

DOI: 10.51790/2712-9942-2023-4-2-07

ЕДИНЫЙ ГОСУДАРСТВЕННЫЙ ФОНД ГИДРОМЕТЕОРОЛОГИЧЕСКИХ ДАННЫХ КАК БОЛЬШИЕ ДАННЫЕ. УНИВЕРСАЛЬНЫЙ ПАРСЕР СТРУКТУРЫ ДАННЫХ В ФОРМАТЕ ЯЗЫКА ОПИСАНИЯ ГИДРОМЕТЕОРОЛОГИЧЕСКИХ ДАННЫХ

Л. О. Перетяtko

Федеральное государственное бюджетное учреждение «Всероссийский научно-исследовательский институт гидрометеорологической информации — Мировой центр данных», г. Обнинск,

Российская Федерация

ORCID: <https://orcid.org/0000-0002-9182-2952>, ✉ peretyatkol@meteo.ru

Аннотация: данная статья относится к серии статей, посвященной работе с данными, размещёнными в Едином государственном фонде данных Федеральной службы России по гидрометеорологии и мониторингу окружающей среды (Росгидромет), а именно, технологиям и инструментам для работы с этими данными как с большими данными.

Статья посвящена универсальному парсеру структуры данных в формате ЯОД (языка описания гидрометеорологических данных). Структура данных в формате ЯОД также называется ЯОД описанием. Разработанный парсер осуществляет анализ, разбор и сохранение структуры ЯОД описания для последующей конвертации данных в этом формате. Для реализации парсера была разработана система классов на основе описания формата ЯОД. Результатом работы парсера является структура ЯОД в виде набора взаимосвязанных объектов разработанных классов для обеспечения возможности чтения или записи данных в формате ЯОД. Универсальный парсер является ключевой составляющей комплекса программных средств взаимной конвертации данных в формате ЯОД в современные распространённые форматы.

В работе представлены диаграмма классов в нотации UML и их описание, особенности универсального парсера, подробное описание этапов его работы.

Универсальный парсер ЯОД описания разработан на базе Всероссийского научно-исследовательского института гидрометеорологической информации — Мирового центра данных.

Ключевые слова: язык описания гидрометеорологических данных, Единый государственный фонд данных, ЕГФД, парсер, технология обработки данных, иерархическая структура данных.

Для цитирования: Перетяtko Л. О. Единый государственный фонд гидрометеорологических данных как большие данные. Универсальный парсер структуры данных в формате языка описания гидрометеорологических данных. *Успехи кибернетики*. 2023;4(2):47–52. DOI: 10.51790/2712-9942-2023-4-2-07.

Поступила в редакцию: 10.04.2023.

В окончательном варианте: 10.04.2023.

THE NATIONAL HYDROMETEOROLOGICAL ARCHIVE AS BIG DATA. A UNIVERSAL PARSER OF THE HYDROMETEOROLOGICAL DATA DESCRIPTION LANGUAGE

L. O. Peretyatko

All-Russia Research Institute of Hydrometeorological Information, World Data Center, Obninsk, Russian Federation

ORCID: <https://orcid.org/0000-0002-9182-2952>, ✉ peretyatkol@meteo.ru

Abstract: this paper refers to a series of papers on handling the big data stored in the National Hydrometeorological Archive (NHA) managed by the Russian Federal Service for Hydrometeorology and Environmental Monitoring.

This study presents a universal parser of the data structure in the HDDL format (hydrometeorological data description language). The parser analyzes, parses, and saves the HDDL document structure for subsequent conversion. We developed a number of classes to handle the HDDL format. The parser produces a data source structure as a set of interconnected objects of the classes to read or write data in the HDDL format. The universal parser is a key component of the software package for HDDL data conversion into other common formats.

The paper presents an UML diagram of the classes and their description, features of the universal parser, and a detailed description of its operation.

The universal parser is developed at the All-Russia Research Institute of Hydrometeorological Information, World Data Center.

Keywords: hydrometeorological data description language, National Hydrometeorological Archive, NHA, parser, data processing, hierarchical data structure.

