are forward and backward chaining. In forward chaining, the rules are repeatedly applied to the known facts to derive new facts until no further ones can be derived. The result is called the *deductive closure* or *fixpoint* of the program, and all questions of interest can be answered simply by examining it. For example, the answer to the query `Ancestor(Alice, x)` ("Who is Alice an ancestor of?") given the rules and facts above is $\{\text{Bob}, \text{Charlie}\}$.

Backward chaining attempts to answer a question by finding facts that match it or rules that have it as their head and facts that match the body, and so on recursively. For example, the query `Ancestor(Alice, Charlie)` does not match any facts, but it matches the rule

$$\text{Ancestor}(x, z) \leftarrow \text{Ancestor}(x, y), \text{Parent}(y, z)$$

and this rule's body matches the facts `Ancestor(Alice, Bob)` and `Parent(Bob, Charlie)`, and therefore the answer is `True`.

Forward and backward chaining in Datalog are *sound* inference procedures, meaning that the answers they give are guaranteed to follow logically from the rules and facts in the program. Logic programs have both *declarative* and *procedural* semantics, meaning a rule can be interpreted both as a statement about the world and as a procedure for computing its head with the given arguments by calling the procedures in the body and combining the results.

The field of inductive logic programming (ILP) is concerned with learning logic programs from data (Lavrač and Džeroski, 1994). For example, an ILP system might induce the rules above from a small database of parent and ancestor relations. Once induced, these rules can answer questions about ancestry chains of any length and involving anyone. Some ILP systems can also do *predicate invention*, i.e., discover relations that do not appear explicitly in the data, akin to hidden variables in neural networks.

## 2.2 Tensor Algebra

A tensor is defined by two properties: its type (real, integer, Boolean, etc.) and its shape (Rabanser et al., 2017). The shape of a tensor consists of its rank (number of indices) and its size (number of elements) along each index. For example, a video can be represented by an integer tensor of shape $(t, x, y, c)$, where $t$ is the number of frames, $x$ and $y$ are a frame's

3

width and height in pixels, and $c$ is the number of color channels (typically 3). A matrix is a rank-2 tensor, a vector a rank-1 tensor, and a scalar a rank-0 tensor. A tensor of rank $r$ and size $n_i$ in the $i$th dimension contains a total of $\prod_{i=1}^{r} n_i$ elements. The element of a tensor $A$ at position $i_1$ along dimension 1, position $i_d$ along dimension $d$, etc., is denoted by $A_{i_1,\dots,i_d,\dots,i_r}$. This generic element of a tensor is often used to represent the tensor itself. The *sum* of two tensors $A$ and $B$ with the same shape is a tensor $C$ such that

$$C_{i_1,\dots,i_d,\dots,i_r} = A_{i_1,\dots,i_d,\dots,i_r} + B_{i_1,\dots,i_d,\dots,i_r}.$$

The *tensor product* of two tensors $A$ and $B$ of rank respectively $r$ and $r'$ is a tensor $C$ of rank $r + r'$ such that

$$C_{i_1,\dots,i_d,\dots,i_r,j_1,\dots,j_{d'},\dots,j_{r'}} = A_{i_1,\dots,i_d,\dots,i_r} B_{j_1,\dots,j_{d'},\dots,j_{r'}}.$$

Einstein notation simplifies tensor equations by omitting all summation signs and implicitly summing over all repeated indices. For example, $A_{ij}B_{jk}$ represents the product of the matrices $A$ and $B$, summing over $j$ and resulting in a matrix with indices $i$ and $k$:

$$C_{ik} = A_{ij}B_{jk} = \sum_j A_{ij}B_{jk}.$$

More generally, the *Einstein sum* (or einsum for short) of two tensors $A$ and $B$ with common indices $\beta$ is a tensor $C$ such that

$$C_{\alpha\gamma} = \sum_\beta A_{\alpha\beta}B_{\beta\gamma},$$

where $\alpha$, $\beta$ and $\gamma$ are sets of indices, $\alpha$ is the subset of $A$'s indices not appearing in $B$, the elements of $\alpha$ and $\beta$ may be interspersed in any order, and similarly for $B$ and $\gamma$. Essentially all linear and multilinear operations in neural networks can be concisely expressed as einsums (Rocktäschel, 2018; Rogozhnikov, 2022).

Like matrices, tensors can be decomposed into products of smaller tensors. In particular, the *Tucker decomposition* decomposes a tensor into a more compact *core tensor* of the same rank and $k$ *factor matrices*, each expanding an index of the core tensor into an index of the original one. For example, if $A$ is a rank-3 tensor, in Einstein notation its Tucker decomposition is

$$A_{ijk} = M_{ip}M'_{jq}M''_{kr}C_{pqr},$$

where $C$ is the core tensor and the $M$'s are the factor matrices.

## 3. Tensor Logic

### 3.1 Representation

Tensor logic is based on the answers to two key questions: What is the relation between tensors and relations? And what is the relation between Datalog rules and einsums?

The answer to the first question is that a relation is a compact representation of a sparse Boolean tensor. For example, a social network can be represented by the neighborhood

matrix $M_{ij}$, where $i$ and $j$ range over individuals and $M_{ij} = 1$ if $i$ and $j$ are neighbors and 0 otherwise. But for large networks this is an inefficient representation, since almost all elements will be 0. The network can be more compactly represented by a relation, with a tuple for each pair of neighbors; pairs not in the relation are assumed to not be neighbors. More generally, a sparse Boolean tensor of rank $n$ can be compactly represented by an $n$-ary relation with a tuple for each nonzero element, and the efficiency gain will typically increase exponentially with $n$.

