# Concept*Mix*:
# Self-Service Analytical Data Integration Based on the Concept-Oriented Model

## Alexandr Savinov
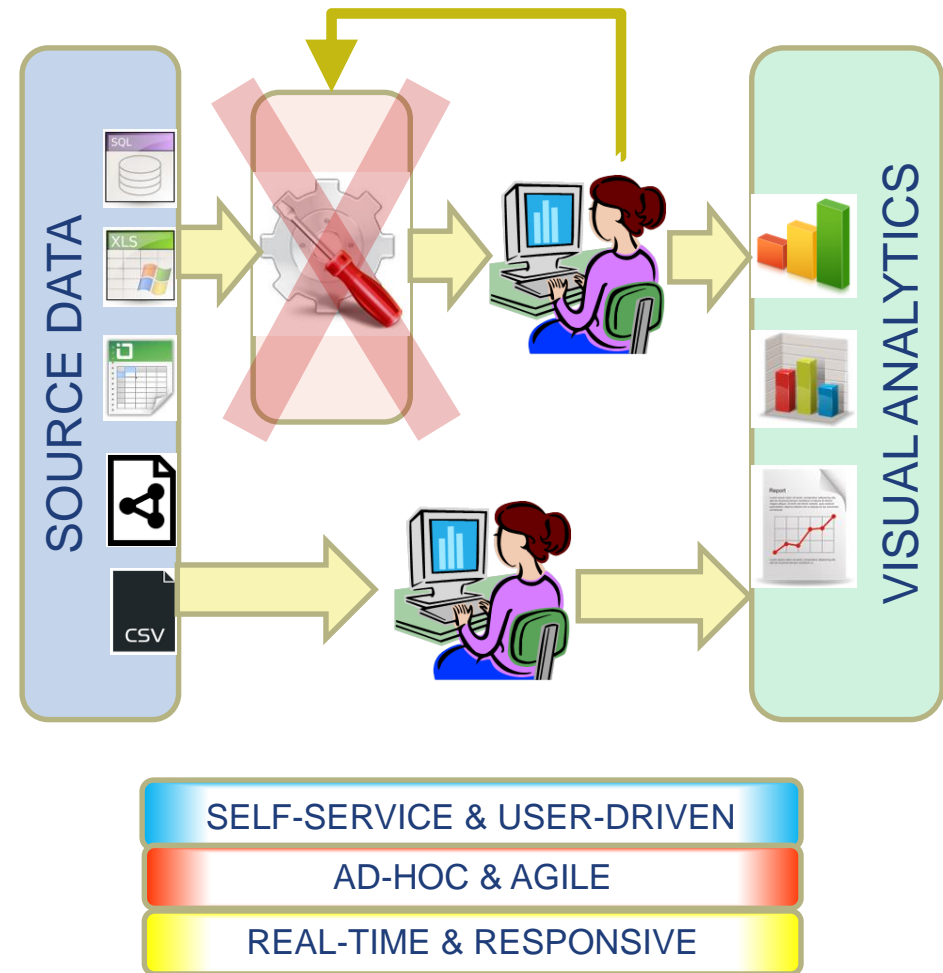
Database Technology Group
Technische Universität Dresden, Germany
Data Commander - http://conceptoriented.com

ESF
Europäischer Sozialfonds
für Deutschland

eXIST
Existenzgründungen aus der Wissenschaft

dresden|exists
WISSEN. GRÜNDEN. UNTERNEHMEN.

# PPROBLEM

- Variety of data sources: one aspect of the big data problem
  - ◆ Integrate: data sources have to be mashed up to produce the desired result
- Data wrangling (curation, munging, scraping) – the most tedious part of the overall analysis process
  - ◆ Transform: refactor the structure of data (schema)
- Original data does not have data the user needs
  - ◆ Analyze: new attributes have to be computed



SELF-SERVICE & USER-DRIVEN

AD-HOC & AGILE

REAL-TIME & RESPONSIVE

Challenge: How to simplify operations with data so that the tool can be used by non-IT users?

