

ITER: Iterative Transformer-based Entity Recognition and Relation Extraction

Anonymous NAACL-HLT 2021 submission

Abstract

Entity Recognition and Relation Extraction are essential components in extracting structured information from text. Recent advances for both tasks generate a structured representation of the information in an autoregressive fashion, a time-intensive and computationally expensive approach. This raises the natural question whether autoregressive methods are necessary in order to achieve comparable results. In this work, we propose ITER, a functionally more expressive, non-autoregressive model, that unifies several improvements to a recent language modeling approach: ITER improves inference throughput by up to $23\times$, is capable of handling all types of nested entities and effectively halves the number of required parameters in comparison. Furthermore, we achieve a SOTA result of 84.30 F1 for the relation extraction dataset ADE and demonstrate comparable performances for both named entity recognition with GENIA and CoNLL03 as well as for relation extraction with CoNLL04 and NYT.

1 Introduction

In recent years, there has been a shift towards using autoregressive methods in many common NLP tasks. Parallel to this development is an increasing focus on approaching NLP tasks such as relation extraction or (nested) named entity recognition as structured prediction problems. Given a sequence of text input, a given model autoregressively generates outputs that encode the structure contained within the input, which offers flexibility since the source and target vocabulary must not share any commonalities.

Flattening the output structure into a single string, preserving the information about the structure in the input and using an autoregressive model to learn to generate this adapted target language (Cabot and Navigli, 2021; Paolini et al., 2021), is an *implicit* approach known to work well across task boundaries (Raffel et al., 2020). In

this case, the target vocabulary typically contains the whole source language vocabulary. However, representing the structured output as a string introduces additional complexity when modeling intra-structure dependencies (Liu et al., 2022). More recently, Liu et al. proposed constraining the autoregressive model to *explicit* generation of the output structure. They define three types of basic actions to be performed at each generation step and use the T5 (Raffel et al., 2020) transformer to autoregressively generate the structure induced by said basic actions.

With this trend of using autoregressive methods for tasks such as relation extraction come however also several problems: As inference time scales quadratically with the output sequence length, language modeling approaches are prone to low inference speed¹ especially with increasing model parameters (Pope et al., 2022). While scaling the model size from hundreds of millions of parameters to billions of parameters yields performance increments for Liu et al., this scaling can become infeasible, both in terms of compute required and the environmental impact when using those large scale models in production.

This raises the natural question whether a *non-autoregressive process* capable of generating such an output structure can achieve similar performance whilst addressing the aforementioned limitations of language modeling approaches. In this paper, we present ITER, an encoder-only transformer-based relation extraction model that addresses the limitations of state-of-the-art architectures and show that the structured prediction problem can be approached without a language modeling objective in mind.

To summarize, our key contributions are the following:

1. We present ITER, a transformer-based encoder-

¹In terms of samples/second

081
082
083
084
085
086
087
088
089
090

091
092
093
094
095
096
097

098
099
100
101
102
103
104
105

106
107

108

109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128

only relation extraction model that addresses the limitations of autoregressive architectures. Instead of employing a language modeling objective, our model generates the structured output in three basic steps. We show that this encoder-based approach achieves performance competitive to language modeling architectures whilst retaining only half of the number of parameters and increasing the inference throughput by a factor of up to $23\times$.

2. We identify several drawbacks with the state-of-the-art architecture ASP (Liu et al., 2022), which limit the model’s expressivity for nested structures and generalization to test data. In our work, we address those limitations and translate the *autoregressive* approach into a three-stage process.
3. In our experiments, we observe that training ITER on NYT, CoNLL04, ADE (RE), or GENIA, CoNLL03 (NER) results in competitive performance across all datasets, while maintaining a significantly smaller size compared to SOTA models. We observe small average improvements of 0.5 F1 points on the ADE dataset.
4. We publish our implementation and checkpoints at GITHUB.COM/ANONYMOUS.

2 Related Work

The goal of relation extraction, sometimes also referred to as *end-to-end* relation extraction or *joint* entity and relation extraction, is to identify the names and types of *named entities*, inside a given text, as well as classify the *relationships* amongst these entities. (Grishman and Sundheim, 1996; Zhao and Grishman, 2005).

First approaches to relation extraction were to decompose the task into *named entity recognition* and *relation classification*, where the named entities are identified first, while the relationships between the found named entities are then classified in a second, separate stage that is being learned independently. This pipeline-based approach is known to be prone to error propagation (Zhong and Chen, 2021; Sui et al., 2020). Because of this known limitation, joint approaches modeling both tasks simultaneously have been introduced and have shown promising results (Gupta et al., 2016; Wang and Lu, 2020).

2.1 Span-based Techniques

Table-filling or span-based strategies were and still are viable approaches to modeling relation extraction and related tasks (Gupta et al., 2016; Wang and Lu, 2020; Joshi et al., 2020; Tang et al., 2022; Urchade et al., 2024). Recent examples of this include DiffusionNER (Shen et al., 2023) and UniRel (Tang et al., 2022). DiffusionNER emits a time-complexity scaling linearly in the number of entities in an input and the number of diffusion steps. Computing the interaction map for UniRel scales quadratically with the sum of the input size and the number of relation types. This quadratic scaling is also an issue for *autoregressive* techniques, where the inference time is scaling quadratically with the length of the output sequence (Shen et al., 2023).

2.2 Autoregressive Techniques

Modeling the task as a *seq2seq* problem has established itself as the state-of-the-art for relation extraction in the last couple of years (Cabot and Navigli, 2021; Wang et al., 2022; Paolini et al., 2021; Liu et al., 2022; Fei et al., 2022; Lu et al., 2022). Enabled by the Transformer (Vaswani et al., 2017), the task is then formulated as a translation objective: Given an example sentence, the model translates the input into a flattened string that encodes the structural information contained within the source text (Liu et al., 2022).

Both (m)REBEL (Cabot and Navigli, 2021; Cabot et al., 2023) and TANL (Paolini et al., 2021) translate the input sequence into a flattened output string, that, in the (m)REBEL case, also no longer resembles natural language. Paolini et al. augment the target output with information about entity types and relations to other named entities. Both models are finetuned to produce a target language specific to the task. A comparison of different model outputs is available in Table 1. Either model can also deal with nested entities, which is crucial when dealing with real-world data, as for example data from the biomedical domain is known to often contain nested entities (Finkel and Manning, 2009).

