

# Structured Tuning for Semantic Role Labeling

**Tao Li**

University of Utah  
tli@cs.utah.edu

**Parth Anand Jawale**

University of Colorado  
Parth.Jawale@colorado.edu

**Martha Palmer**

University of Colorado  
Martha.Palmer@colorado.edu

**Vivek Srikumar**

University of Utah  
svivek@cs.utah.edu

## Abstract

Recent neural network-driven semantic role labeling (SRL) systems have shown impressive improvements in F1 scores. These improvements are due to expressive input representations, which, at least at the surface, are orthogonal to knowledge-rich constrained decoding mechanisms that helped linear SRL models. Introducing the benefits of structure to inform neural models presents a methodological challenge. In this paper, we present a structured tuning framework to improve models using softened constraints only at training time. Our framework leverages the expressiveness of neural networks and provides supervision with structured loss components. We start with a strong baseline (RoBERTa) to validate the impact of our approach, and show that our framework outperforms the baseline by learning to comply with declarative constraints. Additionally, our experiments with smaller training sizes show that we can achieve consistent improvements under low-resource scenarios.

## 1 Introduction

Semantic Role Labeling (SRL, Palmer et al., 2010) is the task of labeling semantic arguments of predicates in sentences to identify who does what to whom. Such representations can come in handy in tasks involving text understanding, such as coreference resolution (Ponzetto and Strube, 2006) and reading comprehension (e.g., Berant et al., 2014; Zhang et al., 2020). This paper focuses on the question of how knowledge can influence modern semantic role labeling models.

Linguistic knowledge can help SRL models in several ways. For example, syntax can drive feature design (e.g., Punyakanok et al., 2005; Toutanova et al., 2005; Kshirsagar et al., 2015; Johansson and Nugues, 2008, and others), and

can also be embedded into neural network architectures (Strubell et al., 2018).

In addition to such influences on input representations, knowledge about the nature of semantic roles can inform structured decoding algorithms used to construct the outputs. The SRL literature is witness to a rich array of techniques for structured inference, including integer linear programs (e.g., Punyakanok et al., 2005, 2008), bespoke inference algorithms (e.g., Täckström et al., 2015), A\* decoding (e.g., He et al., 2017a), greedy heuristics (e.g., Ouchi et al., 2018), or simple Viterbi decoding to ensure that token tags are BIO-consistent.

By virtue of being constrained by the definition of the task, global inference promises semantically meaningful outputs, and could provide valuable signal when models are being trained. However, beyond Viterbi decoding, it may impose prohibitive computational costs, thus ruling out using inference during training. Indeed, optimal inference may be intractable, and inference-driven training may require ignoring certain constraints that render inference difficult.

While global inference was a mainstay of SRL models until recently, today’s end-to-end trained neural architectures have shown remarkable successes without needing decoding. These successes can be attributed to the expressive input and internal representations learned by neural networks. The only structured component used with such models, if at all, involves sequential dependencies between labels that admit efficient decoding.

In this paper, we ask: *Can we train neural network models for semantic roles in the presence of general output constraints, without paying the high computational cost of inference?* We propose a structured tuning approach that exposes a neural SRL model to differentiable constraints during the finetuning step. To do so, we first write the out-

put space constraints as logic rules. Next, we relax such statements into differentiable forms that serve as regularizers to inform the model at training time. Finally, during inference, our structured-tuned models are free to make their own judgments about labels without any inference algorithms beyond a simple linear sequence decoder.

We evaluate our structured tuning on the CoNLL-05 (Carreras and Màrquez, 2005) and CoNLL-12 English SRL (Pradhan et al., 2013) shared task datasets, and show that by learning to comply with declarative constraints, trained models can make more consistent and more accurate predictions. We instantiate our framework on top of a strong baseline system based on the RoBERTa (Liu et al., 2019) encoder, which by itself performs on par with previous best SRL models that are not ensembled. We evaluate the impact of three different types of constraints. Our experiments on the CoNLL-05 data show that our constrained models outperform the baseline system by 0.2 F1 on the WSJ section and 1.2 F1 on the Brown test set. Even with the larger and cleaner CoNLL-12 data, our constrained models show improvements without introducing any additional trainable parameters. Finally, we also evaluate the effectiveness of our approach on low training data scenarios, and show that constraints can be more impactful when we do not have large training sets.

In summary, our contributions are:

1. We present a structured tuning framework for SRL which uses soft constraints to improve models without introducing additional trainable parameters.<sup>1</sup>
2. Our framework outperforms strong baseline systems, and shows especially large improvements in low data regimes.

## 2 Model & Constraints

In this section, we will introduce our structured tuning framework for semantic role labeling. In §2.1, we will briefly cover the baseline system. To that, we will add three constraints, all treated as combinatorial constraints requiring inference algorithms in past work: **Unique Core Roles** in §2.3, **Exclusively Overlapping Roles** in §2.4, and **Frame Core Roles** in §2.5. For each constraint, we will discuss how to use its softened version dur-

<sup>1</sup>Our code to replay our experiments is archived at [https://github.com/utahnlp/structured\\_tuning\\_srl](https://github.com/utahnlp/structured_tuning_srl).

ing training.

We should point out that the specific constraints chosen serve as a proof-of-concept for the general methodology of tuning with declarative knowledge. For simplicity, for all our experiments, we use the ground truth predicates and their senses.

### 2.1 Baseline

We use RoBERTa (Liu et al., 2019) base version to develop our baseline SRL system. The large number of parameters not only allows it to make fast and accurate predictions, but also offers the capacity to learn from the rich output structure, including the constraints from the subsequent sections.

