

Multidimensional design

Training session

This training session is focused on the requirement-driven design of multidimensional schema (MD) in a data warehouse (DW).

Prerequisites:

- All examples are created over the LEARN-SQL system. For better understanding of the system under the study, the schematic/diagrammatic representation of the ontology is showed in Figure 1.
- During the training session as well as in the lab, we will use the *Indyco Builder* tool for conceptual modeling of a DW.
 - Conceptual modeling of a MD schema
 - Uses Dimensional Fact Model (DFM)
 - Intuitive graphical user interface
 - Validation
- We strongly suggest you to check more details and get familiar with the *Indyco Builder* tool following the given links.
 - Dimensional Fact Model (DFM)
 - First publication: (Golfarelli, Maio, & Rizzi, 1998)
 - Wiki: http://en.wikipedia.org/wiki/Dimensional_Fact_Model
 - Extended guidelines: <http://www.indyco.com/kb/dimensional-fact-model/>
 - Indyco Builder tool
 - Official web page: <http://www.indyco.com/builder/>
 - User manual: <https://www.youtube.com/watch?v=haXLU0M2PZQ>
 - Download: <http://www.indyco.com/start-here/request-an-educational-license/> (follow the instructions in Learn SQL on how to download and install Indyco Builder).

This training session has two assignments:

Assignment 1:

Requirement-driven design of a multidimensional schema

Considering the diagrammatic (graphical) ontology representation of the Learn-SQL domain in Figure 1, determine the multidimensional context of the given requirement inside the given ontology and produce valid MD schema solution(s) that satisfy the requirement at hand.

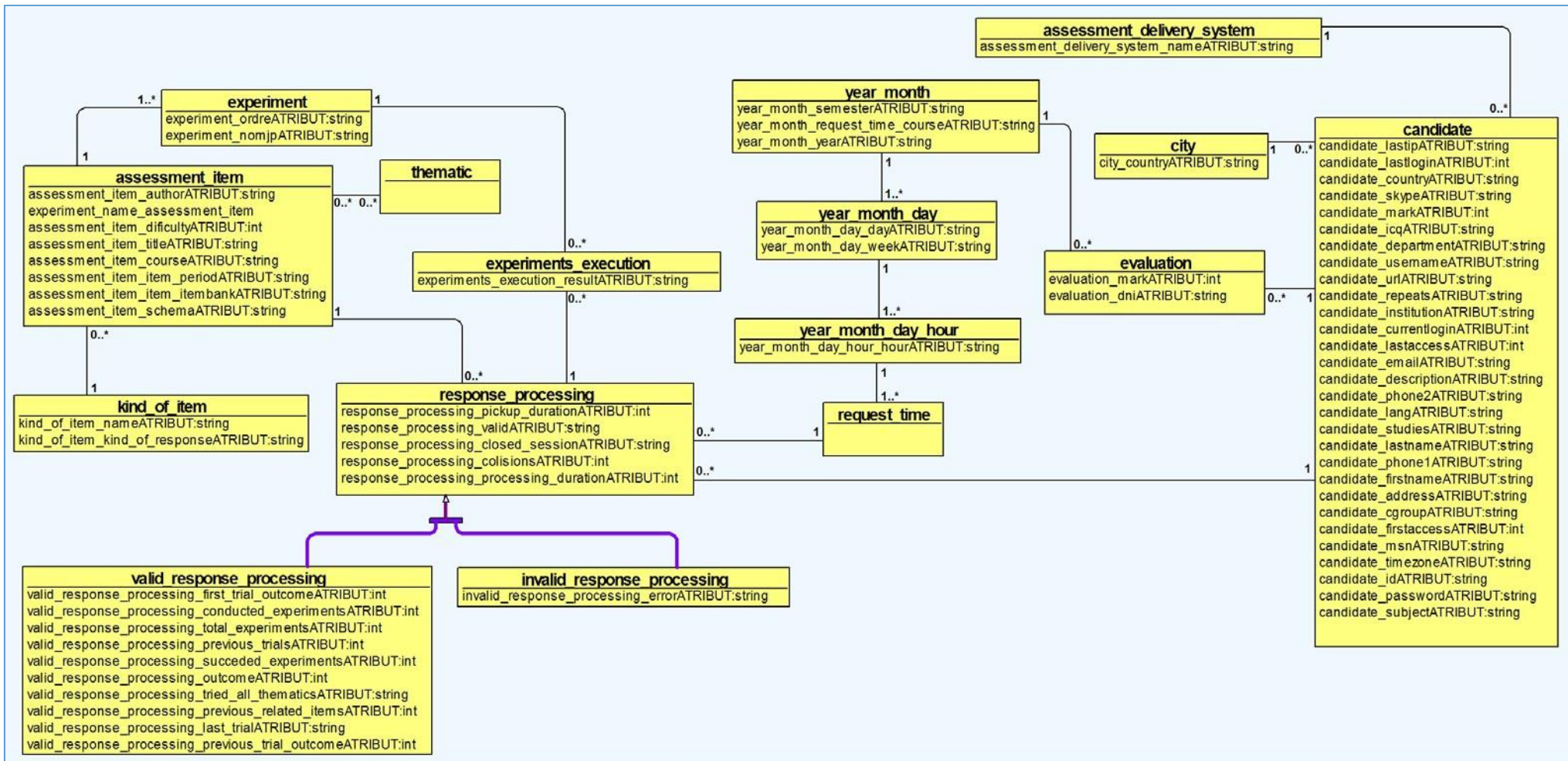


Figure 1: Diagrammatic representation of the domain LEARN-SQL ontology

Assignment 1

Requirement-driven design of a multidimensional schema

New information requirement is posed by the Learn-SQL user:

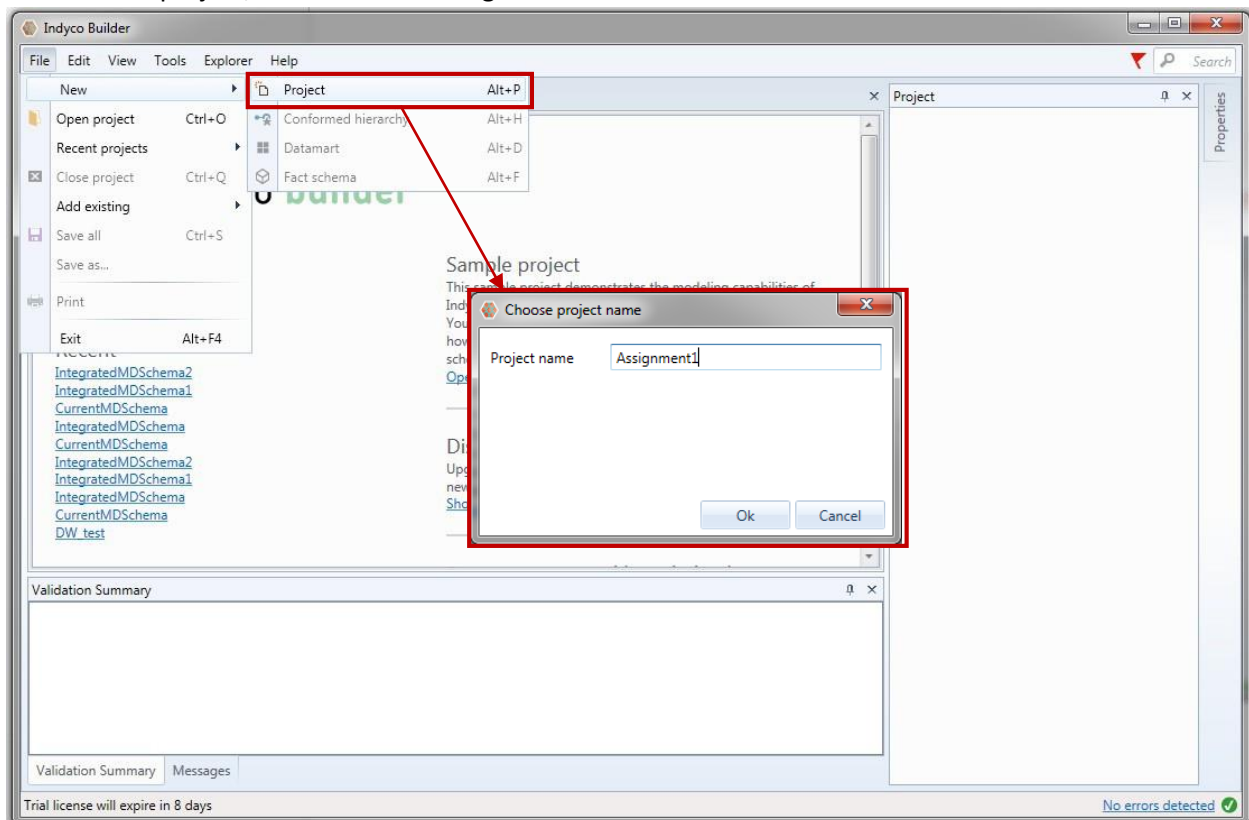
"I want to analyze the average processing duration for each week and candidate where the difficulty level of the assignment item is higher than 7."

Considering the diagrammatic (graphical) ontology representation of the Learn-SQL domain in Figure 1, determine the multidimensional context of the given requirement inside the given ontology and produce valid MD schema solution(s) that satisfy the requirement at hand.

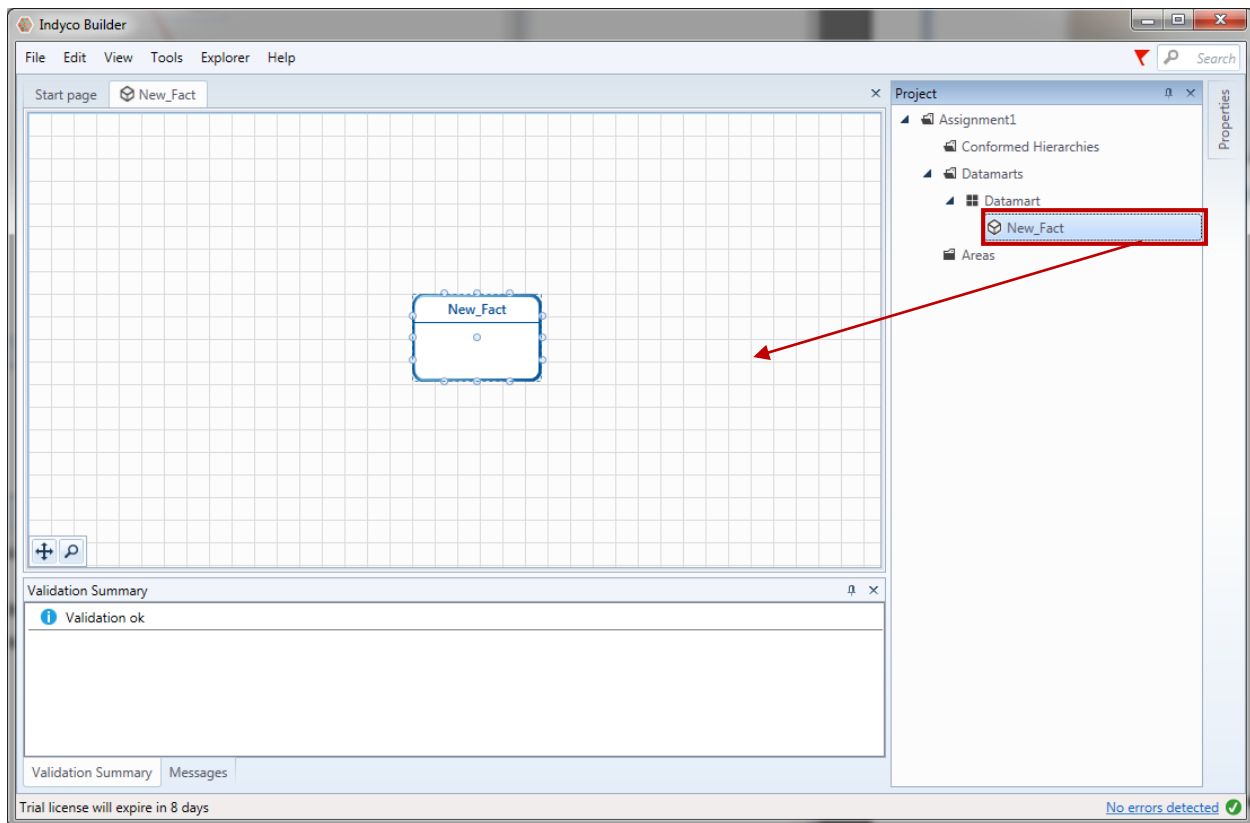
Important note: Different alternative solutions (if more than one) should be provided in different data marts inside the same Indyco project.