2.3 Limitations of Autoregressive Structured Prediction (ASP)

ASP (Liu et al., 2022) has shown that autoregressively generating a sequence of *actions* instead of generating (augmented) natural language yields

129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145

146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172

173
174

175
176
177

| Model | Output |
|-------------|--|
| REBEL | <triplet> Barack Obama <subj> Honolulu, Hawaii <obj> place of birth |
| TANL | [Barack Obama <i>person</i> <i>place of birth</i> = Honolulu, Hawaii] was born in [Honolulu, Hawaii <i>location</i>] |
| ASP | [* Barack Obama] was born in [* Honolulu, Hawaii] |
| ITER (ours) | Barack Obama was born in Honolulu, Hawaii [] [] [] [] |

Table 1: Comparison of different (autoregressive) relation extraction system outputs: REBEL (Cabot and Navigli, 2021), TANL (Paolini et al., 2021) and ASP (Liu et al., 2022). Input into all models was the sentence "Barack Obama was born in Honolulu, Hawaii". Comparing ASP and our method, one can observe that ASP requires [* and] actions to be placed alongside the input to model the same structure of the input.

model performance benefits not only for relation extraction, but for named entity recognition and co-reference resolution as well. At every generation step, their model can perform three distinct types of actions, *structure-building* actions, *bracket-pairing* and *span-labeling* actions. For the *structure-building actions* (Eq. 1), the model can either perform [* or] actions at the current generation position or `copy` the next token from the input into the output.

$$A = \{[*],], \text{copy}\} \quad (1)$$

The generation completes when all input tokens have been copied into the output. *Bracket-pairing* actions (Eq. 2) aim to connect the current position with a previously performed [* action, resulting in a span.

$$B_n = \{m \mid m < n \wedge a_m = [*]\} \quad (2)$$

Span-labeling actions allow both the labeling of individual spans and the linking of the current formed span to an earlier found span in the output sequence, modeling relationships between named entities (Eq. 3). For relation extraction, \mathcal{L} is instantiated as the cartesian product of the named entity and relation types: $\mathcal{L} = \mathbf{T}_E \times \mathbf{T}_R$.

$$Z_n = \{m \mid m < n \wedge a_m =]\} \times \mathcal{L} \quad (3)$$

The authors of ASP employ a conditional language model to learn to produce the optimal output structure $\mathbf{y}^* \in Y_1 \times \dots \times Y_N$ where Y_n is defined as

$$Y_n = A \times B_n \times Z_n \quad (4)$$

At every time-step their model will perform the three basic actions $(a_n, b_n, z_n) \in Y_n$ sourced from their respective sets for *structure-building* actions A , *bracket-pairing* actions B_n and *span-labeling* actions Z_n .

Said approach however is not capable of capturing *nested entities* in all edge cases. At every time-step, ASP can only complete one span with one preceding [* action due to the definition of B_n , hence multiple consecutive] actions would be required to properly model this behaviour². We also hypothesize that the structured prediction process for ASP suffers from suboptimal training due to the nature of the *span-labeling* actions Z_n . Linking the span formed at the current position to another span in the sequence is constrained by the fact that only links to spans that have been completed in the past (i.e. earlier in the sequence) are valid. As it is impossible to link to spans that will be found in the future (i.e. spans that come after the span ending at the current position), the authors of ASP introduce a directionality parameter to counteract the asymmetric property of the relations in the dataset.³

This prevents the two tuples (*Barack Obama*, *work_for*, *the american people*) and (*the american people*, *work_for*, *Barack Obama*) from being indistinguishable.

This is important as those two examples encode drastically different information. The directionality parameter however effectively doubles the number of relations ($\mathbf{T}'_R = \mathbf{T}_R \times \mathbb{B}$), leading to fewer training examples per relation type, and we hypothesize that this yields subpar training results.

Aside from those architectural issues, ASP and similar *seq2seq* Transformer models such as TANL or REBEL all suffer from the quadratically scaling time-complexity of generative architectures, significantly impacting their inference speed (Paolini

²Liu et al. noted issues in generating multiple [*], hence their workaround. We think] faces this issue as well, so we hypothesize that this hurts their models' performance.

³As can be seen in their official implementation: https://github.com/lyutyuh/ASP/blob/master/data/t5minimize_ere.py#L227

et al., 2021). This raises the question whether *eliminating* the requirement for a conditional language model from ASP can retain the same model performance whilst crucially reducing inference time, allowing for simpler identification of all nested entities and addressing the generalization issues.

3 Approach

We base our approach for ITER on the work of Liu et al.. In order to replace the autoregressive component of their approach with our inference process, several modifications to the structured prediction are necessary. We divide the process into three basic steps:

- (1) Firstly, identify all positions n in the input \mathbf{x} where a span is beginning, i.e. the \llbracket action is performed, using *is_left* (Eq. 8).
- (2) Following that, identify all positions $m \geq n$ in the input that pair with any of the previously identified positions n found in step 1, forming named entities. Position m is constrained to come after or at the same position as the beginning of the span. *is_span* (Eq. 9) returns a set of bracket pairings with earlier found \llbracket positions, which in turn can be used to determine the named entities in the input.
- (3) Finally, test for relationships amongst all pairs of named entities found in (2) using *is_link* (Eq. 11). This function will return a vector that, after applying the sigmoid function, contains a probability for each relationship type \mathbf{T}_R . If this probability exceeds 0.5, a relation between the two tested named entities will be predicted.

Each step can individually be computed in parallel, enabling an efficient implementation of our model.

The first necessary change is to transition to a smaller subset of the *structure-building* actions \mathcal{A} :

$$\mathcal{A} = \{ \llbracket, \rrbracket \}$$

Our model must be allowed to perform both \llbracket and \rrbracket actions at the same time, to not lose model expressiveness, otherwise it will not be able to correctly classify single-token spans⁴. Therefore, the

⁴Think of a single-token named entity $x_i = \text{BERLIN}$: the model must be able to determine the span of this entity, since it ends at the same position it started. So a_i must now be a set: $a_i = \{ \llbracket, \rrbracket \}$.

structure-building actions A_n performed at every position n must now be a subset of \mathcal{A} , to allow for this behavior. This is reflected in the definition for our structured output $\mathbf{y}^* \in \mathcal{Y}_1 \times \dots \times \mathcal{Y}_N$:

$$\mathcal{Y}_n = \wp(\mathcal{A}) \times \wp(\mathcal{B}_n) \quad (5)$$

where \wp is the powerset operation. Since we test for relationships only in step 3, we remove the *span-labeling* actions Z_n (Eq. 3).

To be able to properly handle nested entities, or, more specifically, two or more entities ending at the same position, the *bracket-pairing* actions are also present in \mathcal{Y}_n as a subset of all possible such actions \mathcal{B}_n . This change comes in combination with two adjustments to the definition of \mathcal{B}_n itself.

At position m , ITER will be allowed to pair \llbracket actions at positions $n \leq m$, circumventing single-token named entity issues (1, Eq. 6) and each pairing is allowed to hold its own named entity type $t \in \mathbf{T}_E$ (2, Eq. 6).

$$\mathcal{B}_n = \{ n \mid n \leq m \wedge \llbracket \in A_n \} \times \mathbf{T}_E^{(2)} \quad (6)$$

3.1 Identifying Named Entities

Before relationships can be determined, spans need to be uniquely identified by their start and end positions in combination with their type of the named entity in the input sequence. Before any of the three following generation steps, the input \mathbf{x} is given into the encoder of the base model, T5 in our case, which produces a sequence of contextualized vector representations $\mathbf{h} = \langle h_1 \dots h_N \rangle$ for \mathbf{x} with $h_n \in \mathbb{R}^\delta$ where δ corresponds to the hidden dimension size used in the base model. All three stages use gated feed-forward networks of the following form:

$$FFN_\kappa^\psi(\mathbf{h}) = ((\text{GELU}(\mathbf{h}W_a) \otimes \mathbf{h}W_i)W_o) \quad (7)$$

where $W_a, W_i \in \mathbb{R}^{\psi\delta \times 2\psi\delta}$, $W_o \in \mathbb{R}^{2\psi\delta \times \kappa}$ are linear projections learned during training, GELU is the gaussian error linear unit function (Hendrycks and Gimpel, 2016) and κ is the output dimensionality. $\psi = 2$ when two concatenated vectors, indicated by $h_{m,n}$, are input, otherwise $\psi = 1$.

