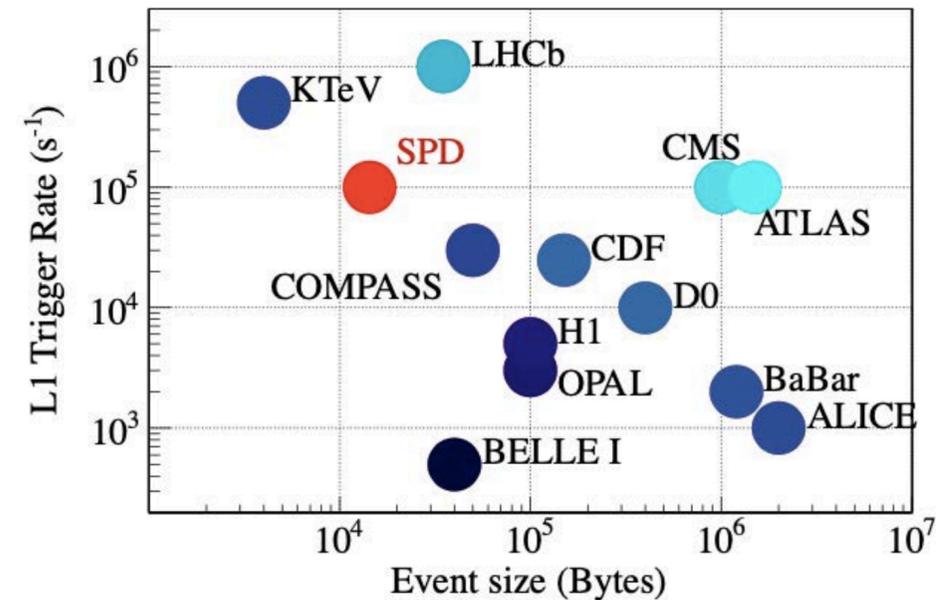


Распределённая система хранения и обработки данных для эксперимента SPD

Артём Петросян, ЛИТ ОИЯИ
Осенняя школа по информационным технологиям ОИЯИ
19 октября 2023

The expected event rate of the SPD experiment is about 3 MHz (pp collisions at $\sqrt{s} = 27$ GeV and 10^{32} cm⁻²s⁻¹ design luminosity). This is equivalent to a **raw data rate** of 20 GB/s or **200 PB/year**, assuming a detector duty cycle is 0.3, while the signal-to-background ratio is expected to be on the order of 10^{-5} . Taking into account the bunch-crossing rate of 12.5 MHz, one may conclude that pile-up probability cannot be neglected.

- SPD TDR



The goal of the **online filter** is at least to decrease the data rate by a factor of 20, so that the **annual growth of data**, including the simulated samples, stays within **10 PB**. Then, data are transferred to the Tier-1 facility, where a full reconstruction takes place and the data is stored permanently. The data analysis and Monte-Carlo simulation will likely run at the remote computing centers (Tier-2s). Given the large data volume, a thorough optimization of the event model and performance of the reconstruction and simulation algorithms are necessary.

- Набор данных ~10 петабайт данных каждый год
- Размер события 10-15 килобайт
- 1 секунда на обработку одного события
- Необходимо контролировать размеры файлов — слишком маленькие создадут нагрузку на каталоги и системы управления нагрузкой (1 файл = 1 запись и 1 задача), слишком большие сложно передавать и долго обрабатывать (оптимум — 6-8 часов)
- Итого: около 60000 ЦПУ для обработки и прирост на 10 ПБ в год

