

Easy access to data with Presto

Iker Martinez de Apellaniz (@mitxino77)
Albert Franzl (@FranziCros)





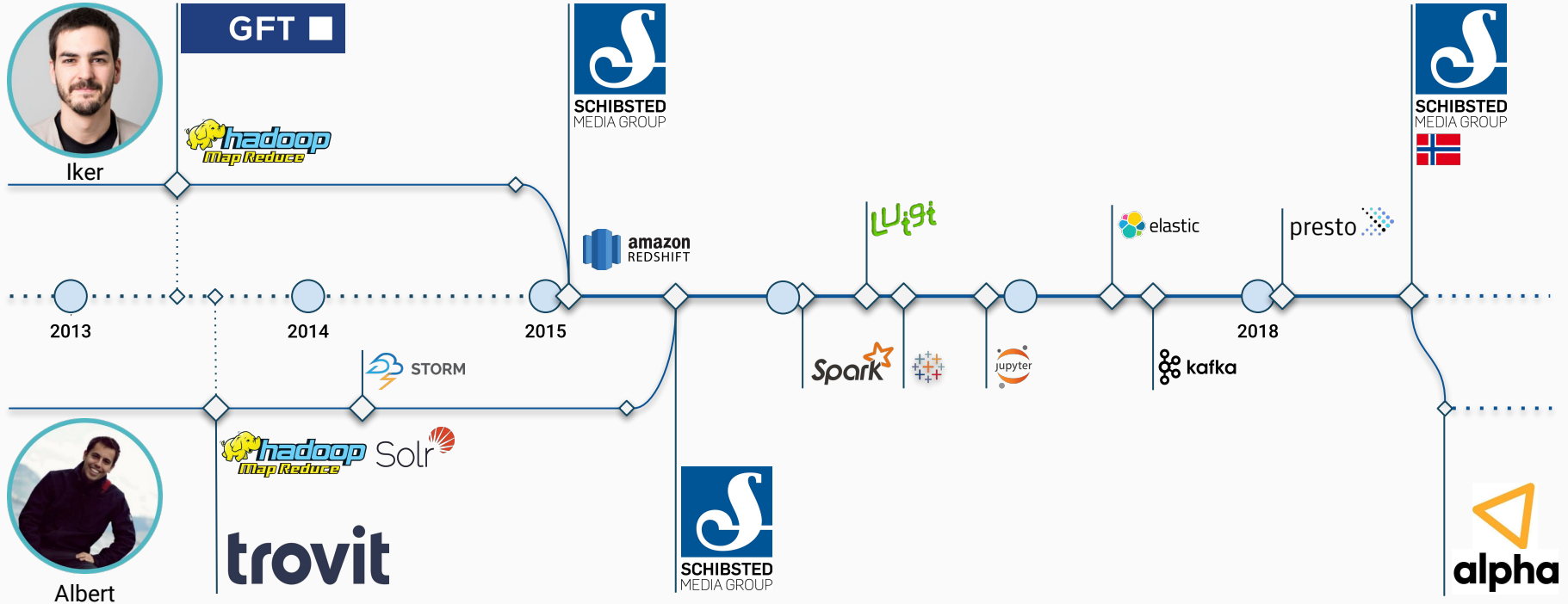
Slides available in:



<< bitly link here >>

About us

euagister:



Our Journey

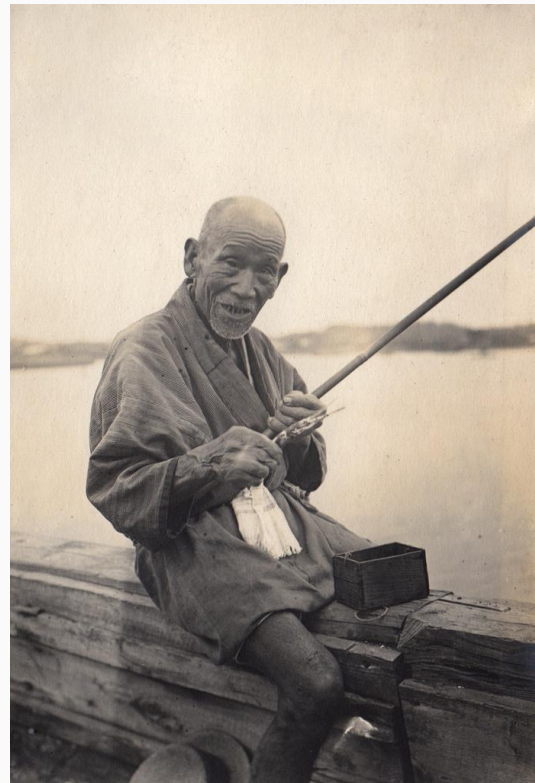
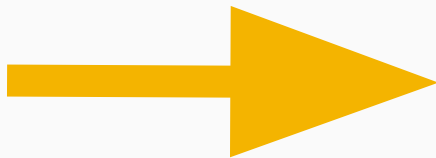
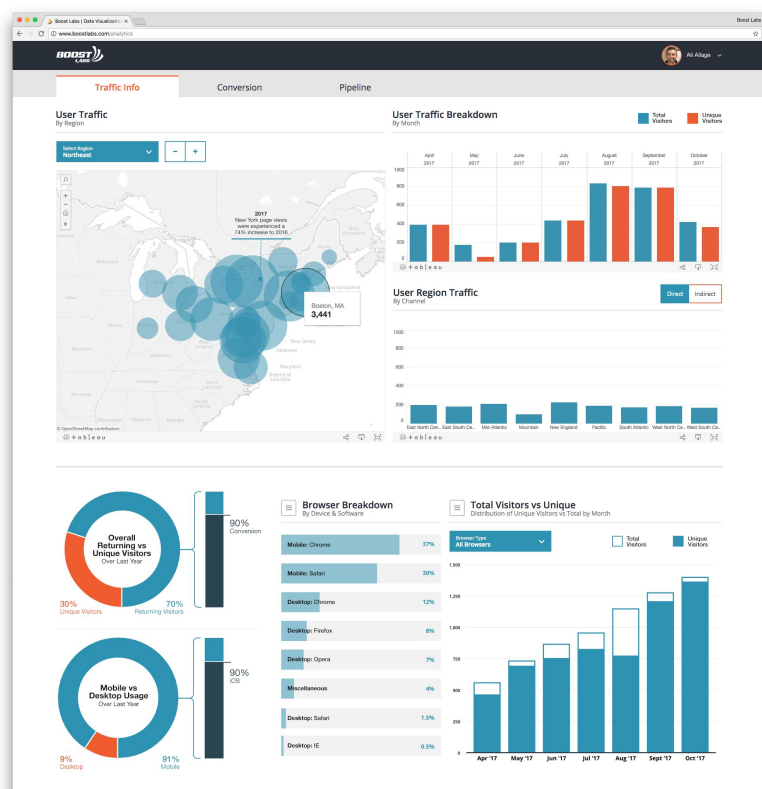