3.1.1 Determining Starting Positions of Named Entities

To identify said spans, the model learns to predict the positions where the spans of named entities in the input \mathbf{x} are beginning. This task is modeled by the function *is_left* (Eq. 8), which receives a

sequence of hidden states \mathbf{h} as input and the output is an equally sized sequence of Boolean values $\langle b_1 \dots b_N \rangle \in \mathbb{B}^N$:

$$is_left(\mathbf{h})_n = FFN_{\kappa=1}^{\psi=1}(h_n) > 0 \quad (8)$$

At all positions where $is_left(\mathbf{h})$ holds true, the left bracket action $\mathbf{1}$ is included in the set of actions A_n performed at position n .

3.1.2 Pairing Left and Right Brackets

After determining where spans of named entities start in the input \mathbf{x} , the next step is to identify which positions \mathbf{x}_m ($m \geq n$) following \mathbf{x}_n in the input form a span of named entity type $t \in \mathbf{T}_E$. Our model learns a projection is_span , that maps the input \mathbf{h} and left bracket positions n , $\mathbf{1} \in A_n$ to a sequence of tuples of indices and entity types (n, t) :

$$is_span : \mathbb{R}^{\delta \times N} \times \mathbb{B}^N \rightarrow \wp(\mathcal{B}_1) \times \dots \times \wp(\mathcal{B}_N) \quad (9)$$

$$is_span(\mathbf{h}, \mathbf{b})_m = \{(n, t) \mid FFN_{\kappa}^{\psi}(h_{m,n})_t > 0\}$$

with $\kappa = |\mathbf{T}_E|$, $\psi = 2$, $\mathbf{b}_n = \top$. For each position m where the output of $B_m = is_span(\mathbf{h})_m$ is non-empty, ITER performs action $\mathbf{1}$ at position m . Each element $(n, t) \in B_m$ determines a pairing from the left bracket at position n with a right bracket at position m of type t , forming a named entity. If a left bracket from step one is left unlinked, no named entity will be identified. The first two stages are visualized in Figure 1.

3.2 Identifying Relationships amongst Named Entities

The third step now tests pairs of identified named entities for their relationship with each other. For the non-nested case, is_link projects two hidden states h_i and h_j to a vector of non-normalized logits, resembling probabilities after applying the sigmoid function (Eq. 10).

$$is_link : \mathbb{R}^{\delta} \times \mathbb{R}^{\delta} \rightarrow \mathbb{R}^{\kappa} \quad (10)$$

$$is_link(h_i, h_j) = \sigma(FFN_{\kappa}^{\psi}(h_{i,j}))$$

where $\kappa = |\mathbf{T}_R|$, $\psi = 2$.

The hidden states are obtained from the positions i, j in the input \mathbf{h} where there is $\mathbf{1} \in A_{i,j}$ and pairings from x_i to x_o ($(o, t_1) \in \mathcal{B}_i$) and similarly for x_j with x_p . is_link outputs vectors containing probabilities for relationships between two entities identified in (1) and (2). During inference, all combinations of found entities are tested for

relationships. The ordering of head and tail entity is important, so $is_link(h_i, h_j) \neq is_link(h_j, h_i)$. The abstraction of using just the hidden state from the last position of the span no longer works when dealing with nested entities, as spans are no longer uniquely identified by their last position, and as such the hidden state of the first position will also be included for such cases:

$$is_link : \mathbb{R}^{\delta} \times \mathbb{R}^{\delta} \times \mathbb{R}^{\delta} \times \mathbb{R}^{\delta} \rightarrow \mathbb{R}^{\kappa}$$

$$is_link(h_i, h_o, h_j, h_p) = \sigma(FFN_{\kappa}^{\psi}(h_{i,o,j,p})) \quad (11)$$

where now $\psi = 4$. This final step is visualized in Figure 2.

3.3 Training

The model of choice for this paper is the T5 (Raffel et al., 2020) Transformer architecture, which can also be used as an encoder, albeit primarily trained for autoregressive applications (Raffel et al., 2020). In order to circumvent error propagation between the three stages of ITER, training will include all three task functions simultaneously: is_left , is_span and is_link . ITER receives as input a sequence of hidden representations (*hidden states*) $\mathbf{h} = \langle h_1, h_2, \dots, h_N \rangle$. The sequence of representations is shared across all three tasks. The loss function used during training is included in the Appendix, in Equations 13, 14 and 15. In a nutshell, in order to minimize the training loss, the model is incentivized to assign weights greater than zero to the respective correct decisions in all three cases, impacting the decisions taken by is_left , is_span and is_link .

3.4 Complexity

In this section, we will briefly discuss the theoretical time complexity of our approach. Steps (1) and (2) can both be parallelized across the sequence dimension. As is_left only uses linear projections and activation functions, its runtime is bound by the length of the input sequence \mathbf{h} , yielding a linear time-complexity $O(N)$. is_span can be optimized to only consider the closest ω left brackets, and in the trivial case we set $\omega = 1$. For nested named entity datasets, ω is another hyperparameter to be adjusted accordingly, but in the trivial case, is_span can be optimized to perform exactly one pass through the FFN per element of the sequence, thereby yielding a time-complexity $O(N)$ linear in the input length N . For $\omega > 1$, we have $O(N * \omega)$, however $\omega \ll N$. Since steps (1) and

| INPUT | | Barack | Obama | was born in Honolulu , | | | | Hawaii | |
|-----------------------------|--|--------|------------|------------------------|---|---|---|---------------------------|---|
| POSITION | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| (1) | $is_left(\mathbf{h}) = \mathbf{b}$ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✓ |
| | \Downarrow | ↓ | {(1, PER)} | | | | ↓ | {(6, LOC), (8, STATE)} | |
| (2) | $is_span(\mathbf{h}, \mathbf{b}) = B_n$ | | | | | | | | |
| | \Downarrow | | | | | | | | |
| $A_n \subseteq \mathcal{A}$ | | [| ↔ |] | | | [| ↔ |] |
| | | | | | | | | [↔] | |

Figure 1: Visualization of stages one (is_left) and two (is_span) of the model. is_left yields three positions where spans are beginning: 1,6 and 8 (Stage 1). is_span then creates pairings of types *person* between position 2 and 1, *location* between position 8 and 6 and *state* for position 8, a 1-token span (indicated by [←]).

