

УДК 004.624

## АНАЛИЗ ТИПИЧНЫХ ОШИБОК ПРИ ОБМЕНЕ ДАННЫМИ МЕЖДУ СИСТЕМАМИ

**С. В. ЛУКАШЕВИЧ, В. Л. ШАБАЛОВ, В. В. ЛУКАШЕВИЧ**  
(Представлено: А. А. СКУКОВСКАЯ)

*В статье рассматриваются типичные ошибки, возникающие при передаче данных между различными программными системами. Особое внимание уделяется проблемам, связанным с кодировками текста, выбором разделителей полей, интерпретацией дат и представлением чисел с плавающей точкой. Анализируются причины возникновения ошибок и предлагаются практические рекомендации по их предотвращению.*

**Ключевые слова:** обмен данными, совместимость форматов, кодировка, CSV, дата, числа с плавающей точкой, межсистемное взаимодействие.

**Введение.** В условиях активного развития цифровых технологий обмен данными между различными системами – от табличных редакторов до облачных баз данных – стал повседневной практикой. Однако, несмотря на стандартизацию многих форматов, на практике часто возникают ошибки, приводящие к искажению или потере информации. Особенно уязвимыми оказываются текстовые данные, содержащие специальные символы, даты и числа, записанные с учётом локальных особенностей. В данной работе анализируются наиболее распространённые источники ошибок и предлагаются пути их минимизации.

**Основная часть.** Одной из наиболее распространённых и коварных проблем при межсистемном обмене данными является некорректная обработка текстовых кодировок. Кодировка представляет собой соглашение о том, как последовательность байтов в цифровом файле интерпретируется как человекочитаемые символы – буквы, цифры, знаки препинания, специальные символы и управляющие последовательности. При несовпадении кодировок между отправляющей и принимающей системами текст может не только искажаться, но и полностью утратить смысл, что приводит к нарушению целостности данных, ошибкам в бизнес-логике и даже сбоям в работе программного обеспечения.

Наиболее ярко эта проблема проявляется при работе с CSV-файлами, экспортируемыми из табличных редакторов, таких как Microsoft Excel. Несмотря на кажущуюся простоту формата, CSV не имеет строгого стандарта, регулирующего кодировку текста. В результате поведение систем при экспорте и импорте может сильно различаться. Например, на русскоязычных версиях операционной системы Windows табличный редактор Excel по умолчанию сохраняет CSV-файлы в кодировке Windows-1251 (также известной как CP1251), которая была разработана Microsoft для поддержки кириллических символов в среде Windows. Эта кодировка эффективна в локальном контексте, однако она не является универсальной.

В то же время большинство современных программных инструментов – включая языки программирования (например, Python, JavaScript, Java), веб-серверы, облачные сервисы и реляционные базы данных – по умолчанию работают с кодировкой UTF-8. Это международный стандарт, закреплённый в спецификациях IETF и W3C, который поддерживает символы практически всех письменностей мира, включая эмодзи и специальные технические символы. UTF-8 также является основной кодировкой для веба: согласно данным W3Techs, более 98 % веб-сайтов используют именно её.

Если CSV-файл, сохранённый в кодировке Windows-1251, будет прочитан как UTF-8, кириллические символы будут интерпретированы неверно. Например, слово «Текст» превратится в последовательность вида РўРµРєСЃС.. Такое искажение возникает потому, что каждый байт (или пара байтов) в исходной кодировке интерпретируется как совершенно другой символ в целевой кодировке. В данном случае символ «Т» (код 0xD2 в CP1251) читается как два байта UTF-8, соответствующих символам «Р» и «ў», и так далее по цепочке. Подобные ошибки особенно опасны, поскольку внешне файл выглядит корректно – текст присутствует, но его содержание бессмысленно.

Ещё одна сложность связана с наличием или отсутствием BOM (Byte Order Mark) – специальной метки в начале текстового файла, указывающей на используемую кодировку. Формат UTF-8 не требует BOM, и многие Unix- и Linux-системы считают его избыточным или даже вредным. Однако Microsoft Excel при сохранении CSV в UTF-8 автоматически добавляет BOM (последовательность байтов EF BB BF). В результате, если такой файл будет обработан парсером, не ожидающим BOM (например, стандартным модулем csv в Python без указания кодировки utf-8-sig), первое поле первой строки может содержать «мусор» в виде невидимого символа или даже отображаться как й». Имя. Это особенно критично при автоматической обработке данных, где имена столбцов используются как ключи.

