

Цель: сравнить работу парсеров, используя учебный корпус для подбора примеров предложений, обладающими различными грамматическими особенностями, а также найти самый удобный и быстрый парсер, отвечающий всем или большинству требований.

Описание работы:

1. Пояснение устройства парсера: сложность алгоритмов и т.д
2. Пояснение логики работы парсера с примерами, которые могли бы показать, как парсер работает с:
 - Подлежащим, если у него есть/нет зависимых слов
 - Сказуемым (точно так же)
 - Однородными членами
 - Разбиением на клаузы
 - Сложными случаями
3. Повторение пункта 2, но с опечатками в предложениях двух видов: **± пунктуация и её отсутствие**
 - случай, когда при замене/пропуске одной буквы у нас существующее слово меняется на существующее (they → the/thy, и что-то похожее)
 - случай, когда существующее слово становится несуществующим (they - > th)
4. Повторение пункта 2, но проверкой влияния наличия/отсутствия знаков препинания на деление клауз, а также определение синтаксических связей внутри предложений.

Ниже можно увидеть таблицу, в которой наглядно представлено сравнение работы парсеров на конкретных примерах с разными грамматическими особенностями. Таблица поделена на условные разделы, которые нам интересны. Она описывает общую логику работы парсеров (пункт 2).

	UDPipe	Stanford CoreNLP	SpaCy
Подлежащее , если у него есть/нет зависимых слов a. It leads to one problem. b. The number of conflicts decreases. c. Sport facilities are a good way to solve a problem.	a. <u>It</u> (<i>nsubj</i>) b. <u>number</u> (<i>nsubj</i>) of <u>conflicts</u> (<i>nmod</i>) c. <u>Sport</u> (<i>compound</i>) <u>facilities</u> (<i>nsubj</i>)	a. <u>It</u> (<i>nsubj</i>) b. <u>number</u> (<i>nsubj</i>) of <u>conflicts</u> (<i>nmod</i>) c. <u>Sport</u> (<i>compound</i>) <u>facilities</u> (<i>nsubj</i>)	a. It (<i>nsubj</i>) b. Justice (for authors - зависит от <i>nsubj</i>) c. Sport (зависит от <i>nsubj</i>) facilities (<i>nsubj</i>)
	<i>Подлежащее -- nsubj, зависят от него modifiers и compounds</i>		
Сказуемое (точно так же) a. In 21 century science <u>goes</u> forward. b. They <u>have caused</u> some huge problems too.	a. <u>demonstrates</u> (<i>root</i>) b. <u>have</u> (<i>aux</i>) <u>caused</u> (<i>root</i>)	a. The figure <u>demonstrates</u> (<i>root</i>) b. They <u>have</u> (<i>aux</i>) <u>caused</u> (<i>root</i>)	a. goes. (<i>science - nsubj</i>) b. have caused (<i>caused - root, they-nsubj</i>)
	<i>Смысловый глагол -- root, вспомогательные элементы -- aux, везде определены верно.</i>		