| FUNCTION | HEAD ENTITY | TAIL ENTITY | OUTPUT |
|--------------------------------|----------------------------|----------------------------|--|
| $is_link(h_1, h_2, h_6, h_8)$ | Barack Obama (1, 2) | Honolulu, Hawaii (6, 8) | <i>live_in</i> <i>born_in</i> <i>work_in</i> (0.05 , 0.98 , 0.01) |
| $is_link(h_1, h_2, h_8, h_8)$ | Barack Obama (1, 2) | Hawaii (8, 8) | <i>live_in</i> <i>born_in</i> <i>work_in</i> (0.17 , 0.45 , 0.02) |
| ⋮ | ⋮ | ⋮ | ⋮ |
| $is_link(h_6, h_8, h_8, h_8)$ | Honolulu, Hawaii (6, 8) | Hawaii (8, 8) | <i>live_in</i> <i>born_in</i> <i>work_in</i> (0.03 , 0.07 , 0.1) |

Figure 2: Visualization of the third stage of the model. The output is already sigmoid-normalized. For illustration purposes, there are three relation types: $\mathbf{T}_R = \langle live_in, born_in, work_in \rangle$. The named entity "Barack Obama" stands in relationship "born_in" with "Honolulu, Hawaii", as $0.98 > 0.5$, which is the criterion defined for this model.

(2) are performed sequentially, combining those remains bound by the sequence length N . Testing for relationships in step (3) requires testing all combinations of entities found and thus yields a quadratic runtime, but not in the sequence length, but in the number of entities E with $E \ll N$. Using ITER thus gives a complexity of $O(N + E^2)$ with E being number of entities.

4 Experimental Results

In this section, we give an overview over the datasets used in our experiments (Section 4.1) and we conclude with our interpretation of the results from said experiments (Section 4.2). Details about the hyperparameter search we performed can be found in the Appendix (Section D)

4.1 Data

To experimentally verify that our model can achieve performances on par or even higher than the baseline from ASP, we used 5 datasets from

two different domains and tasks: CoNLL03 (Sang and Meulder, 2003, NER) (NER), CoNLL04 (Roth and Yih, 2004, RE) and NYT (Riedel et al., 2010, RE) were all annotated from news articles while ADE (Gurulingappa et al., 2012, RE) and GENIA (Kim et al., 2003, NER) contain training examples with biomedical context. Of those five datasets, three contain nested entities (NYT, ADE and GENIA), something that ASP cannot properly model, as shown in Section 2.3, which was another factor for our selection. This portfolio of datasets allows us to verify our claims across a wide range of applications and different levels of data complexity. An overview regarding the datasets can be obtained in Table 8. Following Li et al.; Eberts and Ulges and Cabot and Navigli, we evaluate our model in a *strict* setting: A predicted relation between two entities is only considered correct, if both the span and type of the entity match the gold standard. We report *micro* F1 scores, unless stated otherwise.

| Transformer base model | ASP examples/s | ITER examples/s | Speedup over ASP |
|---------------------------|-------------------|--------------------|---------------------|
| T5 (small) | 44.459 | 1040.55 | $\times 23.41$ |
| T5 (large) | 27.177 | 605.398 | $\times 22.28$ |
| T5 (3b) | 29.427 | 334.843 | $\times 11.38$ |

Table 2: Comparing the inference throughput (as *samples per second*) of ITER versus the autoregressive approach ASP (Liu et al., 2022) on the test set of CoNLL04. Computations were run on a single NVIDIA H100 GPU, a batch size of 64 combined with 10 epochs of training beforehand. Speedups of up to $\times 23.41$ ($\times 18.6$ on avg.) are achieved when using ITER instead of ASP.

| Architecture | Size | NER F1 (strict) | RE F1 (strict) |
|---------------------------------|-------|----------------------|----------------------|
| ITER + FLAN T5 _{large} | 374 M | 89.770 ± 0.51 | 75.175 ± 0.39 |
| ITER + BERT _{large} † | 340 M | ~ 86.02 | ~ 66.922 |
| ASP + FLAN T5 _{base} | 247 M | 89.4 | 73.8 |
| ASP + FLAN T5 _{large} | 783 M | 90.5 | 76.2 |
| TANL | 222 M | 89.8 | 72.6 |
| REBEL (<i>pretrained</i>) | 406 M | - | 75.4 |
| DeepStruct | 10 B | 88.4 | 72.8 |
| DeepStruct (<i>finetuned</i>) | 10 B | 90.70 | 78.3 |
| DiffusionNER | 381M | 92.78 | - |
| TF-MTRNN | - | 93.6 | 72.1 |
| Wang and Lu | - | 90.1 | 73.6 |
| Lu et al. | >770M | - | 75.00 |
| Fei et al. | >770M | - | 75.3 |

Table 3: Final training results for CoNLL04, averaged across five runs for each configuration. † Preliminary test with other Encoder achitecture.

4.2 Results

With the encoder of the FLAN-T5-large model as a base, ITER achieves state-of-the-art results on ADE with on average 0.5 F1 points of improvement (Table 5). Furthermore, it reaches good results when compared to most generative approaches while the number of parameters is significantly smaller (Table 3, 9, 6). Specifically, its performance closely aligns with that of ASP+FLAN T5_{base} and ASP+FLAN T5_{large}, both of which possess a similar parameter count, with the latter having twice the parameters and only being slightly better. Table 2 answers another research question, which was to demonstrate that a higher inference speed can be obtained with a smaller model while reaching comparable results. Especially compared with DeepStruct our model performs well, considering its size and train-

ing time. DIFFUSIONNER performs exceptionally well, and we are not able to match its performance on the NER task, only coming close on CoNLL03. Again, supporting our hypothesis that encoder-only models —like DIFFUSIONNER—can outperform generative models like DeepStruct on structure prediction tasks.

| Architecture | NER F1 (strict) | RE F1 (strict) |
|----------------------------------|--------------------|-------------------|
| ITER + FLAN T5 _{large} | 94.726 \pm 0.16 | 90.707 \pm 0.34 |
| TANL | 94.9 | 90.8 |
| REBEL | - | 92.0 |
| DeepStruct (<i>multi-task</i>) | 95.4 | 93.7 |
| UniRel | - | 93.7 |
| Fei et al. | - | 94.2 |
| Lu et al. | - | 93.54 |

Table 4: Final training results for NYT.

| Architecture | NER F1 (strict) | RE F1 (strict) |
|---------------------------------|----------------------|----------------------|
| ITER + FLAN T5 _{large} | 91.907 ± 0.72 | 84.300 ± 1.52 |
| TANL (<i>multi-task</i>) | 91.2 | 83.8 |
| REBEL (<i>pretrained</i>) | - | 82.2 |
| DeepStruct (<i>finetuned</i>) | 91.1 | 83.8 |
| Wang and Lu | 89.7 | 80.1 |
| Yan et al. | 91.3 | 83.2 |

Table 5: Final training results for ADE with 10-fold cross-validation. F1 metrics are *macro*-averaged.

