## Data Analysis and ETL Process

Document author: arho.virkki@tyks.fi Contributors: anna.hammais@tyks.fi, juhana.valo@medbit.fi, katja.kanerva-leppanen@medbit.fi
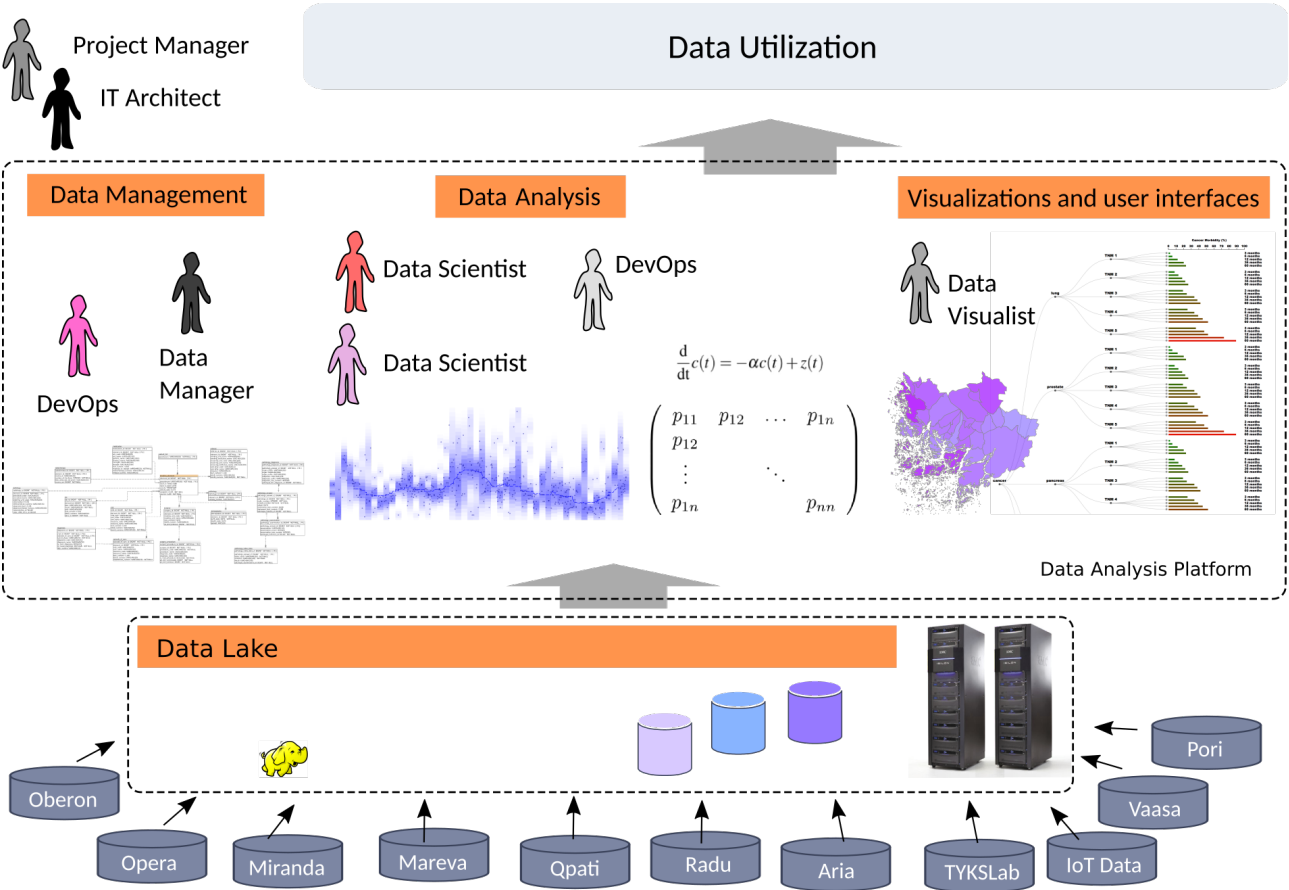
### Data Enrichment and Analysis Steps



