

# EventIndex

каталог событий

Осенняя Школа по информационным технологиям ОИЯИ



- Рассматриваем EventIndex в том виде в каком он был реализован для эксперимента Atlas на LHC, с прицелом на использование в установке SPD на коллайдере NICA.
- Несмотря на то что физические задачи этих экспериментов кардинально различны, системы обработки данных имеют сходные черты, что позволяет использовать похожие решения
- На SPD ожидается объемы данных того же порядка как на ATLAS при даже большем количестве событий
- Сходная модель данных событий и файловая организация
- Основное отличие: отсутствие триггера (но есть Online Filter)



- Группы статистически эквивалентных событий хранятся в файлах на диске или на магнитной ленте.
  - события реальных данных, собранные при тех же условиях детектора
  - смоделированные события, созданные с помощью того же генератора и обработанные теми же версиями ПО
- Каждый файл обычно содержит от 1000 до 10000 событий, в зависимости от формата и балансировки :
  - слишком большое количества маленьких файлов ( $< 1$  ГБ) могут вызвать проблемы с хранением данных
  - слишком большие файлы ( $> 10$  ГБ) могут снизить эффективность передачи.



- Файлы группируются в датасеты (наборы данных)
  - Датасет как правило содержит события относящиеся к одному Run-у
    - Данные для длительных run-ов в процессе обработки как правило разбиваются на несколько датасетов
- Датасеты, в свою очередь могут быть сгруппированы в контейнеры
  - Контейнер может содержать датасеты относящиеся к одному Run-у, или например, к целому году набора данных
  - Один и тот же датасет может быть отнесен к нескольким контейнерам
- Пользователи могут создавать свои датасеты для нужд конкретного анализа



- В зависимости от стадии обработки данные могут иметь различные форматы
  - Изначально данные детектора записываются в формате RAW
  - После реконструкции событий получают AOD датасеты.
  - Derivation (DAOD) - наборы производных данных для анализов
- + на Grid генерируются (MC) аналоги реальных событий
  - EVNT датасеты содержат информацию о смоделированных частицах
  - AOD и DAOD датасеты — аналогично как для реальных событий
- Различные версии датасетов происходящих от одних и тех же событий, полученных с детектора или смоделированных
  - Получены с различными настройками реконструкции и версиями ПО
  - Могут создаваться регулярно либо по мере обнаружения ошибок



- EventIndex – единый каталог всех событий, полученных, обработанных и смоделированных в эксперименте ATLAS
- В рамках проекта была создана инфраструктура, которая
  - предоставляет систему индексирования данных, хранящую указатели на события ATLAS в миллионах файлов разбросанных по сотням серверов глобальной распределенной вычислительной сети LHC GRID, а также основную информацию о событиях.
  - позволяет собирать и сохранять информацию о событиях, автоматически обнаруживать и индексировать новые данные.
  - предоставляет набор инструментов для доступа к этой информации из командной строки, графического интерфейса либо через RESTful API.
  - позволяет быстро и эффективно отбирать искомые события из миллиардов полученных в результате эксперимента и смоделированных событий, используя (в том числе) триггерную информацию.



- Отбор событий (Event Picking)
  - получить событие в нужном формате и обработанное конкретной версией программного обеспечения
- Подсчет и выборка события на основании триггерных решений
- Проверка целостности и согласованности данных
- Создание матриц перекрытия триггерных цепочек
- Создание матриц перекрытия датасетов
- Быстрый просмотр данных
  - Поиск датасетов, составление отчетов, проверка датасетов



- Поиск событий среди сотен миллионов файлов разбросанных по серверам распределенной системы хранения данных
- Получение одного или нескольких событий
  - В нужном формате (AOE, DAOD\_TOP2, RAW) и обработанное конкретной версией программного обеспечения
    - посмотреть на eventDisplay (и вставить в презентацию)
    - Исследовать свойства событий и проверить правильность процедуры восстановления
- Получение большого числа событий
  - Анализ набора событий отобранных согласно какому-то критерию
  - Получение случайного набора событий





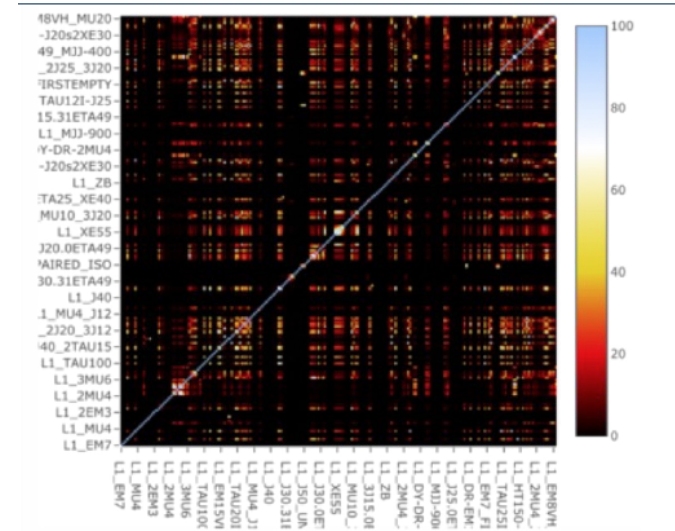
- В процессе обработки данных возможны сбои
  - Одни и те же события могут быть обработаны и/или записаны дважды с одними и тем же либо различными идентификаторами.
  - События может быть пропущены / потеряны
    - из-за ошибки ввода-вывода,
    - некорректной работы алгоритма
    - аварийного выхода программы
  - Особенно часто это происходит в первые годы эксперимента
- События могут повторяться/теряться при передаче данных
- Поломки оборудования также могут вызвать потерю данных
- На начальных этапах эксперимента доля «плохих» данных может быть значительной (десятки процентов)



- **Контроль согласованности данных при обработке**
  - Проверка количества данных во входном и выходном датасете
  - Сравнение идентификаторов событий на входе и на выходе
- **Индексирование новых «официальных» датасетов**
  - Позволяет отлавливать ошибки доступа, каталогов, поврежденные идентификаторы, etc... до начала физического анализа
  - Позволяет выявлять дублирующиеся данные
  - Позволяет выявлять проблемы в триггерной информации
- **Контроль согласованности и целостности данных**
  - НЕ позволяет контролировать физические параметры
  - НЕ гарантирует что обнаруженные ошибки будут исправлены