MD schema design

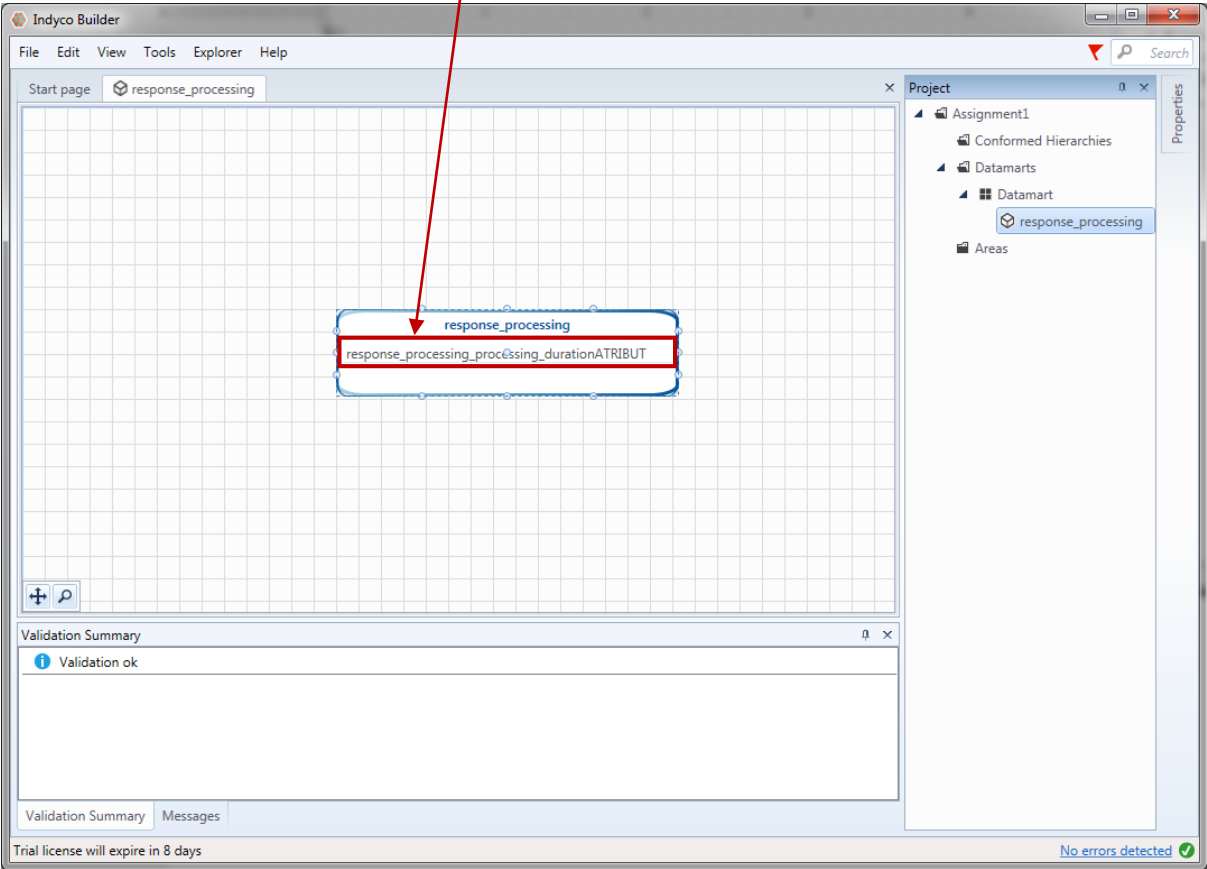
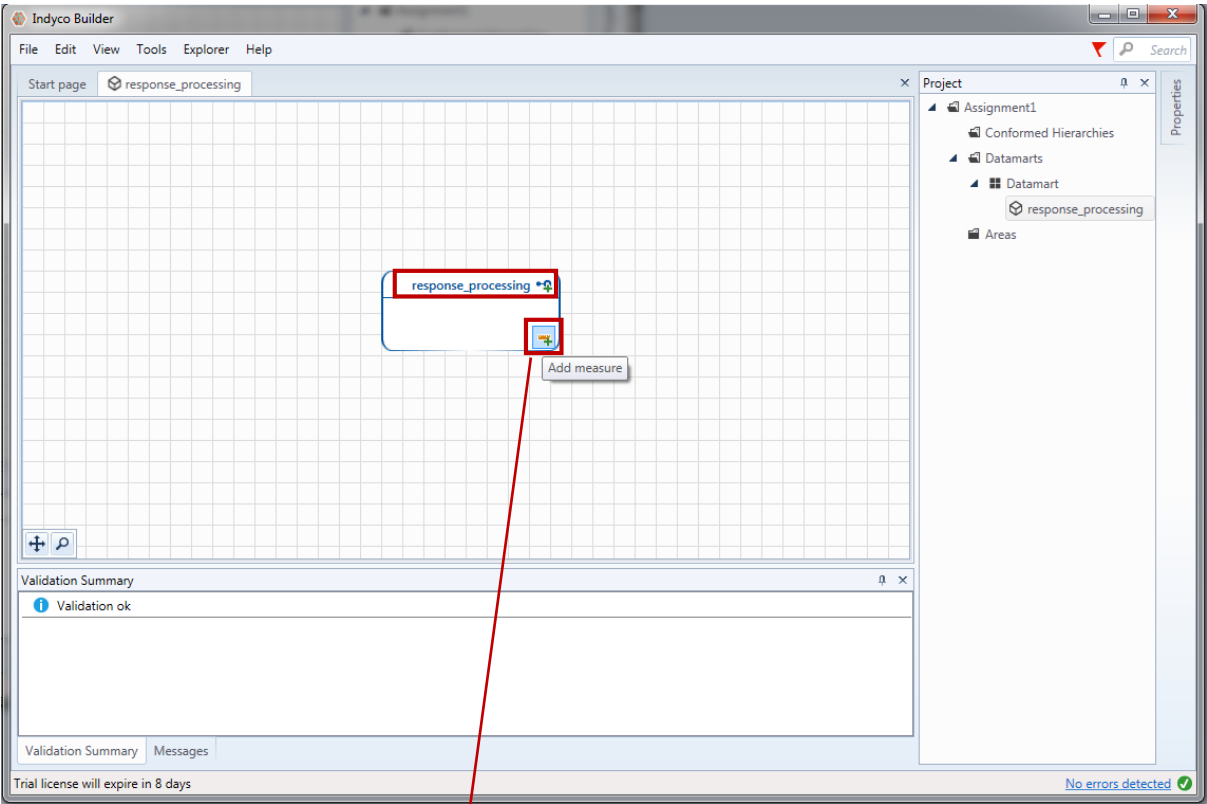
1. **Opening the Indyco Builder tool.** After installing Indyco Builder, you should open it and create new project, as showed in the figure below