The answer to the second question is that a Datalog rule is an einsum over Boolean tensors, with a step function applied elementwise to the result. (Specifically, the Heaviside step function, $H(x) = 1$ if $x > 0$ and 0 otherwise.) For example, consider the rule

$$\texttt{Aunt}(\texttt{x}, \texttt{z}) \leftarrow \texttt{Sister}(\texttt{x}, \texttt{y}), \texttt{Parent}(\texttt{y}, \texttt{z}).$$

Viewing the relations $\texttt{Aunt}(\texttt{x}, \texttt{z})$, $\texttt{Sister}(\texttt{x}, \texttt{y})$ and $\texttt{Parent}(\texttt{y}, \texttt{z})$ as the Boolean matrices $A_{xz}$, $S_{xy}$ and $P_{yz}$, respectively,

$$A_{xz} = H(S_{xy}P_{yz}) = H\left(\sum_y S_{xy}P_{yz}\right)$$

will be 1 iff $S_{xy}$ and $P_{yz}$ are both 1 for at least one $y$. In other words, the einsum $S_{xy}P_{yz}$ implements the join of $\texttt{Sister}(\texttt{x}, \texttt{y})$ and $\texttt{Parent}(\texttt{y}, \texttt{z})$. If $x$ is $z$'s aunt, $y$ is the sibling of $x$ who is also a parent of $z$. The step function is necessary because in general for a given $(x, z)$ pair there may be more than one $y$ for which $S_{xy} = P_{yz} = 1$, leading to a result greater than 1. The step function then reduces this to 1.

Let $U$ and $V$ be arbitrary tensors, and $\alpha$, $\beta$ and $\gamma$ be sets of indices. Then $T_{\alpha\gamma} = H(U_{\alpha\beta}V_{\beta\gamma})$ is a Boolean tensor whose element with indices $\alpha\gamma$ is 1 when there exists some $\beta$ for which $U_{\alpha\beta}V_{\beta\gamma} = 1$. In other words, $T$ represents the join of the relations corresponding to $U$ and $V$.

Since there is a direct correspondence between tensors and relations and between einsums and Datalog rules, there should also be tensor operations that directly correspond to database join and projection. We are thus led to define tensor projection and tensor join as follows.

The *projection* of a tensor $T$ onto a subset of its indices $\alpha$ is

$$\pi_\alpha(T) = \sum_\beta T_{\alpha\beta},$$

where $\beta$ is the set of $T$'s indices not in $\alpha$. ($\beta$'s elements may be interspersed with $\alpha$'s in any order.) In other words, the projection of $T$ onto $\alpha$ is the sum for each value of $\alpha$ of all the elements of $T$ with that value of $\alpha$. For example, a vector may be projected onto a scalar by summing all its elements, a matrix onto a column vector by summing each row into an element of the vector, a cubic tensor onto any one of its faces and then that face onto one of its edges and then onto a corner, etc. If the tensors are Boolean and the projection is followed by a step function, tensor projection reduces to database projection.

The *join* of two tensors $U$ and $V$ along a common set of indices $\beta$ is

$$(U \bowtie V)_{\alpha\beta\gamma} = U_{\alpha\beta}V_{\beta\gamma},$$

where $\alpha$ is the subset of $U$'s dimensions not in $V$ and similarly for $\gamma$ and $V$. (Again, $\alpha$, $\beta$ and $\gamma$ may be interspersed in any order.) In other words, the join of two tensors on a common subset of indices $\beta$ has one element for each pair of elements with the same value of $\beta$, and that element is their product. If $U$ has rank $r$, $V$ has rank $r'$, and $|\beta| = q$, $U \bowtie V$ has rank $r + r' - q$. When two tensors have no indices in common, their join reduces to their tensor product (Kronecker product for matrices). When they have all dimensions in common, it reduces to their elementwise product (Hadamard product for matrices). If the tensors are Boolean, tensor join reduces to database join.

A *tensor logic program* is a set of tensor equations. The left-hand side (LHS) of a tensor equation is the tensor being computed. The right-hand side (RHS) is a series of tensor joins followed by a tensor projection, and an optional univariate nonlinearity applied elementwise to the result. A tensor is denoted by its name followed by a list of indices, comma-separated and enclosed in square brackets. The join signs are left implicit, and the projection is onto the indices on the LHS. For example, a single-layer perceptron is implemented by the tensor equation

$$\texttt{Y} = \texttt{step}(\texttt{W}[\texttt{i}]\,\texttt{X}[\texttt{i}]),$$

where joining on $\texttt{i}$ and projecting it out implements the dot product of $\texttt{W}$ and $\texttt{X}$. Tensors can also be specified by listing their elements, e.g., $\texttt{W} = [0.2, 1.9, -0.7, 3]$ and $\texttt{X} = [0, 1, 1, 0]$. Typing $\texttt{Y}$? then causes $Y$ to be evaluated.

Notice that, like the einsum implementations in NumPy, PyTorch, etc., a tensor equation is more general than the original Einstein notation: the summed-over indices are those that do not appear in the LHS, and thus a repeated index may or may not be summed over. For example, the index $\texttt{i}$ in

$$\texttt{Y}[\texttt{i}] = \texttt{step}(\texttt{W}[\texttt{i}]\,\texttt{X}[\texttt{i}])$$

is not summed over. The implementation of a multilayer perceptron below utilizes this.

