# Evaluating Relaxations of Logic for Neural Networks: A Comprehensive Study

**Mattia Medina Grespan** , **Ashim Gupta** and **Vivek Srikumar**

University of Utah

{mattiamg,ashim,svivek}@cs.utah.edu

## Abstract

Symbolic knowledge can provide crucial inductive bias for training neural models, especially in low data regimes. A successful strategy for incorporating such knowledge involves relaxing logical statements into sub-differentiable losses for optimization. In this paper, we study the question of how best to relax logical expressions that represent labeled examples and knowledge about a problem; we focus on sub-differentiable t-norm relaxations of logic. We present theoretical and empirical criteria for characterizing which relaxation would perform best in various scenarios. In our theoretical study driven by the goal of preserving tautologies, the Łukasiewicz t-norm performs best. However, in our empirical analysis on the text chunking and digit recognition tasks, the product t-norm achieves best predictive performance. We analyze this apparent discrepancy, and conclude with a list of best practices for defining loss functions via logic.

## 1 Introduction

Neural networks are remarkably effective across many domains and tasks; but their usefulness is limited by their data hungriness. A promising direction towards alleviating this concern involves augmenting learning with rules written in first-order logic [Rocktäschel *et al.*, 2015; Li and Srikumar, 2019; Li *et al.*, 2019; Fischer *et al.*, 2019, inter alia]. To make rules amenable with gradient-based learning, this approach calls for relaxing the logical operators to define sub-differentiable loss terms. A systematic method to perform this relaxation uses a well-studied family of binary operators, namely *triangular norms* or *t-norms* [Klement *et al.*, 2013]. Different t-norms define different $[0, 1]$-valued interpretations of the Boolean operators.

There are infinitely many such t-norm logics, but the most commonly used ones are: Product [Rocktäschel *et al.*, 2015; Li *et al.*, 2019; Asai and Hajishirzi, 2020], Gödel [Minervini *et al.*, 2017] and Łukasiewicz [Bach *et al.*, 2017]. While the usefulness of such relaxations is well established, the questions of how they compare against each other, and even how such a comparison should be defined remain open.

This paper is a first step towards answering these important questions by analyzing the three t-norm relaxations and their variants. To do so, we define the criteria for quantifying the goodness of a t-norm based relaxation. On the theoretical side, we rank logic relaxations by their *consistency*, i.e., their ability to preserve the truth of tautologies. On the empirical side, we use a principled approach to construct loss functions using t-norms that subsume traditional loss functions like cross-entropy, and study how well different relaxations compare with standard gradient-based learning. We report the results of experiments on two tasks: jointly recognizing digits and predicting the results of arithmetic operators, and text chunking. Both the theoretical and empirical criteria concur in the recommendation that a variant of the Product t-norm ($\mathcal{R}$-Product) is most suitable for introducing logical rules into neural networks.

In summary, the main contributions of this work are:

- We define theoretical and empirical properties that any relaxation of logic should have to be useful for learning.
- We define the consistency of relaxations to rank them in terms of their ability to preserve tautologies.
- We present empirical comparisons of logic relaxations on two tasks where labeled examples and declaratively stated knowledge together inform neural models.

## 2 Problem Statement and Notation

Several recent efforts have shown the usefulness of declarative knowledge to guide neural network learning towards improving model quality. While different approaches exist for incorporating such rules into neural models [Xu *et al.*, 2018, for example], a prominent strategy involves relaxing logic to the real regime using the well-studied t-norm relaxations.

**Triangular norms** (t-norms) arose in the context of probabilistic metric spaces [Menger, 1942], and for our purposes, represent a relaxation of the Boolean conjunction which agrees with the definition of conjunctions for $\{0, 1\}$-inputs. Given such a relaxation and a relaxation of $\neg X$ as $1 - x$, we have two axiomatic approaches for defining implications. The first (called $\mathcal{S}$-logics) treats implications as disjunctions (i.e., $X \rightarrow Y = \neg X \vee Y$), while the second (called $\mathcal{R}$-logics) defines implications axiomatically. We refer the reader to Klement *et al.* [2013] for a detailed treatment of t-norms.

Table 1 shows the set of t-norms we consider here; these have been used in recent literature to inject knowledge into neural networks. However, despite their increasing prevalence, there is no consensus on which t-norm to employ. A survey of recent papers reveals the use of $\mathcal{S}$-Product [Rocktäschel *et al.*, 2015], $\mathcal{R}$-Product [Li *et al.*, 2019; Asai and Hajishirzi, 2020], Łukasiewicz [Bach *et al.*, 2017], $\mathcal{S}$-Gödel [Minervini *et al.*, 2017] and even a mixture of the $\mathcal{S}$-Gödel and $\mathcal{R}$-Product [Li *et al.*, 2020] t-norms.

*How do these relaxations of logic compare against each other?* In this work, we answer this question from both theoretical and empirical perspectives.

## 2.1 Learning From Logic: The Setup

To compare the different relaxations of logic on an even footing, let us first see a general recipe for converting logic rules to loss functions for a given relaxation[1]. We will use the task of recognizing handwritten digits as a running example.

Given a labeling task, we can represent the fact that the label for an instance $x$ is $y$ as a predicate, say $\texttt{Label}(x,y)$. In our running example, we could define a predicate $\texttt{Digit}(x,y)$ to denote that an image $x$ represents the digit $y$. We could also define a predicate $\texttt{Sum}(x_1, x_2, y)$ to denote the fact that the digits in two images $x_1$ and $x_2$ add up to the digit $y \pmod{10}$. From this perspective, we can treat classifiers as predicting probabilities that the predicates hold.

To train such classifiers, we typically have a training set $D$ of labeled examples $(x, y)$. In our notation, each such example can be represented as the predicate $\texttt{Label}(x,y)$ and the training set is a conjunction of such predicates

$$\bigwedge_{(x,y) \in D} \texttt{Label}(x,y) \qquad (1)$$

Sometimes, instead of a labeled dataset, we may have a constraint written in logic. In our running example, from the definition of the $\texttt{Digit}$ and $\texttt{Sum}$ predicates, we know that

$$\forall x_1, x_2, \bigwedge_{y_1, y_2} \texttt{Digit}(x_1, y_1) \wedge \texttt{Digit}(x_2, y_2) \rightarrow$$
$$\texttt{Sum}(x_1, x_2, (y_1 + y_2) \bmod 10) \qquad (2)$$

Note that such a constraint need not depend on labeled examples, and should hold irrespective of what labels the examples should be assigned. In general, given a large unlabeled set of examples denoted by $x \in U$, we can write constraints

$$\bigwedge_{x \in U} \texttt{C}(x) \qquad (3)$$

These constraints may be composite formulas constructed with predicates as shown in our running example above.

From this standpoint, we can envision the goal of learning as that of ensuring that the formulas representing labeled examples ( equation 1) and constraints ( equation 3) hold. Since we are treating classifiers as predicting probabilities that the atomic predicates hold, we can equivalently state the learning problem as that of finding model parameters that maximize

---

[1]The setup described here is implicitly present in Rocktäschel *et al.* [2015], Li *et al.* [2019], and others.

| | $\mathcal{S}$-Gödel | $\mathcal{R}$-Product | Łukasiewicz |
|---|---|---|---|
| $\wedge$ | $\min(x,y)$ | $x \cdot y$ | $\max(0, x+y-1)$ |
| $\neg$ | $1-x$ | $1-x$ | $1-x$ |
| $\vee$ | $\max(x,y)$ | $x + y - x \cdot y$ | $\min(1, x+y)$ |
| $\rightarrow$ | $\max(1-x, y)$ | $\begin{cases} 1 & \text{if } x \leq y \\ \frac{y}{x} & \text{otherwise} \end{cases}$ | $\min(1, 1-x+y)$ |

Table 1: T-norm relaxations studied in this work. Here, the letters $x$ and $y$ denote the relaxed truth values of the arguments of the formulas. In the implication definitions, $x$ and $y$ denote the antecedent and the consequent respectively. The table does not show $\mathcal{S}$-Product: it agrees with $\mathcal{R}$-Product for all the connectives except the implication, defined as $1 - x + x \cdot y$. We are defining $\mathcal{R}$-Product using its SBL$_\sim$ extension with involutive negation (See Esteva *et al.* [2000]).

the value of a *relaxation* of the conjunction of the formulas representing the data and constraints. In other words, we can use logic to define loss functions.

In this declarative learning setting, we have the choice of using *any* models (e.g., CNNs) for our predicates, and *any* relaxation of logic. If we only have labeled examples, and we use one of the Product relaxations, we recover the widely used cross-entropy loss [Li *et al.*, 2019; Giannini *et al.*, 2019].

**Notation.** We use upper case letters (e.g., P, $\texttt{Digit}$) to represent Booleans, and lower cased letters (e.g., $p$, $digit$) to represent their relaxations. In some places, for clarity, we use square brackets to denote the relaxation of a Boolean formula A (i.e., [A]=a).

## 3 Validity of Relaxed Logic

In this section, we propose three criteria that a logic relaxation should satisfy to be useful for learning.

**Consistency.** The language of logic can declaratively introduce domain knowledge, invariants, or even reasoning skills into neural networks. However, to admit reasoning, tautologies should always hold. That is, the truth value of any tautology should be 1 irrespective of the value of its constituent atomic predicates. Equivalently, the integral of the relaxation of a tautology over the domain of its atomic predicates should be 1. We can formalize this intuition.

**Definition 1.** *Let T be a tautology in predicate logic formed with a set of atomic predicates $\mathcal{T}$, and let L be a logic relaxation. The* consistency *of T in L, denoted as $\kappa^L(\texttt{T})$, is defined as*

$$\kappa^L(\texttt{T}) = \int_0^1 [\texttt{T}] \, \mathrm{d}\mathcal{T} \qquad (4)$$

If the consistency $\kappa^L(\texttt{T}) = 1$, we will say that the tautology T is consistent under the relaxation $L$.

**Self-consistency.** Every Boolean statement implies itself. That is, the statement P $\leftrightarrow$ P is a tautology for any P. This observation gives us the definition of *self-consistency* of a formula under a given relaxation.

