

# Словацкий национальный корпус

Radovan Garabík

Jazykovedný ústav Ľ. Štúra SAV

813 64 Bratislava, Slovakia

[korpus@korpus.juls.savba.sk](mailto:korpus@korpus.juls.savba.sk), <http://korpus.juls.savba.sk>

Príspevok bol prednesený na medzinárodnej konferencii Корпусная лингвистика – 2004 (Sankt Peterburg 2004) a bude publikovaný v zborníku.

**Abstract.** First preliminary work on text corpus of Slovak language started in 1992, with primary intention to provide enough searchable material for lexicographic purposes. In 2002, the Government of the Slovak Republic approved the project to build a representative corpus of contemporary Slovak language to be used in further linguistic research. Currently, Slovak National Corpus contains 192 million lemmatised and morphologically annotated words, with the plans of reaching 200 million words of balanced corpus in 2005. A new morphosyntactic tagset for the Slovak language has been developed and is used to create a manually lemmatised and annotated core, useful for further analysis and NLP tools training. Incoming texts are converted into common text format based on XCES (TEI) standard, with additional linguistic annotation added at this stage, and then converted into internal format for the corpus manager. Access to the corpus can be given to any person in exchange for a simple registration procedure, provided that the corpus is used for educational or personal non-commercial purposes. Texts in the corpus include complete bibliographic and style-genre annotation, described here in detail.

## История

Корпусная (и вообще математическая) лингвистика не пользовалась большим вниманием в словацкой лингвистике, хотя первые шаги были сделаны уже в 1962 – 1970 годах после возникновения отдела Математической лингвистики Института словацкого языка Словацкой академии наук. Отдел занимался преимущественно до тех пор не исследованными аспектами словацкого языка, в том числе главным образом статистическим анализом.

Возобновление интереса к компьютером как к средству обработки лингвистических данных произошло в связи с проникновением компьютеров в научное общество после 1989-ого года[1], но главная польза компьютеров состояла в поддержке лексикографических работ. Лишь в 1991-ом году появились первые мысли[2] о необходимости коллекции достаточного количества текстов в цифровом виде с возможностью сложных поисков. Эта коллекция, без всякой следующей лингвистической обработки, стала называться Корпус текстов словацкого языка[3], и ею пользовались только в институте, главным образом как дополнительный источник данных к подготовке Короткого словаря словацкого языка.

Текстовый корпус словацкого языка продолжал увеличиваться в количестве словоупотреблений, но скоро стало ясно что для подготовки большого, современным стандартом соответствующего корпуса есть необходимо принять более систематический подход. В начале 2002-ого года правительство Словацкой республики одобрило проект создания представительного корпуса современного словацкого языка для будущего лингвистического исследования, вместе с созданием отдела Словацкого национального корпуса в рамках Института языкознания Словацкой академии наук.

## Принципы постройки корпуса

Словацкий национальный корпус, чтобы предоставить желаемую информацию своим пользователям, должен быть не только механический сбор текстов (хотя и это может своим образом использоваться в лингвистических исследованиях, как показывает бывшая версия текстового корпуса словацкого языка), но он должен отвечать нескольким определённым требованиям. [4]

Во первых, Словацкий национальный корпус должен быть корпусом *современного* словацкого письменного языка. Под современным словацким языком подразумевается язык после последней значительной реформы орфографии в 1953 году. Постепенные малейшие поправки не вносили в язык большие изменения, и самая последняя реформа в 1994-ом году, хотя и не так незначительная, не проникла достаточно в общество. В корпус решено включать тексты возникшие в 1955 – 2005 годах. Нижний предел назначен промежутком времени который истёк пока ещё новая орфография не стала общепринятой. Тексты из более давних пор в принципе не исключены из обработки, но они не входят в окончательный корпус – постройка диахронического корпуса Словацкого языка будет отдельным проектом. Принято решение включать в корпус прежде всего современные тексты, поскольку их можно нетрудно приобрести. Но будет сделано тоже сечение по текстах прошлого времени.

Во вторых, корпус должен быть *представительным*, т.е. он должен включать самые разные формы языка и стилей, пропорционально к их распространённости в письменном общении.

Дальше, корпус должен содержать как минимум информацию о лемме лексемы и соответствующих грамматических категориях.