Tensor elements are 0 by default, and equations with the same LHS are implicitly summed. This both preserves the correspondence with logic programming and makes tensor logic programs shorter. Tensor types may be declared or inferred. Setting a tensor equal to a file reads the file into the tensor. Reading a text file results in a Boolean matrix whose $ij$th element is 1 if the $i$th position in the text contains the $j$th word in the vocabulary. (The matrix is not stored in this inefficient form, of course; more on this later.) For example, if the file is the string "Alice loves Bob" and it's read into the matrix $\texttt{M}$, the result is $\texttt{M}[0, \texttt{Alice}] = \texttt{M}[1, \texttt{loves}] = \texttt{M}[2, \texttt{Bob}] = 1$ and $\texttt{M}[\texttt{i}, \texttt{j}] = 0$ for all other $\texttt{i}$, $\texttt{j}$. (Notice how arbitrary constants, not just integers, can be used as indices.) Conversely, setting a file equal to a tensor writes the tensor to the file.

This is the entire definition of tensor logic. There are no keywords, other constructs, etc. However, it is convenient to allow some syntactic sugar that, while not increasing the expressiveness of the language, makes it more convenient to write common programs. For example, we may allow: multiple terms in one equation (e.g., $\texttt{Y} = \texttt{step}(\texttt{W}[\texttt{i}]\,\texttt{X}[\texttt{i}] + \texttt{C})$); index functions (e.g., $\texttt{X}[\texttt{i}, \texttt{t}+1] = \texttt{W}[\texttt{i}, \texttt{j}]\,\texttt{X}[\texttt{j}, \texttt{t}]$); normalization (e.g., $\texttt{Y}[\texttt{i}] = \texttt{softmax}(\texttt{X}[\texttt{i}])$); other tensor functions (e.g., $\texttt{Y}[\texttt{k}] = \texttt{concat}(\texttt{X}[\texttt{i}, \texttt{j}])$); alternate projection operators (e.g., $\texttt{max}=$ or $\texttt{avg}=$ instead of $+=$, which $=$ defaults to); slices (e.g., $\texttt{X}[4 : 8]$); and procedural attachment (predefined or externally defined functions). Tensor logic accepts Datalog syntax; denoting a tensor with parentheses instead of square brackets implies that it's Boolean. In particular,

a sparse Boolean tensor may be written more compactly as a set of facts. For example, the vector $X = [0, 1, 1, 0]$ can also be written as $X(1)$, $X(2)$, with $X(0)$ and $X(3)$ being implicitly 0. Similarly, reading the string "Alice loves Bob" into the matrix $M$ produces the facts $M(0, \texttt{Alice})$, $M(1, \texttt{loves})$ and $M(2, \texttt{Bob})$.)

As another simple example, a multilayer perceptron can be implemented by the equation

$$X[i, j] = sig(W[i, j, k] X[i-1, k]),$$

where $i$ ranges over layers and $j$ and $k$ over units, and $sig()$ is the sigmoid function. Different layers may be of different sizes (and the corresponding weight matrices are implicitly padded with zeros to make up the full tensor). Alternatively, we may use a different equation for each layer.

A basic recursive neural network (RNN) can be implemented by

$$X[i, *t+1] = sig(W[i, j] X[j, *t] + V[i, j] U[j, t]),$$

where $X$ is the state, $U$ is the input, $i$ and $j$ range over units, and $t$ ranges over time steps. The $*t$ notation indicates that $t$ is a virtual index: no memory is allocated for it, and successive values of the $X[i]$ vector are written to the same location. Since RNNs are Turing-complete (Siegelmann and Sontag, 1995), the implementation above implies that so is tensor logic.

### 3.2 Inference

Inference in tensor logic is carried out using tensor generalizations of forward and backward chaining.

In forward chaining, a tensor logic program is treated as linear code. The tensor equations are executed in turn, each one computing the tensor elements for which the necessary inputs are available; this is repeated until no new elements can be computed or a stopping criterion is satisfied.

In backward chaining, each tensor equation is treated as a function. The query is the top-level call, and each equation calls the equations for the tensors on its RHS until all the relevant elements are available in the data or there are no equations for the subqueries. In the latter case (sub)query elements are assigned 0 by default.

The choice of whether to use forward or backward chaining depends on the application.

### 3.3 Learning

Because there is only one type of statement in tensor logic—the tensor equation—automatically differentiating a tensor logic program is particularly simple. Univariate nonlinearity aside, the derivative of the LHS of a tensor equation with respect to a tensor on the RHS is just the product of the other tensors on the RHS. More precisely, if

$$Y[...] = T[...] X_1[...] \ldots X_n[...],$$

then

$$\frac{\partial Y[...]}{\partial T[...]} = X_1[...] \ldots X_n[...].$$

7

Special cases of this include: if $\mathtt{Y} = \mathtt{AX}$, then $\partial\mathtt{Y}/\partial\mathtt{X} = \mathtt{A}$; if $\mathtt{Y} = \mathtt{W[i]}\,\mathtt{X[i]}$, then $\partial\mathtt{Y}/\partial\mathtt{W[i]} = \mathtt{X[i]}$; and if $\mathtt{Y[i,j]} = \mathtt{M[i,k]}\,\mathtt{X[k,j]}$, then $\partial\mathtt{Y[i,j]}/\partial\mathtt{M[i,k]} = \mathtt{X[k,j]}$.

As a result, the gradient of a tensor logic program is also a tensor logic program, with one equation per equation and tensor on its RHS. Omitting indices for brevity, the derivative of the loss $\mathtt{L}$ with respect to a tensor $\mathtt{T}$ is

$$\frac{\partial\mathtt{L}}{\partial\mathtt{T}} = \sum_{\mathtt{E}} \frac{d\mathtt{L}}{d\mathtt{Y}} \frac{d\mathtt{Y}}{d\mathtt{U}} \prod_{\mathtt{U}\backslash\mathtt{T}} \mathtt{X},$$

where $\mathtt{E}$ are the equations whose RHSs $\mathtt{T}$ appears in, $\mathtt{Y}$ is the equation's LHS, $\mathtt{U}$ is its nonlinearity's argument, and $\mathtt{X}$ are the tensors in $\mathtt{U}$.