**Definition 2.** *Let P be any Boolean formula in predicate logic with a set of atomic predicates $\mathcal{P}$, and let L be a logic relaxation. The* self-consistency *of $\mathcal{P}$ in the logic L, denoted as*

$\kappa_S^L(P)$, *is defined as*

$$\kappa_S^L(P) = \kappa^L(P \leftrightarrow P) = \int_0^1 [P \leftrightarrow P] \, \mathrm{d}\mathcal{P} \quad (5)$$

If $\kappa_S^L(P) \neq 1$ we will say that formula P is not self-consistent under a relaxation $L$. Since we consider a dataset to be a conjunction of facts ( equation 1), the self-consistency of large conjunctions allows us to judge whether a dataset implies itself under a relaxation.

**Sub-differentiability.** Since our eventual goal is to relax declaratively stated knowledge to train neural networks, the relaxations should admit training via backpropagation. As a result, the functions defining the relaxed logical operators should at least be sub-differentiable.

In sum, we consider a logic relaxation to be valid if the following properties hold:

(P1) It must be sub-differentiable over the interval $[0, 1]$.

(P2) It must be consistent for any tautology.

(P3) It must be self-consistent for any Boolean formula.

All the relaxations we study here satisfy property P1[2]. Even among sub-differentiable relaxations, some may be easier than others to learn using gradient-based approaches. We consider this empirical question in §5. Properties P2 and P3 do not always hold, and we will prefer relaxations that have higher values of consistency and self-consistency.

## 4 Truth Preservation of Relaxations

In this section, we will assess the relaxations from Table 1 with respect to properties P2 and P3.

### 4.1 Consistency

Property P2 expects *every* tautology to be consistent under a valid logic relaxation. Since predicate logic admits infinitely many tautologies, we will use a representative set of tautologies for our evaluation. This set contains the *Axiom Schemata* of the Hilbert proof system for predicate logic[3], the primitive propositions, and a set of elementary properties defined by Russell and Whitehead [1910].

Table 2 shows the consistency for each tautology for our t-norm relaxations. Comparing results across the columns, we see that Łukasiewicz and $\mathcal{R}$-Product t-norm logics *preserve truth* the most across our representative set. In general, we find $\mathcal{R}$ logics to be better at preserving truth than $\mathcal{S}$ logics.

**Example 1.** *Consistency of the tautology* $A \to A$ *using* $\mathcal{S}$-*Product.*

$$[A \to A] = 1 - a + a^2$$

*By the definition of consistency ( equation 4), we have*

$$\kappa^{\mathcal{S}\text{-}Product}(A \to A) = \int_0^1 1 - a + a^2 \, \mathrm{d}a = \frac{5}{6} \approx 0.83$$

| Tautologies | $\mathcal{S}$-Prod. | $\mathcal{S}$-Gödel | Łuka. | $\mathcal{R}$-Prod. |
|---|---|---|---|---|
| **Axiom Schemata** | | | | |
| P → (Q → P) | 0.92 | 0.79 | 1 | 1 |
| (P → (Q → R)) → ((P → Q) → (P → R)) | 0.88 | 0.75 | 0.96 | 0.93 |
| (¬P → ¬Q) → (Q → P) | 0.86 | 0.75 | 1 | 0.88 |
| **Primitive Propositions** | | | | |
| (P ∨ P) → P | 0.75 | 0.75 | 0.75 | 0.69 |
| Q → (P ∨ Q) | 0.92 | 0.79 | 1 | 1 |
| (P ∨ Q) → (Q ∨ P) | 0.86 | 0.75 | 1 | 1 |
| (P ∨ (Q ∨ R)) → (Q ∨ (P ∨ R)) | 0.91 | 0.78 | 1 | 1 |
| (Q → R) → ((P ∨ Q) → (P ∨ R)) | 0.90 | 0.76 | 1 | 1 |
| **Law of excluded middle** | | | | |
| P ∨ ¬P | 0.83 | 0.75 | 1 | 0.83 |
| **Law of contradiction** | | | | |
| ¬(P ∧ ¬P) | 0.83 | 0.75 | 1 | 0.83 |
| **Law of double negation** | | | | |
| P ↔ ¬(¬P) | 0.70 | 0.75 | 1 | 1 |
| **Principles of transposition** | | | | |
| (P ↔ Q) ↔ (¬P ↔ ¬Q) | 0.61 | 0.67 | 1 | 0.59 |
| **Laws of tautology** | | | | |
| P ↔ (P ∧ P) | 0.69 | 0.75 | 0.75 | 0.50 |
| P ↔ (P ∨ P) | 0.69 | 0.75 | 0.75 | 0.69 |
| **De Morgan's Laws** | | | | |
| (P ∧ Q) ↔ ¬(¬P ∨ ¬Q) | 0.75 | 0.75 | 1 | 1 |
| ¬(P ∧ Q) ↔ ¬(¬P ∨ ¬Q) | 0.75 | 0.75 | 1 | 1 |

Table 2: Consistencies of a (subset) of representative set of tautologies under different logic relaxations. We see here that the $\mathcal{R}$-Product and Łukasiewicz relaxations are generally more consistent, suggesting that they are preferable to the other two relaxations. Other tautologies we examined show the same trends.

### 4.2 Self-consistency

The validity property P3 states that *every* well-formed Boolean formula should be self-consistent for a relaxation of the logic to be valid. It turns out that the definition of implications for any $\mathcal{R}$ logic (including Łukasiewicz) guarantees the self-consistency of every formula.

**Proposition 1.** *Every formula is self-consistent under any $\mathcal{R}$-logic relaxation.*

This follows directly from the definition of the t-norm and the properties of residua. The appendix has a proof sketch.

However, the same is not true for $\mathcal{S}$-logics. For example, for the conjunction P $= A \wedge B$, the self-consistency under the Gödel relaxation $\kappa_S^{\mathcal{S}\text{-Gödel}}(P) = 0.75$ and under the Product relaxation, we have $\kappa_S^{\mathcal{S}\text{-Product}}(P) \approx 0.74$. These results suggest that the $\mathcal{R}$-Product and Łukasiewicz relaxations are preferable from the perspective of property P3 as well.

Intriguingly, we find that the $\mathcal{S}$-Product logic is *eventually* self-consistent for large monotone conjunctions.

**Proposition 2.** *Let* $A = \bigwedge_{i=1}^n A_i$ *be a conjunction consisting of* $n$ *atomic predicates* $A_1, A_2, \ldots, A_n$. *The self-consistency of* $A$ *is given by* $\kappa_S^{\mathcal{S}\text{-Product}}(A) = 1 - \frac{2}{2^n} + \frac{3}{3^n} - \frac{2}{4^n} + \frac{1}{5^n}$.

*Proof.* By induction over the size of the conjunction $n$. □

We see that, as $n \to \infty$, the self-consistency of a monotone conjunction approaches 1. In other words, large conjunctions (e.g., representing large datasets) are essentially self-consistent under the $\mathcal{S}$-Product relaxation.

# 5 Empirical Comparisons of Relaxed Logic

In this section, we empirically study the differences between the different logic relaxations using two tasks: recognizing digits and arithmetic operations, and text chunking. In both tasks, we set up the learning problem in terms of logic, and compare models learned via different logic relaxations.[4]

## 5.1 Recognizing Digits and Arithmetic Operations

These experiments build upon our running example from §2. We seek to categorize handwritten digits; i.e., we learn the predicate `Digit`. In addition, we also seek to predict the sum and product (modulo 10) of two handwritten digit images. These correspond to the predicate `Sum` we have seen, and a new analogous predicate `Product`.

We use the popular MNIST dataset [LeCun, 1998] for our experiments, but *only* to supervise the `Digit` classifier. Rather than directly supervising the other two classifiers, we use coherence constraints over *unlabeled* image pairs that connect them to the `Digit` model. The constraint for the `Sum` classifier is shown in equation 2, and the one for the `Product` classifier is similarly defined.

We set up the learning problem as defined in §2 and compare performance across different relaxations.

### Data and Setup

We partition the 60k MNIST training images into TRAIN and DEV sets, with 50k and 10k images respectively. To supervise the `Digit` model, we sample 1k, 5k and 25k labeled images from TRAIN to form three DIGIT sets. The coherence constraints are grounded in 5k *unlabeled* image pairs consisting of images from TRAIN that are not in any DIGIT set, giving us the PAIR dataset.

For evaluating the `Digit` model, we use the original 10k TEST examples from MNIST. For the development and evaluation of the operator models, we sample random image pairs from DEV and TEST to create the PairDEV and PairTEST sets respectively. The ground truth `Sum` and `Product` labels for these image pairs can be computed by the sum and product modulo 10 of the image labels.

We use CNNs as the `Digit`, `Sum` and `Product` models. For the operator models, we concatenated the two images to get CNN inputs. To jointly train these models using the labeled DIGIT data and the unlabeled PAIR datasets, we define a loss function by relaxing the conjunction of `Digit` predicates over the DIGIT examples and the coherence constraints over the PAIR examples. That is, learning requires minimizing:

$$
-\left[ \left( \bigwedge_{(x,y) \in D} \mathtt{Digit}(x,y) \right) \right.
$$
$$
\wedge \left( \bigwedge_{PAIR} \mathtt{Sum\ Coherence} \right) \qquad (6)
$$
$$
\left. \wedge \left( \bigwedge_{PAIR} \mathtt{Product\ Coherence} \right) \right]
$$

---

[4]Our PyTorch [Paszke *et al.*, 2019] code is archived at https://github.com/utahnlp/neural-logic

|  | 1000 | 5000 | 25000 |
|---|---|---|---|
| $\mathcal{S}$-Gödel | 95.0 (0.3) | 97.4 (0.1) | 97.7 (0.1) |
| $\mathcal{S}$-Product | 95.1 (0.1) | 98.2 (0.0) | 99.0 (0.0) |
| $\mathcal{R}$-Product | **96.3** (0.1) | **98.4** (0.1) | **99.2** (0.0) |
| Łukasiewicz | 95.8 (0.1) | 98.0 (0.1) | 99.1 (0.0) |

Table 3: `Digit` accuracies (and standard deviations) from jointly training `Digit` on DIGIT sizes 1k, 5k, 25k and operators (`Sum`,`Prod`) on PAIR size 5k.

|  | 1000 | 5000 | 25000 |
|---|---|---|---|
| $\mathcal{S}$-Gödel | 87.3 (0.5) | **91.1** (0.1) | 91.5 (0.0) |
| $\mathcal{S}$-Product | 76.9 (0.8) | 88.6 (0.6) | 90.1 (0.0) |
| $\mathcal{R}$-Product | **88.0** (0.3) | 90.8 (0.3) | **91.8** (0.0) |
| Łukasiewicz | 75.9 (3.1) | 84.5 (2.8) | 82.3 (0.3) |

Table 4: Average of `Sum` and `Product` accuracies (and standard deviations) from jointly training `Digit` on DIGIT sizes 1k, 5k, 25k and operators on PAIR size 5k.