Our base system is a standard BIO tagger, briefly outlined below. Given a sentence  $s$ , the goal is to assign a label of the form B-X, I-X or O for each word  $i$  being an argument with label X for a predicate at word  $u$ . These unary decisions are scored as follows:

$$e = \text{map}(\text{RoBERTa}(s)) \quad (1)$$

$$v_u, a_i = f_v(e_u), f_a(e_i) \quad (2)$$

$$\phi_{u,i} = f_{va}([v_u, a_i]) \quad (3)$$

$$y_{u,i} = g(\phi_{u,i}) \quad (4)$$

Here,  $\text{map}$  converts the wordpiece embeddings  $e$  to whole word embeddings by summation,  $f_v$  and  $f_a$  are linear transformations of the predicate and argument embeddings respectively,  $f_{va}$  is a two-layer ReLU with concatenated inputs, and finally  $g$  is a linear layer followed by softmax activation that predicts a probability distribution over labels for each word  $i$  when  $u$  is a predicate. In addition, we also have a standard first-order sequence model over label sequences for each predicate in the form of a CRF layer that is Viterbi decoded. We use the standard cross-entropy loss to train the model.

### 2.2 Designing Constraints

Before looking at the specifics of individual constraints, let us first look at a broad overview of our methodology. We will see concrete examples in the subsequent sections.

Output space constraints serve as prior domain knowledge for the SRL task. We will design our constraints as invariants at the training stage. To do so, we will first define constraints as statements in logic. Then we will systematically relax these Boolean statements into differentiable forms using concepts borrowed from the study of triangular norms (t-norms, Klement et al., 2013). Finally,

we will treat these relaxations as regularizers in addition to the standard cross-entropy loss.

All the constraints we consider are conditional statements of the form:

$$\forall x, L(x) \rightarrow R(x) \quad (5)$$

where the left- and the right-hand sides— $L(x), R(x)$  respectively—can be either disjunctive or conjunctive expressions. The literals that constitute these expressions are associated with classification neurons, *i.e.*, the predicted output probabilities are soft versions of these literals.

What we want is that model predictions satisfy our constraints. To teach a model to do so, we transform conditional statements into regularizers, such that during training, the model receives a penalty if the rule is not satisfied for an example.<sup>2</sup>

To soften logic, we use the conversions shown in Table 1 that combine the product and Gödel t-norms. We use this combination because it offers cleaner derivatives make learning easier. A similar combination of t-norms was also used in prior work (Minervini and Riedel, 2018). Finally, we will transform the derived losses into log space to be consistent with cross-entropy loss. Li et al. (2019) outlines this relationship between the cross-entropy loss and constraint-derived regularizers in more detail.

Logic	$\bigwedge_i a_i$	$\bigvee_i a_i$	$\neg a$	$a \rightarrow b$
Gödel	$\min(a_i)$	$\max(a_i)$	$1 - a$	$-$
Product	$\prod a_i$	$-$	$1 - a$	$\min(1, \frac{b}{a})$

Table 1: Converting logical operations to differentiable forms. For literals inside of  $L(s)$  and  $R(s)$ , we use the Gödel t-norm. For the top-level conditional statement, we use the product t-norm. Operations not used this paper are marked as ‘-’.

### 2.3 Unique Core Roles ( $U$ )

Our first constraint captures the idea that, in a frame, there can be at most one core participant of a given type. Operationally, this means that for every predicate in an input sentence  $s$ , there can be no more than one occurrence of each core argument (*i.e.*,  $\mathcal{A}_{core} = \{A0, A1, A2, A3, A4, A5\}$ ). In

<sup>2</sup>Constraint-derived regularizers are dependent on examples, but not necessarily labeled ones. For simplicity, in this paper, we work with sentences from the labeled corpus. However, the methodology described here can be extended to use unlabeled examples as well.

first-order logic, we have:

$$\forall u, i \in s, X \in \mathcal{A}_{core}, \quad B_X(u, i) \rightarrow \bigwedge_{j \in s, j \neq i} \neg B_X(u, j) \quad (6)$$

which says, for a predicate  $u$ , if a model tags the  $i$ -th word as the beginning of the core argument span, then it should not predict that any other token is the beginning of the same label.

In the above rule, the literal  $B_X$  is associated with the predicted probability for the label  $B-X$ <sup>3</sup>. This association is the cornerstone for deriving constraint-driven regularizers. Using the conversion in Table 1 and taking the natural log of the resulting expression, we can convert the implication in (6) as  $l(u, i, X)$ :

$$\max \left( \log B_X(u, i) - \min_{j \in s, j \neq i} \log(1 - B_X(u, j)) \right).$$

Adding up the terms for all tokens and labels, we get the final regularizer  $L_U(s)$ :

$$L_U(s) = \sum_{(u, i) \in s, X \in \mathcal{A}_{core}} l(u, i, X). \quad (7)$$

Our constraint is universally applied to all words and predicates (*i.e.*,  $i, u$  respectively) in the given sentence  $s$ . Whenever there is a pair of predicted labels for tokens  $i, j$  that violate the rule (6), our loss will yield a positive penalty.

**Error Measurement**  $\rho_u$  To measure the violation rate of this constraint, we will report the percentages of propositions that have duplicate core arguments. We will refer to this error rate as  $\rho_u$ .

### 2.4 Exclusively Overlapping Roles ( $O$ )

We adopt this constraint from Punyakanok et al. (2008) and related work. In any sentence, an argument for one predicate can either be contained in or entirely outside another argument for any other predicate. We illustrate the intuition of this constraint in Table 2, assuming core argument spans are unique and tags are BIO-consistent.

Based on Table 2, we design a constraint that says: if an argument has boundary  $[i, j]$ , then no other argument span can cross the boundary at  $j$ .

<sup>3</sup> We will use  $B_X(u, i)$  to represent both the literal that the token  $i$  is labeled with  $B-X$  for predicate  $u$  and also the probability for this event. We follow a similar convention for the  $I-X$  labels.