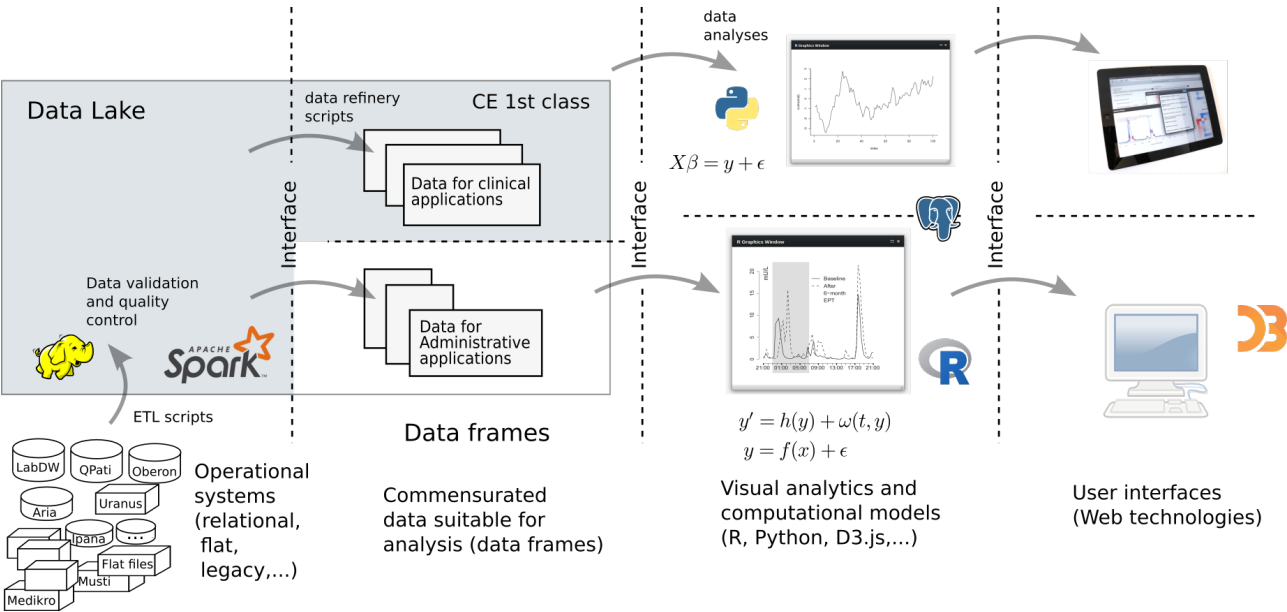**Figure [PDF, SVG].** Process Overview.

