

AMRize, then Parse! Enhancing AMR Parsing with PseudoAMR Data

Anonymous ACL submission

Abstract

As Abstract Meaning Representation (AMR) implicitly involves compound semantic annotations, we hypothesize auxiliary tasks which are semantically or formally related can better enhance AMR parsing. With carefully designed control experiments, we find that 1) Semantic role labeling (SRL) and dependency parsing (DP), would bring much more significant performance gain than unrelated tasks in the text-to-AMR transition. 2) To make a better fit for AMR, data from auxiliary tasks should be properly “AMRized” to PseudoAMR before training. 3) Intermediate-task training paradigm outperforms multitask learning when introducing auxiliary tasks to AMR parsing. From an empirical perspective, we propose a principled method to choose, reform, and train auxiliary tasks to boost AMR parsing. Extensive experiments show that our method achieves new state-of-the-art performance on in-distribution, out-of-distribution benchmarks of AMR parsing. We will release our code upon acceptance.

1 Introduction

Abstract Meaning Representation (AMR) (Banasescu et al., 2013) parsing aims to translate a sentence to a directed acyclic graph, which represents the relations among abstract concepts as shown in Figure 1. AMR can be applied to many downstream tasks, such as information extraction (Rao et al., 2017; Wang et al., 2017; Zhang and Ji, 2021), text summarization, (Liao et al., 2018; Hardy and Vlachos, 2018) question answering (Mitra and Baral, 2016; Sachan and Xing, 2016) and dialogue modeling (Bonial et al., 2020).

Recently, AMR Parsing with the sequence-to-sequence framework achieves most promising results (Xu et al., 2020; Bevilacqua et al., 2021). Comparing with transition-based or graph-based methods, sequence-to-sequence models do not require tedious data processing and is naturally compatible with auxiliary tasks (Xu et al., 2020)

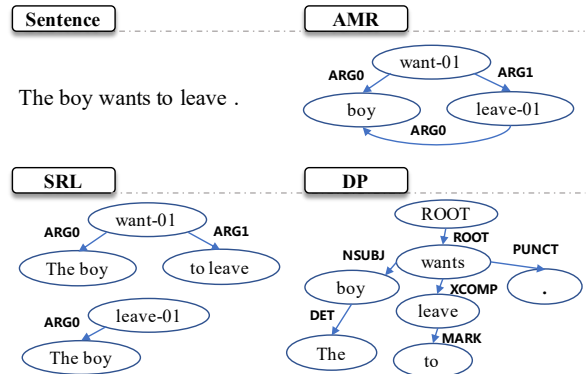


Figure 1: The Abstract Meaning Representation (AMR), Semantic Role Labeling (SRL), and Dependency Parsing (DP) structure of the sentence “The boy wants to leave.”

and powerful pretrained encoder-decoder models (Bevilacqua et al., 2021). Previous work (Xu et al., 2020; Wu et al., 2021) has shown that the performance of AMR parser can be effectively boosted through co-training with certain auxiliary tasks, e.g. Machine Translation or Dependency Parsing.

However, when introducing auxiliary tasks to enhance AMR parsing, we argue that three important issues still remain under-explored in the previous work. **1) How to choose auxiliary task?** The task selection is important since loosely related tasks may even impede the AMR parsing according to Damonte and Monti (2021). However, in literature there are no principles or consensus on how to choose the proper auxiliary tasks for AMR parsing. Though previous work achieves noticeable performance gain through multi-task learning, they do not provide explainable insights on why certain task outperforms others or in which aspects the auxiliary tasks benefit the AMR parser. **2) How to bridge the gap between tasks ?** The form and semantic gaps between AMR parsing and auxiliary tasks are non-negligible. For example, Machine Translation generates text sequence while Dependency Parsing (DP) and Semantic Role Labeling

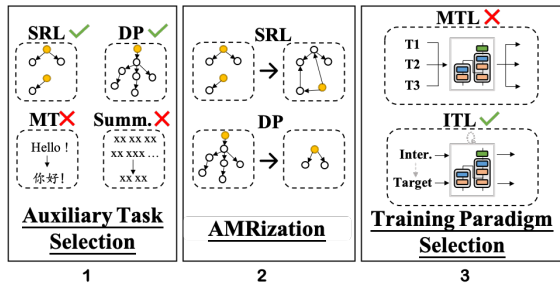


Figure 2: Illustration of methodology in this paper. We proposed a principled method to select, transform and train the auxiliary tasks.

(SRL) produces dependency trees and semantic role forests, respectively. The structural differences between DP, SRL and AMR are visualized in Figure 1. Many prior studies (Xu et al., 2020; Wu et al., 2021; Damonte and Monti, 2021) do not attach particular importance to the gap, which might lead the auxiliary tasks to be what is called *outlier-task* (Zhang and Yang, 2021; Cai et al., 2017) in the Multitask Learning, deteriorating the performance of AMR parsing. **3) How to introduce auxiliary tasks more effectively?** After investigating different training paradigms to combine the auxiliary task training with the major objective (AMR parsing), we figure out that, although all baseline models (Xu et al., 2020; Wu et al., 2021; Damonte and Monti, 2021) choose to jointly train the auxiliary tasks and AMR parsing with Multi-task Learning (MTL), Intermediate-task Learning (ITL) is a more effective way to introduce the auxiliary tasks. Our observation is also consistent with (Pruksachatkun et al., 2020; Poth et al., 2021), which improve other NLP tasks with enhanced pretrained models.

In response to the above three issues, we summarize a principled method to select, transform and train the auxiliary tasks (Figure 2) to enhance AMR parsing from extensive experiments. **1) Auxiliary Task Selection.** We choose auxiliary tasks by estimating their similarities with AMR from the semantics and formality perspectives. AMR is recognized as a deep semantic parsing task which encompasses multiple semantic annotations, e.g. semantic roles, name entities and co-references. As a direct semantic-level sub-task of AMR parsing, we select SRL as one auxiliary task. Traditionally, formal semantics views syntactic parsing a precursor to semantic parsing, leading to the mapping between syntactic and semantic relations. Hence we introduce dependency parsing, a syntactic pars-

ing task as another auxiliary task. **2) AMRization.** Despite being highly related, the output formats of SRL, DP and AMR are distinct from each other. To this end, we introduce transformation rules to “AMRize” SRL and DP to PseudoAMR, intimating the feature of AMR. Specifically, through *Reentrancy Restoration* we transform the structure of SRL to a graph and restore the reentrancy within arguments, which mimics AMR structure. Through *Redundant Relation Removal* we conduct transformation in dependency trees and remove relations that are far from semantic relations in AMR graph. **3) Training Paradigm Selection.** We find that ITL makes a better fit for AMR parsing than MTL since it allows model progressively transit to the target task instead of learning all tasks simultaneously, which benefits knowledge transfer (Zhang and Yang, 2021).