In practice, we found that it is important to use a hyperparameter $\lambda$ that weights the relaxed coherence constraints in the loss. We used the DEV sets for hyperparameter tuning using the average of the accuracy of the `Digit` classifier and the coherences of the other two.

### Results

Table 3 reports accuracies for the `Digit` classifier trained with different sizes of DIGIT, and the coherence constraints instantiated over the 5k PAIR examples. We observe that the $\mathcal{R}$-Prod relaxation dominates across all settings, with higher gains when there are fewer labeled examples. (The bold entries in this and other tables are statistically significantly better than the other relaxations at $p < 0.05$. The accuracies are averages from three runs with different seeds along with the standard deviation.) Table 4 reports the average of `Sum`, and `Prod` accuracies for the same settings, and Table 5 shows the fraction of PairTEST examples where the coherence constraints are satisfied. From these results, we see that the $\mathcal{R}$-Product and $\mathcal{S}$-Gödel relaxations offer the best accuracies.

Interestingly, Łukasiewicz is the least accurate relaxation. However, the losses compiled using Łukasiewicz and $\mathcal{S}$-Gödel t-norms were unstable, and the results shown here were achieved with additional assumptions. We defer this technical discussion to §5.3, and conclude that the stability and accuracy of the $\mathcal{R}$-Product relaxation suggest that it is the most suitable relaxation for this class of problems.

### Joint Learning vs. Pipelines

In the coherence constraints in equation 2, if we knew both the `Digit` terms, we can deterministically compute the value of `Sum`. This suggests a pipeline strategy for training, where we can train the `Digit` classifier alone, and use it to assign (noisy) labels to the unlabeled PAIR data. Subsequently, we can train the `Sum` and `Product` models independently. How does the joint training strategy compare to this pipeline?

Importantly, we should note that both the joint and the pipeline strategies instantiate the declarative learning ap-

|  | 1000 | 5000 | 25000 |
|---|---|---|---|
| $\mathcal{S}$-Gödel | 86.4 (0.6) | 91.4 (0.1) | 91.9 (0.1) |
| $\mathcal{S}$-Product | 79.0 (0.7) | 89.3 (0.5) | 90.3 (0.0) |
| $\mathcal{R}$-Product | **89.1** (0.3) | **91.5** (0.4) | **92.1** (0.1) |
| Łukasiewicz | 77.4 (3.0) | 85.4 (2.7) | 82.6 (0.2) |

Table 5: Average of `Sum` and `Prod` Coherence accuracies (and standard deviations) from jointly training `Digit` on DIGIT sizes 1k, 5k, 25k and operators on PAIR size 5k.

|  | 1000 | 5000 | 25000 |
|---|---|---|---|
| $\mathcal{S}$-Gödel | 95.2 | 97.6 | 97.59 |
| $\mathcal{R}/\mathcal{S}$-Product | **96.7** | **98.4** | **99.12** |
| Łukasiewicz | 96.5 | 98.5 | 98.76 |

Table 6: Pipelined `Digit` accuracies trained on DIGIT sizes 1k, 5k, 25k. Standard deviation is 0.0 for every entry.

proach outlined in §2 where logical facts are compiled into an optimization learning problem via the logic relaxations. However, in the pipeline, the coherence constraints are not explicitly involved in the loss, because the noisy labels already satisfy them. Since conjunctions are identical for the $\mathcal{S}$ and $\mathcal{R}$ relaxations, they give the same results.

Tables 6, 7, and 8 show the results of these pipelined experiments. We observe the same trends as in the joint learning experiments: $\mathcal{R}$-Product achieves the highest performance.

Comparing the joint learning and the pipeline results, we observe that the latter are slightly better across relaxations and data settings. As mentioned above, however, pipelining is only viable here because of the special form of our constraint.

## 5.2 Text Chunking

Our second set of experiments use the NLP task of text chunking using the CoNLL 2000 dataset [Sang and Buchholz, 2000]. This task illustrates how relaxed logic can be used to derive loss functions for sequential inputs and outputs.

Text chunking is a sequence tagging problem, where each word of a sentence is assigned a phrase type. For example, in the sentence '*John is playing in the park.*', the word '*John*' is labeled as `B-NP`, indicating that it starts a noun phrase (`NP`). We use the standard BIO labeling scheme: `B-X` indicates the start of a new phrase labeled `X`, `I-X` indicates the continuation of a phrase labeled `X`, and `O` marks words that do not belong to any of the predefined phrase types.

In the case of text chunking, the input $x$ is a sequence of words and correspondingly, output $y$ is a sequence of labels (phrase types). Consequently, each position in the sequence is associated with a predicate. For an input $x$ with $n$ words, we have $n$ predicates, one per word.

Using our notation from §2 we define the predicate, `Tag`$(x_i, y_i)$, to denote that $i^{th}$ word in input $x$ is assigned the label $y_i$. For our preceding example, *John* being assigned the label `B-NP` corresponds to the predicate `Tag`(John, B-NP).

### Constraints

For each position in the sequence, we have constraints defining pairwise label dependencies. If a word in a sentence has

|  | 1000 | 5000 | 25000 |
|---|---|---|---|
| $\mathcal{S}$-Gödel | 88.3 (0.3) | 90.9 (0.1) | 91.3 (0.1) |
| $\mathcal{R}/\mathcal{S}$-Product | **89.7** (0.1) | **92.2** (0.0) | **93.0** (0.0) |
| Łukasiewicz | 83.0 (1.9) | 86.9 (3.6) | 88.6 (0.4) |

Table 7: Average of Pipelined `Sum` and `Product` accuracies (and standard deviations) trained on PAIR size 5k (noisy) labeled with `Digit` model trained on DIGIT sizes 1k, 5k, 25k.

|  | 1000 | 5000 | 25000 |
|---|---|---|---|
| $\mathcal{S}$-Gödel | 86.7 (0.3) | 90.6 (0.1) | 91.3 (0.0) |
| $\mathcal{R}/\mathcal{S}$-Product | **89.9** (0.2) | **92.5** (0.0) | **92.8** (0.0) |
| Łukasiewicz | 83.5 (1.8) | 87.2 (3.5) | 88.7 (0.4) |

Table 8: Average of Pipelined `Sum` and `Prod` Coherence accuracies (and standard deviations) trained on PAIR size 5k (noisy) labeled with `Digit` model trained on DIGIT sizes 1k, 5k, 25k.

a B/I label of a certain phrase type, then the next word cannot have a I label of a *different* phrase type. For example, we have

$C_{i,1}$:  $\forall i$, `Tag`$(x_i, \text{B-NP}) \rightarrow \neg \text{Tag}(x_{i+1}, \text{I-VP})$
$C_{i,2}$:  $\forall i$, `Tag`$(x_i, \text{I-NP}) \rightarrow \neg \text{Tag}(x_{i+1}, \text{I-VP})$

In general, given a labeled dataset $D$, and a set $C$ of $k$ constraints for each word position, we can write the conjunction of all constraints as:

$$\bigwedge_{x,y \in D} \left\{ \bigwedge_i \left( \text{Tag}(x_i, y_i) \bigwedge_k C_{i,k} \right) \right\} \tag{7}$$

We can now use this Boolean formula to state the goal of learning as finding the model parameters of a neural network that maximizes the value of the relaxation derived for each t-norm. We used a bidirectional LSTM over GloVe embeddings [Pennington *et al.*, 2014] to instantiate `Tag`.

### Experiments and Results

We compare the four t-norms on two settings. First, in purely supervised learning, the model only learns from labeled examples. Second, we augment the labeled dataset with the constraints described above to study the effectiveness of incorporating simple constraints on outputs. For both of these settings, we also study the models in a low data regime with training data restricted to 10%. Following previous work [Sang and Buchholz, 2000], we use the F1 score as our evaluation metric.

We report performances of the models trained with different t-norms in Table 9. We observe that for supervised learning (top rows), both variants of product t-norm outperform the other two t-norms by a significant margin. We again observe that for purely supervised learning, both $\mathcal{R}$-Product and $\mathcal{S}$-Product produce identical results since no implication constraints are used. Further, $\mathcal{S}$-Gödel performs the worst among the four t-norms. Note that this observation is consistent with our analysis from Table 2 where $\mathcal{S}$-Gödel was found to be least consistent with tautologies.

Let us now look at the results of augmenting supervised learning with simple output constraints (Table 9, bottom

| % Train | $\mathcal{S}$-Gödel | $\mathcal{S}$-Prod. | $\mathcal{R}$-Prod. | Łukasiewicz |
|---|---|---|---|---|
| 10% | 70.29 | **79.46** | **79.46** | 76.50 |
| 100% | 85.33 | **89.18** | **89.18** | 87.25 |
| 10% + $C$ | 70.81 | 79.71 | **80.14** | 76.80 |
| 100% + $C$ | 85.49 | 89.19 | **89.75** | 87.59 |

Table 9: Results for Text Chunking. F1 scores on test set of CoNLL 2000 dataset. Product t-norms outperform other t-norms for purely supervised setting (top rows). $\mathcal{R}$-Product performs best among all t-norms when constraints are augmented into the neural models (bottom rows). Rows with $+ C$ include constraints.

rows). First, we observe that incorporating simple constraints into neural models, using the framework discussed in this work, can indeed improve their performance in all cases. Particularly, this improvement is more pronounced in low data regimes, since models need to rely more on external background knowledge when training data is small.