Tableau-aas



Jupyter-aaS



Jupyter as a Service

```
In [1]: import pywikibot

In [6]: site = pywikibot.Site('wikidata', 'wikidata')

In [9]: site
Out[9]: DataSite('wikidata', 'wikidata')

In [13]: item = pywikibot.ItemPage(site, 'Q4115189')

In [14]: item
Out[14]: ItemPage('Q4115189')

In [15]: item_dict = item.get()

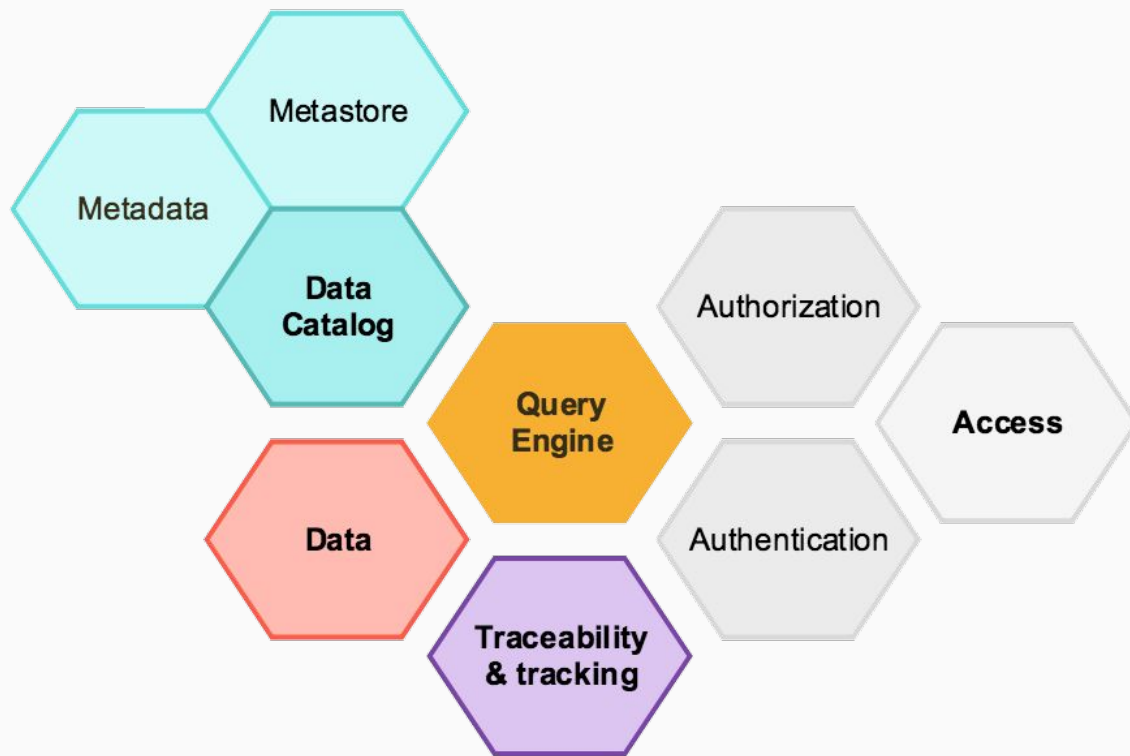
In [16]: item_dict
Out[16]: {'aliases': {'ar': ['مادحة الحريم', 'مادحة اللهب'], 'de': ['Spielewiese', 'Sandbox'], 'de-at': ['Sandkasten', 'Spielplatz'], 'en': ['SD', 'Property test', 'test', 'wikidata SandboxItem', 'wikidata box'], 'it': ['Sandbox di Wikidata'], 'ja': ['Sandbox', 'サンドボックス', '練習用ページ', '練習用項目'], 'nl': ['wikidata-speelruimte'], 'pt-br': ['item para testes', 'teste', 'testes', 'test', 'página de testes'], 'ru': ['rect', 'rect2'], 'zh-hans': ['维基数据沙盘', '维基数据测试']}, 'claims': {'P1110': [{'pywikibot.page.Claim at 0x77f43b30390>'}, 'P1132': [{'pywikibot.page.Claim at 0x77f43b3665b>'}, 'P1302': [{'pywikibot.page.Claim at 0x77f43b3619b>'}, 'pywikibot.page.Claim at 0x77f43b3625b>'}, 'pywikibot.page.Claim at 0x77f43b3665b>'}, 'P1346': [{'pywikibot.page.Claim at 0x77f43b3936b>'}, 'P1350': [{'pywikibot.page.Claim at 0x77f43b394ab>'}, 'P1351': [{'pywikibot.page.Claim at 0x77f43b39a5b>'}, 'P1355': [{'pywikibot.page.Claim at 0x77f43b396f0>'}, 'P1356': [{'pywikibot.page.Claim at 0x77f43b3989b>'}, 'P18': [{'pywikibot.page.Claim at 0x77f43b39e1b>'}, 'P1923': [{'pywikibot.page.Claim at 0x77f43b3965b>'}, 'P2047': [{'pywikibot.page.Claim at 0x77f43b3977b>'}, 'P2630': [{'pywikibot.page.Claim at 0x77f43b3965b>'}, 'P27': [{'pywikibot.page.Claim at 0x77f43b3965b>'}, 'P279': [{'pywikibot.page.Claim at 0x77f43b3965b>'}, 'P31': [{'pywikibot.page.Claim at 0x77f43b39627b>'}, 'pywikibot.page.Claim at 0x77f43b39639b>'}, 'P426': [{'pywikibot.page.Claim at 0x77f43b396f5b>'}, 'P488': [{'pywikibot.page.Claim at 0x77f43b396f5b>'}]}
```



