

Dynamic Named Entity Recognition

Tristan Luiggi*†, Laure Soulier†*, Vincent Guigue ‡, Siwar Jendoubi*, Aurélien Baelde*

*Upskills R&D Choisy-le-roi, France

†Sorbonne Université, CNRS, ISIR, F-75005 Paris, France

‡AgroParisTech, UMR MIA-PS

* Université Paris-Saclay, CNRS, LISN, Orsay, France

firstname.lastname@{upmc.isir.fr;~upskills.ai;~agroparistech.fr}

ABSTRACT

Named Entity Recognition (NER) is a challenging and widely studied task that involves detecting and typing entities in text. So far, NER still approaches entity typing as a task of classification into universal classes (e.g. date, person, or location). Recent advances in natural language processing focus on architectures of increasing complexity that may lead to overfitting and memorization, and thus, underuse of context. Our work targets situations where the type of entities depends on the context and cannot be solved solely by memorization. We hence introduce a new task: Dynamic Named Entity Recognition (DNER), providing a framework to better evaluate the ability of algorithms to extract entities by exploiting the context. The DNER benchmark is based on two datasets, DNER-RotoWire and DNER-IMDb. We evaluate baseline models and present experiments reflecting issues and research axes related to this novel task.

CCS CONCEPTS

• **Computing methodologies** → **Artificial intelligence**; • **Natural language processing** → *Information extraction*;

KEYWORDS

Information extraction, NER, contextualization, datasets

ACM Reference Format:

Tristan Luiggi*†, Laure Soulier†*, Vincent Guigue ‡, Siwar Jendoubi*, Aurélien Baelde*. 2023. Dynamic Named Entity Recognition. In *The 38th ACM/SIGAPP Symposium on Applied Computing (SAC '23)*, March 27-March 31, 2023, Tallinn, Estonia. ACM, New York, NY, USA, Article 4, 8 pages. <https://doi.org/10.1145/3555776.3577603>

1 INTRODUCTION

Information extraction (IE) is a widely studied subfield of natural language processing (NLP) that aims to extract structured information from text. The task of Named Entity Recognition (NER) [16–18] is dedicated to recognize and classify entities (e.g. people, places, dates). NER usually constitutes the foundation of IE systems and strongly impacts the performance of higher-level tasks such as the extraction of relations between entities or the construction of knowledge graphs [19].

Publication rights licensed to ACM. ACM acknowledges that this contribution was authored or co-authored by an employee, contractor or affiliate of a national government. As such, the Government retains a nonexclusive, royalty-free right to publish or reproduce this article, or to allow others to do so, for Government purposes only.

SAC '23, March 27-April 2, 2023, Tallinn, Estonia

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-9517-5/23/03...\$15.00

<https://doi.org/10.1145/3555776.3577603>

Recent advances in artificial intelligence address NER as a sequence modeling task in which each word of an input text is associated with a predefined class. State-of-the-art solutions are formulated through a deep learning framework using specifically designed architectures such as biLSTM-CRF [18] and transformers [1]. These models are tested on reference benchmarks [38, 41].

Some variations of the NER task have been proposed to handle more difficult scenarios and solve specific issues: Entity Disambiguation [33] aims at segmenting entities and linking them to a knowledge base; Entity matching [5] aims at predicting if two detected entities refer to the same singularity.

These different tasks (including NER) share a common drawback: they all consider an entity as a universal concept, linked to a single class, even if it may appear in different surface forms and contexts. This limits the potential of the information extracted which could be useful for more elaborated downstream tasks. As an example, assuming a text relating facts of the *Amazon* company, our objective is to introduce a variability depending on the context, which might be in this case, for instance, its different roles of *seller/buyer*. Indeed, *Amazon* is likely to *sell* a product to an *individual person* but *buy* from another *company*. Moreover, typing an entity to a single class—in the vast majority of cases—makes memorization a competitive approach, as shown in [4, 37].

To address the contextualization of entity type and limit memorization effects, we design in this paper a specific task, called Dynamic Named Entity Recognition (DNER), in which entity labels are sample-dependent. It consists in detecting and categorizing entities whose type varies from one sample to another and addresses more abstract concepts (e.g. winning/losing teams in a basketball match). Related tasks encompass Entity-Aspect Linking (EAL) task [25] in which labels are constrained by fine-grained concepts in a knowledge base, Semantic Role Labeling (SRL) [20] in which labels correspond to roles played by tokens in their context (*verb, subject...*) or Event Extraction (EE) [45] in which a token describing an event (called a trigger) is extracted and linked to associated parameters (defined by the 5W1H *Who, What, Whom, When, Where and How*). Our scenario is different as the class is not specifically constrained by a knowledge base (as in EAL) or need not to be associated to explicit mentions in the input text (as in SRL or EE). The fact that labels are sample dependent ensures that a DNER system will be able to take better advantage of the sentence context and should therefore be more efficient on never seen entities and more transferable to new test data.

With this in mind, we propose two datasets, **DNER-RotoWire** and **DNER-IMDb**. The first one is derived from the academic dataset RotoWire [43], which contains data pairs with NBA match statistics and associated summaries. Our goal is to detect players and classify them as winners or losers based on the summary. The

second dataset is based on the IMDb website from which we extracted movie synopses and the associated ordered list of actors. The objective is to classify the actors according to their credit order (1,2,3,4...). For both datasets, we ensure that entities are classified differently across several samples to make the decision context-dependent. Our contribution is threefold:

- **DNER task formalization** (Section 3): we formalize the task of Dynamic Named Entity Recognition, including also the simplest task of Dynamic Named Entity Typing. We also detail the main associated challenges.

- **DNER Evaluation framework** (Sections 4 and 5): we present the built datasets, **DNER-RotoWire** and **DNER-IMDb**, and introduce a benchmark with metrics and a set of baselines.

