

LB Dev

ALB APIs

ALB Architecture

Architecture files

File Name	Presentation
FY22 Q4 - Architecture Impact Analysis for ALB DataLab - IT11398 - Solvay Digital Labs	
FY22 Q3 - Architecture Impact Analysis for ALB - IT11398	

ALB Data Architecture - Schema/Models

- ELN Schemas
- ELN Spreadsheets design standards
- Lab Booster Data model
 - Context
 - Objective
- Entity-Relationship Diagram (ERD)
 - Entity-Relationship Diagram design
 - Entity dictionary
 - Entity-Relationship Diagram
- BigQuery
 - Batteries
 - Materials
 - Coatings
- Data mapping

This section aims to present the Data Architecture implemented for data flow in Lab Booster.

ELN Schemas

File Name	Data Model File
Agro	
Battery	
Coatings	
Seed Care	
Actizone	
HPC Flocculation	

ELN Spreadsheets design standards

Design is really important for user experience.

The user feedback about the first version of the ELN templates where more about the design than the content. They didn't really enjoy the "black and yellow" spreadsheets.

So for Coatings' Paint Formulation SS V2, we worked on defining design standards using Solvay's color palette as a base.

These standards are going to evolve according to the future needs.

Here are documentations about developed spreadsheet in PROD:

Agro		
Battery	Conductivity	
Battery	Mechanosynthesis	
Coatings	EP	
Seed Care	Formulation	
Seed Care	Results & Requests	

Lab Booster Data model

Overview

A **Data Model** represents the way data is structured in a dataset or a database, such as Lab Booster's data ocean.

The data model defines how the data lake or data ocean is connected to:

- The data input i.e. ELN, LIMS systems, connected instruments etc.
- The data output i.e. the WebApp DataLab in which users can access data

Context

As of mid-2023, each market in Lab Booster has its own data model i.e. its own way to structure data.

At each new project, connections to the data lake must be built again

Objective

Our aim is to have a common data model for all markets, to bring:

- Accelerated delivery of new projects
- Better performance
- Less maintenance

This page is divided two sections

1. Entity-Relationship Diagram (ERD), which served as a basis to design the data model
2. Data model

Entity-Relationship Diagram (ERD)

Data Models are generally based on a diagram or schema called **Entity-Relationship Diagram** defining

- Entities i.e. a definable object or concept within a system

- Relationships i.e. how entities are related to one another

Building the ERD is a preliminary step to designing the actual data model to ensure that all required entities and relationships are accurately defined and represented.