Что это за данные

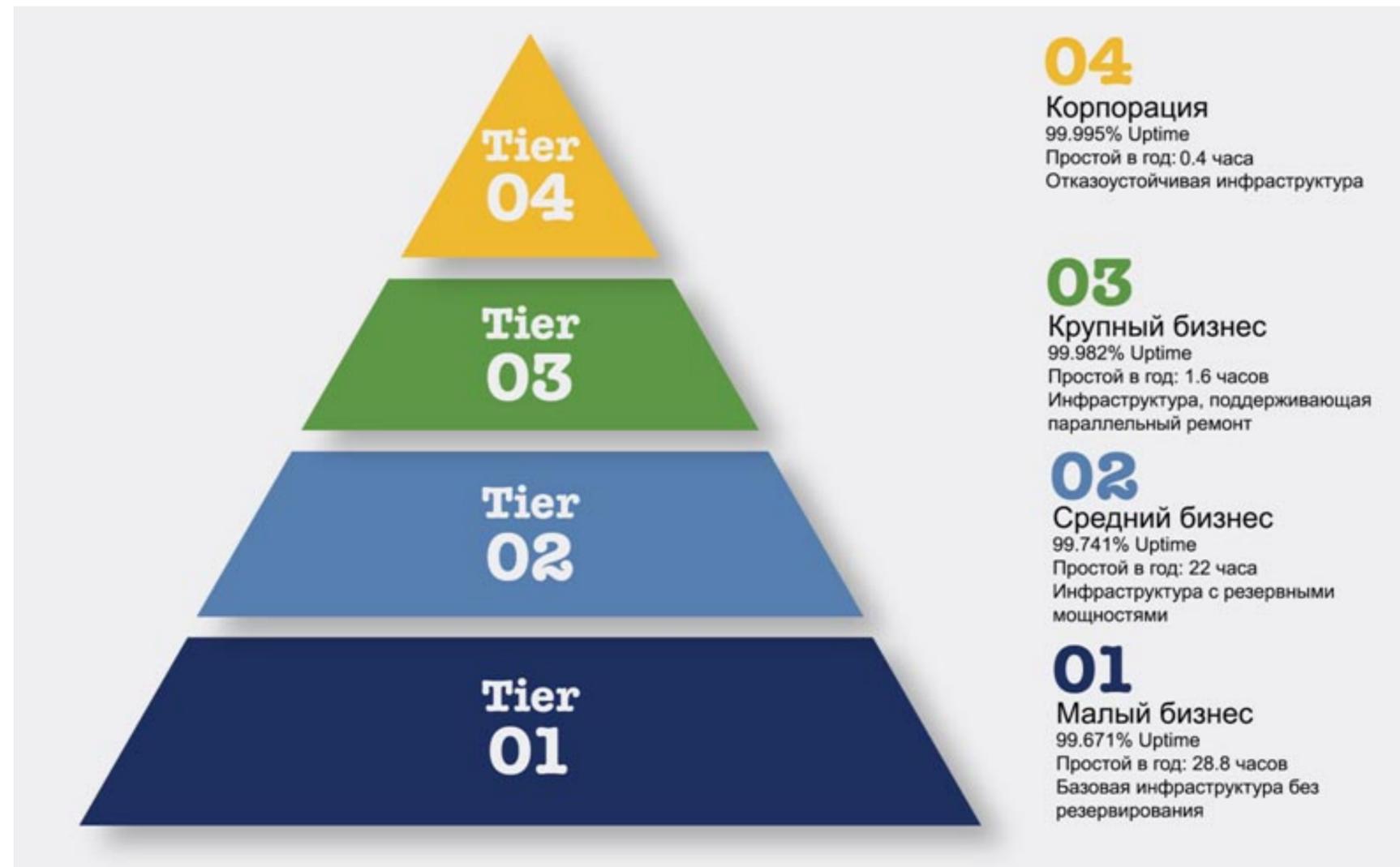
- Результаты моделирования
- Данные с детектора
- Данные различных промежуточных форматов по пути от “сырых” до готовых к анализу физическими группами
- Должны ли они храниться одинаковое количество времени?