- Триггерная информация в EventIndex позволяет
  - Подсчитывать число событий удовлетворяющих определенным наборам триггерных условий (триггерным цепочкам)
  - Строить матрицы перекрытия триггерных цепочек (trigger overlaps)
    - Полезно для оптимизации триггерных таблиц и выбора триггерных условий для анализа.
  - Отбирать события на основании комбинации триггерной информации и других свойств (например мгновенной светимости)
    - Для подсчета частоты появления таких событий
    - Отобранные события можно использования в физических анализах

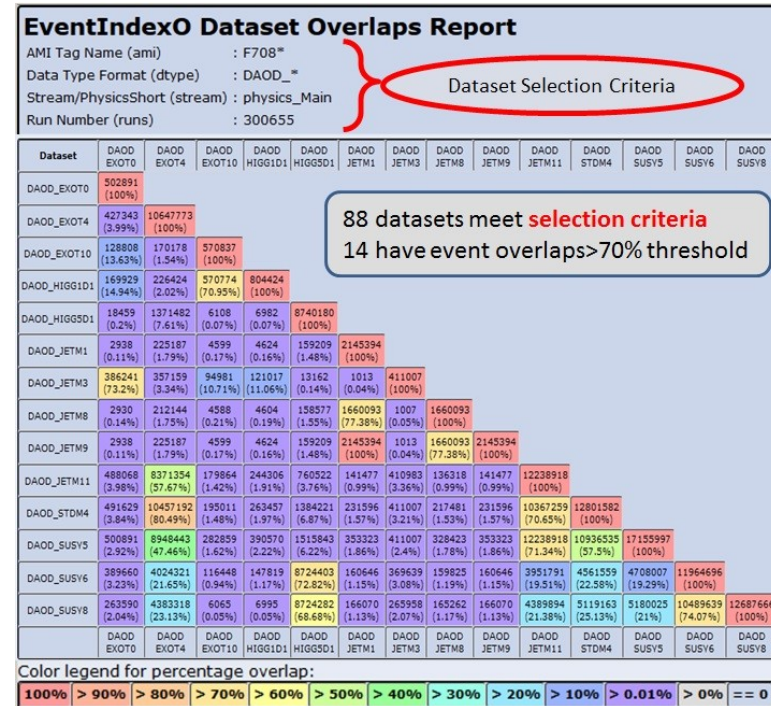




# Перекрытия производных данных



- При создании производных ("derivation") датасетов используются все полностью реконструированные события из, например, AOD датасета, а на выход поступает только часть событий, пригодных для одного или нескольких анализов.



- EventIndex позволяет рассчитывать долю перекрытий наборов событий в DAOD датасетах и таким образом оптимизировать их распределение и использование дискового пространства.

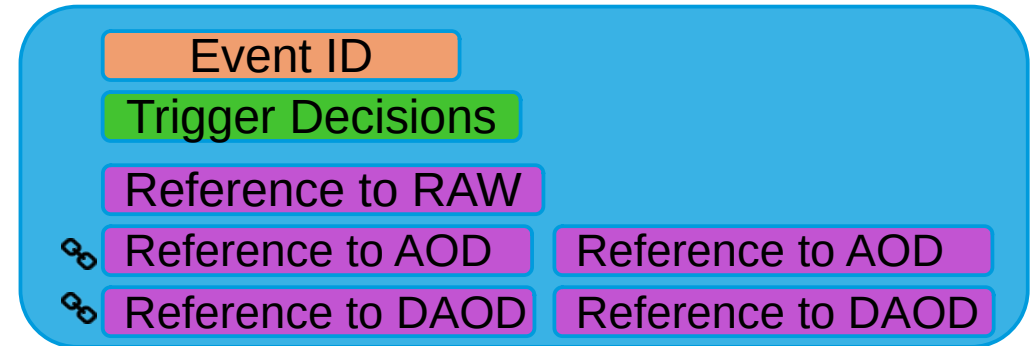
- Речь идет о петабайтах!



- Запись о каждом событии содержит следующие поля:

- Идентификаторы события

- Run and event number
- Trigger stream
- Дополнительная информация
  - Luminosity block
  - Bunch Crossing ID (BCID)



- Информация о триггерных решениях

- Триггерные маски и декодированные триггерные цепочки

- Указатели на файлы, в которых содержится информация о событии, в различных форматах и версиях обработки

- Информация о файлах добавляется по мере накопления данных



# Архитектура EventIndex





- Первичная обработка данных происходит на кластере Tier-0 в CERN. Реконструированные события записываются в AOD
- По готовности файлов, запускаются задачи (*Producer Transformation*) которые извлекают из них информацию для EventIndex.
  - Информация о RAW файлах берется из дочерних AOD
- Producer сохраняет полученную информацию в виде файла
  - проверяется наличие дублирующихся событий
  - выходной файл отправляется во временное хранилище Object Store
  - Сообщение с отчетом о выполненной задаче и ссылкой на файл отправляется программе-диспетчеру EventIndex Supervisor



- EventIndex также собирает информацию о датасетах которые создаются в результате обработки и моделирования данных на серверах распределенной вычислительной системы Grid
- В качестве источника данных о появлении и характеристиках новых наборов данных используются информационные системы AMI и Rucio.
  - Ежедневно проводится опрос системы AMI на предмет появления новых датасетов, их характеристики уточняются с помощью Rucio.
- Полученных списки проверяются с отсевом: dataset может быть поврежден, неполон или быть вовсе пустым
- Отфильтрованные списки передаются в Supervisor