This section is split in two parts

1. Entity-Relationship Diagram design
2. ERD mapping with R&I workflows

Entity-Relationship Diagram design

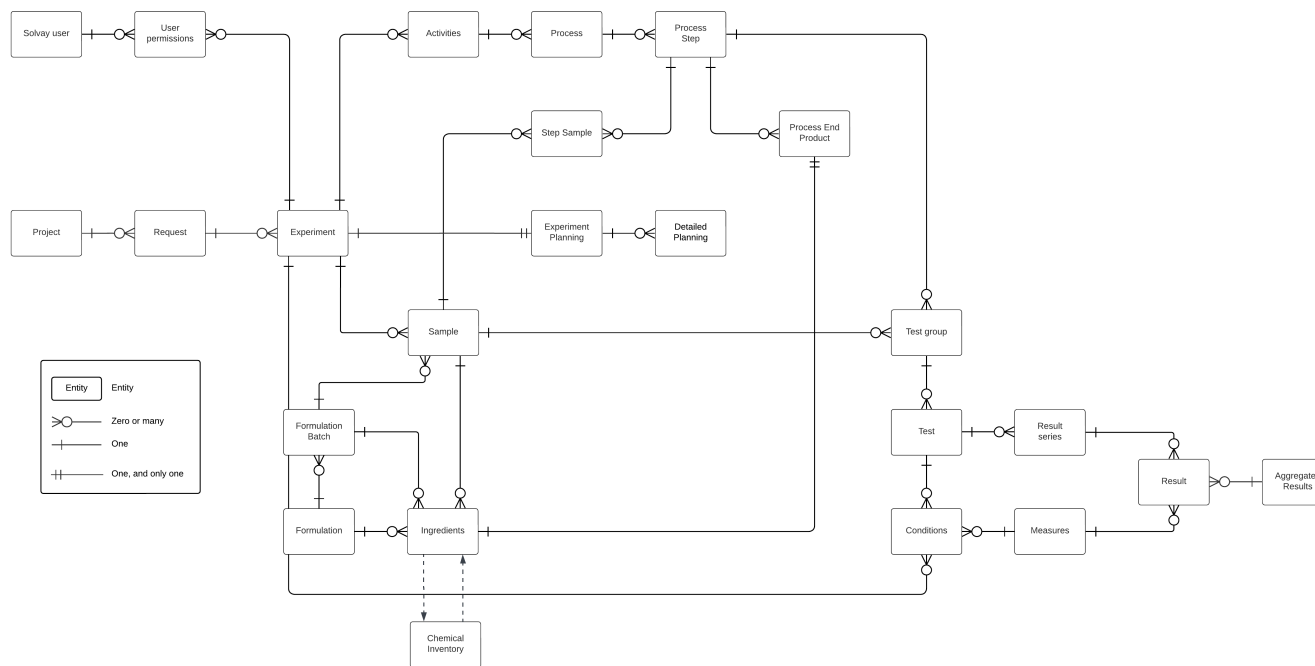
Entity dictionary

Entity	Definition	Example(s)
Experiment	<p>A recording of a workflow performed in the lab by an operator at a given date to achieve an objective</p> <p>An Experiment includes</p> <ul style="list-style-type: none"> • <i>Activities</i> • <i>Samples</i> • <i>Tests</i> • <i>Request</i> • <i>Planning</i> 	<p>Experiments created and recorded in ELN IDBS, LIMS Labware LIMS Agilab...</p>
Solvay User	<p>A recording of the user that created the Experiment, including Solvay ID and email</p>	
User Permissions	<p>A setting determining what application options the user has access to</p>	
Request	<p>A recording of information provided by user requesting an <i>Experiment</i></p> <p>A Request includes</p> <ul style="list-style-type: none"> • Request date • <i>Sample</i> information • Information on user making requests 	<p>Requests for BioMatTech - Biodegradability testing include</p> <ul style="list-style-type: none"> • Request name • Requestor name • Request date • Priority • Status • Test method required • Sample name • Sample ID • Sample status • etc.
Planning	<p>A recording of when the <i>Experiment</i> is supposed to be performed</p> <p>A Planning includes</p> <ul style="list-style-type: none"> • <i>Tests</i> or <i>Activities</i> expected date • <i>Results</i> availability date 	<p>Planning in Novocare - Méréville Request & Results includes</p> <ul style="list-style-type: none"> • Expected application (of slurries & powders on seed) date • Operator performing the application
Activity	<p>A group of <i>Processes</i> performed in the lab in a specific order</p>	<p>In Novocare - Méréville Request & Results, two Activities are found, Application and Testing</p>
Process	<p>A group of <i>Process Steps</i> performed in the lab in a specific order</p>	<p>In BatMat -Mecanosynthesis, the Mecanosynthesis Process is defined by several successive <i>Process Steps</i></p> <ol style="list-style-type: none"> 1. Jar Preparation 2. Milling 3. Drying 4. Calcination 5. Finishing