After we created a new project, the empty project with an initial empty fact will appear on the right as showed in the figure below. After double-click on the fact, the canvas where the MD design is created for that fact should be visible in the screen as shown in the figure below.

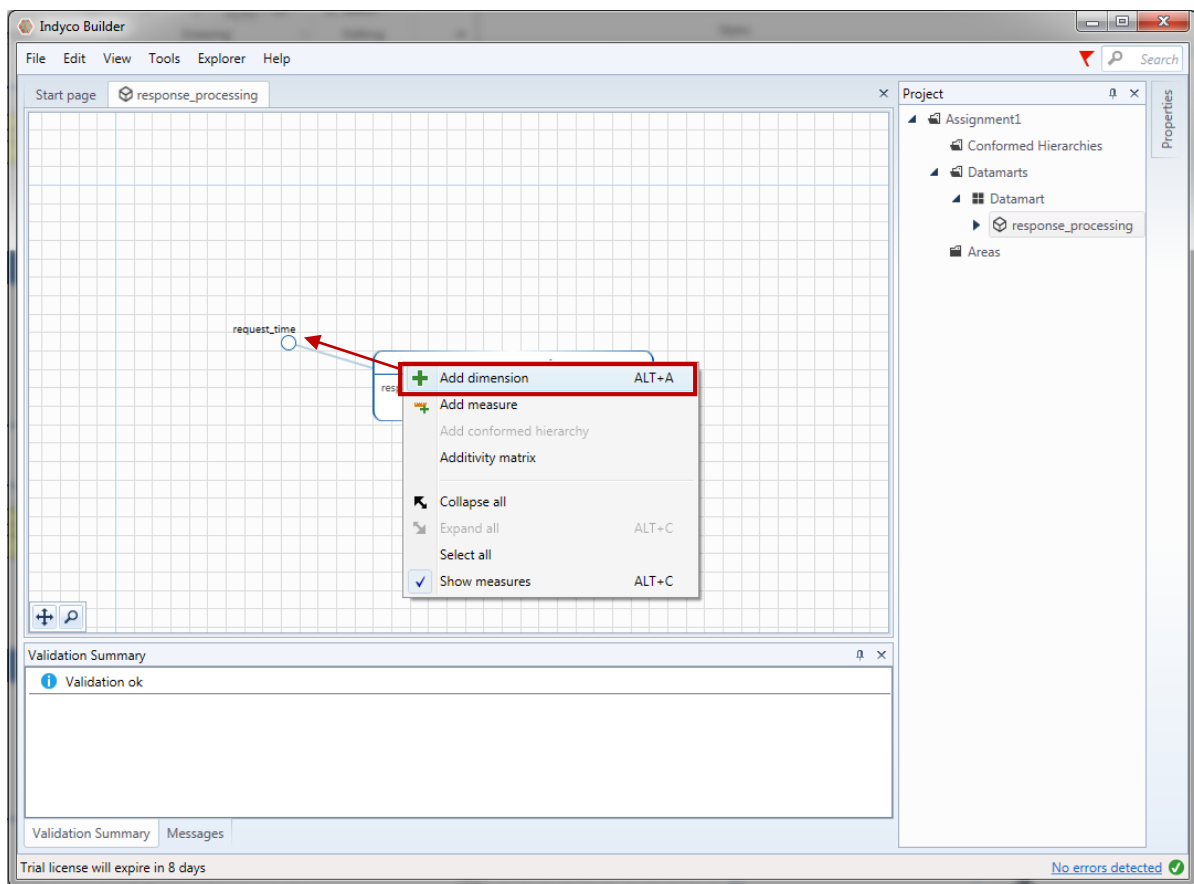


2. **Define the focus of the analysis.** Starting from the above business requirement, we should first identify how the focus of the analysis (a potential fact) maps over the domain ontology (Figure 1). Thus, in our example, we identify that “average processing duration” can map to *response_processing* concept in the LEARN-SQL ontology, more precisely to the attribute *response_processing_processing_durationATRIBUT*. Thus we insert a new (or modify the initial empty) fact, and further add its measure (i.e., *processing_duration*), as showed in the figures below.

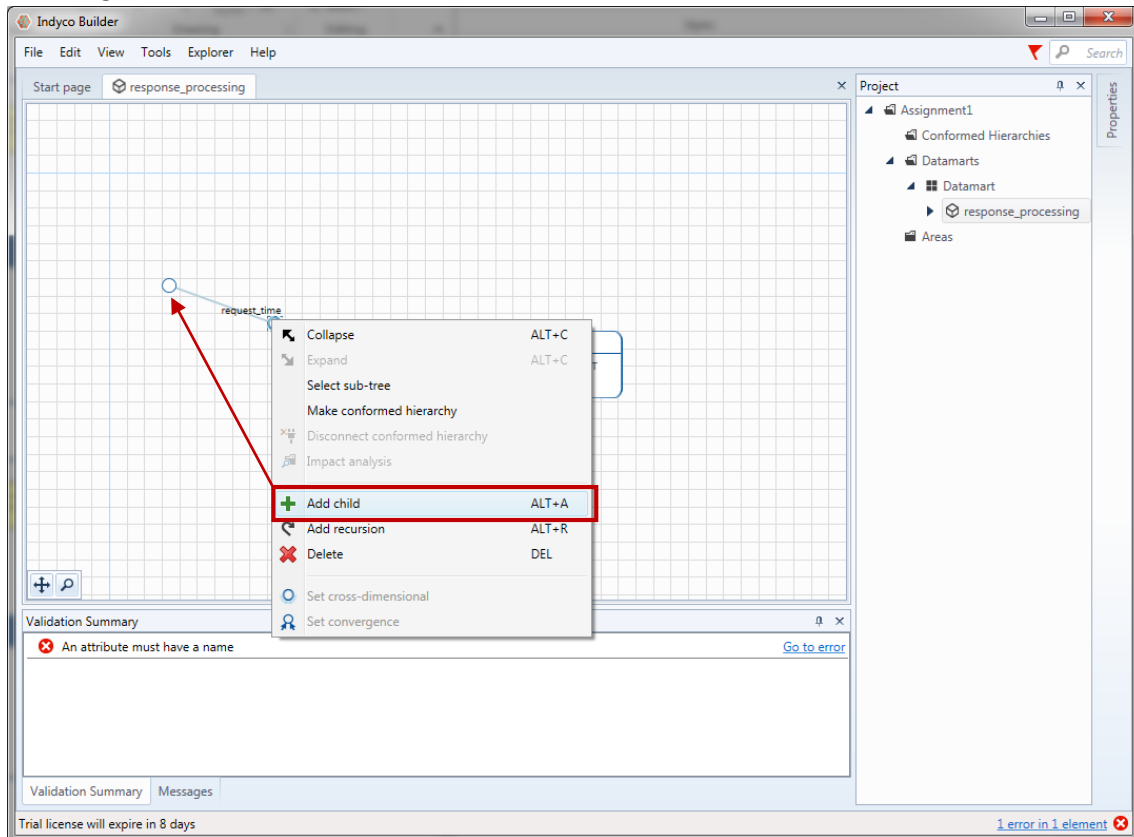


- 3. Adding dimensions.** After adding the focus of the analysis, we need to define the perspectives from which we analyze the added fact (i.e., dimensions).
- Following the given requirement, we identify two dimensional concepts from which a user wants to analyze the *processing duration* measure, i.e., *week* and *candidate*.
 - We can then map them over the given ontology respectively to concepts *year_month_day.year_month_day_week* *ATRIBUT* and *candidate.candidate_id* *ATRIBUT*.
 - Additionally, we need to identify in the ontology, how these dimensional concepts are related to the previously found factual concept (i.e., *response_processing*).
 - If in the ontology the given fact is (by transitivity) related by means of 'to-one' relationship to the given dimensional concept, then we can add the given dimensional concept to the design.
 - Other concepts identified on the path from the factual to the dimensional concept (i.e., *intermediate concepts*), which do not originate from the given requirement, can potentially play the role of both dimensional and factual concepts. This can therefore generate multiple possible alternative solutions. In addition, each solution must respect the MD integrity constraints, explained in Appendix A.

Going from the first solution, where all *intermediate concepts* are interpreted as dimensional concepts, we add them as dimension hierarchies to the previously added fact, as showed in the figures below. Dimension is added by right-click on the fact and choosing *Add dimension*.



Then each further level of the same dimension hierarchy is added by right-click to the last level and choosing *Add child*.



Finally, we obtain the first solution of a MD design as shown in the following figure.