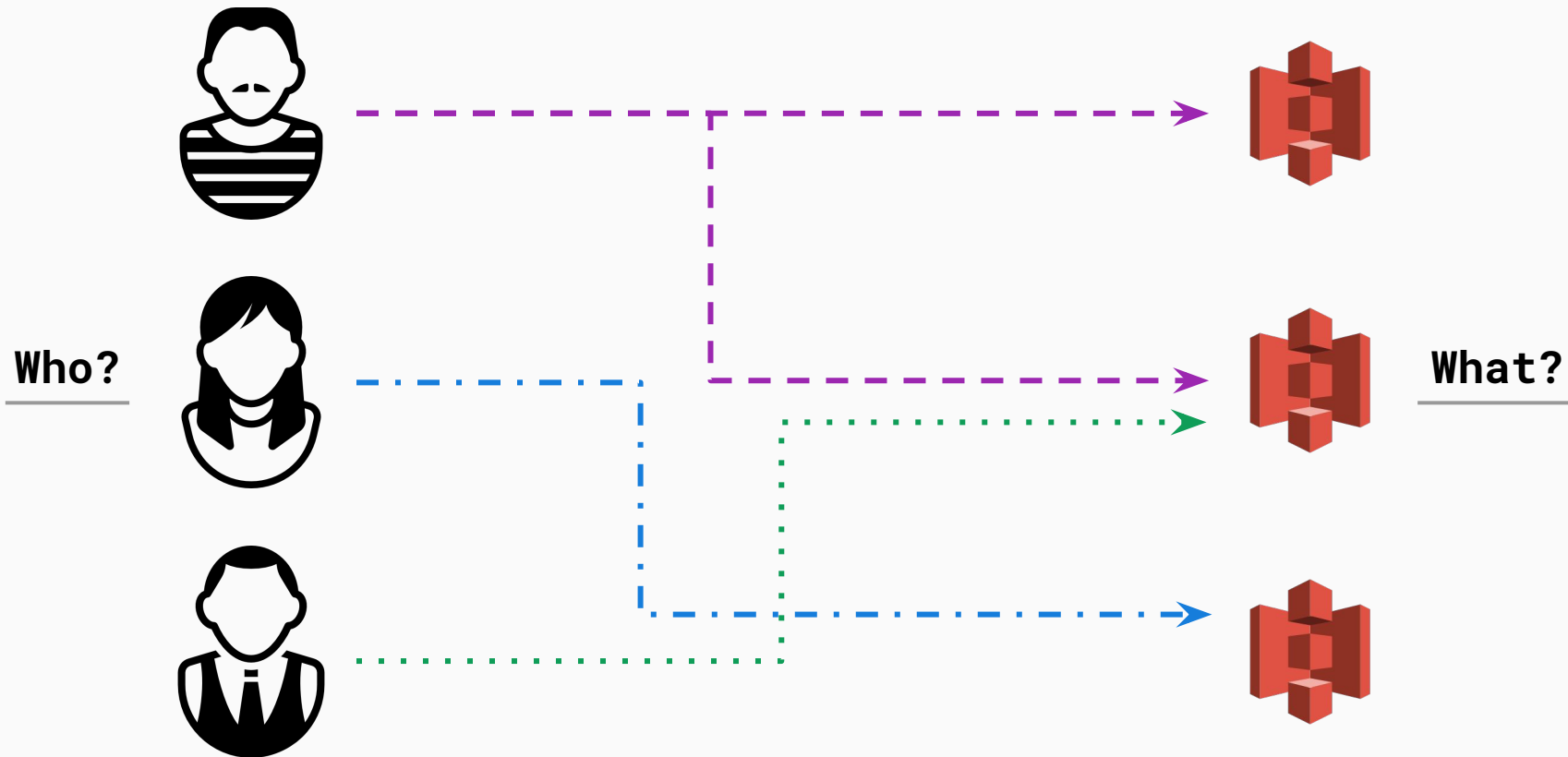

```
SELECT the_data_I_need  
FROM the_table_Im_looking_for  
WHERE my_analysis = 'Makes Sense';
```



Providing Presto as a Service



Auth* (authentication & authorization)



Auth* (authentication & authorization)



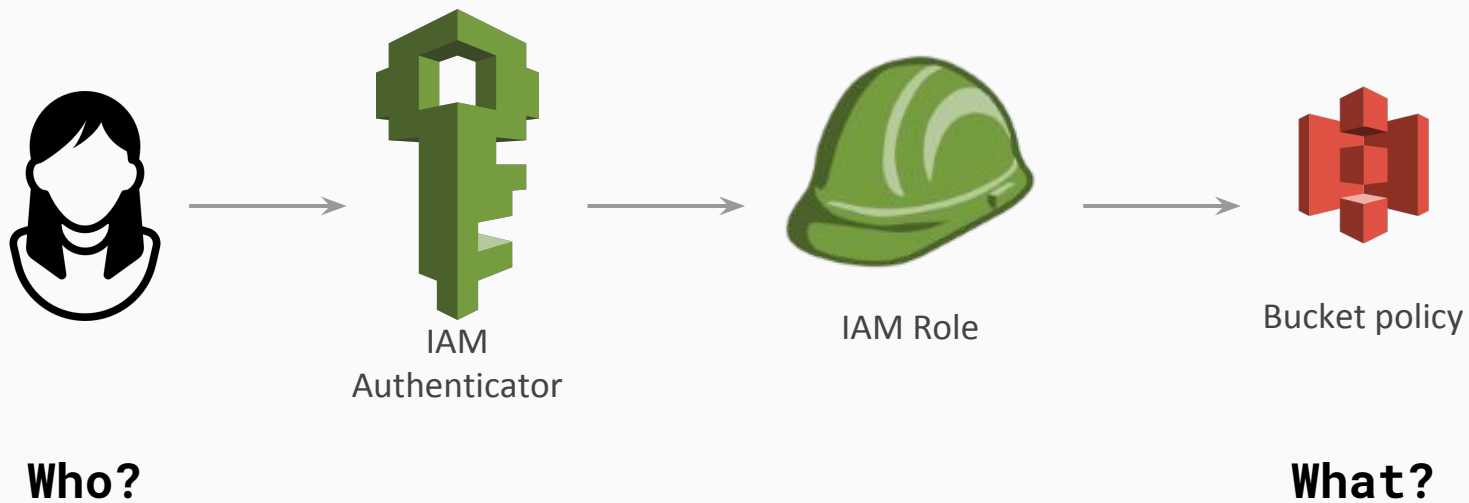
Auth* (authentication & authorization)

```
val authenticationCache: LoadingCache[Credentials, Principal] = CacheBuilder
    .newBuilder()
    .expireAfterWrite(cacheTtl.toMillis, MILLISECONDS)
    .build(new CacheLoader[Credentials, Principal] {
        override def load(key: Credentials): Principal = authenticate(key)
    })
```

Ref: github.com/google/guava/wiki/CachesExplained

Avoid multiple IAM calls by implementing a cache.