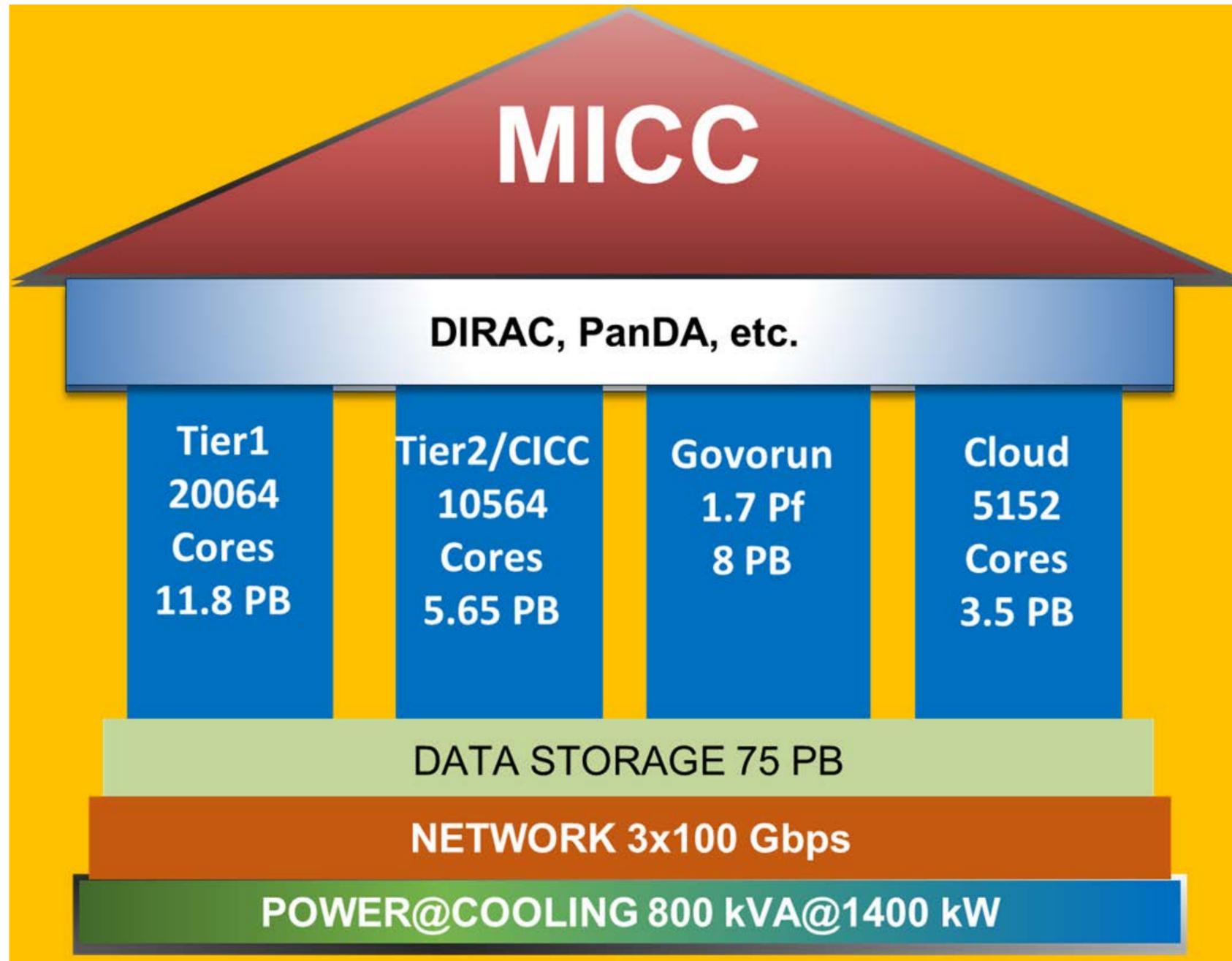
- Данные неоднородные и в зависимости от типа должны храниться различное время
 - Ценные: полученные в результате работы детектора — хранить вечно
 - Их практически невозможно воспроизвести в случае утраты
 - Являются источником физических результатов — как долго сохранять решает коллаборация в каждом конкретном случае
 - Могут быть получены заново, но это достаточно долго и дорого
 - Временные — нужные для проведения какого-то этапа вычислений, после завершения которого могут быть удалены
 - Основная проблема при работе с ними это быстро их (и только их) удалять

- [https://en.wikipedia.org/wiki/Replication_\(computing\)](https://en.wikipedia.org/wiki/Replication_(computing))
- Единственный вариант — хранить несколько копий одного и того же файла в разных местах, подключенных к разным (лучше к нескольким) сетевым каналам и имеющим разное питание (лучше несколько видов), например: одну копию в Москве, а вторую в Санкт-Петербурге
- В рамках одного дата-центра данные тоже могут быть продублированы: любая сетевая файловая система, например Serp FS имеет параметр репликации, задающий, сколько должно быть записано копий каждого файла

Классификация ЦОД

- https://habr.com/ru/company/cloud_mts/blog/332864/
- Tier I — базовая инфраструктура без резервирования
- Tier II — инфраструктура с резервными мощностями
- Tier III — инфраструктура, поддерживающая параллельный ремонт
- Tier IV — отказоустойчивая инфраструктура





4 advanced software and hardware components

- Tier1 grid site
- Tier2/CICC site
- hyperconverged “Govorun” supercomputer
- cloud infrastructure