Token index	$i$	$\dots$	$j$	$j + 1$
$[i-j]$ has label $X$	$B_X$	$\dots$	$I_X$	$\neg I_X$
Not allowed	$-$	$-$	$B_Y$	$I_Y$
Not allowed	$\neg B_Y \wedge \neg I_Y$	$-$	$I_Y$	$I_Y$

Table 2: Formalizing the exclusively overlapping role constraint in terms of the  $B$  and  $I$  literals. For every possible span  $[i-j]$  in a sentence, whenever it has a label  $X$  for some predicate (first row), token labels as in the subsequent rows are not allowed for any other predicate for any other argument  $Y$ . Note that this constraint does not affect the cells marked with a  $-$ .

This constraint applies to all argument labels in the task, denoted by the set  $\mathcal{A}$ .

$$\forall u, i, j \in s \text{ such that } j > i, \text{ and } \forall X \in \mathcal{A},$$

$$P(u, i, j, X) \rightarrow \bigwedge_{\substack{v \in s, Y \in \mathcal{A} \\ (u, X) \neq (v, Y)}} Q(v, i, j, Y) \quad (8)$$

where

$$P(u, i, j, X) = B_X(u, i) \wedge I_X(u, j) \wedge \neg I_X(u, j + 1)$$

$$Q(v, i, j, Y) = Q_1(v, i, j, Y) \wedge Q_2(v, i, j, Y)$$

$$Q_1(v, i, j, Y) = \neg B_Y(v, j) \vee \neg I_Y(v, j + 1)$$

$$Q_2(v, i, j, Y) = B_Y(v, i) \vee I_Y(v, i) \vee \neg I_Y(v, j) \vee \neg I_Y(v, j + 1)$$

Here, the term  $P(u, i, j, X)$  denotes the indicator for the argument span  $[i, j]$  having the label  $X$  for a predicate  $u$  and corresponds to the first row of Table 2. The terms  $Q_1(v, i, j, Y)$  and  $Q_2(v, i, j, Y)$  each correspond to prohibitions of the type described in the second and third rows respectively.

As before, the literals  $B_X$ , etc are relaxed as model probabilities to define the loss. By combining the Gödel and product t-norms, we translate Rule (8) into:

$$L_O(s) = \sum_{\substack{(u, i, j) \in s \\ j > i, X \in \mathcal{A}}} l(u, i, j, X). \quad (9)$$

where,

$$l(u, i, j, X) = \max(0, \log P(u, i, j, X) - \min_{\substack{v \in s, Y \in \mathcal{A} \\ (u, X) \neq (v, Y)}} \log Q(v, i, j, Y))$$

$$P(u, i, j, X) = \min(B_X(u, i), I_X(u, j), 1 - I_X(u, j + 1))$$

$$Q(v, i, j, Y) = \min(Q_1(v, i, j, Y), Q_2(v, i, j, Y))$$

$$Q_1(v, i, j, Y) = 1 - \min(B_Y(v, j), I_Y(v, j + 1))$$

$$Q_2(v, i, j, Y) = \max(B_Y(v, i), I_Y(v, i), 1 - I_Y(v, j), 1 - I_Y(v, j + 1))$$

Again, our constraint applies to all predicted probabilities. However, doing so requires scanning over 6 axes defined by  $(u, v, i, j, X, Y)$ , which is computationally expensive. To get around this, we observe that, since we have a conditional statement, the higher the probability of  $P(u, i, j, X)$ , the more likely it yields non-zero penalty. These cases are precisely the ones we hope the constraint helps. Thus, for faster training and ease of implementation, we modify Equation 8 by squeezing the  $(i, j)$  dimensions using top-k to redefine  $L_O$  above as:

$$\mathcal{T}(u, X) = \arg \text{top-k}_{(i, j) \in s} P(u, i, j, X) \quad (10)$$

$$L_O(s) = \sum_{u \in s, X \in \mathcal{A}} \sum_{(i, j) \in \mathcal{T}(u, X)} l(u, i, j, X). \quad (11)$$

where  $\mathcal{T}$  denotes the set of the top-k span boundaries for predicate  $u$  and argument label  $X$ . This change results in a constraint defined by  $u, v, X, Y$  and the  $k$  elements of  $\mathcal{T}$ .

**Error Measurement  $\rho_o$**  We will refer to the error of the overlap constraint as  $\rho_o$ , which describes the total number of non-exclusively overlapped pairs of arguments. In practice, we found that models rarely make such observed mistakes. In §3, we will see that using this constraint during training helps models generalize better with other constraints. In §4, we will analyze the impact of the parameter  $k$  in the optimization described above.

## 2.5 Frame Core Roles ( $F$ )

The task of semantic role labeling is defined using the PropBank frame definitions. That is, for any predicate lemma of a given sense, PropBank defines which core arguments it can take and what they mean. The definitions allow for natural constraints that can teach models to avoid predicting core arguments outside of the predefined set.

$$\forall u \in s, k \in \mathcal{S}(u),$$

$$\text{Sense}(u, k) \rightarrow \bigwedge_{\substack{i \in s \\ X \notin \mathcal{R}(u, k)}} \neg (B_X(u, i) \wedge I_X(u, i))$$

where  $\mathcal{S}(u)$  denotes the set of senses for a predicate  $u$ , and  $\mathcal{R}(u, k)$  denotes the set of acceptable core arguments when the predicate  $u$  has sense  $k$ .

As noted in §2.2, literals in the above statement can to be associated with classification neurons. Thus the  $\text{Sense}(u, k)$  corresponds to either model prediction or ground truth. Since our focus is to

validate the approach of using relaxed constraints for SRL, we will use the latter.

This constraint can be also converted into regularizer following previous examples, giving us a loss term  $L_F(s)$ .

**Error Measurement**  $\rho_f$  We will use  $\rho_f$  to denote the violation rate. It represents the percentage of propositions that have predicted core arguments outside the role sets of PropBank frames.

**Loss** Our final loss is defined as:

$$L_E(s) + \lambda_U L_U(s) + \lambda_O L_O(s) + \lambda_F L_F(s) \quad (12)$$

Here,  $L_E(s)$  is the standard cross entropy loss over the BIO labels, and the  $\lambda$ 's are hyperparameters.