Cite this article: Peretyatko L. O. The National Hydrometeorological Archive as Big Data. A Universal Parser of the Hydrometeorological Data Description Language. *Russian Journal of Cybernetics*. 2023;4(2):47–52. DOI: 10.51790/2712-9942-2023-4-2-07.

Original article submitted: 10.04.2023.

Revision submitted: 10.04.2023.

Введение

Данная статья относится к серии статей, посвященной работе с данными, размещёнными в Едином государственном фонде данных Федеральной службы России по гидрометеорологии и мониторингу окружающей среды (Росгидромет), а именно, технологиям и инструментам для работы с этими данными как с большими данными [1].

Электронные данные Единого государственного фонда данных о состоянии окружающей среды, ее загрязнении (ЕГФД) хранятся в специальном формате ЯОД (язык описания гидрометеорологических данных). Особенностью этого формата является возможность описывать иерархическую структуру данных, получаемых от наблюдательной сети Росгидромет. Подробное описание формата ЯОД представлено в соответствующей статье [2].

Универсальный парсер структуры данных в формате ЯОД (ЯОД описания) реализован на языке программирования C++ [3], с использованием фреймворка Qt версии 5 [4]. Разработанный парсер является ключевой компонентой системы взаимной конвертации данных различной структуры [5].

Разработана система классов для парсера ЯОД описания, представленная в виде UML диаграммы [6] (см. рис 1).

Класс `StringMassive` представляет отдельные сегменты ЯОД описания, описывающие отдельные части записей или групп. Хранит текстовое описание сегмента, имя вышестоящего по иерархии сегмента и уровень вложенности.

Класс `Element` представляет элементы ЯОД описания. Хранит информацию о элементе: имя, тип, формат хранения, формат преобразования, список допустимых значений, комментариев и прочую вспомогательную информацию.

Класс `Group` представляет группу элементов ЯОД описания. Хранит список элементов, тип, имя, имя вышестоящей по иерархии группы, уровень вложенности, количество экземпляров, описание группы, а также другую вспомогательную информацию.

Класс `Record` представляет записи ЯОД описания. Хранит информацию о записи: имя, список целых групп, список фрагментов групп, номер и список отдельных сегментов описания записи.

Класс `Core` представляет ЯОД описание. Реализует методы для получения и сохранения ЯОД описания и его использования в последующей конвертации. Хранит список записей, заголовков каждой записи, название семейства ЯОД описания, тип данных в формате ЯОД, определяемый по ЯОД описанию.

Особенности парсера

Разработанный парсер ЯОД описания обладает следующими особенностями:

- универсальность — работа парсера была успешно протестирована на примере шестидесяти различных ЯОД описаний метеорологических, аэрологических, океанографических и гидрологических данных;
- использование регулярных выражений [7] для упрощения анализа и разбора на составляющие строк ЯОД описания;
- использование рекурсивных алгоритмов;
- сохранение особенностей иерархической структуры ЯОД, которые учитываются при дальнейшей конвертации;
- сохранение текста ЯОД описания в качестве метаданных.