- **Experiments** (Section 6): we conduct a series of preliminary experiments to evaluate the difficulty of the task. We outline insights reflecting the potential of the task in terms of model design. Our evaluation framework (datasets, metrics, baselines) is available at <https://github.com/Kawatami/DNER>.

2 RELATED WORK

Initial works in NER relied on hand-crafted features and rule-based algorithms [2, 14]. They mainly suffered from maintenance issues, lack of flexibility and thus high adaptation cost [7], leading the community to explore statistical approaches [40]. This marked a turning point in terms of performance, especially with the introduction of the IOB scheme (later extended to IOBES), which allowed the NER task to be treated as a sequence labeling problem [26]. Early proposals were divided into sequence modeling approaches (Hidden Markov Model) [6] and classical discriminators relying on rich contextual features [3, 32]. In this context, the CRF –conditional random field– [17], despite the cost of the Viterbi inference, received much attention [22, 23].

The growing interest in deep learning over the last 10 years [9] had led to significant advances. First deep-NER models [8, 13, 27, 28, 36] exploited the semantics introduced by word representations [9, 24]. Huang et al. [13] introduced the biLSTM-CRF model which quickly became a standard architecture with a powerful bidirectional recurrent neural network preceding a CRF layer to model label dependencies. This backbone has been successively improved by the incorporation of representations at the level of characters [18] and then tokens. Finally efforts have been made toward a better exploitation of additional supervision from external source of information [12, 34, 44].

These progresses have pushed the community to design more challenging NER tasks such as Entity-Aspect Linking, Event Extraction [25] or take over Semantic Role Labeling [20]. These tasks require to exploit the context either by establishing a link to a knowledge base or bound tokens between them with links having a special semantic. These scenarios do not encompass real-life situations where one can deduce entity roles not explicitly described in the input (the role is not part of a knowledge base nor can be represented as a link between two tokens following the predicate-argument structure) and might vary according to the context (such as *winner/loser* in basketball matches, *buyer/seller* in contracts, *etiologic/symptomatic* treatments in medical reports).

Recently, contextualized word representation models, such as BERT [11], have disrupted the NER task. This contextualization allows disambiguation of words while allowing a very efficient fine-tuning on most NLP tasks [15, 42], and leads to better performances in the case of NER [1]. In addition, the very fine modeling of the dependencies in the self-attention layers made it possible to dispense with the costly output CRF layer [37]. Those modern language models offered a lot of opportunities [31]: their extraction abilities improved transfer learning on NER [30], even in the more difficult setting in which the target domain was only associated with distant supervision [21].

Recent studies show that the complexity of those architectures enables them to encode information as a knowledge base [29]. Those capacities also raise new questions: regarding entities, what is the balance between memorization and extraction? From an even more general point of view, is memorization necessary -to integrate prior knowledge- or simply a phenomenon of over-fitting that must be limited? Historical datasets implicitly emphasize memorization capabilities by sharing an important part of the entity set between the training and the test set [4]. This phenomenon can be set aside either by designing a non-overlapping dataset [10] or by investigating transfer between datasets [37]. All these works around the notion of generalization in NER serve as a basis of reflection for this article. The current performances of language models push us to test more and more ambitious problems, this is the position of this article.

3 THE DNER TASK & CHALLENGES

Traditionally NER is formulated as sequence tagging task [13, 18]. Inspired by this formulation, we consider a supervised text T describing a single event (a basketball match or a film synopsis) in which entities (and by extension all their mentions) assume a single class within it. The text itself can be decomposed as a tuple $T = (X, Y)$:

- X corresponds to the raw textual data, split in N tokens: $X = \{x_1, \dots, x_i, \dots, x_N\}$, each token being drawn from a vocabulary \mathcal{X} .

- The entities are not nested, each corresponds to a consistent block of index $I = [i : i + j]$.

- \mathcal{V} which corresponds to the set of possible tags associated to entities. For our two proposed datasets, we define $\mathcal{V} = \{winner, loser\}$ and $\mathcal{V} = \{1, 2, 3, 4\}$, respectively. Note that these labels are not necessarily explicitly associated to any token in the input text X .

- Y stands for the set of IOBES labels associated with tokens mentions within T : $\mathcal{Y} = \{y_1, \dots, y_i, \dots, y_N\}$.

For the datasets proposed, we define $\mathcal{Y} = \{[B, I, E, S] - winner, [B, I, E, S] - loser, \emptyset\}$ or $\mathcal{Y} = \{[B, I, E, S] - 1, [B, I, E, S] - 2, [B, I, E, S] - 3, [B, I, E, S] - 4, \emptyset\}$. Thus, \mathcal{Y} is an extension of the label set \mathcal{V} dedicated to the label sequence tagging task.

3.1 Tasks

Based on these notations, we formalize two tasks to introduce two levels of difficulty for both datasets. We distinguish the DNET task from the DNER one, starting with the simplest. Please note that all metrics will be defined at the entity level.

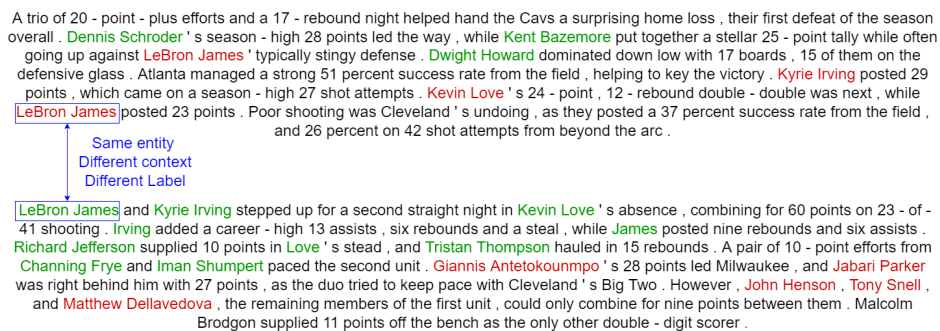


Figure 1: RotoWire preprocessed samples for the DNER task. Players highlighted in green are winners, and those in red are losers. Both samples mention the player "Lebron James" but with different labels.

