

HYPERLANE

Data, Code, Lab, Factory

Assignment

- Cross-Media Audience Measurement System.
- Dozens of data sources. Varying quality.
- Distributed Data Science and Engineering Teams.
- Global Ops Hub.

Hyperlane ...but why?

Sculley et al (2014):

Undeclared Consumers:

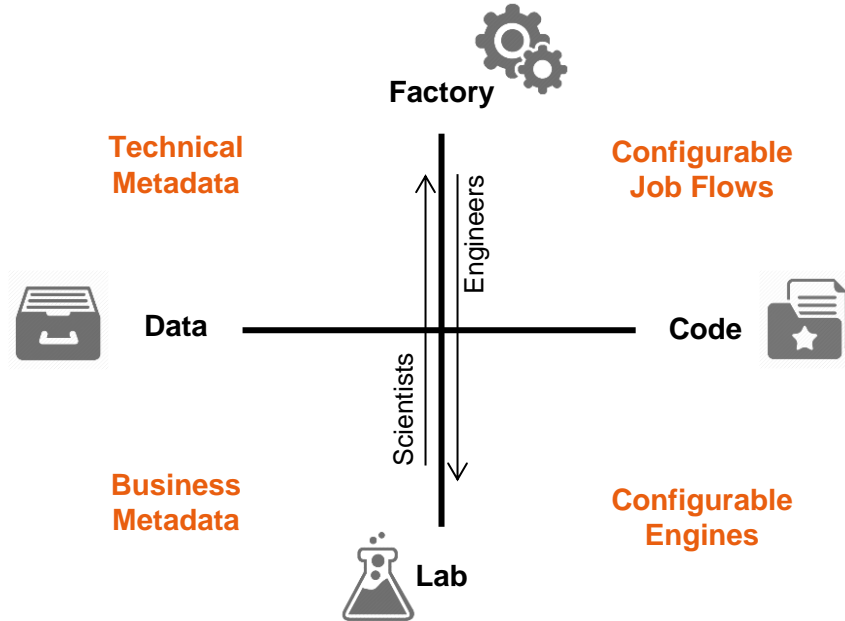
Consuming the output of a given prediction model as an input to another component of the system. Changes will very likely impact these other parts.

Unstable Data Dependencies:

Some input signals are unstable, meaning that they qualitatively change behavior over time.

Static Analysis of Data:

On teams with many engineers, or if there are multiple interacting teams, not everyone knows the status of every single feature.



Sculley et al (2014):

Glue Code:

Using self-contained solutions often results in a glue code system design pattern, in which a massive amount of supporting code is written to get data into and out of general-purpose packages.

Pipeline Jungles:

The system for preparing data in an ML-friendly format may become a jungle of scrapes, joins, and sampling steps, often with intermediate files output.

Configuration Debt:

In a mature system which is being actively developed, the number of lines of configuration can far exceed the number of lines of the code.