We summarize our contributions as follows:

1. Semantically or formally related tasks, e.g., SRL and DP, are better auxiliary tasks for AMR parsing compared with distantly related tasks, e.g. machine translation and machine reading comprehension.
2. We propose task-specific rules to AMRize the structured data to PseudoAMR. SRL and DP with properly transformed output format further improve AMR parsing.
3. ITL outperforms classic MTL methods when introducing auxiliary tasks to AMR Parsing. We show that ITL derives a steadier and better converging process during training.

Extensive experiments show that our method (PseudoAMR + ITL) achieves the new state-of-the-art of single model on in-distribution (85.1 Smatch score on AMR 2.0, 83.9 on AMR 3.0), out-of-distribution benchmarks. Specifically we observe that AMR parser gains larger improvement on the SRL(+3.3), Reentrancy(+3.1) and NER(+2.0) metrics¹, due to higher resemblance with the selected auxiliary tasks.

2 Methodology

As shown in Figure 2, in this paper, we propose a principled method to select auxiliary tasks (Section 2.1), AMRize them into PseudoAMR (Section 2.2) and train PseudoAMR and AMR effectively

¹Computed on AMR 2.0 and 3.0 dataset.

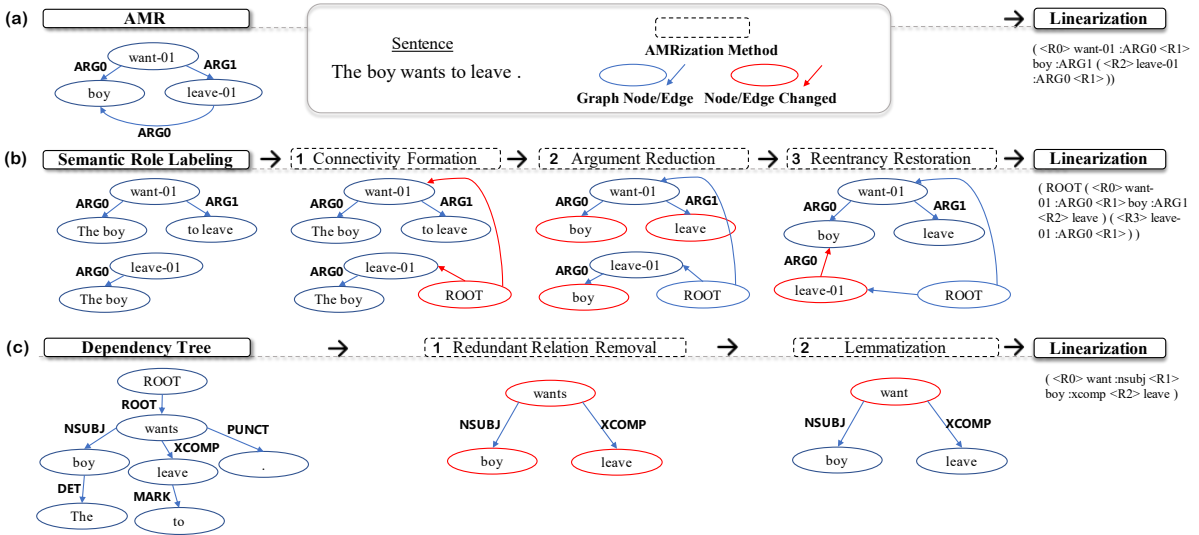


Figure 3: Illustration of AMRization methods and Graph Linearization. The source sentence is “The boy wants to leave.”

(Section 2.3) to boost AMR parsing. We formulate both PseudoAMR and AMR parsing as the sequence-to-sequence generation problem. Given a sentence $x = [x_i]_{1 \leq i \leq N}$, the model aims to generate a linearized PseudoAMR or AMR graph $y = [y_i]_{1 \leq i \leq M}$ (the right part of Figure 3) with a product of conditional probability:

$$P(y) = \prod_{i=1}^M p(y_i | (y_1, y_2, \dots, y_{i-1}))$$

2.1 Auxiliary Task Selection

When introducing auxiliary tasks for AMR parsing, the selected tasks should be formally or semantically related to AMR, thus the knowledge contained in them can be transferred to AMR parsing. Based on this principle of relevance, we choose semantic role labeling (SRL) and dependency parsing (DP) as our auxiliary tasks. We involve Translation and Summarization tasks for comparison.

Semantic Role Labeling SRL aims to recover the predicate-argument structure of a sentence, which can enhance AMR parsing, because: (1) Recovering the predicate-argument structure is also a sub-task of AMR parsing. As illustrated in Figure 3(a,b), both AMR and SRL locate the predicates (‘want’, ‘leave’) of the sentence and conduct word-sense disambiguation. Then they both capture the multiple arguments of center predicate. (2) SRL and AMR are known as shallow and deep semantic parsing, respectively. It is reasonable to think that the shallow level of semantic knowledge in SRL is useful for deep semantic parsing.

Dependency Parsing DP aims to parse a sentence into a tree structure, which represents the dependency relation among tokens. The knowledge of DP is useful for AMR parsing, since: (1) Linguistically, DP (syntax parsing task) can be the precursor task of AMR (semantic parsing). (2) The dependency relation of DP is also related to semantic relation of AMR, e.g., as illustrated in Figure 1(c), ‘NSUBJ’ in DP usually represents ‘:ARG0’ in AMR. Actually, they both correspond to the agent-patient relations in the sentence. (3) DP is similar to AMR parsing from the perspective of edge prediction, because both of them need to capture the relation of nodes (tokens/concepts) in the sentence.

2.2 AMRization

Although SRL and DP are highly related to AMR parsing, there still exists gaps between them, e.g., SRL annotations may be disconnected, while AMR is always a connected graph. To bridge these gaps, we transform them into PseudoAMR, which we call AMRization.

2.2.1 Transform SRL to PseudoAMR

We summarize typical gaps between SRL and AMR as: (1) *Connectivity*. AMR is a connected directed graph while the structure of SRL is a forest. (2) *Span-Concept Gap*. Nodes in AMR graph represent concepts (e.g., “boy”) while that of SRL are token spans (e.g., “the boy”, “that boy”). Actually all the mentioned token spans correspond to the same concept. (3) *Reentrancy*. Reentrancy is an important feature of AMR as shown in Figure 3(a),

the instance boy is referenced twice as ARG0. The feature can be applied to conduct coreference resolution. However, there is no reentrancy in SRL. To bridge such gaps, we propose **Connectivity Formation, Argument Reduction** and **Reentrancy Restoration** to transform SRL into PseudoAMR.