Process Step	<p>A recording of tasks performed in the lab, defined by its name and date</p> <p>A Process Step includes</p> <ul style="list-style-type: none"> • <i>Conditions</i> in which it is carried out • Input and output <i>Step Samples</i> • <i>Tests</i> performed during Process Step • <i>Process End Product</i> <p>A Process Step follows a Standard Operating Procedure (SOP)</p>	<p>In Aroma - Fermentation the Growth Process Step is defined by the date on which it is performed and includes</p> <ul style="list-style-type: none"> • <i>Conditions</i> - Scale, Temperature, pH... • Input <i>Step Samples</i> - Starter media and Substrate • Output <i>Step Samples</i> - Sample #, Date and Time • <i>Tests</i> - Optical Density and Glucose analysis • <i>Process End Product</i> - Growth media
Process End Product	<p>The chemical output of a <i>Process</i>, defined by its name and date</p> <p>Process End Product characteristics include composition, aspect, mass and/or volume...</p> <p>A Process End Product can be registered as a new <i>Ingredient</i> for other <i>Formulation (Batch)</i> or <i>Process Steps</i></p>	<p>In Aroma - Fermentation, the Process End Product of the <i>Process Step</i> "Bioconversion" is vanillin</p> <p>In Novacare - Méréville Formulation Recipe, the Process End Product of the <i>Formulation Process Step</i> is a formulation</p> <p>In BatMat - Mecanosynthesis Jar Slurries, Amorphous Precursors and Raw Calcined Products are Process End Products</p>
Ingredient	<p>A chemical product, defined by its name and unique ID and recorded in an inventory</p> <p>Ingredient characteristics include date, batch number, supplier, physical state (liquid/solid), density, color...</p> <p>An Ingredient can be:</p> <ul style="list-style-type: none"> • A <i>Formulation Batch</i> • A <i>Sample</i> • A <i>Process End-Product</i> 	<p>In Aroma - Fermentation, the substrate Ferulic acid is an Ingredient</p> <p>In Novacare - Méréville Request & Results, Slurries and Powders are Ingredients</p> <p>In BatMat - Mecanosynthesis Jar Precursors, Slurries, Amorphous Precursors and Raw Calcined Products are Ingredients</p>
Formulation	<p>A combination of chemical products defined by the <i>Ingredients</i>, the <i>Ingredients</i> target proportions and its name</p> <p>Formulation characteristics include total number of chemical products, target concentration, target volume, calculated density...</p>	<p>In Novacare - Méréville Request & Results, a Recipe is a Formulation and is defined by name, ID and label.</p> <p>Characteristics include Number of products, Products, Recipe unit, Recipe Price, Calculated Recipe Density...</p>
Formulation Batch	<p>A combination of chemical products defined by the <i>Ingredients</i> actual proportions, its name, unique ID and date</p> <p>Formulation Batch characteristics include total number of chemical products, actual concentration, total volume, density, container (vessel, jar, bottle)...</p> <p>A Formulation Batch is a <i>Formulation</i> that has been created in the lab</p>	<p>In Novacare - Méréville Request & Results, a Batch of Recipe is a Formulation Batch and is defined by name, ID and label</p> <p>Characteristics include Recipe selection, Actual Weight (of Products)</p>
Sample	<p>A part of a substance or component that is taken from the whole substance or component, defined by its name, unique ID and date</p> <p>A Sample can come from</p> <ul style="list-style-type: none"> • An <i>Ingredient</i> • A <i>Formulation Batch</i> • A <i>Process End-Product</i> • A <i>Request</i> <p>A Sample can be used for</p> <ul style="list-style-type: none"> • A <i>Test</i> • A <i>Process Step</i> <p>See <i>Step Sample</i> for Samples taken during a <i>Process Step</i></p>	<p>Samples come from</p> <ul style="list-style-type: none"> • An <i>Ingredient</i> : Inoculum in BioMatTech - Biodegradability • A <i>Formulation Batch</i> : Batch of Recipe in Novacare - Méréville Formulation • A <i>Process End-Product</i> : Finished Product in BatMat - Mecanosynthesis <p>Samples are used for</p> <ul style="list-style-type: none"> • A <i>Test</i>: Batch of Recipe to characterize at t0 in Novacare - Méréville Formulation • A <i>Process Step</i>: Growth mass used in Bioconversion <i>Process Step</i> in Aroma - Fermentation