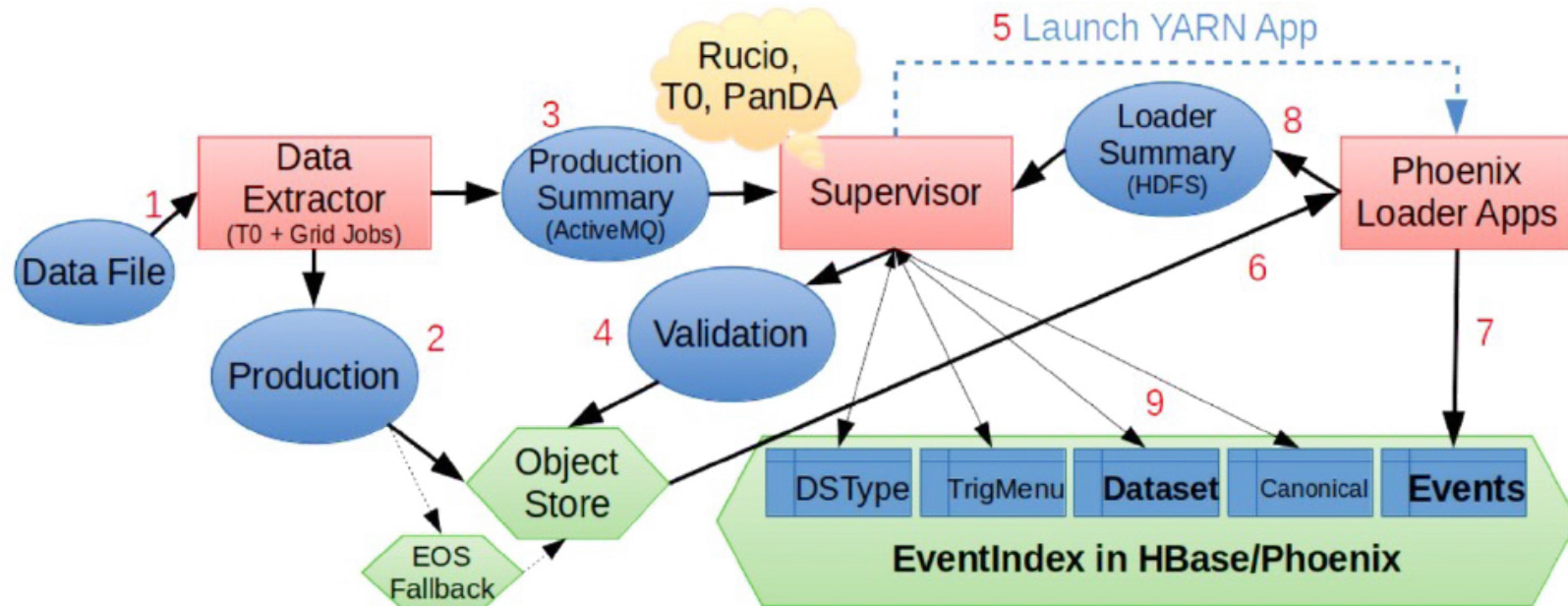




# EventIndex Supervisor



- EI Data Collection Supervisor – программа которая контролирует получение данных EventIndex от producer-ов и организует их передачу в постоянное хранилище





- Supervisor получает сообщения которые посылают задачи извлекающий информацию для EventIndex
  - Имя датасета, идентификаторы задачи, местоположение объекта с файлами eventIndex, id обработанных файлов и число событий
- Проверяется соответствие этих данных с данными Rucio
- Supervisor запускает программу Loader, которая
  - считывает данные из Object Store ,
  - записывает их в таблицы событий EventIndex
  - Записывает отчет о выполнении в HDFS для Supervisor
- Supervisor записывает данные в таблицу датасетов
- Supervisor можно контролировать через web интерфейс



- EventIndex разрабатывался в 2012-2013 гг для нужд Run 2 в качестве замены системы на основе Oracle, не справлявшейся с возросшими объемами данных.
- Для ядра системы было решено использовать нереляционные (NoSQL) базы данных из семейства Hadoop.
- При этом часть информации копировалась в таблицы Oracle для реализации быстрых запросов EventPicking
- Использование самых передовых на тот момент технологий обработки и хранения данных позволило обеспечить удовлетворительную производительность в течении всего периода работы начиная с запуска в середине 2015 года.





- Однако ожидаемое существенное увеличение потока и темпа поступления данных в Run 3 и Run 4 потребовало соответствующей модернизации системы EventIndex с использованием новой платформы хранения данных
- Создана новая версия использующая Apache HBase для хранения данных и Apache Phoenix для организации запросов и записи данных
- Это позволяет объединить функции хранилища больших данных в HBase с возможностью выполнения SQL-запросов через интерфейс Phoenix



- Apache HBase принадлежит к семейству HADOOP.
  - Нереляционная (NoSQL) база данных с открытым кодом
  - Распределенное и масштабируемое хранилище данных
- HBase организует данные в виде таблиц
  - Ряды таблицы имеют уникальный ключ-идентификатор
  - Каждый ряд может иметь свою схему
  - Данные в каждом ряду могут быть сгруппированы в заранее заданное семейства колонок
  - Доступ к значениям возможен по ключу и имени колонки



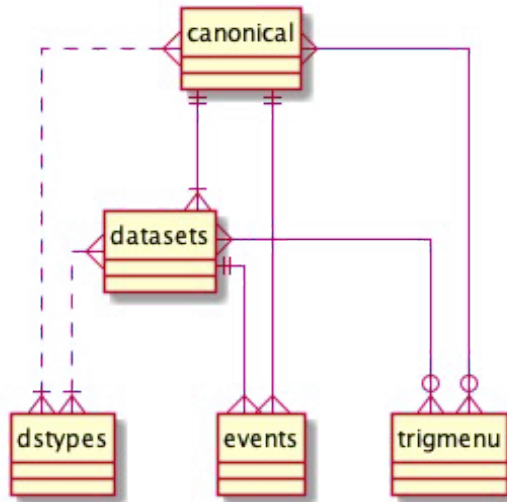
- Свойства такой организации (RowKey):
  - Относительно небольшой размер
  - Позволяет уникально идентифицировать события
  - Позволяет производить поиск по интервалам
  - При росте таблиц используется гомогенное пространство RowKey
- Для тривиальной модели достаточно простой организации данных в виде пар ключ:значение, характерных для NoSQL
- Добавление функциональности SQL к HBase дает дополнительные преимущества:
  - Структурированные данные проще понимать и поддерживать
  - Стандартная логика запросов, пригодная для сложных запросов



- Транзакционная система и операционная аналитика для HBase
  - Принимает SQL запросы
  - Преобразует их в серии сканирований в HBase
  - Напрямую использует API HBase, а также сопроцессоры и специальные фильтры
  - Выдает результаты запросов в стандарте JDBC
- Дизайн RowKey адаптируется к типам и размерам Phoenix
  - за счет небольшого уменьшения производительности
- Phoenix позволяет использовать поля RowKey в запросах, при этом они сохраняются как один объект в HBase



