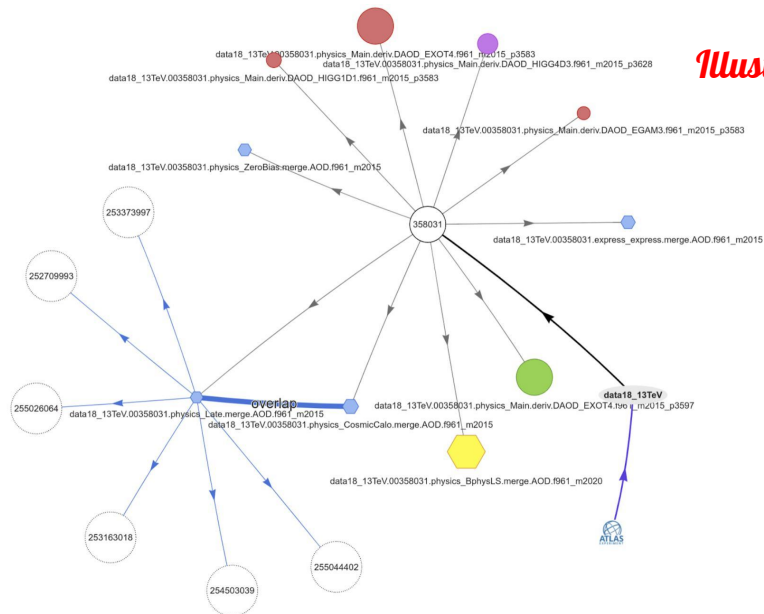


Julius Hrivnac (ATLAS Orsay) on behalf of the ATLAS Collaboration

Illustrating Graph Databases potential on new ATLAS Event Index prototypes.



- *H&P storage*
- *Graph Databases*
- *Atlas Event Index*
 - *Graph Database Prototypes*
- *Graph Databases for H&P*
 - *What they can do for us*

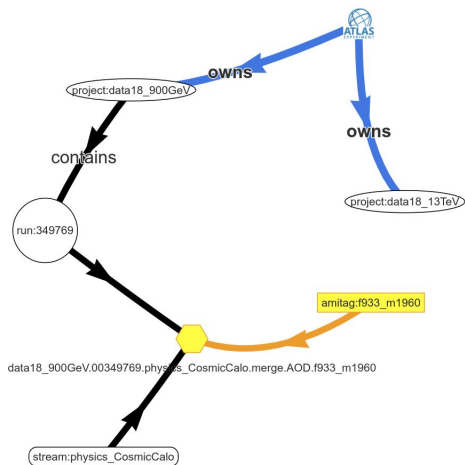


HEP Storage

- Traditional data structures in HEP:
 - tuples (tables)
 - trees
 - nested tuples (trees of tuples)
 - relational (SQL-like)
- Schema-based or schema-less
- **But many of HEP data are graph-like & schema-less**
 - **Entities with relations**
- Not handled by standard tree-ntuple storage
 - Relations should be added and interpreted outside storage
- Not well covered by relational (SQL) databases
 - We need to add new relations, not covered by schema
- Difficult to manage by Object Oriented (OO) databases or serialisation
 - Problem to distinguish essential relations from volatile ones

Graph Databases

- Storing Graphs in a database
- **Graph = (Vertexes, Edges), $G = (V, E)$**
- Vertices and Edges have properties



- Graph databases have existed for a long time
 - Matured only recently thanks to Big Data & AI (Graph NN)
 - Very good implementations & (de-facto) standards available
 - Rapid evolution
- **Moving essential structure from code to data**
 - Together with migration from imperative to declarative semantics
 - Things don't **happen**, but **exist**
 - Structured data with relations facilitates **Declarative Analyses**
- **Data elements appear in a Context**
 - Which simplifies understanding, analyses and processing
- The difference between SQL and Graph database is similar as between Fortran and C++/Java
 - On one side, a rigid system, which can be very optimized
 - On the other side, a flexible dynamical system, which allows expressing of complex structures
- Graph database is a synthesis of OO and SQL databases
 - Expressing web of objects without fragility of OO world
 - Capturing only essential relations, not an object dump



Graph DB: Languages

➤ Direct manipulation of Vertices and Edges

- Always available from all languages
- Doesn't use full graph expression power

➤ Cypher (and GQL)

- Pure declarative
- Inspired by SQL and OQL
 - But applied to schema-less database
- Available to all languages via JDBC-like API
 - Semantic mismatch, passed as String
 - There is a wall between coder and database, with a thin tunnel, only Strings can pass
- Coming from Neo4J
 - Accepted as a standard
 - Neo4J can be also used with Gremlin

```
MATCH (a:run)-[:has]->(b:dataset)
WHERE a.rnumber = 98765
RETURN b.name
```

*All three methods can be used to access the same database
(e.g. using Cypher for query and Gremlin for traversal).*

➤ Gremlin

- Functional syntax
- Originated from *Groovy*, but available to all languages supporting functional programming
- Integrated in the language

```
g.V().has('run', 'rnumber', 98765)
    .out('has')
    .values('name')
```



Atlas Event Index

- ATLAS Event Index Service keeps references to all real and simulated ATLAS events. Hadoop Map files and HBase tables are used to store the Event Index data, a subset of data is also stored in the Oracle database.
 - Contains information about **events, datasets, streams, runs**,... - and their relations.
- Several user interfaces are currently used to access and search the data. From the simple command line interface, through programmatical API to sophisticated Graphical Web Services.
- History
 - Original Event Index (EI) in Oracle
 - Too rigid (can't easily add columns, relations), other problems
 - Migrated to Hadoop
 - Map files in HDFS
 - Flexible
 - Too slow for searching (ok for processing)
 - Typeless
 - Partially migrated to HBase
 - Two tables: Catalog + Events
 - Tables contain a lot of ad-hoc relations (references to other entries)
 - **We have in fact implemented a poor-man's GraphDB on top of HBase**
 - Several prototypes developed to study next generation Event Index (to be fully deployed for Run-3, end 2020), using Graph Database in different ways:
 - Prototype storing all EI data directly in JanusGraph database over HBase storage
 - Prototype storing data in a HBase table with Phoenix SQL interface, Graph structure added via another auxiliary HBase table

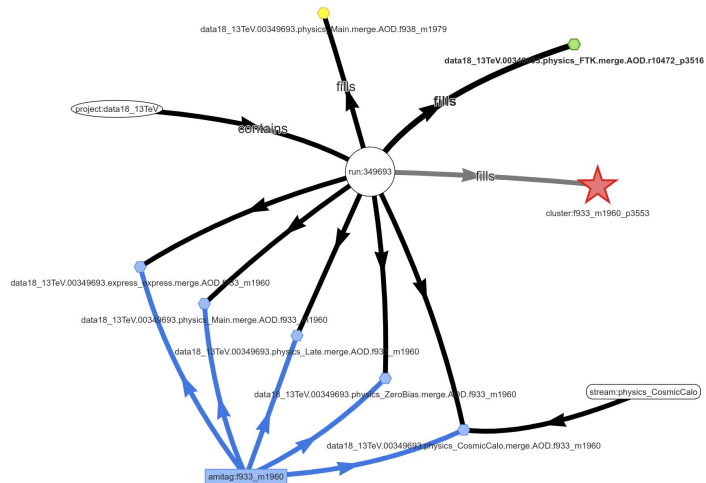
The ATLAS Event Index for LHC Run 3 - Track 4 on Th at 14:45

Prototype 1:

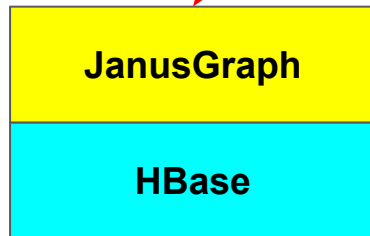
JanusGraph on top of HBase

Prototype 1

- Subset of data imported into full **JanusGraph** Database storing data in an **HBase** table
- Part of existing functionality implemented
 - Large part of code (handling relations and properties) moved to structure of the graph
- Most of the graphical part implemented
 - By standalone JS implementation



Gremlin API



Implementations Choices

Functional syntax with additional navigational semantics !