Тексты в корпусе также должны сохранять информацию о их происхождении и включать остальные метатекстовые данные о стиле, жанре, авторе, полу автора, источнике текста итд. Эта информация очень важна для эвентуальной постройки разных субкорпусов.

И наконец, корпус должен быть общедоступный для всех интересующихся, как профессиональных лингвистов, так и для любителей языка или учащихся словацкий язык..

Желанная версия представительного Словацкого национального корпуса должна содержать приблизительно 200 миллионов словоупотреблений с полной морфологической разметкой и с добавленной библиографической информацией.

## Сбор текстов

Учитывая проблемы которые встретили остальные корпуса, тексты в Словацкий национальный корпус входят только когда удастся приобрести разрешение включить тексты в состав корпуса. Хотя в словацком авторском законе есть возможность использовать тексты для научных целей, следующее распространение текстов не позволено ни в каком случае, и такой корпус не мог бы быть общедоступен.

Главными источниками текстов являются издательства, газеты и журналы, но довольно часто тексты (или разрешение для включения текстов в корпус)

приобретены прямо от авторов.

## Обработка текстов

Тексты в корпус входят в четырёхуровневую иерархию. На первом уровне – архив, в котором данные сохраняются в оригинальном виде, т.е. так как они были получены.

Следующий уровень – банк, где тексты переведены в общий формат вроде простого XML, содержавшего нужную информацию о типографии в оригинальном формате. От действительного XML этот формат отличается возможностью перекося тэгов, который иногда очень трудно удалить. На этом уровне тексты аннотированы библиографической разметкой. Один текст в банке соответствует одной логической единицы документа, например один роман, одна статья газетного текста или один закон из сбора. Из этого следует что одному документу в архиве часто соответствует большое количество единиц в банке.

Следует корпусоид, в котором тексты находятся в формате XCES[5], содержавшие уже тексты с полной морфосинтаксической разметкой. Иначе структура корпусоида совпадает структуре банка.

Наконец тексты поступают в секцию «дата», в цифровом формате конкретного корпусного менеджера. Здесь тексты уже распределены по разным субкорпусам.

## Метатекстовая разметка

Каждая единица текста в банке аннотирована с учётом библиографического источника, стиля и жанра текста. Для каждого текста, аннотация находится в самостоятельном текстовом файле, где в каждой строке записана комбинация ключ-значение. Значение может быть либо свободное, либо избранное из множества возможных значений. Аннотация в архиве ведётся подобным образом, но немного проще.

У каждого ключа возможны эти специальные значения:

- пустое место: неопределённое значение – о значении не можно сказать ничего точного
- XXX: неизвестное значение – его не удалось открыть, напр. имя автора в газетной статье
- YYY: неопределимое значение – в принципе невозможно определит его, напр. пол переводчика, если произведение не является переводом
- MIX: смешение других значений
- MSC: другое значение
- TTT: неизвестное значение которое надо определить позже

Полное описание аннотации архива и банка в приложении.

## Лингвистическая разметка

Основным исходным пунктом лингвистической обработки корпуса является токенизация. В корпусе принято решение разделять текст на самые малейшие

логические единицы которые можно определить в письменных документах. Это прежде всего значит, что в корпусе ни каким образом не выделены составные названия – например *Sankt-Peterburg* будут три токена, *Peterburg* с полной морфосинтаксической пометкой, *Sankt* без морфологической информации и „-“ как пунктуационный знак. Главная причина таково решения была невозможность точного определения границы когда составное название уже не можно считать одной лексической единицей текста.

Но с другой стороны, соединённые слова не разделяются, так например местоимения с агглютинативными предлогами (*prej, oň, zaňho...*) оставлены как особый вид местоимений. В словацком языке, в отличие от чешского, они встречаются достаточно редко чтобы такой подход был обоснованным.