Connectivity Formation To address the connectivity gap, we need to merge all SRL trees into a connective graph. Note that the merging doesn't guarantee correctness in semantic level. As shown in Figure 3(b-1), we first add a virtual root node, then generating a directed edge from the virtual root to each root of SRL trees, thus the SRL annotation becomes a connected graph.

Argument Reduction To address the Span-Concept Gap, as shown in Figure 3(b-2), if the argument of current predicate is a span with more than one token, we will replace this span with its head token in its dependency structure. Thus token spans "the boy", "that boy" will be transformed to "boy", more similar to the corresponding concept. Similar method has been applied by (Zhang et al., 2021) to find the head of token spans of argument.

Reentrancy Restoration For the reentrancy gap, we design a heuristic algorithm based on DFS to restore reentrancy in SRL. As shown in Figure 3(b-3), the core idea of the restoration is that we create a variable when the algorithm first sees a node. If the DFS procedure meets node with the same name, the destination of current edge will be redirected to the variable we have created at first. Please refer to Appendix A for the pseudo code of the reentrancy restoration. Note that the algorithm can not guarantee the merging of nodes is unbiased since there are might be different variables with same name in real AMR and SRL itself doesn't provide information to rebuild all correct reentrancies. However, this simple transformation leads to significant improvements in terms of SRL and Reentrancy fine-grained scores of AMR parser as shown in Table 7, revealing the transform can enhance knowledge transfer between SRL and AMR Parsing.

2.2.2 Transform Dependency Structure to PseudoAMR

We summarize the gaps between Dependency Tree and AMR as: (1) *Redundant Relation*. Some relations in dependency parsing focus on syntax, e.g., ":PUNCT" and ":DET", which are far from semantic relations in AMR. (2) *Token-Concept Gap*. The

basic element of dependency structure is token while that of AMR is the concept, which captures deeper syntax-independent semantics. We use **Redundant Relation Removal** and **Token Lemmatization** to transform the dependency structure to PseudoAMR to handle the gaps.

Redundant Relation Removal For the Redundant Relation Gap, we remove some relations which are far from the sentence's semantics most of the time, such as "PUNCT" and "DET". As illustrated in Figure 3(c-1), by removing some relations of the dependence, the parsing result become more compact compared with original DP tree, forcing the model to ignore some semantics-unrelated tokens during seq2seq training.

Token Lemmatization As shown in Figure 3(c-2), for Token-Concept Gap, we conduct lemmatization on the node of dependency tree based on the observation that the affixes of single word do not affect the concept it corresponds to. Together with the smart-initialization (Bevilacqua et al., 2021) by setting the concept token's embedding as the average of the subword constituents, the embedding vector of lemmatized token ('want') becomes closer to the vector concept ('want-01') in the embedding matrix, therefore requiring the model to capture deeper semantic when conducting DP task.

2.2.3 Linearization

After all AMRization steps, the graph structure of SRL/DP also should be linearized before doing seq2seq training. As depicted in the right part of Figure 3, we linearize the graph by the DFS-based travel, and use special tokens <R0>, ..., <Rk> to indicate variables, and parentheses to mark the depth, which is the best AMR linearization method of Bevilacqua et al. (2021).

2.3 Training Paradigm Selection

After task selection and AMRization, we still need to choose an appropriate training paradigm to train PseudoAMR and AMR effectively. We explore three training paradigms as follows:

Multitask training Following Xu et al. (2020); Damonte and Monti (2021), we use classic schema in sequence-to-sequence multitask training by adding special task tag at the beginning of input sentence and training all tasks simultaneously. The validation of best model is conducted only on the AMR parsing sub-task.

Intermediate training Similar to Pruksachatkun et al. (2020), we first fine-tune the pretrained model on the intermediate task (PseudoAMR parsing), followed by fine-tuning on the target AMR parsing task under same training setting.

Multitask & Intermediate training We apply a joint paradigm to further explore how different paradigms affect AMR parsing. We first conduct multitask training, followed by fine-tuning on AMR parsing. Under this circumstance, Multitask training plays the role as the intermediate task.

3 Experiments

3.1 Datasets

AMR Datasets We conducted out experiment on two AMR benchmark datasets, AMR 2.0 and AMR 3.0. AMR2.0 contains 36521, 1368 and 1371 sentence-AMR pairs in training, validation and testing sets, respectively. AMR 3.0 has 55635, 1722 and 1898 sentence-AMR pairs for training validation and testing set, respectively. We also conducted experiments in out-of-distribution datasets (BIO,TLP,News3) and low-resources setting.

Auxiliary Task Datasets Apart from DP/SRL, we choose NLG tasks including summarization and translation to evaluate the contributions of auxiliary tasks. Description of datasets is listed Appendix C.

3.2 Evaluation Metrics

We use the Smatch scores (Cai and Knight, 2013) and further the break down scores (Damonte et al., 2017) to evaluate the performance.

To fully understand the aspects where auxiliary tasks improve AMR parsing, we divide the fine-grained scores to two categories: **1) Concept-Related** including Concept, NER and Negation scores, which care more about concept centered prediction. **2) Topology-Related** including Unlabeled, Reentrancy and SRL scores, which focus on edge and relation prediction. NoWSD and Wikification are listed as isolated scores because NoWSD is highly correlated with Smatch score and wikification relies on external entity linker system.

3.3 Experiment Setups

Model Setting We use current state-of-the-art Seq2Seq AMR Paring model SPRING (Bevilacqua et al., 2021) as our main baseline model and apply BART-Large (Lewis et al., 2020) as our pretrained model. Blink (Li et al., 2020) is used to add wiki

tags to the predicted AMR graphs. We do not apply re-category methods and other post-processing methods are the same with Bevilacqua et al. (2021) to restore AMR from token sequence. The hyperparameters tuning details is listed in Appendix D

AMRization Setting For SRL, we explore four AMRization settings. 1) Trivial. Concept :multi-sentence and relation :snt are used to represent the virtual root and its edges to each of the SRL trees. 2) With Argument Reduction. We use dependency parser from Stanford CoreNLP Toolkit (Manning et al., 2014) to do the argument reduction. 3) With Reentrancy Restoration 4) All techniques.