We also find that both variants of product t-norms outperform $\mathcal{S}$-Gödel and Łukasiewicz. However, in both data regimes, $\mathcal{R}$-Product outperforms $\mathcal{S}$-Product when incorporating constraints into the neural model. We observe that these results are consistent with the ones obtained in §5.1.

### 5.3 Theory vs. Experiments

From our analysis of t-norms in §4, we found that Łukasiewicz and $\mathcal{R}$-Product are the most preferable. Empirically, $\mathcal{R}$-Product outperforms the other relaxations.

The consistency analysis suggests that since $\mathcal{S}$-Gödel is least consistent, we should also expect $\mathcal{S}$-Gödel to empirically perform the worst. This argument is reinforced by the results across all tasks and settings, if we naively applied the Gödel t-norm to define the loss. The $\mathcal{S}$-Gödel results shown in this work were obtained by warm starting the learning using the $\mathcal{R}$-Product relaxation. Without the warm start, we found that not only is learning unstable, the resulting accuracies were also low.

One other discrepancy bears attention: although Łukasiewicz is most preferable in terms of consistency on tautologies, the experiments suggest otherwise. From an empirical perspective, a valid relaxation of logic should provide a sufficient gradient signal to make learning feasible. We discovered that loss functions defined by the Łukasiewicz t-norm are not amenable to gradient-based learning. Here, we briefly explain this phenomenon.

Consider the Łukasiewicz conjunction over $n$ atoms:

$$\left[ \bigwedge_{i}^{n} \mathtt{A}_i \right] = \max \left( 0, \sum_{i}^{n} a_i - (n-1) \right)$$

For this operation to provide a non-zero output, we need $\sum_{i}^{n} a_i > (n-1)$, or on average, each of the conjuncts $a_i$ should exceed $\frac{n-1}{n}$.

That is, to provide a useful signal, Łukasiewicz t-norm requires the model to assign high probabilities to correct labels. This, of course, is not true for a randomly initialized model when learning starts, contributing to near-zero gradients, and no model updates at all.

To make learning feasible with Łukasiewicz t-norm, we implemented a less strict definition of conjunction inspired by the MAX-SAT relaxation of Bach *et al.* [2017]. This transforms our learning objective to the form $\max_\theta \left( \sum_{i}^{n} a_i \right)$. An important consequence of this approach is that Łukasiewicz t-norm in our experiments is merely an approximation of the original Łukasiewicz t-norm.

## 6 Related Work and Discussion

The use of t-norm relaxations to encode knowledge into learning is increasingly prevalent in recent years [Wang *et al.*, 2020; Minervini *et al.*, 2017; Donadello *et al.*, 2017, for example]. Additionally, Grefenstette [2013], and Nandwani *et al.* [2019] also implicitly use t-norms to similar ends. These works do not frame the learning problem as a declarative statement encoded using a single predefined logical language, as in the case of [Sikka *et al.*, 2020] and [Giannini *et al.*, 2019].

Our framework is closer to Li *et al.* [2019], Asai and Hajishirzi [2020], and Wang *et al.* [2020] where both data and background knowledge are declaratively stated and encoded in one single t-norm relaxation that defines the loss. Our experiments for the extended MNIST digit classification are inspired by Manhaeve *et al.* [2018], who employ constraints similar to our coherence constraints. Of course, the goal of this work is to theoretically, and empirically investigate which relaxation is best suited for use in such problems. To our knowledge, none of the previously mentioned works analyze the choice of the relaxation they use.

Similar studies to this paper are that of Evans and Grefenstette [2018], and van Krieken *et al.* [2020]. The former, includes comparison in performance of Łukasiewicz, Product, and Gödel t-norm operators used to induce differentiable functions from definitive clauses in neural program synthesis. The latter, perhaps closer to our approach, provides a general theoretical and empirical analysis for different t-norm relaxations. However, there are two key differences with our work: (a) they do not treat labeled data as part of the declarative problem specification, and the constraints are added into a standard cross entropy loss, (b) their analysis involves the derivatives of the losses and does not discuss the consistency and self-consistency properties. We also studied the importance of the loss gradient signal in gradient-based learning, and hypothesise that a less consistent t-norm would perform worse at characterizing the truth of a Boolean statement than a more consistent one, and as a result, would correspond to poorer empirical performance.

## 7 Conclusions

In this work, we studied the question of how best to relax declarative knowledge to define loss functions. To this end, we define a set of criteria that characterizes which relaxations of logic would be most amenable for preserving tautologies and offer support for gradient-based learning. We also present empirical studies on two tasks using the paradigm of formulating entire learning problems via logic. All our analyzes concur that the $\mathcal{R}$-Product relaxation is best suited for learning in this paradigm.

## Acknowledgments

## References

[Asai and Hajishirzi, 2020] Akari Asai and Hannaneh Hajishirzi. Logic-Guided Data Augmentation and Regularization for Consistent Question Answering. In *ACL*, 2020.

[Bach *et al.*, 2017] Stephen H Bach, Matthias Broecheler, Bert Huang, and Lise Getoor. Hinge-loss Markov Random Fields and Probabilistic Soft Logic. *The Journal of Machine Learning Research*, 18, 2017.

[Donadello *et al.*, 2017] Ivan Donadello, Luciano Serafini, and Artur d'Avila Garcez. Logic Tensor Networks for Semantic Image Interpretation. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, 2017.

[Esteva *et al.*, 2000] Francesc Esteva, Lluís Godo, Petr Hájek, and Mirko Navara. Residuated Fuzzy Logics with an Involutive Negation. *Archive for Mathematical Logic*, 39, 2000.

[Evans and Grefenstette, 2018] Richard Evans and Edward Grefenstette. Learning Explanatory Rules from Noisy Data. *Journal of Artificial Intelligence Research (JAIR)*, 61, 2018.

[Fischer *et al.*, 2019] Marc Fischer, Mislav Balunovic, Dana Drachsler-Cohen, Timon Gehr, Ce Zhang, and Martin Vechev. DL2: Training and Querying Neural Networks with Logic. In *ICML*, 2019.

[Giannini *et al.*, 2019] Francesco Giannini, Giuseppe Marra, Michelangelo Diligenti, Marco Maggini, and Marco Gori. On the Relation Between Loss Functions and T-Norms. In *ILP*, 2019.

[Grefenstette, 2013] Edward Grefenstette. Towards a Formal Distributional Semantics: Simulating Logical Calculi with Tensors. In *\*SEM@NAACL*, 2013.

[Klement *et al.*, 2013] Erich Peter Klement, Radko Mesiar, and Endre Pap. *Triangular Norms*, volume 8. Springer Science & Business Media, 2013.

[LeCun, 1998] Yann LeCun. The MNIST Database of Handwritten Digits. http://yann.lecun.com/exdb/mnist/, 1998. Accessed: 2019-11-01.

[Li and Srikumar, 2019] Tao Li and Vivek Srikumar. Augmenting Neural Networks with First-order Logic. In *ACL*, 2019.

[Li *et al.*, 2019] Tao Li, Vivek Gupta, Maitrey Mehta, and Vivek Srikumar. A Logic-Driven Framework for Consistency of Neural Models. In *EMNLP*, 2019.

[Li *et al.*, 2020] Tao Li, Parth Anand Jawale, Martha Palmer, and Vivek Srikumar. Structured Tuning for Semantic Role Labeling. In *ACL*, 2020.

[Loshchilov and Hutter, 2017] Ilya Loshchilov and Frank Hutter. SGDR: stochastic gradient descent with restarts. In *ICLR*, 2017.

[Manhaeve *et al.*, 2018] Robin Manhaeve, Sebastijan Dumancic, Angelika Kimmig, Thomas Demeester, and Luc De Raedt. DeepProbLog: Neural Probabilistic Logic Programming. In *NIPS*, 2018.

[Menger, 1942] Karl Menger. Statistical Metrics. *Proceedings of the National Academy of Sciences of the United States of America*, 28, 1942.

[Minervini *et al.*, 2017] Pasquale Minervini, Thomas Demeester, Tim Rocktäschel, and Sebastian Riedel. Adversarial sets for regularising neural link predictors. In *UAI*, 2017.

[Nandwani *et al.*, 2019] Yatin Nandwani, Abhishek Pathak, and Parag Singla. A Primal-Dual Formulation for Deep Learning with Constraints. In *NeurIPS*, 2019.

[Paszke *et al.*, 2019] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *NeurIPS*, 2019.

[Pennington *et al.*, 2014] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *EMNLP*, 2014.

[Rocktäschel *et al.*, 2015] Tim Rocktäschel, Sameer Singh, and Sebastian Riedel. Injecting Logical Background Knowledge into Embeddings for Relation Extraction. In *NAACL*, 2015.

[Russell and Whitehead, 1910] Bertrand Russell and Alfred North Whitehead. *Principia Mathematica Vol. I*. Cambridge University Press, 1910.

[Sang and Buchholz, 2000] Erik Tjong Kim Sang and Sabine Buchholz. Introduction to the CoNLL-2000 Shared Task Chunking. In *CoNLL*, 2000.

[Sikka *et al.*, 2020] Karan Sikka, Andrew Silberfarb, John Byrnes, Indranil Sur, Ed Chow, Ajay Divakaran, and Richard Rohwer. Deep Adaptive Semantic Logic (DASL): Compiling Declarative Knowledge into Deep Neural Networks. *arXiv preprint arXiv:2003.07344*, 2020.

[van Krieken *et al.*, 2020] Emile van Krieken, Erman Acar, and Frank van Harmelen. Analyzing Differentiable Fuzzy Implications. In *KR*, 2020.

[Wang *et al.*, 2020] Haoyu Wang, Muhao Chen, Hongming Zhang, and Dan Roth. Joint Constrained Learning for Event-Event Relation Extraction. In *EMNLP*, 2020.

[Xu *et al.*, 2018] Jingyi Xu, Zilu Zhang, Tal Friedman, Yitao Liang, and Guy Van den Broeck. A Semantic Loss Function for Deep Learning with Symbolic Knowledge. In *ICML*, 2018.