Distributed multi-layer data storage system

- Disks
- Robotized tape library

Network

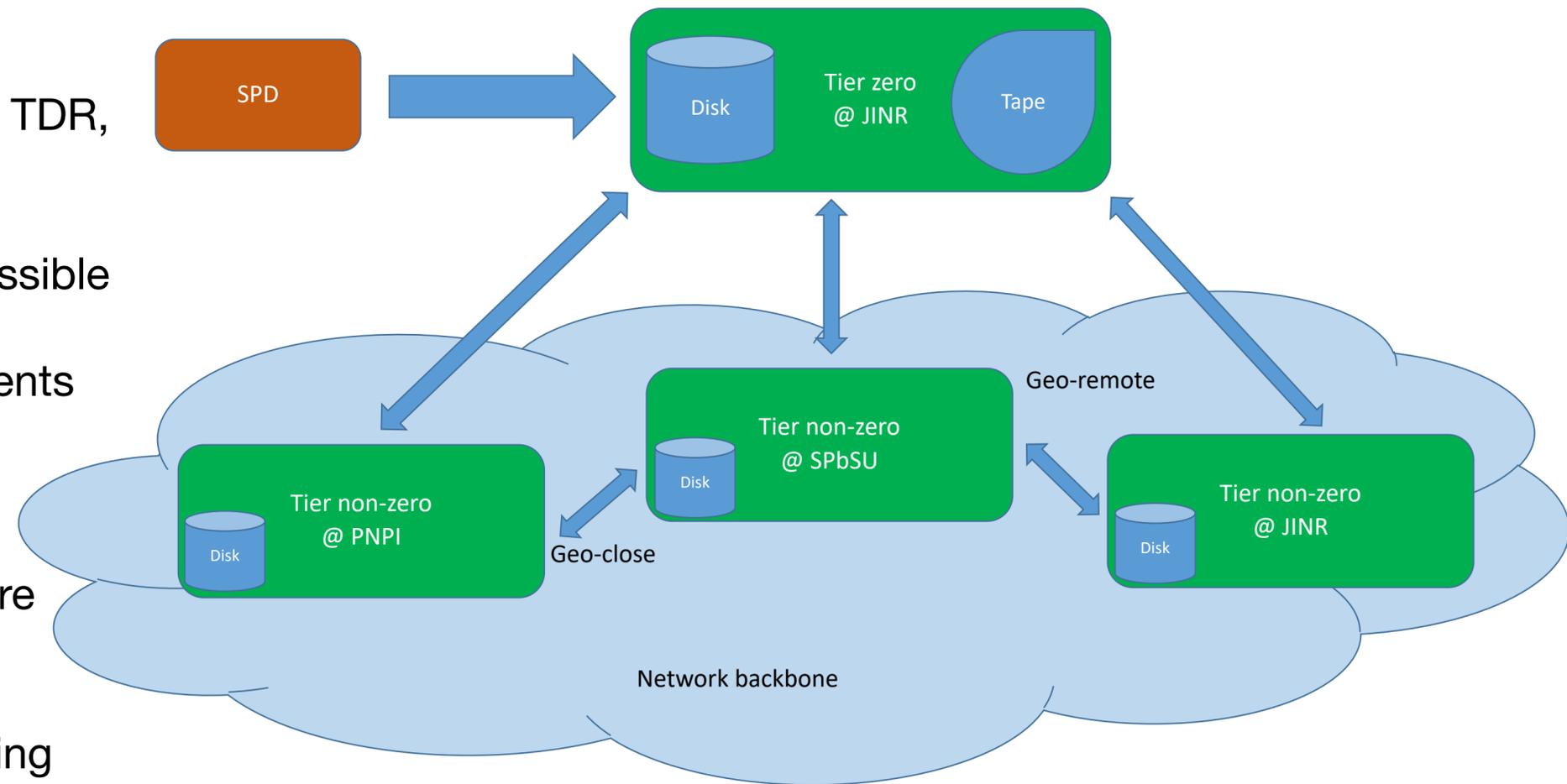
- Wide Area Network
- Local Area Network

