

Logical Approach in Software and Data Design

软件和数据设计中的逻辑方法

Evgeny Cherkashin

Matrosov Institute for System Dynamics and Control Theory
of Siberian Branch of Russian Academy of Sciences, Irkutsk,
Russia

俄罗斯伊尔库茨克, Matrosov
俄罗斯科学院西伯利亚分院系统动力学与控制理论研究所
eugeneai@icc.ru

2024, March, Shandong, China

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

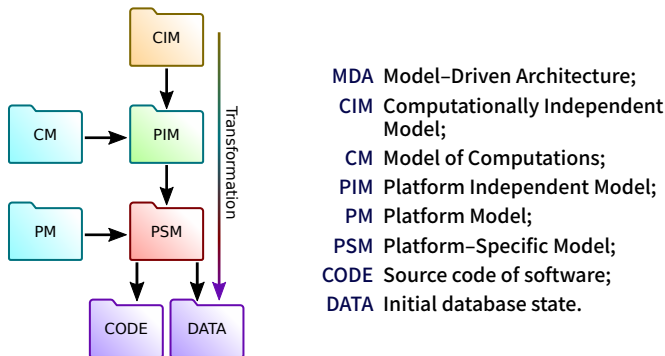
Model-Driven Architecture: Research objectives

Main objective of the research is to construct a MDA technology based on nowadays system modeling visual languages (SysML, UML, BPMN, CMMN) and existing Semantic Web vocabularies and technologies. The following techniques and software are under development:

1. CIM representation with SysML, BPMN, CMMN, and results of source code processing,
2. CIM, PIM, PSM representation in UML, RDF with existing vocabularies,
3. transformation implementation with logical language Logtalk,
4. usage of LOD sources in transformations for obtaining additional semantic data,
5. generation of documents and user interfaces with LOD markup.

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

Model-Driven Architecture



E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

Logtalk as transformation definition language

We have chosen Logtalk as it

- inherits widely known Prolog language syntax and runtime;
- is implemented as macro package, performance penalties are about 1.5%;
- has flexible semantics: we can define transformations and constraints within the same syntax;
- implement object-oriented knowledge (rules) structuring, encapsulation and replacement;
- compositional way of transformation implementation;
- powerful engine to post constraints on object-to-object messages (events);
- has implementation for various Prolog engines.

The «regular» language allow us to use its libraries not directly related to MDA transformations.

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

QR-Code of the presentation



E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

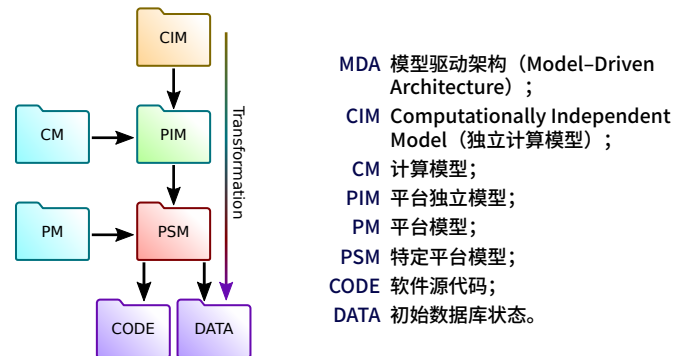
模型驱动架构：研究目标

研究的主要目标是基于当今的系统建模可视化语言（SysML、UML、BPMN、CMMN）和现有的语义网词汇表和技术构建一种 MDA 技术 以下技术和软件正在开发中

1. 用 SysML、BPMN、CMMN 和源代码处理结果表示 CIM、
2. 用 UML、RDF 和现有词汇表表示 CIM、PIM 和 PSM、
3. 使用逻辑语言 Logtalk 实现转换、
4. 在转换中使用 LOD 源以获取额外的语义数据、
5. 使用 LOD 标记生成文档和用户界面。

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

模型驱动架构（Model-Driven Architecture）



E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

Logtalk 作为转换定义语言

我们选择Logtalk是因为它

- 继承了广为人知的 Prolog 语言的语法和运行时；
- 作为宏包实现，性能损失约为1.5
- 具有灵活的语义：我们可以在相同的语法中定义转换和约束；
- 实现了面向对象的知识（规则）结构化、封装和替换；
- 实现转换的组合方式；
- Logtalk是一个强大的引擎，可以对对象到对象的消息（事件）发布约束；
- 为各种Prolog引擎提供实现。