# Схемы Phoenix для EventIndex



dstypes
<ul style="list-style-type: none"> <li>type : tinyint</li> <li>name : varchar(20)</li> </ul>
id : smallint

trigmenu
<ul style="list-style-type: none"> <li>smk : integer</li> <li>type : tinyint</li> <li>name : varchar(200)</li> </ul>
id : smallint

events
<ul style="list-style-type: none"> <li>dspid : integer</li> <li>dstypeid : smallint</li> <li>eventno : bigint</li> <li>seq : smallint</li> </ul>
a.tid : integer a.sr : binary(32) a.mcc : integer a.mcw : float
b.pv : binary(34) ARRAY[]
c.lb : integer c.bcid : integer c.lpsk : integer c.etime : timestamp c.id : integer c.tbp : smallint[] c.tap : smallint[] c.tav : smallint[]
d.lb1 : integer d.bcid1 : integer d.hpsk : integer d.lph : smallint[] d.ph : smallint[]

datasets
<ul style="list-style-type: none"> <li>runno : integer</li> <li>project : varchar(200)</li> <li>datatype : varchar(200)</li> <li>streamname : varchar(200)</li> <li>prodstep : varchar(200)</li> <li>version : varchar(200)</li> <li>tid : integer</li> <li>dspid : integer</li> <li>dstypeid : smallint</li> </ul>
smk : integer events_rucio : bigint files : integer events : bigint events_uniq : bigint events_dup : bigint files_dup : integer rank : smallint status : varchar(200) rucio_at : timestamp updated_at : timestamp dups_at : timestamp trigger_at : timestamp is_open : boolean is_deleted : boolean has_raw : boolean has_trigger : boolean prov_seen : smallint ARRAY[] sr_cnt : varchar(200) sr_clid : varchar(200) sr_tech : varchar(200)

canonical
<ul style="list-style-type: none"> <li>runno : integer</li> <li>project : varchar(200)</li> <li>datatype : varchar(200)</li> <li>streamname : varchar(200)</li> <li>prodstep : varchar(200)</li> <li>version : varchar(200)</li> <li>dspid : integer</li> <li>dstypeid : smallint</li> </ul>
smk : integer events_rucio : bigint files : integer events : bigint events_uniq : bigint events_dup : bigint files_dup : bigint rank : smallint status : varchar(200) rucio_at : timestamp updated_at : timestamp dups_at : timestamp trigger_at : timestamp is_open : boolean is_deleted : boolean has_raw : boolean has_trigger : boolean prov_seen : smallint ARRAY[] is_rucio_eq : boolean sr_cnt : varchar(200) sr_clid : varchar(200) sr_tech : varchar(200)





- Ключи строк (Row keys). Они должны :
  - Содержать самую необходимую информацию
  - Иметь минимально возможный размер
- В EI выбран формат: `dspid.dstypeid.eventno.seq`
  - ID датасета → `dspid`
  - тип данных → `dstypeid`
  - номер события → `eventno`
  - номер дубля события → `seq`
- Семества данных:
  - A: расположение события (и MC info)
  - B: происхождение события
  - C: информация триггера 1 уровня (L1).
  - D: данные триггера высокого уровня
    - (EF of HLT) и L2 для Run 1

<code>dspid</code>	<code>integer</code>	<code>NOT NULL</code>	<code>,</code>
<code>dstypeid</code>	<code>smallint</code>	<code>NOT NULL</code>	<code>,</code>
<code>eventno</code>	<code>bigint</code>	<code>NOT NULL</code>	<code>,</code>
<code>seq</code>	<code>smallint</code>	<code>NOT NULL</code>	<code>,</code>
<code>a.tid</code>	<code>integer</code>		<code>,</code>
<code>a.sr</code>	<code>binary(24)</code>		<code>,</code>
<code>a.mcc</code>	<code>integer,</code>		
<code>a.mcw</code>	<code>float,</code>		
<code>b.pv</code>	<code>binary(26) array</code>		<code>,</code>
<code>c.lb</code>	<code>integer</code>		<code>,</code>
<code>c.bcid</code>	<code>integer</code>		<code>,</code>
<code>c.lpsk</code>	<code>integer</code>		<code>,</code>
<code>c.etime</code>	<code>timestamp</code>		<code>,</code>
<code>c.id</code>	<code>bigint</code>		<code>,</code>
<code>c.tbp</code>	<code>smallint array</code>		<code>,</code>
<code>c.tap</code>	<code>smallint array</code>		<code>,</code>
<code>c.tav</code>	<code>smallint array</code>		<code>,</code>
<code>d.lb1</code>	<code>integer,</code>		
<code>d.bcid1</code>	<code>integer,</code>		
<code>d.hpsk</code>	<code>integer,</code>		
<code>d.lph</code>	<code>smallint array,</code>		
<code>d.lpt</code>	<code>smallint array,</code>		
<code>d.lrs</code>	<code>smallint array,</code>		
<code>d.ph</code>	<code>smallint array,</code>		
<code>d.pt</code>	<code>smallint array,</code>		
<code>d.rs</code>	<code>smallint array,</code>		
<code>CONSTRAINT events_pk PRIMARY KEY</code>			
<code>(dspid, dstypeid, eventno, seq)</code>			



- Вспомогательные таблицы для хранения идентификаторов датасетов и данных учета
- Таблицы датасетов:
  - Расположение датасетов
  - Сгенерированные идентификаторы (dspid)
  - Статус импорта
  - Некоторая информация о метаданных: количество всех событий, уникальные события, дублирование
- Таблица типов данных:
  - типы данных
    - (RAW, EVENT, AOD, DAOD)
  - подтипы для производных

dstypes
<ul style="list-style-type: none"><li>• type : tinyint</li><li>• name : varchar(20)</li></ul>
id : smallint

trigmenu
<ul style="list-style-type: none"><li>• smk : integer</li><li>• type : tinyint</li><li>• name : varchar(200)</li></ul>
id : smallint

datasets
<ul style="list-style-type: none"><li>• runno : integer</li><li>• project : varchar(200)</li><li>• datatype : varchar(200)</li><li>• streamname : varchar(200)</li><li>• prodstep : varchar(200)</li><li>• version : varchar(200)</li><li>• tid : integer</li><li>• dspid : integer</li><li>• dstypeid : smallint</li></ul>
smk : integer events_rucio : bigint files : integer events : bigint events_uniq : bigint events_dup : bigint files_dup : integer rank : smallint status : varchar(200) rucio_at : timestamp updated_at : timestamp dups_at : timestamp trigger_at : timestamp is_open : boolean is_deleted : boolean has_raw : boolean has_trigger : boolean prov_seen : smallint ARRAY[] sr_cnt : varchar(200) sr_clid : varchar(200) sr_tech : varchar(200) name : varchar(250)