Engineering infrastructure

- Power
- Cooling

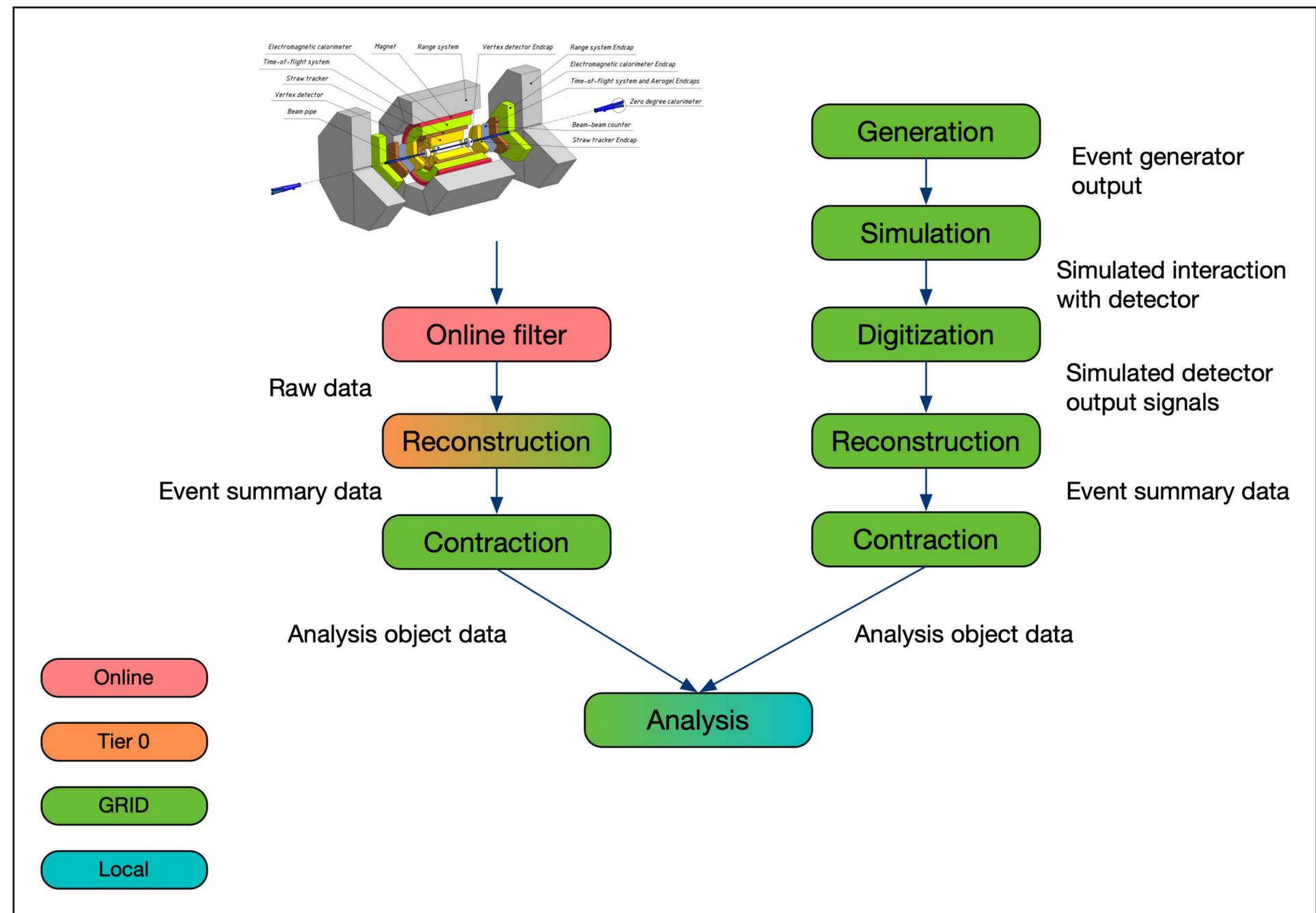
The main objective of the project is to ensure multifunctionality, scalability, high performance, reliability and availability in 24x7x365 mode for different user groups that carry out scientific studies within the JINR Topical Plan.

- Data volume mandates some baselines
 - >10 Gbps network per site (from TDR)
 - >500 TB storage capacity per site (not from TDR, but might be added to the next version)
- Try to use existing free software as much as possible
 - Experience comes from large LCG experiments
- Optimize management and operation effort
 - Do not deploy home-grown solutions that are different from site to site
 - Provide a reasonable guidelines for interfacing physical resources with central data management services