Auth* (authentication & authorization)



Traceability



IamAuthenticator

EventListener (Logs): Report to Datadog all created and completed queries.

SystemAccessControl: Report to Datadog all user interactions with Presto.

Traceability - SPI plugins < SystemAccessControl >

```
import com.facebook.presto.spi.security.AccessDeniedException._
import com.facebook.presto.spi.security.{Identity, SystemAccessControl}
import com.facebook.presto.spi.CatalogSchemaTableName

trait TableAccessControl extends SystemAccessControl {
  self: Logging with DatadogReporter =>

  override def checkCanCreateTable(identity: Identity, table: CatalogSchemaTableName): Unit = {
    reportTableEvent("can_create_table", identity, table)
  }

  override def checkCanDropTable(identity: Identity, table: CatalogSchemaTableName): Unit = {
    reportTableEvent("can_drop_table", identity, table)
  }

  override def checkCanRenameTable(identity: Identity, table: CatalogSchemaTableName,
    newTable: CatalogSchemaTableName): Unit = {
    reportTableEvent("can_rename_table", identity, table, allowed = false)
    denyRenameTable(table.toString, newTable.toString)
  }

  ...
}
```

Traceability - SPI plugins < SystemAccessControl >

```
class SaasSystemAccessControl(val statsd: DatadogMetricsClient) extends SystemAccessControl
  with UserAccessControl
  with CatalogAccessControl
  with SchemaAccessControl
  with TableAccessControl
  with ViewAccessControl
  with DatadogReporter
  with Logging
```