- Клиент для поиска событий Atlas EventIndex (event lookup):
  - Пользователь запрашивает службу поиска событий, она отправляет запрос в EventIndex и возвращает `htpекmfnfn` пользователю
- Пользователь предоставляет список событий для поиска
- Результат может быть представлен в различных форматах:
  - простой список
  - подробный список с дополнительными параметрами
  - JSON (полезно для вызова из программ, например из `panDA`)
- Пример:

```
$ event-lookup -c name 348895:3613825 -d AOD  
348895 3613825 0d1c27ff-c5b2-6b46-a0fa-edbdc51c6b77 AOD data18_13TeV.00348895.physics_Main.merge.AOD.f937_m1972  
348895 3613825 78c6837f-f448-e811-ae3c-44a8420a7621 RAW data18_13TeV.00348895.physics_Main.merge.AOD.f937_m1972
```



## Event Lookup

List of 'runnumber evtnumber' (-e)

AMI tag (-p) (substring match)

Stream name (-s)

Data type (-d)

GUID type (-f)  AOD  ESD  RAW  all

-api  simple  rich  indexer  mc (indexer)

-details  event  type  id  dataset  rich output

-email  (implies asynchronous execution)

<  >

-e 00278880 558085589, 00278880 210318172

-details typenullnullnull

-api indexer

```
6 guids found for 1 runs with 2 events, 0 guids missing, 0s spend
148C774B-B5BB-BF41-9346-F82F213EC4A6 StreamDAOD_HIGG1D1
6C2A1CB8-8F8F-4B42-9FAF-AC55B0957CFD StreamAOD
886C7D3A-3E56-E511-9D00-44A8420A5EB7 StreamRAW
99ACB5FB-11FB-464F-B60F-E496213CABF1 StreamDAOD_HIGG1D1
EBC919C0-BAE8-9F42-BA63-21D782C0FAB7 StreamAOD
FC28C1E6-1D56-E511-AEDE-44A8420A8576 StreamRAW
```



- Trigger Tool из предыдущего EventIndex:
  - Пользователь может искать события, которые удовлетворяют некоторым условиям запуска (прошли некоторые цепочки триггеров)
    - Пользователь может указать список цепочек триггеров для поиска
    - Также может быть предоставлен список вето
- Выполнение занимало существенное время, так как EventIndex должен был выполнять Map/reduce Job для огромного объема данных
  - Результатом является список событий
- Использовался для расчета матриц перекрытий триггеров
  - Выполнялся автоматически на новых наборах данных



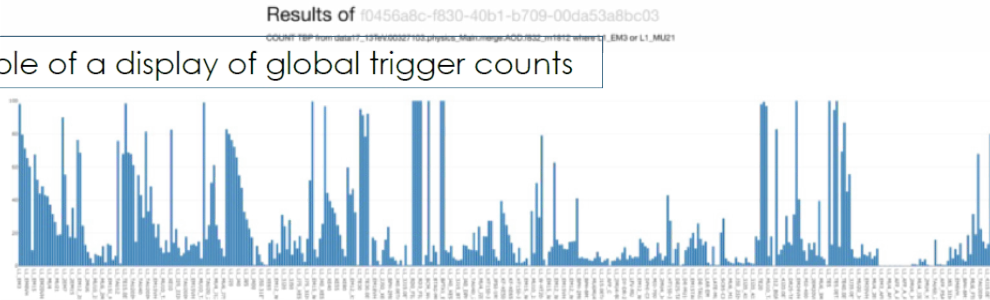
# Web-интерфейс: счетчик триггеров



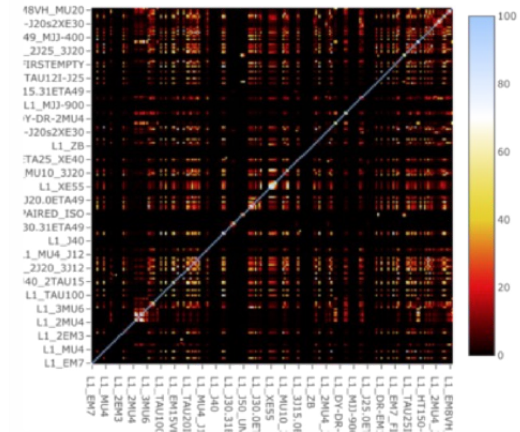
The web form to use the Trigger Counter

Run Number: 341312  
Stream: physics\_BphysLS  
Dataset: data17\_13TeV.00341312.physics\_BphysLS.merge.AOD.f903\_m1917  
Operation: COUNT  
Level: HLT  
Trigger: PH  
Expression:   
Group by: No  
Limit:   
  
[Show/Hide Triggers](#) +

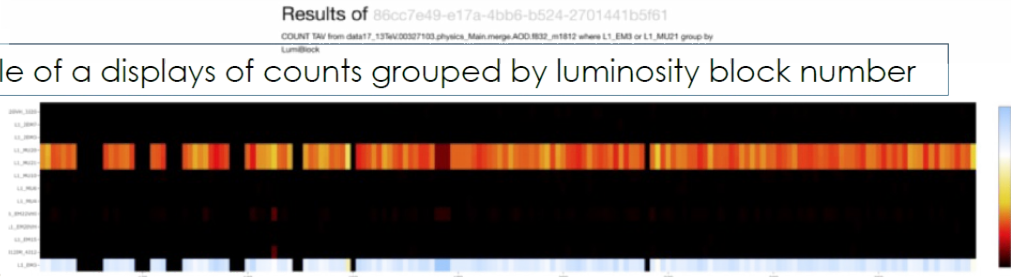
Example of a display of global trigger counts