<p>Однородные члены</p> <p>a. Non-metallic minerals and chemicals have approximately 10 percent of rail transportation.</p> <p>b. Non-metallic minerals, chemicals have approximately 10 percent of rail transportation.</p> <p>c. It also develops and strengthens its economy.</p>	<p>a. <u>minerals</u> (<i>nsubj</i>) <u>and</u> (<i>cc</i>) <u>chemicals</u> (<i>conj</i>)</p> <p>b. <u>minerals</u> (<i>obl</i>) , (<i>punct</i>) <u>chemicals</u> (<i>nsubj</i>)</p>	<p>a. <u>minerals</u> (<i>nsubj</i>) <u>and</u> (<i>cc</i>) <u>chemicals</u> (<i>conj</i>)</p> <p>b. <u>minerals</u> (<i>nsubj</i>) , (<i>punct</i>) <u>chemicals</u> (<i>appos</i>)</p>	<p>a. minerals (<i>nsubj</i>) and(<i>cc</i>) chemicals(<i>conj</i>) (have - root)</p> <p>b. Non-metallic minerals, chemicals (<i>разделяет на две разные клаузы по запятой</i>) - нужно проверить на других примерах похожих</p>
	<p>c. <u>develops</u> (<i>root</i>) <u>and</u> (<i>cc</i>) <u>strengthens</u> (<i>conj</i>)</p>	<p>c. <u>develops</u> (<i>root</i>) <u>and</u> (<i>cc</i>) <u>strengthens</u> (<i>conj</i>)</p>	<p>c. develops(<i>root</i>) and (<i>cc</i>) strengthens (<i>conj</i>)</p>
	<p><i>У Spacy и Stanford однородный член присоединяется к первому на том же уровне, что и соединительный союз и обозначается как conj (получается два уровня), если есть запятая, Spacy не справляется, так как делит сразу предложение на разные клаузы. UDPipe присоединяет союз на уровень ниже второго однородного слова (получается три уровня). Только Stanford правильно работает с однородными членами, разделенными запятой, без союза.</i></p>		
<p>Разбиение на клаузы</p> <p>a. relative clauses The metals which took the first place in the previous chart take only the 5th with 11%.</p> <p>b. subordinate clauses A lot of people think that it is ugly and unimpressive.</p> <p>c. coordinate clauses Children are our future, and they should be healthy from birth.</p> <p>d. adverbial clause modifier Maybe I would have gone to France for several months and would speak French as fluently as the French do.</p>	<p>a. <u>which</u> (<i>nsubj</i>) <u>took</u> (<i>acl:relcl</i>) <u>the first place in the previous chart</u></p>	<p>a. <u>which</u> (<i>nsubj</i>) <u>took</u> (<i>acl:relcl</i>) <u>the first place in the previous chart</u></p>	<p>a. <u>which</u>(<i>nsubj</i>) <u>took</u> (<i>relcl</i>, <i>зависит от metals-nsubj</i>, <i>зависит от took-root</i>) <u>the first place in the previous chart</u> .</p>
	<p><i>Relative clauses всеми парсерами распознаются, вершина зависимой клаузы зависит от слова в главной клаузе, к которому семантически присоединено.</i></p>		
	<p>b. <u>that</u> (<i>mark</i>) <u>it</u> (<i>nsubj</i>) <u>is</u> (<i>cop</i>) <u>ugly</u> (<i>ccomp</i>) <u>and</u> (<i>cc</i>) <u>unimpressive</u> (<i>conj</i>)</p>	<p>b. <u>that</u> (<i>mark</i>) <u>it</u> (<i>nsubj</i>) <u>is</u> (<i>cop</i>) <u>ugly</u> (<i>ccomp</i>) <u>and</u> (<i>cc</i>) <u>unimpressive</u> (<i>conj</i>)</p>	<p>b. <u>that</u>(<i>mark</i>) <u>it</u> (<i>nsubj</i>) <u>is</u>(<i>ccomp</i>, <i>зависит от root - think</i>) <u>ugly</u>(<i>acomp</i>) and <u>unimpressive</u> (<i>conj</i> - <i>зависит от ugly</i>).</p>
	<p><i>Subordinate clauses всеми парсерами распознаются, вершина зависимой клаузы зависит от вершины главной.</i></p>		
	<p>c. <u>and</u> (<i>cc</i>) <u>they</u> (<i>nsubj</i>) <u>should</u> (<i>aux</i>) <u>be</u> (<i>cop</i>) <u>healthy</u> (<i>conj</i>)</p>	<p>c. <u>and</u> (<i>cc</i>) <u>they</u> (<i>nsubj</i>) <u>should</u> (<i>aux</i>) <u>be</u> (<i>cop</i>) <u>healthy</u> (<i>conj</i>) <u>from birth</u></p>	<p>c. <u>and</u> (<i>cc</i>) <u>they</u> <u>should</u>(<i>aux</i>) <u>be</u>(<i>conj</i>, <i>зависит от root - are</i>) <u>healthy</u> from birth.</p>
<p><i>Coordinate clauses всеми парсерами распознаются, вершина зависимой клаузы зависит от вершины главной.</i></p>			
<p>d. <u>as</u> (<i>case</i>) <u>the</u> (<i>det</i>) <u>French</u> (<i>obl</i>) <u>do</u>(<i>advmod</i>)</p> <p><i>Не распознает такое как отдельную клаузу: считает все</i></p>	<p>d. <u>as</u> (<i>mark</i>) <u>the</u> <u>French</u> (<i>nsubj</i>) <u>do</u> (<i>advcl</i>)</p>	<p>d. <u>as</u> the <u>French</u>(<i>nsubj</i>) <u>do</u>(<i>advcl</i> - <i>зависит от fluently</i>)</p>	

