

Тамара & Леннарт Лённгрен

## ПРИНЦИПЫ ЛЕММАТИЗАЦИИ СЛОВОФОРМ, НАПИСАННЫХ С ПРОПИСНОЙ БУКВЫ

В Институте славянских языков при Уппсальском университете составлен машинный фонд русских текстов объемом в 1 миллион словоупотреблений. Необходимые экономические средства для проведения работы были предоставлены шведским исследовательским органом "Humanistiska forskningsrådet". Основную работу над проектом осуществляли Леннарт Лённгрен и Людмила Ферм<sup>1</sup>. В последнее время в этой работе также принимала участие соавтор данной статьи, Тамара Лённгрен.

Материал фонда был подвергнут ряду обработок, касающихся прежде всего его словарного состава. В частности, была проведена "лемматизация", то есть приведение каждой словоформы к соответствующей "лемме". Под термином "лемма" в данном случае мы подразумеваем словарную (исходную) форму данной словоформы. Такой анализ всего лексического материала — в фонде содержится свыше 132.000 словоформ — очень трудоемкий, так как он только в ограниченной степени подлежит автоматизации. Однако лемматизация является необходимым условием для получения двух основных результатов настоящего проекта: конкорданса и частотного словаря, представленных совместно на уровне словоформы и леммы.

Лемматизация большого лексического материала является не только трудоемкой механической работой, но и ставит много сложных лингвистических задач, требующих нетривиальных решений. Одной из таких задач является лемматизация словоформ, написанных

---

<sup>1</sup>См. *Slovo* 36, с. 77-83 и 39, с. 125-141.

с большой буквы. Здесь мы не могли, как во многих других случаях обращаться к толковым словарям, в которых, как известно, сло этой категории вообще не приводятся.

В обычном представлении, слова, написанные с большой буквы это прежде всего имена собственные. Тогда возможность возникновения проблем при лемматизации имени собственного может показаться по меньшей мере странной: лемма словоформы *Москвой* будет только *Москва*, а *Ольге* — только *Ольга*. На первый взгляд — все довольно понятно, но это только на первый взгляд...

Не углубляясь в проблему имени собственного — объекта исследования самостоятельной лингвистической науки ономастики, в котором существует огромное количество литературы, — остановимся на предельно лаконичном его определении в "Русской грамматике"<sup>1</sup> "собственные имена в узком смысле этого слова — это географические и астрономические названия и имена людей и животных"<sup>2</sup>. Но ведь прописной буквы пишутся и многие другие названия, не являющиеся именами собственными в узком смысле, и для которых существует довольно пестрая терминология, например, таксоны, номены, хремотонимы, прагматонимы, документонимы, эргонимы, фалеронимы<sup>3</sup>.

В многословных наименованиях с прописной буквы пишется, как правило, первое слово, например *Тихий океан* (исключение составляют случаи управления, типа *проспект Мира*, а также единичные случаи согласования, ср. *юбилейные Игры*). Написание остальных нарицательных имен в составе названия регулируется определенными правилами правописания<sup>4</sup>, ср. *Золотой Рог* (бухта), но *Черное море*. Не

<sup>2</sup>Русская грамматика, М., 1980, с. 480.

<sup>3</sup>А. В. Суперанская, *Аппеллятив — онома, Имя нарицательное и собственное* (под ред. А. В. Суперанской), М., 1978, с. 5-33.

<sup>4</sup>Д. Э. Розенталь, *Вопросы русского правописания*, М., 1970, с. 26 и сл.

редко одно название входит в состав другого, ср. *Институт русского языка Академии Наук СССР*.

Что касается имен собственных в узком смысле, то они до настоящего времени в словарной практике лемматизации не подвергались. При составлении имеющегося русского частотного словаря под редакцией Л. Н. Засориной (ЧС)<sup>5</sup> собственные имена (очевидно, авторы отождествляют "имя собственное" и "антропоним") просто не выписывались из текста, не фиксировались клички животных и т. д. (ЧС, с. 14). В Упсальском машинном фонде принято другое решение — ввести все без исключения словоформы, содержащихся в подобранных текстах. И здесь мы полностью разделяем мнение А. В. Суперанской о том, что "Собственные имена несправедливо игнорируются многими исследователями как маловажные, не занимающие должного места в лексическом составе языка. Однако, как бы ни был беден словарный запас человека, собственные имена неизбежно присутствуют в нем, а некоторые наиболее распространенные имена употребляются едва ли не чаще, чем многие нарицательные."<sup>6</sup>