The heat map of trigger overlaps within a given dataset



Example of a displays of counts grouped by luminosity block number



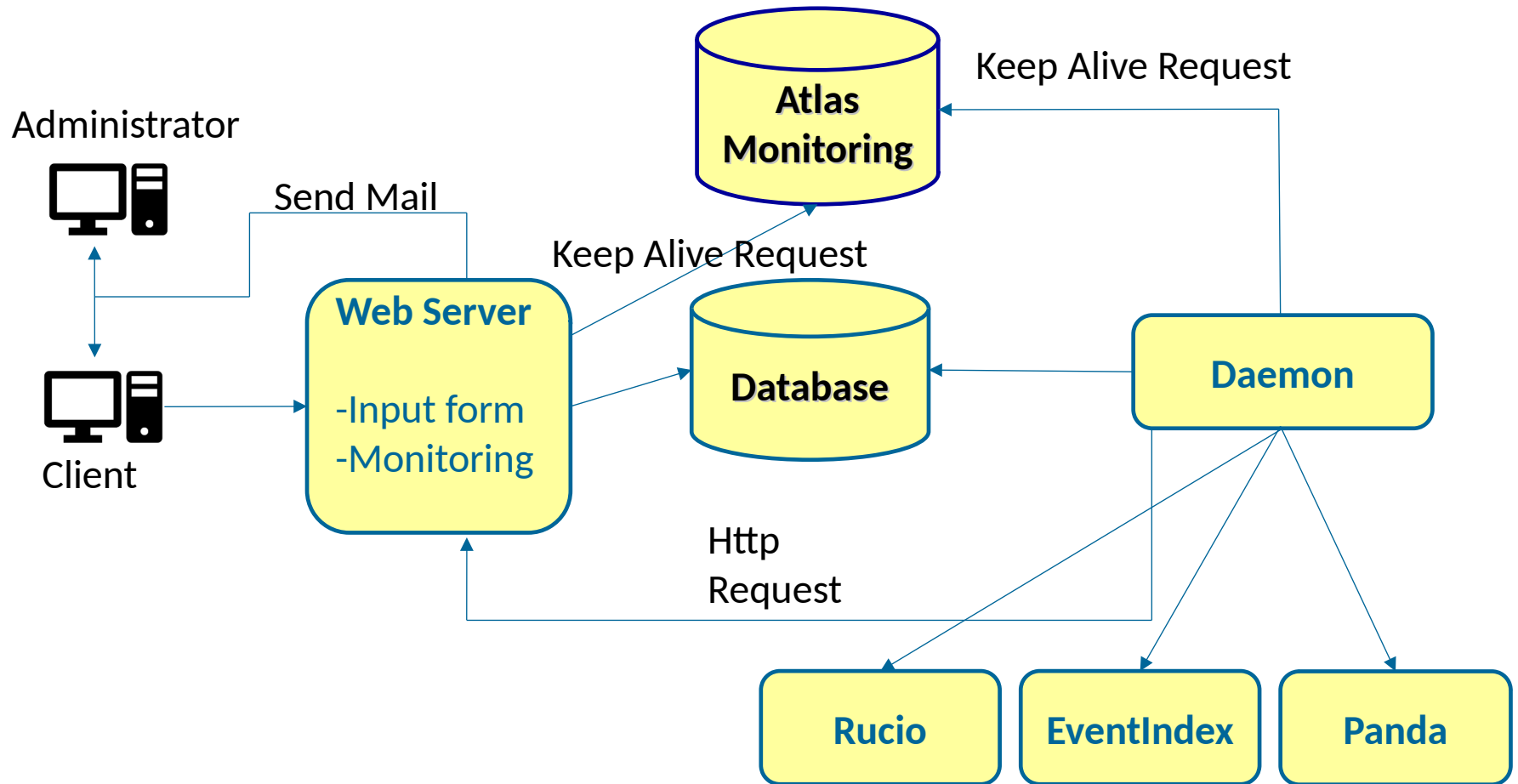
5 July 20



- Основная цель системы — автоматизировать операции, необходимые для массового выбора событий (event picking):
  - Графический интерфейс (GUI) получает от пользователя список событий и дополнительные параметры (тип данных для поиска, тип данных для извлечения, поток, тег AMI и т. д.).
  - EPS разбивает список (списки) по run number, запрашивает интерфейс командной строки (EventIndex Query CLI), отправляет задания PanDA по сбору событий, повторяет их, если время ожидания истекло (например если события находятся на ленте), проверяет выходные данные, информирует пользователя о завершении.
  - Пользователь может следить за прогрессом с помощью графического интерфейса.



# Архитектура Event Picking Service







# Архитектура Event Picking Service



← → ↻ atlas-event-picking.cern.ch/ep\_debug/service?rt=view&sid=110

★ Bookmarks 📁 Темы диссертаций... 🌐 Еврокубки в цифр... 🗂 Флибуста | Книжно... VS | VegaSat.ru 📄 Open Source softw... 🌐 Directory Bind -... 🔄 Hybrid MPI/OpenM... 📄 api.lib

Event Picking Service v. 0.9.01

Common Info New Request Requests Monitoring Data Tables

* Data format:	RAW	▼
Project:	Select Data format	Carlo
Trigger stream:	RAW	
	EVNT	
	AOD	
AMI tag:		
File containing run and event numbers :	Выбрать файлы	Файл не выбран

* User (client) name:	<input type="text"/>
* User e-mail (identifier):	<input type="text"/>

Submit request Clear form

**Data format** - required field - choose one of several formats.

**Project** - enables (if real data) or disables (if Monte Carlo) the field "**Trigger stream**".

**Trigger stream** - the field is active if real data.

**AMI tag** - the field is active if the data format is not "RAW".

**File containing run and event numbers** - required field - a text file containing strings like "XXX YYYY", first column - **run number**; second column - **event number**.