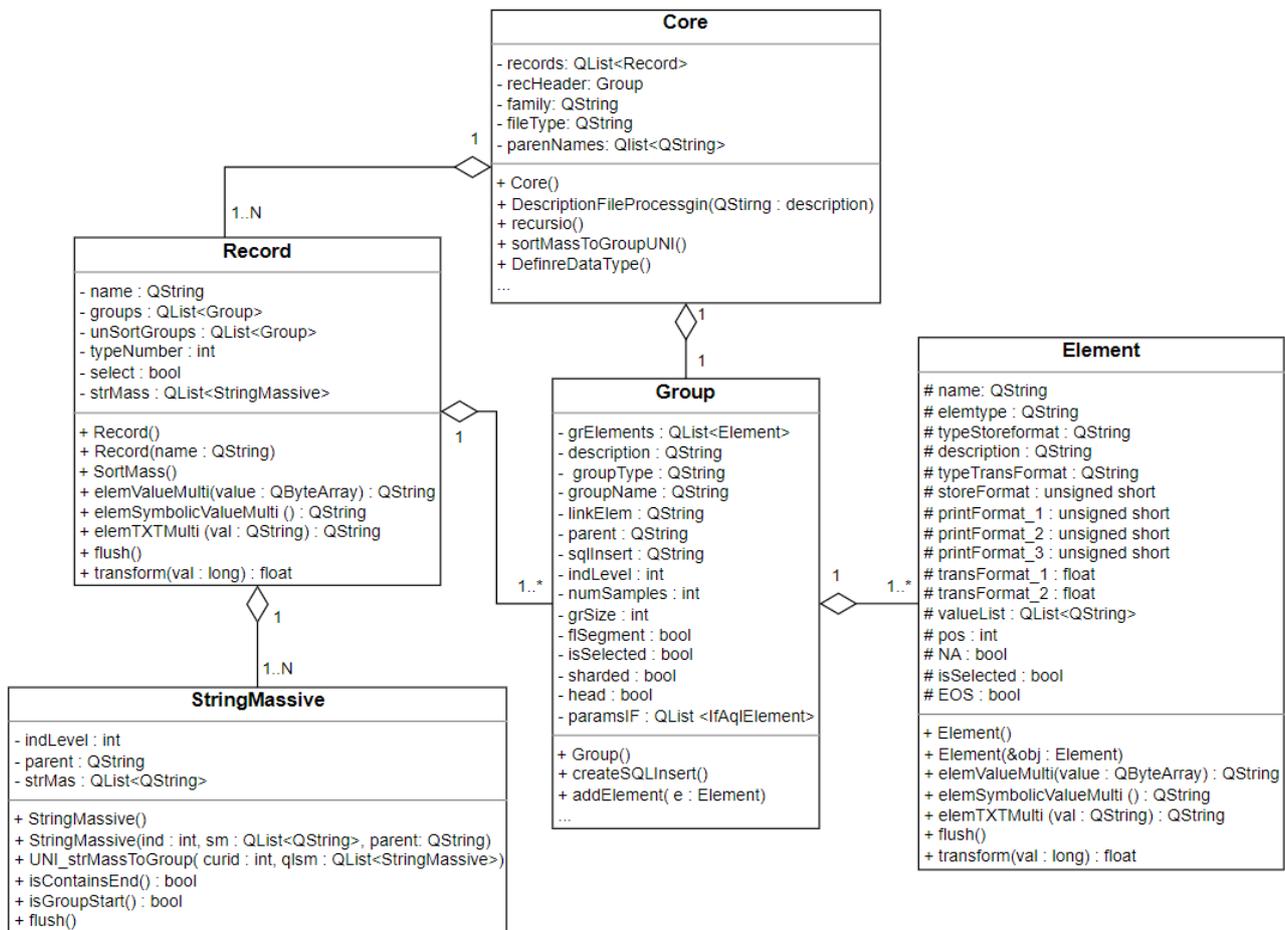


Рис. 1. UML диаграмма классов парсера ЯОД описания

Алгоритм работы парсера

Файл с ЯОД описанием считывается построчно. После чтения каждой строки из неё удаляются лишние пробелы в начале и в конце, множественные пробелы или табуляции заменяются на одинарные пробелы.

Работу парсера ЯОД описания можно разделить на четыре основных этапа.

На первом этапе ведётся поиск строки с ключевым словом FAMILY, в которой выделяется и сохраняется имя семейства файлов, указанного после ключевого слова (SI — имя семейства, см. рис 2). Затем выполняется поиск одного из двух возможных ключевых слов — RECORD или RECORDS. Наличие ключевого слова RECORDS означает наличие описания заголовка для множества типов записей в ЯОД описании (см. рис 2) и дальнейшее появление ключевого слова RBODY, обозначающего начало описания каждой отдельной записи. Заголовок каждой записи имеют одинаковую структуру и содержит данные, которые позволяют идентифицировать каждую запись в файле с данными. При обнаружении слова RECORDS выполняется сохранение описания заголовка записей для дальнейшей обработки. Ключевое слово RECORD означает, что описывается только одна запись, а описание заголовка отсутствует. При обнаружении слова RECORD выполняется чтение и сохранение описания всей записи.