Step Sample	<p>A part of a substance or component that is taken from the whole substance or component in relation to a <i>Process Step</i>, defined by its name and date</p> <p>A Step Sample can be</p> <ul style="list-style-type: none"> • An input for the <i>Process Step</i> • An output of the <i>Process Step</i> 	<p>In Aroma - Fermentation, Step Samples are taken throughout the three <i>Process Steps</i> to monitor the chemical reactions</p>
Sample Test Plan	<p>A <i>Planning</i> defined for a set of <i>Samples</i>, defined by its name and the timing</p> <p>The Sample Test Plan characteristics include total number of <i>Samples</i>, <i>Tests</i> to perform ...</p> <p>A Sample Test Plan can apply in the context of</p> <ul style="list-style-type: none"> • A <i>Process Step</i> • A <i>Request</i> • A <i>Planning</i> 	<p>In Novocare - Méréville Formulation the Sample Test Plan defines when <i>Samples</i> should be taken during an ageing <i>Process Step</i></p> <p>It is defined by</p> <ul style="list-style-type: none"> • Protocol name • Initial storage date • Number of <i>Samples</i>
Test Group	<p>A group of <i>Tests</i> performed on the same <i>Sample</i></p>	<p>Characterization tests (OD manual, OD dencytee and Glucose) performed during the Growth <i>Process Step</i> in Aroma - Fermentation for a Test Group</p>
Test	<p>A measure of <i>Sample</i> behavior when a procedure is carried out</p>	<p>Tests performed in BatMat - Mecanosynthesis include Particle size test, SEM test, Lumisizer test, H NMR test, P31 NMR test, Li7 NMR test, Discrete value test</p>
Measure	<p>A property that can be measured</p> <p>A Measure can serve both a <i>Condition</i> and/or a <i>Result</i></p>	<p>pH is a <i>Condition</i> in Aroma - Fermentation and a <i>Result</i> in BioMatTech - Biodegradability</p>
Conditions	<p>A variable or setting defined by the operator for</p> <ul style="list-style-type: none"> • A <i>Test</i> and affecting its <i>Result</i> • A <i>Process Step</i> 	<p>In BioMatTech - Biodegradability, Conditions for the Dry matter <i>Test</i> include Empty aluminium cup weight</p> <p>In Aroma - Fermentation, Conditions of the Growth Process Step include Scale, Temperature, pH...</p>
Results	<p>The outcome of a <i>Test</i> performed on a <i>Sample</i> in specified <i>Conditions</i></p> <p>Results can take the form of</p> <ul style="list-style-type: none"> • A numerical value • A set of numerical values (i.e. curve) • A non numerical value (i.e. observations) 	<p>A pH value is a Result of a biodegradability <i>Test</i> in BioMatTech - Biodegradability</p> <p>A conductivity curve is a Result of a conductivity <i>Test</i> in BatMat - Conductivity</p> <p>Observations are a Result of a Look after Attrition <i>Test</i> in Novocare - Méréville Request & Results</p>
Results Series	<p>A set of <i>Results</i>, obtained at different time intervals, for a <i>Test</i> performed in the same <i>Conditions</i> on the same <i>Sample</i></p>	
Aggregated Result	<p>A <i>Result</i> obtained by aggregating <i>Results</i> from several <i>Tests</i></p>	<p>In Aroma - Fermentation, the maximum amount of vanillin produced during the Bioconversion <i>Process Step</i> is an Aggregated Result as it aggregates several vanillin concentration measure <i>Results</i></p> <p>In Novocare - Méréville Request & Results, averages calculated from two different <i>Test Results</i> are Aggregated Results</p>

Entity-Relationship Diagram



ERD mapping with R&I workflows (WIP)

Three types of R&I workflows were identified

- Formulation workflows
- Synthesis workflows
- Analysis workflows

This was done in order to ensure that the ERD defined accomodates all types of R&I workflows.

The mapping done for different workflows is summarized in the table below.

GBU/F- R&I	Workflow name	Workflow type	Mapping status	Link to mapping	Documentation - Data capture
Novecare GBU	Seed Care Formulation	Formulation	Done	Seed Care mapping	ELN template
Novecare GBU	Seed Care Request & Results	Formulation	Done	Seed Care mapping	ELN template
Battery Platform	Mecanosynthesis	Synthesis	Done	Mecanosynthesis mapping	ELN template
Aroma Performance GBU	Fermentation	Synthesis	Done	Fermentation mapping	ELN spreadsheet mockup
BioMatTech Platform	Biodegradability	Analysis	Done	Biodegradability mapping	LIMS spreadsheet mockup
Specialty Polymers GBU	Aging, Mechanical, Thermal	Analysis	Ongoing		
Specialty Polymers GBU		Synthesis	To do		
Novecare GBU	Agro	Formulation	To do		
Novecare GBU	EP Coatings	Synthesis	To do		
Novecare GBU	Paint Coatings	Formulation	To do		
Corporate R&I	Solvent platform - Solubilization		To do		
Corporate R&I		Analysis	To do		

Green Hydrogen Platform	Conductivity	Analysis	To do		
-------------------------	--------------	----------	-------	--	--

BigQuery

New Data Model of ALB Data Mart (Exposition layer): https://app.genmymodel.com/api/projects/_k07o4lBOEe29ie0vpi-P5A/diagrams/_k07o4oBOEe29ie0vpi-P5A/svg

Data Mapping to Data Mart:

The following BigQuery datasets are all staging as per the data convention explained previously.