For DP, we apply four AMRization settings 1) Trivial. Extra relations in dependency tree are added to the vocabulary of BART 2) With Lemmatization. We use NLTK (Bird, 2006) to conduct token lemmatization 3) With Redundant Relation Removal. We remove PUNCT, DET, MARK and ROOT relations. 4) All techniques.

3.4 Main Results

We report the result (ITL + All AMRization Techniques) on benchmark AMR 2.0 and 3.0 in Table 1. On AMR 2.0, our models with DP or SRL as intermediate task gains consistent improvement over the SPRING model by a large margin (1.1 Smatch) and reach new state-of-the-art for single model (85.1 Smatch). Compared with SPRING with 200k extra data, our models achieve higher performance with much less extra data (40k v.s. 200k), suggesting the effectiveness of our intermediate tasks. We also compare our models with contemporary work (Lam et al., 2021; Zhou et al., 2021b). It turns out that our ensemble model beats its counterpart with less extra data, reaching a higher performance (85.3 Smatch). In fact, even without ensembling, our model still performs better than those ensembling models, showing the effectiveness of our methods.

On AMR 3.0, Our models consistently outperform other models under both single model (83.9 Smatch) and ensembling setting (84.0 Smatch). Same as AMR 2.0, our single model reaches higher Smatch score than those ensembling models, revealing the effectiveness of our proposed methods.

Fine-grained Performance To better analyse how the AMR parser benefits from the intermediate training and how different intermediate tasks affect the overall performance. We report the fine-grained score as shown in Table 1. We can tell

Model	Extra Data	SMATCH	NoWSD	Wiki	Concept-related			Topology-related			
					Conc.	NER	Neg.	Unl.	Reen.	SRL	
AMR 2.0	Cai and Lam (2020)	N	78.7	79.2	81.3	88.1	87.1	66.1	81.5	63.8	74.5
	Fernandez Astudillo et al. (2020)	N	80.2	80.7	78.8	88.1	87.5	64.5	84.2	70.3	78.2
	Zhou et al. (2021a)	70k	81.8	82.3	78.8	88.7	88.5	69.7	85.5	71.1	80.8
	SPRING (Bevilacqua et al., 2021)	N	83.8	84.4	84.3	90.2	90.6	74.4	86.1	70.8	79.6
	SPRING (w/ silver) (Bevilacqua et al., 2021)	200k	84.3	84.8	83.1	90.8	90.5	73.6	86.7	72.4	80.5
	SPRING (Ours)	N	84.0	84.3	83.5	89.9	91.8	75.1	87.1	71.3	81.3
	Ours (w/ DP)	40k	85.0	85.4	84.1	90.4	92.5	74.7	88.2	74.7	83.1
	Ours (w/ SRL)	40k	85.1	85.6	83.6	90.4	91.4	75.7	88.2	75.0	83.5
	*Graphene 4S ^E (Lam et al., 2021)	200k	84.8	85.3	83.9	90.6	92.2	75.2	88.0	71.4	83.5
	*Structure-aware ^E (Zhou et al., 2021b)	47k	84.9	-	-	-	-	-	-	-	-
Ours (w/ SRL) ^E	40k	85.3	85.7	83.9	90.7	92.2	75.0	88.4	75.0	83.6	
AMR 3.0	Bevilacqua et al. (2021) (w/ silver)	200k	83.0	83.5	82.7	89.8	87.2	73.0	85.4	70.4	78.9
	Ours (w/ DP)	40k	83.9	84.3	81.6	89.7	89.2	73.0	87.0	73.7	82.3
	Ours (w/ SRL)	40k	83.9	84.3	81.0	89.7	88.4	73.9	87.0	73.9	82.5
	*Graphene 4S ^E (Lam et al., 2021)	200k	83.8	84.2	81.9	90.1	88.3	74.6	86.9	70.2	82.5
	*Structure-aware ^E (Zhou et al., 2021b)	47k	83.1	-	-	-	-	-	-	-	-
	Ours (w/ SRL) ^E	40k	84.0	84.5	80.7	90.0	88.9	73.1	87.1	73.9	82.6

Table 1: SMATCH and fine-grained F1 scores on AMR 2.0 and 3.0. ^E denotes result with model ensemble (the details of the ensembling models are described in Appendix B). We conduct ensembling by averaging the parameters of models from three random seeds following Zhou et al. (2021b). Model with * denotes contemporary work.

that by incorporating intermediate tasks, considerable increases on most sub-metrics, especially on the Topology-related terms, are observed. On both AMR 2.0 and 3.0 our single model with SRL as intermediate task achieves the highest score in Unlabeled, Reentrancy and SRL metrics, suggesting that SRL intermediate task improves our parser’s capability in Coreference and SRL.

DP leads to consistent improvement in topology-related metrics, which also derives the best result on NER sub-task (92.5 on AMR 2.0, 89.2 on AMR 3.0). We suppose that the ":nn" relation which signifies multi-word name entities in dependency parsing helps the AMR parser recognize multi-word name entities. Generally speaking, AMR parser gains large improvement in Topology-related sub-tasks and NER by incorporating our intermediate tasks in terms of the Smatch scores.

3.5 Exploration in Auxiliary Task Selection

We explore how different tasks affect AMR parsing apart from DP and SRL. We involve two classic conditional NLG tasks, Summarization and Translation for comparison as shown in Table 2.

The comparison implies that SRL and DP are better auxiliary tasks for AMR Parsing even under the circumstance where their counterparts exploit far more data (40k v.s. 400k). In fact, the performance of MT drops while introducing more data, which contradicts with Xu et al. (2020)’s findings that more MT data can lead to better result when pretraining the *raw Transformer model*. However, this is not surprising under the back-

Model	Extra	SMATCH	Conc.	Topo.
Ours (w/ NLG)				
- w/ DialogSum	13k	84.5	85.5	81.5
- w/ CNNDM	40k	84.4	85.5	81.7
- w/ CNNDM	80k	84.2	85.1	81.4
- w/ EN-DE	40k	84.4	85.3	81.5
- w/ EN-DE	80k	84.4	85.4	81.4
- w/ EN-DE	200k	84.2	84.6	81.2
- w/ EN-DE	400k	83.6	84.9	80.6
Ours (w/ Parsing)				
- w/ DP	40k	85.0	85.9	82.0
- w/ SRL	40k	85.1	85.8	82.2

Table 2: Result of Task Selection. We first train BART on different auxiliary tasks for 10 epochs before AMR Parsing. We also report the average scores of Concept-related (Conc.) and Topology-related metrics (Topo.)

ground of Intermediate-task Learning where we already have a pretrained model with large-scale pretraining. Whether the intermediate tasks’ form fits for the target task is far more important than the amount of data in the intermediate-task as also revealed by Poth et al. (2021). According to their observation, tasks with the most data (QQP 363k, MNLI 392k) perform far worse (-97.4% relative performance degradation at most) on some target tasks compared with tasks having much smaller datasets (CommonsenseQA 9k, SciTail 23k) which on the contrary give a positive influence.