To further improve our understanding of ITER and its shortcomings, we analyse ADE using confusion matrices. Figure 3 shows that our model does not struggle with the span-prediction task. The model also learned to predict the actions [and] at the same step where appropriate. The main challenge seems to be the named entity type *Adverse-Effekt*, which is falsely predicted and missed several times. Furthermore, when considering the confusion matrices and the correlation between NER and RE+ score for ADE, it becomes apparent that a good NER score mostly leads to better RE+ results. Our confusion matrices show that ITER never confuses the two entity types and when predicting actions it also almost never confuses: [,] and the one token case with each other, which in turn leads to better results.

| Architecture | Dataset | F1 (<i>strict</i>) |
|----------------------------------|---------|----------------------|
| ITER + FLAN T5 _{large} | CoNLL03 | 91.593 \pm 0.39 |
| ASP + T5 _{base} | | 91.8 |
| ASP + T5 _{large} | | 92.8 |
| DeepStruct (<i>multi-task</i>) | | 93.10 |
| DiffusionNER | | 92.78 |
| TF-MTRNN | | 86.80 |
| Lu et al. | | 92.99 |
| Fei et al. | | 93.2 |
| ITER + FLAN T5 _{large} | GENIA | 80.153 \pm 0.25 |
| TANL | | 76.4 |
| DeepStruct (<i>finetuned</i>) | | 80.8 |
| DiffusionNER | | 81.53 |

Table 6: Final training results for CoNLL03 and GENIA. For CoNLL03, we trained ITER with a linear learning rate schedule, instead of the SMAC3 prediction to use a constant schedule, as the final performance did significantly degrade using a constant schedule.

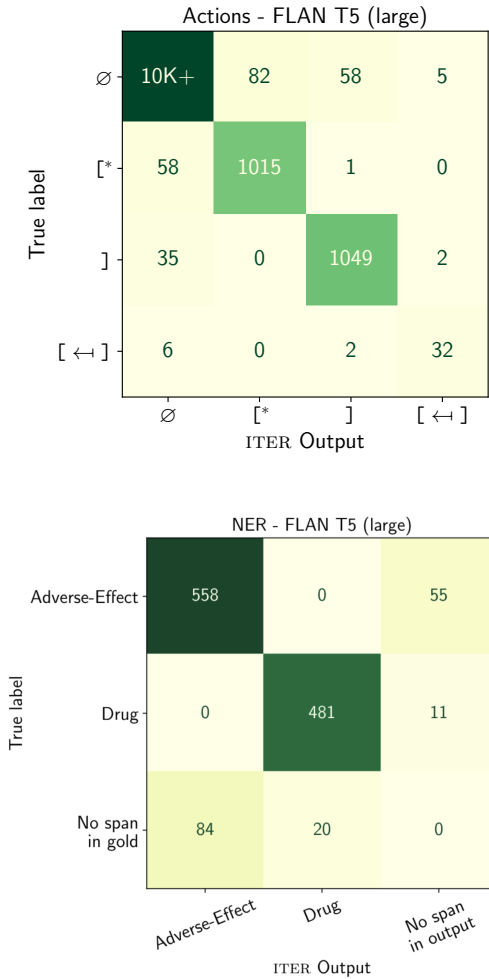


Figure 3: Confusion matrix for actions \mathcal{A} and NER on ADE.

5 Conclusion

In this paper, we identified several key drawbacks in a fairly recent state-of-the-art method for relation extraction, ASP, and proposed several improvements to counteract those issues. We investigated whether it is possible to translate the autoregressive process into a constant-in-time-complexity approach, whilst maintaining an equal level of performance.

We unify the aforementioned proposed improvements together with a new three-stage process in ITER, an encoder-based *non-autoregressive* relation extraction model. Our model achieves performances on par with state-of-the-art methods on all datasets and sets a new state-of-the-art on ADE of 84.3 F1, whilst being functionally more expressive and reducing inference time significantly, when compared to ASP. In our experiments, we highlight the time saving benefits of encoder-based models over autoregressive *seq2seq* approaches, suggesting that they are just as viable in a structure prediction task.

6 Future Work

While our model is currently built on top of a T5 encoder stack, it might be insightful to explore the performance of this architecture with other pre-trained (encoder-only) language models such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019) or the Nyströmformer (Xiong et al., 2021)

One area of future work might be exploring an even larger set of datasets. Unfortunately, there exist no benchmark suites for relation extraction, which itself might be an area of future work. While most datasets are open-source, there exist proprietary datasets, preventing the democratization of research in NLP and in machine learning in general. Evaluating architectures on a diverse portfolio of datasets instead of a limited amount of selected or hand-picked datasets should also allow to gain more significant insights into the performance, capabilities and limitations of relation extraction systems in general.

7 Limitations

One of ITER’s limitations are named entities that are not directly contained in the input text. This issue can arise when combining the NER stage of ITER with tasks such as entity linking, where the task is then to not only identify the named entity

and its type, but also to link said entity to a knowledge base entity, something particularly interesting when using a relation extraction pipeline to create knowledge graphs. The severity of this limitation strongly depends on the datasets used, we focussed on experiments with datasets where this issue cannot surface.

Furthermore, we are not able to fully compare ITER with ASP since we are not able to evaluate it on ACE05. If it is the case that the input has not been preprocessed, our model also requires a very tedious preprocessing-step that requires the programmer to correctly align the input string with the tokens that the model will be trained on. This is a limitation of the *sentencepiece* (Kudo and Richardson, 2018) tokenizer used in our experiments, as the tokenization process does not guarantee entity-level boundaries being respected during tokenization, meaning that a token spanning the characters i to j might contain the beginning of a span k ($i < k < j$). While generative approaches can circumvent this problem by introducing additional tokens into the target language text, encoder-based approaches such as our work are limited to dealing with this issue pre-tokenization.

Another limitation of ITER would be the strong task-dependent design of the functions *is_left*, *is_span* and *is_link*. This prevents a few-shot task transfer without finetuning for new relations or entity types.