Таким образом, в процессе обработки большого материала, положенного в основу создаваемого частотного словаря, вполне закономерно встал вопрос о выработке соответствующих принципов лемматизации словоформ, написанных с большой буквы. Здесь важно отметить, что мы не имеем в виду все слова, написанные с большой буквы в самих текстах. Уже при вводе текстового материала различались "синтаксическое" употребление прописной буквы, то есть в начале предложения, в начале цитаты и т. д., и "лексическое" ее употребление, в частности, в именах собственных. Только в последнем случае

---

<sup>5</sup>Частотный словарь русского языка (под ред. Л. Н. Засориной), М., 1977 (далее: ЧС).

<sup>6</sup>А. В. Суперанская, *Ударение в собственных именах в современном русском языке*. М., 1966, с. 24.

большая буква является признаком данной словоформы как лексической единицы. При дальнейшей обработке словарного состава фон было важно, чтобы большая буква "сохранялась" только в последнем случае. Это осуществлялось таким образом: специальным знаком были отмечены все словоформы последнего — "лексического" — типа. В остальных словоформах информация о прописной букве не сохранялась.

Однако и эта работа не была полностью лишена проблем. Прописные буквы сохранялись в основном только там, где данное название обладает хотя бы каким-то постоянством (хотя, может быть, полная последовательность здесь не была достигнута). Кроме того, сохранялись прописные буквы только в словоформах знаменательной части речи (так, например, не была обозначена звездочкой (\*) буква "Н" в случае: картина "На реке").

Самым простым решением проблемы было бы оставить прописную букву в соответствующих леммах. Но это привело бы к нежелательным результатам особенно в таких случаях, когда одна и та же словоформа выступает в двух вариантах: со строчной и с прописной начальной буквой. На уровне леммы мы бы получили, например, два разные "академии": *академия* и *Академия* (ср. *учиться в академии Академия Наук СССР*), два разных "тихий": *тихий* и *Тихий* (ср. *тихий день : Тихий океан*) и т. д. В результате произошло бы лингвистически неоправданное расщепление частотности подобных словоформ по двум леммам, что, в свою очередь, привело бы к получению неверных статистических данных. Неправомерность такой трактовки словоформ применительно к словосочетаниям, являющимся названиями, признавалась, по-видимому, и авторами ЧС, так как, например, *Черное море* они анализируют как одно употребление прилагательного *черный* и одно употребление существительного *море* (ЧС, с. 15). Чтобы избежать каких-либо неточностей в статистических подсчетах:

мы выработали определенные принципы лемматизации тех 11.677 словоформ, в которых сохранилась начальная прописная буква.

Дифференцированная трактовка словоформ, написанных с прописной буквы, обусловлена и такими соображениями: лексическая обработка машинного фонда сводится к ступенчатому обобщению языковых единиц, к их редукции. Первый шаг обобщения — это переход от текстов к списку словоформ (token → type), второй — от словоформы к лексеме. Нам представлялось естественным, чтобы этот шаг включал некоторую редукцию и в области правописания.

Таким образом, "судьба" прописной буквы определяется в нашей системе на трех различных уровнях: на уровне текста, на уровне словоформы и на уровне леммы. Одни прописные буквы существуют только в самих текстах (случай синтаксического употребления), вторые сохраняются как лексическая характеристика словоформы, но потом "сливаются" со строчными (*Тихого [океана]* → *тихий*), третьи же сохраняются вплоть до последнего уровня.

Теоретически не исключена возможность "обратной" лемматизации, то есть от словоформы со строчной буквой к лемме с прописной (например, в случае *рождество*, где наблюдается колебание, можно считать написание с прописной буквой ближе к новой складывающейся норме). Но, во избежание сложностей, мы отказались от этого направления лемматизации.

