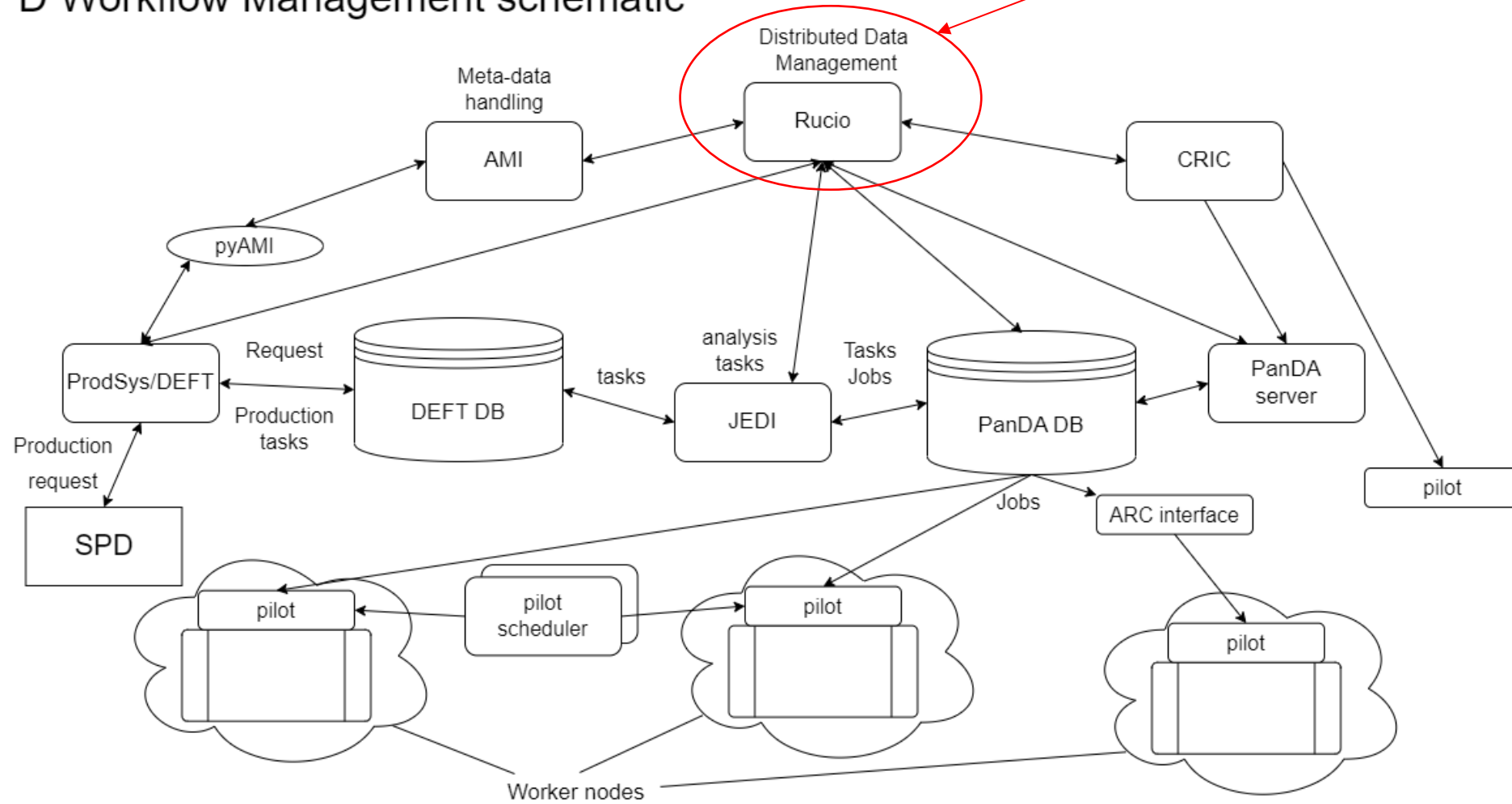


# Система управления научными данными Rusio

Конак А.С. (ТулГУ), Петросян А.Ш. (ОИЯИ)

# SPD offline data processing

SPD Workflow Management schematic



# Что такое Rucio?

Rucio — это программный комплекс с открытым исходным кодом, которая обеспечивает масштабируемый функционал для организации, управления и доступа к данным. Данные могут быть распределены по географически распределенным центрам обработки данных.

Первоначально Rucio был разработан для удовлетворения требований эксперимента ATLAS, и теперь он постоянно расширяется для поддержки как экспериментов LHC так и других научных сообществ.

# Зачем нужна система Rucio?

Эксперимент SPD будет генерировать большой поток данных. Все эти данные будут собраны в файлы, которые будут храниться в разных местах.