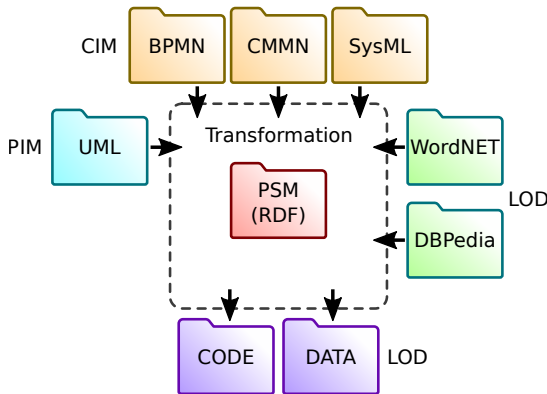
这种常规语言允许我们使用它与 MDA 转换没有直接关系的库。

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

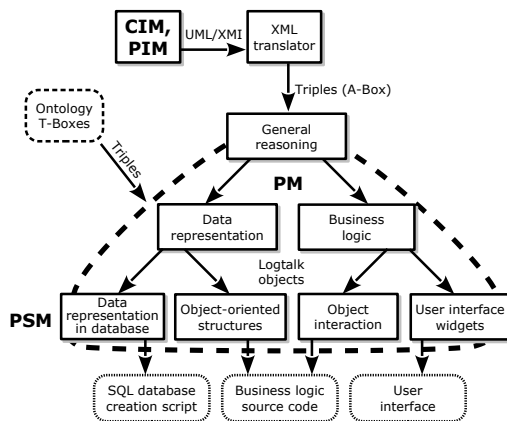
1. Information is published in Internet with open access license;
 2. It is represented in a machine-readable form, e.g., Excel table instead of a bitmap picture;
 3. An open format used, e.g., CSV instead of Excel;
 4. The format is based on W3C recommended standards, allowing RDF and SPARQL reference;
 5. Published data refer to objects, forming context.
- Thus, applications publish data as relations of objects (entities).

Model Driven Architecture and Linked Open Data

模型驱动架构和关联开放数据



Architecture of transformation modules, 转换模块的结构



Implementation of Query object, 实现查询对象

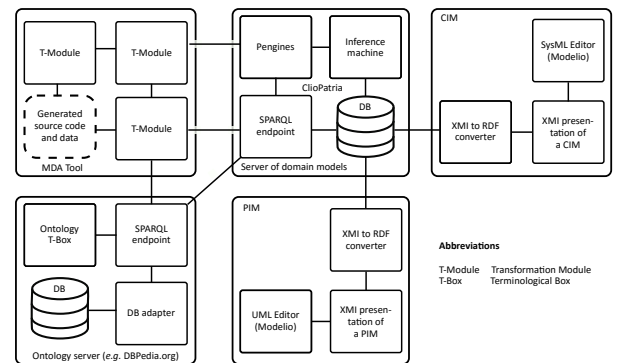
```

:- object(query(_XMI)).
:- protected(xmi/1).
:- public([class/2, attribute/3, method/3]).
xmi(XMI) :- parameter(1, XMI).
class(Name, ID):-
    ::xmi(XMI),
    XMI::rdf(ID, rdf:type, uml:'Class'),
    XMI::rdf(ID, rdfs:label, literal(Name)).
attribute(Name, ClassID, ID):-
    ::xmi(XMI),
    XMI::rdf(ClassID, xmi:ownedAttribute, ID),
    XMI::rdf(ID, rdfs:label, literal(Name)).
method(Name, ClassID, ID):-
    ::xmi(XMI),
    XMI::rdf(ClassID, xmi:ownedOperation, ID),
    XMI::rdf(ID, rdfs:label, literal(Name)).
% .....
:- end_object.