Словоформы с прописной буквой разделяются, в первую очередь, на две группы: слова, которые никогда не пишутся со строчной буквы, и слова, для которых мыслимы оба написания. Слова первого типа, конечно, никогда не приводятся к леммам со строчной буквой, независимо от структуры данного названия, ср. *Азия*, *Средняя Азия*. В связи с этим следует обратить внимание на своеобразную омографию типа *Дзержинский* (фамилия) : *Дзержинский район*, где по форме совпадают имя существительное и имя прилагательное; ср. также

*Петропавловск-Камчатский* (город), но: *Петропавловск-Камчатск* училище. В таких случаях прописная буква сохраняется только леммах имен существительных.

Для лемматизации словоформ второго типа — с двояким возможным написанием — очень существенно разделение названий на однословные (простые) и многословные (составные). Последние, как правило, содержат по крайней мере одну словоформу, которая в другом контексте пишется со строчной буквы. В случае однословных названий границы слова и названия совпадают. Если данное слово по происхождению имя нарицательное или принадлежит к другой части речи, чем существительное, оно неизбежно значительно преобразуется. Имя нарицательное превращается в имя собственное и теряет способность склоняться по категории числа, например (река) *Канаэ* (деревня) *Кулички*. Имя прилагательное, кроме того, переходит другую часть речи, в класс существительных, и в связи с этим теряет способность склоняться по родам, например (поселок) *Дружно*. Исходя из этих соображений мы решили привести все однословные названия к леммам с прописной буквой.

Исключением из названного принципа является более или менее окказиональное употребление прописной буквы в однословных "наименованиях". Известно, что в особом стилистическом употреблении с прописной буквы пишутся слова типа *Родина*, *Отчизна*, *Человек*; ср. также следующие контексты: "принадлежность Себе и Делу"; "полюбить Тень, Звук имени..."; "услыхала Его голос". Все такие случаи приведены к леммам с соответствующей строчной буквой. О окказиональных личных именах можно говорить в случаях типа *Папаня*, *Маманя*, а об окказиональном топониме — в случае "*Где Лошади Плачут*" (Ф. Искандер: "Широколобый"); ср. также: "... сам Россия на перегоне от станции Вчера до станции Завтра" (Б. Васильев: "Вы чье, старичье?"). В леммах подобных слов прописная буква

тоже не сохраняется. Такое решение оправдано тем, что "собственность" этих слов сохраняется только в определенном контексте, при особой стилистической обусловленности.

Еще одним исключением являются случаи, где наблюдается колебание. Например, не сразу было найдено однозначное решение лемматизации таких слов как *Запад*, *Восток*, (в смысле "западные, восточные страны"), *Север*, *Юг* (в смысле "территория"). С одной стороны, было бы закономерным оставить в леммах их словоформ прописную букву, но в связи с тем, что в текстовых файлах встречается двойное написание этих слов в одном и том же значении, они были лемматизированы со строчной буквы (но с большой буквой остаются, конечно, употребления типа *папирсы "Север"*).

Разделение между однословными и многословными названиями было проведено не случайно. В многословных названиях, типа *Тихий океан*, *Средняя Азия*, *Набережные Челны*, входящие в состав данных словосочетаний нарицательные имена сами по себе не переходят в разряд имен собственных. В частности, они сохраняют свои морфологические и синтаксические свойства. Поэтому такие словоформы, написанные с прописной буквы, приводятся к леммам с соответствующей строчной буквой.

В результате вышеупомянутой "фильтрации" уже на уровне текстов из многословных названий остались, в основном, только словосочетания, состоящие из независимого существительного и разного рода определений. Названия другой синтаксической структуры — речь идет, прежде всего, о названиях произведений искусства, типа *картина "На реке"* — не отмечались специально в текстовых файлах и тем самым не сохранялись. Исключения составляют единичные случаи типа журнал *"Знание — сила"*, журнал *"Информатика и образование"*.

Приняв упомянутые принципы лемматизации, мы столкнулись с

особой трудностью: не всегда ясно, является ли данное название словосочетанием, или только составной частью словосочетания. Другими словами, иногда бывает трудно определить линейные (синтаксические) границы названия, что, разумеется, необходимо для правильной его классификации по отношению к признаку однословности/многословности. Какими критериями руководствоваться в случаях типа *курган Садовый* и *Садовая улица*? *проспект Мира* и *улица Магросская Тишина*? Что следует включать в данное название и что — нет?