For more ETL (extraction, transformation, loading) details, please refer to: [App Lab Booster \(ALB\) - Data](#)

Batteries

Staging

- instr_bat_electro_bol_bcs_conso
- instr_bat_electro_bol_bcs_delta
- instr_bat_electro_bol_maccor_conso
- instr_bat_electro_bol_maccor_delta
- instr_bat_electro_bol_vmp_conso
- instr_bat_electro_bol_vmp_delta
- instr_conduct_conductivity_conso
- instr_conduct_parameter_conso
- instr_mechano_platine_z24_conso
- instr_mechano_rotative_oven_conso
- instr_rawmat_critical_current_density_conso
- instr_rawmat_lithium_metal_compatibility_...
- instr_rawmat_tortuosity_conso

Materials

Staging

- labware_conso
- labware_delta
- labware_testcomponent_conso
- labware_testcomponent_delta
- method_description_conso
- method_description_conso_(1)
- method_description_delta
- raw_data_conso
- raw_data_conso_(1)
- raw_data_delta
- sample_information_conso
- sample_information_conso_(1)
- sample_information_delta
- summary_results_conso
- summary_results_conso_(1)
- summary_results_delta

Coatings


```

1 SELECT Count(*) FROM `prj-labbooster-materials-prod.ODS.raw_data_conso`
2 WHERE test_number = '26888881'

```

Press Alt+F1 for Accessibility Options

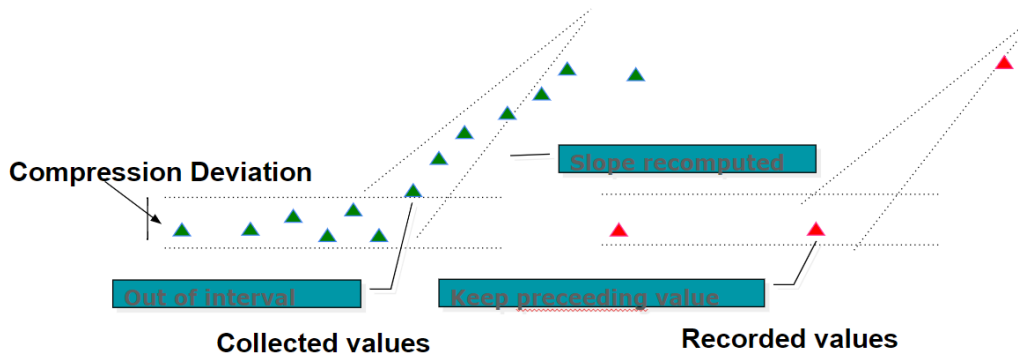
Query results SAVE RESULTS EXPLORE DATA

JOB INFORMATION RESULTS JSON EXECUTION DETAILS

Row	f0_
1	9132

The chart, as plotted on Tableau, is this one (540 values for each trend):

The idea consisted in applying a MES domain' compression algorithm (named Swinging Doors) for minimizing the storage needed on hard disks, and preserving the trend' characteristics!



Normally, this algorithm is applied to only one trend, but with a little modification it's possible to extend it to multiple series at the same time!

Accepting an error of 0.1% (evaluated on the data span of each series), the rows decrease from 540 to... only 36, and the trend is preserved!