### 3 Experiments & Results

In this section, we study the question: *In what scenarios can we inform an end-to-end trained neural model with declarative knowledge?* To this end, we experiment with the CoNLL-05 and CoNLL-12 datasets, using standard splits and the official evaluation script for measuring performance. To empirically verify our framework in various data regimes, we consider scenarios ranging from where only limited training data is available, to ones where large amounts of clean data are available.

#### 3.1 Experiment Setup

Our baseline (described in §2.1) is based on RoBERTa. We used the pre-trained base version released by Wolf et al. (2019). Before the final linear layer, we added a dropout layer (Srivastava et al., 2014) with probability 0.5. To capture the sequential dependencies between labels, we added a standard CRF layer. At testing time, Viterbi decoding with hard transition constraints was employed across all settings. In all experiments, we used the gold predicate and gold frame senses.

Model training proceeded in two stages:

1. We use the finetuned the pre-trained RoBERTa model on SRL with *only* cross-entropy loss for 30 epochs with learning rate  $3 \times 10^{-5}$ .
2. Then we continued finetuning with the combined loss in Equation 12 for another 5 epochs with a lowered learning rate of  $1 \times 10^{-5}$ .

During both stages, learning rates were warmed up linearly for the first 10% updates.

For fair comparison, we finetuned our baseline twice (as with the constrained models); we found that it consistently outperformed the singly finetuned baseline in terms of both error rates and role F1. We grid-searched the  $\lambda$ 's by incrementally adding regularizers. The combination of  $\lambda$ 's with good balance between F1 and error  $\rho$ 's on the dev set were selected for testing. We refer readers to the appendix for the values of  $\lambda$ 's.

For models trained on the CoNLL-05 data, we report performance on the dev set, and the WSJ and Brown test sets. For CoNLL-12 models, we report performance on the dev and the test splits.

#### 3.2 Scenario 1: Low Training Data

Creating SRL datasets requires expert annotation, which is expensive. While there are some efforts on semi-automatic annotation targeting low-resource languages (e.g., Akbik et al., 2016), achieving high neural network performance with small or unlabeled datasets remains a challenge (e.g., Fürstenu and Lapata, 2009, 2012; Titov and Klementiev, 2012; Gormley et al., 2014; Abend et al., 2009).

In this paper, we study the scenario where we have small amounts of fully labeled training data. We sample 3% of the training data and an equivalent amount of development examples. The same training/dev subsets are used across all models.

Table 3 reports the performances of using 3% training data from CoNLL-05 and CoNLL-12 (top and bottom respectively). We compare our strong baseline model with structure-tuned models using all three constraints. Note that for all these evaluations, while we use subsamples of the dev set for model selection, the evaluations are reported using the full dev and test sets.

We see that training with constraints greatly improves precision with low training data, while recall reduces. This trade-off is accompanied by a reduction in the violation rates  $\rho_u$  and  $\rho_f$ . As noted in §2.4, models rarely predict label sequences that violate the exclusively overlapping roles constraint. As a result, the error rate  $\rho_o$  (the number of violations) only slightly fluctuates.

#### 3.3 Scenario 2: Large Training Data

Table 4 reports the performance of models trained with our framework using the full training set of

CoNLL-05 (3%, 1.1k)								
Dev	P	R	F1	$\delta$ F1	$\rho_u$	$\rho_o$	$\rho_f$	
RoBERTa <sup>2</sup>	67.79	<b>72.69</b>	70.15		14.56	23	6.19	
+U,F,O	<b>70.40</b>	71.91	<b>71.15</b>	1.0	8.56	20	5.82	
WSJ	P	R	F1	$\delta$ F1	$\rho_u$	$\rho_o$	$\rho_f$	
RoBERTa <sup>2</sup>	70.48	<b>74.96</b>	72.65		13.35	37	NA	
+U,F,O	<b>72.60</b>	74.13	<b>73.36</b>	0.7	7.46	49	NA	
Brown	P	R	F1	$\delta$ F1	$\rho_u$	$\rho_o$	$\rho_f$	
RoBERTa <sup>2</sup>	62.16	<b>66.93</b>	64.45		12.94	6	NA	
+U,F,O	<b>64.31</b>	65.64	<b>64.97</b>	0.5	5.47	6	NA	
CoNLL-12 (3%, 2.7k)								
Dev	P	R	F1	$\delta$ F1	$\rho_u$	$\rho_o$	$\rho_f$	
RoBERTa <sup>2</sup>	74.39	<b>76.88</b>	75.62		7.43	294	3.23	
+U,F,O	<b>75.99</b>	76.80	<b>76.39</b>	0.8	4.37	245	3.01	
Test	P	R	F1	$\delta$ F1	$\rho_u$	$\rho_o$	$\rho_f$	
RoBERTa <sup>2</sup>	74.79	<b>77.17</b>	75.96		6.92	156	2.67	
+U,F,O	<b>76.31</b>	76.88	<b>76.59</b>	0.6	4.12	171	2.41	

Table 3: Results on low training data (3% of CoNLL-05 and CoNLL-12). RoBERTa<sup>2</sup>: Baseline finetuned twice. U: Unique core roles. F: Frame core roles. O: Exclusively overlapping roles.  $\delta$ F1: improvement over baseline.  $\rho_f$  is marked NA for the CoNLL-05 test results because ground truth sense is unavailable on the CoNLL-05 shared task page.