User Stories

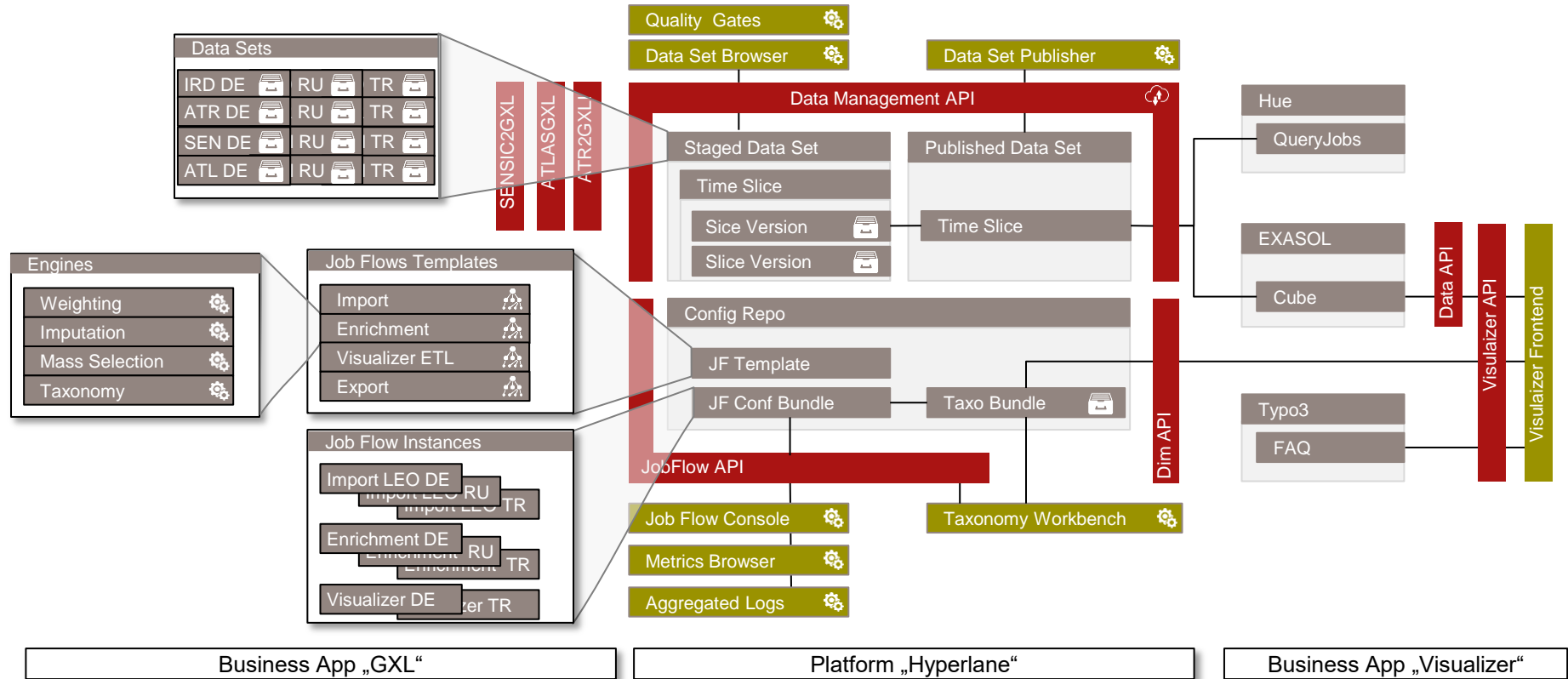
- As a Data Provider..

- I want to provide either (time-based) increments or full snapshots of my data in order to avoid expensive transformations.
- I want to provide data using the most recent schema even if this schema evolves over time in order to avoid expensive schema mappings.
- I want to provide new versions of data that has already been delivered in order to reflect reworks.
- I want to tag and comment my data (-increments, -versions) in order to retain qualitative and business meta-data that has an impact on how the data can/should (not) be used. This can include lineage data encoded in tags.

- As a Data Consumer..

- I want to find and browse (sampled) data sets based on their business meta data, tags, comments etc. in order to find relevant data and understand the status this data has.
- I want to access data using consistent schema and format whenever possible (no matter if the the data was written using evolving schemas or particular technical representations or not) in order to avoid writing expensive and error prone mapping code.
- I want to have (automated) access to growing data sets (Feeds) filtered by tags / business meta data in order to feed it into data pipelines.
- I want to be exposed to immutable data while I'm working on it in order to avoid breaking processing jobs due to concurrent write/update operations (potentially using a different schema) during the run time of my job.
- I want to be able to use standard tooling (Hive, PIG, JDBC, etc.) to access data in order to avoid lock-in and to benefit from the momentum in the engineering and data science community.
- I want to understand at what data center a particular data set is physically stored in order to adhere to regulation and to optimize the processing (e.g. avoid remote data access).