На втором этапе, если обнаружено ключевое слово RECORDS:

- выполняется обработка текстового описания заголовка записей и сохранение его в виде объекта класса Group (поле recHeader класса Core);
- выполняется поиск ключевого слова RBODY и считывание описания каждой записи до конца файла. Описание каждой записи заканчивается комбинацией ключевого слова END и именем записи, указанного после RBODY;
- выполнение рекурсивной обработки полученного описания записи — разделение на отдельные текстовые сегменты (см. рис 3). Для каждого сегмента определяется уровень вложенности (начинается

```

FAMILY SI V(4044);

RECORDS;
LNG  ДЛЗАП  В(2) PC(4); //
MIT  НУЛИ    В(2) PC(4); //
KEY(I,СПГОД) ГОД      В(2) PC(4); // Год
KEY(I,СПМЕСЯЦ ) МЕСЯЦ В(1) PC(2); // Месяц
KEY(U) СТАНЦИЯ В(4) PC(7); // Станция
MRC(I) ТИПЗАП В(1) PC(2); // Тип записи (1-12)
    
```

Рис. 2. Описание семейства ЯОД файлов и заголовка записей в ЯОД описании

с нуля) и имя вышестоящего сегмента. Последние два шага выполняются для каждой следующей записи.

Если обнаружено ключевое слово RECORD, выполняется рекурсивная обработка и преобразование записи в объект.