In conclusion, our findings suggest that the selection of intermediate task is important and should be closely related to AMR parsing in form, otherwise it would even lead to a performance drop for AMR parsing.

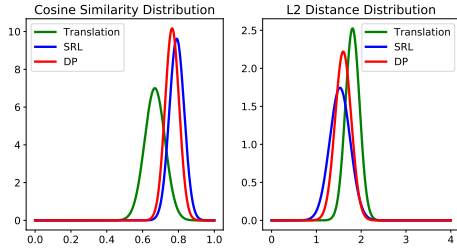


Figure 4: The distance distribution of sentences representation. SRL and DP consistently provide more similar sentence representation to AMR than Translation. The computation is illustrated in Figure 6 in appendix.

4 Analysis

4.1 More Similar Sentence Representation

To examine how different auxiliary tasks affect AMR parsing, we collect the sentences’ representation from different tasks’ trained encoders. We use the average hidden state of the encoder’s output as the sentence representation. We compute the Cosine Similarity and L2 distance between auxiliary tasks’ representation and AMR’s representation for same sentence. The test split of AMR 2.0 is used for evaluation. Finally, We apply Gaussian distribution to fit the distribution of distances and draw the probability distribution function curves as shown in Figure 4. It turns out that under both distance metrics, SRL/DP consistently provide more similar sentence representation to AMR than Translation and SRL/DP are more similar to AMR parsing. It empirically justifies our hypothesis that semantically or formally related tasks can lead to a better initialization for AMR parsing.

4.2 Ablation Study on AMRization Methods

As shown in Table 3, we conduct ablation study on how different AMRization methods affect the performance AMR parsing. For both SRL and DP, jointly adopting our AMRization techniques can further improve the performance of AMR parsing significantly, comparing to trivial linearization. The imperfect reentrancy restoration method leads to a significant improvement in terms of both the Topology and Concept related scores. It reveals that transformation of structure to mimic the feature of AMR can better the knowledge transfer between shallow and deep semantics.

As shown in Table 7, compared with jointly using the two techniques, it is worth noting that model with solely Reentrancy Restoration reaches highest fine-grained scores in especially on Reentrancy

Model	SMATCH	Conc.	Topo.
Ours (w/ Semantic Role Labeling)	84.5	85.5	81.6
- w/ Arg. Reduction(AR)	84.8	85.6	81.9
- w/ Reen. Restoration(RR)	85.0	86.1	82.5
- w/ AR+RR	85.1	85.8	82.2
Ours (w/ Dependency Parsing)	84.4	84.7	81.7
- w/ Redundant Relation Removal (RRR)	84.5	85.2	81.8
- w/ Lemmatization (Lemma)	84.7	85.5	81.7
- w/ RRR + Lemma	85.0	85.9	82.0

Table 3: We report the average scores of Concept-related scores and Topology-related scores. The full scores are listed in Table 7. The improvement of involving all techniques against trivial linearization is significant with $p < 0.005$ for both SRL and DP.

Model	Extra	SMATCH
Ours (w/ Intermediate)		
- w/ DP	40k	85.0
- w/ SRL	40k	85.1
- w/ DP,SRL	80k	84.7
Ours (w/ Multitask)		
- w/ DP	40k	83.7
- w/ SRL	40k	83.6
- w/ DP,SRL	80k	83.5
Ours (w/ Multi. + Inter.)		
- w/ DP	40k	84.1
- w/ SRL	40k	84.1
- w/ DP,SRL	80k	83.9

Table 4: Analysis on Training Paradigms. Intermediate-task training is more suitable for AMR parsing than Multitask training

and SRL scores. To explore the reason why it outperform adopting both techniques, we analyse the number of restored reentrancy. The result shows that about 10k more reentrancies are added when Argument Reduction (AR) is previously executed. It’s expected since AR replaces the token span to the root token. Compared with token span, single token is more likely to be recognized as the coreference variable according to the Reentrancy Restoration (RR) algorithm, thus generating more reentrancy, which might include bias to the model. This explains why solely using RR can lead to better results on SRL and Reen.

4.3 ITL Outweighs MTL

We report the result of different fine-tuning paradigms in Table 4. It justifies our assumption that classic multitask learning with task tag as previously applied in Xu et al. (2020); Damonte and Monti (2021) does not compare with intermediate training paradigm for AMR Parsing task.

As shown in Figure 5, Intermediate-task training provides a faster and better converging process

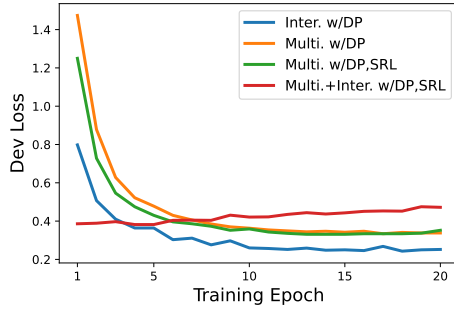


Figure 5: The loss curve on development set of AMR 2.0 for different training paradigms.

Model	BIO	TLP	News3
SPRING	59.7	77.3	73.7
SPRING + silver	59.5	77.5	71.8
SPRING ^E	60.5	77.9	74.7
Ours	61.2	78.9	75.4

Table 5: Analysis on OOD data. ^E denotes result given by the ensembling of models. Our model exploits SRL as the intermediate task.

than MTL. We assume this is due to the huge gap between AMR parsing and auxiliary tasks which may harm the optimization process of MTL. The process of optimizing all auxiliary tasks simultaneously may introduce noise to AMR Parsing.

We also find that under the setting of ITL, sequentially training SRL and DP tasks did not bring further improvement to AMR parsing. We guess this is due to the catastrophic forgetting problem. Further regularization during training might help the model progressively learn from different auxiliary tasks and relieve catastrophic forgetting.

4.4 Exploration in Out-of-Distribution Generalization

Following Bevilacqua et al. (2021); Lam et al. (2021), we assess the performance of our models when trained on out-of-distribution (OOD) data. The models trained solely on AMR 2.0 training data are used to evaluate out-of-distribution performance on the BIO, the TLP and the News3 dataset.

Table 6 shows the result of our out-of-distribution experiments. Our model surpass other models even the ensembled one(Lam et al., 2021), creating new state-of-the-art for single model.