## A  Relevant Proofs

**Proof for Proposition 1.**

*Proof.* The t-norms we consider (Gödel, Łukasiewicz and Product) are left-continuous binary operators. For a left-continuous t-norm $T$, we can uniquely define the residuum (denoted by $\rightarrow$) as the binary operator that satisfies

$$\forall x, y, z \in [0,1], \quad T(z,x) \le y \text{ iff. } z \le (x \rightarrow y). \quad (8)$$

The residuum generalizes the logical implication. From the definition, we see that $x \rightarrow y$ is greater than or equal to every $z$ such that $T(z,x) \le y$. This allows us to write the residuum in the following equivalent way:

$$(x \rightarrow y) = \sup\{z \mid T(z,x) \le y\}. \quad (9)$$

Next, we use the monotonicity property of t-norms. Every t-norm is a non-decreasing function of both its arguments. In particular, we have:

$$T(z,x) \le T(1,x) = x \quad (10)$$

Consequently, for any $x, y \in [0,1]$ such that $x \le y$, we have $T(z,x) \le y$. For such $x, y$ pairs, the supremum in equation 9 is identical to 1. Indeed, we can show that whenever the residuum $x \rightarrow y$ evaluates to 1, we have $x \le y$. In short,

$$(x \rightarrow y) = 1 \text{ iff. } x \le y. \quad (11)$$

Now, consider the self-consistency of a well-formed Boolean formula P and an $\mathcal{R}$-logic relaxation L, denoted by $\kappa_S^L(\mathbf{P})$. From the definition of self-consistency, we have:

$$\kappa_S^L(\mathbf{P}) = \kappa^L(\mathbf{P} \leftrightarrow \mathbf{P})$$
$$= \kappa^L((\mathbf{P} \leftarrow \mathbf{P}) \wedge (\mathbf{P} \rightarrow \mathbf{P}))$$

When we relax the conjunction in the final expression, we will end up with a t-norm, both of whose arguments are one from equation 11. Consequently, the t-norm itself evaluates to one. Applying the definition of consistency (the integral of the constant 1 over the domain) gives us the required result. $\square$

**Proof Sketch for Proposition 2.**

*Proof.* Let $F = A \leftrightarrow A = \bigwedge_{i=1}^{n} \mathbf{A}_i \leftrightarrow \bigwedge_{i=1}^{n} \mathbf{A}_i$. By the definition of the $\mathcal{S}$-Product logic connectives we have,

$$[F] = \left(1 + \left(\prod_{i=1}^{n} a_i\right)^2 - \left(\prod_{i=1}^{n} a_i\right)\right)^2$$

and by the definition of self-consistency, we have

$$\kappa_S^{\mathcal{S}\text{-Product}}(A) = \int_0^1 [F] \, \mathrm{d}\mathcal{F}$$

Let us simplify the large product corresponding to the n-conjunction by introducing an auxiliary variable that includes all but the final element in the conjunction $A$. We will call this $n-1$-conjunction $b_{n-1}$. That is,

$$b_{n-1} = \prod_{i=1}^{n-1} a_i$$

This allows us to write $[F]$ as

$$[F] = (1 + a_n^2 b_{n-1}^2 - a_n b_{n-1})^2.$$

Integrating out the $n^{th}$ variable $a_n$ by expanding out the polynomial, we obtain

$$\kappa_n = \int_0^1 \left(1 + a_n^2 b_{n-1}^2 - a_n b_{n-1}\right)^2 \mathrm{d}a_n$$
$$= 1 - b_{n-1} + b_{n-1}^2 - \frac{1}{2}b_{n-1}^3 + \frac{1}{5}b_{n-1}^4.$$

Applying the same strategy as before, we can define

$$b_{n-2} = \prod_{i=1}^{n-2} a_i$$

then,

$$b_{n-1} = a_{n-1}b_{n-2}.$$

Substituting $b_{n-1}$ in $\kappa_n$ and taking the integral over the $(n-1)^{th}$ variable we obtain,

$$\kappa_{n-1} = \int_0^1 \kappa_n \, \mathrm{d}a_{n-1}$$
$$= \int_0^1 1 - a_{n-1}b_{n-2} + a_{n-1}^2 b_{n-2}^2$$
$$\quad - \frac{1}{2}a_{n-1}^3 b_{n-2}^3 + \frac{1}{5}a_{n-1}^4 b_{n-2}^4 \, \mathrm{d}a_{n-1}$$
$$= 1 - \frac{2}{2^2}b_{n-2} + \frac{3}{3^2}b_{n-2}^2 - \frac{2}{4^2}b_{n-2}^3 + \frac{1}{5^2}b_{n-2}^4.$$

Repeating this procedure, every subsequent integral will add one to the degree of the denominators in each term. This allows us to generalize the integral as

$$\kappa_{n-k} = 1 - \frac{2}{2^k}b_{n-k} + \frac{3}{3^k}b_{n-k}^k - \frac{2}{4^k}b_{n-k}^3 + \frac{1}{5^k}b_{n-k}^4$$

where,

$$b_{n-k} = \prod_{i=1}^{n-k} a_i.$$

In particular, when $k = n$, the empty product becomes 1. Therefore, substituting $b_{n-k}$ in $\kappa_{n-k}$ we obtain,

$$\kappa_S^{\mathcal{S}\text{-Product}}(A) = \kappa_0 = 1 - \frac{2}{2^n} + \frac{3}{3^n} - \frac{2}{4^n} + \frac{1}{5^n}$$

$\square$

## B  $\mathcal{R}$-logics and the Connective for Negation

In table 1, we abuse the nomenclature for $\mathcal{R}$-Product and $\mathcal{R}$-Gödel logics. In these relaxations, the negation and conjunction connectives are respectively defined as

$$n_\neg(x) = \begin{cases} 1 & \text{if x=0} \\ 0 & \text{otherwise} \end{cases}$$

and,

$$x \wedge y = \begin{cases} 0 & \text{if x=y=0} \\ 1 & \text{otherwise} \end{cases}$$

Unfortunately, these functions are not sub-differentiable. Instead, we introduce the corresponding $SBL_\sim$ extensions for residuated $\mathcal{R}$ t-norm logics [Esteva *et al.*, 2000] which are constructed using the involutive negation $n_\neg(x) = 1 - x$ and sub-differentiable in all of their connectives. We keep the same names $\mathcal{R}$-Product and $\mathcal{R}$-Gödel for simplicity.

## C  Reproducibility

### C.1  Digit Arithmetic experiments

**Model Architecture.**  For `Digit`, `Sum` and `Product` models we use 2-layer Convolutional Neural Networks with 3-by-3 padded filters, with ReLU activation and 2-by-2 Maxpool layers with stride 2 in between, with 2 fully connected layers with ReLU activation and dropout ($p = 0.5$) in between.

**Seeds and multiple runs.**  We use a random seed 20 to create all the dataset used for training, evaluation, and testing. The same seed (20) is used for hyper-parameter tuning. We perform three runs of every experiment using seeds 0, 20, and 50. The standard deviation from the three runs with different seeds was high (up to 5%) for some of the relaxations over the scarce data settings (DIGIT size 1000 and/or PAIR size 1000). In these cases, we perform extra runs with seeds 1, 10, and 60, reporting the average of the best three.

**Hyper-parameter tuning Details**  Validation sets are used for tuning hyper-parameters: learning rate, batch size, optimizer, lambda coefficient for the coherence constraints in the loss. Specifically we perform grid search with:

- learning rates: $\{10^{-1}, 5 \times 10^{-2}, 10^{-2}, 5 \times 10^{-3}, 10^{-3}, 5 \times 10^{-4}, 10^{-4}, 5 \times 10^{-5}, 10^{-5}\}$
- $\lambda$'s: $\{0.05, 0.1, 0.5, 1, 1.5, 2\}$
- batch sizes: $\{8, 16, 32, 64\}$.

We also treat the optimization method as a hyper-parameter, using the one resulting in best performance between standard Stochastic Gradient Descent and Adam optimizers.

**Number of epochs for training.**  The number of epochs for convergence varies between 100 to 1000 epochs across t-norm relaxations and data settings. All the experiments were run for a maximum of 1600 epochs.

**Significance.**  We performed t-tests and observed a statistical significance of at least $p < 0.05$ in all the model comparisons we discuss in the main paper.

**Warm-restarts.**  Using the Stochastic Gradient Descent with Restarts techniques [Loshchilov and Hutter, 2017], brings, considerably, more stability to the experiments. We use the best number iterations for the first restart ($T_0$) among 50, 100 and 200, with a factor of increase of 1 ($T_{mult}$)

### C.2  Experimental Details for Text Chunking

**Model Architecture.**  For all our models, we use bidirectional RNN based models, experimenting with LSTMs and GRUs. We use GloVe word embeddings of 300 dimension at input along with 100 dimensional character embeddings. We fix the dropout rate to 0.5 for all our experiments.

**Hyper-parameter Tuning**  Since the original CoNLL 2000 dataset does not accompany a validation set, for each experiment, we select 10% of the total training data for validation. We perform hyper-parameter tuning using grid search on the following parameters: hidden unit size of RNN in the set $\{200, 256, 300, 400, 500, 512\}$, learning rate in the set $\{0.001, 0.005, 0.0001\}$.

### C.3  Weighted Constrained Loss

We observe that in order to stabilize training, we have to use a $\lambda$ (hyper-parameter) to relatively weigh the loss from both `Sum` and `Product` coherence constraints. We find that product based t-norms are less sensitive to different values of $\lambda$.

For chunking, for example, for Łukasiewicz t-norm, values around 0.0001 works best, while for $\mathcal{R}$-Product, values greater than 1 provides best results on validation set.

### C.4  Code and Computing Infrastructure

We implement all our experiments in the deep learning library PyTorch [Paszke *et al.*, 2019] on a server with the following configuration:

- CPU: AMD EPYC 7601, 32 cores, 2.2 GHz ,
- Memory: 512 gb,
- GPU: Titan RTX,
- Disk Space: 8 TB.

## D  Learning behaviour of neural models

### D.1  Product relaxations

$\mathcal{S}$-Product and $\mathcal{R}$-Product logic are the most stable and less brittle t-norm logics across different data settings.