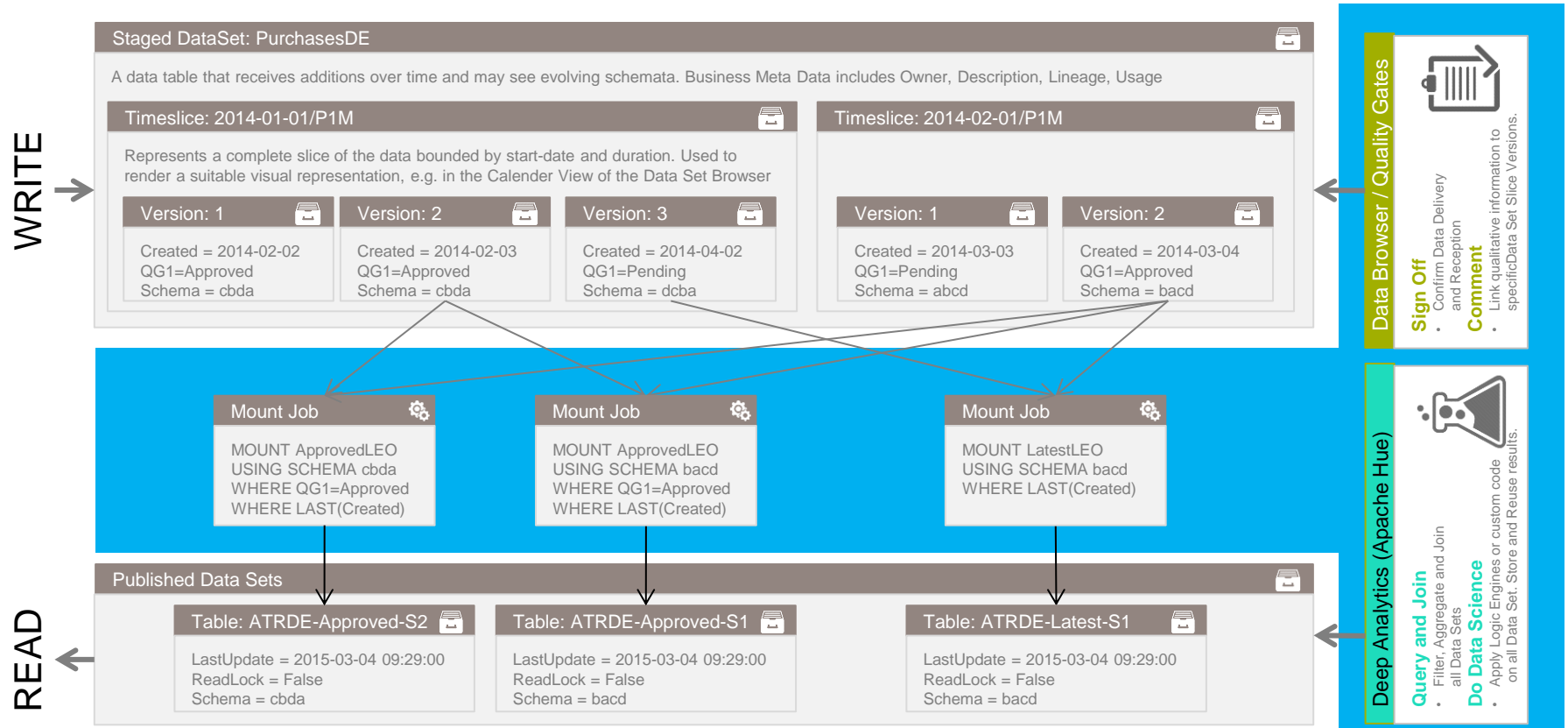
Hyperlane – Functional Architecture



HYPERLANE DATA

Berlin | 16 March 2015 | Get Together

Hyerlane - Data Set Management



Meta Data Management

Challenges:

- Versioned Data, Quality Gates, Comments as Meta Data, Logging
- Wording: Dataset, Instance, Partition, Timeslice? Dataset & Timeslice



Make or Buy?