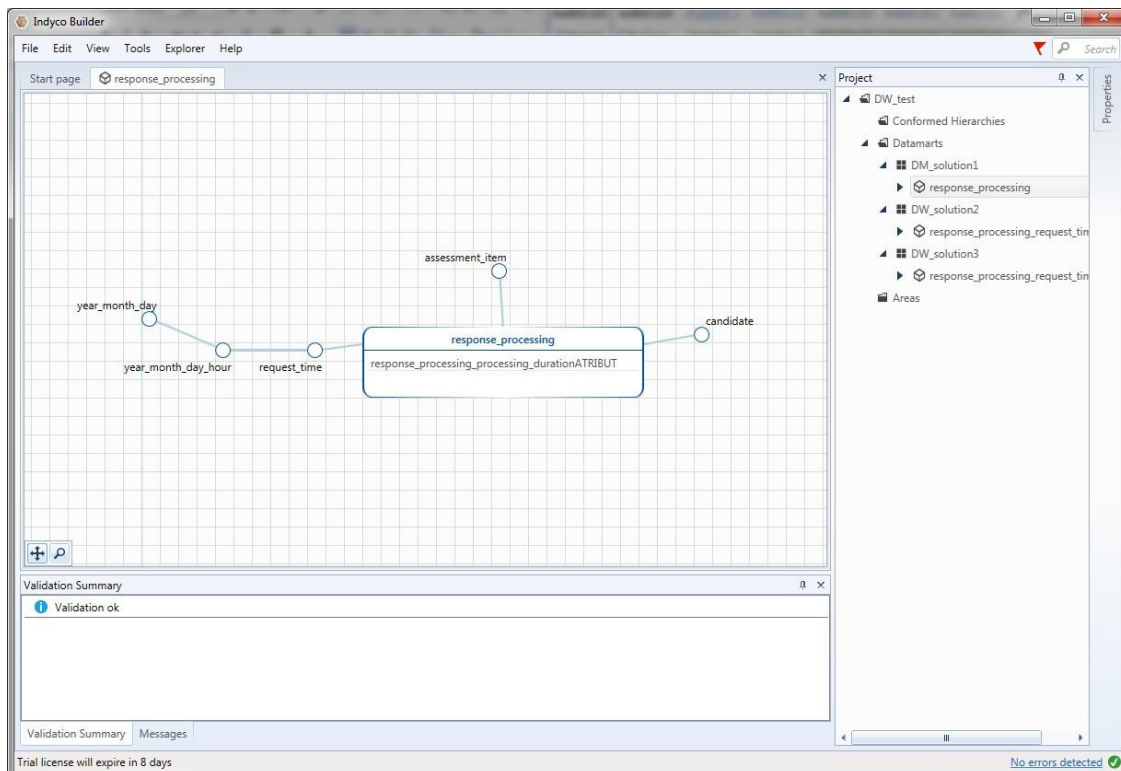


Figure 2: Solution 1

Futhermore, we can obtain two more potential solutions where the *intermediate concepts*, i.e., *request_time*, or *request_time* and *year_month_day_hour* play factual roles. This in fact means that the factual concept will be at different levels of granularity (*year_month_day_hour* or *year_month_day*). Notice that the solution where *request_time* plays a dimensional role and *year_month_day_hour* plays a factual role is incorrect as it does not respect the MD integrity constraints (see Appendix A).