**Caution !!** If the file format does not meet this requirement, the server will return an error.

**User name** - required field

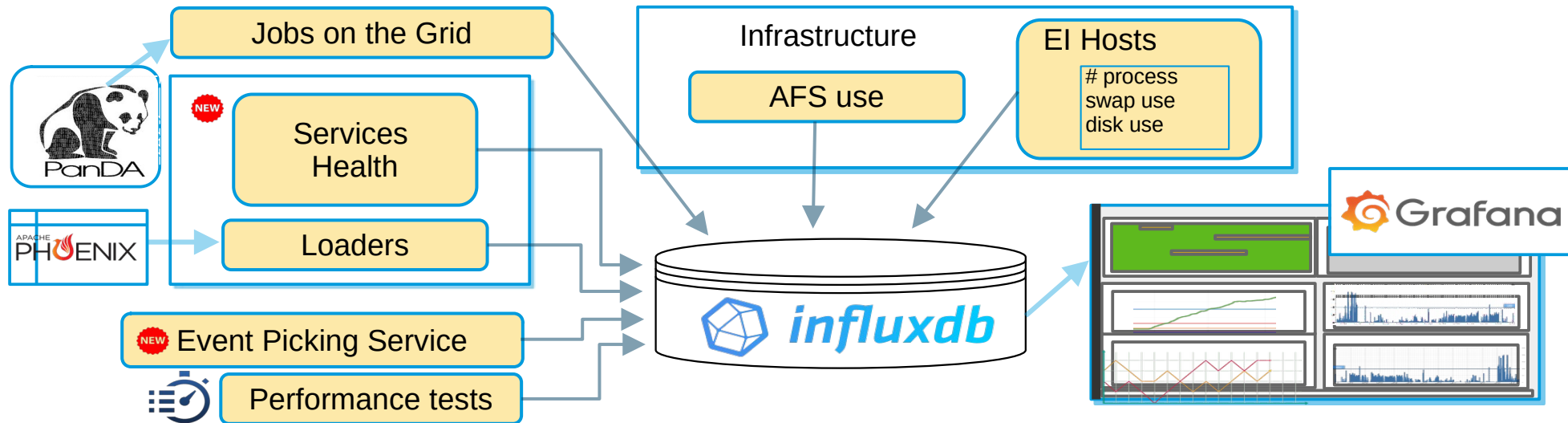
**User e-mail** - required field - e-mail, to which the results of the request processing will be sent.



# Мониторинг системы

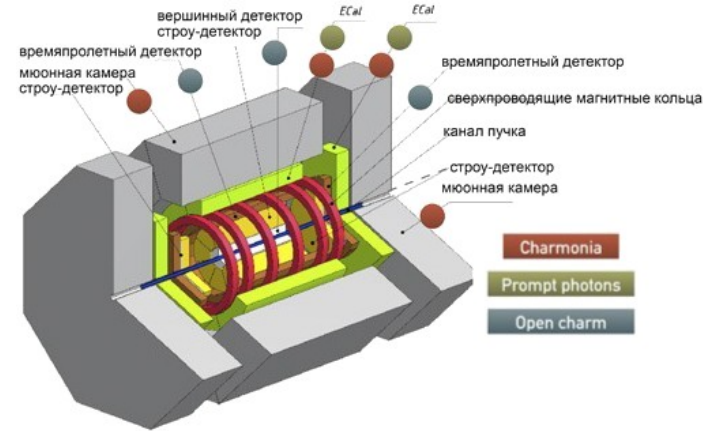
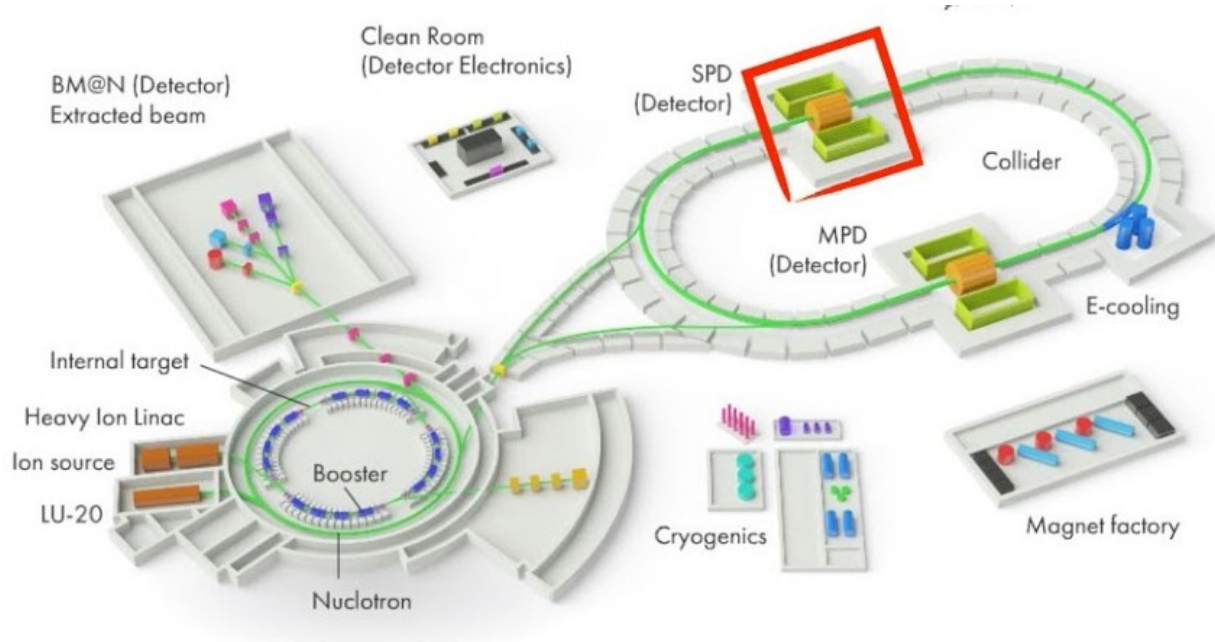


- Данные о состоянии и производительности системы (~15k значений в день) собираются разными способами (сканирование журналов и HDFS, REST запросы и анализ WEB-страниц)
- Сбор и обработка информации управляется планировщиком
- Результаты загружаются в InfluxDB / Grafana через REST





## Spin Physics Detector





- Предполагаемые объемы данных с SPD достаточно велики
  - 20 ГБ/с => ~ или 200 ПБ/год (необработанные (RAW) данные)
  - ~30 миллиардов событий в год
  - это величины того же порядка или больше, чем в ATLAS
- После реконструкции и обработки для использования в физическом анализе добавятся еще несколько экземпляров каждого события
- Дополнительно ожидается сравнимый объем данных моделирования (MC)
- Данные о событиях в различных форматах записываются в файлы, организованные в датасеты различных типов



- Данные временных слайсов, записанные в формате RAW
  - Временные данные, нет нужды индексировать
- После онлайн-фильтра → формат filtered RAW:
  - Результаты быстрой реконструкции событий,
  - распакованные RAW хиты, блоки RAW данных
- После оффлайн реконструкции → датасеты AOD и ESD:
  - AOD: Физические объекты: частицы, треки, кластеры ...
  - ESD: + необработанные хиты
- Деривации для использования в конкретных анализах (DAOD)
  - Десятки типов
  - Как правило с меньшим числом событий и размером события



- Генератор событий создает датасеты EVGEN
  - Содержат "истинные" параметры частиц (MC truth)
- После моделирования → RAW MC события
  - MC truth + распакованные RAW хиты
- После оффлайн реконструкции → датасеты AOD и ESD
  - То же, что и для реальных данных + MC truth
- Как и для реальных данных, деривации (DAOD)
- Данные моделирования появятся задолго до запуска установки, что позволит провести тестирование системы.