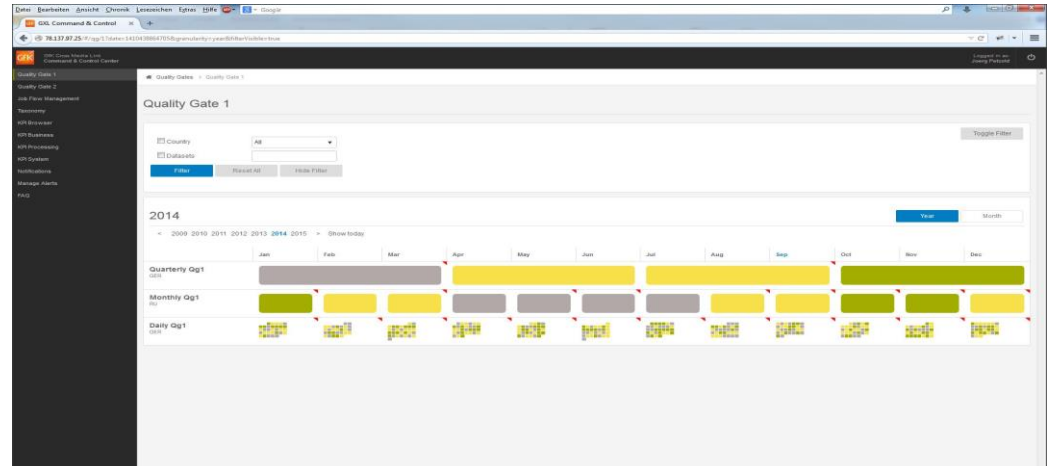
- Cloudera Navigator
+ Policy Engine, - Vendor-Lock-In, Inflexible
- Waterline Data
+ very good for unstructured data, + profiling, + tag propagation, +lineage, - no qg, no key value tags, - product fairly new

Metadata:

1. Semantics: What does it actually mean?
2. Quality: Checked and Approved
3. Validity: What can/can't we do with this data?

Tool: Data Set Browser + QGs

DEMO TIME!



HYPERLANE CODE

Berlin | 16 March 2015 | Get Together

Configurable Engines

- Business Logic is packaged as Pig Scripts
- Takes Parameters and invokes R or JAVA code
- Easy to re-use through Apache Hue and Apache Oozie
- No UI

Weighting Engine

Given the mass (a subset of panelists) and a description (constraints) of the real world online population, this engine computes weights, such that weighted (for the mass) can be considered representative for the actual online population.

This engine applies two phases:

1. Compute input data for the weighting R script
2. Run the R script on the data produces in step 1.

<https://stash.gfk.com/projects/GO/repos/marketing-science/browse/weighting/specification/>

Input Data

Name	Type	Description
Mass selected and Poss/Socio Imputed data (output of Device Possession Imputation)	Semicolon separated file	Output of Device Possession Imputation
setting-constraints	Comma separated file	Settings for the weighting
weighting-constraints	Comma separated file	Constraints for the weighting

Parameters

Name	Description	Example
inputLocation	Location of the input file (i.e. Device Possession Imputation engine output)	/user/hdfs/GO-262/output/3020_mass_completed_201407
outputLocation	Location of the output produced by the engine	Output of Device Possession Imputation
queueName	Name of the queue	root.default
appPath	Base path of the oozie workflow	/user/hyptane/engines/weighting/