Для решения этой проблемы необходимо остановиться на вопросе о приложении, в роли которого выступает название, и о видах синтаксической связи в таких конструкциях. В данном случае приложение — это определяющий член словосочетания, а в качестве определяемого члена выступает нарицательное слово (редко — словосочетание). Между определяемым словом и приложением существует два вида синтаксической связи: согласование или примыкание, причем постпозитивную позицию всегда занимает приложение, например, *река Нева* (ср. *на реке Нева*), *река Дон*, (ср. *на реке Дон*). Следует обратить внимание на принципиальное различие между внешне похожими словосочетаниями *река Быстрая* и *улица Детская*. В последнем случае *улица* — это составная часть многословного названия. Порядок слов здесь является вторичным, маркированным, и возникает в контекстах типа *Я живу не на Садовой улице, а на улице Детской*.

В приведенных случаях всегда возможно опустить определяемое слово, ср. *на озере Байкал* — *на Байкале*. Если же данное имя нарицательное присоединяется с последующим зависимым именем собственным с помощью управления, то ни о каком приложении не может быть речи, например *проспект Мира*. Здесь в принципе невозможно опустить независимое слово (только в разговорной речи можно услышать конструкции типа *Я живу на Мира*), так что границы словосочетания

и названия совпадают.

Изложенное понимание приложения несколько расходится со мнениями советских грамматистов. Так, в "Краткой русской грамматике"<sup>7</sup> приложение характеризуется как вид согласования (с. 352), но, по-видимому, одновременно воспринимается как определяющее слово, так как говорится о "приложениях-топонимах" (с. 354). Отмечается, что случаи, где нет согласования, "по существу выходят за рамки приложения" (с. 354). В. А. Белошапкина<sup>8</sup>, не называя подобные случаи приложением, совершенно справедливо указывает на то, что связь примыкания здесь выражается "контактным постпозитивным расположением примыкающего компонента" (с. 54).

Ниже проиллюстрируем указанные принципы лемматизации на примерах, взятых из машинного фонда. Сначала рассмотрим однословные названия, потом — многословные. В обоих случаях данные названия могут быть или "самостоятельными", или выступать в роли приложения, и в соответствии с этим проводится главное разделение. Многословные названия, кроме того, разделяются по синтаксической связи. В случае "вложения" одного названия в другое данный пример классифицируется по связи между наиболее крупными составными.

### 1.1 Однословные самостоятельные названия

а) антропонимы (имена, отчества, фамилии, прозвища): Ирина, Искра, Зубр, Капа, Машка, Петька; Кузьмич, Петровна, Данилыч; Буренков, Гончаров, Тимофеев-Ресовский, Петров-Воткин; Стамеска, Рыжий, Пипка, Светка-Пипетка;

б) клички животных: Матера, Канитель, Купчиха, Красуля, Лохма-

<sup>7</sup>Краткая русская грамматика (под ред. Н. Ю. Шведовой & В. В. Лопатина), М., 1989.

<sup>8</sup>В. А. Белошапкина, Современный русский язык. Синтаксис, М., 1977.

тый;

в) астрономические названия: Сириус, Венера;

г) топонимы: Байкал, Крым, Серебрянка (река), Швеция, Минь Кислицы, Грязное (два последних – населенные пункты), Невск (проспект /эллипсис/);

д) административные названия, организации, предприятия: Интури, Аэрофлот, Госкомиздат, Минстанкопром, КамАЗ, СССР, НИИ;

е) прочие названия: "Известия" (газета), "Мир" (спутник), "Север" (холодильник).

## 1.2 Однословные названия-приложения

а) антропонимы: старуха Дарья, директор Кутаков;

б) клички животных: корова Поляна, собака Тузик, кот Рыжко;

в) астрономические названия: планета Марс, звезда Солнце;

г) топонимы: река Волга, река Лазурь, курган Садовый, полуостров Рыбачий, озеро Лама, озеро Великое, высота "Плоская", роща "Светлая", канал Ростов, город Бахмач, город Озеры, город Улан-Удэ, деревня Кулички, деревня Степановская, село Крылос, станция "Минский", метро "Сокол";

д) названия предприятий: совхоз "Комсомолец", совхоз "Погорелоский", объединение "Швейтекстильтрикотаж", комплекс "Восточный", трест "Железобетон";

е) прочие названия: кинотеатр "Чайка", пансионат "Ласточка", общество "Искусство", журнал "Природа", спутник "Экран" (ср. также, определяемым словом в множественном числе: спутники "Метеор").

