# *REST-based Data Integration Services for Software Engineering Domain*

Fridolin Koch, Bachelor's Thesis – Kickoff Presentation

Software Engineering for Business Information Systems (sebis)
Department of Informatics
Technische Universität München, Germany

wwwmatthes.in.tum.de

# Outline

1. **Motivation**
   - Problem statement
   - Existing ETL solutions
2. **Research Questions**
3. **Solution Approach**
   - UI Prototype
   - Framework Workflow
   - Technology Stack
   - Current Architecture
4. **Next Steps**
5. **Timeline**

# Problem Statement

- Existing barrier in the adoption of knowledge management systems in software engineering domain

  - Many different software architecture life cycle tools produce data in different formats (Enterprise Architect, Excel, Jira, etc.)

  - Repeatedly integrating this data into such a system can be a challenging and tedious task

- In general the task of data integration is addressed by **E**xtract-**T**ransform-**L**oad-Tool (ETL-Tool)

  - Wide range of commercial and open source ETL-Tool available

  - But: Mostly tailored to generic use cases → Difficult to embedded in existing domain specific tools

- **Potential Solution**: Analyze popular ETL-Tools and create an easily extendable framework

# Existing ETL-Tools: Overview

| Tool | OpenSource | Technology | Mode | Domain |
|------|-----------|-----------|------|--------|
| Apatar | Yes | Java | Standalone | Generic |
| CloverETL | Core only | Java | Standalone, Embedded | Generic |
| Talend Open Studio for Data Integration | Yes | Java | Designer / Script-Generator | Generic |
| Pentaho | Yes, but less functionality | Java | Standalone | Generic |
| RhinoETL | Yes | C# .net | Framework | Generic |
| UnifiedViews | Yes | Java | Standalone | Linked-Data (RDF) |

**sebis**


apatar
connecting data


Flexible Deployment Options

Desktop Application · Server Engine · Embedded

No constraints around deployment to build scalable solutions.


No Coding! Visual Job Designer

Distinct · Transform · Validate · Connectors · TextFile · Salesforce.com · Oracle · Ldap · MsSQL · MySQL

Salesforce.com · Join · Transform · Oracle · MySQL

Graphical tools enable non-developers to connect applications on the spot.


No Coding! Transformation Mapper

Custom Table.ID · Lookup · CONTACT_ID

Custom Table.FIRSTNAME · Append · FNAME

Visual mapping tool to link data, and to create and modify complex transfomations.

- Java-Based
- Open-Source
- Visual job designer
- Generic usage domain

Source: www.apatar.com

CloverETL

- Java-Based ETL-Tool
- Open-Source (Core only)
- Visual job designer (Community ad Commercial Edition)
- Standalone and embedded
- Generic usage domain
- Custom Domain-Specific-Language to define business logic ("CTL")
- Clusterable
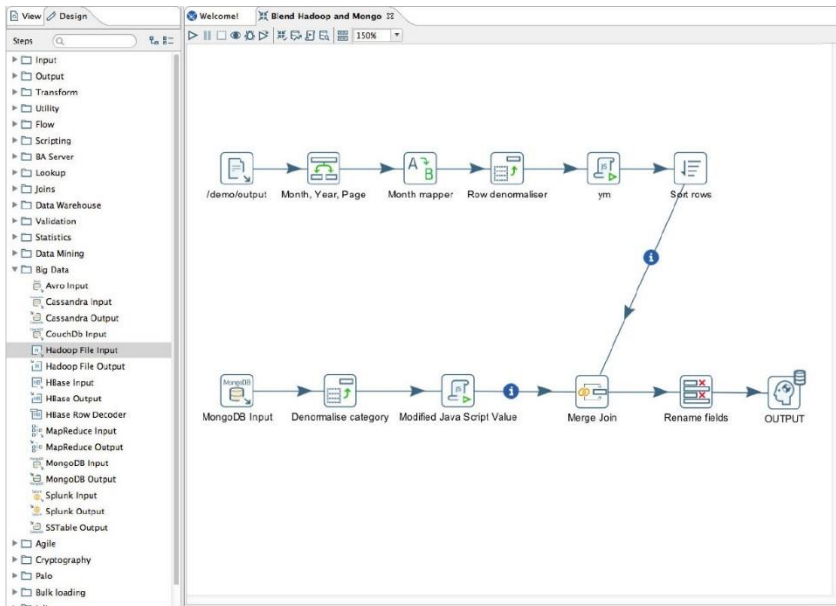- Many data connectors