При подготовке формального описания грамматических категорий словацкого языка подходящего для автоматической обработки мы опирались преимущественно на два уже существующие проекта автоматической морфологической разметки словацкого языка. Первый из них – автоматический анализатор построен на Институте формальной и прикладной лингвистики Карлова университета в Праге[6, 7]. Структура помет использована в этом анализаторе (и именно и анализатор) идентична с пометами использованными в Чешском национальном корпусе, с некоторыми мелкими отличиями совпадающими с отсутствием редких грамматических явлений чешского языка в словацком. Пометы характеризованы позиционной системой, каждой позиции ж помете соответствует одна грамматическая категория. Категории не имеющее смысл для указанной части речи обозначены знаком „-“. Анализатор включает модуль дисамбигуации и он хорошо проверен входя в многие проекты нуждающиеся в морфологическом анализе, в том числе и коммерческие. Словарь корневых морфем и парадигм словацкого языка достаточно разработан, хотя и не достигает размер чешского словаря.

Второй анализатор *Ajka* построен на Факультете информатики в Брно[8]. Он независимый на языке и в принципе подходит для всех языков которые используют главным образом аффиксы для реализации грамматических категорий. Вопрос только в существовании словарей. Структура помет состоит из пар символов, где первый символ пары определяет грамматическую категорию и второй – значение категории. Если категория не имеет смысл для указанной части речи, соответствующий символ просто не включён в помету. *Ajka* не включает модуль дисамбигуации, и в настоящее время проводятся работы над созданием достаточного словаря словацких морфем в сотрудничестве отдела Словацкого национального корпуса и Факультета информатики.

Есть ещё третья система помет, которая использована в проекте *Multext-East*[9]. Это система была разработана для общей разметки славянских языков. Словацкая версия не существует, но система достаточно универсальная чтобы было можно ей воспользоваться.

При проектировке собственной системы было необходимо учесть несколько противоречивых требований, преимущественно основанных на непригодности формальных теоретических описаний морфологии словацкого языка к машинной обработке. Главный принятый принцип – стремление к описанию чисто морфологических категорий, без учёта синтаксических или даже семантических ролей словоупотреблений. К сожалению, не всегда было возможно соблюдать этот

принцип.

Наша система помет имеет черты как позиционных, так и ключевых систем.. Пометка начинается знаком определяющим часть речи, следует знак соответствующий парадигму (если это флективная часть речи), и следуют знаки грамматических категорий, причём значение каждой грамматической категории соответствует одно-однозначно одному знаку, знаки не повторяются и следуют в определённом порядке, и в пометке находятся только знаки релевантных категорий.

Строго морфологический принцип ведёт к тому что прошедшее время глаголов не выделено как особая категория, вместо того выделяется L-причастие. Будущее время назначено только в совершенных глаголах, и будущее время глаголов несовершенного вида назначено только в помете вспомогательного глагола „byť“

Вопрос звательного падежа – в словацком языке звательный падеж не осуществлён морфологией, он почти всегда совпадает с именительным. Но, есть несколько десятков слов, у которых сохранилась древняя морфология в звательном падеже, и некоторые из этих слов встречаются очень часто, чтобы ими было можно пренебречь. Сверх того, в разговорном языке появляется „новый звательный падеж“, возникший из сокращённых или уменьшительных вариантов прежде всего собственных имён или названий семейных отношений – эволюцией звательного падежа словацкий язык близко напоминает русский. Из-за этого, было решено включить в систему помет тоже звательный падеж, хотя это вносит довольно большую омонимию с именительным.

Причастия, хотя своей морфологией не отличаются от прилагательных имён, введены как собственная часть речи.

Омонимы, которых морфология та же самая, не отличаются (у них та же самая лемма).

Описание системы помет в Таб.1.

Для автоматической разметки в корпусе используется сейчас словацкая версия анализатора сделанного в Праге. Кроме того есть „ядро“ корпуса, состоящее из текстов с ручной морфосинтаксической разметкой, полезное для будущего исследования и тренировки программных средств языковой обработки.

## Поиск в корпусе

В главную очередь, корпус разделяется на три субкорпуса, в зависимости от лицензии текстов. Первый, доступный для широкого общества, состоит из текстов для которых есть разрешение использования в корпусе без добавленных ограничений. Однако, некоторые текстодатели разрешили включить тексты в корпус при условии что ими могут пользоваться только сотрудники института языкознания, и эти тексты вместе с первой группой создают второй субкорпус.

Третий субкорпус состоит из совершенно всех текстов, даже из тех для которых не существует разрешение предоставления другим лицам (потому что переговоры о приобретении разрешения ещё не завершились, или есть другие препятствия). Этот субкорпус доступен только членам отдела словацкого национального корпуса.