References

- Pere-Lluís Huguet Cabot and Roberto Navigli. 2021. [REBEL: relation extraction by end-to-end language generation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 16-20 November, 2021*, pages 2370–2381. Association for Computational Linguistics.
- Pere-Lluís Huguet Cabot, Simone Tedeschi, Axel-Cyrille Ngonga Ngomo, and Roberto Navigli. 2023. [Red^{fm}: a filtered and multilingual relation extraction dataset](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2023, Toronto, Canada, July 9-14, 2023, pages 4326–4343. Association for Computational Linguistics.
- 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, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.
- Markus Eberts and Adrian Ulges. 2020. [Span-based joint entity and relation extraction with transformer pre-training](#). In *ECAI 2020 - 24th European Conference on Artificial Intelligence, 29 August-8 September 2020, Santiago de Compostela, Spain, August 29 - September 8, 2020 - Including 10th Conference on Prestigious Applications of Artificial Intelligence (PAIS 2020)*, pages 2006–2013.
- Hao Fei, Shengqiong Wu, Jingye Li, Bobo Li, Fei Li, Libo Qin, Meishan Zhang, Min Zhang, and Tat-Seng Chua. 2022. [Lasuie: Unifying information extraction with latent adaptive structure-aware generative language model](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- Jenny Rose Finkel and Christopher D. Manning. 2009. [Nested named entity recognition](#). In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing, EMNLP 2009, 6-7 August 2009, Singapore, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 141–150. ACL.
- Ralph Grishman and Beth Sundheim. 1996. [Message understanding conference- 6: A brief history](#). In *16th International Conference on Computational Linguistics, Proceedings of the Conference, COLING 1996, Center for Sprogteknologi, Copenhagen, Denmark, August 5-9, 1996*, pages 466–471.
- Pankaj Gupta, Hinrich Schütze, and Bernt Andrassy. 2016. [Table filling multi-task recurrent neural network for joint entity and relation extraction](#). In *COLING 2016, 26th International Conference on Computational Linguistics, Proceedings of the Conference: Technical Papers, December 11-16, 2016, Osaka, Japan*, pages 2537–2547.
- Harsha Gurulingappa, Abdul Mateen Rajput, Angus Roberts, Julianne Fluck, Martin Hofmann-Apitius, and Luca Toldo. 2012. [Development of a benchmark corpus to support the automatic extraction of drug-related adverse effects from medical case reports](#). *J. Biomed. Informatics*, 45(5):885–892.
- Dan Hendrycks and Kevin Gimpel. 2016. [Bridging nonlinearities and stochastic regularizers with gaussian error linear units](#). *CoRR*, abs/1606.08415.
- Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. 2020. [Spanbert: Improving pre-training by representing and predicting spans](#). *Trans. Assoc. Comput. Linguistics*, 8:64–77.

| | | |
|-----|---|-----|
| 659 | Jin-Dong Kim, Tomoko Ohta, Yuka Tateisi, and Jun'ichi | 716 |
| 660 | Tsujii. 2003. GENIA corpus - a semantically anno- | 717 |
| 661 | tated corpus for bio-textmining . In <i>Proceedings of</i> | 718 |
| 662 | <i>the Eleventh International Conference on Intelligent</i> | 719 |
| 663 | <i>Systems for Molecular Biology, June 29 - July 3, 2003,</i> | 720 |
| 664 | <i>Brisbane, Australia</i> , pages 180–182. | 721 |
| 665 | Taku Kudo and John Richardson. 2018. Sentencepiece: | 722 |
| 666 | A simple and language independent subword tok- | 723 |
| 667 | enizer and detokenizer for neural text processing . In | |
| 668 | <i>Proceedings of the 2018 Conference on Empirical</i> | |
| 669 | <i>Methods in Natural Language Processing, EMNLP</i> | |
| 670 | <i>2018: System Demonstrations, Brussels, Belgium,</i> | |
| 671 | <i>October 31 - November 4, 2018</i> , pages 66–71. | |
| 672 | Fei Li, Meishan Zhang, Guohong Fu, and Donghong | |
| 673 | Ji. 2017. A neural joint model for entity and relation | |
| 674 | extraction from biomedical text . <i>BMC Bioinform.</i> , | |
| 675 | 18(1):198:1–198:11. | |
| 676 | Marius Lindauer, Katharina Eggersperger, Matthias | |
| 677 | Feurer, André Biedenkapp, Difan Deng, Carolin Ben- | |
| 678 | jamins, Tim Ruhkopf, René Sass, and Frank Hutter. | |
| 679 | 2022. SMAC3: A versatile bayesian optimization | |
| 680 | package for hyperparameter optimization . <i>J. Mach.</i> | |
| 681 | <i>Learn. Res.</i> , 23:54:1–54:9. | |
| 682 | Tianyu Liu, Yuchen Eleanor Jiang, Nicholas Monath, | |
| 683 | Ryan Cotterell, and Mrinmaya Sachan. 2022. Autore- | |
| 684 | gressive structured prediction with language models . | |
| 685 | In <i>Findings of the Association for Computational</i> | |
| 686 | <i>Linguistics: EMNLP 2022, Abu Dhabi, United Arab</i> | |
| 687 | <i>Emirates, December 7-11, 2022</i> , pages 993–1005. | |
| 688 | Association for Computational Linguistics. | |
| 689 | Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man- | |
| 690 | dar Joshi, Danqi Chen, Omer Levy, Mike Lewis, | |
| 691 | Luke Zettlemoyer, and Veselin Stoyanov. 2019. | |
| 692 | Roberta: A robustly optimized BERT pretraining | |
| 693 | approach . <i>CoRR</i> , abs/1907.11692. | |
| 694 | Yaojie Lu, Qing Liu, Dai Dai, Xinyan Xiao, Hongyu | |
| 695 | Lin, Xianpei Han, Le Sun, and Hua Wu. 2022. Uni- | |
| 696 | fied structure generation for universal information | |
| 697 | extraction . <i>CoRR</i> , abs/2203.12277. | |
| 698 | Giovanni Paolini, Ben Athiwaratkun, Jason Krone, | |
| 699 | Jie Ma, Alessandro Achille, Rishita Anubhai, | |
| 700 | Cícero Nogueira dos Santos, Bing Xiang, and Ste- | |
| 701 | fano Soatto. 2021. Structured prediction as transla- | |
| 702 | tion between augmented natural languages . In <i>9th</i> | |
| 703 | <i>International Conference on Learning Representa-</i> | |
| 704 | <i>tions, ICLR 2021, Virtual Event, Austria, May 3-7,</i> | |
| 705 | <i>2021</i> . OpenReview.net. | |
| 706 | Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, | |
| 707 | Jacob Devlin, James Bradbury, Anselm Levskaya, | |
| 708 | Jonathan Heek, Kefan Xiao, Shivani Agrawal, and | |
| 709 | Jeff Dean. 2022. Efficiently scaling transformer in- | |
| 710 | ference . <i>CoRR</i> , abs/2211.05102. | |
| 711 | Colin Raffel, Noam Shazeer, Adam Roberts, Katherine | |
| 712 | Lee, Sharan Narang, Michael Matena, Yanqi Zhou, | |
| 713 | Wei Li, and Peter J. Liu. 2020. Exploring the limits | |
| 714 | of transfer learning with a unified text-to-text trans- | |
| 715 | former . <i>J. Mach. Learn. Res.</i> , 21:140:1–140:67. | |
| | Sebastian Riedel, Limin Yao, and Andrew McCallum. | |
| | 2010. Modeling relations and their mentions with- | |
| | out labeled text . In <i>Machine Learning and Knowl-</i> | |
| | <i>edge Discovery in Databases, European Conference,</i> | |
| | <i>ECML PKDD 2010, Barcelona, Spain, September</i> | |
| | <i>20-24, 2010, Proceedings, Part III</i> , volume 6323 of | |
| | <i>Lecture Notes in Computer Science</i> , pages 148–163. | |
| | Springer. | |
| | Dan Roth and Wen-tau Yih. 2004. A linear program- | |
| | ming formulation for global inference in natural lan- | |
| | guage tasks . In <i>Proceedings of the Eighth Confer-</i> | |
| | <i>ence on Computational Natural Language Learning,</i> | |
| | <i>CoNLL 2004, Held in cooperation with HLT-NAACL</i> | |
| | <i>2004, Boston, Massachusetts, USA, May 6-7, 2004,</i> | |
| | pages 1–8. ACL. | |
| | Santosh Tokala Yaswanth Sri Sai, Prantika Chakraborty, | |
| | Sudakshina Dutta, Debarshi Kumar Sanyal, and | |
| | Partha Pratim Das. 2021. Joint entity and relation | |
| | extraction from scientific documents: Role of lin- | |
| | guistic information and entity types . In <i>Proceed-</i> | |
| | <i>ings of the 2nd Workshop on Extraction and Evalua-</i> | |
| | <i>tion of Knowledge Entities from Scientific Documents</i> | |
| | <i>(EEKE 2021) co-located with JCDL 2021, Virtual</i> | |
| | <i>Event, September 30th, 2021</i> , pages 15–19. | |
| | Erik F. Tjong Kim Sang and Fien De Meulder. 2003. | |
| | Introduction to the conll-2003 shared task: Language- | |
| | independent named entity recognition . In <i>Proceed-</i> | |
| | <i>ings of the Seventh Conference on Natural Language</i> | |
| | <i>Learning, CoNLL 2003, Held in cooperation with</i> | |
| | <i>HLT-NAACL 2003, Edmonton, Canada, May 31 -</i> | |
| | <i>June 1, 2003</i> , pages 142–147. ACL. | |
| | Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, | |
| | Weiming Lu, and Yueting Zhuang. 2023. Diffusion- | |
| | ner: Boundary diffusion for named entity recognition . | |
| | In <i>Proceedings of the 61st Annual Meeting of the As-</i> | |
| | <i>sociation for Computational Linguistics (Volume 1:</i> | |
| | <i>Long Papers)</i> , <i>ACL 2023, Toronto, Canada, July 9-14,</i> | |
| | <i>2023</i> , pages 3875–3890. | |
| | Dianbo Sui, Yubo Chen, Kang Liu, Jun Zhao, Xian- | |
| | grong Zeng, and Shengping Liu. 2020. Joint entity | |
| | and relation extraction with set prediction networks . | |
| | <i>CoRR</i> , abs/2011.01675. | |
| | Wei Tang, Benfeng Xu, Yuyue Zhao, Zhendong Mao, | |
| | Yifeng Liu, Yong Liao, and Haiyong Xie. 2022. Unirel: | |
| | Unified representation and interaction for | |
| | joint relational triple extraction . In <i>Proceedings of</i> | |
| | <i>the 2022 Conference on Empirical Methods in Natu-</i> | |
| | <i>ral Language Processing, EMNLP 2022, Abu Dhabi,</i> | |
| | <i>United Arab Emirates, December 7-11, 2022</i> , pages | |
| | 7087–7099. | |
| | Zaratiana Uurchade, Nadi Tomeh, Pierre Holat, and | |
| | Thierry Charnois. 2024. An autoregressive text-to- | |
| | graph framework for joint entity and relation extrac- | |
| | tion . | |
| | Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob | |
| | Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz | |
| | Kaiser, and Illia Polosukhin. 2017. Attention is all | |