# PRODUCT VISION

**Data sources**

**Formula bar**

Product Categories
- Id
- Name

Orders
- Id
- Amount

Customers
- Id
- Country

= COUNT( this <- (Orders) -> (Customers) )    ④

| Category | Totoal Amount | Customers |
|---|---|---|
| Drinks | 50.000 | 876 |
| Electronics | 10.543 | 356 |
| Garden | 3.826 | 84 |
| Toys | 23.82 | 1.539 |

**Mash-up**

① ② ③

- ◆ ConceptMix: <u>self-service</u> data integration, transformation and analysis tool
- ◆ Concept*Mix* is <u>column-oriented</u> rather than cell-oriented
- ◆ Data is defined by <u>column formulas</u> (4) rather than cell-formulas
- ◆ Drag-n-drop a source column (1-3) with automatic recommendations

# TECHNOLOGY

- ● Key enabler: concept-orientation:
  - ◆ Concept-oriented model of data (COM)
    - ▶ Unified model: simple and natural representation
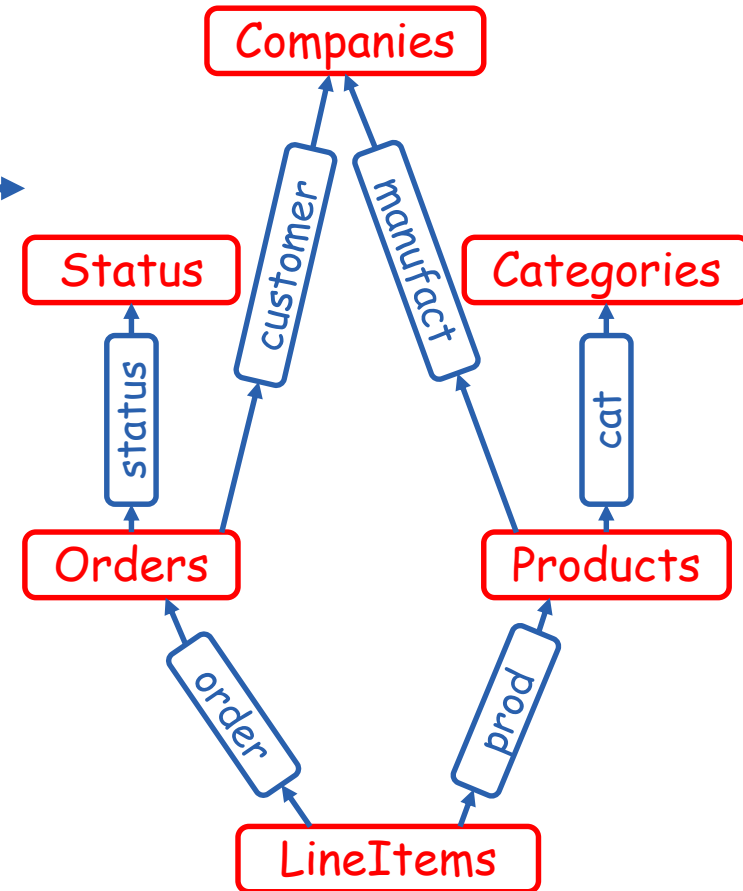    - ▶ Partially ordered set
    - ▶ Functional approach
  - ◆ Concept-oriented expression language (COEL)
    - ▶ No joins, no group-bys, no formal logic
    - ▶ Simple and expressive analytical operations
    - ▶ Algebra of functions
  - ◆ Column-based data processing model
    - ▶ Fast analytical operations with data (analytical database)
    - ▶ Column is a function
- ● More info: http://conceptoriented.org