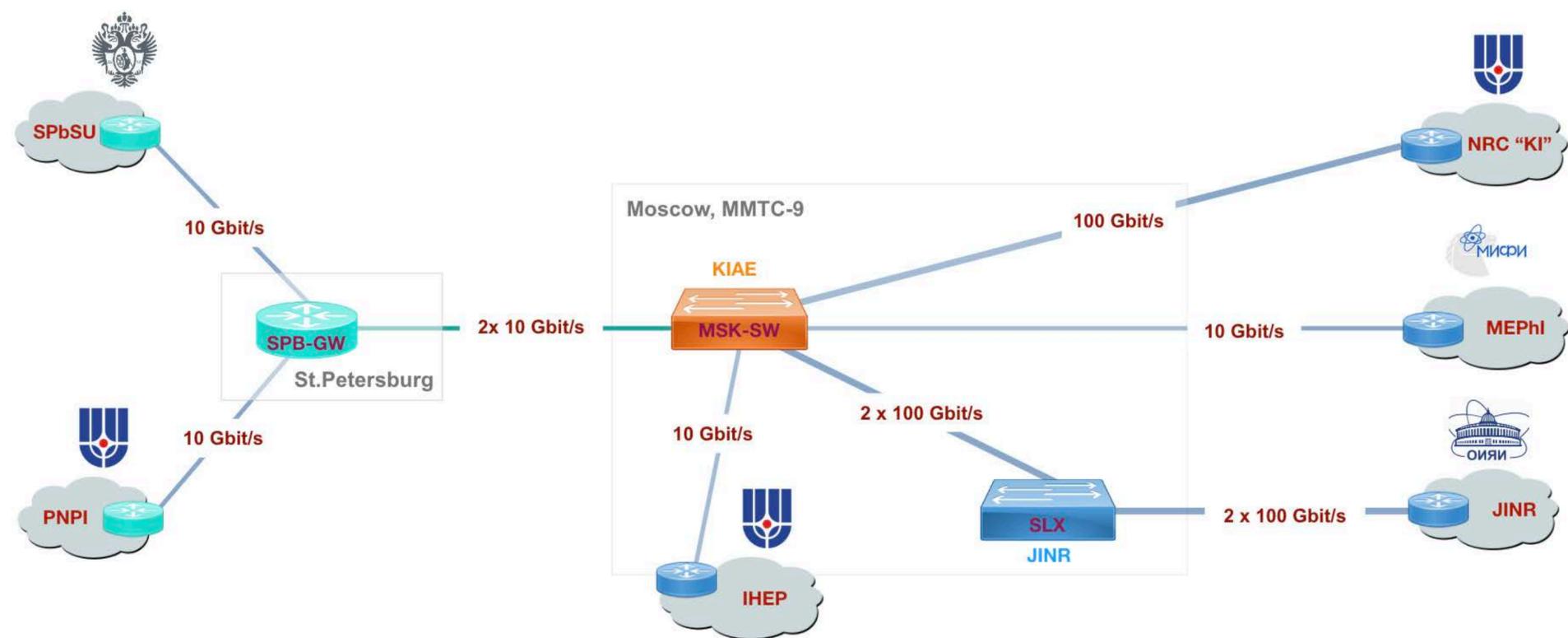
Processing steps distribution over computing resource types

- Execution of events reconstruction and reprocessing jobs is accompanied by intensive I/O operations and will be done mostly on the dedicated farms on JINR site as Tier 0 component of the distributed computing system
- The use of Tier 0 is dictated by huge amount of initial data, gathered by the physics facility — data must be reduced as much as possible in order to be ready for distribution
- Less I/O intensive steps, especially Monte-Carlo production, can be performed on the remote computing centers
- User analysis can be run on every close to user resource