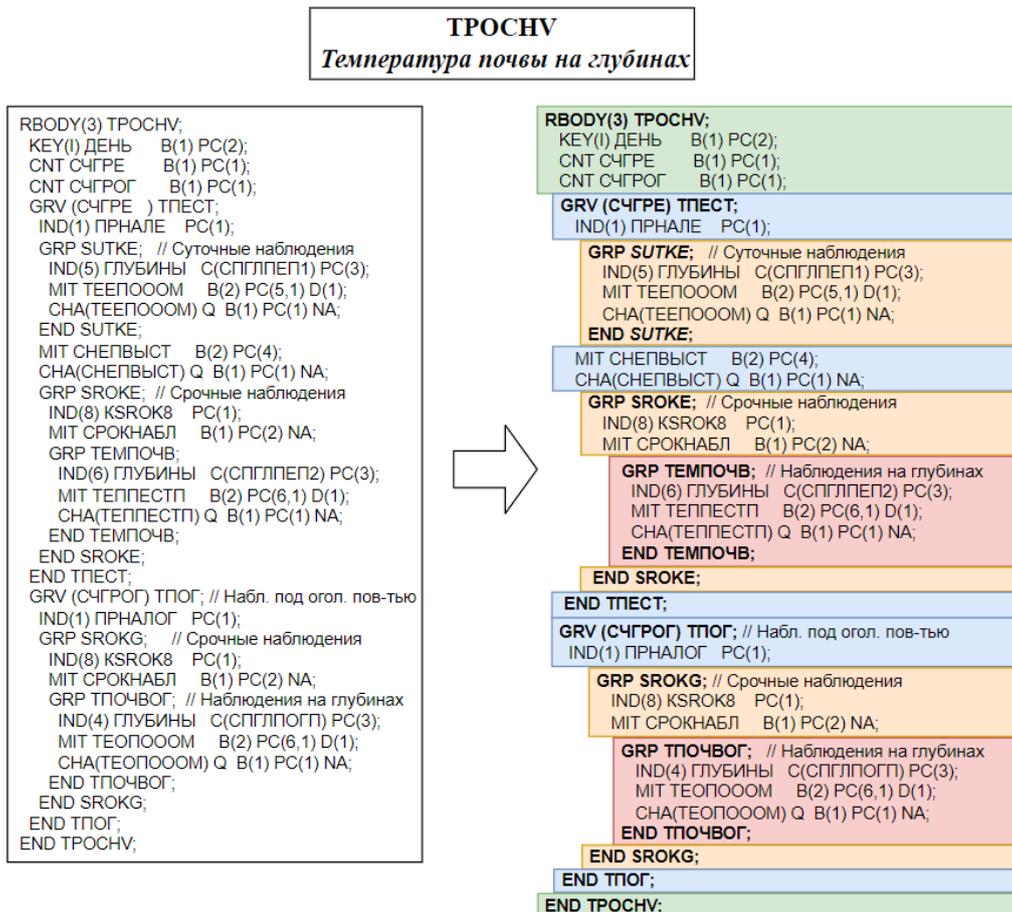


Рис. 3. Разделение ЯОД описания записи ТРОСНУ на отдельные тестовые сегменты

На третьем этапе в записи выполняется объединение текстовых сегментов, описывающих одну и ту же группу или элементы записи на одном и том же уровне. Таким образом получается список объединённых сегментов, который можно представить в виде иерархической структуры. На рисунке 4 представлен пример результата работы третьего этапа.

На четвертом этапе выполняется преобразование списка сегментов записи в список объектов класса Group, каждый из которых содержит полную информацию о группе. Причём преобразованию