	сравнительным оборотом без подлежащего и сказуемого		
	С последним пунктом проблема возникла у UDPipe, остальные парсеры распознали конструкцию и верно выделили подлежащее и сказуемое в интересующей клаузе.		
<p>5. сложные случаи</p> <p>a. инверсия The bachelor's degree have a 45 per cent of males and 55 per cent of females.</p> <p>b. gerund as a subject Having thoughts of committing crimes comes from the period of socialization.</p> <p>c. there is/are There are several opinions on this issue that I am going to discuss.</p> <p>d. предложения с именным сказуемым Sport facilities are a good way to solve a problem.</p> <p>e. пассив The building was destroyed by the demolition crew.</p> <p>f. have-иметь и have-вспомогательный Our society always has crime. I have done my homework.</p>	a. <u>degree</u> (nsubj), <u>have</u> (root)	a. <u>degree</u> (nsubj), <u>have</u> (root)	a. <u>degree</u> (nsubj), <u>have</u> (root)
	Все парсеры не отслеживают нарушение порядка слов SVO вне зависимости от согласования, следовательно, не улавливают инверсию, даже если определенные подлежащее со сказуемым не согласовываются		
	b. <u>Having</u> (csubj) thoughts of committing crimes <u>comes</u> (root)	b. <u>Having</u> (csubj) thoughts of committing crimes <u>comes</u> (root)	b. <u>Having</u> (csubj) thoughts of committing crimes <u>comes</u> (root)
	Верно определены подлежащее-герундий и сказуемое всеми парсерами.		
	c. <u>There</u> (expl) <u>are</u> (root) several <u>opinions</u> (nsubj)	c. <u>There</u> (expl) <u>are</u> (root) several <u>opinions</u> (nsubj)	c. <u>There</u> (expl) <u>are</u> (root) <u>opinions</u> (attr) Именная часть подлежащего выделяется тегом attr
	Верно определяются подлежащее и сказуемое клауз в конструкциях такого типа всеми парсерами.		
	d. Sport <u>facilities</u> (nsubj) <u>are</u> (cop) a good <u>way</u> (root)	d. Sport <u>facilities</u> (nsubj) <u>are</u> (cop) a good <u>way</u> (root)	d. Sport <u>facilities</u> (nsubj) <u>are</u> (root) a good <u>way</u> (attr)
	Из-за разной логики парсеров по-разному выделяются вершины клауз и, следовательно, зависимости внутри них. UDPipe и Stanford считают сказуемым именную часть, а глагол-связку обозначают как copula. Spacy же считает сказуемым глагол-связку, а именную часть обозначает как attribute.		
	e. <u>The building</u> (nsubjpass) <u>was</u> (auxpass) <u>destroyed</u> (root)	e. <u>The building</u> (nsubjpass) <u>was</u> (auxpass) <u>destroyed</u> (root)	e. <u>The building</u> (nsubjpass) <u>was</u> (auxpass) <u>destroyed</u> (root)
	Конструкция пассива верно определяется всеми парсерами, предусмотрены отдельные теги для подлежащего и сказуемого в пассивных конструкциях.		
f. <u>society</u> (nsubj) <u>has</u> (aux, AUX) <u>crime</u> (root, VERB) I (nsubj) <u>have</u> (aux, AUX) <u>done</u> (root) В идеале, должен различать have-aux и	f. <u>society</u> (nsubj) <u>has</u> (root, VERB) I (nsubj) <u>have</u> (aux, AUX) <u>done</u> (root)	f. has(root, но в частеречном разборе считает have aux в любом случае) society (nsubj)	