## Data Flow



**Figure [PDF, SVG].** Data Analysis Workflow.

## Research Process Steps
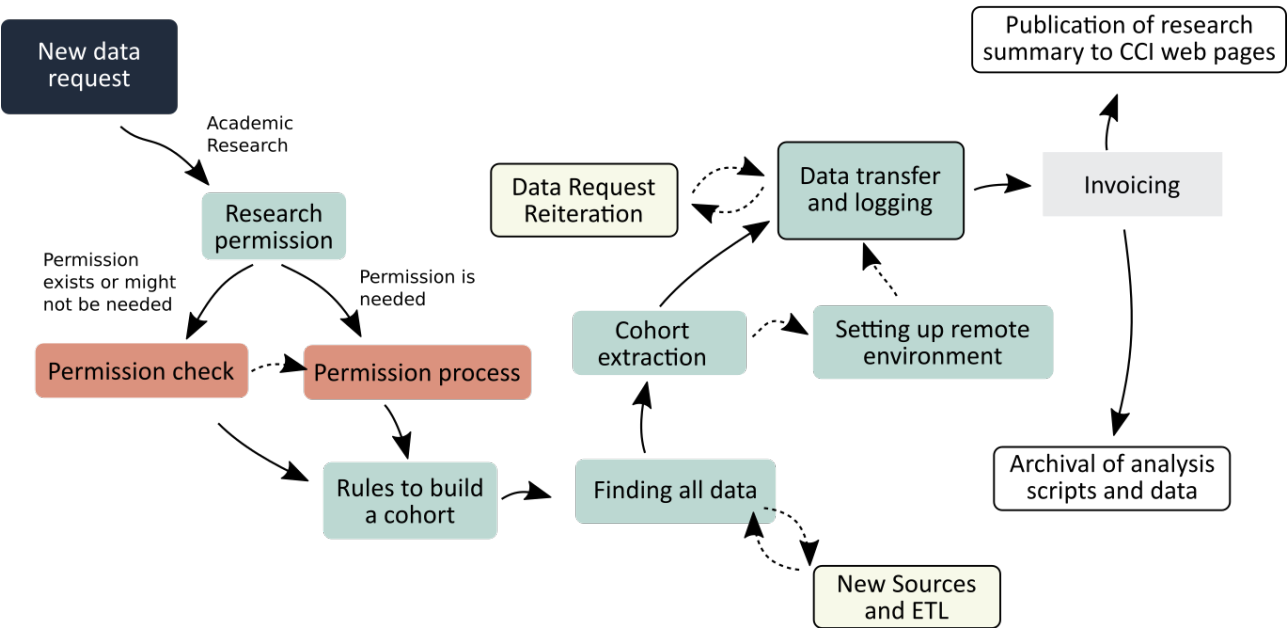
### Academic Research Process



**Figure [PDF, SVG].** Academic research process.
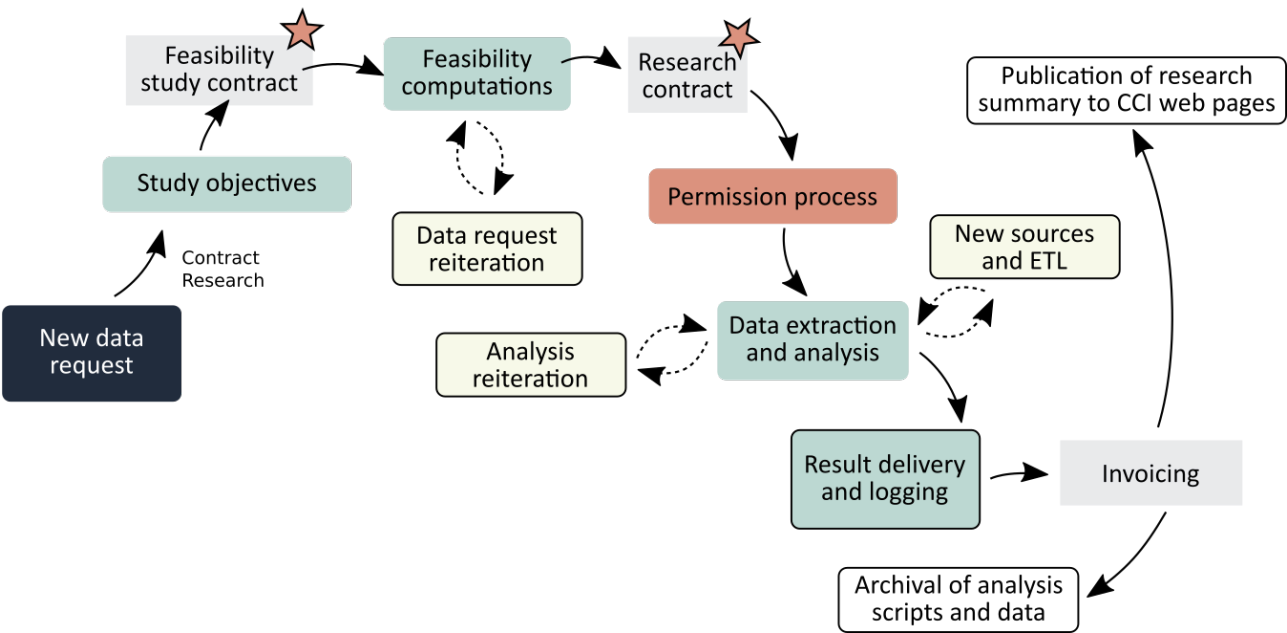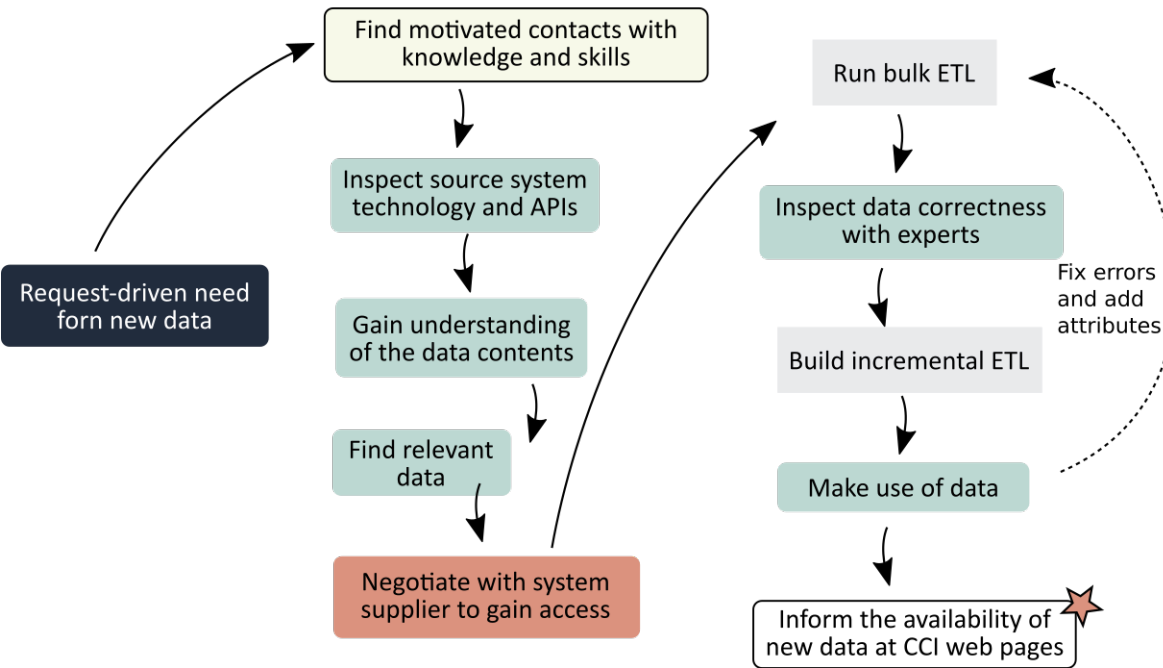
**Contract Research Process**