- De-facto standard language/api: **Gremlin**
 - Gremlin is a functional, data-flow language to **traverse a property graph**. Every Gremlin traversal is composed of a sequence of (potentially nested) steps. A step performs an atomic operation on the data stream. Every step is either a *map*-step (transforming the objects in the stream), a *filter*-step (removing objects from the stream), or a *sideEffect*-step (computing statistics about the stream).
 - Gremlin supports **transactional** & **non-transactional** processing in **declarative** or **imperative** manner.
 - Gremlin can be expressed in all languages supporting function composition & nesting.
 - Supported languages: **Java**, **Groovy**, **Scala**, **Python**, Ruby, Go, ...
- Commonly used framework: **TinkerPop**
- Leading implementation: **JanusGraph**
 - Supported storage backends: Cassandra, **HBase**, Google Cloud, Oracle BerkeleyDB
 - Supported graph data analytics: Spark, Giraph, Hadoop
 - Supported searches: Elastic Search, Solr, Lucene
 - Growing popularity of Neo4J
- Chosen visualisation: **visj.org**
 - Many others exist



Gremlin Syntax

Prototype 1

- Functional syntax
- **Functional & navigational** semantics
- Very intuitive, **no special syntax needed** (using existing functional syntax), easy integration.
- Database just accessed as objects with structure and relations.
 - Nested collections with links.
- Can use functional API (streams) and Lambda.
- **No semantic mismatch.**
 - Using one language.
- Came from **Groovy**
 - (Almost) identical for other supported languages (Python, Scala, Go,...).
- Both **search** and **traversal** steps.
- Search steps can be boosted by indexes.
- Functions can be loaded on server for faster execution.

```
# add a vertex 'experiment' with the name 'ATLAS'
g.addV('experiment').property('ename', 'ATLAS')
# add edges 'owns' from all vertices 'project' to vertex 'experiment' 'ATLAS'
g.V().hasLabel('project')
  .addE('owns')
  .from(g.V()
    .hasLabel('experiment')
    .has('ename', 'ATLAS'))
```

```
# show datasets with more events or number of events in an interval
g.V().has("run", "number", 358031)
  .out()
  .has("nevents", gt(7180136))
  .values("name", "nevents")
g.V().has("run", "number", 358031)
  .out()
  .has("nevents", inside(7180136, 90026772))
  .values("name", "nevents")
```

Event-Lookup function (server side UDF)

```
def el(run, event, g) {
  e = g.V().hasLabel('run')
    .has('rnumber', run)
    .out('fills')
    .out('keeps')
    .has('enumber', event)
    .values('guid')
  # all runs
  # selected run
  # all datasets filling that run
  # all events kept in that dataset
  # selected event
  # its guid
}
```

CLI command

```
curl -XPOST -d '{"gremlin":"el(run, event)"}' http://ei-gremlin-server.cern.ch:8182
```

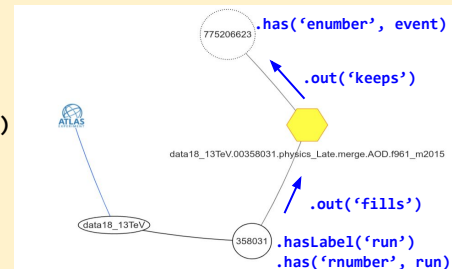
or using standard gremlin client

```
gremlin << EOF
```

```
:remote connect tinkerpops.server $janusgraph_home/conf/remote.yaml
```

```
el(run, event)
```

```
EOF
```





Prototype 2: Phoenix Project

Prototype 2

➤ Aim

- **Simple and flexible NoSQL storage (HBase)**
 - Good experience
- **Compatibility with other SQL database (Phoenix API)**
 - SQL used in other ATLAS components
- **Advantages of Graph Database**

➤ Solution

- Extend Phoenix/HBase storage with pure HBase storage
 - Sharing the same keys
- So keeping Phoenix advantages (speed, SQL interface) for RO data
- While opening for new possibilities, adaptability to changing environment
- Phoenix/HBase for static data, HBase for dynamic data

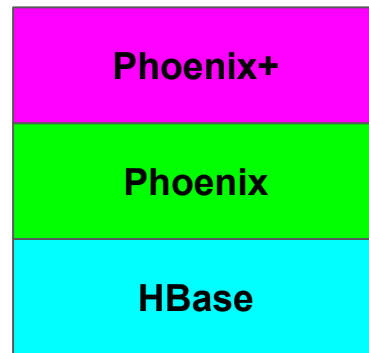
➤ Data stored in pure HBase table with Phoenix SQL API

➤ Graph structure added via aux HBase table

- 'Lazy Graphs': only created when needed

➤ Home-grown 'object API'

- With Gremlin subset API



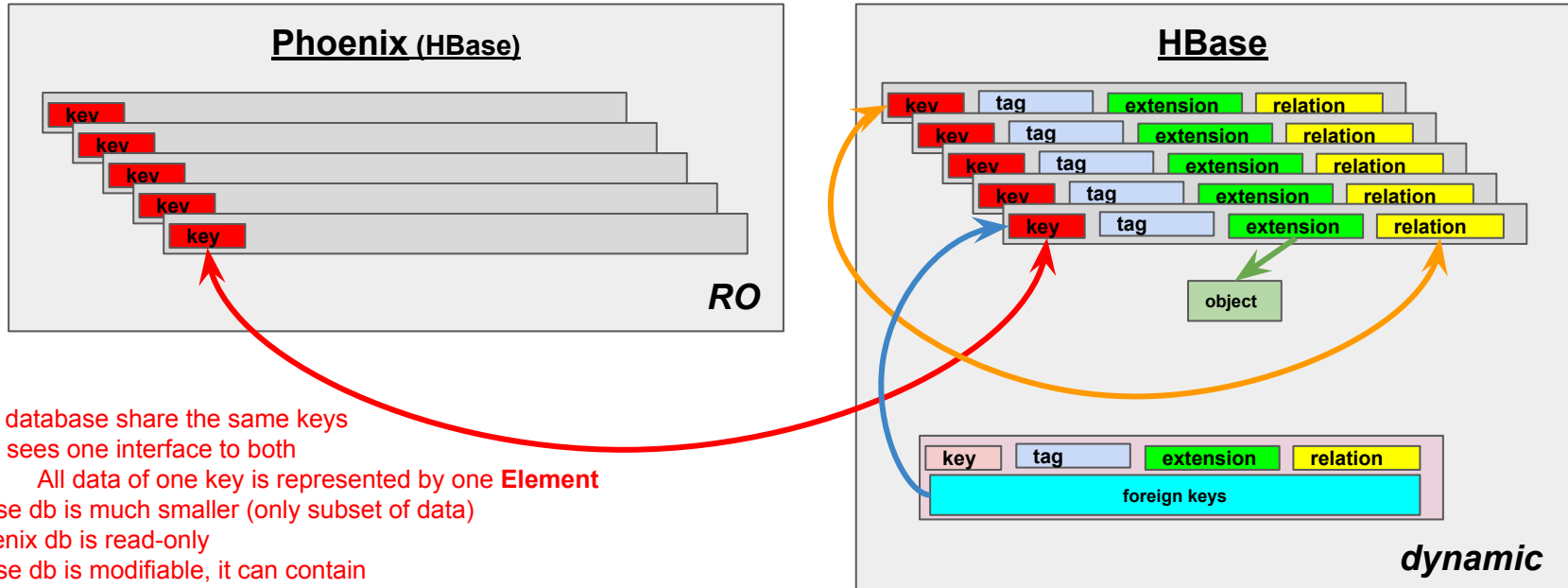
**Object API
& Gremlin**

SQL

Architecture

Prototype 2

Creating Graph structure on top of an SQL database - a user sees one database.