Learning a tensor logic program requires specifying the loss function and the tensors it applies to by means of one or more tensor equations. For example, to learn an MLP by minimizing squared loss on the last layer's outputs we can use the equation

$$\mathtt{Loss} = (\mathtt{Y[e]} - \mathtt{X[*e,N,j]})^2,$$

where $\mathtt{e}$ ranges over training examples and $\mathtt{j}$ over units, $\mathtt{Y}$ contains the target values, $\mathtt{X}$ is the MLP as defined above extended with a virtual index for examples, and $\mathtt{N}$ is the number of layers. By default, all tensors that are not supplied as training data will be learned, but the user can specify if any should remain constant (e.g., hyperparameters). The optimizer itself can be encoded in tensor logic, but typically a pre-supplied one will be used.

While backpropagation in traditional neural networks is applied to the same architecture for all training examples, in tensor logic the structure may effectively vary from example to example, since different equations may apply to different examples, and backpropagating through the union of all possible derivations of the example would be wasteful. Fortunately, a solution to this problem is already available in the form of backpropagation through structure, which for each example updates each equation's parameters once for each time it appears in the example's derivation (Goller and Küchler, 1996). Applying this to RNNs yields the special case of backpropagation through time (Werbos, 1990).

Learning a tensor logic program consisting of a fixed set of equations is quite flexible, since an equation can represent any set of rules with the same join structure. (E.g., an MLP can represent any set of propositional rules.) Further, tensor decomposition in tensor logic is effectively a generalization of predicate invention. For example, if the program to be learned is the equation

$$\mathtt{A]i,j,k]} = \mathtt{M[i,p]}\,\mathtt{M'[j,q]}\,\mathtt{M''[k,r]}\,\mathtt{C[p,q,r]}$$

and $\mathtt{A}$ is the sole data tensor, the learned $\mathtt{M}$, $\mathtt{M'}$, $\mathtt{M''}$ and $\mathtt{C}$ form a Tucker decomposition of $\mathtt{A}$; and thresholding them into Booleans turns them into invented predicates.

## 4. Implementing AI Paradigms

The implementations below use forward chaining unless otherwise specified.

## 4.1 Neural Networks

A convolutional neural network is an MLP with convolutional and pooling layers (LeCun et al., 1998). A convolutional layer applies a filter at every location in an image, and can be implemented by a tensor equation of the form

$$\texttt{Features}[\texttt{x},\texttt{y}] = \texttt{relu}(\texttt{Filter}[\texttt{dx},\texttt{dy},\texttt{ch}]\,\texttt{Image}[\texttt{x}+\texttt{dx},\texttt{y}+\texttt{dy},\texttt{ch}]),$$

where x and y are pixel coordinates, dx and dy are filter coordinates, and ch is the RGB channel. A pooling layer combines a block of nearby filters into one, and can be implemented by

$$\texttt{Pooled}[\texttt{x}/\texttt{S},\texttt{y}/\texttt{S}] = \texttt{Features}[\texttt{x},\texttt{y}],$$

where / is integer division and S is the stride. This results in the filter outputs at S successive positions in each dimension being summed into one. This implements sum-pooling; max-pooling would replace = with max=, etc. A convolutional and pooling layer can be combined into one with the equation $\texttt{Pooled}[\texttt{x}/\texttt{S},\texttt{y}/\texttt{S}] = \texttt{relu}(\ldots)$.

Graph neural networks (GNNs) apply deep learning to graph-structured data (e.g., social networks, molecules, metabolic networks, the Web) (Liu and Zhou, 2022). Table 1 shows the implementation of a simple GNN. The network's graph structure is defined by the $\texttt{Neig}(\texttt{x},\texttt{y})$ relation, with one fact for each adjacent $(\texttt{x},\texttt{y})$ pair; or equivalently, by the Boolean tensor $\texttt{Neig}[\texttt{x},\texttt{y}] = 1$ if x and x are adjacent and 0 otherwise. The main tensor is $\texttt{Emb}[\texttt{n},\texttt{l},\texttt{d}]$, containing the d-dimensional embedding of each node n in each layer l. Initialization sets each node's 0th-layer embeddings to its features (externally defined or learned). The network then carries out L iterations of message passing, one per layer. Each iteration starts by applying one or more perceptron layers to each node. (Table 1 shows one. To preserve permutation invariance, the weights $\texttt{W}_{\texttt{P}}$ do not depend on the node. Although there are no sub/superscripts in tensor logic, I will use them here for brevity.) The GNN then aggregates each node's neighbors' new features Z by joining the tensors $\texttt{Neig}(\texttt{n},\texttt{n}')$ and $\texttt{Z}[\texttt{n}',\texttt{l},\texttt{d}]$. For each node, this zeroes out the contributions of all non-neighbors; the result is the sum of the neighbors' features. (Internally, this can be done efficiently by iterating over the node's neighbors or by other methods; see Section 6.) The aggregated features may then be passed through another MLP (not shown), after which they are combined with the node's features using weights $\texttt{W}_{\texttt{Agg}}$ and $\texttt{W}_{\texttt{Self}}$ to produce the next layer's embeddings.

The most common applications of GNNs are node classification, edge prediction and graph classification. For two-class problems, each node is classified by doing the dot product of its final embedding with a weight vector, and passing the result through a sigmoid to yield the class probability. For multiclass problems (not shown), each node's final embedding is dotted with a weight vector for each class c, $\texttt{W}_{\texttt{Out}}[\texttt{c},\texttt{d}]$, yielding a vector of logits that is then passed through a softmax to yield the class probabilities $\texttt{Y}[\texttt{n},\texttt{c}]$. Edge prediction predicts whether there is an edge between each pair of nodes by dotting their embeddings and passing the result through a sigmoid. Graph classification produces a class prediction for the entire graph, and is identical to node classification save for the result being a scalar Y instead of a vector $\texttt{Y}[\texttt{n}]$.