### 2.1.1 Многословные самостоятельные названия с согласованием

а) антропонимы (имя + отчество + фамилия; прозвища): Римма Львовна Львовская; Сухарь Сухарыч, Нюська Светлячок (с точки зрения современного языка здесь вряд ли можно говорить о словосочетании).

четании, так как данная синтаксическая связь скорее сочинительного, чем подчинительного характера); Володька Ржавый, Александр Македонский, Александр Невский (с несомненной подчинительной связью; в последних двух случаях прописная буква словоформ-прилагательных сохраняется и в соответствующих леммах, так как эти определения воспринимаются как имена собственные);

б) астрономические названия: Млечный путь, Большой Пес;

в) топонимы: Лисий Нос, Кривой Рог, Набережные Челны, Северная бухта, Азовское море, Кольский полуостров, Серебряный бор, Кузякин лог, Садовое кольцо, Малая Садовая (улица /эллипсис/);

г) административные названия, предприятия: Таджикская ССР, Ясногорский район, Дзержинский район, Усть-Лабинский район, Краснодарский край, Херсонская область; Октябрьский народный суд, Копейский горком партии, Калевальский райком КПСС, Воронежский облисполком, Селигеровский сельсовет; Светлогорский целлюлозно-бумажный комбинат, Селенгинская фабрика, Чебоксарская АЭС; Советские Вооруженные Силы;

д) прочие названия: Варшавский Договор, Основной Закон; Генеральный Секретарь, Верховный Главнокомандующий, Нобелевские лауреаты; Золотая медаль; "Литературная газета", "Красная Звезда" (газета); юбилейные Игры.

### 2.1.2 Многословные самостоятельные названия с управлением

а) астрономические названия: созвездие Дракона, альфа Большого Пса;

б) топонимы: море Лаптевых (лемматизируется: море, Лаптев), проспект Мира, площадь Революции, улица Свободы;

в) административные названия: Союз Советских Социалистических Республик, Страна Советов, Совет Министров СССР, Верховный Совет СССР, Верховный Суд СССР, Ревизионная Комиссия КПСС,

Комитет Мира, Институт имени И. В. Курчатова, коммуна имени Ф Дзержинского (в двух последних примерах имеет место "вложен одного словосочетания в другое, причем подчиненное словосочета состоит из названия-приложения и определяемого слова *имя*);

г) прочие названия: Итоговый документ Венской встречи, Заключительный акт общеевропейского совещания, Соборное уложение ца Закон Эстонской ССР, Февральский Пленум ЦК КПСС; Председат Совета Министров, Герой Социалистического Труда; Орден Крас Звезды; Игры I Олимпиады, велогонка Мира; могила Неизвестн Солдата.

### 2.2.1 Многословные названия-приложения с согласованием

а) антропонимы: школьник Сева Гребин, некий старик Евсей Быков;

б) топонимы: улица Матросская Тишина, улица Первый Смоленск Ручей, деревня Белый Бор, поселок Старые Ляды;

в) названия предприятий: колхоз "Ленинский путь", объединен "Краснодарский чай";

г) прочие названия: газета "Советский спорт", журнал "Всемирн следопыт".

### 2.2.2 Многословные названия-приложения с управлением

а) топонимы: село Годы Турка;

б) названия предприятий, прочие названия: совхоз "Заветы Ильича газета "Огни Ангары", журнал "Вопросы философии".

Своеобразно ведут себя слова типа *имя*, *название*, которые мог выступать в качестве определяемого слова приложения: *имя Паве название Тверь*, или же просто подчинять себе название, не явля щееся приложением: *имя Сталина, название Солодовой улицы*. первом случае имеет место примыкание или согласование (ср. *назв*

ние "Разбойниково", но: фамилия Серпилин — фамилию Серпилина), во втором — управление. Словосочетания этого типа — то же касается словосочетания с определяемым словом "слово" (ср. слово Тверь, слово "Канны") — ввиду их метаязыкового характера не были включены в изложенную классификацию.