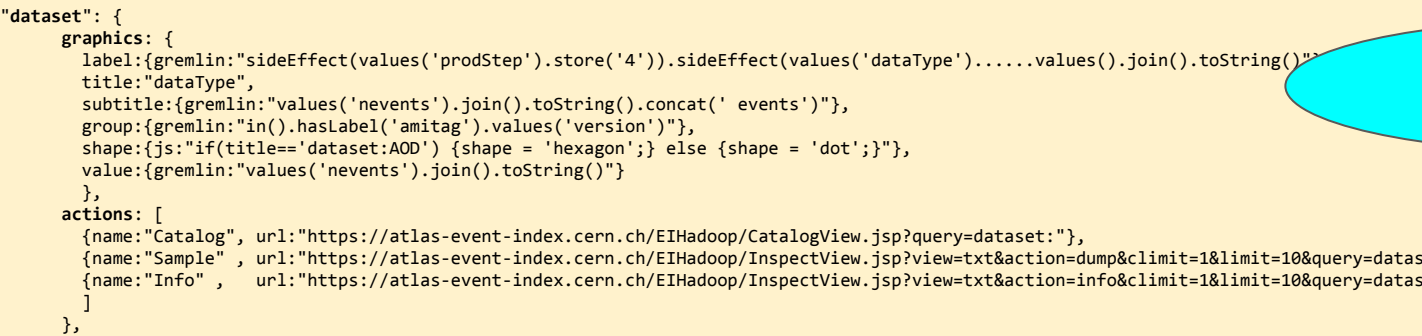


- Both database share the same keys
- User sees one interface to both
 - All data of one key is represented by one **Element**
- HBase db is much smaller (only subset of data)
- Phoenix db is read-only
- HBase db is modifiable, it can contain
 - **Simple Tags**
 - They can be also used in search filter
 - **Extensions** with any object
 - E.g. Trigger statistics and overlap, duplicated events list,...
 - **Relations** to other elements (like Graph DB)
 - E.g. overlaps between datasets

- HBase can also contain Elements without Phoenix partner: **Hubs**
 - They represent pre-defined **virtual collections** of Elements
 - E.g. Amitag, Stream, Run, Project,...
 - They can be extended and searched in the same way as other Elements
- Ad-hoc virtual collections can be build also using Tags



- Can display **any Gremlin-compatible** database
 - Visualisation can be customised via Stylesheet (JSON)
- Implemented **completely in JavaScript**
 - So doesn't need server-side application
 - Connects to standard Gremlin server to get JSON view of data



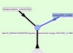
Direct API and REST Web Service also available.

WS GUI

Prototype 2

- Very generic
- Very thin layer on top of data
 - Structure is already there
- Hierarchical navigation
 - A'la Google Earth
- Show all available data, their relations and properties
 - And all available actions for them

1. Select period, run, ami tag, stream, dataset or event
 - a. Or any other Hub (virtual collection)
2. Get it, together with all related entities and information stored in EI
3. Each entity has a set of possible actions to perform
4. Each entity can be annotated
5. Detailed search is possible too



ATLAS Event Index

Problems or Questions ? - Ask [service manager](#) !

[ATLAS](#)

data09_2TeV

data10_7TeV

data09_900GeV

1

* 140541 140571 140737 140747 140836 140842 140953 140955 141266 141270 141374 141387 141689 141691 141695 141700 141748 141749 141755 141811 141818 141841 141994 141998 141999 142042 142065 142183 142144 142149 142154 142155 142157 142159 142161 142165 142166 142171 142174 142189 142190 142191 142192 142193 142194 142195 142199 142203 142205 142210 142259 142309 142310 142383 142390 142391 142392 142394 142395 142397 142400 142404 142405 142406

data10_900GeV

AMI Tags

f174_m268 f175_m273 f176_m273 f176_m278 f177_m278 f178_m283 f179_m283 f170_m258 f185_m298 f186_m298 f187_m304 f188_m304 f189_m304 f180_m288 f181_m293 f181_m298 f183_m298 f196_m325 f190_m310 f190_m315 f190_m320 f191_m315 f191_m320 f192_m320 f193_m320 f193_m325 f215_m377 f215_m382 f215_m387 f216_m382 f216_m387 f217_m387 f217_m392 f218_m392 f218_m392 f212_m377 f234_m422 f236_m427 f238_m427 f239_m427 f231_m422 f232_m422

Streams

[physics_L1Calo](#) [express](#) [express](#) [physics_MinBias](#) [physics_BPTX](#) [physics_RNDM](#) [physics_L1CaloEM](#) [physics_MuonswBeam](#) [physics_CosmicCalo](#) [physics_CosmicMuons](#) [physics_CosmicCaloEM](#)

Individual Element Searches

Dataset data10_900 149751 physics_RN merge AOD f217_m387

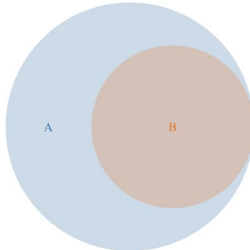
Event 510799

5

A = data09_900GeV.140571.physics_BPTX.merge.AOD.f175_m273
B = data09_900GeV.140571.physics_RNDM.merge.AOD.f175_m273

A 3253420028
B 2084349450
A∩B 4132819449
A∪B 1204950029 (29.16%)
A-A∩B 2048469999 (62.96%)
B-A∩B 879399421 (42.19%)

select



DOOverlap Remove Describe

Venn

Customize the interactions with the graph.

Cluster by group type Cluster by group size Expand all clusters

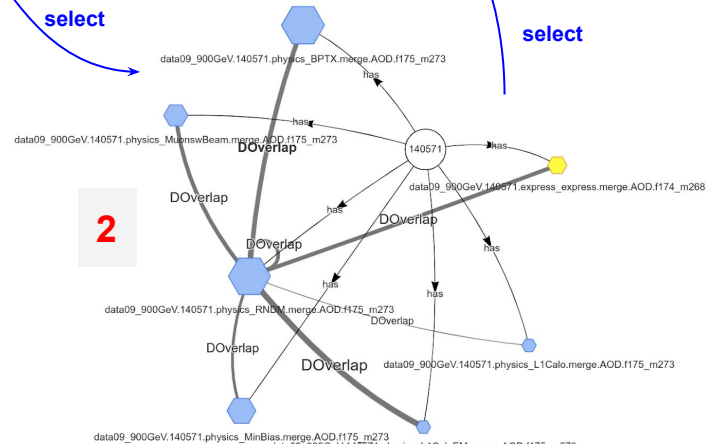
live remove old

filter:

operation feedback --- Loading Run 140571

select

2



select



Added Values & New Possibilities for Event Index

- Big part of the application code absorbed in Graph Database structure
 - Delegate implementation/optimisation details to suitable database framework
 - Database carries information about data structure which would otherwise have to be handled by the application code
- Simple graphical access
 - JavaScript client connects directly do Gremlin server
- Using standards
 - So components can be replaced
 - JanusGraph with Neo4J,...
 - Hadoop with Cassandra,...
- Possibility of Creation of Virtual Entities
 - Virtual Collections
 - 'Whiteboard' functionality
 - Can create API allowing users creating their own (or group-wise) Collections
 - Persistent Requests results
 - Results can be stored as new objects with relations (cache functionality)
 - Query Spaces
 - Abstract space on top of data to allow faster searching and navigation
- Finding Higher level correlations
 - Multiple trigger overlaps between events, trigger pattern in data derivations,...