```

1. 信息在互联网上发布，采用开放式获取许可；
 2. 它以机器可读的形式表示，例如 Excel 表格，而不是位图图片；
 3. 使用的开放格式，如 CSV 而非 Excel；
 4. 该格式基于 W3C 推荐的标准，允许引用 RDF 和 SPARQL；
 5. 发布的数据指代对象，形成上下文。
- 因此，应用程序以对象（实体）关系的形式发布数据。

MDA infrastructure, 基础设施



PSM: Scenario of a Class synthesis, PSM: 课堂综合情景

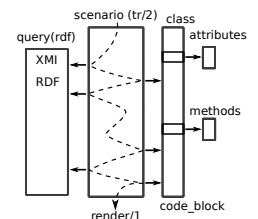
```

:- object(direct(_Package, _LocalProf, _CodeProf)). % 转换驱动程序对象
:- public([tr/4, tr/3]). % 类综合方案的公共接口
% .....
tr(class, Class, ClassID):-package(Package), % 合成一个类
query(Package):class(Name, ClassID), % XMI 中的查询包结构
create_object(Class, % ..... % 创建类对象
create_object(Attributes, % ..... % 创建属性对象
create_object(Methods, % ..... % 创建方法对象
Class:name(Name), % ..... % 为类级命名。
% Generate attributes of the class,
% organizing them in a local database.
% .....methods...
Class:attributes(Attributes), % 为类设置属性。
Class:methods(Methods), % .....方法。

tr(tr(attribute, Attribute, ClassID, AttributeID):- % 属性转换
package(Package),
query(Package):attribute(Name, ClassID, AttrID),
create_object(Attribute, % ..... % 属性命名。
Attribute:name(Name), % .....方法名称。

tr(tr(method, Method, ClassID, MethodID):- % 方法的转变
package(Package),
query(Package):method(Name, ClassID, MethodID),
create_object(Method, % ..... % 方法名称。
Method:name(Name), % .....方法名称。
:- end_object.

```

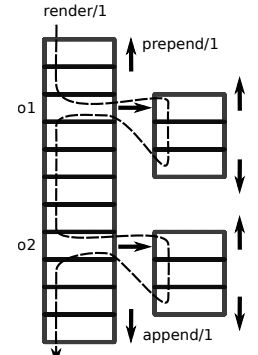


Code Block (idea is taken from llvmlite*, 代码块，创意来自“llvmlite”图书馆)

```

:- object(code_block, specializes(root)).
% Public interface of the object
:- public([append/1, prepend/1, clear/0,
render/1, render_to/1, remove/1,
item/1, items/1]).
% Code block items
:- dynamic([item/_1]).
:- private([item/_1]).
% Methods specialized during inheritance
:- protected([renderitem/2, render_to/2]).
% .....
% Delegate rendering to object itself
renderitem(Object, String):-
current_object(Object), !,
Object::render(String).
% Convert a literal to its string
% representation
renderitem(literal(Item), String):-!,
atom_string(Item, String).
% Just print the item (debugging).
renderitem(Item, String):-
root::iswritef(String, '%q', [Item]).
:- end_object.

```



*) <https://github.com/numba/llvmlite>

E. Cherkashin, et al. Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

E. Cherkashin, et al. Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

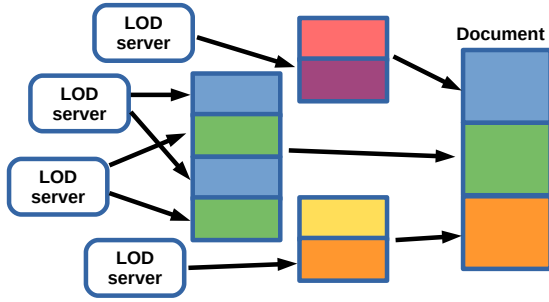
E. Cherkashin, et al. Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