Dynamic Named Entity Typing - DNET.

This task consists in classifying an already identified span of tokens indexed by I . Thus, we design a function f_{dnet} that makes a decision for a single entity mention within the given a context X , $f_{dnet} : \mathcal{X}^N \times I \rightarrow \mathcal{V}$.

Dynamic Named Entity Recognition - DNER.

This task corresponds to the complete NER, including both span identification and label assignment. The task is thus formalized as a sequence tagging: $f_{dner} : \mathcal{X}^N \rightarrow \mathcal{Y}^N$.

3.2 Challenges

Having the task formalization in mind, we can outline the different challenges of DNET and DNER:

- *Label variability*: an entity may have different labels depending on the sample, making the context influential and decreasing the informativeness of the entity's surface form. Typically, two basketball teams play each other several times, possibly with different results. The challenge is then to focus on the language elements designating the winners and losers but not on the team membership, which would lead to overfitting. It is worth noting that while NER may also assume such variability in principle it is not constrained by design and is empirically rarely found (e.g., 97.49% of entities in OntoNote are associated to a single label).

- *Label consistency*: an entity may be mentioned several times in a text X , possibly in slightly varying forms. It is important to maintain label consistency per entity during the inference process. For basketball matches, if a player is one of the winners, all his mentions must be labeled accordingly in the same sample. This challenge is also found for the NER task to a lesser extent: in the case of DNER it become crucial as the label variability takes a greater importance.

- *Out-of-scope entities and distant supervision*: the set Y of labels does not necessarily provide supervision for all of the named entities. For example, in basketball matches, we focus only on the players while other entities are likely to appear, in particular the coaches of both teams. This raises the question for the task and the metrics: do we need to detect this type of entity? If so, how do we label them? In this article, we focus primarily on the two aforementioned challenges and use common NER metrics. We leave the in-depth analysis of this challenge for future work.

4 DNER DATASETS

In this section, we describe the construction process and bias analysis of the two introduced datasets.

4.1 DNER-RotoWire

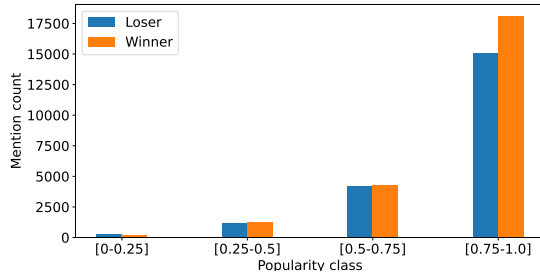
Construction and statistics. The Rotowire dataset [43] consists of pairs of tabular data (match statistics) and English summaries written by sport reporters. The primary goal of this dataset is to provide a benchmark for data-to-text generation models. To fit with the DNET and DNER tasks, we reprocessed the RotoWire dataset by identifying players as entities and denoting whether they belong to the winning or the losing team (a sample is provided in Figure 1). The procedure is done via regexp following the distant supervision paradigm which may introduce noise in the datasets, particularly, when names vary between tabular data and summaries; for instance, our script is designed to handle partial mentions (mentioning only the last name) but shows limitations when dealing with nicknames which would require an external knowledge base to be handled correctly.

Finally, to allow fair comparison between models, particularly Transformer-based approaches that are mostly limited to 512 tokens, we truncate summaries at this size limit. This implies that only 1.48 entities are removed on average per summary¹.

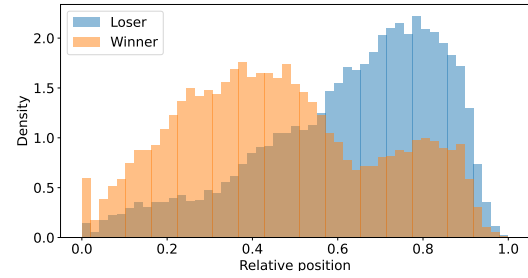
To measure the impact of the context memorization, we design a specific pipeline to split the dataset into train/test sets inspired by [37]. The goal is to separate test samples according to the increasing level of difficulty. Samples might share common properties (such as the teams involved in basketball matches) that a model could overfit and thus, bias its performance. To measure this phenomenon we divide the tests to measure performances in situations where the context (e.g., the team) have been seen and not seen during training. For this dataset we define the test set *Seen* as samples in which both involved teams are seen during training, *Unseen* when both teams are unseen during training, and *Seen/Unseen* for which only one team is seen during training. For more details about the splitting procedure, see section 5.

Dataset statistics are provided in Table 1, in the resulting sets we observe an imbalance in the class distribution towards winners (55% versus 45%), thus better performances are expected for this class.

¹Both the truncated and the full version of the datasets will be provided.



(a) Winning statistics w.r.t. popularity. The quartile of most popular players tends to win more frequently.



(b) Player relative position distribution according to their labels. Winning players tend to be mentioned earlier.

Figure 2: Analysis of popularity and relative position bias in the DNER-Rotowire dataset.

Set	Samples	Entities	Entity tokens	Winner	Loser
Train	1532	14202	24940	0.54	0.46
Validation	511	4615	8086	0.53	0.47
Seen (test)	511	4748	8360	0.54	0.46
Seen/Unseen (test)	1996	18293	31776	0.53	0.47
Unseen (test)	303	2721	4674	0.54	0.46

Table 1: DNER-RotoWire statistics. *Entities* refers to the number of entity mentions, *Entity tokens* to the number of tokens associated to entity mentions. Columns *Winner* and *loser* mention the proportion of each label category.