Performance

- Requests in general in three phases
 - First search of the initial entry point (event, dataset, run,...)
 - Could be optimised
 - Natural order
 - Indexes
 - Elastic Search
 - Spark
 - More hierarchical navigation
 - Then navigation on the graph
 - Very fast
 - And finally accumulation of results
- Data can still be accessed directly, without Graph Database API
 - So with the same performance as non-GraphDB
 - Navigational step (instead of sub-search) can only speed it up
- In general:
 - Very fast retrieval
 - Slower import
 - Because import creates structures
 - Which are used in retrieval (simpler & faster)



Graph Databases for Functional Programming

*General
Comments*

- Relations (edges) can be considered as functions
 - Navigation as a function execution
 - **From the user point of view, there is no difference in creating new object or navigating to it**
 - Both operation can be 'lazy'
- Functional processing and graph navigation ("Graph Oriented Programming") can work very well together
 - Using the same functional syntax
 - Both are realisation of Categories
 - Vertex == object, Edge == morphism
 - Functional program can be modeled as a Graphs
 - Graph data can be navigated using functions
 - Data **ready for parallel access**
- Very well implemented by Gremlin

Extending parallel-ready functional model from code to data !



Graph Databases for Deep Learning

- Neural Network itself is a Graph
 - Using Graph Database to describe NN itself
- In many cases, Neural Network handles Graph data (objects with relations)
 - They can operate either on individual nodes (Node-focused tasks)
 - Or on the whole graph (Graph-focused tasks)
- GraphNN can be seen as a generalisation of ConvolutionalNN
 - Non-geometric
- Possibility to impose constraints/knowledge to NN
 - Inductive Bias
 - Semantic Induction

Graph Neural Networks create a Natural environment for Deep Learning !

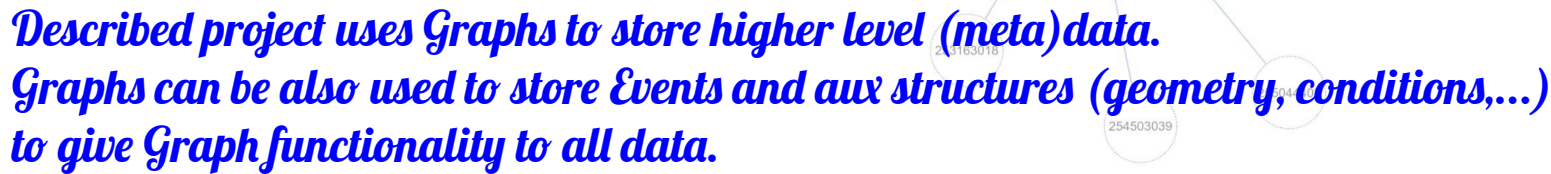


- Parallel programming
- Functional programming

- More structured data => simpler and faster access

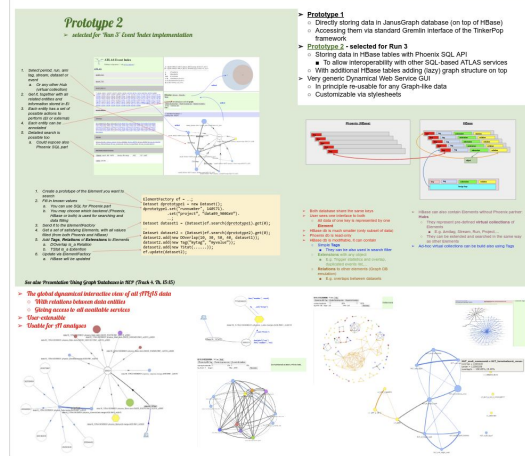
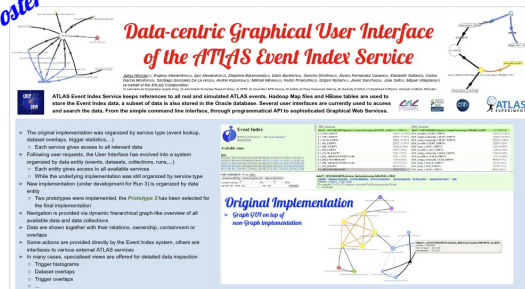
- More transparent code
 - Stable data structure is handled in the storage layer
- Suitable for **Functional Style** and **Parallelism**
- Suitable for **Deep Learning**
- Suitable for **Declarative Analyses**
- Can help with **Analysis Preservation**
- Language & Framework neutral

- Store data in a real Graph database
- Build a Graph layer on top of the existing storage
 - Close to DB layer
 - In the application layer





See also Poster 211



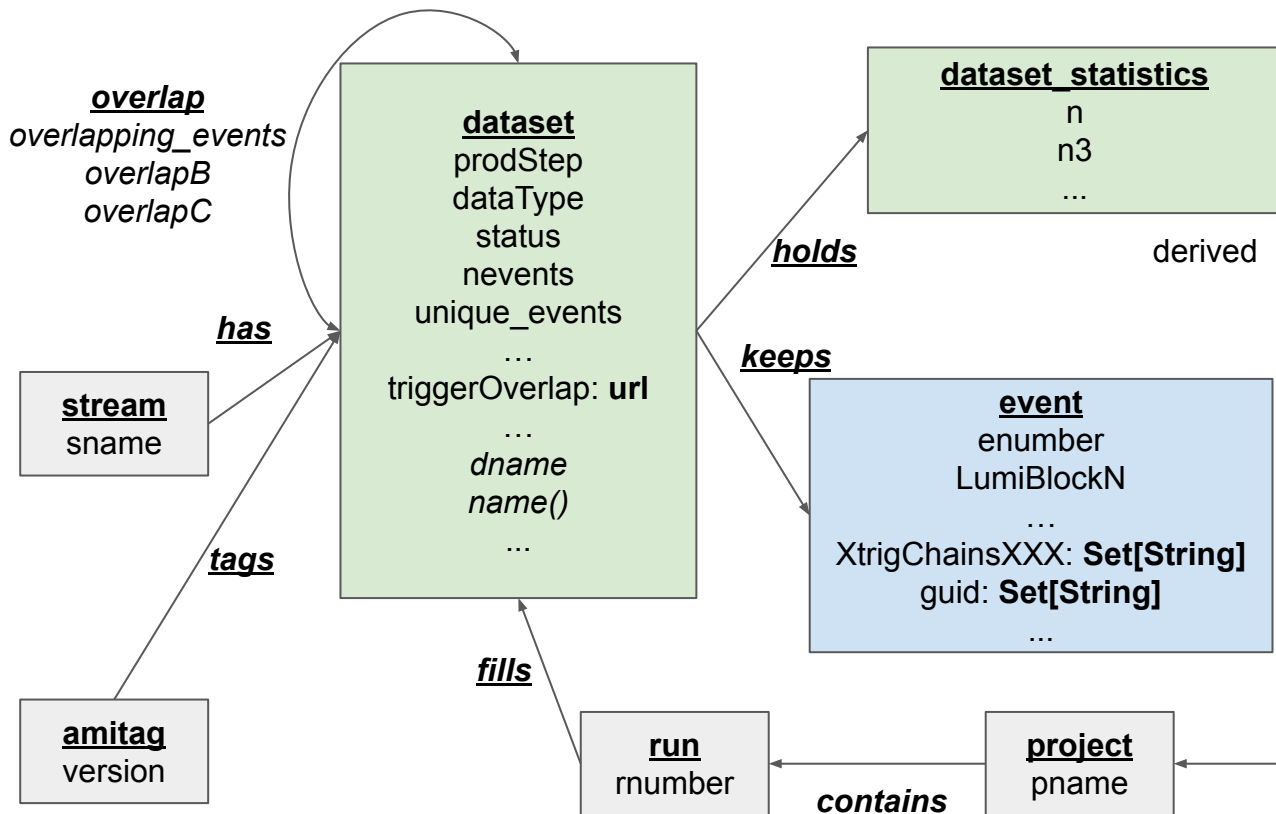


Schema may be specified to speed up searching and navigation.

Other nodes & vertexes can be added.