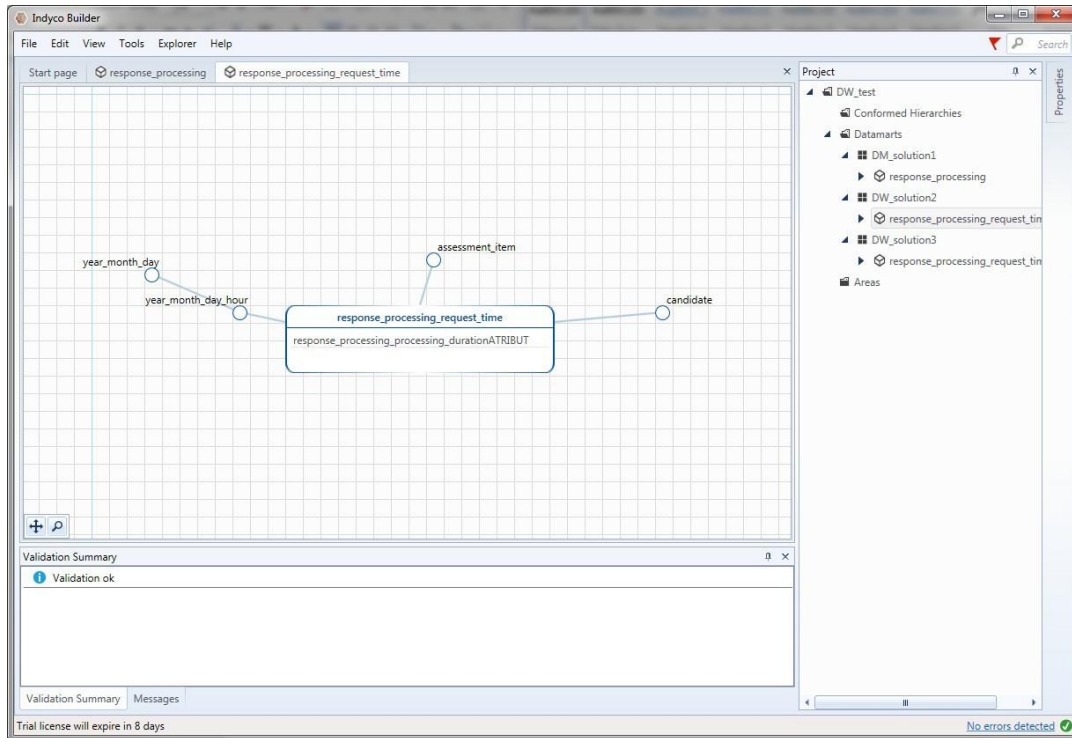


Figure 3: Solution 2

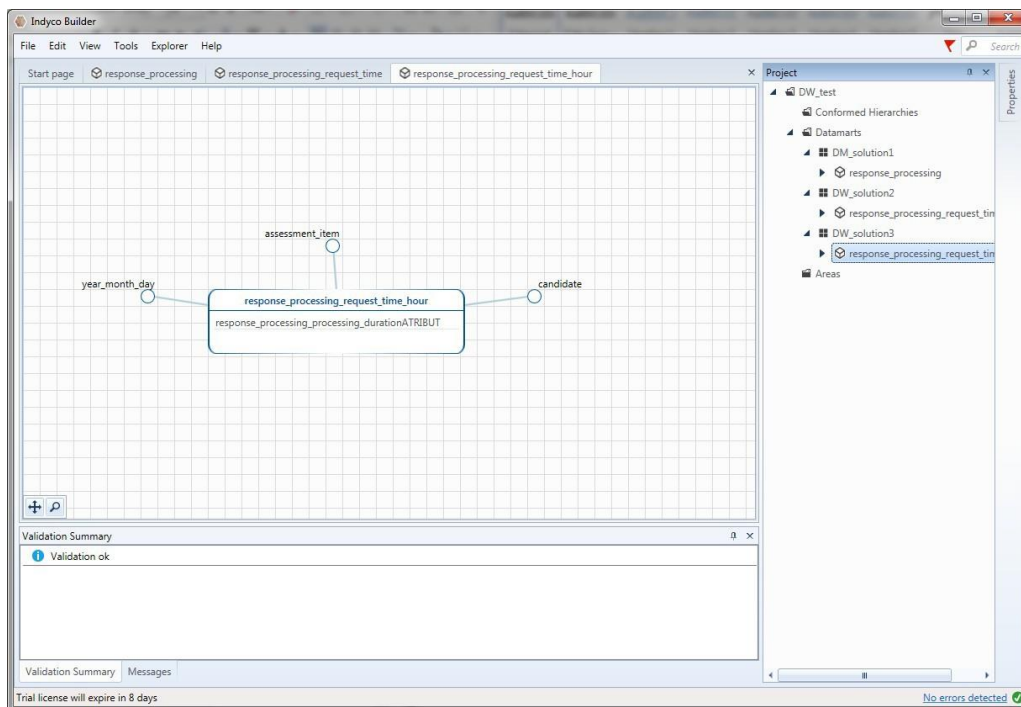
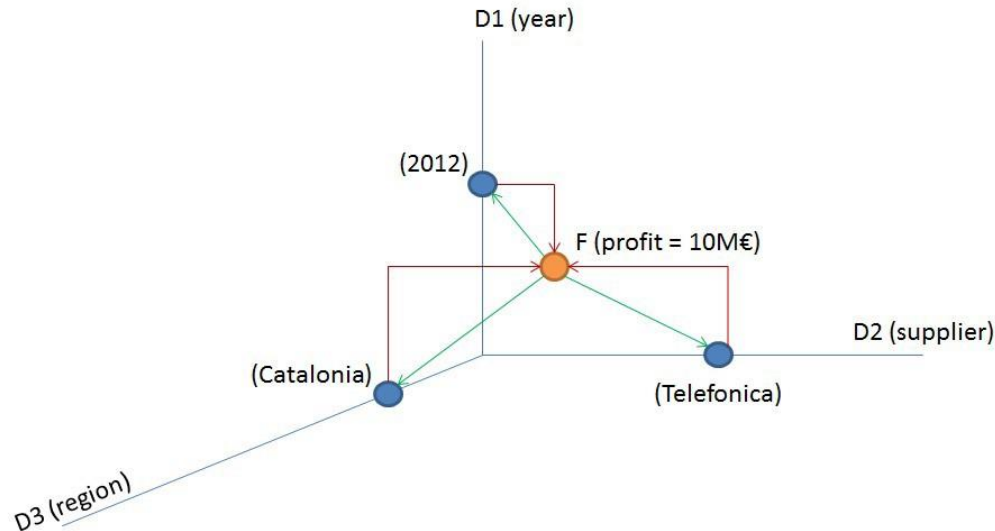


Figure 4: Solution 3

Appendix A

Multidimensional Design - integrity constraints

- MD space: the dimensions over the fact must arrange the correct MD space. On the one side, for each fact instance (e.g., profit = 10M€), there must one and only one dimensional concept in each of the given dimensions (e.g., year = 2012, region = Catalonia, supplier = Telefonica). On the other side, for each set of instances (one from each dimension), there can be one and only one fact instance. See the figure below.
- Summarizability: Data summarization must be correct which is ensured by applying necessary conditions for summarization (see more details in (Mazón, Lechtenbörger, & Trujillo, 2009)):