Для управления всеми файлами необходима специальная система.

# Зачем нужна система Rucio?

В настоящее время система Rucio используется для:

- каталогизации данных представленных в виде файлов и наборов файлов (датасетов);
- единое взаимодействие с разнородными инфраструктурами хранения;
- восстановление данных;
- адаптивная репликация.

# Компоненты Rucio

Уровень клиентов связывает пользователя с системой и состоит из таких компонентов, как клиенты командной строки (CLI), клиенты Python, а также пользовательский веб-интерфейс.

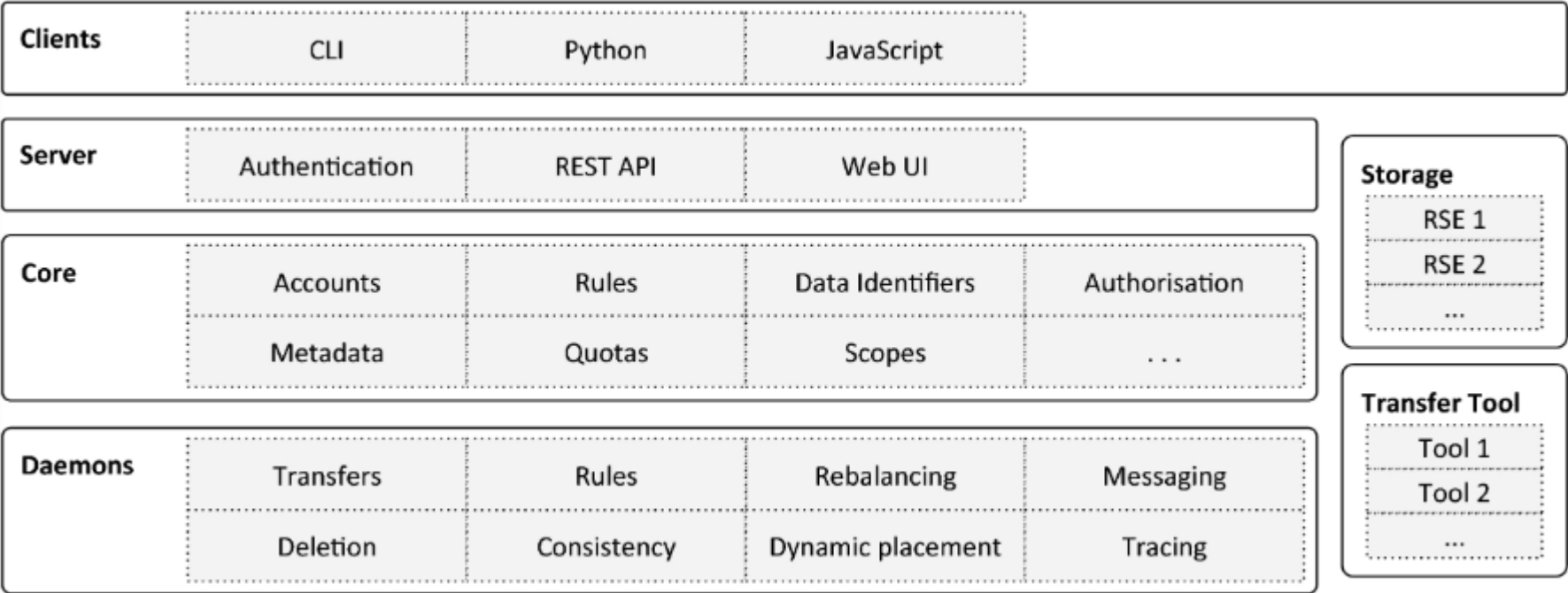
Уровень сервера служит для аутентификации и предоставляет пользователю привилегии для взаимодействия с клиентами и другими внешними приложениями.

Ядро - основной уровень, определяющий все взаимодействия пользователя с данными.

Уровень демонов заботится обо всех асинхронных и непрерывных рабочих процессах в фоновом режиме.

Уровень хранилища - взаимодействие с различными промежуточными инструментами и системами хранения.

Инструменты передачи позволяют демонам отправлять, запрашивать и отменять передачи независимо от фактической используемой службы передачи.



# Концепция Rucio

Основные концепции, поддерживаемые Rucio, охватывают

- Namespace - обеспечивает единое пространство имен объединяя все системы хранения в одно пространство.
- Accounts - обрабатывает авторизацию, аутентификацию и выдает права пользователю.
- Storage - обеспечивает унифицированный интерфейс с распределенными центрами обработки данных. Rucio Storage Element (RSE) – минимальная единица адресуемого хранилища, содержащее описание всех атрибутов для доступа к пространству хранения.
- Subscriptions - используются для крупномасштабных политик потоков данных.
- Replication rules - обеспечивают согласованное распределенное состояние пространства имен в хранилище, также на правилах репликации основано управление репликами.