you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.

Chenguang Wang, Xiao Liu, Zui Chen, Haoyun Hong, Jie Tang, and Dawn Song. 2022. [Deepstruct: Pre-training of language models for structure prediction](#). In *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 803–823.

Jue Wang and Wei Lu. 2020. [Two are better than one: Joint entity and relation extraction with table-sequence encoders](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 1706–1721.

Yijun Wang, Changzhi Sun, Yuanbin Wu, Hao Zhou, Lei Li, and Junchi Yan. 2021. [Unire: A unified label space for entity relation extraction](#). *CoRR*, abs/2107.04292.

Yunyang Xiong, Zhanpeng Zeng, Rudrasis Chakraborty, Mingxing Tan, Glenn Fung, Yin Li, and Vikas Singh. 2021. [Nyströmformer: A nyström-based algorithm for approximating self-attention](#). In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 14138–14148. AAAI Press.

Zhiheng Yan, Chong Zhang, Jinlan Fu, Qi Zhang, and Zhongyu Wei. 2021. [A partition filter network for joint entity and relation extraction](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 185–197.

Deming Ye, Yankai Lin, Peng Li, and Maosong Sun. 2022. [Packed levitated marker for entity and relation extraction](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 4904–4917.

Shubin Zhao and Ralph Grishman. 2005. [Extracting relations with integrated information using kernel methods](#). In *ACL 2005, 43rd Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference, 25-30 June 2005, University of Michigan, USA*, pages 419–426. The Association for Computer Linguistics.

Zexuan Zhong and Danqi Chen. 2021. [A frustratingly easy approach for entity and relation extraction](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 50–61.

A Appendix

Definitions. The LSE operation used in the following equations is defined as:

$$\text{LSE}_{n=1}^N(\mathbf{x}) = \log \sum_{n=1}^N \exp(\mathbf{x}_n) \quad (12)$$

where $\mathbf{x} \in \mathbb{R}^N$ be a vector of reals.

During training, the model will learn to minimize the following loss function:

$$L_{\text{ITER}} = \sum_{i=1}^N \sum_{j=1}^3 \left[\begin{matrix} L_{is_left}(i) \\ L_{is_span}(i) \\ L_{is_link}(i) \end{matrix} \right]_j, \quad (13)$$

a summation of three loss values for each position in the input sequence \mathbf{x} .

The loss function L_{is_left} is defined as:

$$L_{\text{left}}(n) = \text{LSE}_{j=1}^2 \left[\begin{matrix} \mathbf{h}_n \\ 0 \end{matrix} \right]_j - \text{LSE}_{j=1}^2(\Gamma) \quad (13)$$

$$\Gamma = \left[\begin{matrix} \mathbf{h}_n + (1 - \alpha) * -M \\ \alpha * -M \end{matrix} \right]$$

where $\mathbf{h}_n \in \mathbb{R}$ is the real-valued output of $FFN_{\kappa=1}^{\psi=1}(h_n)$, $M \rightarrow \infty$.

$$\alpha = \begin{cases} 1 & \text{iff. } \boxed{\text{I}} \in \mathcal{A}_n \\ 0 & \text{otherwise} \end{cases}$$

equals to one if the model should perform a $\boxed{\text{I}}$ action at time-step n , effectively cancelling out one of the terms in the above equation. Accordingly, we define L_{lr} :

$$L_{\text{lr}}(n) = \text{LSE}_{j=1}^2(\pi) - \text{LSE}_{j=1}^2(\Pi) \quad (14)$$

where

$$\pi = \left[\begin{matrix} (\text{LSE}_{i=1}^{\eta} \mathbf{h}_{n,m,i}) \\ 0 \end{matrix} \right]$$

$$\Pi = \left[\begin{matrix} (\text{LSE}_{i=1}^{\eta} \mathbf{h}_{n,m,i} + \Delta_{n,m,i}) + (1 - \beta) * (-\infty) \\ 0 + \beta * (-\infty) \end{matrix} \right]$$