### D.2  Discussion about Gradients in Łukasiewicz relaxation

From an empirical perspective, a valid relaxation of logic should provide enough gradient signal to make the gradient-based learning process possible. One relevant case of a logic relaxation that does not satisfy this property is the Łukasiewicz t-norm. By inspecting the Łukasiewicz conjunction connective definition:

$$\left[ \bigwedge_i^n \mathtt{A}_i \right] = \max\left( 0, \sum_i^n a_i - (n-1) \right) \quad (12)$$

we can see that it requires almost absolute certainty for each of its conjuncts to result in a non-zero truth value. This means that the gradient signal provided from an objective that is the conjunction of the dataset term (equation 1) and the constraint term (equation 3) is almost always zero.

To be able to study the Łukasiewicz logic, for our experiments we implemented a less strict definition of the conjunction inspired by the MAX-SAT relaxation of Bach *et al.* [2017] assuming the best case scenario $\sum_i^n a_i \geq n - 1$. Therefore, the learning problem using Łukasiewicz logic, defined in §2.1 as

$$\max_\theta \left( \max\left( 0, \sum_i^n a_i - (n-1) \right) \right)$$

is restated for our experiments as

$$\max_\theta \left( \sum_i^n a_i \right) \qquad (13)$$

### D.3 Practical considerations for $\mathcal{S}$-Gödel logic

In practice, we find that optimizing the loss function derived from the $\mathcal{S}$-Gödel relaxation was difficult due to learning instability. To make learning possible, we apply the following techniques:

- By definition of $\mathcal{S}$-Gödel conjunction, for every epoch, we make the update with respect to the example with the minimum output value in the data. Instead, we apply mini-batch descent over the minimum of each batch, and proceed with the standard SGD process. This means that we use the $\mathcal{S}$-Gödel logic over the batches but the $\mathcal{S}$-Product t-norm (which agrees with the cross-entropy loss) over the epochs. While this is not strictly a Gödel conjunction, we found that this was necessary for the optimizer to work.

- Given that at the early stages of learning the system is uncertain about every example, and $\mathcal{S}$-Gödel conjunction operates over the "worst" example, it is hard to get any meaningful signal from the updates. To break ties, we "warm-up" the system by running 2 to 10 epochs of learning using the $\mathcal{R}$-Product logic. Interestingly, even though this strategy is needed to learn with $\mathcal{S}$-Gödel on both joint and pipeline learning frameworks, the joint system needs only one or two epochs to warm-up, while the pipelined approach needs between 6 to 10 epochs to start learning on its own. These warm-up processes represent on average 40% of the final accuracies obtained across our experiments.

- In the joint model, we noticed that the accuracies for the `Digit` classifier decrease after a few epochs hindering training. This is caused because the model is trying to reach the trivial solution to satisfy the coherence constraints (the `Digit` classifier predicting low scores for all digits to make the value of the conjunction in the antecedent of the constraints as low as possible). In order to avoid this, we froze the parameters corresponding to the `Digit` model right after the warm-up stage; in other words, we train `Digit` classifier using $\mathcal{R}$-Product relaxation for a few epochs; enough to successfully continue the training.

## E  Examples of Consistency Calculation

Here we calculate consistency under different relaxations for the axiom schema tautology $P \to (Q \to P)$ which we will denote by $A1$.

- $\mathcal{S}$-Gödel t-norm.

$$[A1] = \max\left(1 - [P], [Q \to P]\right)$$
$$= \max\left(1 - p, \max(1 - q, p)\right)$$

The consistency is given by

$$\kappa^{\mathcal{S}\text{-Gödel}}(A1)$$
$$= \int_0^1 \int_0^1 \max\left(1 - p, \max(1 - q, p)\right) \, \mathrm{d}p \, \mathrm{d}q$$
$$= \frac{19}{24} \approx 0.79$$

- $\mathcal{S}$-Product t-norm.

$$[A1] = 1 - [P] + ([P] \cdot [Q \to P])$$
$$= 1 - p + (p \cdot (1 - q + q \cdot p))$$

The consistency is given by

$$\kappa^{\mathcal{S}\text{-Product}}(A1)$$
$$= \int_0^1 \int_0^1 1 - p + (p \cdot (1 - q + q \cdot p)) \, \mathrm{d}p \, \mathrm{d}q$$
$$= \frac{11}{12} \approx 0.92$$

- $\mathcal{R}$-Product t-norm.

$$[A1] = \begin{cases} 1 & \text{if} \quad [P] \le [Q \to P] \\ \frac{[Q \to P]}{[P]} & \text{else} \end{cases}$$

We have two cases:

1. $[P] < [Q]$

   By the definition of $\mathcal{R}$-Product implication,

   $$[Q \to P] = \frac{p}{q}$$

   then,

   $$[A1] = \begin{cases} 1 & \text{if} \quad p \le \frac{p}{q} \\ \frac{1}{q} & \text{else} \end{cases}$$

   but notice that $p \not> \frac{p}{q}$ because $0 < q < 1$ by definition. Therefore, $[A1] = 1$.

2. $[P] \ge [Q]$

   By the definition of $\mathcal{R}$-Product implication,

   $$[Q \to P] = 1$$

   then, by the definition of $[A1]$ under $\mathcal{R}$-Product and the fact that $0 < [P] < 1$ we obtain that $[A1] = 1$.

- Łukasiewicz t-norm.

$$[A1] = \min\left(1, 1 - [P] + \min\left(1, 1 - [Q] + [P]\right)\right)$$
$$= \min\left(1, 1 - p + \min(1 - q + p)\right)$$

The consistency is given by

$$\kappa^{\mathcal{S}\text{-Łukasiewicz}}(A1)$$
$$= \int_0^1 \int_0^1 \min\left(1, 1 - p + \min(1 - q + p)\right) \, dp \, dq$$
$$= 1$$

# F  Additional Experiments and Results

## F.1  Recognizing Digits and Arithmetic Operations

For the same set of experiments described in section 5.1, tables 10, 11, 12, 13, 14, 15 report the results from instantiating the sum and product coherence constraints over 1k and 25k PAIR examples, respectively, to obtain the training signal for the `Sum` and `Prod` classifiers.

|  | 1000 | 5000 | 25000 |
| --- | --- | --- | --- |
| $\mathcal{S}$-Gödel | 94.5 (0.3) | 97.0 (0.2) | 96.7 (0.1) |
| $\mathcal{S}$-Product | 92.5 (0.1) | 97.8 (0.0) | **99.0** (0.1) |
| $\mathcal{R}$-Product | **95.5** (0.1) | **97.9** (0.0) | 96.0 (0.0) |
| Łukasiewicz | 95.1 (0.4) | 96.2 (0.2) | 97.9 (0.3) |

Table 10: `Digit` accuracies (and standard deviations) from jointly training `Digit` on DIGIT sizes 1k, 5k, 25k and operators (`Sum`,`Prod`) on PAIR size 1k.

|  | 1000 | 5000 | 25000 |
| --- | --- | --- | --- |
| $\mathcal{S}$-Gödel | **62.7** (0.9) | **62.0** (0.3) | 64.1 (0.2) |
| $\mathcal{S}$-Product | 49.6 (1.1) | 59.1 (1.0) | 60.6 (0.3) |
| $\mathcal{R}$-Product | 53.3 (0.7) | 59.3 (0.5) | **66.5** (0.2) |
| Łukasiewicz | 44.5 (4.2) | 45.3 (3.1) | 50.8 (0.7) |

Table 11: Average of `Sum` and `Product` accuracies (and standard deviations) from jointly training `Digit` on DIGIT sizes 1k, 5k, 25k and operators on PAIR size 1k.

|  | 1000 | 5000 | 25000 |
| --- | --- | --- | --- |
| $\mathcal{S}$-Gödel | **64.1** (0.9) | **62.1** (0.2) | 63.9 (0.2) |
| $\mathcal{S}$-Product | 52.1 (1.3) | 59.4 (1.2) | 60.7 (0.5) |
| $\mathcal{R}$-Product | 54.4 (0.5) | 59.1 (0.2) | **67.0** (0.0) |
| Łukasiewicz | 45.4 (4.5) | 45.8 (2.9) | 50.9 (1.0) |

Table 12: Average of `Sum` and `Prod` Coherence accuracies (and standard deviations) from jointly training `Digit` on DIGIT sizes 1k, 5k, 25k and operators on PAIR size 1k.

Similarly, tables 16, 17, 18, 19 report the results from assigning noisy labels to 1k and 25k unlabeled PAIR data examples to train the `Sum` and `Prod` classifiers independently following the pipeline strategy we use in section 5.1.

For all these different data configurations, we observed a similar trend in the results to the one observed in section 5.1 with $\mathcal{R}$-Product achieving the highest performance among the different relaxations.

|  | 1000 | 5000 | 25000 |
| --- | --- | --- | --- |
| $\mathcal{S}$-Gödel | 93.5 (0.2) | 97.7 (0.1) | 98.0 (0.0) |
| $\mathcal{S}$-Product | 95.2 (0.0) | 98.2 (0.0) | 99.0 (0.0) |
| $\mathcal{R}$-Product | **96.2** (0.0) | **98.4** (0.0) | 99.0 (0.0) |
| Łukasiewicz | 94.8 (0.1) | 98.2 (0.1) | **99.0** (0.0) |

Table 13: `Digit` accuracies (and standard deviations) from jointly training `Digit` on DIGIT sizes 1k, 5k, 25k and operators (`Sum`,`Prod`) on PAIR size 25k.

|  | 1000 | 5000 | 25000 |
| --- | --- | --- | --- |
| $\mathcal{S}$-Gödel | 88.3 (0.5) | 94.9 (0.2) | 95.9 (0.0) |
| $\mathcal{S}$-Product | 83.6 (0.5) | 90.9 (0.5) | 95.0 (0.0) |
| $\mathcal{R}$-Product | **89.6** (0.2) | **95.7** (0.1) | **96.1** (0.1) |
| Łukasiewicz | 70.8 (3.0) | 89.4 (1.4) | 90.4 (0.4) |

Table 14: Average of `Sum` and `Product` accuracies (and standard deviations) from jointly training `Digit` on DIGIT sizes 1k, 5k, 25k and operators on PAIR size 25k.