Доступ к первому, публичному корпусу может приобрести любое частное лицо, подписавшее договор о пользовании корпусом. Право пользоваться корпусом

личное, не переносное и действует в течении одного года, после чего его надо возобновить. Доступ организациям или другим группам не разрешается, с исключением студентов на уроках связанных с корпусной лингвистикой, которым предоставлен общий доступ, действителен обыкновенно в течении одного семестра. Доступ позволен исключительно для научно-исследовательских или частных целей, и ни в каком случае не может быт предоставлен для коммерческой пользы.

Для поиска в корпусе пользуется система Manatee/Bonito[10], которая состоит из сервера (Manatee) и клиента (Bonito). Пользователи корпуса подключаются к серверу и задают вопросы, которые проводятся на сервере. Сервер позволяет простой поиск одного слова, или фразы (нескольких слов в очередном порядке), или произвольных регулярных выражений. К каждому словоупотреблению может быть добавлена информация в виде атрибутов, по которым тоже можно проводить поиск. В Словацком корпусе добавлены атрибуты *lemma* и *tag*, соответствующие лемме и морфосинтаксической помете.

Manatee дальше позволяет проводить разные статистические исследования в рамках колокаций, как MI-score, T-score, и также частотный анализ разных форм отысканных единиц.

К каждому тексту может быть получена полная информация о метатекстовой (библиографической) разметке, которая находится в банке, и в контексте можно включить изображение структурных помет текста, показывавших или разные типографские черты текста (как например заголовки, курсив), или внутреннюю структуру текста (конец абзаца, конец предложения). Поиск можно ограничить только в рамках желающих структур. Manatee позволяет работать с несколькими корпусами, и каждый пользователь имеет доступ к версии корпуса содержащей все тексты подходящей лицензии, и к версии с исключёнными текстами, которые по разным причинам мешали бы в лингвистическом анализе – например тексты из специализированных лингвистических журналов или тексты с отсутствующими диакритическими знаками.

## Заключение

В настоящее время Словацкий национальный корпус состоит из 192 миллионов словоупотреблений. Публично доступный субкорпус содержит 183 миллионов словоупотреблений. Корпусы включают метатекстовую библиографическую и стилистическо-жанровую разметку, и лингвистическую разметку с показанием леммы и морфосинтаксических категорий. В корпус постоянно добавляются новые тексты и дополнительная лингвистическая обработка.

Отдел Словацкого национального корпуса принимает участие в разных проектах, например в подстройке словацкой версии WordNet (вместе с Университетом в Прешове), синтаксически аннотированного корпуса словацкого языка (вместе с Университетом в Прешове и Институтом формальной и прикладной лингвистики Карлова университета в Праге), параллельных Словацко-Хорватского (вместе с Хорватским национальным корпусом) и Словацко-Русского (вместе с Кафедрой математической лингвистики Филологического факультета Санкт-Петербургского государственного университета) корпусов.

## Приложение

Описание аннотации в архиве (перечень возможных ключей):

- Name: Название документа
- Acqdate: Число когда текст попал в корпус
- Id: Однозначный идентификационный номер текста
- Source: Источник текста (напр. название издательства)
- License: Тип лицензии текста
- Comment: Произвольный комментарий

Описание аннотации в банке:

- Name: Название текста
- Origname: Оригинальное название текста (если это перевод)
- Author: Имя автора (авторов), как напечатано в произведении
- Origauthor: Оригинальное имя автора (на оригинальном языке, если оно отличается от предыдущего значения)
- Translator: Имя переводчика
- Translation:

<i>trn</i>	перевод
<i>org</i>	оригинальный текст
<i>fir</i>	свободный перевод

- ISBN: ISBN номер
- ISSN: ISSN номер
- SourceId: Идентификационный номер документа архива, который является источником этого текста
- Rhyme: рифм текста

<i>nrh</i>	не рифмический
<i>rhy</i>	рифмический
<i>MIX</i>	частично рифмический

- Type: тип текста, подробнее определён ключом Subtype:

	Type		Subtype
<i>img</i>	художественный		
		<i>poe</i>	поэзия
		<i>pro</i>	проза
		<i>dra</i>	драма
<i>inf</i>	информативный текст		
		<i>pub</i>	публицистика
		<i>adv</i>	реклама, объявления
		<i>adm</i>	администрация
<i>prf</i>	профессиональный текст		
		<i>sci</i>	научные тексты, статьи, профессиональные журналы, учебники для ВУЗ-ов
		<i>pop</i>	популярно-научные тексты
		<i>txb</i>	учебники для основных и средних школ

	Type		Subtype
		<i>enc</i>	энциклопедии
		<i>man</i>	справочники
<i>liv</i>	коммуникация		
		<i>spk</i>	разговорная речь
		<i>wri</i>	письменная коммуникация

● Genre: жанр текста

	Genre
<i>ver</i>	стихи
<i>son</i>	песня
<i>scd</i>	театральный сценарий
<i>scf</i>	киносценарий
<i>scr</i>	перезапись радиопередач
<i>nov</i>	роман
<i>col</i>	рассказ
<i>ess</i>	очерк
<i>mem</i>	воспоминания
<i>let</i>	письмо
<i>chr</i>	хроника
<i>dis</i>	диссертация
<i>std</i>	студия
<i>tcl</i>	статья
<i>rfl</i>	рассуждение
<i>ref</i>	реферат
<i>lct</i>	доклад

● Subgenre: точнее определяет жанр текста для жанров *nov*, *col*, *ess*

	Subgenre
<i>crm</i>	детективный роман
<i>scf</i>	фантастика
<i>bel</i>	художественная литература
<i>jun</i>	литература для молодёжи
<i>trv</i>	путёвые записки
<i>fac</i>	литература факта

● Domain

	Domain		Subdomain
<i>ars</i>	художественная наука		
		<i>mus</i>	музыка, балет, опера
		<i>cin</i>	фильм
		<i>arc</i>	архитектура
		<i>art</i>	изобразительные искусства, фотография, скульптура
		<i>the</i>	театр, театральные науки, критика
		<i>lit</i>	литература
<i>hum</i>	гуманитарная наука		
		<i>his</i>	история, археология
		<i>psy</i>	психология
		<i>edu</i>	педагогика
		<i>soc</i>	социология, коммуникация
		<i>phi</i>	философия
		<i>inf</i>	библиотекведение
		<i>pol</i>	политология
		<i>lin</i>	языкознание
		<i>eth</i>	этнология, этнография



	Domain		Subdomain
		<i>cul</i>	культурология
<i>law</i>	право		
		<i>jur</i>	правосудие
<i>nat</i>	естественные науки		
		<i>agr</i>	сельское хозяйство
		<i>med</i>	медицина
		<i>zoo</i>	зоология
		<i>bot</i>	ботаника
		<i>bio</i>	биология
		<i>che</i>	химия
		<i>mat</i>	математика
		<i>ggr</i>	география
		<i>phy</i>	физика
		<i>met</i>	метеорология
		<i>geo</i>	геология
		<i>env</i>	экология
<i>tec</i>	техника		
		<i>tra</i>	транспорт, средства связи
		<i>ene</i>	энергетика
		<i>ind</i>	промышленность
		<i>com</i>	информатика
		<i>bui</i>	строительство
		<i>sta</i>	нормализация
<i>ecn</i>	экономика, управление		
		<i>eco</i>	экономика, банковское дело, торговля
		<i>mng</i>	управление, менеджмент
		<i>mer</i>	товароведение
<i>blf</i>	вера, сверхъестественное		
		<i>rel</i>	религия, вера, секты
		<i>teo</i>	теология
		<i>exc</i>	магия, сверхъестественное
<i>lif</i>	быт		
		<i>hou</i>	домашнее хозяйство
		<i>fsh</i>	мода
		<i>spo</i>	спорт, физкультура
		<i>sct</i>	общественная жизнь
		<i>amu</i>	хобби, свободное время
		<i>min</i>	национальности
		<i>reg</i>	регионы
		<i>cnl</i>	консультация
		<i>cli</i>	культура

● Medium: происхождение текста

	Medium
<i>lib</i>	книга
<i>nws</i>	газета
<i>jou</i>	журнал
<i>ste</i>	лекции
<i>net</i>	интернет и другие компьютерные сети
<i>for</i>	формуляр
<i>occ</i>	документы вышедшее по случаю
<i>npu</i>	рукопись
<i>tvf</i>	телевидение
<i>rad</i>	радио

● Authsex: пол автора

<b>Authsex</b>	
<i>msc</i>	мужской
<i>fem</i>	женский

- Lang: язык документа  
Всегда „slk“ (словацкий)

- Varieta: вариант языка

<b>Varieta</b>	
<i>std</i>	стандартный
<i>nst</i>	нестандартный (диалект, архаичный, сленг...)