- Регулярно будет производиться поворотная обработка данных, каждый раз добавляя новые экземпляры датасетов
- Данные будут распределенный по вычислительным серверам, их число будет меньше чем на LHC
- Для эффективного доступа к событиям, необходима база данных база данных, содержащая ссылки на экземпляры событий в распределенном хранилище данных
- То есть, **EventIndex**
- При разработке EventIndex предполагается опираться на опыт работающей на Atlas системы, адаптируя его к реалиям SPD



- Отбор событий (Event Picking)
  - получить событие в нужном формате и обработанное конкретной версией программного обеспечения
- Подсчет и выборка события
  - на основании параметров событий
  - на основании данных онлайн триггера
- Проверка целостности и согласованности данных
- Создание матриц перекрытия датасетов
- Быстрый просмотр данных
  - Поиск датасетов, составление отчетов, проверка датасетов
- Создание виртуальных датасетов **OPTIONAL**





- Идентификаторы события
  - Run number
  - Event number
  - Frame ID
  - Bunch Crossing ID (BCID)
- Информация онлайн триггера
- Указатели на **постоянные** файлы, содержащие информация о событии, в различных форматах и версиях обработки
  - Информация о файлах добавляется по мере накопления данных
- Параметры события (например поляризация) 



- Первичные поля (при создании записи)

- 'RunNumber' integer YMMDDhh[0..9] 280308060
- 'EventNumber' integer [-2147483648 .. +2147483647] 474 836
- 'FrameID' integer [-2147483648 .. +2147483647] 5 000
- 'OFStream' varchar(128) 'Main'

**OPTIONAL**

- 'OFData' JSONB ...
- 'GUIDRaw' UUID '78c6837f-f448-e811-ae3c-44a8420a7621'

- Дополнительные (по мере обработки и создания датасетов)

- 'GUID\_AOD' UUID '0d1c27ff-c5b2-6b46-a0fa-edbdc51c6b77'
- 'ImportantValue' real 3.14
- 'GUID\_DAOD\_PHYS' UUID 'aa321c54-ca5e-b748-845d-02f30074f8de'

**OPTIONAL**

**OPTIONAL**

- 'ImportantNumber' smallint 42



- Разработку EventIndex естественно начать с выбора платформы для хранения и управления данными.
- Нужна Open Source СУБД позволяющая быстро и надежно оперировать терабайтными объемами данных
- Рассматриваются несколько вариантов



- ClickHouse от Яндекс  ClickHouse

- столбцовая СУБД для онлайн обработки аналитических запросов (OLAP)

- Apache HBase



- распределенная, колоночно-ориентированная, мультиверсионная база типа «ключ-значение».

- PostgreSQL



- объектно-реляционная система управления базами данных (СУБД)



- Столбцовое хранение данных — данные считываются только из нужных колонок, и однотипная информация эффективно сжимается
- Поддержка приближённых вычислений на части выборки — снижается число обращений к жёсткому диску, что ещё больше повышает скорость обработки данных;
- распараллеливание операций как в пределах одного сервера на несколько процессорных ядер, так и в рамках распределённых вычислений на кластере за счёт механизма шардирования;
- линейная масштабируемость — есть возможность построить кластер очень большого размера



- Разреженный индекс делает ClickHouse плохо пригодным для точечных чтений одиночных строк по своим ключам.
- Возможность изменять или удалять ранее записанные данные с низкими задержками и высокой частотой запросов не предоставляется. Есть массовое удаление и изменение данных для очистки более не нужного
- Отсутствие уникальности первичных ключей, что вместе с предыдущим пунктом в нашем случае осложняет добавление новых полей в EventRecord



- Данные организованы в таблицы, проиндексированные первичным ключом, который в Hbase называется RowKey.
- Для каждого RowKey ключа может храниться неограниченный набор атрибутов (или колонок).
- Колонки организованы в группы колонок, называемые Column Family.
  - Как правило в одну Column Family объединяют колонки, для которых одинаковы паттерн использования и хранения.
- Для каждого атрибута может храниться несколько различных версий. Разные версии имеют разный timestamp.



- Записи физически хранятся в отсортированном по RowKey порядке. При этом данные соответствующие разным Column Family хранятся отдельно, что позволяет при необходимости читать данные только из нужного семейства колонок.
- Атрибуты, принадлежащие одной группе колонок и соответствующие одному ключу физически хранятся как отсортированный список. Любой атрибут может отсутствовать или присутствовать для каждого ключа, при этом если атрибут отсутствует — это не вызывает накладных расходов на хранение пустых значений.
- Список и названия групп колонок фиксирован и имеет четкую схему.



- PostgreSQL это а объектно-реляционная СУБД с открытым исходным кодом.
- Поддержка многочисленных типов данных
  - Численные, булевые, символьные, составные, сетевые типы данных, перечисление, типы «дата/время», массивы, etc.
- Поддержка JSON, что позволяет использовать schema-less данные
  - Встроенные специализированные JSON операторы и функции
- Поддержка пользовательских объектов и их поведения, включая типы данных, функции, операции и индексы.
- Индексирование: частичное, функциональное, GiST, GIN, BRIN
- Функции виртуальных таблиц, Материализованные представления
- Возможность добавления/изменения столбцов





- Обеспечивает высокую надежность и производительность:
  - соответствие принципам ACID (атомарность, изолированность, непротиворечивость, сохранность данных)
  - Многоверсионный контроль конкурентных транзакций и изоляция транзакций
    - возможность изменения баз данных одновременно разными пользователями.
    - минимизирует блокировки данных и позволяет увеличить производительность
  - Журналы опережающей записи (**Write Ahead Logging**)
    - фиксирующая все изменения до их фактического применения.
  - Резервное копирование и восстановление
  - Возможность восстановления базы данных Point in Time Recovery
    - Откат состояние базы данных к предыдущему стоянию используя WAL.



- Основные ограничения

Максимальный размер БД	Неограничен
Максимальный размер таблицы	32 TB
Максимальный размер строки	1.6 TB
Максимальный размер поля	1 GB
Максимальное количество строк в таблице	Неограниченно
Максимальное количество столбцов в таблице	250-1600
Максимальное количество индексов в таблице	Неограничено