Source: http://www.cloveretl.com/

- Java-Based ETL-Tool
- Open-Source
- Code generator for data transformation scripts (Java)
- Based on Eclipse
- +900 connector and components
- Generic usage domain

Source: https://www.talend.com/download/talend-open-studio

# Existing ETL-Tools: Pentaho

sebis



- Java-Based
- Community and Enterprise edition
- Visual Designer
- Rich library of pre-built ETL components
- Generic usage domain

Source: http://www.pentaho.com/product/data-integration

# Existing ETL-Tools: Rhino ETL

- C# .Net Framework
- Open source
- Hello-World application available on GitHub
- Pure framework no additional connectors

Source: https://hibernatingrhinos.com/oss/rhino-etl

**sebis**

- Java based
- Open Source
- Specialized on RDF-Data (Linked data)
- Visual Designer to build Job (Web-Based)
- Extendable through plugins
- Developed at *Charles University, Prague*

Source: http://unifiedviews.eu/

sebis

- Almost all tool have a generic use case domain, but are manly advertised for Business Intelligence and Big Data Integration / Analysis
    - Tools have thousands of adapters, transformers and settings → High entry barrier
    - Heavy duty tools for "Big Data" → Higher configuration and maintenance effort
- SyncPipes is lightweight an quick to integrate into your infrastructure
    - TypeScript / JavaScript provides an ecosystem that is easily extensible
    - ~260.000 Packages available through *npm* to speed up development
    - RESTful API assures easy integration into existing system architecture
    - Rule of thumb: Create new adapters and the corresponding Workflow within a day
    - Docker-Support out of the Box ("Zero configuration")

**sebis**

**Research Objective:** Create a REST-based Data Integration Framework to enable developers to implement adapters for ETL-Workflows easily. Facilitate the End-User to visualize the source and target system's domain model in the conjunction with creating new Data Integration Workflows.

**Research Questions**

*Q1*: "What are the key features that must be supported by data integration framework?"

*Q2*: "How does the framework's architecture support its extensibility with new adapters?"