- Paragraphs: сообщает сохранены-ли в тексте правильно абзацы

<i>tru</i>	да, абзацы сохранены
<i>fls</i>	нет, информация о абзацах отсутствует

- Emphasis: сообщает сохранена-ли в тексте информация о изменении начертания текста (напр. жирный, курсив, другой шрифт)

<i>tru</i>	да, информация сохранена
<i>fls</i>	нет, информация отсутствует

- Diacritics: сообщает сохранены-ли в тексте диакритические знаки

<i>tru</i>	да, диакритические знаки сохранены
<i>fls</i>	нет, диакритические знаки отсутствуют

- Transsex: пол переводчика. Смотри Authsex.

- Origlang: оригинальный язык документа. Следует норму ISO 639-2.

- Date: число издания

- Dateorig: число первого издания (возникновения документа)

- Conglomerate: множество документов к которому этот принадлежит (например один месяц газеты)

- Vogocong: сокращённое название Conglomerate.

- Comment: Произвольный комментарий

Position:	1.	2.	3.	4.	5.	6.	7.	
<b>S</b> Substantives	S A F U	<b>Paradigm</b> substantive adjective mixed incomplete	<b>Gender</b> m masc. animate i masc. inanim. f feminine n neutrum	s p o	<b>Number</b> singular plural unknown	<b>Case</b> 1 nominative 2 genitive 3 dative 4 accusative 5 vocative 6 locative 7 instrumental o unspecified		
<b>A</b> Adjectives	A F U	<b>Paradigm</b> adjective mixed incomplete	<b>Gender</b> m masc. animate i masc. inanim. f feminine n neutrum o unspecified	s p o	<b>Number</b> singular plural unknown	<b>Case</b> 1 nominative 2 genitive 3 dative 4 accusative 5 vocative 6 locative 7 instrumental o unspecified	x y z <b>Grade</b> positive comparative superlative	
<b>P</b> Pronouns	S A P F U D	<b>Paradigm</b> substantive adjective pronoun mixed incomplete adverbial	<b>Gender</b> m masc. animate i masc. inanim. f feminine n neutrum o unspecified h general	s p o	<b>Number</b> singular plural unknown	<b>Case</b> 1 nominative 2 genitive 3 dative 4 accusative 5 vocative 6 locative 7 instrumental o unspecified	g <b>Agglutinated</b> agglutinated	
<b>N</b> Numerals	S A N F U D X	<b>Paradigm</b> substantive adjective numeral mixed incomplete adverbial solitaire use	<b>Gender</b> m masc. animate i masc. inanim. f feminine n neutrum o unspecified	s p o	<b>Number</b> singular plural unknown	<b>Case</b> 1 nominative 2 genitive 3 dative 4 accusative 5 vocative 6 locative 7 instrumental o unspecified		
<b>V</b> Verbs	I K M H L B	<b>Form</b> infinitive indicative imperative transgressive /l-participle futurum	<b>Aspect</b> d perfective e imperfective j ambivalent	s p	<b>Number</b> singular plural	<b>Person</b> a first b second c third	m i f n o h <b>Gender</b> masc. animate masc. inanimate feminine neutrum unspecified general	+ - <b>Negation</b> affirmative negative
<b>G</b> Participles	k t	<b>Type</b> active passive	<b>Gender</b> m masc. animate i masc. inanim. f feminine n neutrum o unspecified	s p o	<b>Number</b> singular plural unknown	<b>Case</b> 1 nominative 2 genitive 3 dative 4 accusative 5 vocative 6 locative 7 instrumental o unspecified	x y z <b>Grade</b> positive comparative superlative	
<b>D</b> Adverbs	x y z	<b>Grade</b> positive comparative superlative						
<b>E</b> Prepositions	v u	<b>Form</b> vocalised non-vocalised	2 3 4 6 7 o <b>Binds with</b> genitive dative accusative locative instrumental unspecified					
<b>O</b> Conjunctions	Y	contains conditional morpheme <i>by</i>						
<b>T</b> Particles	Y	contains conditional morpheme <i>by</i>						