5 Related Work

AMR Parsing AMR parsing is a challenging semantic parsing task, since AMR is a deep semantic representation and consists of many sep-

arate annotations (Banarescu et al., 2013) (e.g., semantic relations, named entities, co-reference and so on). There are four major methods to do AMR Parsing currently, sequence-to-sequence approaches (Ge et al., 2019; Xu et al., 2020; Bevilacqua et al., 2021), tree-based approaches (Zhang et al., 2019b,a), graph-based approaches (Lyu and Titov, 2018; Cai and Lam, 2020) and transition-based approaches (Naseem et al., 2019; Lee et al., 2020; Zhou et al., 2021a).

There are two ways to incorporate other tasks to AMR Parsing. Goodman et al. (2016) builds AMR graph directly from dependency trees while (Ge et al., 2019) parse directly from linearized syntactic tree. Xu et al. (2020) introduces Machine Translation, Constituency Parsing as pretraining tasks for Seq2Seq AMR parsing and Wu et al. (2021) introduces Dependency Parsing for transition-based AMR parsing. However all of them do not take care of the semantic and formal gap between the auxiliary tasks and AMR parsing.

Multitask & Intermediate-task Learning

Multi-task Learning (MTL) (Caruana, 1997) aims to jointly train multiple related tasks to improve the performance of all tasks. Different from MTL, Intermediate-task Learning (ITL) is proposed to enhance pretrained models e.g. BERT by training on intermediate task before fine-tuning on the target task. Recent studies(Pruksachatkun et al., 2020; Poth et al., 2021) on ITL expose that choosing right intermediate tasks is important. Tasks that don't match might even bring negative effect to the target even if it has far more data.

Xu et al. (2020); Damonte and Monti (2021); Procopio et al. (2021) utilize auxiliary tasks in a MTL fashion with specific task tags. Bevilacqua et al. (2021); Zhou et al. (2021b) adopt sliver training data in a ITL paradigm. However, there is no work comparing ITL and MTL when introducing auxiliary tasks to enhance PTM-based AMR parser.

6 Conclusion

In this paper, We find that semantically or formally related tasks, e.g. SRL and DP are better auxiliary tasks for AMR parsing and can further improve the performance by proper AMRization methods to bridge the gap between tasks. And Intermediate-task Learning is more effective in introducing auxiliary tasks compared with Multitask Learning. Extensive experiments and analyses show the effectiveness and priority of our proposed methods.

References

- 591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
- Laura Banarescu, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider. 2013. Abstract meaning representation for sembanking. In *Proceedings of the 7th linguistic annotation workshop and interoperability with discourse*, pages 178–186.
- Michele Bevilacqua, Rexhina Blloshmi, and Roberto Navigli. 2021. One spring to rule them both: Symmetric amr semantic parsing and generation without a complex pipeline. In *Proceedings of the Thirty-Fifth AAAI Conference on Artificial Intelligence*.
- Steven Bird. 2006. [NLTK: The Natural Language Toolkit](#). In *Proceedings of the COLING/ACL 2006 Interactive Presentation Sessions*, pages 69–72, Sydney, Australia. Association for Computational Linguistics.
- Claire Bonial, L. Donatelli, Mitchell Abrams, Stephanie M. Lukin, Stephen Tratz, Matthew Marge, Ron Artstein, David R. Traum, and Clare R. Voss. 2020. Dialogue-amr: Abstract meaning representation for dialogue. In *LREC*.
- Deng Cai and Wai Lam. 2020. [AMR parsing via graph-sequence iterative inference](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1290–1301, Online. Association for Computational Linguistics.
- Shu Cai and Kevin Knight. 2013. [Smatch: an evaluation metric for semantic feature structures](#). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 748–752, Sofia, Bulgaria. Association for Computational Linguistics.
- Sirui Cai, Yuchun Fang, and Zhengyan Ma. 2017. Will outlier tasks deteriorate multitask deep learning? In *Neural Information Processing*, pages 246–255. Springer International Publishing.
- Rich Caruana. 1997. Multitask learning. *Machine learning*, 28(1):41–75.
- Yulong Chen, Yang Liu, Liang Chen, and Yue Zhang. 2021. [DialogSum: A real-life scenario dialogue summarization dataset](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 5062–5074, Online. Association for Computational Linguistics.
- Marco Damonte, Shay B. Cohen, and Giorgio Satta. 2017. [An incremental parser for Abstract Meaning Representation](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 536–546, Valencia, Spain. Association for Computational Linguistics.
- Marco Damonte and Emilio Monti. 2021. [One semantic parser to parse them all: Sequence to sequence multi-task learning on semantic parsing datasets](#). In *Proceedings of *SEM 2021: The Tenth Joint Conference on Lexical and Computational Semantics*, pages 173–184, Online. Association for Computational Linguistics.
- Ramón Fernandez Astudillo, Miguel Ballesteros, Tahira Naseem, Austin Blodgett, and Radu Florian. 2020. [Transition-based parsing with stack-transformers](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1001–1007, Online. Association for Computational Linguistics.
- DongLai Ge, Junhui Li, Muhua Zhu, and Shoushan Li. 2019. Modeling source syntax and semantics for neural amr parsing. In *IJCAI*, pages 4975–4981.
- James Goodman, Andreas Vlachos, and Jason Naradowsky. 2016. [Noise reduction and targeted exploration in imitation learning for Abstract Meaning Representation parsing](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1–11, Berlin, Germany. Association for Computational Linguistics.
- Hardy Hardy and Andreas Vlachos. 2018. Guided neural language generation for abstractive summarization using abstract meaning representation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 768–773.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *NIPS*.
- Hoang Thanh Lam, Gabriele Picco, Yufang Hou, Young-Suk Lee, Lam M. Nguyen, Dzung T. Phan, Vanessa Lopez, and Ramon Fernandez Astudillo. 2021. [Ensembling graph predictions for amr parsing](#).
- Young-Suk Lee, Ramon Fernandez Astudillo, Tahira Naseem, Revanth Gangi Reddy, Radu Florian, and Salim Roukos. 2020. Pushing the limits of amr parsing with self-learning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 3208–3214.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Belinda Z. Li, Sewon Min, Srinivasan Iyer, Yashar Mehdad, and Wen-tau Yih. 2020. [Efficient one-pass end-to-end entity linking for questions](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6433–6441, Online. Association for Computational Linguistics.