Attention, the basis of large language models, is also straightforward to implement in tensor logic (Vaswani et al., 2017). Given an embedding matrix $\texttt{X}[\texttt{p},\texttt{d}]$, where p ranges over

Table 1: Graph neural networks in tensor logic

| Component | Equation |
|---|---|
| Graph structure | $\text{Neig}(x, y)$ |
| Initialization | $\text{Emb}[n, 0, d] = X[n, d]$ |
| MLP | $Z[n, l, d'] = \text{relu}(W_P[l, d', d] \, \text{Emb}[n, l, d])$, etc. |
| Aggregation | $\text{Agg}[n, l, d] = \text{Neig}(n, n') \, Z[n', l, d]$ |
| Update | $\text{Emb}[n, l+1, d] = \text{relu}(W_{\text{Agg}} \, \text{Agg}[n, l, d] + W_{\text{Self}} \, \text{Emb}[n, l, d])$ |
| Node classification | $Y[n] = \text{sig}(W_{\text{Out}}[d] \, \text{Emb}[n, L, d])$ |
| Edge prediction | $Y[n, n'] = \text{sig}(\text{Emb}[n, L, d] \, \text{Emb}[n', L, d])$ |
| Graph classification | $Y = \text{sig}(W_{\text{Out}}[d] \, \text{Emb}[n, L, d])$ |

items (e.g., positions in a text) and $d$ over embedding dimensions, the query, key and value matrices are obtained by multiplying $X$ by the corresponding weight matrices:

$$\text{Query}[p, d_k] = W_Q[d_k, d] \, X[p, d]$$
$$\text{Key}[p, d_k] = W_K[d_k, d] \, X[p, d]$$
$$\text{Val}[p, d_v] = W_V[d_v, d] \, X[p, d]$$

Attention can then be computed in two steps, the first of which compares the query at each position with each key:

$$\text{Comp}[p, p'.] = \text{softmax}(\text{Query}[p, d_k]) \, \text{Key}[p', d_k] \, / \, \text{sqrt}(D_k)),$$

where $\text{sqrt}(D_k)$ scales the dot products by the square root of the keys' dimension. The notation $p'.$ indicates that $p'$ is the index to be normalized (i.e., for each $p$, softmax is applied to the vector indexed by $p'$). The attention head then returns the sum of the value vectors weighted by the corresponding comparisons:

$$\text{Attn}[p, d_v] = \text{Comp}[p, p'] \, \text{Val}[p', d_v].$$

We can now implement an entire transformer with just a dozen tensor equations (Table 2). As we saw in Subsection 3.1, a text can be represented by the relation $X(p, t)$, stating that the $p$th position in the text contains the $t$th token. (Tokenization rules are easily expressed in Datalog, and are not shown.) The text's embedding $\text{EmbX}[p, d]$ is then obtained by multiplying $X(p, t)$ by the embedding matrix $\text{Emb}[t, d]$. The next equation implements positional encoding as in the original paper (Vaswani et al., 2017); other options are possible. (Incidentally, this equation also shows how conditionals and case statements can be implemented in tensor logic: by joining each expression with the corresponding condition.) The residual stream is then initialized to the sum of the text's embedding and the positional encoding.

Attention is implemented as described above, with two additional indices for each tensor: $b$ for the attention block and $h$ for the attention head. The attention heads' outputs are then concatenated, added to the residual stream and layer-normalized. MLP layers are implemented as before, with additional indices for block and position, and their outputs

Table 2: Transformers in tensor logic

| Component | Equation(s) |
|---|---|
| Input | $\mathtt{X(p, t)}$ |
| Embedding | $\mathtt{EmbX[p, d] = X(p, t)\, Emb[t, d]}$ |
| Pos. encoding | $\mathtt{PosEnc[p, d] = Even(d)\, sin(p/L^{d/D_e}) + Odd(d)\, cos(p/L^{d-1/D_e})}$ |
| Residual stream | $\mathtt{Stream[0, p, d] = EmbX[p, d] + PosEnc[p, d]}$ |
| Attention | $\mathtt{Query[b, h, p, d_k] = W_Q[b, h, d_k, d]\, Stream[b, p, d]}$, etc. |
|  | $\mathtt{Comp[b, h, p, p'.] = softmax(Query[b, h, p, d_k]\, Key[b, h, p', d_k]/sqrt(D_k))}$ |
|  | $\mathtt{Attn[b, h, p, d_v] = Comp[b, h, p, p']\, Val[b, h, p', d_v]}$ |
| Merge and | $\mathtt{Merge[b, p, d_m] = concat(Attn[b, h, p, d_v])}$ |
| layer norm | $\mathtt{Stream[b, p, d.] = lnorm(W_S[b, d, d_m]\, Merge[b, p, d_m] + Stream[b, p, d])}$ |
| MLP | $\mathtt{MLP[b, p] = relu(W_P[p, d]\, Stream[b, p, d])}$, etc. |
| Output | $\mathtt{Y[p, t.] = softmax(W_O[t, d]\, Stream[B, p, d])}$ |

are also normalized and added to the stream (not shown). Finally, the output (token probabilities) is obtained by dotting the stream with an output weight vector for each token and passing through a softmax.

### 4.2 Symbolic AI

A Datalog program is a valid tensor logic program. Therefore anything that can be done in Datalog can be done in tensor logic. This suffices to implement many symbolic systems, including reasoning and planning in function-free domains. Accommodating functions (as in Prolog) requires implementing unification in tensor logic (Lloyd, 1987).