<b>J</b> Interjection	<b>#</b> Not a word
<b>R</b> Reflexive particle/pronoun <i>sa, si</i>	<b>%</b> Citation element (e.g. foreign language word)
<b>Y</b> Morpheme <i>by</i>	<b>0</b> Digits
<b>Z</b> Punctuation	
<b>W</b> Abbreviation	
<b>Q</b> Unknown POS type	

Таб.1) Полное описание морфосинтаксических помет в Словацком национальном корпусе

```

<tok>
  <orth>meč</orth>
  <disamb>
    <base>meč</base>
    <ctag>SSis1</ctag>
  </disamb>
  <lex>
    <base>meč</base>
    <ctag>SSis1</ctag>
  </lex>
  <lex>
    <base>mečať</base>
    <ctag>VMesb+</ctag>
  </lex>
</tok>

```

Фрагмент текста в корпусе в XCES формате

## Литература

1. *Metódy výskumu a opisu lexiky slovanských jazykov*. Materiály zo sympózia konaného v rámci 7. zasadnutia Lexikologicko-lexikografickej komisie pri Medzinárodnom komitáte slavistov (Nové Vozokany 24. – 26. April 1989). Edited by V. Blanár. Bratislava: Jazykovedný ústav Ľ. Štúra SAV (1990)
2. Benko, V.: Slovak Language Lexical Database. In: Computational Lexicography. Research Institute for Linguistics, Hungarian Academy of Sciences, Budapest, Hungary (1991) 5–8
3. Jarošová, A.: *Korpus textov slovenského jazyka*. In: Slovenská reč 2 (1993) 89–95
4. Šimková, M.: *Počítačové spracovanie prirodzeného jazyka a Slovenský národný korpus*. Budmerice: Počítačová podpora prekladu (2003)
5. Ide, N., Bonhome, P., Romary, L.: *XCES: An XML-based Encoding Standard for Linguistic Corpora*. In: Proceedings of the Second International Language Resources and Evaluation conference. Paris: European Language Resources Association (2000)
6. Hajič, J., Hladká, B.: *Czech Language Processing - POS Tagging*. In: Proceedings of the First International Conference on Language Resources and Evaluation. Granada,

- Spain (1998) 931–936
7. Hajič, J., Hric, J., Kuboň, V.: *Machine Translation of Very Close Languages*. In: Proceedings of the ANLP 2000. Seattle, U.S.A. (2000) 7–12
  8. Sedláček, R.: *Morfologický analyzátor češtiny*. PhD. thesis. Faculty of Informatics, Masaryk University Brno, (1999)
  9. Erjavec, T. et al.: *The MULTEXT-East Morphosyntactic Specifications for Slavic Languages*. In: Proceedings of the EACL 2003 Workshop on the Morphological Processing of Slavic Languages.
  10. <http://www.textforge.cz>
  11. Przepiórkowski, A.: *The IPI PAN Corpus preliminary version*. Instytut Podstaw Informatyki PAN, Warsaw, Poland (2004)
  12. Zakharov, V., Volkov, V.: *Morphological Tagging of Russian Texts of the XIX<sup>th</sup> Century*. In: Text, Speech and Dialogue. Proceedings of the 7<sup>th</sup> International Conference TSD 2004. Brno, Czech Republic (2004) 235–242
  13. Forróová, M., Horák, A.: *Morfologická anotácia korpusu*. In: Slovenčina na začiatku 21. storočia. Prešovská univerzita, Prešov, Slovakia (2004) 174–183
  14. Gianitsová, L.: *Dynamika skloňovania substantív (na materiáli SNK)*. In: „Tzv. základní výzkum v lingvistice – desideratum, nebo realis?“ V. mezinárodní setkání mladých lingvistů (Olomouc 17 – 19 May 2004). To be published.