**sebis**

# SyncPipes

**①** ——————— ② ——————— ③ ——————— ④

Start

---

**≡✓** **Select existing**
Edit an existing integration workflow

Select workflow ▼

NEXT

---

**✎** **Create**
Create a new integration workflow

Workflow name

**Select input connector**
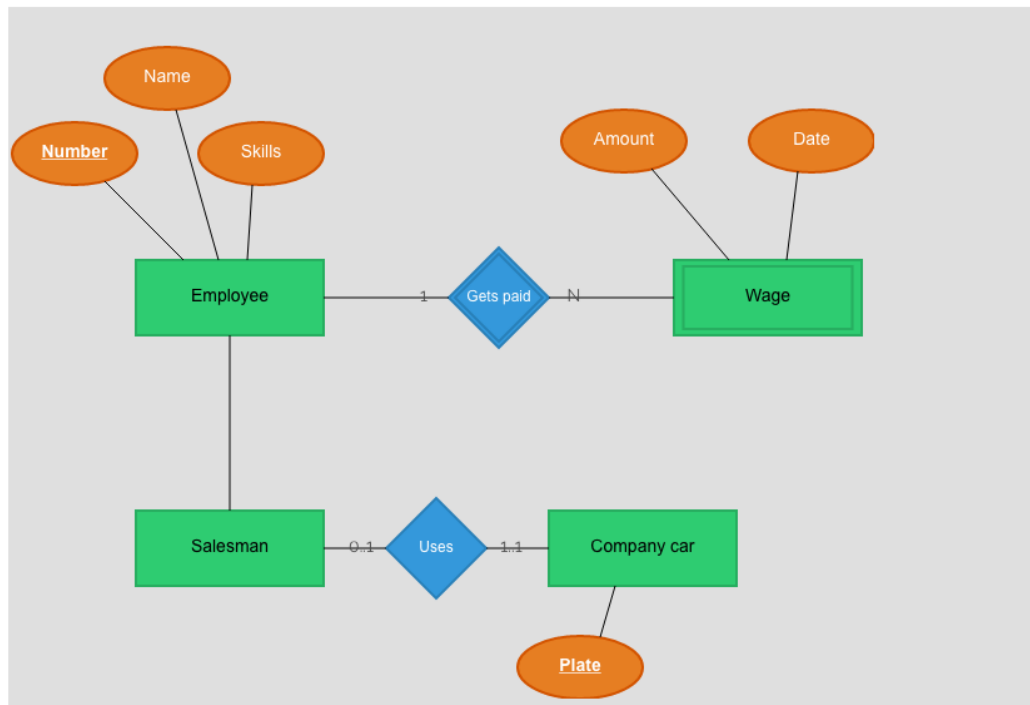
Select connector ▼

**Select output connector**

# UI Prototype (II)

Mapping

Entity–relationship model view

Click on attributes to see values

**SOURCE SYSTEM**    TARGET SYSTEM

Name

Number    Skills

Amount    Date

Employee — 1 — Gets paid — N — Wage

Salesman — 0.1 — Uses — 1.1 — Company car

Plate

Click an attribute to view the values

Date

Sat May 03 1997 14:32:36 GMT+0000 (UTC)
Tue Jun 12 2001 06:02:40 GMT+0000 (UTC)
Sun Jan 26 1975 12:22:34 GMT+0000 (UTC)
Wed Oct 08 2008 03:45:25 GMT+0000 (UTC)
Fri Aug 02 1974 15:04:37 GMT+0000 (UTC)
Thu Mar 08 2012 13:17:46 GMT+0000 (UTC)
Mon Nov 19 1984 09:56:47 GMT+0000 (UTC)
Thu Apr 25 1991 20:12:19 GMT+0000 (UTC)
Tue Nov 03 1992 18:06:46 GMT+0000 (UTC)
Mon Mar 09 1987 05:17:47 GMT+0000 (UTC)
Fri Jan 04 1980 03:21:28 GMT+0000 (UTC)
Mon Aug 01 1988 19:11:29 GMT+0000 (UTC)
Fri Sep 04 2009 22:02:47 GMT+0000 (UTC)
Sun Feb 22 1970 08:23:15 GMT+0000 (UTC)
Wed Dec 09 2009 22:10:16 GMT+0000 (UTC)