### 4.3 Kernel Machines

A kernel machine can be implemented by the equation

$$\mathtt{Y[Q] = f(A[i]\, Y[i]\, K[Q, i] + B)},$$

where $\mathtt{Q}$ is the query example, $\mathtt{i}$ ranges over support vectors, and $\mathtt{f}()$ is the output nonlinearity (e.g., a sigmoid) (Schölkopf and Smola, 2002). The kernel $\mathtt{K}$ is then implemented by its own equation. For example, a polynomial kernel is

$$\mathtt{K[i, i'] = (X[i, j]\, X[i', j])^n},$$

where $\mathtt{i}$ and $\mathtt{i'}$ range over examples, $\mathtt{j}$ ranges over features, and $\mathtt{n}$ is the degree of the polynomial. A Gaussian kernel is

$$\mathtt{K[i, i'] = exp(-(X[i, j] - X[i', j])^2\, /\, Var)}.$$

(More precisely, $\mathtt{K}$ is the Gram matrix of the kernel with respect to the examples.) Structured prediction, where the output consists of multiple interrelated elements (Bakr et al., 2007), can be implemented by an output vector $\mathtt{Y[Q, k]}$ and equations stating the interactions among outputs and between outputs and inputs.

Table 3: Graphical models in tensor logic

| Component | Implementation |
|---|---|
| Factor | Tensor |
| Marginalization | Projection |
| Pointwise product | Join |
| Join tree | Tree-like program |
| P(Query|Evidence) | Prog(Q,E)/Prog(E) |
| Belief propagation | Forward chaining |
| Sampling | Selective projection |

### 4.4 Probabilistic Graphical Models

A graphical model represents the joint probability distribution of a set of random variables as a normalized product of factors,

$$P(X{=}x) = \frac{1}{Z} \prod_k \phi_k(x_{\{k\}}),$$

where each factor $\phi_k$ is a non-negative function of a subset of the variables $x_{\{k\}}$ and $Z = \sum_x \prod_k \phi_k(x_{\{k\}})$ (Koller and Friedman, 2009). If each factor is the conditional probability of a variable given its parents (predecessors in some partial ordering), the model is a Bayesian network and $Z = 1$.

Table 3 shows how the constructs and operations in discrete graphical models map directly onto those in tensor logic. A factor is a tensor of non-negative real values, with one index per variable and one value of the index per value of the variable. The unnormalized probability of a state $x$ is the product of the element in each tensor corresponding to $x$. A Bayesian network can thus be encoded in tensor logic using one equation per variable, stating the variable's distribution in terms of its conditional probability table (CPT) and the parents' distributions:

$$\mathtt{P_X[x] = CPT_X[x, par_1, ..., par_n]\, P_1[par_1]\ldots P_n[par_n]}.$$

Inference in graphical models is the computation of marginal and conditional probabilities, and consists of combinations of two operations: marginalization and pointwise product. The marginalization of a subset of the variables $Y$ in a factor $\phi$ sums them out, leaving a factor over the remaining variables $X$:

$$\phi'(X) = \sum_Y \phi(X, Y).$$

Marginalization is just tensor projection. The pointwise product of two potentials over subsets of variables $X$ and $Y$ combines them into a single potential over $X \cup Y$, and is the join of the corresponding tensors.

Every graphical model can be expressed as a join tree, a tree of factors where each factor is a join of factors in the original model. All marginal and conditional queries can

be answered in time linear in the size of the tree by successively marginalizing factors and pointwise-multiplying them with the parent's factor. A join tree is a tree-like tensor logic program, i.e., one in which no tensor appears in more than one RHS. As a result, linear-time inference can be carried out by backward chaining over this program. Specifically: the partition function $Z$ can be computed by adding the equation $\texttt{Z} = \texttt{T}[...]$ to the program, where $\texttt{T}[...]$ is the LHS of the root factor's equation, and querying $\texttt{Z}$; the marginal probability of evidence $P(E)$ can be computed by adding $E$ to the program as a set of facts, querying $\texttt{Z}$, and dividing by the original one; and the conditional probability of a query given evidence can be computed as $P(E) = P(Q, E)/P(E)$.

However, the join tree may be exponentially larger than the original model, necessitating approximate inference. The two most popular methods are loopy belief propagation and Monte Carlo sampling. Loopy belief propagation is forward chaining on the tensor logic program representing the model. Sampling can be implemented by backward chaining with selective projection (i.e., replacing a projection by a random subset of its terms).

## 5. Reasoning in Embedding Space

The most interesting feature of tensor logic is the new models it suggests. In this section I show how to perform knowledge representation and reasoning in embedding space, and point out the reliability and transparency of this approach.

Consider first the case where an object's embedding is a random unit vector. The embeddings can be stored in a matrix $\texttt{Emb}[\texttt{x}, \texttt{d}]$, where $\texttt{x}$ ranges over objects and $\texttt{d}$ over embedding dimensions. Multiplying $\texttt{Emb}[\texttt{x}, \texttt{d}]$ by a one-hot vector $\texttt{V}[\texttt{x}]$ then retrieves the corresponding object's embedding. If $\texttt{V}[\texttt{x}]$ is a multi-hot vector representing a set,

$$\texttt{S}[\texttt{d}] = \texttt{V}[\texttt{x}]\,\texttt{Emb}[\texttt{x}, \texttt{d}]$$

is the superposition of the embeddings of the objects in the set. The dot product

$$\texttt{D}[\texttt{A}] = \texttt{S}[\texttt{d}]\,\texttt{Emb}[\texttt{A}, \texttt{d}]$$

for some object $\texttt{A}$ is then approximately 1 if $\texttt{A}$ is in the set and approximately 0 otherwise (with standard deviation $\sqrt{N/D}$, where $N$ is the cardinality of the set and $D$ is the embedding dimension). Thresholding this at $\frac{1}{2}$ then tells us if $\texttt{A}$ is in the set with an error probability that decreases with the embedding dimension. This is similar to a Bloom filter (Bloom, 1970).