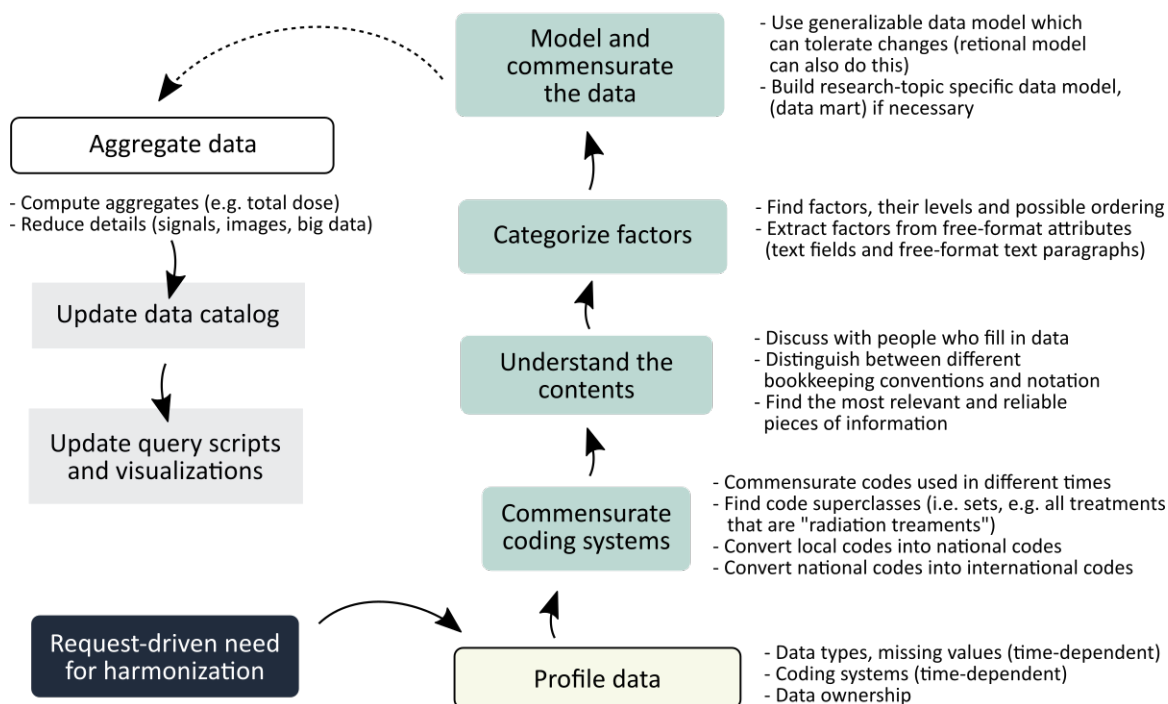**Figure [PDF, SVG].** Contract research process.

## New Data Sources and Data Harmonization

**New Data Source Process**



**[PDF, SVG].** New data source process.**

**Data harmonization Process**



**[PDF, SVG].** Data harmonization process.

## ETL Steps

1. Data extraction from source
2. File upload
3. Format conversion
4. Type conversion
5. Data integration
6. Semantic unification

## ETL Script Repository

The ETL scripts are saved to the Git-repository `ktp@ktpgit:/opt/git/ETL.git`. Two working copies are used in production:

```
ktphadoop.vsshp.net:/var/lib/hive/ETL/
ktpanalytics.vsshp.net:/opt/ktp/ETL/
```

The first is for pre-processing the flat text files at the Hadoop machine, and the second is used with Pentaho Kettle in the data integration phase. The automated etl scripts are run by the *ktp* user. For details, see *crontab -e* as *ktp* at *ktpanalytics*.

## Details

Hadoop Environment
PostgreSQL Setup