	<p><i>have-verb, однако в значительной части случаев(см f.1) не справляется с этой задачей</i></p>		
<p><i>Для Spacy have в частеречном разборе всегда AUX, но это не влияет на проверку правильного/неправильного согласования, так как, когда Spacy выделяет have как root - have - модальный или вспомогательный, а иначе have - вспомогательный, потому что логика парсера такова, что в аналитических формах root является смысловой глагол.</i></p> <p><i>Для остальных парсеров ситуация иная: в UDPipe есть разделения на have-aux и have-verb, однако в значительной части случаев(см f.1) оно не является верным. У Stanford NLP также есть разделения на have-aux и have-verb. Но стоит отметить, что это разделение нужно для определения вершины клаузы, а это все парсеры делают верно.</i></p>			

English version:

	UDPipe	Stanford CoreNLP	SpaCy
<p>Subject if it has / doesn't have dependent words</p> <p>a. It leads to one problem.</p> <p>b. The number of conflicts decreases.</p> <p>c. Sport facilities are a good way to solve a problem.</p>	<p>a. <u>It</u> (<i>nsubj</i>)</p> <p>b. <u>number</u> (<i>nsubj</i>) of <u>conflicts</u> (<i>nmod</i>)</p> <p>c. <u>Sport</u> (<i>compound</i>) <u>facilities</u> (<i>nsubj</i>)</p>	<p>d. <u>It</u> (<i>nsubj</i>)</p> <p>e. <u>number</u> (<i>nsubj</i>) of <u>conflicts</u> (<i>nmod</i>)</p> <p>f. <u>Sport</u> (<i>compound</i>) <u>facilities</u> (<i>nsubj</i>)</p>	<p>d. <u>It</u> (<i>nsubj</i>)</p> <p>e. Justice (for authors - зависит от <i>nsubj</i>)</p> <p>f. Sport (зависит от <i>nsubj</i>) <u>facilities</u> (<i>nsubj</i>)</p>
<p><i>Subject is marked as 'nsubj', dependent words are modifiers or compounds.</i></p>			
<p>Predicate</p> <p>a. In 21 century science goes forward.</p> <p>b. They <u>have caused</u> some huge problems too.</p>	<p>a. <u>demonstrates</u> (<i>root</i>)</p> <p>b. <u>have</u> (<i>aux</i>) <u>caused</u> (<i>root</i>)</p>	<p>a. The figure <u>demonstrates</u> (<i>root</i>)</p> <p>b. They <u>have</u> (<i>aux</i>) <u>caused</u> (<i>root</i>)</p>	<p>a. goes. (<i>science - nsubj</i>)</p> <p>b. have caused (<i>caused - root, they-nsubj</i>)</p>
<p><i>Main verb is marked as 'root'; auxiliary verbs are marked as 'aux'. All 3 parsers marked them correctly</i></p>			
<p>Homogeneous parts</p> <p>a. Non-metallic minerals and chemicals have approximately 10 percent of rail transportation.</p>	<p>a. <u>minerals</u> (<i>nsubj</i>) <u>and</u> (<i>cc</i>) <u>chemicals</u> (<i>conj</i>)</p> <p>b. <u>minerals</u> (obl) , (<i>punct</i>) <u>chemicals</u> (nsubj)</p>	<p>a. <u>minerals</u> (<i>nsubj</i>) <u>and</u> (<i>cc</i>) <u>chemicals</u> (<i>conj</i>)</p> <p>b. <u>minerals</u> (nsubj) , (<i>punct</i>) <u>chemicals</u> (appos)</p>	<p>a. minerals (<i>nsubj</i>) and(cc) chemicals(<i>conj</i>) (<i>have - root</i>)</p> <p>b. Non-metallic minerals, chemicals</p>