CoNLL-05 (100%, 36k)								
Dev	P	R	F1	$\delta$ F1	$\rho_u$	$\rho_o$	$\rho_f$	
RoBERTa <sup>2</sup>	86.74	87.24	86.99		1.97	3.23		
+U,F,O	<b>87.24</b>	<b>87.26</b>	<b>87.25</b>	0.3	1.35	2.99		
Oracle					0.40	2.34		
WSJ	P	R	F1	$\delta$ F1	$\rho_u$	$\rho_f$		
RoBERTa <sup>2</sup>	87.75	87.94	87.85		1.71	NA		
+U,F,O	<b>88.05</b>	<b>88.00</b>	<b>88.03</b>	0.2	0.85	NA		
Oracle					0.30	NA		
Brown	P	R	F1	$\delta$ F1	$\rho_u$	$\rho_f$		
RoBERTa <sup>2</sup>	79.38	78.92	78.64		3.36	NA		
+U,F,O	<b>80.04</b>	<b>79.56</b>	<b>79.80</b>	1.2	1.24	NA		
Oracle					0.30	NA		

Table 4: Results on the *full* CoNLL-05 data. Oracle: Errors of oracle.  $\rho_o$  is in [0,6] across all settings.

the CoNLL-05 dataset which consists of 35k sentences with 91k propositions. Again, we compare RoBERTa (twice finetuned) with our structured models. We see that the constrained models consistently outperform baselines on the dev, WSJ, and Brown sets. With all three constraints, the constrained model reaches 88 F1 on the WSJ. It also

generalizes well on new domain by outperforming the baseline by 1.2 points on the Brown test set.

As in the low training data experiments, we observe improved precision due to the constraints. This suggests that even with large training data, direct label supervision might not be enough for neural models to pick up the rich output space structure. Our framework helps neural networks, even as strong as RoBERTa, to make more correct predictions from differentiable constraints.

Surprisingly, the development ground truth has a 2.34% error rate on the frame role constraint, and 0.40% on the unique role constraint. Similar percentages of unique role errors also appear in WSJ and Brown test sets. For  $\rho_o$ , the oracle has no violations on the CoNLL-05 dataset.

The exclusively overlapping constraint (*i.e.*  $\rho_o$ ) is omitted as we found models rarely make such prediction errors. After adding constraints, the error rate of our model approached the lower bound. Note that our framework focuses on the learning stage without any specialized decoding algorithms in the prediction phase except the Viterbi algorithm to guarantee that there will be no BIO violations.

### What about even larger and cleaner data?

The ideal scenario, of course, is when we have the luxury of massive and clean data to power neural network training. In Table 5, we present results on CoNLL-12 which is about 3 times as large as CoNLL-05. It consists of 90k sentences and 253k propositions. The dataset is also less noisy with respect to the constraints. For instance, the oracle development set has no violations for both the unique core and the exclusively overlapping constraints.

We see that, while adding constraints reduced error rates of  $\rho_u$  and  $\rho_f$ , the improvements on label consistency do not affect F1 much. As a result, our best constrained model performs on a par with the baseline on the dev set, and is slightly better than the baseline (by 0.1) on the test set. Thus we believe when we have the luxury of data, learning with constraints would become optional. This observation is in line with recent results in Li and Srikumar (2019) and Li et al. (2019).

### But is it due to the large data or the strong baseline?

To investigate whether the seemingly saturated performance is from data or from the model, we also evaluate our framework on the original

BERT (Devlin et al., 2019) which is relatively less powerful. We follow the same model setup for experiments and report the performances in Table 5 and Table 9. We see that compared to RoBERTa, BERT obtains similar F1 gains on the test set, suggesting performance ceiling is due to the train size.

CoNLL-12 (100%, 90k)						
Dev	P	R	F1	$\delta F1$	$\rho_u$	$\rho_f$
RoBERTa <sup>2</sup>	<b>86.62</b>	<b>86.91</b>	<b>86.76</b>		0.86	1.18
+U,F,O	86.60	86.89	86.74	0	0.59	1.04
Oracle					0	0.38
Test	P	R	F1	$\delta F1$	$\rho_u$	$\rho_f$
RoBERTa <sup>2</sup>	86.28	86.67	86.47		0.91	0.97
+U,F,O	<b>86.40</b>	<b>86.83</b>	<b>86.61</b>	0.1	0.50	0.93
Oracle					0	0.42
Dev	P	R	F1	$\delta F1$	$\rho_u$	$\rho_f$
BERT <sup>2</sup>	85.62	86.22	85.92		1.41	1.12
+U,F,O	<b>85.97</b>	<b>86.38</b>	<b>86.18</b>	0.3	0.78	1.07
Test	P	R	F1	$\delta F1$	$\rho_u$	$\rho_f$
BERT <sup>2</sup>	85.52	86.24	85.88		1.32	0.94
+U,F,O	<b>85.82</b>	<b>86.36</b>	<b>86.09</b>	0.2	0.79	0.90

Table 5: Results on CoNLL-12. BERT<sup>2</sup>: The original BERT finetuned twice.  $\rho_o$  is around 50 across all settings. With the luxury of large and clean data, constrained learning becomes less effective.

## 4 Ablations & Analysis

In §3, we saw that constraints not just improve model performance, but also make outputs more structurally consistent. In this section, we will show the results of an ablation study that adds one constraint at a time. Then, we will examine the sources of improved F-score by looking at individual labels, and also the effect of the top-k relaxation for the constraint  $O$ . Furthermore, we will examine the robustness of our method against randomness involved during training. We will end this section with a discussion about the ability of constrained neural models to handle structured outputs.

**Constraint Ablations** We present the ablation analysis on our constraints in Table 6. We see that as models become more constrained, precision improves. Furthermore, one class of constraints do not necessarily reduce the violation rate for the others. Combining all three constraints offers a balance between precision, recall, and constraint violation.

One interesting observation that adding the  $O$  constraints improve F-scores even though the  $\rho_o$  values were already close to zero. As noted in §2.4, our constraints apply to the predicted scores of all labels for a given argument, while the actual decoded label sequence is just the highest scoring sequence using the Viterbi algorithm. Seen this way, our regularizers increase the decision margins on affected labels. As a result, the model predicts scores that help Viterbi decoding, and, also generalizes better to new domains *i.e.*, the Brown set.