**sebis**

## Mapping

Map source entities & attributes to your target system

Entites ⊕

employee ▼

   Attributes ⊕

number ▼

skills ▼

Entites ⊕

salesman ▼

   Attributes ⊕

name ▼

name ▼

NEXT

# UI Prototype (IV)

① ——————— ② ——————— ③ ——————— ④

## Schedule

## Scheduling

🔘 Schedule Workflow to run periodically

| MINUTE | HOUR | DAY OF THE MONTH | MONTH | DAY OF WEEK |

Repeat

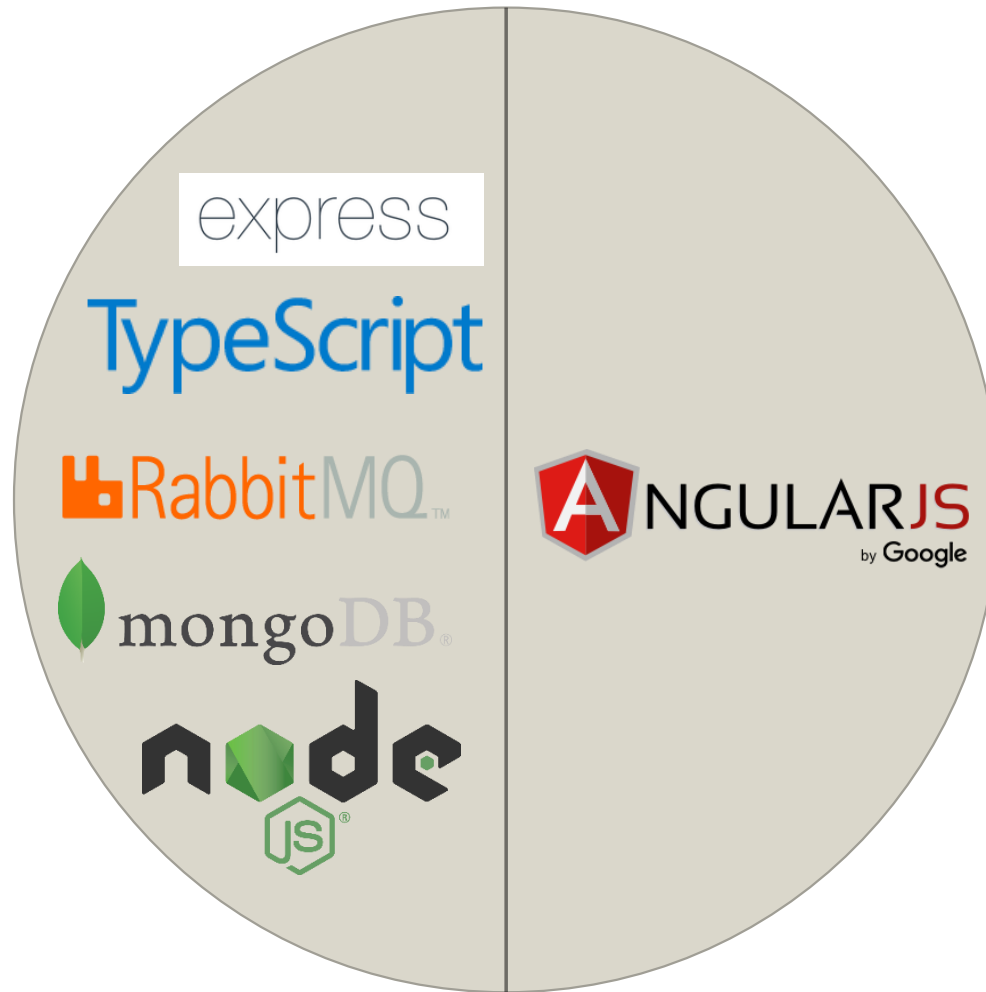every minute ▼

NEXT

# UI Prototype (V)

sebis

① — ② — ③ — ④

Execute

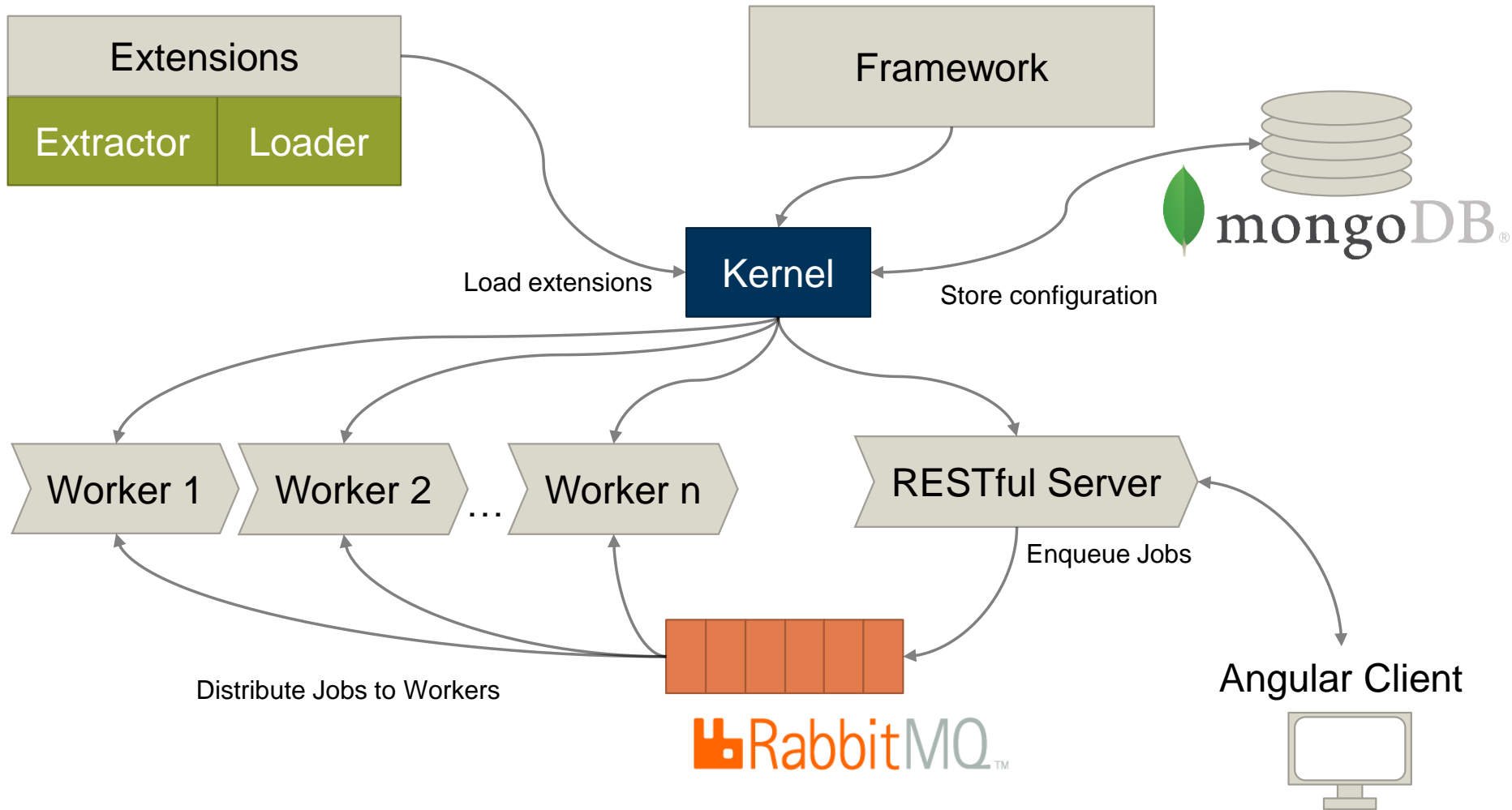## Execute your Workflow
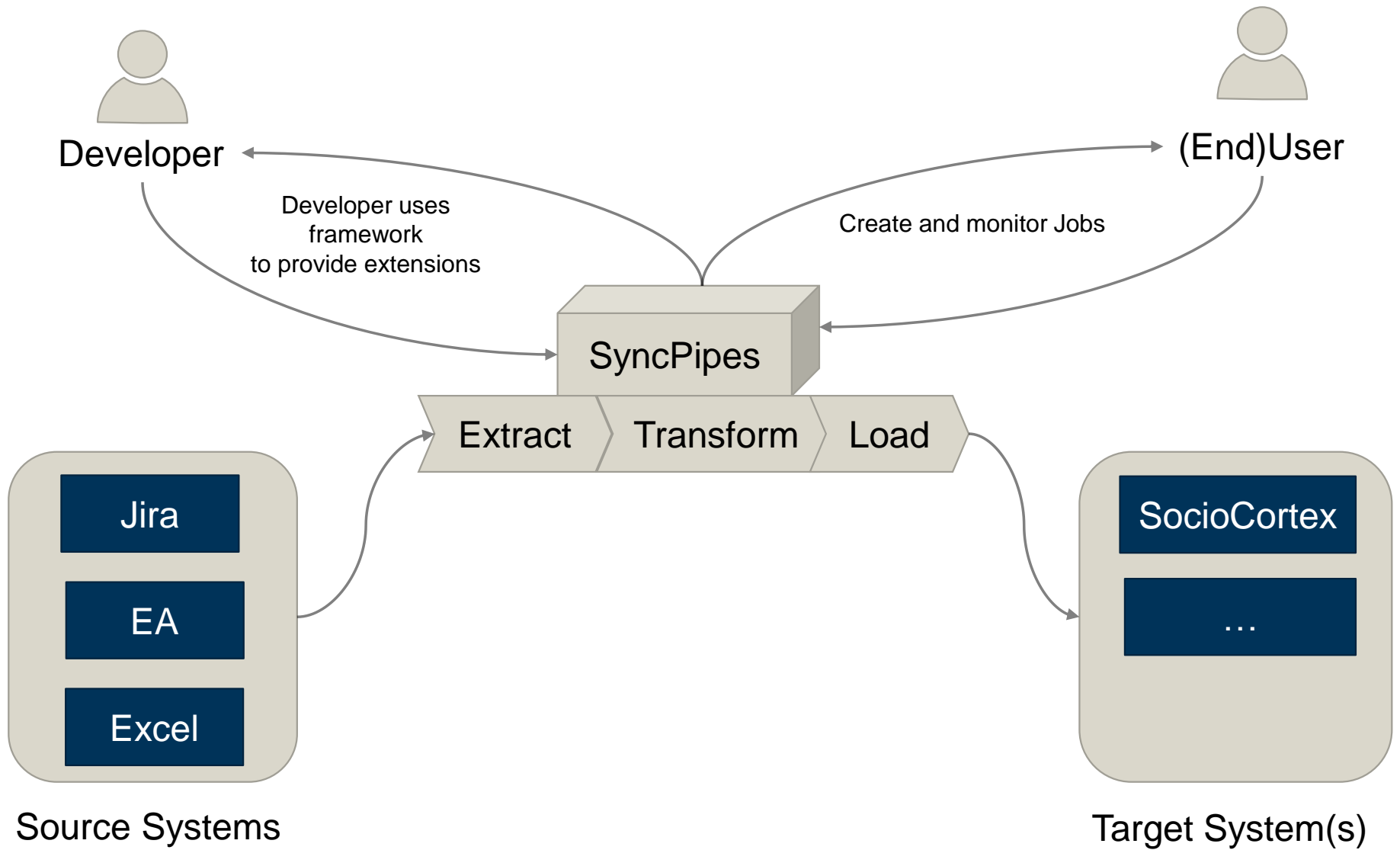
TEST WORKFLOW EXECUTION    SAVE WORKFLOW

```
Copying dataset 0 from Jira to SocioCortex
Copying dataset 1 from Jira to SocioCortex
Copying dataset 2 from Jira to SocioCortex
Copying dataset 3 from Jira to SocioCortex
Copying dataset 4 from Jira to SocioCortex
Copying dataset 5 from Jira to SocioCortex
Copying dataset 6 from Jira to SocioCortex
Copying dataset 7 from Jira to SocioCortex
Copying dataset 8 from Jira to SocioCortex
Copying dataset 9 from Jira to SocioCortex
Copying dataset 10 from Jira to SocioCortex
Copying dataset 11 from Jira to SocioCortex
Copying dataset 12 from Jira to SocioCortex
Copying dataset 13 from Jira to SocioCortex
Copying dataset 14 from Jira to SocioCortex
Copying dataset 15 from Jira to SocioCortex
Copying dataset 16 from Jira to SocioCortex
Copying dataset 17 from Jira to SocioCortex
Copying dataset 18 from Jira to SocioCortex
```