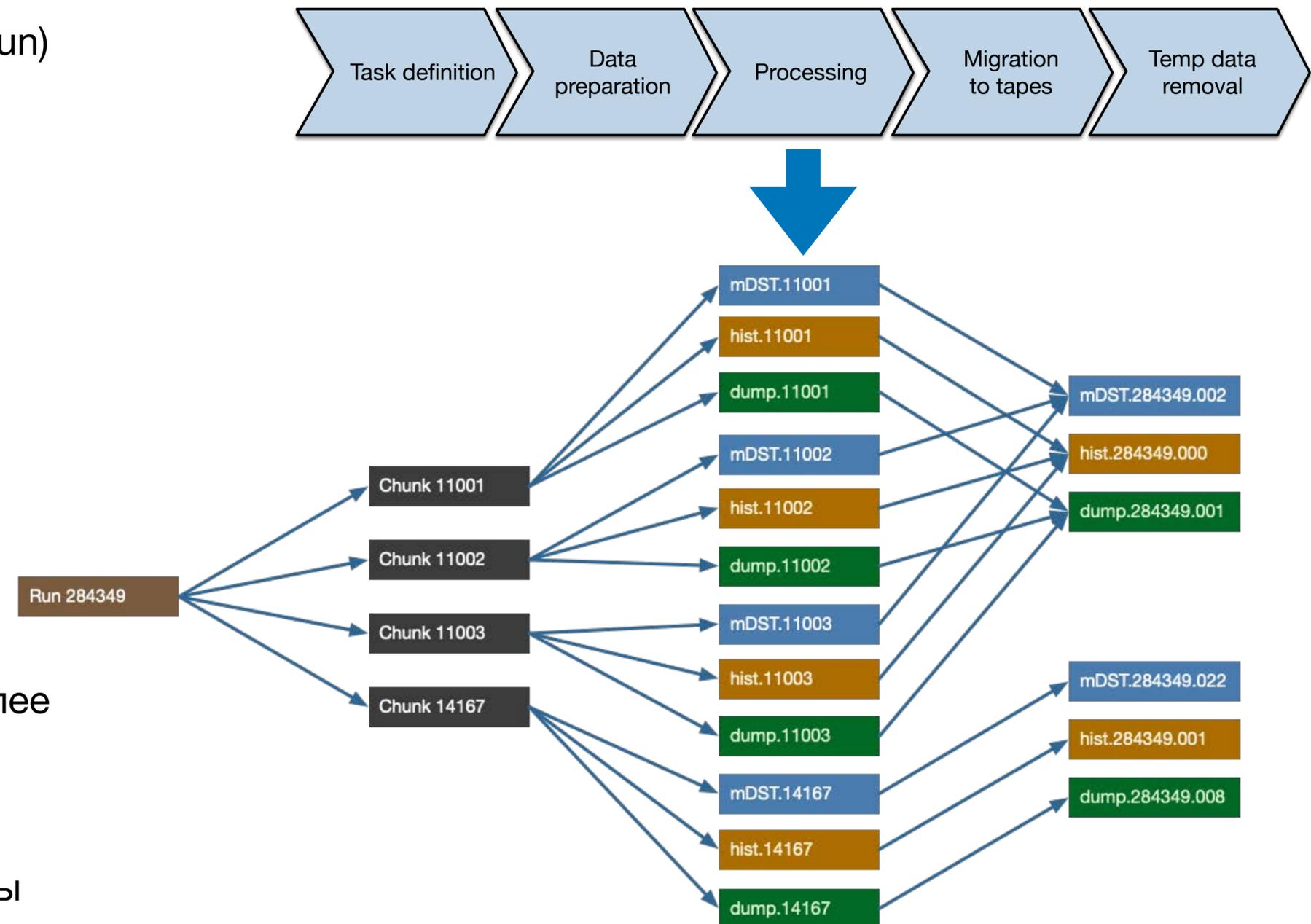


Уже подтвердившие участие в обработке данных в рамках проекта SPD вычислительные центры

- Участники коллаборации, уже предоставляющие вычислительные ресурсы: СПбГУ, ПИЯФ, ИЯФ БУ
- Работаем над тем, чтоб расширить список участников: СамГУ готов предоставить кластер “Королёв”
- Мы должны быть готовы в любой момент подключить коммерческие облака и суперкомпьютеры