Traceability - SPI plugins < Event Listener >

```
import com.facebook.presto.spi.eventlistener._

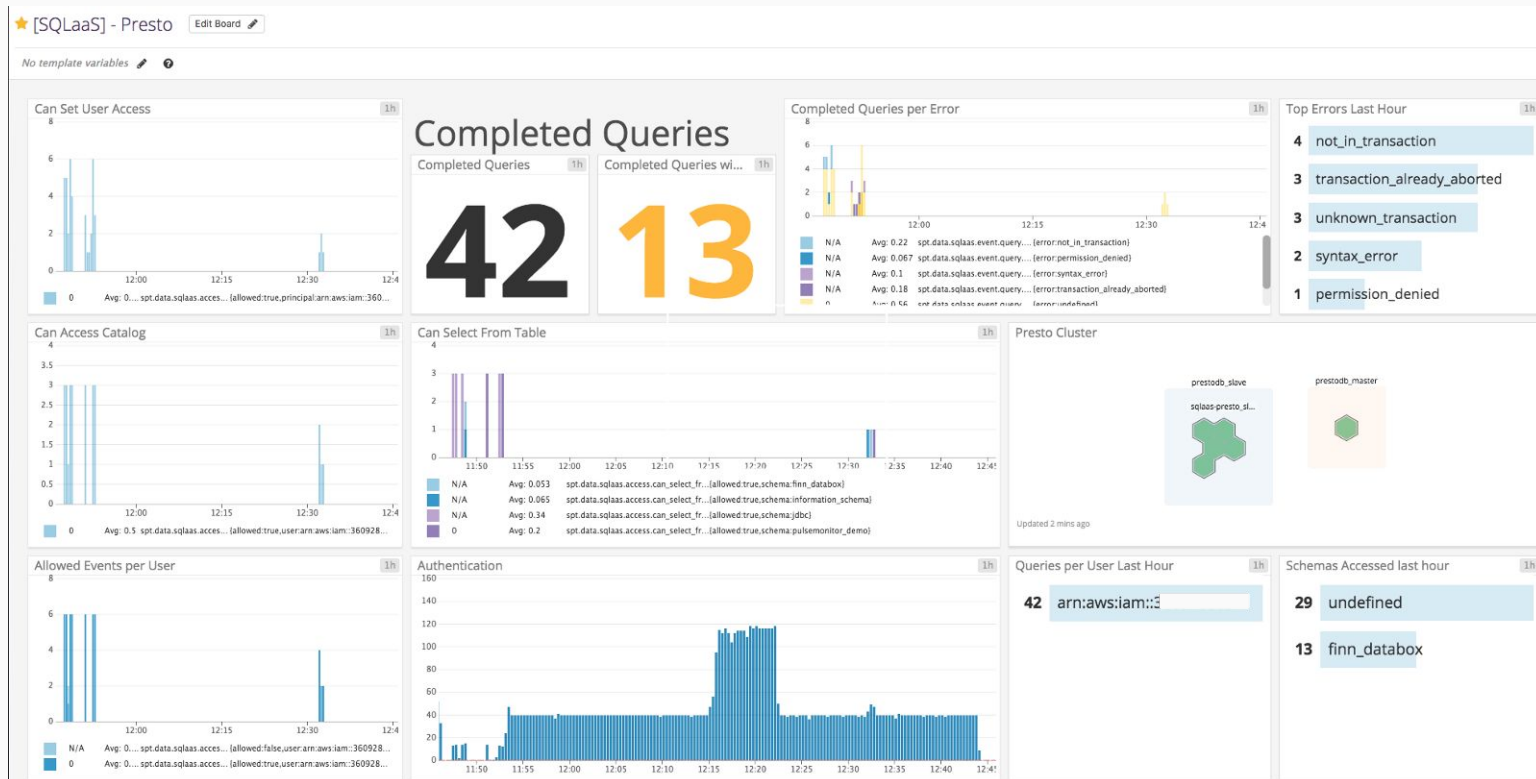
class SaasEventListener(statsd: DatadogMetricsClient) extends EventListener with Logging {

  override def queryCreated(queryCreatedEvent: QueryCreatedEvent): Unit = {
    val context = queryCreatedEvent.getContext

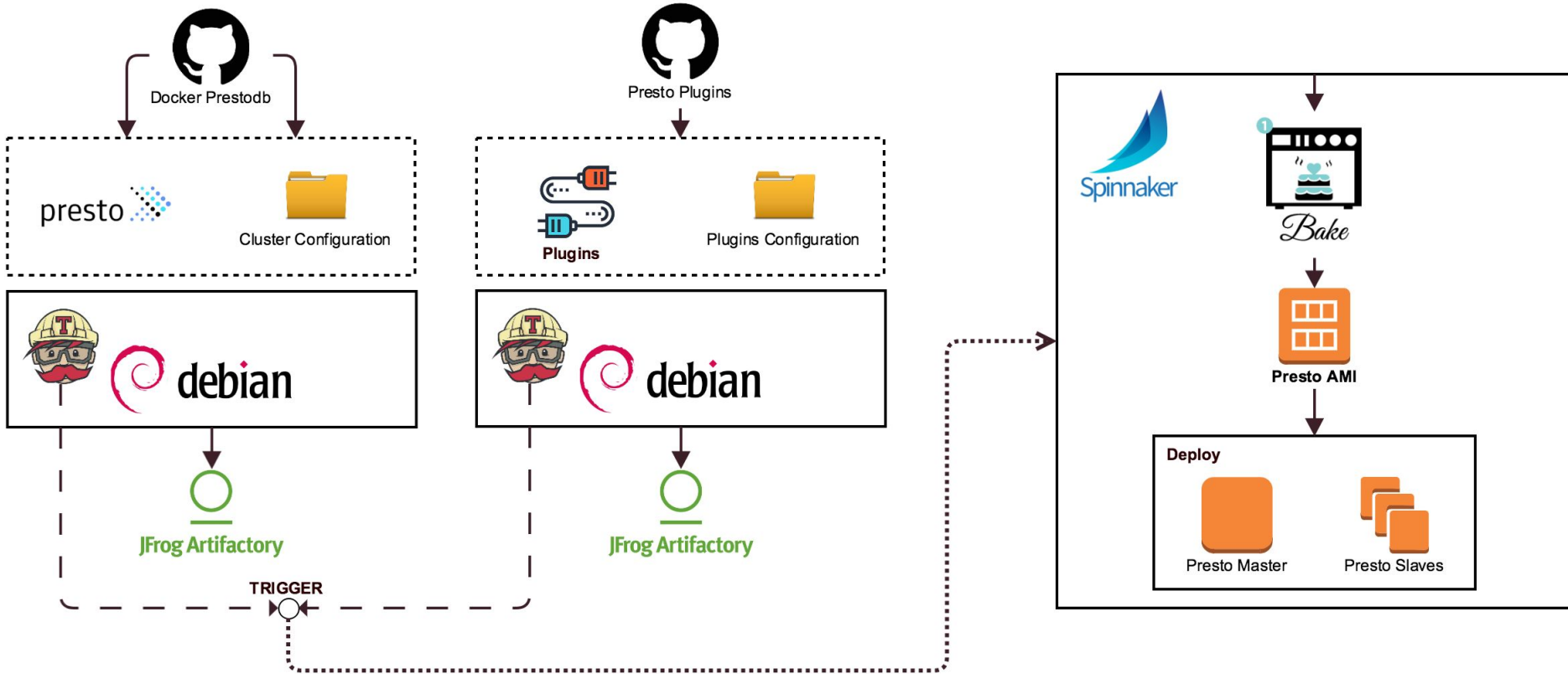
    val tags = Map(
      User -> context.getUser,
      PrincipalName -> context.getPrincipal.getOrElse(Undefined),
      Source -> context.getSource.getOrElse(Undefined),
      Catalog -> context.getCatalog.getOrElse(Undefined),
      Schema -> context.getSchema.getOrElse(Undefined),
      Resource -> context.getResourceGroupName.getOrElse(Undefined)
    )

    statsd.increment("query.created", tags)
  }
  ...
}
```