 1. *Disjointness* - the sets of objects to be aggregated must be disjoint;
 2. *Completeness* - the union of subsets must constitute the entire set; and
 3. *Compatibility* of the dimension - the type of measure being aggregated and the aggregation function.



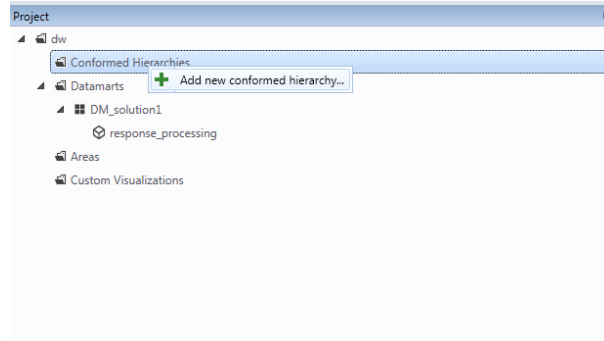
MD schema solutions

- The produced MD schemas should:
 - **Respect the above MD integrity constraints**, and also
 - Have the **minimal structural complexity** (schema-wise), e.g., number of facts, number of dimensions, attributes, associations, etc.
- To achieve the minimal design we need to look for the parts of the schema that can be:
 - **reused** in answering new requirements
 - correctly **collapsed** with the adjacent parts to build single MD concept