In this case the data compression factor is the following

$$\text{Ratio} = 1 - \frac{[\text{Final rows}]}{[\text{Original rows}]} = 93.33\%$$

With an error of 0.5%, we could drop to 17 rows, but the trend is not so smooth (Comp Ratio = 96.8%):

You can find here the test I've made, on Google Sheet:

Raw Data Compression

- on the tab "Original Data", there are the values, as stored in GBQ
- on the tab "Data Table", you can find the "real" table, pivoted by the column "component_name", aggregated by time, and filtered by a specimen (cell C1, white background)
- the maximum error on each variable is evaluated starting from the relative error (cell C5)

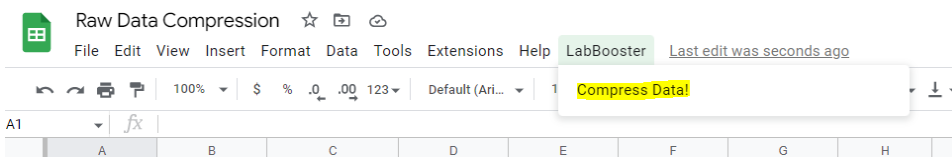
Raw Data Compression ☆ 📄 🌐

File Edit View Insert Format Data Tools Extensions Help LabBooster1 Last edit v

The screenshot shows a spreadsheet with the following data:

A	B	C	D	E	F
	Specimen	2			
	Values	5810			
	Trends	5			
	Exc [%]	0.1			
	Min	-0.06899	18.53328	-0.00018	0.46333
	Max	3.29991	9297.761	2.01725	232.44402
	Exc [-]	0.0033689	9.27922772	0.00201743	0.23198069

- In order to launch the algorithm, it's possible clicking on the custom menu "LabBooster":



- The results are plotted in the "Results" tab;
- The Compression ratio is evaluated in the "KPIs" tab.

More insights

For example, for Lab Booster - Materials, we have 572 tests within our database:

Its raw data is stored in a single table, that has 32 millions rows!

Row	test_number	specimen_n...	time	time_unit
1	27107970	1.0	2.20000	(s)
2	27107970	2.0	51.46300	(s)
3	27107970	6.0	50.41800	(s)
4	27107970	2.0	21.50000	(s)
5	27107970	2.0	6.95000	(s)
6	27107970	5.0	2.45000	(s)

Results per page: 50 1 - 50 of 32264374

In order to present this result to Tableau, we're facing some issues, since Tableau has to copy the entire dataset from Google Cloud Platform to the Tableau Server itself, and this process takes almost 2 hours to complete, preventing us to work with the same update frequency of Talend (Talend - refresh scheduled every hour; Tableau - refresh scheduled every 4 hours).

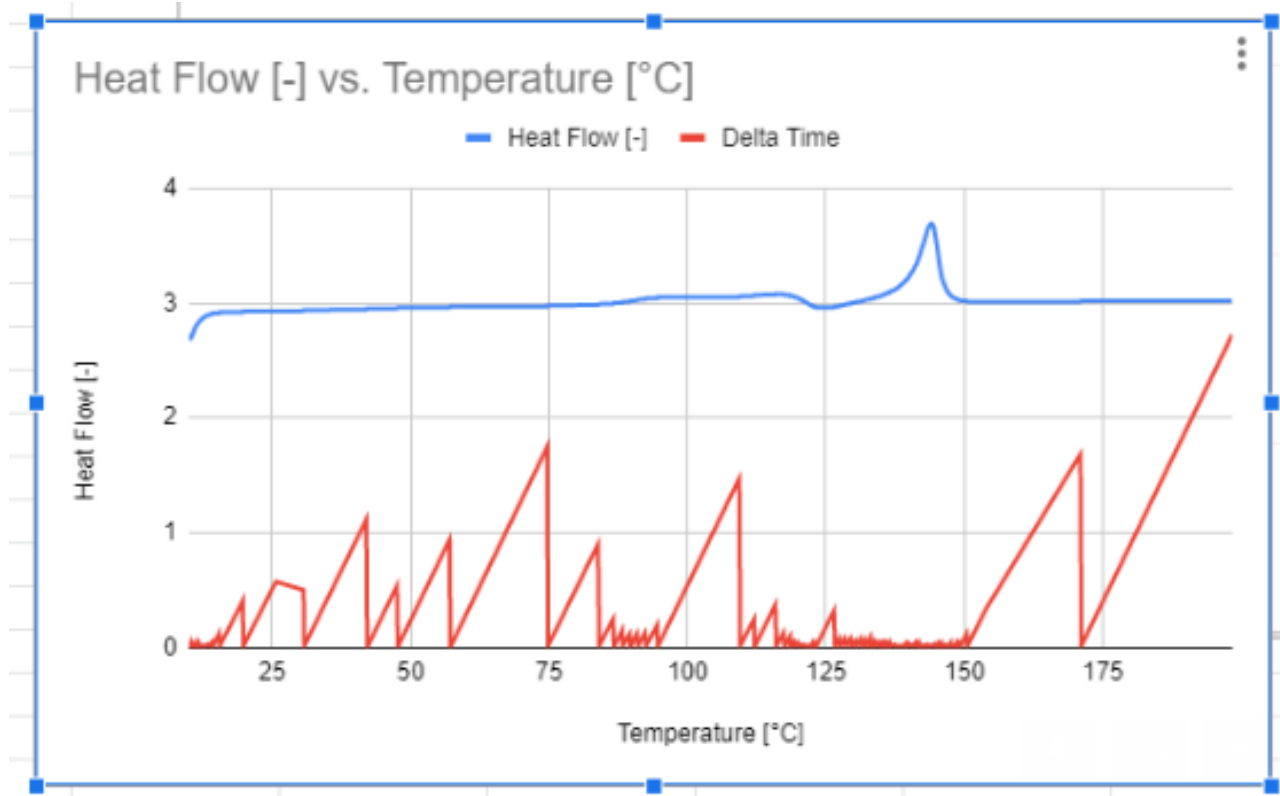
This will be the case also of DSC Thermal analyses, for example: ADC/Labs database is not designed to host very large tables of raw data (as far as I understood from the tech lead of ACD). For a single DCS analysis, we should save the data of three ramps, for having the info of time/temperature/heat flow. In each ramp, the time is recorded every 0.01 s, and we could have several values (I analyzed the 1st heating of a test in Bollate):