sebis



**M**ongoDB **E**xpress.js **A**ngular.js **N**ode.js

sebis

**Extensions**

| Extractor | Loader |
|-----------|--------|

**Framework**

Load extensions

**Kernel**

Store configuration

mongoDB

Worker 1    Worker 2    ...    Worker n

**RESTful Server**

Enqueue Jobs

Distribute Jobs to Workers

RabbitMQ™

**Angular Client**

**sebis**

Developer

(End)User

Developer uses
framework
to provide extensions

Create and monitor Jobs

SyncPipes

Extract  Transform  Load

Jira

EA

Excel

SocioCortex

…

Source Systems

Target System(s)

- Prototype evaluation
  1. Present prototypical implementation to 2-3 developers which a familiar with the target domain (e.g. researchers at the SEBIS chair)
  2. Ask developers to implement extractor and loader extensions
  3. Gather feedback through interviews
  4. Improve prototype based on the provided feedback
  5. Ask developers to implement similar adapters again
  6. Gather feedback
- Improve UI / Frontend to work with RESTful backend
- Write thesis

**sebis**



| February | March | April | May | June | July |
|----------|-------|-------|-----|------|------|

Literature research

Analyze ETL-Tools

UI Prototype

UI + REST

Framework Implementation

Evaluation

Improve Framework

Evaluation

Writing + Buffer

# Thank you for your attention.

**Fridolin Koch**

sebis

Technische Universität München
Department of Informatics
Chair of Software Engineering for
Business Information Systems

Boltzmannstraße 3
85748 Garching bei München

frido.koch@tum.de
wwwmatthes.in.tum.de