Schema

Prototype 1



- Should decide what is a feature and what is a relation
- Arrows have only logical meaning, they can be navigated equally from both sides
- Relations and features have defined multiplicities (not shown here) and types (int, string, date,..., Set[...],...)
- Defined entities (and combinations of them) can be indexed for fast search
- Other features and relations can be freely added to any entity (vertex or edge)

Object API

Prototype 2

- *REST and GUI WS are build on top of this API*
- *Subset of Gremlin API also available*

1. Create a prototype of the Element you want to search

2. Fill in known values

- You can use SQL for Phoenix part
- You may choose which backend (Phoenix, HBase or both) is used for searching and data filling

3. Send it to the ElementFactory

4. Get a set of satisfying Elements, with all values filled (from both Phoenix and HBase)

5. Add Tags, Relations of Extensions to Elements

- DOverlap is_a Relation
- TStat is_a Extension

6. Update via ElementFactory

- HBase will be updated

```
ElementFactory ef = ...;
Dataset dprototype1 = new Dataset();
dprototype1.set("runnumber", 140571).
    .set("project", "data09_900GeV").
    ...;
Dataset dataset1 = (Dataset)ef.search(dprototype1).get(0);
...
Dataset dataset2 = (Dataset)ef.search(dprototype2).get(0);
dataset2.add(new DOverlap(10, 30, 50, 40, dataset1));
dataset2.add(new Tag("mytag", "myvalue"));
dataset2.add(new TStat(...));
ef.update(dataset2);
```

GUI WS

Operations on View.

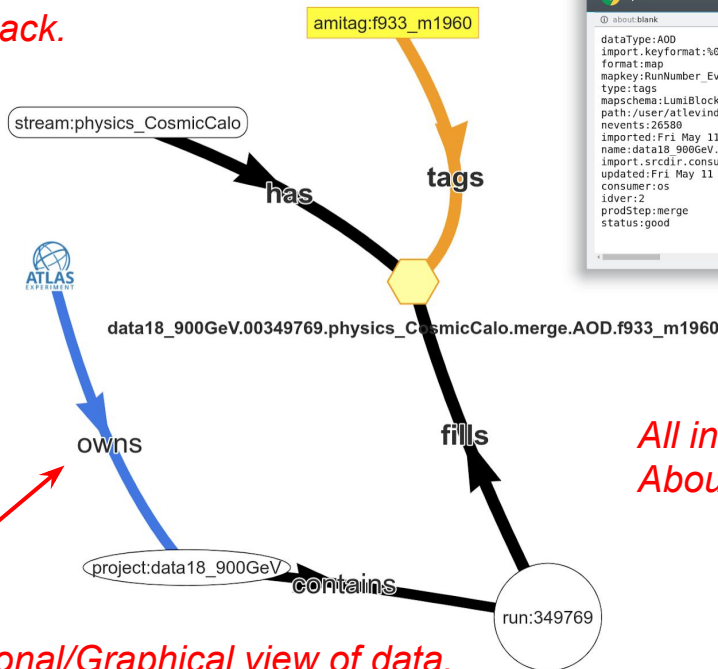
Context-sensitive action
on selected element(s).

Customize the interactions with the graph.
Cluster by group type | Cluster by group size | Expand all clusters
☐ live ☒ expand children ☒ expand parents ☐ remove old
Expanding data18_900GeV.00349769.physics_CosmicCalo.merge.AOD.f933_m1960 # Showing 0 new elements # Showing 3 new elements # Showing 3 new elements #

data18_900GeV.00349769.physics_CosmicCalo.merge.AOD.f933_m1960 Remove Describe
Catalog - Sample - Info -

--- commands output ---

Operation feedback.

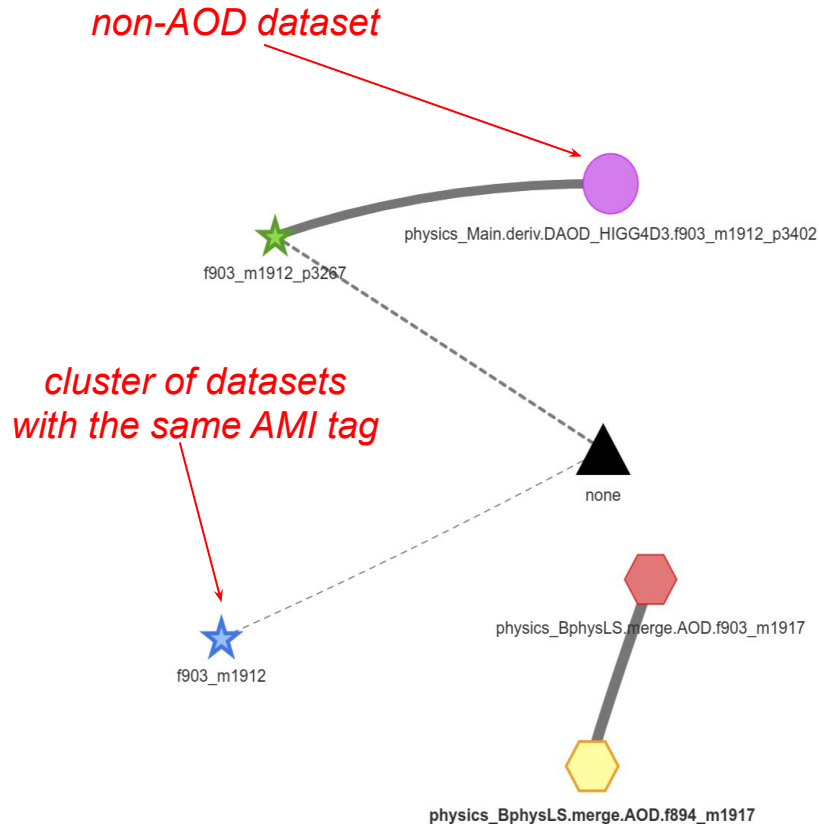


```
17637376 - Google Chrome
dataType:AOD
import.keyformat:%08d-%011d
format:map
mapkey:RunNumber_EventNumber=String
type:tags
mapschema:LumiBlock=Int BunchId=Int EventTime=Int EventTimeNanoSec=Int EventWeight=float
path:/user/atleind/EI18.1/data18_900GeV.00349769.physics_CosmicCalo.merge.AOD.f933_m1960
nevents:26580
imported:Fri May 11 23:00:00 CEST 2018
name:data18_900GeV.00349769.physics_CosmicCalo.merge.AOD.f933_m1960
import.srcdir.consumer:/user/atleind/ObjectStoreConsumerData/datasets/data18_900GeV.003
updated:Fri May 11 23:01:38 CEST 2018
consumer:os
idver:2
prodStep:merge
status:good
```

Context-sensitive action
command output.