# Представление данных в Rucio

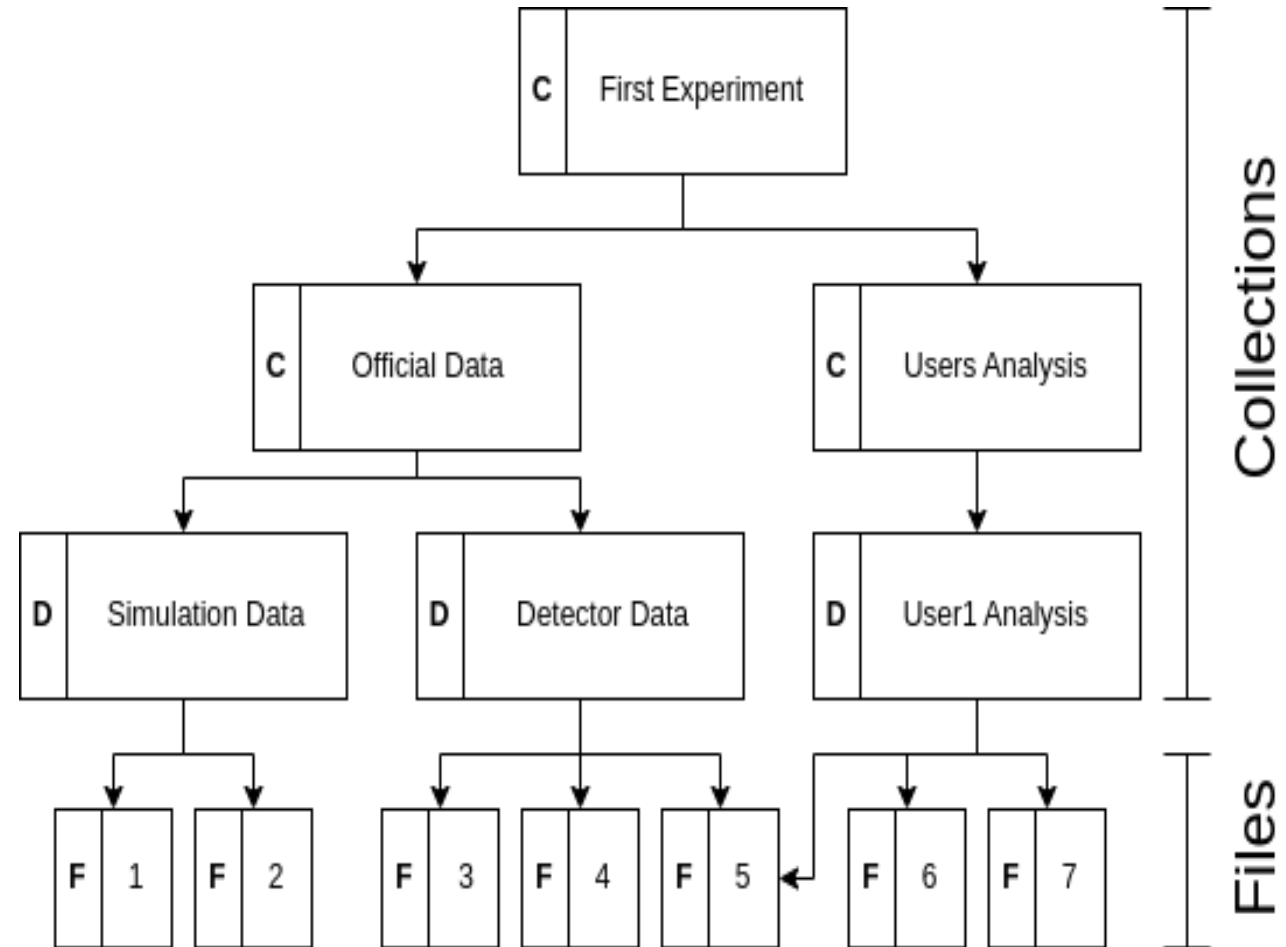
Данные физически хранятся в файлах. Файл - наименьшая операционная единица данных в Rucio. Файлы могут быть сгруппированы в наборы данных, а наборы данных могут быть сгруппированы в контейнеры.

Для работы со всеми типами данных используются Data Identifiers (DID). Идентификатор данных — это имя отдельного файла, набора данных или контейнера (data2023:mysearch1).

У каждого объекта есть статусы.

Статус файла (доступен, утерян, удален)

Статусы коллекций( открытый; монотонный).