CoNLL-05 (100%, 36k)					
Dev	P	R	F1	$\rho_u$	$\rho_f$
RoBERTa <sup>2</sup>	86.74	87.24	86.99	1.97	3.23
+U	87.21	87.32	87.27	1.29	3.23
+U,F	87.19	<b>87.54</b>	<b>87.37</b>	<b>1.20</b>	3.11
+U,F,O	<b>87.24</b>	87.26	87.25	1.35	<b>2.99</b>
WSJ	P	R	F1	$\rho_u$	$\rho_f$
RoBERTa <sup>2</sup>	87.75	87.94	87.85	1.71	NA
+U	87.88	88.01	87.95	1.18	NA
+U,F	<b>88.05</b>	<b>88.09</b>	<b>88.07</b>	0.89	NA
+U,F,O	<b>88.05</b>	88.00	88.03	<b>0.85</b>	NA
Brown	P	R	F1	$\rho_u$	$\rho_f$
RoBERTa <sup>2</sup>	79.38	78.92	78.64	3.36	NA
+U	79.36	79.15	79.25	1.74	NA
+U,F	79.60	79.24	79.42	<b>1.00</b>	NA
+U,F,O	<b>80.04</b>	<b>79.56</b>	<b>79.80</b>	1.24	NA

Table 6: Ablation tests on CoNLL-05.

**Sources of Improvement** Table 7 shows label-wise F1 scores for each argument. Under low training data conditions, our constrained models gained improvements primarily from the frequent labels, e.g., A0-A2. On CoNLL-05 dataset, we found the location modifier (AM-LOC) posed challenges to our constrained models which significantly performed worse than the baseline. Another challenge is the negation modifier (AM-NEG), where our models underperformed on both datasets, particularly with small training data. When using the CoNLL-12 training set, our models performed on par with the baseline even on frequent labels, confirming that the performance of soft-structured learning is nearly saturated on the larger, cleaner dataset.

**Impact of Top- $k$  Beam Size** As noted in §2.4, we used the top- $k$  strategy to implement the constraint  $O$ . As a result, there is a certain chance for predicted label sequences to have non-exclusive

	CoNLL-05 3%		CoNLL-05 100%		CoNLL-12 3%		CoNLL-12 100%	
	RoBERTa <sup>2</sup>	+U,F,O	RoBERTa <sup>2</sup>	+U,F,O	RoBERTa <sup>2</sup>	+U,F,O	RoBERTa <sup>2</sup>	+U,F,O
A0	81.28	82.11	93.43	93.52	84.99	85.73	92.78	92.81
A1	72.12	73.59	89.23	89.80	78.36	79.67	89.88	89.75
A2	46.50	47.52	79.53	79.73	68.24	69.20	84.93	84.90
A3	39.58	42.11	81.45	81.86	33.26	34.47	72.96	73.24
A4	51.61	51.56	74.60	75.59	56.29	58.38	80.80	80.33
AM-ADV	44.07	47.56	66.67	66.91	55.26	54.93	66.37	66.92
AM-DIR	16.39	18.92	55.26	55.56	36.51	35.81	64.92	64.95
AM-DIS	71.07	70.84	80.20	80.50	76.35	76.40	82.86	82.71
AM-LOC	53.08	51.60	69.02	66.50	59.74	59.94	72.74	73.21
AM-MNR	44.30	44.18	68.63	69.87	56.14	55.67	70.89	71.13
AM-MOD	91.88	91.60	98.27	98.60	95.50	95.76	97.88	98.04
AM-NEG	91.18	88.35	94.06	93.60	93.29	93.05	95.93	95.83
AM-TMP	74.05	74.13	88.24	88.08	79.00	78.78	87.58	87.56
Overall	70.48	71.55	87.33	87.61	76.66	77.45	87.60	87.58

Table 7: Label-wise F1 scores for the CoNLL-05 and CoNLL-12 development sets.

overlap without our regularizer penalizing them. What we want instead is a good balance between coverage and runtime cost. To this end, we analyze the CoNLL-12 development set using the baseline trained on 3% of CoNLL-12 data. Specifically, we count the examples which have such overlap but the regularization loss is  $\leq 0.001$ . In Table 8, we see that  $k = 4$  yields good coverage.

k	1	2	4	6
# Ex.	10	8	3	2

Table 8: Impact of  $k$  for the top- $k$  strategy, showing the number of missed examples for different  $k$ . We set  $k = 4$  across all experiments.

**Robustness to random initialization** We observed that model performance with structured tuning is generally robust to random initialization. As an illustration, we show the performance of models trained on the full CoNLL-12 dataset with different random initializations in Table 9.

**Can Constrained Networks Handle Structured Prediction?** Larger, cleaner data may presumably be better for training constrained neural models. But it is not that simple. We will approach the above question by looking at how good the transformer models are at dealing with two classes of constraints, namely: 1) structural constraints that rely *only* on available decisions (constraint  $U$ ), 2) constraints involving external knowledge (constraint  $F$ ).

For the former, we expected neural models to perform very well since the constraint  $U$  repre-

CoNLL-12 (100%, 90k)				
Test F1	Seed1	Seed2	Seed3	avg $\delta$ F1
BERT <sup>2</sup>	85.88	85.91	86.13	0.1
+U,F,O	<b>86.09</b>	<b>86.07</b>	<b>86.19</b>	
Test F1	Seed1	Seed2	Seed3	avg $\delta$ F1
RoBERTa <sup>2</sup>	86.47	86.33	86.45	0.1
+U,F,O	<b>86.61</b>	<b>86.48</b>	<b>86.57</b>	

Table 9: F1 scores models trained on the CoNLL-12 data with different random seeds. The randomness affects the initialization of the classification layers and the batch ordering during training.

sents a simple local pattern. From Tables 4 and 5, we see that the constrained models indeed reduced violations  $\rho_u$  substantially. However, when the training data is limited, *i.e.*, comparing CoNLL-05 3% and 100%, the constrained models, while reducing the number of errors, still make many invalid predictions. We conjecture this is because networks learn with constraints mostly by memorization. Thus the ability to generalize learned patterns on unseen examples relies on training size.