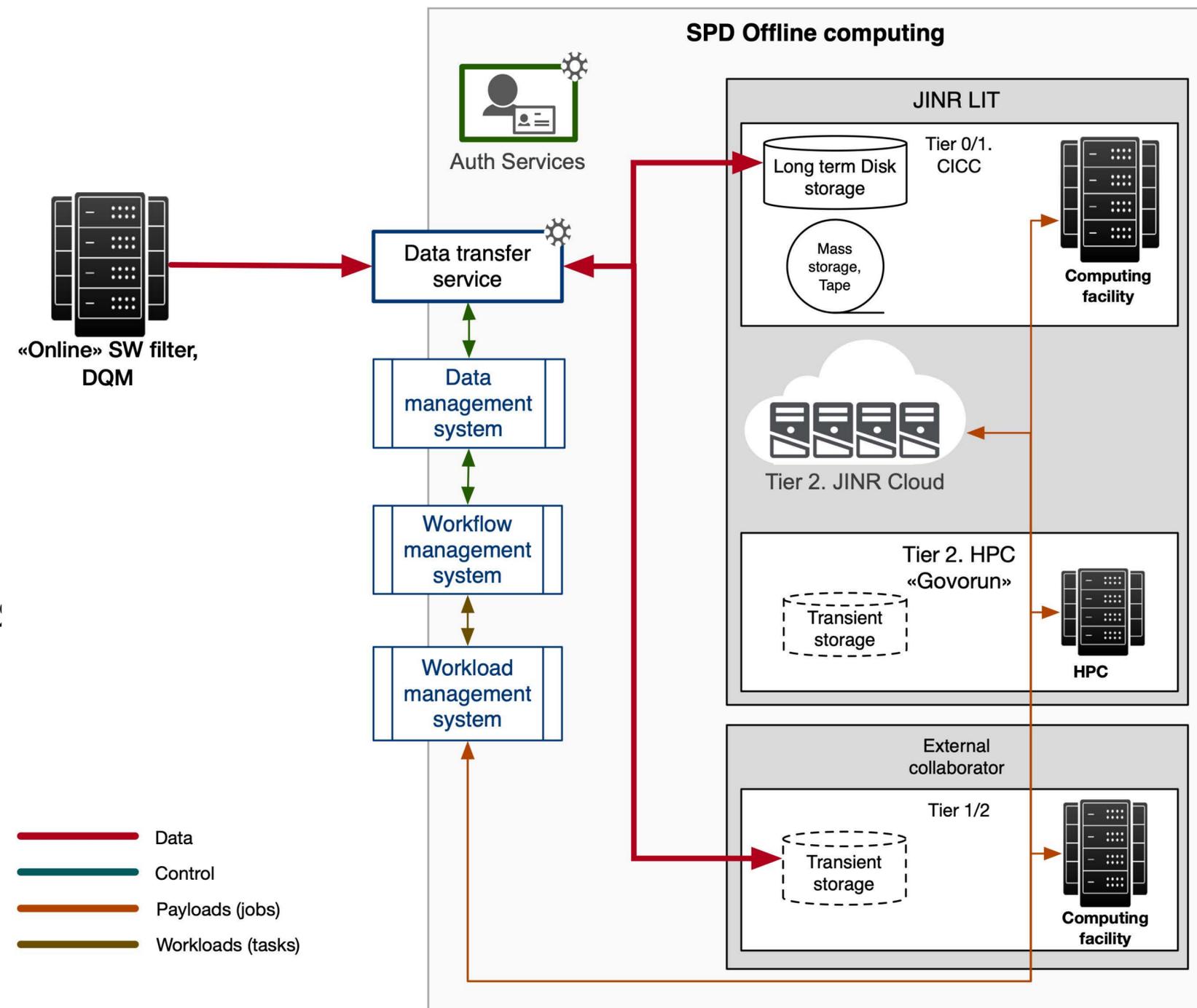


- Набор файлов из одного сброса детектора (run) может состоять из различного количества файлов необработанных данных (chunks)
- Система управления процессом обработки данных создает по задаче на каждый файл
- Каждая завершившаяся задача создает выходные файлы разных типов: дерево событий (mini Data Summary Tree, mDST), гистограмму, и файл, содержащий только данные о значимых для данного задания событиях (event dump)
- Задание на обработку может состоять из более чем 200 000 файлов, а с промежуточными, временными и финальными их количество может превышать 1 000 000 – необходимо максимально автоматизировать все процессы



Как управлять подобными масштабными инфраструктурами?

- Основными системами и сервисами подобных распределенных инфраструктур являются: система аутентификации, система авторизации, информационная система, система управления распределенными данными, система управления нагрузкой, система управления процессами обработки, сервис передачи данных, сервис кэширования программного обеспечения.
- Каждый из этих компонентов это большой (некоторым проектам уже более 20 лет) и развиваемый группой разработчиков программный продукт



- Создание системы мониторинга подобной инфраструктуры это отдельная крупная задача
- Система мониторинга должна включать не только мониторинг движения данных и задач, но и работу каждого из сервисов, в том числе и мониторинг состояния сетевых каналов между центрами обработки и хранения
- Должна позволять как отслеживать обработку данных на самом высоком уровне, так и предоставлять возможности анализа деталей выполнения каждой отдельной задачи
- Должна формировать сводный пульт наблюдения и управления, куда будут выводиться только самые важные события, требующие немедленного реагирования в случае сбоев

- Работа с подобным метавычислителем и метохранилищем требует большого внимания к вопросам безопасности
- Каждый пользователь должен получить сертификат стандарта X.509 (в будущем JSON Web Token), и получить соответствующую своим обязанностям роль в системе
- Каждое действие подписывается проху-сертификатом, время жизни которого не превышает нескольких суток
- Загрузка прикладного ПО в систему доступна нескольким ответственным разработчикам после проведения набора тестов на уровне Gitlab и в процессе интеграции, и в дальнейшем размещается на файловой системе, доступной только для чтения
- Любые подозрительные действия или действия, приводящие к перегрузке каких-то частей системы, приводят к отключению пользователя, если же действия носили намеренный характер, то к отзыву сертификата
- Действия отслеживаются как на уровне системы управления, так и на уровне вычислительных центров — по цифровым сигнатурам использования ЦПУ и работы с памятью и сетью

- Строящийся детектор SPD будет генерировать потоки данных, с хранением и обработкой которых в рамках одного ЦОД справиться невозможно
- В рамках решения задач по управлению данными мы создаем распределенное хранилище с контролем реплик, управлением временем жизни данных и их целостностью
- Для обработки мы привлекаем все возможные вычислительные ресурсы, и стараемся использовать их наиболее оптимальным образом, распределяя задачи по наиболее подходящим центрам
- Для управления процессами обработки данных эксперимента SPD мы строим высокоавтоматизированную систему, учитывающую при работе разнообразные параметры: размеры файлов, их местонахождение, подходящие для каждого этапа обработки процессоры и память, состояние сетевых соединений и тд.

Спасибо за внимание!