Source \ Target	Target				
	Train	Validation	Seen	Seen/Unseen	Unseen
Train	100	96.71	97.05	98.82	20.484
Validation	99.28	100	97.50	98.84	19.60
Seen (test)	99.22	96.80	100	98.79	21.89
Seen/Unseen (test)	70.72	67.97	66.68	100	59.43
Unseen (test)	40.94	37.11	35.51	99.52	100

Table 2: Proportion of common players between sets in DNER-RotoWire. From the source (rows) that appear in the target (column).

We checked the label variability of entities (main hypothesis of the proposed DNER task): we found that over 44579 total mentions, 44212 (99.17%) belonged to players with variable labels and 367 (0.82%) with constant ones. Complementary statistics about the sets are available in Table 2.

Bias analysis. We investigate here the potential bias behind the entities regarding our DNER task. Specifically, we consider two features:

- *Popularity:* some players are more popular than others, mainly due to their performances. This might impact the results of matches in which they play and lead to an over-citation ratio in summaries.
- *Position in the narrative of summaries:* It seems that sport journalists tend to present the facts/players of the winning team first, then those of the losing team.

In Figure 2 (a), we provide a visual representation of the popularity bias regarding players' labels. The popularity is estimated by the ratio of a player's mentions over the total mentions in the dataset. Quartiles are then extracted to group players within four groups ranging from low to high popularity. We can observe a relative

equal balancing between losing and winning mentions when players are not very popular (three first quartiles) and an unbalancing toward winning players for the most popular players.

Figure 2 (b) depicts the label distributions according to their relative position. Positions are normalized to range within $[0, 1]$ (beginning/end of the text), and grouped within 50 bins. Each distribution mode reflects the position bias: the winning players tend to be mentioned earlier on average.

These analyses show the importance for future models to leverage the textual context to deal with label variability and avoid any bias towards popularity and position.

4.2 DNER-IMDb

We provide the DNER-IMDb dataset with the two goals: 1) adding more variability in terms of vocabulary and 2) increasing the difficulty of the task with more output classes and uncertain labels.

Construction and statistics. IMDb²³ is an online database related to media content. We focus on movies, characterized by two main pieces of information that we use for the DNER task:

- *Movie synopses:* these are short English descriptions of movies. Synopses mention fictional characters, their importance and behavior in the films, and the relationships between the characters. A movie may have several synopses.
- *Character meta-data:* character information, including the actors playing those roles and their credit rank.

To fit with our objective of dynamic labeling, we replace fictional characters with actor names. Indeed, an actor will appear in several movies, probably with different labels. For instance, *Bruce Willis* played as first credited character in "Die Hard" but fourth in "Pulp Fiction" (see Figure 3). Once the substitution is done, the designed task is to identify the credit order of all actors given a synopsis.

In practice, we built the DNER-IMDb dataset with a raw database of 245,404 movie synopses. The labels range from the first credited actor to the eighth, but less than 1% of the samples have more than 4 credited actors. We only retain films with no more than 4 credited actors, produced between 1970 and 2021, and with a minimum of 600 views on the IMDb site. This ensures that synopses are written in a modern style with sufficient metadata about the

²<https://www.imdb.com/interfaces/>

³<https://rapidapi.com/apidojo/api/imdb8> (as of 27/01/2022)

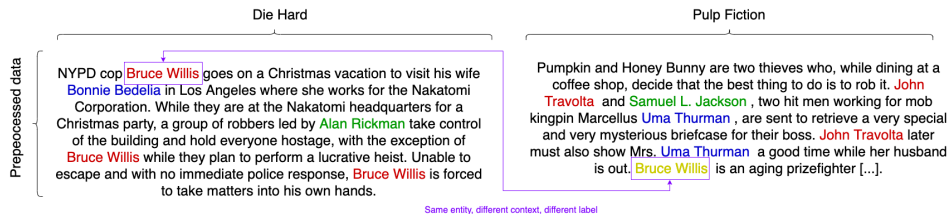
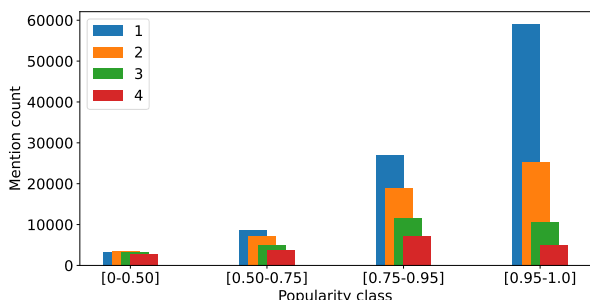


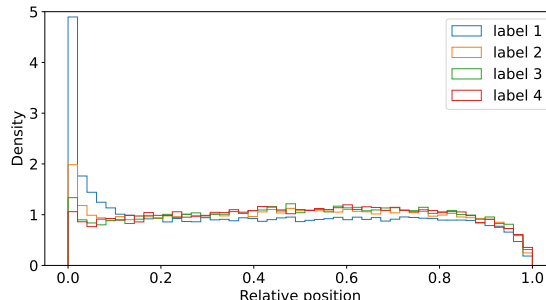
Figure 3: IMDb preprocessed samples for the DNER task. In movies *Die Hard* (left) and *Pulp Fiction* (right), actors are colored w.r.t. their labels. Red, blue, green, and yellow labels are resp. for first, second, third and fourth actors.

Set	# samples	# Entity	# Tokens	credit 1	credit 2	credit 3	credit 4
Train	13328	90726	185850	39.50%	24.76%	21.97%	13.75%
Validation	1725	12008	24557	37.40%	24.14%	23.55%	14.89%
Seen (test)	805	5450	11176	39.68%	23.92%	23.43%	12.95%
Unseen (test)	4030	27346	56038	39.38%	24.11%	22.60%	13.89%

Table 3: Sample, entity and token counts for DNER-IMDb. Entity label distribution for every set.