The constraint  $F$  requires external knowledge from the PropBank frames. We see that even with large training data, constrained models were only able to reduce error rate  $\rho_f$  by a small margin. In our development experiments, having larger  $\lambda_F$  tends to strongly sacrifice argument F1, yet still does not to improve development error rate substantially. Without additional training signal in the form of such background knowledge, constrained inference becomes a necessity, even with strong neural network models.



## 5 Discussion & Conclusion

**Semantic Role Labeling & Constraints** The SRL task is inherently knowledge rich; the outputs are defined in terms of an external ontology of frames. The work presented here can be generalized to several different flavors of the task, and indeed, constraints could be used to model the interplay between them. For example, we could revisit the analysis of Yi et al. (2007), who showed that the PropBank A2 label takes on multiple meanings, but by mapping them to VerbNet, they can be disambiguated. Such mappings naturally define constraints that link semantic ontologies.

Constraints have long been a cornerstone in the SRL models. Several early linear models for SRL (e.g. Punyakanok et al., 2004, 2008; Surdeanu et al., 2007) modeled inference for PropBank SRL using integer linear programming. Riedel and Meza-Ruiz (2008) used Markov Logic Networks to learn and predict semantic roles with declarative constraints. The work of (Täckström et al., 2015) showed that certain SRL constraints admit efficient decoding, leading to a neural model that used this framework (FitzGerald et al., 2015). Learning with constraints has also been widely adopted in semi-supervised SRL (e.g., Fürstenaу and Lapata, 2012).

With the increasing influence of neural networks in NLP, however, the role of declarative constraints seem to have decreased in favor of fully end-to-end training (e.g., He et al., 2017b; Strubell et al., 2018, and others). In this paper, we show that even in the world of neural networks with contextual embeddings, there is still room for systematically introducing knowledge in the form of constraints, without sacrificing the benefits of end-to-end learning.

**Structured Losses** Chang et al. (2012) and Ganchev et al. (2010) developed models for structured learning with declarative constraints. Our work is in the same spirit of training models that attempts to maintain output consistency.

There are some recent works on the design of models and loss functions by relaxing Boolean formulas. Kimmig et al. (2012) used the Łukasiewicz t-norm for probabilistic soft logic. Li and Srikumar (2019) augment the neural network architecture itself using such soft logic. Xu et al. (2018) present a general framework for

loss design that does not rely on soft logic. Introducing extra regularization terms to a downstream task have been shown to be beneficial in terms of both output structure consistency and prediction accuracy (e.g., Minervini and Riedel, 2018; Hsu et al., 2018; Mehta et al., 2018; Du et al., 2019; Li et al., 2019).

**Final words** In this work, we have presented a framework that seeks to predict structurally consistent outputs without extensive model redesign, or any expensive decoding at prediction time. Our experiments on the semantic role labeling task show that such an approach can be especially helpful in scenarios where we do not have the luxury of massive annotated datasets.

### Acknowledgements

We thank members of the NLP group at the University of Utah for their valuable insights and suggestions; and reviewers for pointers to related works, corrections, and helpful comments. We also acknowledge the support of NSF Cyberlearning-1822877, SaTC-1801446, U.S. DARPA KAIROS Program No. FA8750-19-2-1004, DARPA Communicating with Computers DARPA 15-18-CwC-FP-032, HDTRA1-16-1-0002, and gifts from Google and NVIDIA.

The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of DARPA or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright annotation therein.

### References

- Omri Abend, Roi Reichart, and Ari Rappoport. 2009. Unsupervised argument identification for semantic role labeling. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*.
- Alan Akbik, Vishwajeet Kumar, and Yunyao Li. 2016. Towards semi-automatic generation of proposition Banks for low-resource languages. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*.
- Jonathan Berant, Vivek Srikumar, Pei-Chun Chen, Abby Vander Linden, Brittany Harding, Brad Huang, Peter Clark, and Christopher D. Manning. 2014.