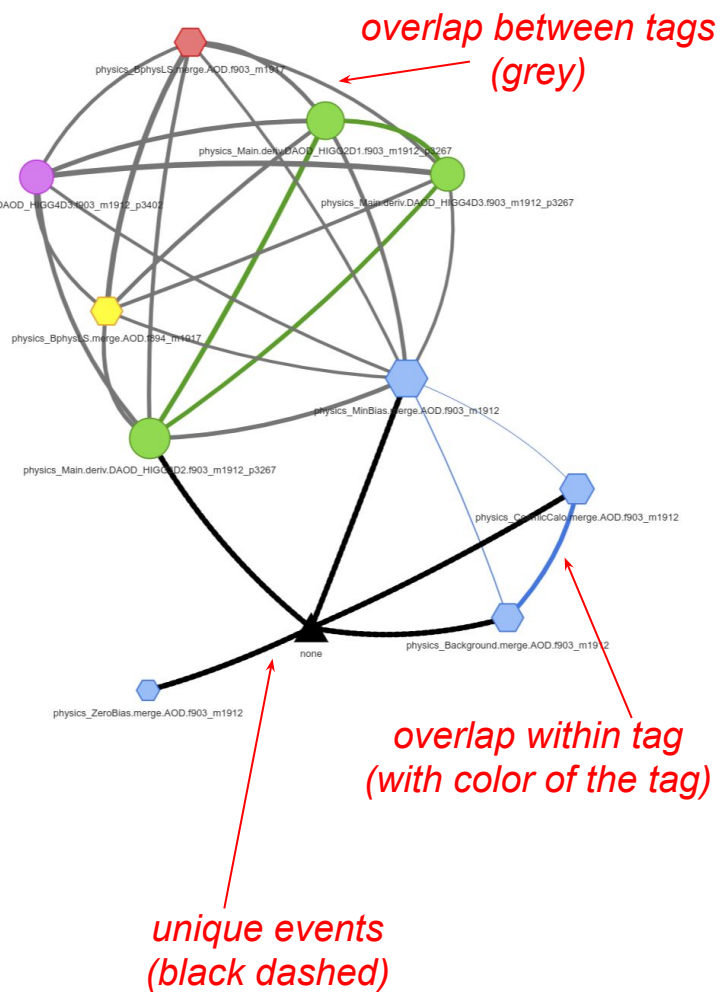
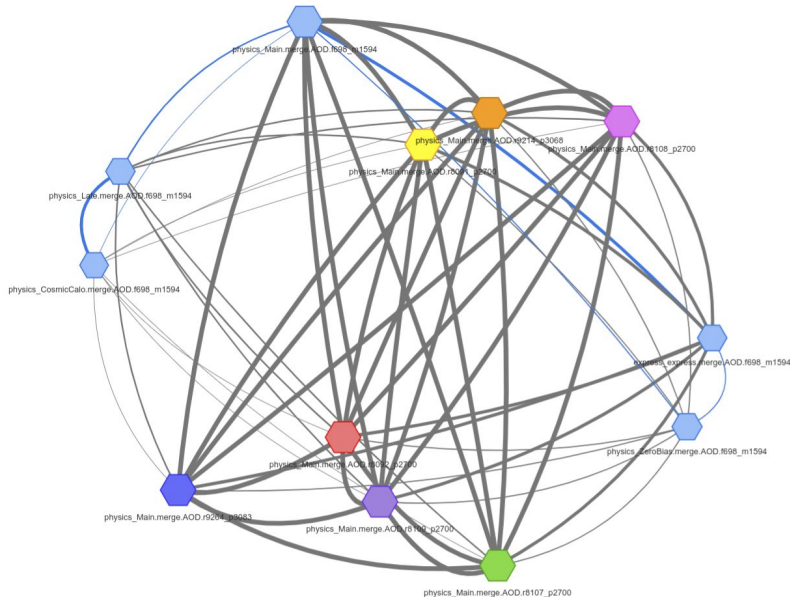
All information
About selected element.

InterActive Relational/Graphical view of data.



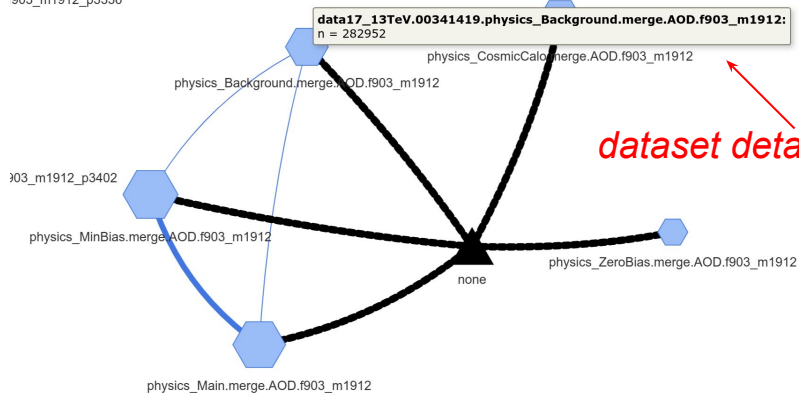


Context-sensitive menu will be here.



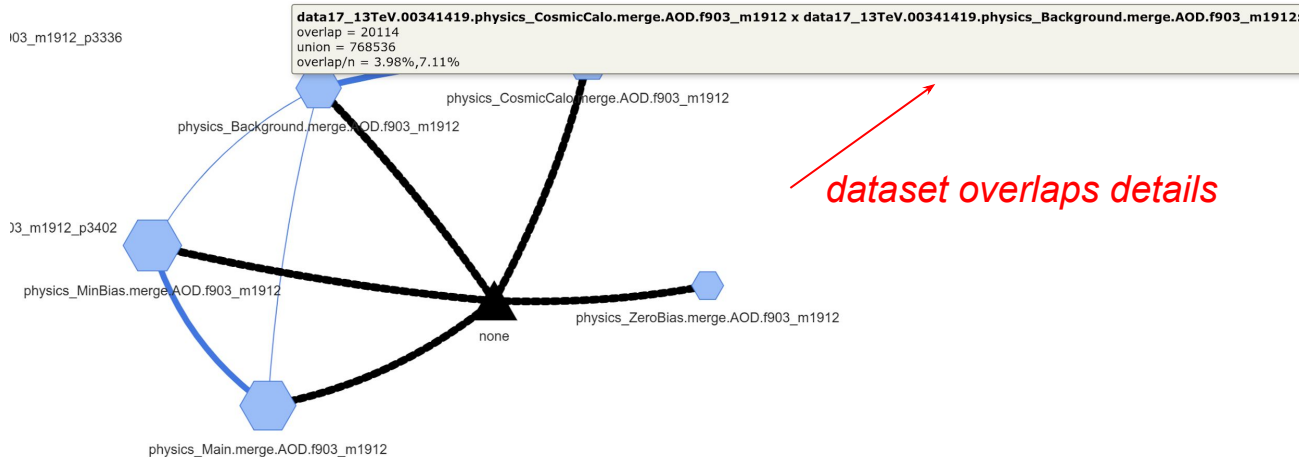
Details

03_m1912_p3336



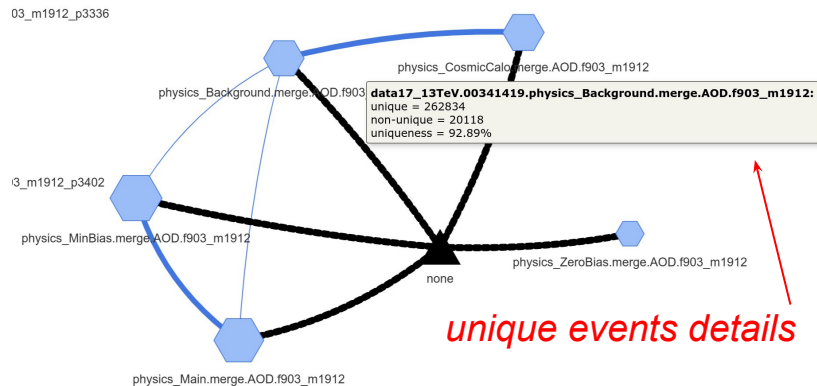
dataset details

03_m1912_p3336



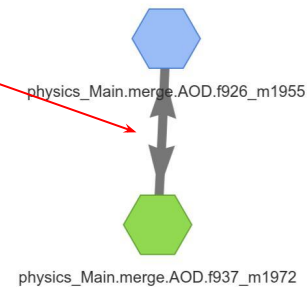
dataset overlaps details

03_m1912_p3336

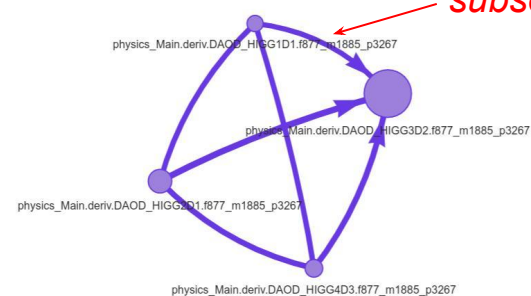


unique events details

identity
(mutual subset)