(a) Credit order statistics w.r.t. popularity. Popular actors tend to take first and second place in the credit while being less expected to be in third and fourth.



(b) Actor relative position distribution according to their labels. First credited actors tend to be cited more often earlier in the summary.

Figure 4: Analysis of popularity and relative position bias in the DNER-IMDb dataset.

actors. Moreover, actor names explicitly mentioned in the synopsis are removed (e.g. "F.B.I. trainee Clarice Starling (**Jodie Foster**) works hard [...] - *The Silence of the Lambs*). Then, we replace the names of the characters with the names of the actors using regular expressions. There are still a few improperly formatted samples; the main errors being (1) the mismatch between the characters' surface forms provided by the IMDb database and that found in the synopsis, resulting in inconsistent or partial permutations, and (2) mentions of characters without associated data. Similarly as in DNER-RotoWire, we restricted the synopsis length to 512 tokens. We observed that synopses are very short on average therefore the size limit of 512 tokens has almost no impact here.

The construction of the training and test sets follows a similar procedure for the DNER-RotoWire dataset (see Section 5). As samples are only described by unique movie titles on the synopsis level (and not two teams for a match summary in the DNER-RotoWire dataset), the procedure is simplified as no *seen/unseen* set is produced.

The resulting dataset consists of 44,189 samples (a sample is provided in Figure 3) with synopses averaging 106.89 words in length and 4.59 actor mentions. Dataset statistics are given in Table 3 and complementary statistics are available in Table 4. To assess the label variability hypothesis, we estimate the distribution of actors'

Source \ Target	Target			
	Train	Validation	Seen	Unseen
Train	100	19.13	20.76	20.76
Validation	95.76	100	30.57	43.83
Seen (test)	97.47	52.37	100	50.34
Unseen (test)	54.84	23.12	15.50	100

Table 4: Proportion of common actors between sets in DNER-IMDb. From the source (rows) that appear in the target (column).

mentions w.r.t. their number of different associated labels in the ground truth. Although there is an important number of mentions related to actors with the same label (23.9%), we measure that most mentions are associated with entities with multiple labels (18.7% have 2 labels, 26% have 3 labels, and 31.4% have 4 labels).

Bias analysis. This dataset shares common similarities with DNER-RotoWire regarding biases. We hypothesize that an actor also has a popularity factor and a position in synopses that might influence the decision process as well. Many movie synopses start by mentioning the first character (e.g. "**Mr. Cobb** a unique con artist can enter anyone's dreams [...] - *Inception*) which constitutes a potential bias. Actors' relative positions have been analyzed in Figure 4 (a). We notice that first credited actors are usually mentioned earlier, specifically at the very beginning of synopses. This effect

occurs for all classes, with less impact on other credit orders. Except for the beginning of synopses, we observe a homogeneous distribution of the position across labels. Similarly to DNER-RotoWire, we analyzed the popularity bias in Figure 4 (b). The popularity of each actor is estimated by the ratio of her/his mentions over the number of mentions in the ground truth. We can observe that the number of mentions for classes 1 and 2 increases with the popularity level while it remains stable or decreases for classes 3 and 4. This means that popular actors are more likely to be assigned to first or second credited roles than to be assigned to last credited ones. However, when actors are not really popular, they can be uniformly assigned to any roles in the credit order.

As for the DNER-RotoWire, this bias analysis reinforces our intuition that models need to focus on language elements to avoid learning by heart bias and not being robust to label variability.

5 TEST SET CONSTRUCTION PROCEDURE

To measure the impact of context memorization, we design a specific pipeline for splitting the datasets into train/test sets. We make the hypothesis that current architecture such as LSTM or transformer are complex enough to retain entities/labels association via overfitting and not solely on context analysis. To measure this phenomenon we divide the tests to measure performances in situations where the context is known and unknown. The first set hold samples with context seen during training, we expect the best performances as this set share most of its entities with the training set (see Table 2 and 4) which facilitates the segmentation process, in addition, a model may overfit on datasets biases previously mentioned in section 4.1 and 4.2. The second set holds data without any common context with the training data, we expect a decrease in performances as a model cannot overfit on the context in this case in addition to be exposed to never seen entities. We, therefore, build sets to dispose of seen and unseen contexts in the test set regarding those which compose the train set. The splitting pipeline is illustrated in Figure 5 and includes the following main steps:

- The set of context (teams and movie titles) is split into TrainV0 and TestV0u (75%/25%). TestV0u is considered as the set including "unseen context".
- The TrainV0 is split into two new sets TrainV1 and TestV0s (75%/25%). TestV0s represents the set of "seen context" while TrainV1 corresponds to the training contexts.
- A final step is necessary to aggregate these sets at the sample level. To do so, we process the list of samples and assign to the training fold if the current context belong to the set TrainV1 and to the test fold otherwise. From the TrainV1, we split into final train/validation sets on the basis of 80%/20%. For context in the test set, we distinguish two situations : "unseen" if the context belong to the set TestV0u and "seen" if the context belong to TestV0s. Therefore, the "seen" test set is guaranteed to hold sample sharing context with ones found in the train set.

Please note that depending on the data at hand a context might be defined by one property such as unique movie title in the case of DNER-IMDb. But several properties might take place in for other data such as DNER-RotoWire.

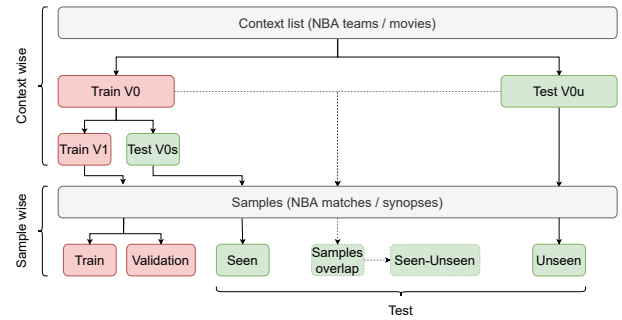


Figure 5: Data splitting procedure for set creation.