<p>b. Non-metallic minerals, chemicals have approximately 10 percent of rail transportation.</p> <p>c. It also develops and strengthens its economy.</p>	<p>c. <u>develops</u> (root) <u>and</u> (cc) <u>strengthens</u> (conj)</p>	<p>c. <u>develops</u> (root) <u>and</u> (cc) <u>strengthens</u> (conj)</p>	<p>c. <u>develops</u>(root) and (cc) <u>strengthens</u> (conj)</p>
<p><i>Spacy and Stanford CoreNLP place the head of the second homogeneous part on the same level as the coordinating conjunction. Both words are connected to the head of the first homogeneous part, which is placed on one level above. Thus two levels are made.</i></p> <p><i>UDPipe connects the coordinating conjunction to the second homogeneous part (which itself is connected to the first homogeneous part), thus making three levels.</i></p> <p><i>Spacy doesn't parse well sentences in which homogeneous parts are divided by a comma because it often separates clauses and creates two trees, one for each clause.</i></p> <p><i>Stanford CoreNLP handles well homogeneous parts that are divided by a comma.</i></p> <p><i>UDPipe does not recognise two homogenous parts divided by commas. The parser marks one of the parts correctly, while the second one depends from it either as oblique (for nouns) or as adverbial clause (for verbs). In case with three or more comma-separated parts UDPipe performs well.</i></p>			
<p>Dependent clauses</p> <p>a. relative clauses The metals which took the first place in the previous chart take only the 5th with 11%.</p> <p>b. subordinate clauses A lot of people think that it is ugly and unimpressive.</p> <p>c. coordinate clauses Children are our future, and they should be healthy from birth.</p> <p>d. adverbial clause modifier Maybe I would have gone to France for several months and would speak French as fluently as the French do.</p>	<p>a. <u>which</u> (nsubj) <u>took</u> (acl:relcl) <u>the first place in the previous chart</u></p>	<p>a. <u>which</u> (nsubj) <u>took</u> (acl:relcl) <u>the first place in the previous chart</u></p>	<p>a. <u>which</u>(nsubj) <u>took</u> (relcl, <i>завучум ом metals-nsubj, завучум ом took-root</i>) <u>the first place in the previous chart</u>.</p>
<p><i>All parsers correctly recognize relative clauses. The head of the relative clause is dependent from the word of the main clause to which it is semantically related to.</i></p>			
<p><i>Subordinate clauses are recognized correctly by all parsers. The head of the subordinate clause is dependent from the head of the main clause.</i></p>			
<p><i>Coordinate clauses are recognized correctly by all parsers. The head of the coordinate clause is dependent from the head of the main clause.</i></p>			
<p><i>Coordinate clauses are recognized correctly by all parsers. The head of the coordinate clause is dependent from the head of the main clause.</i></p>			
<p><i>UDPipe failed to recognize this clause, marking it as an adverbial modifier instead.</i></p>			