To store this data within a database' table, we need 3420 rows (more than 10k for storing the three ramps). So, if we have 100 analyses to store within the database, we need more than 1 million rows.

The [compression approach](#) is data analytics-friendly: instead of storing all the data, this approach simplifies the data itself, in order to not store data that is "correlated", with an adaptive sampling step!

- the first step consists in defining a maximum error, named deviation. A reasonable value is 0.1% of the trend span (*i.e.* the difference between the maximum and the minimum of the trend itself)
- At every step, the algorithm evaluates the derivative between the previous and the current point.
- An isosceles parallelogram is defined:

and this one for the compressed values (with this approach, we have the maximum definition on the "curve portions", and a "less-dense" point definition when the derivative is constant): the red line shows the time span between the actual compressed value, and the previous one!



The compressed trend is created with 6 times less the original points, and it preserve all the information for the 1st/2nd grade derivatives:

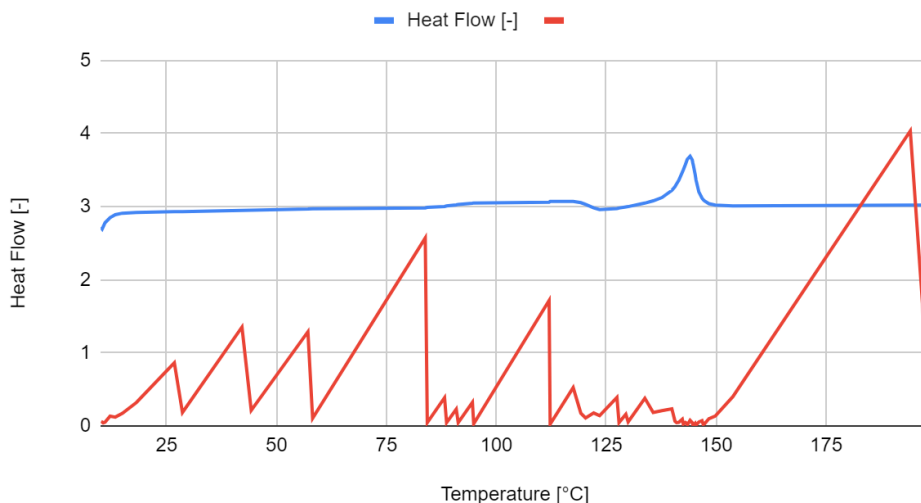
Original Values	3420
Compressed Values	552
Reduction	83.86%

The "missing" points could be reconstructed through a linear interpolation approach, committing an error smaller than the 0.1% of the series' span!