6 EXPERIMENT PROTOCOL

Baselines. We propose a set of baselines inspired by state-of-the-art NER architectures, which rely on Transformers supplemented either by a classical discriminator or a CRF layer for difficult cases [35].

For the simpler DNET task, we design a lightweight architecture that first encodes the texts at the token level. Since the entities are composed of a variable number of tokens, an aggregation problem must be solved before the classification stage. In their work, [39] study several approaches to compute such representation ranging from token selection, pooling, or attention. Their experiments highlight that the best is task-dependent, but the max-pooling procedure is very robust across a wide range of tasks. Thus, we use it to compute the representations of spans associated to the entities.

For both DNET and DNER, we also consider a supplementary feature modeling the general context. By exploiting the advantages of the BERT architecture, we integrate the special *CLS* token into the classifier features. We believe that this additional information can potentially be useful for our task-specific challenges such as label consistency or business knowledge modeling.

As a result, we consider 4 baselines for DNER and 2 for DNET:

- **BERT-Linear:** token representations are contextualized through BERT and then classified using a linear layer.
- **BERT-CLS:** as BERT-Linear with the additional *CLS* token concatenated to each word representation before classification.
- **BERT-CRF:** following the SOTA in transfer NER, we propose to add a CRF output layer to explicitly model label dependencies. This architecture is dedicated to the sequence labeling task and therefore not suited for the simple CNET task.
- **BERT-CLS-CRF:** as BERT-CRF with the additional *CLS* token. A visual representation of baseline models is provided in Figure 6.

Metrics. To measure the quality of the entity classification, we make use of $\mu F1$. In the case of the DNER task, the entities are extracted via the IOBES scheme in the case where it is well formatted (BEGIN token followed by INSIDE and finally END or just a SINGLE token), the entity labels are extracted via the associated class from \mathcal{V} .

To measure our ability to detect entities (whatever their classes), we provide a span quality metric, **Entity**. For each entity belonging to the associated reference summary, we compare each entity boundaries in the reference summary with boundaries from the *begin* and *end* tokens. If boundaries match, the entity has been correctly detected. We then estimate the *F1* score.

Models	Set	RotoWire			IMDb		
		DNET	DNER	Entity	DNET	DNER	Entity
BERT-Linear	Seen	0.81	0.66	0.86	0.67	0.36	0.58
	Seen/Unseen	0.81	0.65	0.85	-	-	-
	Unseen	0.80	0.63	0.81	0.45	0.31	0.56
BERT-CLS	Seen	0.81	0.67	0.88	0.69	0.37	0.60
	Seen/Unseen	0.81	0.68	0.87	-	-	-
	Unseen	0.80	0.67	0.85	0.46	0.32	0.58
BERT-CRF	Seen	-	0.67	0.90	-	0.60	0.94
	Seen/Unseen	-	0.67	0.88	-	-	-
	Unseen	-	0.66	0.87	-	0.52	0.92
BERT-CLS-CRF	Seen	-	0.61	0.82	-	0.56	0.90
	Seen/Unseen	-	0.61	0.81	-	-	-
	Unseen	-	0.60	0.79	-	0.48	0.88

Table 5: Experiment results. The $\mu F1$ score is reported for both datasets and tasks.

Model	Set	RotoWire				IMDb					
		GT	All	W	L	GT	All	1	2	3	4
Bert-Linear	S	5.44%	17.69%	14.61%	21.87%	0%	7.24%	4.47%	6.27%	11.11%	9.92%
	S/U	4.64%	20.88%	17.52%	26.15%	-	-	-	-	-	-
	U	2.63%	20.54%	17.96%	26.31%	0.31%	8.64%	5.59%	10.05%	12.02%	9.20%
Bert-CLS	S	5.44%	13.27%	9.23%	18.75%	0%	5.01%	3.68%	3.76%	7.20%	7.14%
	S/U	4.64%	18.27%	13.40%	25.92%	-	-	-	-	-	-
	U	2.63%	10.81%	7.81%	17.54%	0.31%	5.68%	4.58%	6.41%	6.43%	6.25%

Table 6: Inconsistency analysis statistics. Entity with a single mention are ignored. S stands for the *seen* test set, S/U for the *seen/unseen* test set and U for the *unseen*

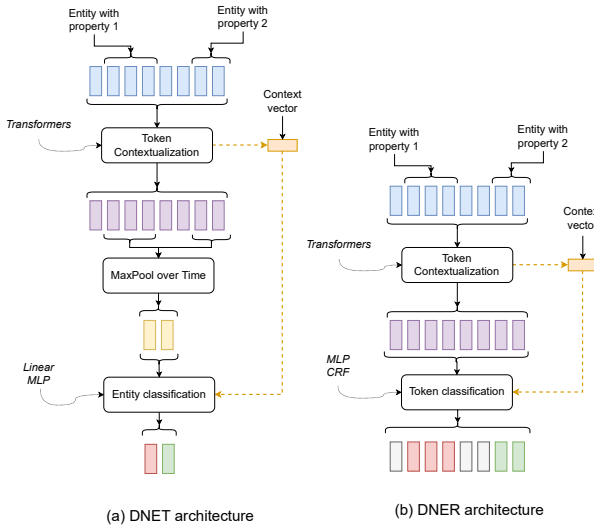


Figure 6: Baseline architectures

To measure the entity consistency (challenges in Section 3.2), we design the **inconsistency metric**. It compares the labels of all mentions of the same entity within a sample. If an entity obtains the same label for all its mentions, the incoherence metric is equal to 0. Otherwise, its value is 1. This metric is then aggregated over all multi-mentioned entities⁴ of all samples.