	The only parser that failed to recognize comparative clause was UDPipe. Spacy and StanfordNLP performed the task correctly.		
<p>5. difficult cases</p> <p>a. inversion The bachelor's degree have a 45 per cent of males and 55 per cent of females.</p> <p>b. gerund as a subject Having thoughts of committing crimes comes from the period of socialization.</p> <p>c. there is/are There are several opinions on this issue that I am going to discuss.</p> <p>d. sentences with a noun predicate Sport facilities are a good way to solve a problem.</p> <p>e. passive The building was destroyed by the demolition crew.</p> <p>f. semantic 'have' and auxiliary 'have' Our society always has crime. I have done my homework.</p>	a. <u>degree</u> (nsubj), <u>have</u> (root)	a. <u>degree</u> (nsubj), <u>have</u> (root)	a. <u>degree</u> (nsubj), <u>have</u> (root)
	<i>All three parsers do not recognize the disruption of the standard SVO order regardless of whether the agreement is right or wrong. Therefore, they are unable to recognize inversion even if certain subjects and predicate verbs do not have correct agreement.</i>		
	b. <u>Having</u> (csubj) thoughts of committing crimes <u>comes</u> (root)	b. <u>Having</u> (csubj) thoughts of committing crimes <u>comes</u> (root)	b. <u>Having</u> (csubj) thoughts of committing crimes <u>comes</u> (root)
	<i>Gerund subject is marked as 'csubj'. All three parsers recognize it correctly.</i>		
	c. <u>There</u> (expl) <u>are</u> (root) several <u>opinions</u> (nsubj)	c. <u>There</u> (expl) <u>are</u> (root) several <u>opinions</u> (nsubj)	c. <u>There</u> (expl) <u>are</u> (root) <u>opinions</u> (attr) Именная часть подлежащего выделяется тегом attr
	<i>All three parsers are able to correctly parse this type of sentence.</i>		
	<i>All three parsers mark 'there' as 'expl', 'are' or 'is' as 'root'. UDPipe and Stanford CoreNLP mark the noun as 'nsubj'; Spacy marks the noun as 'attr'.</i>		
	d. Sport <u>facilities</u> (nsubj) <u>are</u> (cop) a good <u>way</u> (root)	d. Sport <u>facilities</u> (nsubj) <u>are</u> (cop) a good <u>way</u> (root)	d. Sport <u>facilities</u> (nsubj) <u>are</u> (root) a good <u>way</u> (attr)
	<i>Spacy differently marks the head of the main clause because its parsing logic is different from that of UDPipe and Stanford CoreNLP. The latter two mark the noun predicate as 'root', and the copula verb as 'cop'. Spacy, on the other hand, marks the noun predicate as 'attr', and the copula verb as 'root'.</i>		
	e. <u>The building</u> (nsubjpass) <u>was</u> (auxpass) <u>destroyed</u> (root)	e. <u>The building</u> (nsubjpass) <u>was</u> (auxpass) <u>destroyed</u> (root)	e. <u>The building</u> (nsubjpass) <u>was</u> (auxpass) <u>destroyed</u> (root)
<i>Passive clauses are recognized correctly by all three parsers. There are distinct tags for auxiliary verbs and subjects in the passive construction – 'auxpass', 'nsubjpass'.</i>			
f. <u>society</u> (nsubj) <u>has</u> (aux, AUX) <u>crime</u> (root, VERB)	f. <u>society</u> (nsubj) <u>has</u> (root, VERB)	f. has(root, AUX) society (nsubj)	
I (nsubj) <u>have</u> (aux, AUX) <u>done</u> (root)	I (nsubj) <u>have</u> (aux, AUX) <u>done</u> (root)		
<i>Ideally, UDPipe should distinguish between semantic 'have' and auxiliary 'have', but it often (see f.1) fails to do so.</i>			

	<p><i>Both UDPipe and Stanford CoreNLP have a distinction between semantic 'have' and auxiliary 'have': their POS-taggers mark the former as 'VERB' (while also syntactically marking it as 'root') and the latter as 'AUX' (while also syntactically marking it as 'aux');. However, in most cases UDPipe marks the verb (see f.1) incorrectly. Stanford CoreNLP usually handles the distinction well.</i></p> <p><i>Spacy's POS-tagger always marks 'have' as 'AUX', but it does not disrupt the work of the parser. In most cases, when the verb is semantic it is marked as 'root', and when the verb is auxiliary it is marked as 'aux'.</i></p>		

It also develops, strengthens its economy.

Asia, Europe, Africa didn't change during two periods.