Traceability - Datadog Dashboards



CI / CD



CI / CD - Env Properties

```
# etc.template/config.properties.template.https
{% - if HTTPS_ENABLED is defined and HTTPS_ENABLED %}
http-server.https.enabled=true
http-server.https.port={{HTTPS_PORT}}
http-server.https.keystore.path=/opt/keystore/keystore.jks
http-server.https.keystore.key=presto
{% endif %}

{% - if AUTHENTICATION_TYPE is defined %}
http-server.authentication.type={{AUTHENTICATION_TYPE}}
{% endif %}
```

```
# etc.template/config.properties.template
presto.version={{PRESTO_VERSION}}.schibsted
{% include "config.properties.template.coordinator" %}
{% include "config.properties.template.http" %}
{% include "config.properties.template.https" ignore missing %}
...
```



Env variables defined in the
Spinnaker Pipeline

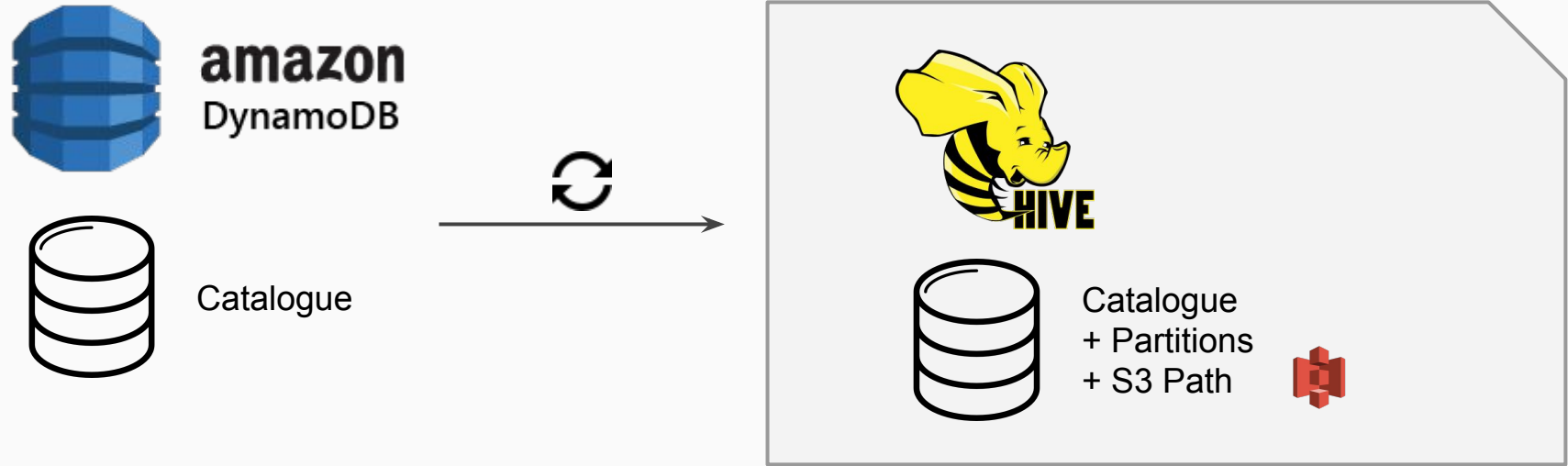
Metastore - Aws Glue

```
# etc/catalog/hive.properties
hive.metastore=glue
hive.metastore.glue.region=eu-west-1
hive.metastore.glue.assume-role=arn:aws:iam::0000000000:role/presto
```