7 BENCHMARK RESULTS

RotoWire - DNET. In this experiment, the goal is to classify player contextualized representations within two categories: *winner* and *loser*. Results are shown in Table 5 (left).

We logically observe a decrease in performance when the difficulty increases from *Seen* to *Unseen*. Even if the difference is limited,

⁴For clarity, all entities that are mentioned only once are removed from the calculations.

it is easier for a model to decide if some properties of the context have already been seen during training. Our two baselines perform similarly. We can observe in Table 6 (left) the label inconsistency metric. The effect of remote supervision is visible on the ground truth with an inconsistency that varies from 2.63% to 5.44%. This error is amplified by the model whose inconsistencies rise to 20.54% (*unseen* test set) for BERT-Linear. This indicates that the consistency challenge is difficult to meet without explicit modeling of team membership constraints. It is interesting to note that the introduction of a general context (CLS) enables the model to significantly reduce inconsistencies (10.81% for the *unseen* test set).

RotoWire - DNER. All DNER results are shown in table 5 (right). The first conclusion from this table is that the CLS token provides a performance gain compared to the baseline *Bert-Linear* architecture. This is consistent with the experiments with *DNET*, where the CLS token exhibited better robustness to inconsistency. This effect could be of greater magnitude due to the IOBES scheme, which requires the classification of a larger number of classes. *BERT-CRF* performs better in entity recognition, which is easily explained by the CRF layer effectively maintaining label coherence. This suggests that the CRF helps in maintaining such a factor, but is not able to correctly analyze a context. The combination of the CLS token and the CRF layer generally performs well, but the added complexity could trigger an overfitting effect. This confirms our suspicions about the added difficulty of our proposed task.

The best baseline (*BERT-CLS*) proposes an average F1 performance of 0.67, mainly due to errors in typing and precision issues; although interesting, it is clear that the many challenges mentioned in Section 3.2 must be addressed in a specific way to cope with the difficulty of the DNER task.

IMDb - DNET. On these data, the scores are globally worse than on RotoWire (Table 5 (right)). This is easily explained by a change from 2 to 4 classes. The loss of performance by going from seen

to unseen is much more marked than on RotoWire with more than 22 points of $\mu F1$ on average: the memorization effect brings information for all the entities with a single class.

In terms of labeling inconsistencies (Table 6 (right)), the problem is virtually absent from the ground truth. On the test data, the inconsistencies rise to 8.64%. While this figure again shows the need to specifically address this issue, it is still far below that of RotoWire. This difference is easily explained: the number of entities per document is much lower on IMDb (4 against 9 on average on unseen) and this intrinsically reduces the risk of inconsistencies.

IMDb - DNER. Entity detection is better on IMDb than on RotoWire, but it relies heavily on the CRF layer. The overall improvement is probably due to the artificial aspect of the dataset, where entities always have the same surface forms and to the overlap rate between learning and testing (54.84%, even on unseen movies).

We then return to the conclusion of the previous section: once detected, entities are hard to categorize. The difference between seen and unseen films is very large (6 $\mu F1$ points on average) and the overall performance tops out at 0.60 of $\mu F1$.

8 CONCLUSION

This paper introduces the Dynamic Named Entity Recognition (DNER) task which aims at detecting entities and classifying them in a frame where labels are dynamic. This task raises several challenges such as label variability, label consistency and taking into account entity position or popularity bias. We provide benchmarks in the form of two supervised datasets associated with test sets of increasing difficulty. These benchmarks are provided with metrics and reference models to ensure reproducibility and to encourage the emergence of new models to address the specific challenges of the task. Indeed, despite a reference architecture based on transformers, our analyses show that the DNER task is particularly difficult and the results obtained can be improved. The presented datasets were designed for experimental purposes and might not be relevant for real world applications.