E. Cherkashin et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

E. Cherkashin et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

Structure of a document, 文件结构



E. Cherkashin, et al. Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

Representation, 代表性

```
<html lang=" ru" xmlns=http://www.w3.org/1999/xhtml
xmlns:taa=http://irnok.net/engine/rdfa-manipulation
xml:lang=" ru" metal:define-macro=" page" >
<head> ... </head>
<body prefix=" rdf: http://www.w3.org/1999/...-ns# foaf: http://xmlns.com/foaf/...
imei: imei.html# course: https://irnok.net/college/plan/01.16-... \
%D0%BA_PB-SM.plm.xml.xls-...2.3.1.html# " resource=" #post"
typeof=" schema:CreativeWork sioc:Post prov:Entity" >
<!-- The application control panel -->

<main lang=" ru" resource=" #annotation" typeof=" oa:Annotation" id=" main-doc-cnt" >
<div property=" oa:hasTarget" resource=" #course-work-prog" ></div>
<article property=" oa:hasBody" typeof=" foaf:Document curr:WorkingProgram"
resource=" #course-work-program" id=" main-document" >
<div taa:content=" imei:title-page" ></div>
<div taa:content=" imei:neg-UMK" ></div>
<section id=" TOC" class=" break-after" ><h2>Table of Contents</h2>
<div id=" tableOfContents" ></div>
</section>
<section id=" course-description" resource=" #description"
property=" schema:hasPart" typeof=" schema:CreativeWork" >
<div property=" schema:hasPart" resource=" #purpose"
typeof=" dc:Text cnt:ContentAsText" >
<div property=" cnt:chars" datatype=" xsd:string" >
<h2 property=" dc:title" datatype=" xsd:string" >
Aims and objectives of the discipline (module)</h2>
<p>The aim of teaching the discipline ...</p>
</div>
```

E. Cherkashin, et al. Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

Generated list of title page preambles, 生成扉页序言列表



E. Cherkashin, et al. Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

Imported time distribution for lecture, seminars, ..., 导入讲座和研讨会的时间分配、

загрузка...

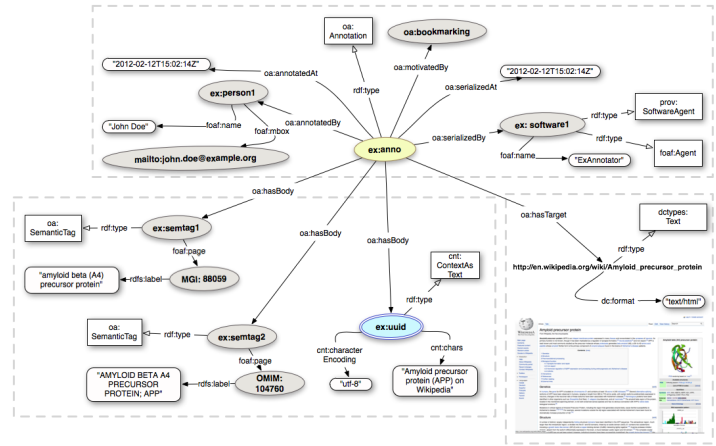
- методиками экстремального и agile-программирования.

4. Объем дисциплины (модуля) и виды учебной работы (разделяется по формам обучения)

Вид учебной работы	Всего часов / зачетных единиц	Семестры	
		3	4
Аудиторные занятия (всего)	108	33	75
в том числе:			
Лекции	36		36
Практические занятия (ПЗ)			
Семинары (С)			
Лабораторные работы (ЛР)	66	30	36
КСР	6	3	3
Самостоятельная работа обучающихся (СР-ОБ)	45	20	25

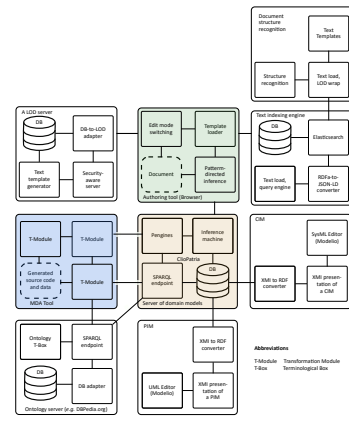
E. Cherkashin, et al. Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

Open Annotation ontology (oa), 开放注释本体论



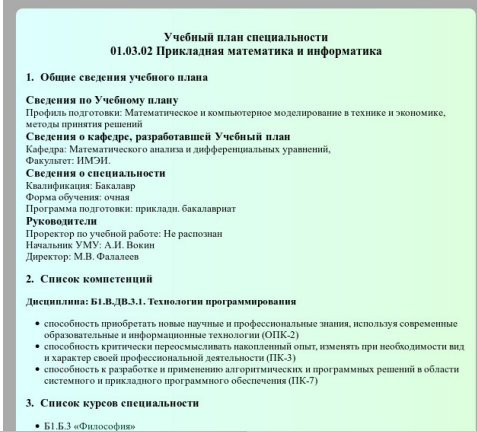
E. Cherkashin, et al. Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