PR: github.com/prestodb/presto/pull/10864

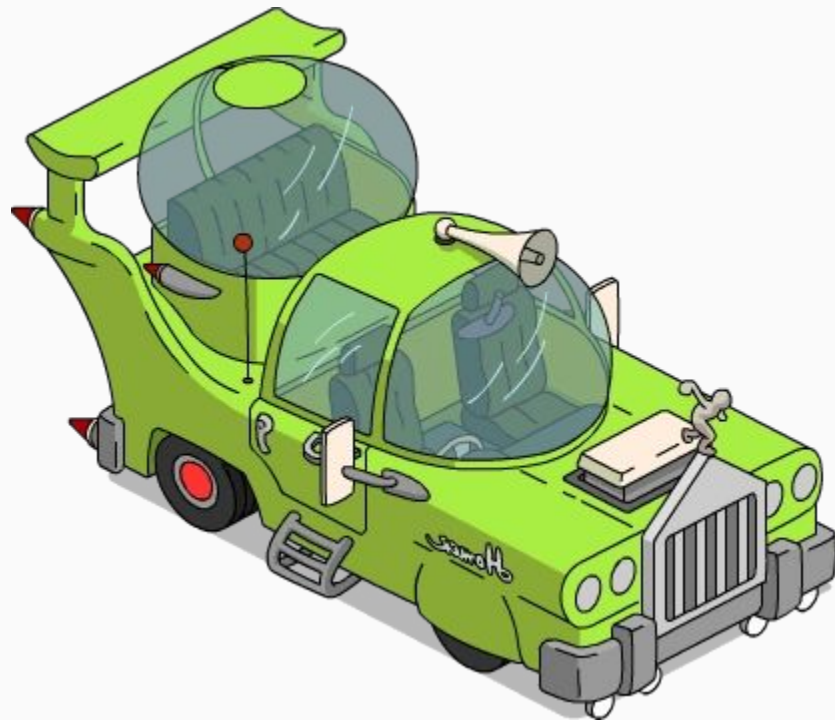
Metastore - Aws Glue - Catalogue Synchronizer



Metastore - Aws Glue - Catalogue Synchronizer

```
class Synchronizer(dynamoSynchronizer: DynamoSynchronizer, glueSynchronizer: GlueSynchronizer) {  
  
    def execute(): Unit = {  
        //- 1. Get Status from DynamoDB  
        val metaCatalogue: MetaCatalogue = dynamoSynchronizer.getDynamoDbMetastore  
  
        //- 2. Get Status from Glue (Schemas, Tables & Partitions)  
        val glueCatalogue: GlueCatalogue = glueSynchronizer.getGlueMetastore(metaCatalogue.schemas.map(_.id))  
  
        //- 3. Update Glue Metastore  
        //- 3.1. Create Missing Schemas  
        glueSynchronizer.createMissingSchemas(metaCatalogue, glueCatalogue)  
        //- 3.2. Create Missing Tables  
        glueSynchronizer.createMissingTables(metaCatalogue, glueCatalogue)  
        //- 3.3. Update outdated Tables  
        glueSynchronizer.syncExistingTables(metaCatalogue, glueCatalogue)  
        //- 3.4. Update table partitions  
        glueSynchronizer.syncTablePartitions(metaCatalogue)  
    }  
}
```

Data



Data



Data: Schema Hell

```
"object": {  
  "id": 123,  
  "location": "barcelona"  
},  
"sch:custom": "{ `this`: `json` }",  
"@type": "event",  
"ref-v2": "something"
```

```
"object": {  
  "ID": "123",  
    
},  
"sch:custom": { "this": "json" },  
"@type": "event"
```

```
"object": {  
  "id": "123",  
  "location": {"lat", "long"}  
},  
"sch:custom": { "this": "json" },  
"@type": "event",  
"ref-V2": "something"
```

```
"schema-hell": {  
  "id": "123",  
  "location": "barcelona"  
},
```

Data: Hard Rules / Soft Rules

```
"event_id": "ABC-1234",  
"provider": "sdrn:fotocasa",  
"timestamp": "2018-09-14T01:02Z",  
"type": "view",  
"object": {  
  "id": 123,  
  "type": "classified_ad",  
  "category": "loft"  
},  
"tracker": {  
  "type": "JS",  
  "version": "1.2.3.4-alpha",  
},  
"location": {  
  "lat": 10.0,  
  "long": 10.0  
},  
"ref_v2": "something"
```

Ad Views x Category

NumOf Events x version

```
"provider": "sdrn:fotocasa",  
"timestamp": "2018-09-14T01:02Z",  
"type": "view",  
"object": {  
  "id": 123,  
  "type": "classified_ad",  
  "category": "loft"  
}
```

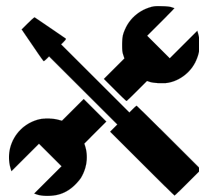
```
"event_id": "ABC-1234",  
"provider": "sdrn:fotocasa",  
"timestamp": "2018-09-14T01:02Z",  
"tracker": {  
  "type": "JS",  
  "version": "1.2.3.4-alpha"  
}
```

Summary

SQL on S3 Data
is Super Cool and
gives value from
minute 1



Configuring Presto
to your needs is
possible, if you
have capacity




No magic tool will
help you get gold
of your crappy
data



Links of Interest

Accessing S3 data through SQL with Presto

 bit.ly/access-data-presto

BigData SQL Query Engine benchmark

 bit.ly/query-engine-benchmark

Latest Presto News & Articles

 starburstdata.com/newsletter