The same scheme can be extended to embedding a relation. Consider a binary relation $\texttt{R}(\texttt{x}, \texttt{y})$ for simplicity. Then

$$\texttt{EmbR}[\texttt{i}, \texttt{j}] = \texttt{R}(\texttt{x}, \texttt{y})\,\texttt{Emb}[\texttt{x}, \texttt{i}]\,\texttt{Emb}[\texttt{y}, \texttt{j}]$$

is the superposition of the embeddings of the tuples in the relation, where the embedding of a tuple is the tensor product of the embeddings of its arguments. This is a type of tensor product representation (Smolensky, 1990). It can be computed in time linear in $|\texttt{R}|$ by iterating over the tuples adding the corresponding tensor product to the result. The equation

$$\texttt{D}[\texttt{A}, \texttt{B}] = \texttt{EmbR}[\texttt{i}, \texttt{j}]\,\texttt{Emb}[\texttt{A}, \texttt{i}]\,\texttt{Emb}[\texttt{B}, \texttt{j}]$$

retrieves $R(A, B)$, i.e., $D[A, B]$ is approximately 1 if the tuple $(A, B)$ is in the relation and 0 otherwise, since

$$
\begin{aligned}
D[A, B] & = \text{EmbR}[i, j]\,\text{Emb}[A, i]\,\text{Emb}[B, j] \\
& = (R(x, y)\,\text{Emb}[x, i]\,\text{Emb}[y, j])\,\text{Emb}[A, i]\,\text{Emb}[B, j] \\
& = R(x, y)\,(\text{Emb}[x, i]\,\text{Emb}[A, i])\,(\text{Emb}[y, j]\,\text{Emb}[B, j]) \\
& \simeq R(A, B).
\end{aligned}
$$

The penultimate step is valid because einsums are commutative and associative. (In particular, the result does not depend on the order the tensors appear in, only on their index structure.) The last step is valid because the dot product of two random unit vectors is approximately 0.

By the same reasoning, the equation

$$
D[A, y] = \text{EmbR}[i, j]\,\text{Emb}[A, i]\,\text{Emb}[y, j]
$$

returns the superposition of the embeddings of the objects that are in relation $R$ to $A$, and

$$
D[x, y] = \text{EmbR}[i, j]\,\text{Emb}[x, i]\,\text{Emb}[y, j]
$$

returns the entire relation $R(x, y)$. $\text{EmbR}[i, j]$, $\text{Emb}[x, i]$ and $\text{Emb}[y, j]$ form a Tucker decomposition of the data tensor $D[x, y]$, with $\text{EmbR}[i, j]$ as the core tensor and $\text{Emb}[x, i]$ and $\text{Emb}[y, j]$ as the factor matrices.

The relation symbols themselves may be embedded. (E.g., $R$, $A$ and $B$ in $R(A, B)$ may all be embedded.) This results in a rank-3 tensor. Relations of arbitrary arity can be reduced to sets of *(relation, argument, value)* triples. Thus an entire database can be embedded as a single rank-3 tensor.

The next step is to embed rules. We can embed a Datalog rule by replacing its antecedents and consequents by their embeddings: if the rule is

$$
\text{Cons}(...) \leftarrow \text{Ant}_1(...), \ldots, \text{Ant}_n(...),
$$

its embedding is

$$
\text{EmbCons}[...] = \text{EmbAnt}_1[...]\ldots\text{EmbAnt}_n[...],
$$

where

$$
\text{EmbAnt}_1[...] = \text{Ant}_1(...)\,\text{Emb}[...]\ldots\text{Emb}[...],
$$

etc. Reasoning in embedding space can now be carried out by forward or backward chaining over the embedded rules. The answer to a query can be extracted by joining its tensor with its arguments' embeddings, as shown above for any relation. This gives approximately the correct result because each inferred tensor can be expressed as a sum of projections of joins of embedded relations, and the product $\text{Emb}[x, i]\,\text{Emb}[x', i]$ for each of its arguments is approximately the identity matrix. The error probability decreases with the embedding dimension, as before. To further reduce it, we can extract, threshold and re-embed the inferred tensors at regular intervals (in the limit, after each rule application).

The most interesting case, however, is when objects' embeddings are learned. The product of the embedding matrix and its transpose,

$$\mathtt{Sim}[\mathbf{x}, \mathbf{x}'] = \mathtt{Emb}[\mathbf{x}, \mathbf{d}]\,\mathtt{Emb}[\mathbf{x}', \mathbf{d}],$$

is now the Gram matrix measuring the similarity of each pair of objects by the dot product of their embeddings. Similar objects "borrow" inferences from each other, with weight proportional to their similarity. This leads to a powerful form of analogical reasoning that explicitly combines similarity and compositionality in a deep architecture.

If we apply a sigmoid to each equation,

$$\sigma(x, T) = \frac{1}{(1 + e^{-x/T})},$$

setting its temperature parameter $T$ to 0 effectively reduces the Gram matrix to the identity matrix, making the program's reasoning purely deductive. This contrasts with LLMs, which may hallucinate even at $T = 0$. It's also exponentially more powerful than retrieval-augmented generation (Jiang et al., 2025), since it effectively retrieves the deductive closure of the facts under the rules rather than just the facts.

Increasing the temperature makes reasoning increasingly analogical, with examples that are less and less similar borrowing inferences from each other. The optimal $T$ will depend on the application, and can be different for different rules (e.g., some rules may be mathematical truths and have $T = 0$, while others may serve to accumulate weak evidence and have a high $T$).