If we vary the derivative error limit, here you can find some results (in red the step size):

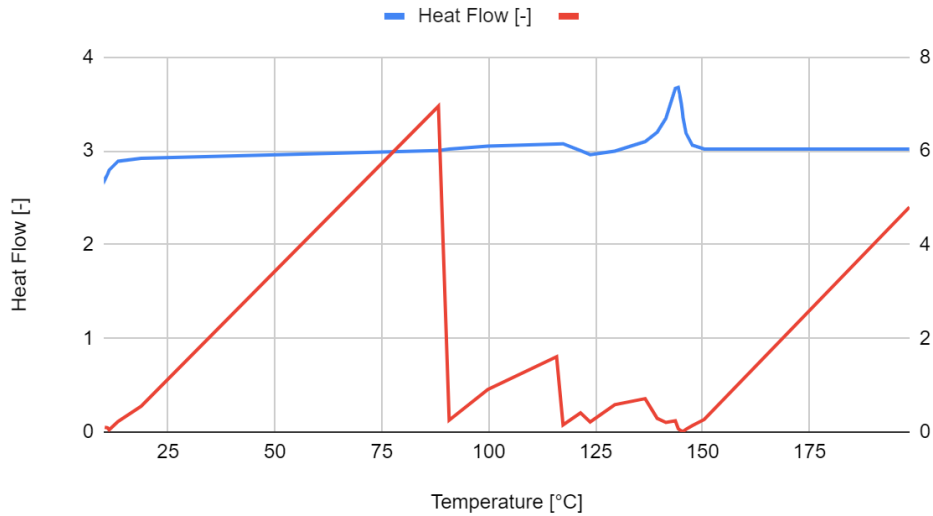
- Error = 0.5% --> 243 Values - Compr. Ratio = 93%

Heat Flow [-] vs. Temperature [°C]



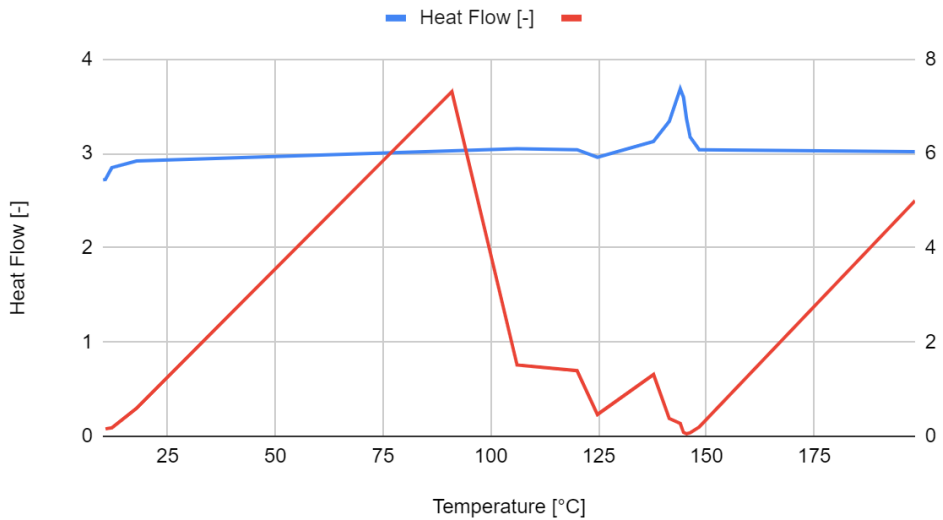
- Error = 1% --> 107 Values - Compr. Ratio = 97% (we start losing some information; the axis for the time step is on the right)

Heat Flow [-] vs. Temperature [°C]



-
- Error = 2% --> 63 Values - Compr. Ratio > 98% (we've lost the info regarding 2nd grade derivative)

Heat Flow [-] vs. Temperature [°C]



-

In my opinion, we could have two possibilities:

- applying this algorithm for capturing data, defining a small derivative error (= 0.1%), and saving between 80% and 90% of table rows
- applying this algorithm for visualization purposes only: there will be two "raw data tables" - original and compressed - and we expose to Tableau only the compressed ones.

The implemented algorithm

Input: a plain table, that has in 1st column the independent variable, and in the others the dependent variables;

Output: the compressed table, with the same structure, but with much less rows!