# SETS

- Goal: define a new set in terms of existing sets and functions
- Two operations
  - Product: all combinations of greater sets
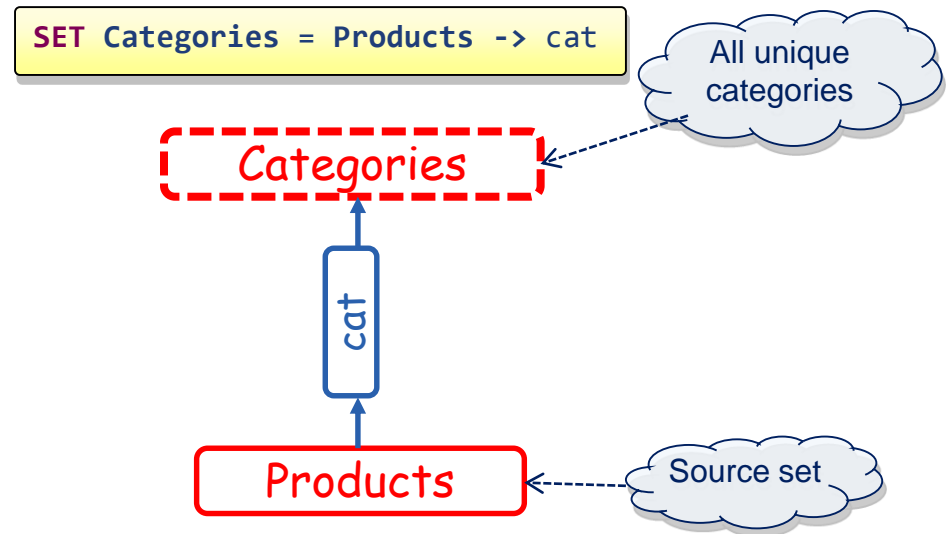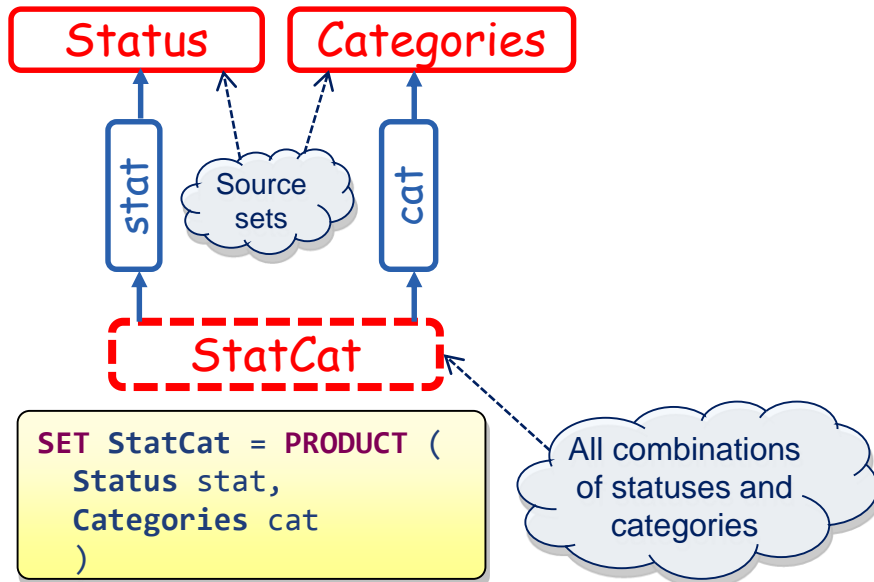  - Project: all outputs of some function

**Extraction dialog**

| | | |
|---|---|---|
| Table Name: | Categories | |
| Column Name: | Category | |
| Source Table: | Products | |

Column mapping:

| Select | Source | Type |
|---|---|---|
| ☐ | List Price | Double |
| ☑ | Reorder Level | Integer |
| ☐ | Target Level | Integer |
| ☑ | Category | String |

OK   Cancel



```
SET StatCat = PRODUCT (
    Status stat,
    Categories cat
)
```

All combinations of statuses and categories

```
SET Categories = Products -> cat
```

All unique categories

Source set

# LINKS

- ## Goal: link to sets using existing functions

# AGGREGATION

- **Parameters:**
  - ◆ Fact set
  - ◆ Grouping function
  - ◆ Measure function
  - ◆ Aggregation function



```
Double TotalAmount = AGGREGATE (
    LineItems,
    prod.cat,
    amount,
    SUM
    )
```

# CONCLUSION

- **Novelties:**
  - ◆ <u>Unified data model</u> and expression language are used
  - ◆ <u>Column formulas</u> as opposed to cell formulas for derived data
- **Advantages of ConceptMix (Data Commander):**
  - ◆ <u>Ease of use</u>: radically simplifies analytical data integration; kills complexities when manipulating data
  - ◆ <u>Fast time-to-value</u>: from months to minutes
  - ◆ <u>Lower IT costs</u>: move the burden of authoring BI contents to the end users
  - ◆ Increase <u>motivation</u>; more convenient consumption of data
- **Future work:**
  - ◆ <u>Assistance engine</u>: recommending mappings, relationships, sources
  - ◆ <u>Selection propagation</u> and inference for interactive analysis
- **More info:** <u>http://conceptoriented.org</u>