Output Data

Pig Wrapper

```
1 SET mapred.min.split.size 1000000000;
2 SET pig.maxCombinedSplitSize 1000000000;
3 SET pig-exec.reducers.max 1;
4 SET pig.splitCombination 'true';
5 SET mapreduce.map.java.opts '-Xmx8G';
6 SET mapred.child.java.opts '-Xmx8G';
7
8 DEFINE R_SCRIPT 'Rscript $r_script_name --min-nsize=100M --min-nsize=10M 1>62'
9 INPUT(stdin USING PigStreaming(';'))
10 OUTPUT('output.csv' USING PigStreaming(';'))
11 SHIP('$r_script_name', '$r_resources_name', '$weighting_constraints', '$weighting_settings');
12
13 Data = LOAD '$input_location' USING PigStorage(';');
14
15 Data = STREAM Data THROUGH R_SCRIPT;
16
17 STORE Data INTO '$output_location' USING PigStorage(';');
```

Pig UDF

```
1 package com.gfk.hyperlane.engine.taxonomy;
2
3 import java.util.Arrays;
4 import java.util.List;
5
6 import org.apache.pig.backend.executionengine.ExecException;
7 import org.apache.pig.data.Tuple;
8
9 /**
10  * @author thevis
11  *
12  */
13 public class AddTaxonomy extends EnhanceOrUpdate {
14
15     static String[] COLUMNS = {
16         "TAXONOMY_COPYRAW_LABEL",
17         "TAXONOMY_COPYRAW_ID",
18         "TAXONOMY_BRAND_LABEL",
19         "TAXONOMY_BRAND_ID",
20         "TAXONOMY_PRODUCT_LABEL",
21         "TAXONOMY_PRODUCT_ID",
22         "TAXONOMY_ACTIVITY_LABEL",
23         "TAXONOMY_ACTIVITY_ID",
24         "TAXONOMY_CONTENT_LABEL",
25         "TAXONOMY_CONTENT_ID",
26         "TAXONOMY_CONTEXT_LABEL",
27         "TAXONOMY_CONTEXT_ID"
28     };
29
30     public AddTaxonomy(String taxonomy, String a2c) {
31         super(taxonomy, a2c, AddTaxonomy.class.getSimpleName());
32     }
33
34     /* (non-Javadoc)
35      * @see com.gfk.hyperlane.engine.taxonomy.EnhanceOrUpdate#getSpecificOutputCalu
36      */
37 }
```

R Code

```
1 luergui <- function(path.in, path.out, path.constraints, path.settings, path.logs, jobstamp,
2   ## This Tool is a GfK SE Marketing Sciences product. ##
3   ## Contact Philipp Gaffert for questions and details.##
4   ## Key developers are: Philipp Gaffert (lead, core)
5   ## Markus Ziegler (GUL logic, IP)
6   ## Dr. Markus Illenthal (outer core)
7   ## Dr. Harek Jiruse (cleaning, error management)
8   ## Dr. Volker Bosch (theory and former versions)
9   ## June, 12th 2013
10  ##
11  # weighted.minmax=FALSE : minimum, maximum bounds are set relative to the design weight
12  # dropcons=FALSE : weighting ignores the constant, i.e. sample / population size is
13  # jobstamp : Integer Job ID
14  # path.func : Location of R sources
15  # path.in : Input data
16  # path.out : output folder
17  # path.constraints : path to constraints file
18  # path.settings : path to settings file
19
20
21  #-----#
22  # Set input file definitions #
23  #-----#
24
25  #path for R functions
26  if (is.null(path.func)) (path.func <- getwd())
27
28  #make important parameters global
29  jobstamp<-jobstamp
30  path.out<-path.out
31  path.logs<-path.logs
32  options(scipen=12)
```



layout.json

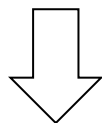
workflow.xml

```
{
  "id": 4,
  "type": "circle",
  "configuration": {
    "name": "first data 4",
    "configuration": []
  }
},
{
  "id": 5,
  "type": "rectangle",
  "name": "Engine 1",
  "configuration": {
    {
      "name": "config-a",
      "type": "string"
    }
  }
},
{
  "id": 6,
  "type": "rectangle",
  "name": "Engine 2",
  "configuration": [
    {
      "name": "config-a",
      "type": "string"
    },
    {
      "name": "config-b",
      "type": "string"
    }
  ]
},
{
  "id": 7
```