The inferred tensors can be extracted at any point during inference. This makes reasoning highly transparent, in contrast with LLM-based reasoning models. It's also highly reliable and immune to hallucinations at sufficiently low temperature, again in contrast with LLM-based models. At the same time, it has the generalization and analogical abilities of reasoning in embedding space. This could make it ideal for a wide range of applications.

## 6. Scaling Up

For large-scale learning and inference, equations involving dense tensors can be directly implemented on GPUs. Operations on sparse and mixed tensors can be implemented using (at least) one of two approaches.

The first is separation of concerns: operations on dense (sub)tensors are implemented on GPUs, and operations on sparse (sub)tensors are implemented using a database query engine, by treating (sub)tensors as relations. The full panoply of query optimization can then be applied to combining these sparse (sub)tensors. An entire dense subtensor may be treated as single tuple by the database engine, with an argument pointing to the subtensor. Dense subtensors are then joined and projected using GPUs.

The second and more interesting approach is to carry out all operations on GPUs, first converting the sparse tensors into dense ones via Tucker decomposition. This is exponentially more efficient than operating directly on the sparse tensors, and as we saw in the previous section, even a random decomposition will suffice. The price is that there will be a small probability of error, but this can be controlled by appropriately setting the embedding dimension and denoising results by passing them through step functions. Scaling up via

Tucker decompositions has the significant advantage that it combines seamlessly with the learning and reasoning algorithms described in previous sections.

## 7. Discussion

Tensor logic is likely to be useful beyond AI. Scientific computing consists essentially of translating equations into code, and with tensor logic this translation is more direct than with previous languages, often with a one-to-one correspondence between symbols on paper and symbols in code. In scientific computing the relevant equations are then wrapped in logical statements that control their execution. Tensor logic makes this control structure automatically learnable by relaxing the corresponding Boolean tensors to numeric ones, and optionally thresholding the results back into logic. The same approach is applicable in principle to making any program learnable.

Any new programming language faces a steep climb to wide adoption. What are tensor logic's chances of succeeding? AI programming is no longer a niche; tensor logic can ride the AI wave to wide adoption in the same way that Java rode the Internet wave. Backward compatibility with Python is key, and tensor logic lends itself well to it: it can initially be used as a more elegant implementation of einsum and extension of Python to reasoning tasks, and as it develops it can absorb more and more features of NumPy, PyTorch, etc., until it supersedes them.

Above all, adoption of new languages is driven by the big pains they cure and the killer apps they support, and tensor logic very much has these: e.g., it potentially cures the hallucinations and opacity of LLMs, and is the ideal language for reasoning, mathematical and coding models.

Fostering an open-source community around tensor logic will be front and center. Tensor logic lends itself to IDEs that tightly integrate coding, data wrangling, modeling and evaluation, and if it takes off vendors will compete to support it. It is also ideally suited to teaching and learning AI, and this is another vector by which it can spread.

Next steps include implementing tensor logic directly in CUDA, using it in a wide range of applications, developing libraries and extensions, and pursuing the new research directions it makes possible.

For more information on tensor logic, visit tensor-logic.org.

## Acknowledgments

## References

G. Bakr, T. Hofmann, B. Schölkopf, A. Smola, B. Taskar, and S. Vishwanathan, editors. *Predicting Structured Data*. MIT Press, Cambridge, MA, 2007.

B. Bloom. Space/time tradeoffs in hash coding with allowable errors. *Comm. ACM*, 13: 422–426, 1970.

C. Goller and A. Küchler. Learning task-dependent distributed representations by back-propagation through structure. In *Proc. Int. Conf. Neural Networks*, pp. 347–352, 1996.

S. Greco and C. Molinaro. *Datalog and Logic Databases*. Morgan & Claypool, San Rafael, CA, 2016.

P. Jiang, S. Ouyang, Y. Jiao, M. Zhong, R. Tian, and J. Han. Retrieval and structuring augmented generation with large language models. In *Proc. Int. Conf. Knowl. Disc. & Data Mining*, pp. 6032–6042, 2025.

D. Koller and N. Friedman. *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, Cambridge, MA, 2009.

N. Lavrač and S. Džeroski. *Inductive Logic Programming: Techniques and Applications*. Ellis Horwood, Chichester, UK, 1994.

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proc. IEEE*, 86:2278–2324, 1998.

Z. Liu and J. Zhou. *Introduction to Graph Neural Networks*. Morgan & Claypool, San Rafael, CA, 2022.

J. W. Lloyd. *Foundations of Logic Programming* (2nd ed.). Springer, Berlin, Germany, 1987.

S. Rabanser, O. Shchur, and S. Günnemann. Introduction to tensor decompositions and their applications in machine learning. arXiv:1711.1078, 2017.

T. Rocktäschel. Einsum is all you need – Einstein summation in deep learning. https://rockt.ai/2018/04/30/einsum, 2018.

A. Rogozhnikov. Einops: Clear and reliable tensor manipulations with Einstein-like notation. In *Proc. Int. Conf. Learn. Repr.*, 2022.

B. Schölkopf and A. J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, Cambridge, MA, 2002.

H. Siegelmann and E. Sontag. On the computational power of neural nets. *J. Comp. & Sys. Sci.*, 50:132–150, 1995.

P. Smolensky. Tensor product variable binding and the representation of symbolic structures in connectionist systems. *Artif. Intel.*, 46:159–216, 1990.

A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. *Adv. Neural Inf. Proc. Sys.*, 30:5998–6008, 2017.

P. Werbos. Backpropagation through time: What it does and how to do it. *Proc. IEEE*, 78:1550–1560, 1990.