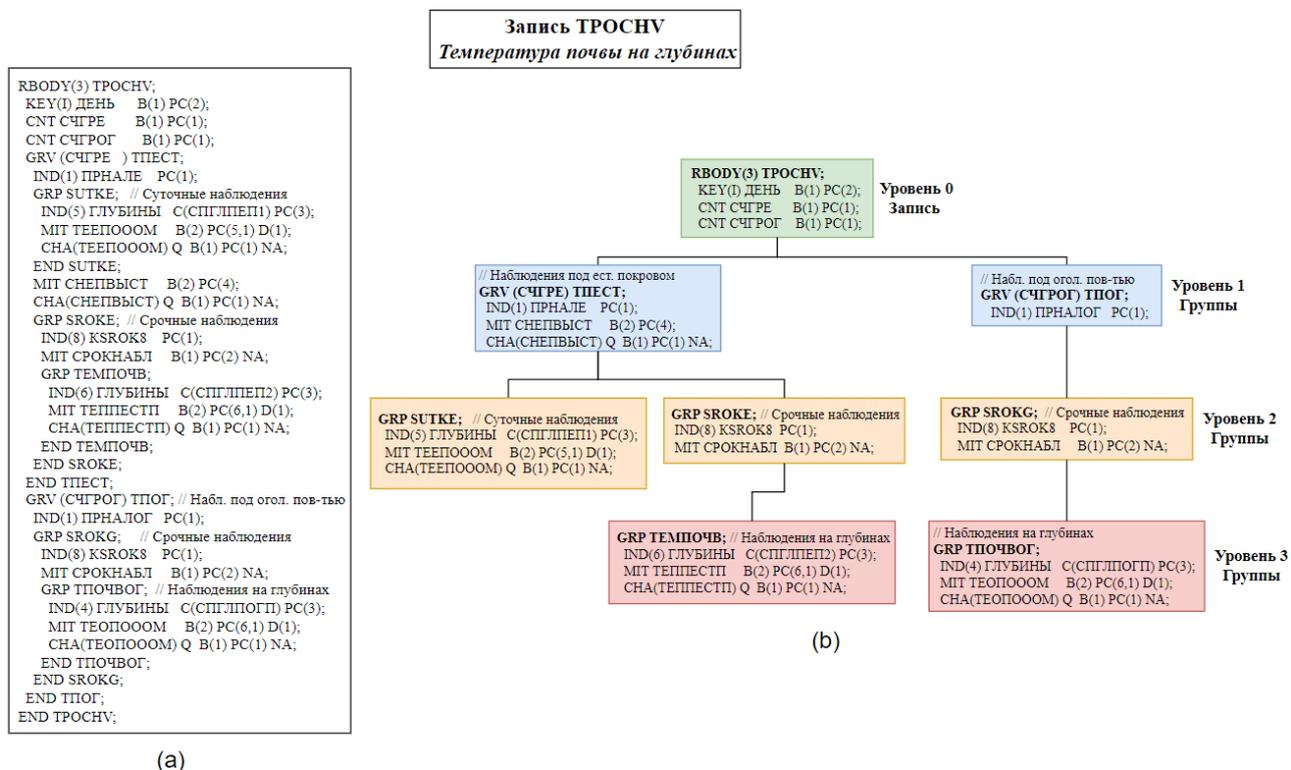


Рис. 4. Запись ТРОСНУ в формате ЯОД: (а) описание записи ТРОСНУ в формате ЯОД; (b) иерархическая структура записи ТРОСНУ из объединённых сегментов

Список фрагментов групп		Список целых групп	
RBODY(3) ТРОСНУ; KEY(I) ДЕНЬ В(1) PC(2); CNT СЧГРЕ В(1) PC(1); CNT СЧГРОГ В(1) PC(1);	Уровень 0; Предок ТРОСНУ	RBODY(3) ТРОСНУ; KEY(I) ДЕНЬ В(1) PC(2); CNT СЧГРЕ В(1) PC(1); CNT СЧГРОГ В(1) PC(1);	Уровень 0; Предок ТРОСНУ
GRV (СЧГРЕ) ТПЕСТ; // Наблюдения под ест. покровом IND(1) ПРНАЛЕ PC(1);	Уровень 1; Предок ТРОСНУ	GRV (СЧГРЕ) ТПЕСТ; // Наблюдения под ест. покровом IND(1) ПРНАЛЕ PC(1); MIT СНЕПВЫСТ В(2) PC(4); СНА(СНЕПВЫСТ) Q В(1) PC(1) NA; END ТПЕСТ;	Уровень 1; Предок ТРОСНУ
GRP SUTKE; // Суточные наблюдения IND(5) ГЛУБИНЫ С(СПГЛПЕП1) PC(3); MIT ТЕЕПОООМ В(2) PC(5,1) D(1); СНА(ТЕЕПОООМ) Q В(1) PC(1) NA; END SUTKE;	Уровень 2; Предок ТПЕСТ	GRP SUTKE; // Суточные наблюдения IND(5) ГЛУБИНЫ С(СПГЛПЕП1) PC(3); MIT ТЕЕПОООМ В(2) PC(5,1) D(1); СНА(ТЕЕПОООМ) Q В(1) PC(1) NA; END SUTKE;	Уровень 2; Предок ТПЕСТ
MIT СНЕПВЫСТ В(2) PC(4); СНА(СНЕПВЫСТ) Q В(1) PC(1) NA;	Уровень 1; Часть ТПЕСТ; Предок ТРОСНУ	GRP SROKE; // Срочные наблюдения IND(8) КSROK8 PC(1); MIT СРОКНАБЛ В(1) PC(2) NA; END SROKE;	Уровень 2; Предок ТПЕСТ
GRP SROKE; // Срочные наблюдения IND(8) КSROK8 PC(1); MIT СРОКНАБЛ В(1) PC(2) NA;	Уровень 2; Предок ТПЕСТ	GRP SROKE; // Срочные наблюдения IND(8) КSROK8 PC(1); MIT СРОКНАБЛ В(1) PC(2) NA; END SROKE;	Уровень 2; Предок ТПЕСТ
GRP ТЕМПОЧВ; // Наблюдения на глубинах IND(6) ГЛУБИНЫ С(СПГЛПЕП2) PC(3); MIT ТЕПЕСТП В(2) PC(6,1) D(1); СНА(ТЕПЕСТП) Q В(1) PC(1) NA; END ТЕМПОЧВ;	Уровень 3; Предок SROKE	GRP ТЕМПОЧВ; // Наблюдения на глубинах IND(6) ГЛУБИНЫ С(СПГЛПЕП2) PC(3); MIT ТЕПЕСТП В(2) PC(6,1) D(1); СНА(ТЕПЕСТП) Q В(1) PC(1) NA; END ТЕМПОЧВ;	Уровень 3; Предок SROKE
GRV (СЧГРОГ) ТПОГ; // Набл. под огол. пов-тью IND(1) ПРНАЛОГ PC(1);	Уровень 1; Предок ТРОСНУ	GRV (СЧГРОГ) ТПОГ; // Набл. под огол. пов-тью IND(1) ПРНАЛОГ PC(1); END ТПОГ;	Уровень 1; Предок ТРОСНУ
GRP SROK8; // Срочные наблюдения IND(8) КSROK8 PC(1); MIT СРОКНАБЛ В(1) PC(2) NA;	Уровень 2; Предок ТПОГ	GRP SROK8; // Срочные наблюдения IND(8) КSROK8 PC(1); MIT СРОКНАБЛ В(1) PC(2) NA; END SROK8;	Уровень 2; Предок ТПОГ
GRP ТПОЧВОГ; // Наблюдения на глубинах IND(4) ГЛУБИНЫ С(СПГЛПОГП) PC(3); MIT ТЕОПОООМ В(2) PC(6,1) D(1); СНА(ТЕОПОООМ) Q В(1) PC(1) NA; END ТПОЧВОГ;	Уровень 3; Предок SROK8	GRP ТПОЧВОГ; // Наблюдения на глубинах IND(4) ГЛУБИНЫ С(СПГЛПОГП) PC(3); MIT ТЕОПОООМ В(2) PC(6,1) D(1); СНА(ТЕОПОООМ) Q В(1) PC(1) NA; END ТПОЧВОГ;	Уровень 3; Предок SROK8