$\eta = |\mathbf{T}_E|$ is the number of entity types and $\mathbf{h}_{n,m} = is_span(\mathbf{x}, n, m) \in \mathbb{R}^{\eta}$ is a vector containing one logit per such entity type.

$$\beta = \begin{cases} 1 & \text{iff. } \boxed{\text{J}} \in \mathcal{A}_n \\ 0 & \text{otherwise} \end{cases}$$

equals one iff. the performing $\boxed{\text{J}}$ is a correct action at time-step n . We also define $m = \max\{m \mid m \leq n \wedge \boxed{\text{I}} \in \mathcal{A}_m\}$, $m \leq n$, the largest index of the

| Dataset | Learning Rate | | -Schedule | | Warmup | | Weight Decay | | Batch -size | Activation -function |
|---------|---------------|--------|-------------------------|-------------------------|--------|------|--------------|--------|----------------|-------------------------|
| | T5 | ITER | T5 | ITER | T5 | ITER | T5 | ITER | | |
| CoNLL04 | 1e-4 | 4e-4 | linear with warmup | constant with warmup | 0.2 | 0.01 | 0.070 | 0.133 | 8 | GELU |
| ADE | 2e-4 | 2.8e-4 | constant with warmup | linear with warmup | 0.1 | 0.01 | 0.028 | 0.027 | 32 | ReLU |
| NYT | 2.5e-4 | 1e-4 | linear with warmup | linear | 0.2 | 0 | 0.016 | 0.07 | 8 | ReLU |
| GENIA | 2.6e-4 | 8e-4 | linear with warmup | linear with warmup | 0.2 | 0.1 | 0.045 | 0.056 | 16 | ReLU |
| CoNLL03 | 2e-4 | 9.7e-5 | constant | constant | 0 | 0 | 0.096 | 0.0098 | 8 | ReLU |

Table 7: Hyperparameter search results obtained using SMAC3 (Lindauer et al., 2022). For all datasets, the search was performed for 8 GPU hours using a single NVIDIA H100 GPU per dataset. The single best incumbent configuration has been selected for final training on the respective datasets.

preceding positions where $[\in \mathcal{A}_m$. Finally, we define

$$\Delta_{n,m,i} = \begin{cases} 0 & \text{iff. } (m, t_i) \in \mathcal{B}_n, t_i \in \mathbf{T}_E \\ -\infty & \text{otherwise} \end{cases}$$

to equal zero iff. there is a bracket pairing between the positions m and n of type $t_i \in \mathbf{T}_E$, and a large negative value otherwise. In order to minimize the loss function, the model is hereby incentivized to assign negative values to not existing interactions between two positions m and n of a certain type t_i . Lastly, L_{is_link} is defined as the binary cross entropy loss function:

$$L_{is_link}(n) = \sum_m \sum_{i=1}^{|\mathbf{T}_R|} \begin{cases} \mu & \text{iff. } [\in \mathcal{A}_{n,m} \\ 0 & \text{otherwise} \end{cases} \quad (15)$$

where

$$\mu = \sum \left[\frac{\theta_{n,m,i} * \log(\mathbf{h}_{n,m,i})}{(1 - \theta_{n,m,i}) * \log(1 - \mathbf{h}_{n,m,i})} \right]$$

with

$$\mathbf{h}_{n,m,i} = is_link(\mathbf{x}, n, m)$$

and $\theta_{n,m,i} = 1$ iff. the spans ending at positions n and m are in relationship i , $\theta_{n,m,i} = 0$ otherwise.

B Dataset Statistics

| Dataset | TRAIN | DEV | TEST | Nested Entities |
|---------|--------|-------|-------|-----------------|
| ADE | 4,272 | 10%* | 10%* | ✓ |
| NYT | 56,196 | 5,000 | 5,000 | ✓ |
| CoNLL03 | 954 | 216 | 231 | ✗ |
| CoNLL04 | 922 | 231 | 288 | ✗ |
| GENIA | 16,692 | † | 1,854 | ✓ |

Table 8: Number of samples per dataset split. * No official dataset split exists for ADE so we employ 10-fold cross-validation with 10% of the total examples following . † GENIA comes with only two files.

C Proofs

Theorem 1. Let $\mathbf{x} \in \mathcal{V}^N$ be a sequence of tokens with $x_N = \text{EOS}$. If $\mathbf{y} \in \mathcal{Y}_1 \times \dots \mathcal{Y}_M$ is the decoded sequence of actions, then $M \geq N$ holds for all $\mathbf{x} \in \mathcal{V}^N$.

Proof. Let a_m be the action chosen at step m , $\#_{\text{copy}}(m) = \sum_{i=1}^m \mathbb{1}_{[a_i = \text{copy}]}$ be the number of tokens x_n that have been copied until generation step m . Recall: generation completes at step m when $x_{\#_{\text{copy}}(m)} = \text{EOS} \wedge a_m = \text{copy}$ (1), i.e. the EOS token has been copied into the output.

Let $\#A(m) = m$ be the number of actions performed up until a certain point m in the output sequence \mathbf{y} of length M . It holds that

$$\#A(m) = \sum_{i=1}^m \mathbb{1}_{a_i = \text{copy}} + \sum_{i=1}^m \mathbb{1}_{a_i \neq \text{copy}} \cdot \begin{matrix} \geq 0 & \geq 0 \end{matrix}$$

With that, it follows that $\#_{\text{copy}}(m) \leq \#A(m)$ (2). Using (1) we get $\#_{\text{copy}}(M) = N$ and with (2) we then get $N \leq \#A(M) = M \implies N \leq M \Leftrightarrow M \geq N$ \square

D Hyperparameter search

Before training all of our models, we perform a hyperparameter search for all datasets using SMAC3 (Lindauer et al., 2022). For all datasets, we search for 8 hours, optimizing for high RE+ or NER F1, depending on the task. The search space consists of learning rates $lr \in [1e-3, 2e-5]$, learning rate schedules (*constant* or *linear*), warmup ratio $r \in \{0.0, 0.05, 0.1, 0.2\}$ and weight decay rate $wd \in [0, 0.1]$ for both the parameters of the base model (T5 in our case) and the parameters on top that are responsible for modeling the functions *is_left*, *is_span* and *is_link*, combined with batch size $bs \in \{8, 16, 32, 64\}$ and choice of activation

function $act \in \{\text{GELU}, \text{ReLU}, \text{tanh}\}$. The results of the hyperparameter search can be obtained in Table 7.

| Architecture | NER F1 (strict) | RE F1 (strict) |
|-----------------------------------|--------------------|-------------------|
| ITER + FLAN T5 _{large} † | 66.82 | 35.03 |
| (Ye et al., 2022) | 69.9 | 41.6 |
| (Sai et al., 2021) | 70.53 | 39.41 |
| SpERT | 67.62 | 46.44 ‡ |
| (Urchade et al., 2024) | 69.7 | 38.6 |
| (Wang et al., 2021) | 68.4 | 36.9 |

Table 9: Final training results for SciERC. † Preliminary tests. ‡ *BoundaryScore* reported.