703	Kexin Liao, Logan Lebanoff, and Fei Liu. 2018. Abstract meaning representation for multi-document summarization. In <i>Proceedings of the 27th International Conference on Computational Linguistics</i> , pages 1178–1190.	758
704		759
705		760
706		761
707		
708	Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. 2019. On the variance of the adaptive learning rate and beyond. <i>arXiv preprint arXiv:1908.03265</i> .	762
709		763
710		764
711		765
712	Chunchuan Lyu and Ivan Titov. 2018. AMR parsing as graph prediction with latent alignment . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 397–407, Melbourne, Australia. Association for Computational Linguistics.	766
713		767
714		768
715		769
716		770
717		771
718	Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Rose Finkel, Steven Bethard, and David McClosky. 2014. The stanford corenlp natural language processing toolkit. In <i>ACL</i> .	772
719		773
720		774
721		775
722	Mitchell Marcus, Beatrice Santorini, Mary Marcinkiewicz, and Ann Taylor. 1999. Penn treebank 3.	776
723		777
724		778
725	Arindam Mitra and Chitta Baral. 2016. Addressing a question answering challenge by combining statistical methods with inductive rule learning and reasoning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 30.	779
726		780
727		781
728		782
729		783
730	Tahira Naseem, Abhishek Shah, Hui Wan, Radu Florian, Salim Roukos, and Miguel Ballesteros. 2019. Rewarding smatch: Transition-based amr parsing with reinforcement learning. <i>arXiv preprint arXiv:1905.13370</i> .	784
731		785
732		786
733		787
734		788
735	Clifton Poth, Jonas Pfeiffer, Andreas Rücklé, and Iryna Gurevych. 2021. What to pre-train on? Efficient intermediate task selection . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 10585–10605, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	789
736		790
737		791
738		792
739		793
740		794
741		795
742	Luigi Procopio, Rocco Tripodi, and Roberto Navigli. 2021. SGL: Speaking the graph languages of semantic parsing via multilingual translation . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 325–337, Online. Association for Computational Linguistics.	796
743		797
744		798
745		799
746		800
747		801
748		802
749		803
750	Yada Pruksachatkun, Jason Phang, Haokun Liu, Phu Mon Htut, Xiaoyi Zhang, Richard Yuanzhe Pang, Clara Vania, Katharina Kann, and Samuel Bowman. 2020. Intermediate-task transfer learning with pre-trained language models: When and why does it work? In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5231–5247.	804
751		805
752		806
753		807
754		808
755		809
756		810
757		811
	Sudha Rao, Daniel Marcu, Kevin Knight, and Hal Daumé III. 2017. Biomedical event extraction using abstract meaning representation. In <i>BioNLP 2017</i> , pages 126–135.	812
		813
	Mrinmaya Sachan and Eric Xing. 2016. Machine comprehension using rich semantic representations. In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 486–492.	
	Yanshan Wang, Sijia Liu, Majid Rastegar-Mojarad, Liwei Wang, Feichen Shen, Fei Liu, and Hongfang Liu. 2017. Dependency and amr embeddings for drug-drug interaction extraction from biomedical literature. In <i>Proceedings of the 8th acm international conference on bioinformatics, computational biology, and health informatics</i> , pages 36–43.	
	Ralph M. Weischedel, Eduard H. Hovy, Mitchell P. Marcus, and Martha Palmer. 2017. Ontonotes : A large training corpus for enhanced processing.	
	Taizhong Wu, Junsheng Zhou, Weiguang Qu, Yanhui Gu, Bin Li, Huilin Zhong, and Yunfei Long. 2021. Improving amr parsing by exploiting the dependency parsing as an auxiliary task. <i>Multim. Tools Appl.</i> , 80:30827–30838.	
	Dongqin Xu, Junhui Li, Muhua Zhu, Min Zhang, and Guodong Zhou. 2020. Improving amr parsing with sequence-to-sequence pre-training. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2501–2511.	
	Sheng Zhang, Xutai Ma, Kevin Duh, and Benjamin Van Durme. 2019a. AMR parsing as sequence-to-graph transduction . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 80–94, Florence, Italy. Association for Computational Linguistics.	
	Sheng Zhang, Xutai Ma, Kevin Duh, and Benjamin Van Durme. 2019b. Broad-coverage semantic parsing as transduction . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3786–3798, Hong Kong, China. Association for Computational Linguistics.	
	Yu Zhang and Qiang Yang. 2021. A survey on multi-task learning . <i>IEEE Transactions on Knowledge and Data Engineering</i> , pages 1–1.	
	Zhisong Zhang, Emma Strubell, and Eduard Hovy. 2021. Comparing span extraction methods for semantic role labeling . In <i>Proceedings of the 5th Workshop on Structured Prediction for NLP (SPNLP 2021)</i> , pages 67–77, Online. Association for Computational Linguistics.	
	Zixuan Zhang and Heng Ji. 2021. Abstract meaning representation guided graph encoding and decoding for joint information extraction. In <i>Proceedings of</i>	

814 *the 2021 Conference of the North American Chapter*
815 *of the Association for Computational Linguistics:*
816 *Human Language Technologies*, pages 39–49.

817 Jiawei Zhou, Tahira Naseem, Ramón Fernandez As-
818 tudillo, and Radu Florian. 2021a. [AMR parsing with](#)
819 [action-pointer transformer](#). In *Proceedings of the*
820 *2021 Conference of the North American Chapter of*
821 *the Association for Computational Linguistics: Hu-*
822 *man Language Technologies*, pages 5585–5598, On-
823 line. Association for Computational Linguistics.

824 Jiawei Zhou, Tahira Naseem, Ramón Fernandez As-
825 tudillo, Young-Suk Lee, Radu Florian, and Salim
826 Roukos. 2021b. [Structure-aware fine-tuning of](#)
827 [sequence-to-sequence transformers for transition-](#)
828 [based AMR parsing](#). In *Proceedings of the 2021*
829 *Conference on Empirical Methods in Natural Lan-*
830 *guage Processing*, pages 6279–6290, Online and
831 Punta Cana, Dominican Republic. Association for
832 Computational Linguistics.

A Algorithms

Algorithm 1 Reentrancy Restoration for SRL

Input: Treenode:T

Output: Graph:G

Description: T is root node of the original SRL after node ROOT is added to form tree structure. G is the output graph with possible reentrancy restored.