|  | 1000 | 5000 | 25000 |
| --- | --- | --- | --- |
| $\mathcal{S}$-Gödel | 86.0 (0.4) | 94.6 (0.2) | 95.4 (0.1) |
| $\mathcal{S}$-Product | 85.8 (0.5) | 91.7 (0.4) | 95.6 (0.1) |
| $\mathcal{R}$-Product | **90.8** (0.1) | **95.9** (0.0) | **96.3** (0.0) |
| Łukasiewicz | 72.9 (3.1) | 90.2 (1.2) | 90.7 (0.4) |

Table 15: Average of `Sum` and `Prod` Coherence accuracies (and standard deviations) from jointly training `Digit` on DIGIT sizes 1k, 5k, 25k and operators on PAIR size 25k.

|  | 1000 | 5000 | 25000 |
| --- | --- | --- | --- |
| $\mathcal{S}$-Gödel | 60.7 (0.6) | **63.0** (0.4) | 62.9 (0.2) |
| $\mathcal{R}/\mathcal{S}$-Product | **61.5** (0.2) | 60.9 (0.3) | **62.5** (0.1) |
| Łukasiewicz | 54.9 (3.1) | 55.5 (6.7) | 55.1 (0.3) |

Table 16: Average of Pipelined `Sum` and `Product` accuracies (and standard deviations) trained on PAIR size 1k (noisy) labeled with `Digit` model trained on DIGIT sizes 1k, 5k, 25k.

|  | 1000 | 5000 | 25000 |
| --- | --- | --- | --- |
| $\mathcal{S}$-Gödel | 60.7 (0.4) | **63.5** (0.3) | 62.3 (0.4) |
| $\mathcal{R}/\mathcal{S}$-Product | **62.1** (0.4) | 61.2 (0.3) | **62.6** (0.3) |
| Łukasiewicz | 55.6 (2.4) | 55.7 (3.4) | 55.6 (1.0) |

Table 17: Average of Pipelined `Sum` and `Prod` Coherence accuracies (and standard deviations) trained on PAIR size 1k (noisy) labeled with `Digit` model trained on DIGIT sizes 1k, 5k, 25k.

**Arithmetic Properties**

We also evaluate the resulting models from both of the training regimes (joint and pipelines) on arithmetic properties they are not being trained for. Specifically, we create test sets to evaluate if the commutativity, associativity, and distributivity properties are satisfied by the operator classifiers.

|  | 1000 | 5000 | 25000 |
|---|---|---|---|
| $\mathcal{S}$-Gödel | 90.5 (0.1) | 94.5 (0.0) | 95.4 (0.0) |
| $\mathcal{R}/\mathcal{S}$-Product | **91.1** (0.2) | **95.3** (0.1) | **96.9** (0.0) |
| Łukasiewicz | 85.9 (1.4) | 93.3 (1.9) | 95.4 (0.3) |

Table 18: Average of Pipelined `Sum` and `Product` accuracies (and standard deviations) trained on PAIR size 25k (noisy) labeled with `Digit` model trained on DIGIT sizes 1k, 5k, 25k.

.

|  | 1000 | 5000 | 25000 |
|---|---|---|---|
| $\mathcal{S}$-Gödel | 88.5 (0.2) | 93.8 (0.2) | 94.8 (0.0) |
| $\mathcal{R}/\mathcal{S}$-Product | **91.0** (0.1) | **95.3** (0.1) | **96.9** (0.0) |
| Łukasiewicz | 85.8 (1.5) | 93.3 (1.5) | 95.4 (0.3) |

Table 19: Average of Pipelined `Sum` and `Prod` Coherence accuracies (and standard deviations) trained on PAIR size 25k (noisy) labeled with `Digit` model trained on DIGIT sizes 1k, 5k, 25k.

To evaluate commutativity, we check how many times an operator classifier predicts the same digit for pairs of TEST images and their reverse order.

To evaluate the associativity property

$$\forall x_1, x_2, x_3 \quad (x_1 + x_2) + x_3 = x_1 + (x_2 + x_3)$$

in a operator classifier, say the `Sum` (mod 10) operator, we use triples of different images from TEST. For each triple $x_1$, $x_2$, $x_3$, we consider the predicted digits $y_{1,2}$ and $y_{2,3}$ corresponding to the `Sum` of $x_1$ and $x_2$ ($x_1 + x_2$), and the `Sum` of $x_2$ and $x_3$ ($x_2 + x_3$) respectively. Then, we randomly sample an image $x_{1,2}$ of a $y_{1,2}$, and an image $x_{2,3}$ of a $y_{2,3}$ from TEST, and check whether the respective predicted `Sum` of $x_{1,2}$ and $x_3$ ($x_{1,2} + x_3$) and `Sum` of $x_1$ and $x_{2,3}$ ($x_1 + x_{2,3}$) agree. To compensate to the random choice of the image in TEST representing the predicted associations in the property, we repeat the process for 6 iterations and consider the most frequent resulting digit on each side of the equation respectively to make the final comparison.

We follow an analogous approach to evaluate the distributivity property which involves predictions from both `Sum` and `Prod` classifiers.

In tables 20, 21, 22, 23, 24, 25, the first two columns respectively show the average for `Sum` and `Prod` classifiers of the fraction of test instances where the commutativity and associativity properties are satisfied. Similarly, the last column shows the average of the fraction of examples satisfying the left and right distributivity properties. In these experiments, we observe that Gödel and $\mathcal{R}$-Product are the best relaxations. On the other hand, Łukasiewicz t-norm performance was considerably poor with high standard deviations.

## G   Full Tautologies Table and Consistencies

Table 26 shows the entire set of tautologies we considered for the evaluation of consistencies across different relaxations. We observe that the top three relaxations preserving truth (entries equal to 1) are $\mathcal{R}$-Gödel, Łukasiewicz and $\mathcal{R}$-Product, however, $\mathcal{R}$-Gödel does not satisfy P1 as its definition of implication is not a sub-differentiable function (see §4.1).

|  |  | Commut. | Assoc. | Dist. |
|---|---|---|---|---|
| 1000 | $\mathcal{S}$-Gödel | **59.1** (3.4) | **47.6** (2.4) | **43.3** (4.1) |
|  | $\mathcal{S}$-Product | 55.8 (1.4) | 45.6 (1.9) | 39.7 (2.5) |
|  | $\mathcal{R}$-Product | 49.4 (1.8) | 41.4 (1.2) | 32.1 (1.9) |
|  | Łukasiewicz | 43.0 (5.4) | 37.6 (5.2) | 27.2 (4.8) |
| 5000 | $\mathcal{S}$-Gödel | **60.0** (2.6) | 48.1 (2.0) | **41.8** (2.9) |
|  | $\mathcal{S}$-Product | 54.2 (1.2) | 43.5 (1.2) | 37.2 (1.9) |
|  | $\mathcal{R}$-Product | 54.0 (1.1) | **48.9** (1.0) | 41.4 (1.2) |
|  | Łukasiewicz | 43.8 (6.0) | 39.9 (5.7) | 31.1 (3.7) |
| 25000 | $\mathcal{S}$-Gödel | 62.7 (0.5) | 50.0 (1.0) | 42.7 (0.9) |
|  | $\mathcal{S}$-Product | 52.8 (0.4) | 44.7 (0.3) | 39.0 (0.4) |
|  | $\mathcal{R}$-Product | **63.0** (0.2) | **51.8** (0.3) | **44.7** (0.3) |
|  | Łukasiewicz | 52.2 (4.2) | 44.9 (4.6) | 36.8 (3.6) |

Table 20: Arithmetic properties accuracies (and standard deviations) from jointly training `Digit` on DIGIT sizes 1k, 5k, 25k and operators (Sum,Prod) on PAIR size 1k.

|  |  | Commut. | Assoc. | Dist. |
|---|---|---|---|---|
| 1000 | $\mathcal{S}$-Gödel | 90.0 (0.3) | 85.8 (0.3) | **82.6** (0.3) |
|  | $\mathcal{S}$-Product | 79.5 (1.0) | 69.7 (0.7) | 66.5 (1.2) |
|  | $\mathcal{R}$-Product | **91.0** (0.8) | **87.0** (0.9) | 82.0 (0.5) |
|  | Łukasiewicz | 76.5 (2.6) | 65.7 (4.1) | 64.3 (5.7) |
| 5000 | $\mathcal{S}$-Gödel | **91.8** (0.0) | **89.2** (0.2) | **86.8** (0.1) |
|  | $\mathcal{S}$-Product | 89.4 (0.5) | 85.6 (0.4) | 83.0 (0.5) |
|  | $\mathcal{R}$-Product | 90.9 (0.5) | 88.4 (0.6) | 83.9 (0.1) |
|  | Łukasiewicz | 85.3 (1.8) | 80.2 (2.9) | 74.7 (4.5) |
| 25000 | $\mathcal{S}$-Gödel | **91.9** (0.0) | **89.4** (0.1) | **87.2** (0.0) |
|  | $\mathcal{S}$-Product | 90.1 (0.1) | 86.8 (0.0) | 84.3 (0.1) |
|  | $\mathcal{R}$-Product | 91.4 (0.1) | 89.3 (0.2) | 84.1 (0.0) |
|  | Łukasiewicz | 86.2 (0.6) | 81.2 (0.7) | 67.7 (4.0) |

Table 21: Arithmetic properties accuracies (and standard deviations) from jointly training `Digit` on DIGIT sizes 1k, 5k, 25k and operators (Sum,Prod) on PAIR size 5k.