# Текущий прогресс

На данный момент к тестированию готово стандартное окружение, развернутое на виртуальной машине на базе облачного сервиса ЛИТ.

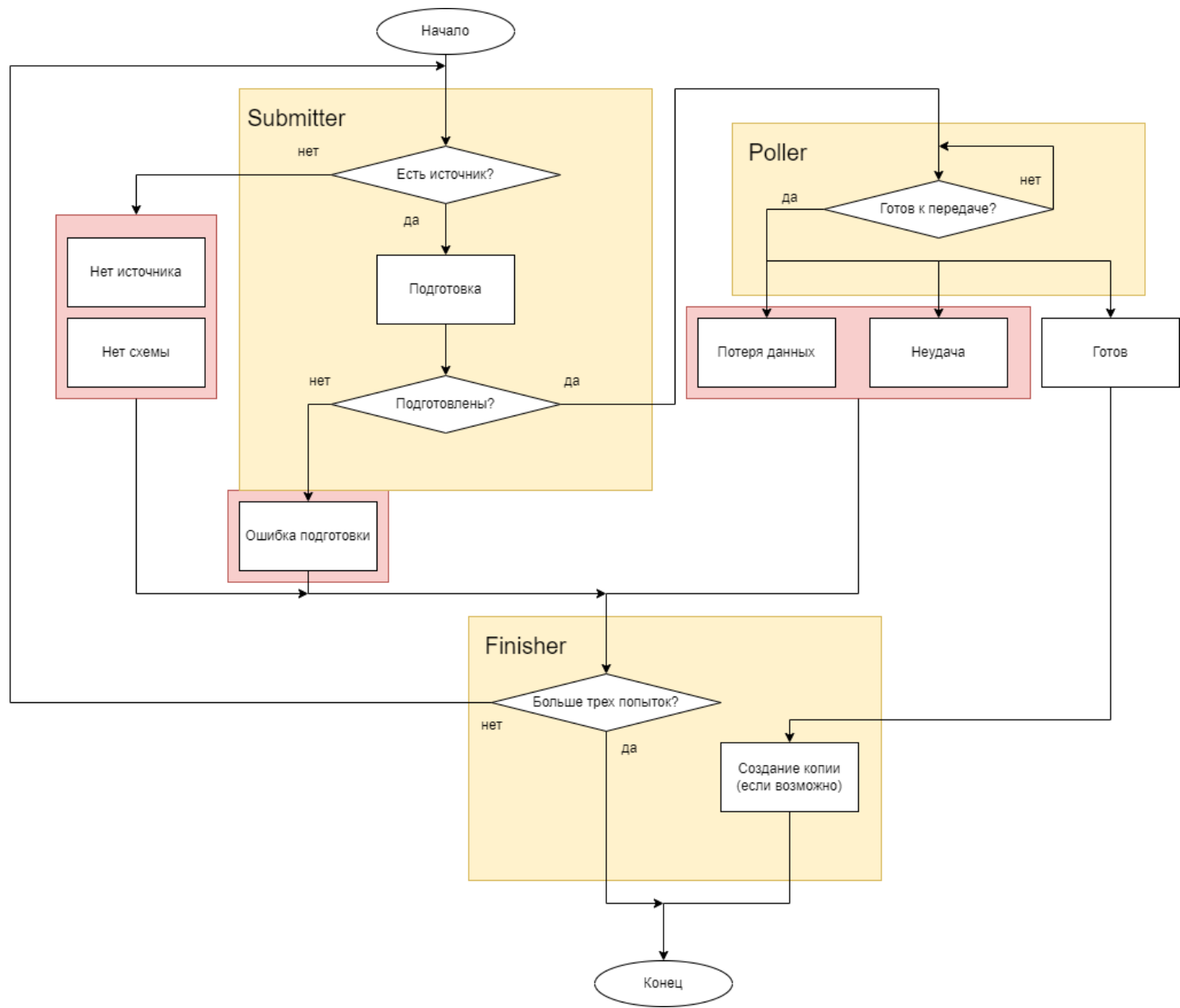
- Rucio-server;
- PostgreSQL;
- Daemons.

# Структура системы демонов

Submitter - выполняет выбор источника, вычисление пути и фактическую отправку переводов во внешний инструмент передачи.

Poller - регулярно опрашивает внешний инструмент передачи на предмет статуса ожидающих переводов и помечает их как успешные/неудачные.

Finisher - действует на успешные или неудачные передачи.



# Планы

- Тестирование взаимодействия с системой управления нагрузкой PanDA в рамках обработки задач.
- Развертывание рабочей системы.
- Перевод системы на аутентификацию и авторизацию на основе технологии JWT.

**Спасибо за внимание!**