Global Variables: Dict: V={}. Here Dict is the official data structure of Python’s dictionary.

```

1: for predicate in T.sons do
2:   for son in predicate.sons() do
3:     if son.name in V.keys() then
4:       son = V[son.name]
5:       # restore reentrancy
6:     else
7:       V[son.name] = son
8: return T

```

B Ensemble Models’ Methods

Graphene-4S^E Lam et al. (2021) make use of 4 SPRING models from different random seeds and their proposed graph ensemble algorithm to do the ensembling. They also include another ensemble model named Graphene All which includes four checkpoints from models of different architectures, SPRING(Bevilacqua et al., 2021), APT(Zhou et al., 2021a), T5, and Cai&Lam(Cai and Lam, 2020). We do not report the score of Graphene All since it aggregates models with different inductive bias while our ensemble model only use models from one structure. It is out of the scope for fair comparison.

Structure-aware^E Zhou et al. (2021b) use ensemble results from 3 models’ combination to generate the ensemble model.

Ours (w/ SRL)^E We use the setting the same as Zhou et al. (2021b), we use the average of three models’ parameters as the ensemble model.

C Auxiliary Datasets Description

C.1 Summarization

CNN/DM(Hermann et al., 2015) The CNN / DailyMail Dataset is an English-language dataset containing news articles as written by journalists at CNN and the Daily Mail. The dataset is widely

accepted as benchmark to test models’ performance of summarizing .

DIALOGSUM(Chen et al., 2021) The Real-Life Scenario Dialogue Summarization (DIALOGSUM), is a large-scale summarization dataset for dialogues. Unlike CNN/DM which focuses on monologue news summarization, DIALOGSUM covers a wide range of daily-life topics in the form of spoken dialogue. We use all the training data (13k) to conduct the intermediate training.

C.2 Translation

WMT14 EN-DE We select the first 40k,80k,200k and 400k training examples from WMT14 EN-DE training set to form EN-DE translation intermediate tasks.

C.3 Dependency Parsing

PENN TREEBANK(Marcus et al., 1999) The Penn Treebank (PTB) project selected 2,499 stories from a three year Wall Street Journal (WSJ) collection of 98,732 stories for syntactic annotation. We only utilize the dependency structure annotations to form our intermediate dependency parsing task. There are 39,832 (~40k) sentences.

C.4 Semantic Role Labeling

ONTONOTES(Weischedel et al., 2017) The OntoNotes project is built on two resources, following the PENN TREEBANK(Marcus et al., 1999) for syntax and the PENN PROPBANK for predicate-argument structure. We select 40k sentences with SRL annotations to form intermediate task.

D Model Tuning Description

We tune the hyper-parameters on the SPRING baseline, and then adding the auxiliary data using just those hyper-parameters without any changing.

We use RAdam (Liu et al., 2019) as our optimizer, and the learning rate is $3e^{-5}$. Batch-size is set to 2048 tokens with 10 steps accumulation. The dropout rate is set to 0.3.

Parameter	Searching Space
Learning rate	1e-5, 3e-5, 5e-5, 1e-4
Batch-size	256, 512, 1024, 2048, 4096
Grad. accu.	10
Dropout	0.1, 0.2, 0.3

Table 6: Hyper-parameters searching space

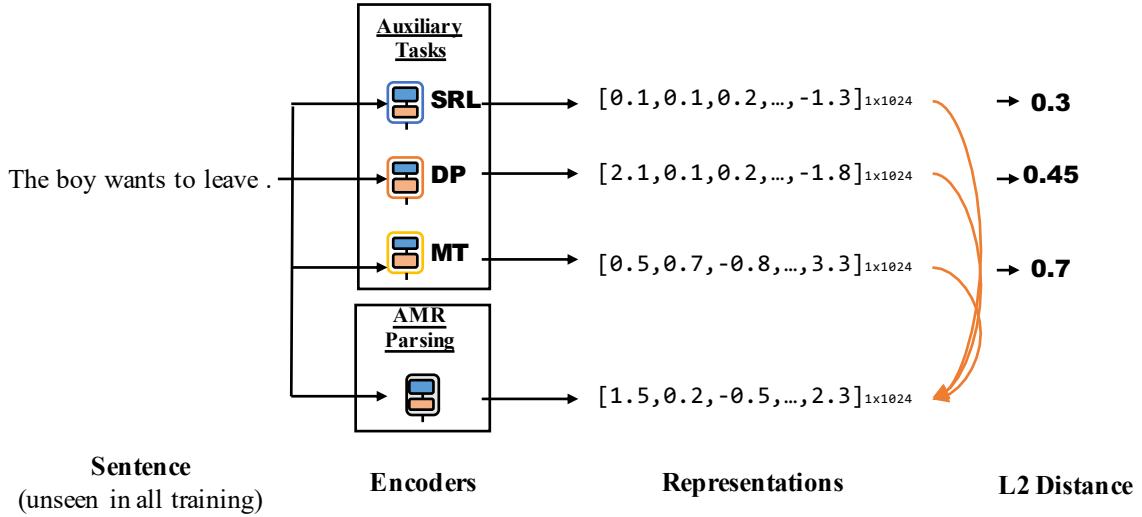


Figure 6: Illustration of how to compute sentence representation distance of different tasks. The sentences used for evaluate are never seen in the training of AMR Parsing and other auxiliary tasks. Cosine Similarity is computed the same way. We collect all sentences’ distance of one encoder to draw the Gaussian distribution curve.

Model	Extra Data	SMATCH	NoWSD	Wiki	Concept-related			Topology-related		
					Conc.	NER	Neg.	Unll.	Reen.	SRL
SPRING (w/ silver) (Bevilacqua et al., 2021)	200k	84.3	84.8	83.1	90.8	90.5	73.6	86.7	72.4	80.5
Ours (w/ Semantic Role Labeling)	40k	84.5	84.9	84.0	90.2	91.8	74.6	87.7	74.2	82.8
- w/ Arg. Reduction(AR)	40k	84.8	85.2	83.9	90.4	92.2	74.2	88.1	74.5	83.0
- w/ Reen. Restoration(RR)	40k	85.0	85.4	83.5	90.6	92.1	75.6	88.2	75.5	83.7
- w/ AR+RR	40k	85.1	85.6	83.6	90.4	91.4	75.7	88.2	75.0	83.5
Ours (w/ Dependency Parsing)	40k	84.4	84.9	82.9	90.1	90.5	73.5	87.8	74.3	82.9
- w/ Redundant Relation Removal (RRR)	40k	84.5	85.0	83.5	90.2	91.2	74.3	88.0	74.5	82.9
- w/ Lemmatization (Lemma)	40k	84.7	85.2	83.8	90.2	91.2	75.0	88.0	74.1	83.0
- w/ RRR + Lemma	40k	85.0	85.4	84.1	90.4	92.5	74.7	88.2	74.7	83.1

Table 7: Full scores of ablation on AMRization methods.