|  |  | Commut. | Assoc. | Dist. |
|---|---|---|---|---|
| 1000 | $\mathcal{S}$-Gödel | 93.9 (0.1) | 89.5 (0.1) | 85.7 (0.0) |
|  | $\mathcal{S}$-Product | 89.8 (0.4) | 82.2 (0.2) | 79.0 (0.6) |
|  | $\mathcal{R}$-Product | **95.3** (0.3) | **91.8** (0.4) | **86.3** (0.3) |
|  | Łukasiewicz | 79.9 (2.5) | 63.3 (4.9) | 55.7 (6.9) |
| 5000 | $\mathcal{S}$-Gödel | 96.6 (0.0) | 95.0 (0.0) | 92.2 (0.0) |
|  | $\mathcal{S}$-Product | 92.3 (0.1) | 88.8 (0.2) | 86.5 (0.2) |
|  | $\mathcal{R}$-Product | **97.5** (0.2) | **96.4** (0.0) | **93.8** (0.0) |
|  | Łukasiewicz | 95.5 (0.3) | 93.7 (0.3) | 83.6 (0.4) |
| 25000 | $\mathcal{S}$-Gödel | 96.7 (0.0) | 95.6 (0.0) | **94.5** (0.0) |
|  | $\mathcal{S}$-Product | 96.4 (0.2) | 94.9 (0.1) | 93.6 (0.2) |
|  | $\mathcal{R}$-Product | **97.5** (0.0) | **96.5** (0.2) | 94.1 (0.1) |
|  | Łukasiewicz | 94.9 (0.3) | 91.5 (0.1) | 83.2 (0.3) |

Table 22: Arithmetic properties accuracies (and standard deviations) from jointly training `Digit` on DIGIT sizes 1k, 5k, 25k and operators (Sum,Prod) on PAIR size 25k.

|       |                      | Commut.      | Assoc.       | Dist.        |
| ----- | -------------------- | ------------ | ------------ | ------------ |
| 1000  | $\mathcal{S}$-Gödel  | 56.4 (0.9)   | 44.7 (0.9)   | 33.9 (1.4)   |
|       | $\mathcal{R}/\mathcal{S}$-Product | **56.6** (1.2) | **46.6** (0.8) | **39.1** (0.9) |
|       | Łukasiewicz          | 55.0 (2.8)   | 43.8 (9.5)   | 21.5 (8.5)   |
| 5000  | $\mathcal{S}$-Gödel  | **55.9** (0.5) | **46.1** (0.4) | **39.9** (0.6) |
|       | $\mathcal{R}/\mathcal{S}$-Product | 54.3 (0.3) | 44.6 (0.4)   | 36.8 (0.3)   |
|       | Łukasiewicz          | 48.7 (2.3)   | 74.1 (9.1)   | 27.8 (8.8)   |
| 25000 | $\mathcal{S}$-Gödel  | 54.4 (0.2)   | **47.0** (0.3) | **47.2** (0.3) |
|       | $\mathcal{R}/\mathcal{S}$-Product | **54.9** (0.3) | 45.2 (0.3) | 35.2 (0.2)   |
|       | Łukasiewicz          | 52.9 (1.3)   | 44.5 (7.1)   | 30.3 (7.5)   |

Table 23: Sum and Product arithmetic properties accuracies averages (and standard deviations) from pipeline training on PAIR size 1k (noisy) labeled with Digit model trained on DIGIT sizes 1k, 5k, 25k.

|       |                      | Commut.      | Assoc.       | Dist.        |
| ----- | -------------------- | ------------ | ------------ | ------------ |
| 1000  | $\mathcal{S}$-Gödel  | 90.0 (0.2)   | 86.0 (0.2)   | 82.4 (0.5)   |
|       | $\mathcal{R}/\mathcal{S}$-Product | **92.4** (0.3) | **88.9** (0.3) | **85.1** (0.2) |
|       | Łukasiewicz          | 85.4 (1.9)   | 74.1 (6.9)   | 58.6 (7.5)   |
| 5000  | $\mathcal{S}$-Gödel  | 90.6 (0.0)   | 88.0 (0.1)   | 85.7 (0.3)   |
|       | $\mathcal{R}/\mathcal{S}$-Product | **92.4** (0.0) | **90.1** (0.0) | **87.5** (0.1) |
|       | Łukasiewicz          | 88.8 (1.2)   | 80.1 (5.3)   | 76.5 (6.1)   |
| 25000 | $\mathcal{S}$-Gödel  | 90.5 (0.0)   | 88.1 (0.0)   | 85.6 (0.2)   |
|       | $\mathcal{R}/\mathcal{S}$-Product | **92.4** (0.0) | **90.4** (0.0) | **87.9** (0.0) |
|       | Łukasiewicz          | 87.5 (0.8)   | 80.5 (6.0)   | 78.4 (5.4)   |

Table 24: Sum and Product arithmetic properties accuracies averages (and standard deviations) from pipeline training on PAIR size 5k (noisy) labeled with Digit model trained on DIGIT sizes 1k, 5k, 25k.

|       |                      | Commut.      | Assoc.       | Dist.        |
| ----- | -------------------- | ------------ | ------------ | ------------ |
| 1000  | $\mathcal{S}$-Gödel  | 94.4 (0.0)   | 91.0 (0.0)   | 87.2 (0.0)   |
|       | $\mathcal{R}/\mathcal{S}$-Product | **95.5** (0.0) | **91.9** (0.0) | **87.6** (0.0) |
|       | Łukasiewicz          | 92.9 (0.5)   | 79.9 (6.6)   | 74.7 (5.0)   |
| 5000  | $\mathcal{S}$-Gödel  | 95.9 (0.0)   | 94.2 (0.0)   | 93.0 (0.0)   |
|       | $\mathcal{R}/\mathcal{S}$-Product | **96.8** (0.0) | **95.5** (0.0) | **93.3** (0.0) |
|       | Łukasiewicz          | 94.3 (1.0)   | 90.0 (3.1)   | 89.7 (2.8)   |
| 25000 | $\mathcal{S}$-Gödel  | 96.0 (0.0)   | 94.7 (0.0)   | 93.5 (0.0)   |
|       | $\mathcal{R}/\mathcal{S}$-Product | **97.4** (0.0) | **96.8** (0.0) | **95.3** (0.0) |
|       | Łukasiewicz          | 96.0 (0.2)   | 94.7 (0.9)   | 92.9 (3.1)   |

Table 25: Sum and Product arithmetic properties accuracies averages (and standard deviations) from pipeline training on PAIR size 25k (noisy) labeled with Digit model trained on DIGIT sizes 1k, 5k, 25k.

| Tautologies | $\mathcal{S}$-**Prod.** | $\mathcal{S}$-**Gödel** | **Łuka.** | $\mathcal{R}$-**Pro.** | $\mathcal{R}$-**Gödel** |
|---|---|---|---|---|---|
| **Axiom Schemata** | | | | | |
| $P \rightarrow (Q \rightarrow P)$ | 0.92 | 0.79 | 1 | 1 | 1 |
| $(P \rightarrow (Q \rightarrow R)) \rightarrow ((P \rightarrow Q) \rightarrow (P \rightarrow R))$ | 0.88 | 0.75 | 0.96 | 0.93 | 1 |
| $(\neg P \rightarrow \neg Q) \rightarrow (Q \rightarrow P)$ | 0.86 | 0.75 | 1 | 0.88 | 0.79 |
| **Primitive Propositions** | | | | | |
| $(P \vee P) \rightarrow P$ | 0.75 | 0.75 | 0.75 | 0.69 | 1 |
| $Q \rightarrow (P \vee Q)$ | 0.92 | 0.79 | 1 | 1 | 1 |
| $(P \vee Q) \rightarrow (Q \vee P)$ | 0.86 | 0.75 | 1 | 1 | 1 |
| $(P \vee (Q \vee R)) \rightarrow (Q \vee (P \vee R))$ | 0.91 | 0.78 | 1 | 1 | 1 |
| $(Q \rightarrow R) \rightarrow ((P \vee Q) \rightarrow (P \vee R))$ | 0.90 | 0.76 | 1 | 1 | 1 |
| **Law of excluded middle** | | | | | |
| $P \vee \neg P$ | 0.83 | 0.75 | 1 | 0.83 | 0.75 |
| **Law of contradiction** | | | | | |
| $\neg(P \wedge \neg P)$ | 0.83 | 0.75 | 1 | 0.83 | 0.75 |
| **Law of double negation** | | | | | |
| $P \leftrightarrow \neg(\neg P)$ | 0.70 | 0.75 | 1 | 1 | 1 |
| **Principles of transposition** | | | | | |
| $(P \leftrightarrow Q) \leftrightarrow (\neg P \leftrightarrow \neg Q)$ | 0.61 | 0.67 | 1 | 0.59 | 0.17 |
| $((P \wedge Q) \rightarrow R) \leftrightarrow ((P \wedge \neg R) \rightarrow \neg Q)$ | 0.84 | 0.78 | 1 | 0.86 | 0.65 |
| **Laws of tautology** | | | | | |
| $P \leftrightarrow (P \wedge P)$ | 0.69 | 0.75 | 0.75 | 0.5 | 1 |
| $P \leftrightarrow (P \vee P)$ | 0.69 | 0.75 | 0.75 | 0.69 | 1 |
| **Laws of absorption** | | | | | |
| $(P \rightarrow Q) \leftrightarrow (P \leftrightarrow (P \wedge Q))$ | 0.66 | 0.71 | 0.83 | 0.67 | 1 |
| $Q \rightarrow (P \leftrightarrow (P \wedge Q))$ | 0.82 | 0.75 | 1 | 1 | 1 |
| **Assoc., Comm., Dist. laws** | | | | | |
| $(P \wedge (Q \vee R)) \leftrightarrow ((P \wedge Q) \vee (P \wedge R))$ | 0.69 | 0.72 | 0.90 | 0.89 | 1 |
| $(P \vee (Q \wedge R)) \leftrightarrow ((P \vee Q) \wedge (P \vee R))$ | 0.69 | 0.72 | 0.90 | 0.94 | 1 |
| **De Morgans Laws** | | | | | |
| $(P \wedge Q) \leftrightarrow \neg(\neg P \vee \neg Q)$ | 0.75 | 0.75 | 1 | 1 | 1 |
| $\neg(P \wedge Q) \leftrightarrow \neg(\neg P \vee \neg Q)$ | 0.75 | 0.75 | 1 | 1 | 1 |
| **Material excluded middle** | | | | | |
| $(P \rightarrow Q) \vee (Q \rightarrow P)$ | 0.97 | 0.83 | 1 | 1 | 1 |

Table 26: Degrees of consistency given by the different logic relaxations under consideration over a representative set of tautologies.