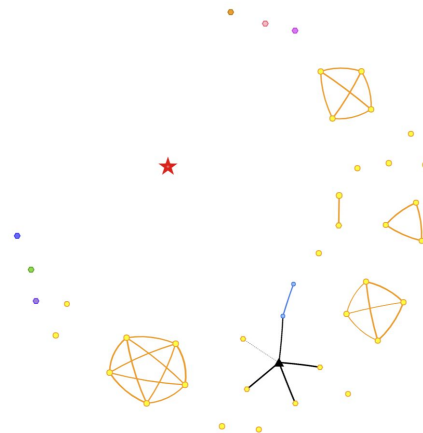
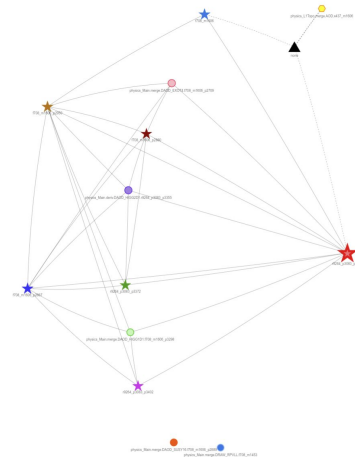
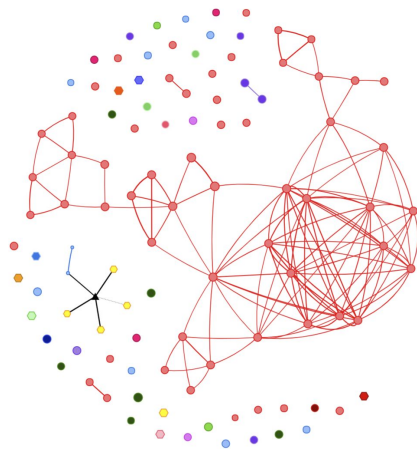
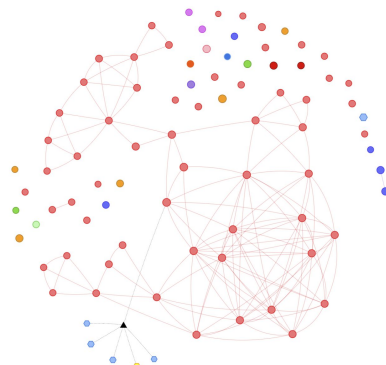
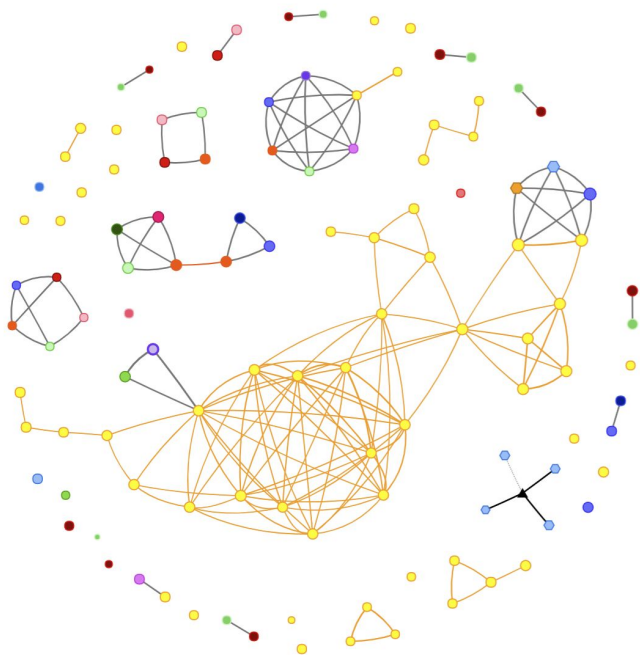
subset



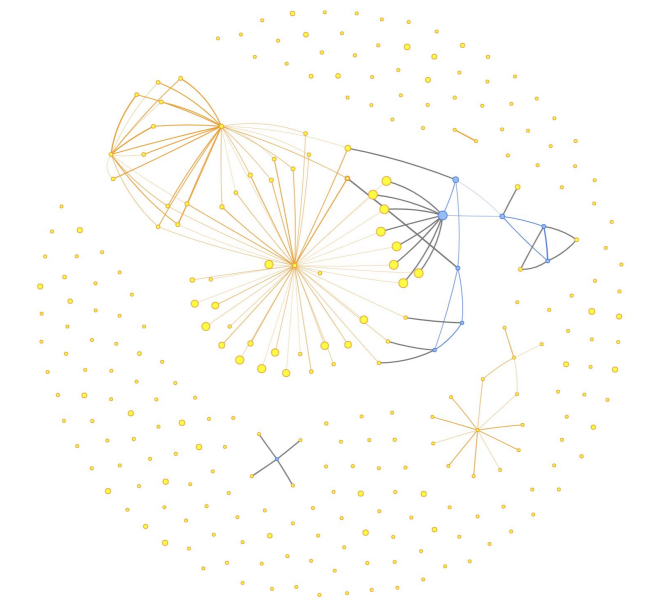
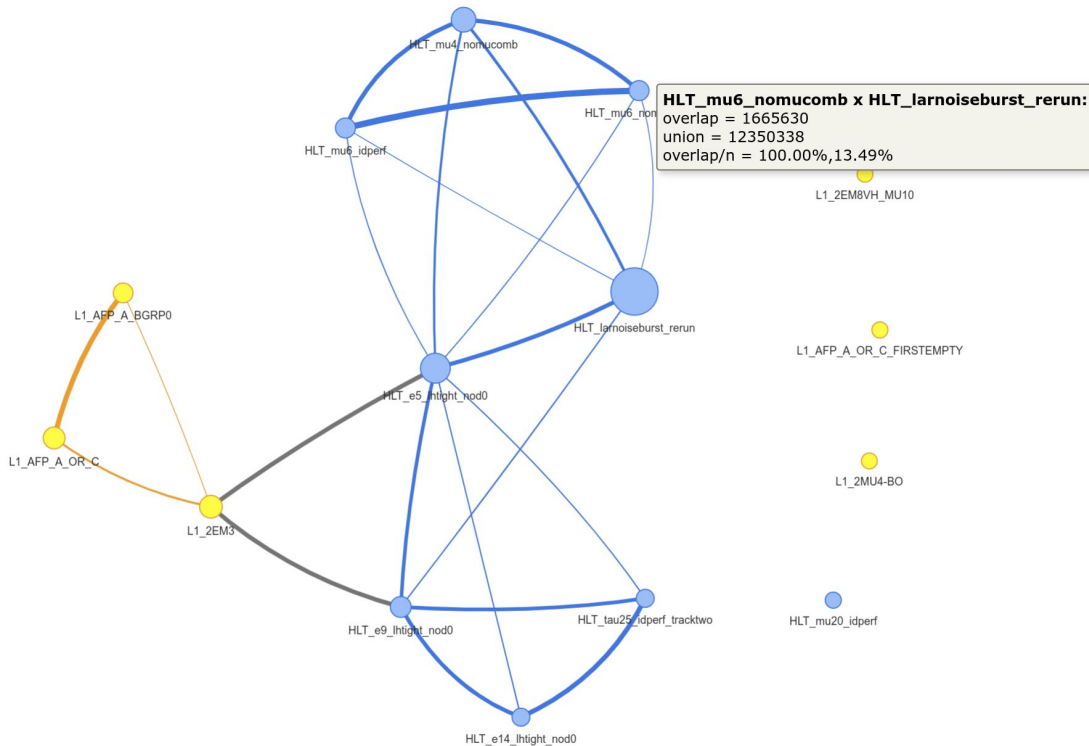
Clusters

El16.1/00298690: ☒ live [Help](#)
Cluster by AMI Tag
overlap thresholds:
tag level: 99 target: null filter: null