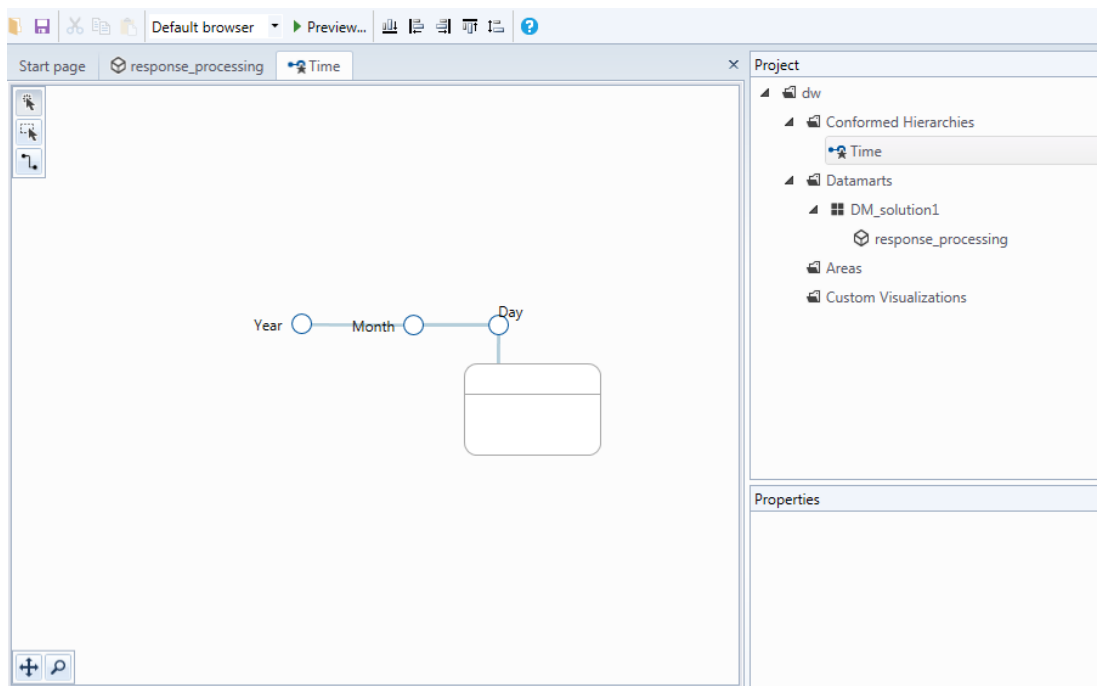
Appendix B

Shared dimensions (Conformed Hierarchies)

When you create multiple datamarts, it might be the case that you want to reuse a dimension previously defined (for instance the *Time* dimension). *Conformed Hierarchies* allow you to define once a dimension and reuse it many times in different datamarts. Create a conformed hierarchy as depicted in the following figure:



Then, you will be able to define the different levels composing it as usual.

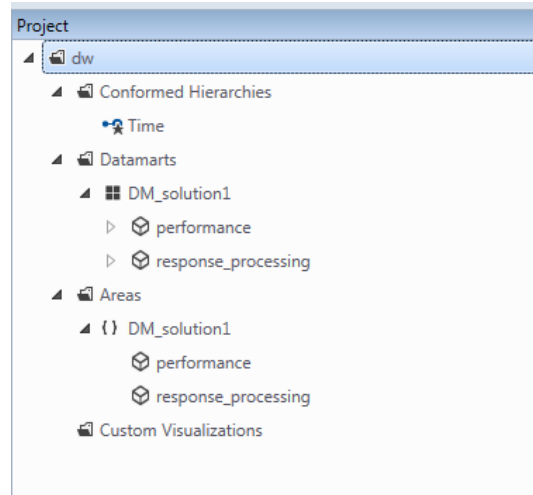


To reuse it in different data marts, you only need to drag and drop the conformed hierarchy to the fact that we want to connect with. Indico will request the *root* of the conformed hierarchy, which is the lowest level of granularity that you want to connect the fact with.

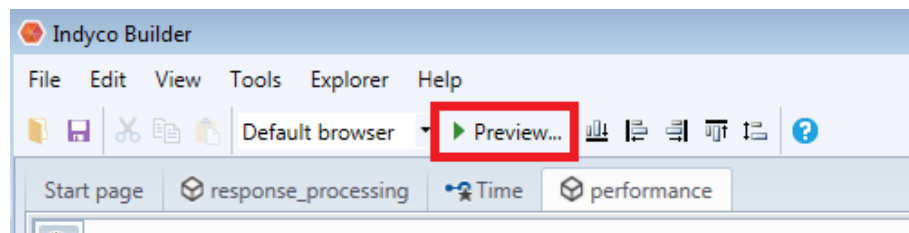
Appendix C

Visualizing all fact schemas at once

Indyco offers the possibility to visualize the different fact schemas from one datamart at once. First, you will need to create an *Area* element and copy (drag and drop) the fact schemas you want to visualize.



When you click “Preview”, Indyco will open a web browser with the visualization of the different fact schemas.



Bibliography

Golfarelli, M., Maio, D., & Rizzi, S. (1998). The Dimensional Fact Model: A Conceptual Model for Data Warehouses. *Int. J. Cooperative Inf. Syst.*, 7(2-3), 215–247. doi:10.1142/S0218843098000118

Mazón, J.-N., Lechtenbörger, J., & Trujillo, J. (2009). A survey on summarizability issues in multidimensional modeling. *Data Knowl. Eng.*, 68(12), 1452–1469.