- Modeling biological processes for reading comprehension. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Xavier Carreras and Lluís Màrquez. 2005. Introduction to the CoNLL-2005 shared task: Semantic role labeling. In *Proceedings of the Ninth Conference on Computational Natural Language Learning (CoNLL-2005)*.
- Ming-Wei Chang, Lev Ratinov, and Dan Roth. 2012. Structured Learning with Constrained Conditional Models. *Machine learning*, 88(3):399–431.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*.
- Xinya Du, Bhavana Dalvi, Niket Tandon, Antoine Bosselut, Wen tau Yih, Peter Clark, and Claire Cardie. 2019. Be Consistent! Improving Procedural Text Comprehension using Label Consistency. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- Nicholas FitzGerald, Oscar Täckström, Kuzman Ganchev, and Dipanjan Das. 2015. Semantic role labeling with neural network factors. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 960–970.
- Hagen Fürstenau and Mirella Lapata. 2009. Graph alignment for semi-supervised semantic role labeling. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*.
- Hagen Fürstenau and Mirella Lapata. 2012. Semi-supervised semantic role labeling via structural alignment. *Computational Linguistics*, 38(1):135–171.
- Kuzman Ganchev, Jennifer Gillenwater, Ben Taskar, et al. 2010. Posterior Regularization for Structured Latent Variable Models. *Journal of Machine Learning Research*.
- Matthew R. Gormley, Margaret Mitchell, Benjamin Van Durme, and Mark Dredze. 2014. Low-resource semantic role labeling. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.
- Luheng He, Kenton Lee, Mike Lewis, and Luke Zettlemoyer. 2017a. Deep semantic role labeling: What works and what’s next. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.
- Luheng He, Kenton Lee, Mike Lewis, and Luke Zettlemoyer. 2017b. Deep semantic role labeling: What works and what’s next. In *ACL*, volume 1, pages 473–483.
- Wan-Ting Hsu, Chieh-Kai Lin, Ming-Ying Lee, Kerui Min, Jing Tang, and Min Sun. 2018. A Unified Model for Extractive and Abstractive Summarization using Inconsistency Loss. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*.
- Richard Johansson and Pierre Nugues. 2008. Dependency-based semantic role labeling of PropBank. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*.
- Angelika Kimmig, Stephen Bach, Matthias Broecheler, Bert Huang, and Lise Getoor. 2012. A short Introduction to Probabilistic Soft Logic. In *Proceedings of the NIPS Workshop on Probabilistic Programming: Foundations and Applications*.
- Erich Peter Klement, Radko Mesiar, and Endre Pap. 2013. *Triangular Norms*. Springer Science & Business Media.
- Meghana Kshirsagar, Sam Thomson, Nathan Schneider, Jaime Carbonell, Noah A. Smith, and Chris Dyer. 2015. Frame-semantic role labeling with heterogeneous annotations. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*.
- Tao Li, Vivek Gupta, Maitrey Mehta, and Vivek Srikumar. 2019. A logic-driven framework for consistency of neural models. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*.
- Tao Li and Vivek Srikumar. 2019. Augmenting Neural Networks with First-order Logic. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Sanket Vaibhav Mehta, Jay Yoon Lee, and Jaime Carbonell. 2018. Towards Semi-Supervised Learning for Deep Semantic Role Labeling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.
- Pasquale Minervini and Sebastian Riedel. 2018. Adversarially Regularising Neural NLI Models to Integrate Logical Background Knowledge. In *Proceedings of the 22nd Conference on Computational Natural Language Learning*.

- Hiroki Ouchi, Hiroyuki Shindo, and Yuji Matsumoto. 2018. A span selection model for semantic role labeling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.
- Martha Palmer, Daniel Gildea, and Nianwen Xue. 2010. Semantic role labeling. *Synthesis Lectures on Human Language Technologies*, 3(1):1–103.
- Simone Paolo Ponzetto and Michael Strube. 2006. Exploiting semantic role labeling, WordNet and Wikipedia for coreference resolution. In *Proceedings of the Human Language Technology Conference of the NAACL, Main Conference*.
- Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Hwee Tou Ng, Anders Björkelund, Olga Uryupina, Yuchen Zhang, and Zhi Zhong. 2013. Towards robust linguistic analysis using OntoNotes. In *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*.
- Vasin Punyakanok, Dan Roth, and Wen-tau Yih. 2005. The necessity of syntactic parsing for semantic role labeling. In *IJCAI*.
- Vasin Punyakanok, Dan Roth, and Wen-tau Yih. 2008. The importance of syntactic parsing and inference in semantic role labeling. *Computational Linguistics*.
- Vasin Punyakanok, Dan Roth, Wen-tau Yih, and Dav Zimak. 2004. Semantic role labeling via integer linear programming inference. In *COLING 2004: Proceedings of the 20th International Conference on Computational Linguistics*.
- Sebastian Riedel and Ivan Meza-Ruiz. 2008. Collective semantic role labelling with markov logic. In *CoNLL 2008: Proceedings of the Twelfth Conference on Computational Natural Language Learning*.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*.
- Emma Strubell, Patrick Verga, Daniel Andor, David Weiss, and Andrew McCallum. 2018. Linguistically-Informed Self-Attention for Semantic Role Labeling. *EMNLP*.
- Mihai Surdeanu, Lluís Màrquez, Xavier Carreras, and Pere R Comas. 2007. Combination strategies for semantic role labeling. *Journal of Artificial Intelligence Research*, 29:105–151.
- Oscar Täckström, Kuzman Ganchev, and Dipanjan Das. 2015. Efficient inference and structured learning for semantic role labeling. *Transactions of the Association for Computational Linguistics*.
- Ivan Titov and Alexandre Klementiev. 2012. Semi-supervised semantic role labeling: Approaching from an unsupervised perspective. In *Proceedings of COLING 2012*.
- Kristina Toutanova, Aria Haghighi, and Christopher Manning. 2005. Joint learning improves semantic role labeling. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL’05)*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*.
- Jingyi Xu, Zilu Zhang, Tal Friedman, Yitao Liang, and Guy Van den Broeck. 2018. A Semantic Loss Function for Deep Learning with Symbolic Knowledge. In *International Conference on Machine Learning*.
- Szu-ting Yi, Edward Loper, and Martha Palmer. 2007. [Can Semantic Roles Generalize Across Genres?](#) In *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Proceedings of the Main Conference*, pages 548–555, Rochester, New York. Association for Computational Linguistics.
- Zhuosheng Zhang, Yuwei Wu, Hai Zhao, Zuchao Li, Shuailiang Zhang, Xi Zhou, and Xiang Zhou. 2020. Semantics-aware bert for language understanding. *AAAI Conference on Artificial Intelligence (AAAI)*.

## A Appendices

### A.1 Hyperparameters

We show the hyperparameters of  $\lambda$ 's in Table 10. We conducted grid search on the combinations of  $\lambda$ 's for each setting and the best one on development set is selected for reporting.

Model	$\lambda_U$	$\lambda_O$	$\lambda_F$
RoBERTa CoNLL-05 (3%) +U,F,O	2	0.5	0.5
RoBERTa CoNLL-2012 (3%) +U,F,O	1	2	1
RoBERTa CoNLL-05 (100%) +U +U,F +U,F,O	1 1 1	 0.5 0.5	  0.1
RoBERTa CoNLL-2012 (100%) +U,F,O	1	1	0.1
BERT CoNLL-2012 (100%) +U,F,O	0.5	1	0.1

Table 10: Values of hyperparameter  $\lambda$ 's.