Помимо различий между Windows-1251 и UTF-8, существуют и другие кодировки, применявшиеся в разные эпохи и в разных экосистемах. Например, KOI8-R широко использовалась в раннем русскоязычном интернете и Unix-системах, ISO-8859-5 – в международных стандартах, а MacCyrillic – на компьютерах Apple до перехода на Unicode. Хотя сегодня эти кодировки встречаются редко, они всё ещё

могут присутствовать в архивных данных, устаревших информационных системах или при миграции данных из легаси-систем. Отсутствие единства в выборе кодировки между участниками обмена многократно увеличивает риск ошибок, особенно в условиях глобализированных проектов, где данные генерируются в разных странах и обрабатываются разнородными системами.

Второй распространённой проблемой при обмене табличными данными является выбор разделителя полей. В формате CSV (Comma-Separated Values) предполагается использование запятой (,) в качестве разделителя. Однако на практике это правило часто нарушается из-за региональных настроек операционной системы. В англоязычных странах действительно используется запятая, тогда как в странах СНГ, Германии, Франции и других европейских государствах Excel автоматически выбирает точку с запятой (;) как разделитель, поскольку запятая там традиционно применяется как десятичный разделитель в числах. Если принимающая система (например, скрипт на Python или веб-сервис) не учитывает эту особенность и пытается разобрать файл с разделителем; как стандартный CSV с ,, поля будут неправильно разграничены. Это приведёт к смешению данных: значение из столбца «Email» может попасть в поле «Телефон», а числовые данные – интерпретироваться как текст.

Особую сложность представляет работа с датами. Формат записи даты сильно зависит от региональных и культурных норм. В США принят формат MM/DD/YYYY, в большинстве европейских стран – DD/MM/YYYY, а в Японии и Китае – YYYY/MM/DD. При этом технические системы, особенно веб-API и базы данных, предпочитают использовать международный стандарт ISO 8601 – YYYY-MM-DD (например, 2024-03-04). При передаче даты 03/04/2024 из одной системы в другую возможна двусмысленность: для американца это 3 апреля, для европейца – 4 марта. В худшем случае дата может быть вообще проигнорирована, преобразована в строку или интерпретирована как число. Например, Excel хранит даты как количество дней, прошедших с 1 января 1900 года (или 1904 года на macOS), поэтому дата может превратиться в число вроде 45352, что сделает её бессмысленной при импорте в другую систему без дополнительной обработки.

Не менее критичной является проблема представления чисел с плавающей точкой. В русскоязычных и многих европейских локациях десятичным разделителем служит запятая (3,14), тогда как в англоязычной традиции и во всех основных форматах обмена данными (JSON, XML, SQL, YAML) требуется точка (3.14). При импорте числа 2,718 в систему, ожидающую точку, оно может быть интерпретировано как строка, как два отдельных числа (2 и 718), или вызвать ошибку парсинга. Это особенно опасно в финансовых, научных и инженерных приложениях, где даже небольшая ошибка в интерпретации числа может привести к катастрофическим последствиям.

Для решения указанных проблем необходимо придерживаться следующих рекомендаций:

Всегда явно указывать кодировку при сохранении и чтении текстовых файлов. Для CSV-файлов, содержащих не-ASCII символы, предпочтительно использовать UTF-8 с BOM при работе с Excel и UTF-8 без BOM при интеграции с веб-сервисами.

Устанавливать единые соглашения по формату данных между всеми участниками обмена – в виде технических спецификаций, контрактов API или схем данных.

Применять международные стандарты: ISO 8601 для дат и времени, точку как десятичный разделитель, RFC 4180 для CSV.

Валидировать данные на каждом этапе передачи с помощью строгих парсеров, схем (например, JSON Schema, XML Schema) или библиотек с поддержкой типизации (например, Pydantic в Python).

При работе с CSV указывать разделитель явно, например, используя параметр `sep=';'` в библиотеке pandas или `delimiter=';'` в модуле csv языка Python.

Таким образом, проблемы совместимости форматов данных – особенно в части кодировок, разделителей, дат и чисел – требуют не только технического понимания, но и дисциплинированного подхода к проектированию межсистемного взаимодействия. Только при соблюдении единых стандартов и тщательной валидации можно обеспечить надёжную, точную и безопасную передачу информации в современных распределённых системах.

**Заключение.** Типичные ошибки при обмене данными – не просто технические нюансы, а серьёзные риски для целостности и достоверности информации. Кодировки, разделители, даты и числа требуют особого внимания при проектировании межсистемного взаимодействия. Соблюдение единых стандартов, явное указание параметров форматов и тщательная валидация данных позволяют значительно снизить вероятность ошибок и обеспечить надёжную передачу информации между разнородными системами.

## ЛИТЕРАТУРА

1. McKinney, W. Python for Data Analysis. – O'Reilly Media, 2022.
2. Bray, T. (Ed.). The JavaScript Object Notation (JSON) Data Interchange Format. – RFC 8259, 2017.
3. Microsoft. Excel Specifications and Limits. – Документация Microsoft, 2023.
4. PostgreSQL Global Development Group. PostgreSQL Documentation. – <https://www.postgresql.org/docs/>