Рис. 5. Списки описания групп: (а) список частей групп; (b) список целых групп

подвергаются как список отдельных сегментов записи, так и список объединённых сегментов. В итоге получается список целых групп (см. рис. 5.a) и список частей групп (см. рис. 5.b).

Список целых групп используется для создания структуры, например, в реляционной базе данных, и выстраивания связей между таблицами в ней. Список частей группы отражает порядок ЯОД и используется при чтении данных в формате ЯОД.

Заключение

В настоящей статье были представлены: диаграмма классов в нотации UML с подробным описанием каждого класса, особенности разработанного парсера, поэтапное описание алгоритма работы парсера на примере описания записи ТРОСНУ (температура почв на глубине), относящейся к восьмисрочным метеорологическим данным. В результате тестирования работы парсера на примере шестидесяти различных ЯОД описаний была доказана универсальность разработанного парсера.

ЛИТЕРАТУРА

1. Перетяцько Л. О. Единый государственный фонд гидрометеорологических данных как большие данные. Технологии и инструменты для работы с ним. *Успехи кибернетики*. 2022;3(4):98–101. DOI: 10.51790/2712-9942-2022-3-4-11.
2. Кофтан Ю. Р., Перетяцько Л. О. К построению технологии взаимной конвертации баз данных различной структуры для пополнения и верификации данных ЕГФД, а также для обслуживания потребителей. *Труды ВНИИГМИ-МЦД*. 2018;181:162–174.
3. Шлее М. *Qt 5.10. Профессиональное программирование на C++*. СПб.: БХВ-Петербург; 2018. 1072 с.
4. *Документация Qt версии 5*. Режим доступа: <https://doc.qt.io/qt-5/qstring.html>.
5. Перетяцько Л. О., Кофтан Ю. Р. Система взаимной конвертации данных различной структуры для обслуживания потребителей ЕГФД. *Труды ВНИИГМИ-МЦД*. 2020;186:163–175.
6. Арлоу Д., Нейштадт И. *UML 2 и Унифицированный процесс, Практический объектно-ориентированный анализ и проектирование*, 2-е изд. / пер. с англ. СПб: Символ-Плюс; 2007. 624 с.
7. Фицджеральд М. *Регулярные выражения: основы* / пер. с англ. М.: ООО «И.Д. Вильямс»; 2015. 144 с.