[illegible]

JobFlow Templates

Scientists & Engineers



de/config.json

uk/config.json

tur/config.json

```
{
  "tags": {
    "allowUnknownTags": true,
    "dictionaries": ["jsdoc", "closure"]
  },
  "source": {
    "includePattern": "\\.+\\.\\.js(doc)?$",
    "excludePattern": "<^/\\|/\\\\\\\\_>"
  },
  "plugins": [],
  "templates": {
    "cleverlinks": false,
    "monospaceLinks": false
  }
}
```

```
{
  "tags": {
    "allowUnknownTags": true,
    "dictionaries": ["jsdoc", "closure"]
  },
  "source": {
    "includePattern": ".+\\.js(doc)?$",
    "excludePattern": "(^|\\/|\\\\|_)_"
  },
  "plugins": [],
  "templates": {
    "cleverlinks": false,
    "monospaceLinks": false
  }
}
```

```
{
  "tags": {
    "allowUnknownTags": true,
    "dictionaries": ["jsdoc", "closure"]
  },
  "source": {
    "includePattern": "-+\\.js(doc)?$",
    "excludePattern": "(^|\\/|\\\\|_)_"
  },
  "plugins": [],
  "templates": {
    "cleverLinks": false,
    "monospaceLinks": false
  }
}
```

JobFlow Instances

Operations Hub

HYPERLANE WRAP UP

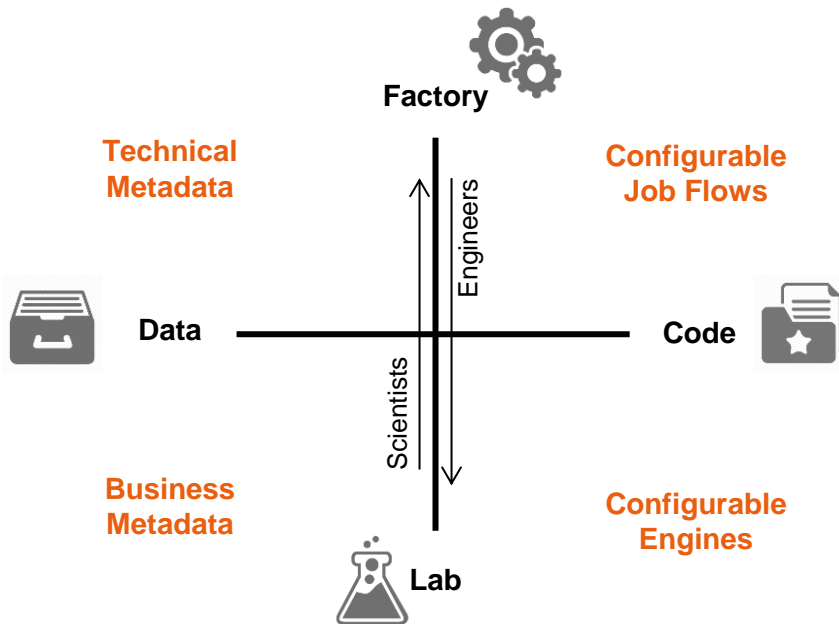
Berlin | 16 March 2015 | Get Together

Hyperlane. Wrap Up.



Improved findability
of relevant data sets using a data catalog that features rich meta data (context, lineage, evolving schemas).

Better collaboration
and strict quality controls enforced through sign-offs and clear accountabilities for providers and consumers of data sets.



More transparency for **productionized Job Flows** visually represented as interactive DAGs used for documentation, configuration and monitoring.

Easier fusion, imputation and weighting of data sets through **reusable and extensible "Engines"** implemented in PIG, JAVA and R.