REFERENCES

- [1] Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. Contextual String Embeddings for Sequence Labeling. In *COLING*.
- [2] Douglas E. Appelt, J.R. Hobbs, John Bear, David Israel, Megumi Kameyama, Andy Kehler, David Martin, Karen Myers, and Mabry Tyson. 1995. SRI International FASTUS System MUC-6 Test Results and Analysis. In *MUC-6*.
- [3] Masayuki Asahara and Yuji Matsumoto. 2003. Japanese Named Entity Extraction with Redundant Morphological Analysis. In *HLT-NAACL*.
- [4] I Augenstein, L Derczynski, and K Bontcheva. 2017. Generalisation in named entity recognition: A quantitative analysis. *Computer Speech & Language* (2017).
- [5] Nils Barlaug and Jon Atle Gulla. 2020. Neural Networks for Entity Matching. *CoRR abs/2010.11075* (2020).
- [6] Daniel M. Bikel, Scott Miller, Richard Schwartz, and Ralph Weischedel. 1997. Nymble: a High-Performance Learning Name-finder. In *ACL*.
- [7] E Brill. 1995. Transformation-based error-driven learning and natural language processing: A case study in part-of-speech tagging. *Computational linguistics* (1995).
- [8] Jason P. C. Chiu and Eric Nichols. 2015. Named Entity Recognition with Bidirectional LSTM-CNNs. *ACL* (2015).
- [9] R Collobert, J Weston, L Bottou, M Karlen, K Kavukcuoglu, and P Kuksa. 2011. Natural Language Processing (Almost) from Scratch. *J. Mach. Learn. Res.* (2011). <http://dl.acm.org/citation.cfm?id=2078183.2078186>
- [10] L Derczynski, E Nichols, M van Erp, and N Limsopatham. 2017. Results of the WNUT 2017 Shared Task on Novel and Emerging Entity Recognition. In *ACL*.
- [11] J Devlin, M.W Chang, K Lee, and K Toutanova. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *CoRR* (2018).
- [12] A Ghaddar, P Langlais, A Rashid, and M Rezagholizadeh. 2021. Context-aware Adversarial Training for Name Regularity Bias in Named Entity Recognition. *CoRR* (2021). arXiv:2107.11610
- [13] Zhiheng Huang, Wei Xu, and Kai Yu. 2015. Bidirectional LSTM-CRF Models for Sequence Tagging. *PACLIC* (2015).
- [14] S B Huffman. 1995. Learning information extraction patterns from examples. In *IJCAI*.
- [15] Mandar Joshi, Omer Levy, Daniel S Weld, and Luke Zettlemoyer. 2019. BERT for coreference resolution: Baselines and analysis. *ACL* (2019).
- [16] Onur Kuru, Ozan Arkan Can, and Deniz Yuret. 2016. CharNER: Character-Level Named Entity Recognition. In *COLING*.
- [17] J Lafferty, A McCallum, and F Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. (2001).
- [18] G Lample, M Ballesteros, S Subramanian, K Kawakami, and Cs Dyer. 2016. Neural Architectures for Named Entity Recognition. In *ACL*.
- [19] Jing Li, Aixun Sun, Jianglei Han, and Chenliang Li. 2018. A Survey on Deep Learning for Named Entity Recognition. *IEEE* (2018).
- [20] Zuchao Li, Hai Zhao, Shexia He, and Jiaxun Cai. 2020. Syntax Role for Neural Semantic Role Labeling. *ACM* (2020).
- [21] Chen Liang, Yue Yu, Haoming Jiang, Siawpeng Er, Ruijia Wang, Tuo Zhao, and Chao Zhang. 2020. Bond: Bert-assisted open-domain named entity recognition with distant supervision. In *SIGKDD*.
- [22] Xiaohua Liu, Shaodian Zhang, Furu Wei, and Ming Zhou. 2011. Recognizing named entities in tweets. In *ACL*.
- [23] A McCallum and W Li. 2003. Early results for Named Entity Recognition with Conditional Random Fields, Feature Induction and Web-Enhanced Lexicons. In *HLT-NAACL*.
- [24] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient Estimation of Word Representations in Vector Space. *ICLR* (2013).
- [25] Federico Nanni, Simone Paolo Ponzetto, and Laura Dietz. 2018. Entity-Aspect Linking: Providing Fine-Grained Semantics of Entities in Context. *ACM/IEEE*.
- [26] Lance Ramshaw and Mitch Marcus. 1995. Text Chunking using Transformation-Based Learning. In *ACL*.
- [27] Marek Rei. 2017. Semi-supervised Multitask Learning for Sequence Labeling. *CoRR abs/1704.07156* (2017).
- [28] Marek Rei, Gamal K. O. Crichton, and Sampo Pyysalo. 2016. Attending to Characters in Neural Sequence Labeling Models. *ACL* (2016).
- [29] Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How much knowledge can you pack into the parameters of a language model? *CoRR* (2020).
- [30] J.D Rodriguez, A Caldwell, and A Liu. 2018. Transfer learning for entity recognition of novel classes. In *COLING*.
- [31] Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in bertology: What we know about how bert works. *ACL* (2020).
- [32] Satoshi Sekine, Ralph Grishman, and Hiroyuki Shinnou. 1998. A Decision Tree Method for Finding and Classifying Names in Japanese Texts. In *ACL*.
- [33] Ö Sevgili, A Shelmanov, M Y. Arkipov, A Panchenko, and C Biemann. 2020. Neural Entity Linking: A Survey of Models based on Deep Learning. *SWJ* (2020).
- [34] Shreyam Shah and Jyoti Ramteke. 2021. Context aware Named Entity Recognition with Pooling. *ICCICT 2021* (2021).
- [35] Fábio Souza, Rodrigo Nogueira, and Roberto Lotufo. 2019. Portuguese named entity recognition using BERT-CRF. *CoRR* (2019).
- [36] E Strubell, P Verga, D Belanger, and A McCallum. 2017. Fast and Accurate Sequence Labeling with Iterated Dilated Convolutions. *EMNLP* (2017).
- [37] B Taillé, V Guigie, and P Gallinari. 2020. Contextualized embeddings in named-entity recognition: An empirical study on generalization. In *ECIR*.
- [38] Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the CoNLL-2003 Shared Task: Language-Independent Named Entity Recognition. In *NAACL*.
- [39] Shubham Toshniwal, Haoyue Shi, Bowen Shi, Lingyu Gao, Karen Livescu, and Kevin Gimpel. 2020. A Cross-Task Analysis of Text Span Representations. *ACL*.
- [40] Kristina Toutanova, Dan Klein, Christopher D Manning, and Yoram Singer. 2003. Feature-rich part-of-speech tagging with a cyclic dependency network. In *ACL*.
- [41] Ralph M. Weischedel, Eduard H. Hovy, Mitchell P. Marcus, and Martha Palmer. 2017. OntoNotes : A Large Training Corpus for Enhanced Processing.
- [42] G Wiedemann, S Remus, A Chawla, and C Biemann. 2019. Does BERT make any sense? Interpretable word sense disambiguation with contextualized embeddings. *KONVENS* (2019).
- [43] Sam Wiseman, Stuart Shieber, and Alexander Rush. 2017. Challenges in Data-to-Document Generation. In *EMNLP*.
- [44] Chuhan Wu, Fangzhao Wu, Tao Qi, and Yongfeng Huang. 2020. Named Entity Recognition with Context-Aware Dictionary Knowledge. <https://aclanthology.org/2020.ccl-1.85>
- [45] Wei Xiang and Bang Wang. 2019. A Survey of Event Extraction From Text. *IEEE Access* (2019).