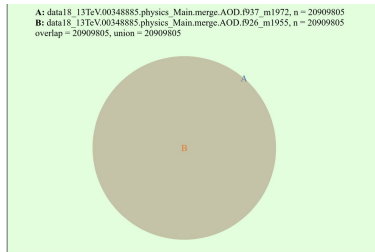
Context-sensi



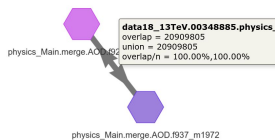
Trigger Overlaps



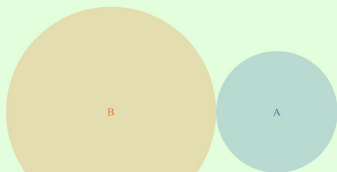
Overlaps Venn Diagrams



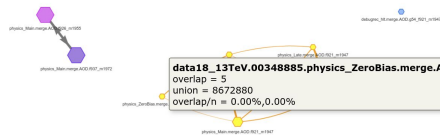
[data18_13TeV.00348885.physics_Main.merge.AOD.f937_m1972 x data18_13TeV.00348885.physics_Main.merge.AOD.f926_m1955 Venn Diagram](#)



A: data18_13TeV.00348885.physics_ZeroBias.merge.AOD.f921_m1947, n = 186395
B: data18_13TeV.00348885.express_express.merge.AOD.f921_m1947, n = 556987
overlap = 1, union = 743381

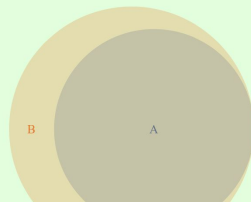


[data18_13TeV.00348885.physics_ZeroBias.merge.AOD.f921_m1947 x data18_13TeV.00348885.express_express.merge.AOD.f921_m1947 Venn Diagram](#)



data18_13TeV.00348885.physics_ZeroBias.merge.AOD.f921_m1947 x data18_13TeV.00348885.express_express.merge.AOD.f921_m1947 Venn Diagram

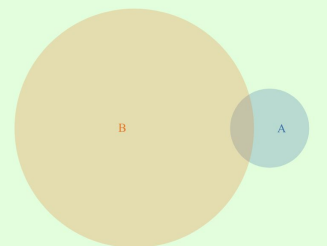
A: HLT_10j40_L16J15, n = 21866
B: HLT_10j40_L14J150ETA25, n = 32799
overlap = 21866, union = 32799



[HLT_10j40_L16J15 x HLT_10j40_L14J150ETA25 Venn Diagram](#)



A: data18_13TeV.00348885.physics_Main.deriv.DAOD_HIGG1D1.f926_m1955_p3544, n = 267984
B: data18_13TeV.00348885.physics_Main.deriv.DAOD_EXOT4.f926_m1955_p3544, n = 2452676
overlap = 61468, union = 2659192



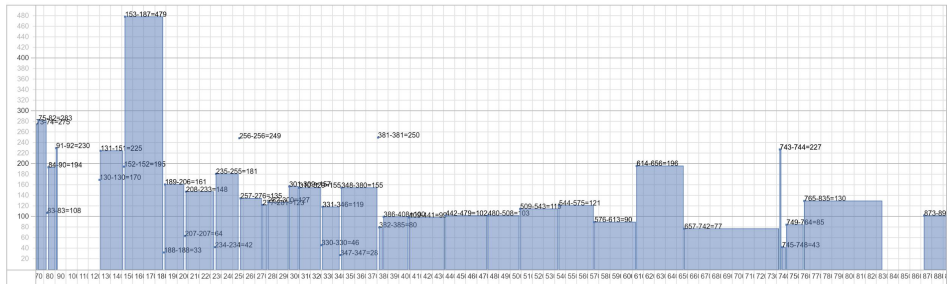
[data18_13TeV.00348885.physics_Main.deriv.DAOD_HIGG1D1.f926_m1955_p3544 x data18_13TeV.00348885.physics_Main.deriv.DAOD_EXOT4.f926_m1955_p3544 Venn Diagram](#)

data18_13TeV.00348885.physics_Main.deriv.DAOD_HIGG1D1.f926_m1955_p3544 x data18_13TeV.00348885.physics_Main.deriv.DAOD_EXOT4.f926_m1955_p3544 Venn Diagram

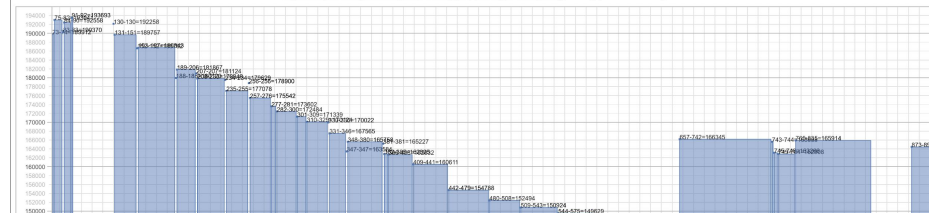


Trigger Overlaps for LB Ranges

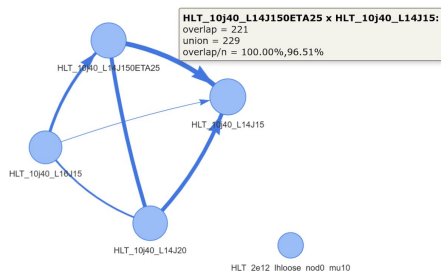
HLT_1040_L16/15 x HLT_1040_L14/150ETA25 overlap for data19_13TeV.0031364.physics_Main.merge.AOD.PS7_m1972 (per LB ranges, normalised to reverts)



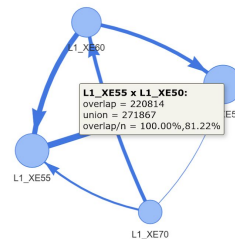
L1_XE55 x L1_XE50 overlap for data19_13TeV.0031364.physics_Main.merge.AOD.PS7_m1972 (per LB ranges, normalised to reverts)



HLT_1040_L16/15 x HLT_1040_L14/150ETA25
Venn Diagram - LB Graph



L1_XE55 x L1_XE50
Venn Diagram - LB Graph(*)



Performance Example

Event Lookup

```
gremlin> el(358031, 775206623, g).profile()
```

```
==>Traversal Metrics
```

Step	Count	Traversers	Time (ms)	% Dur
JanusGraphStep([],[~label1.eq(event), number.eq... _condition=(~label = event AND number = 775206623) _isFitted=true _query=multiKSQ[1]@2147483647 _index=event:number:u _orders=[] _isOrdered=true optimization optimization backend-query _query=event:number:u:multiKSQ[1]@2147483647	1	1	204.805	75.74
JanusGraphVertexStep(IN,[keeps],vertex) _condition=type[keeps] _isFitted=true _vertices=1 _query=org.janusgraph.diskstorage.keycolumnvalue.SliceQuery@b3a55b7f _orders=[] _isOrdered=true optimization backend-query _query=org.janusgraph.diskstorage.keycolumnvalue.SliceQuery@b3a55b7f	1	1	25.560	9.45
JanusGraphVertexStep(IN,[fills],vertex) _condition=type[fills] _isFitted=true _vertices=1 _query=org.janusgraph.diskstorage.keycolumnvalue.SliceQuery@b3a605c1 _orders=[] _isOrdered=true optimization backend-query _query=org.janusgraph.diskstorage.keycolumnvalue.SliceQuery@b3a605c1	1	1	10.388	3.84
HasStep([rnumber.eq(358031)])	1	1	13.129	4.86
SelectOneStep(last,e)	1	1	0.993	0.37
NoOpBarrierStep(2500)	1	1	0.159	0.06
JanusGraphPropertiesStep([guid],value) _condition=type[guid] _isFitted=true _vertices=1 _query=org.janusgraph.diskstorage.keycolumnvalue.SliceQuery@b11f98a7 _orders=[] _isOrdered=true optimization	2	2	14.800	5.47
NoOpBarrierStep(2500)	2	2	0.568	0.21
>TOTAL	-	-	270.406	-

*75% of the time is spend by the entry point search,
following graph traversal is very fast.*

*This was the first request,
the second one will be cca 10x faster
(even on different event).*