Architecture, 建筑学



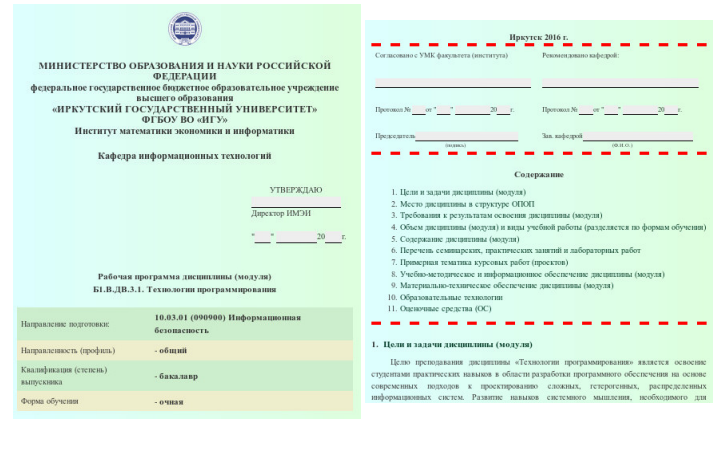
E. Cherkashin, et al. Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

Generated part of a course description, 生成课程说明的一部分



E. Cherkashin, et al. Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

Complete document, 完整文件



E. Cherkashin, et al. Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

- Friend-of-a-friend (foaf) - agent information: individuals, legal entities, program agents.
- Provenance (prov) - references between documents.
- Dublin Core (dc) - edited annotation mark up.
- DBPedia resource (dbr) – references to instant objects and classes.
- Schema.org (schema) - Google, Yandex, Yahoo, etc. searchable objects, structural elements.
- The Bibliographic Ontology (bibo) - literature reference mark up.

Conclusion

A tools (components) for digital archive implementation, which allows to device information systems and document processing services with the following features:

- load LOD marked up document, extract, store in a graph and index RDF data;
- retrieve RDF data as triples or as a result of full-text search query;
- combine existing LOD data and its content in new documents dynamically with browser based context inference machine;
- use server-site inference machine (Prolog) to process RDF data upon request from browser's part of the system;
- convert created RDFa marked up HTML5 documents into Excel and Word formats.

Applications

- Document authoring automation;
- Context-depended editing;
- Self-organizing global document flows;
- Documents as data sources for information systems.

TabbyXL

Software Platform for Rule-Based Spreadsheet Data Extraction and Transformation

基于规则的电子表格数据提取和转换软件平台

Alexey Shigarov, Vasiliy Khristyuk, et al.

shigarov@icc.ru

动机

- 关于任意电子表格表格
 - ▶ 为科学和商业应用提供大量宝贵数据
 - ▶ 丰富多样的布局、风格和内容功能
 - ▶ 以人为本（结构不正确，内容杂乱无章）
 - ▶ 没有明确的语义供计算机解释
- 挑战
 - ▶ 如何从工作表中提取表格
 - ▶ 如何识别和纠正细胞结构异常
 - ▶ 如何恢复自动解释所需的语义
 - ▶ 如何使用外部词汇表将提取的数据概念化

- Friend-of-a-friend (foaf): 代理人信息：个人、法人实体、计划代理人
- Provenance (prov): 文件之间的引用
- Dublin Core (dc): 编辑的注释标记
- DBPedia resource (dbr): 即时对象和类的引用
- Schema.org (schema): Google、Yandex、Yahoo 等搜索对象、结构元素
- The Bibliographic Ontology (bibo): 文献参考标注

结论

实施数字档案的工具（组件），可以 设备信息系统和文件处理服务具有以下功能:

- 加载 LOD 标注文件、提取、存储在图表中并为 RDF 数据编制索引
- 以三元组或全文搜索查询结果的形式检索 RDF 数据
- 利用基于浏览器的上下文推理机，将现有 LOD 数据及其内容动态结合到新文件中
- 使用服务器端推理机（Prolog）处理 RDF 数据。处理 RDF 数据
- 将创建的 RDFa 标记 HTML5 文档转换为 Excel 和 Word 格式。

Applications

- 文件编写自动化
- 根据上下文进行编辑
- 自组织全球文件流
- 文件作为信息系统的数据源

Motivation

- About arbitrary spreadsheet tables
 - ▶ A large volume of valuable data for science and business applications
 - ▶ A big variety of layout, style, and content features
 - ▶ Human-centeredness (incorrect structure and messy content)
 - ▶ No explicit semantics for interpretation by computers
- Challenges
 - ▶ How to extract tables from worksheets
 - ▶ How to recognize and correct cell structure anomalies
 - ▶ How to recover semantics needed for the automatic interpretation
 - ▶ How to conceptualize extracted data by using external vocabularies

Background

Table understanding includes the following tasks

1. Extraction, detecting a table and recognizing the physical structure of its cells
2. Role analysis, extracting functional data items from cell content
3. Structural analysis, recovering internal relationships between extracted functional data items
4. Interpretation, linking extracted functional data items with external vocabularies (general-purpose or domain-specific ontologies)

背景介绍

表格理解 包括以下任务

1. 提取: 检测表格并识别其单元格的物理结构
2. 角色分析: 从单元格内容中提取功能数据项
3. 结构分析: 恢复提取的功能数据项之间的内部关系
4. 口译: 将提取的功能数据项与外部词汇表 (通用或特定领域本体) 连接起来

E. Cherkashin, et al.

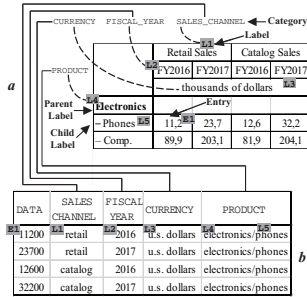
Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

捐款

TabbyXL 是一个软件平台, 旨在开发和执行基于规则的程序, 用于电子表格数据提取和从任意表 (a) 到关系表 (b) 的转换。

Novelty

- ❑ 为数据项而非单元格分配角色的表格对象模型
- ❑ CRL, 特定领域语言, 用于表达用户定义的表格分析和解释规则
- ❑ CRL 到 Java 翻译器, 用于合成电子表格数据转换的可执行程序



E. Cherkashin, et al.

Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

表格对象模型

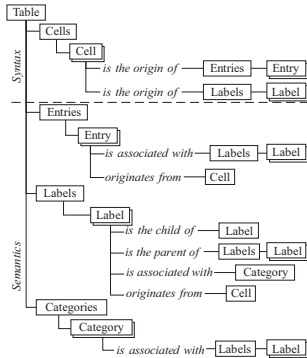
物理层

以布局、样式和内容特征的单元格

逻辑层

功能数据项及其关系:

- ❑ 条目 (数值)
- ❑ 标签 (键)
- ❑ 类别 (概念)
- ❑ 条目标签对
- ❑ 标签-标签对
- ❑ 标签-类别对



E. Cherkashin, et al.

Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

Cell Cleansing

The actions correct an inaccurate layout and content of a hand-coded table

- ❑ **<merge>** combines two adjacent cells when they share one border
- ❑ **<split>** divides a merged cell that spans n -tiles (row-column intersections) into n -cells
- ❑ **<set text>** modifies a textual content of a cell
- ❑ **<set indent>** modifies a text indentation of a cell

Example

```
when
cell corner: cl == 1, rt == 1, blank
cell c: cl > corner.cr, rt > corner.rb
then
split c
```

E. Cherkashin, et al.

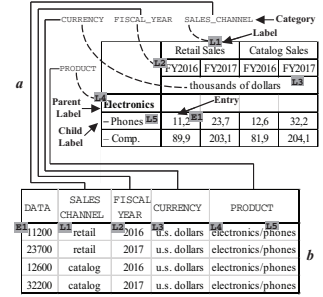
Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

Contribution

TabbyXL is a software platform aiming at the development and execution of rule-based programs for spreadsheet data extraction and transformation from arbitrary (a) to relational tables (b)

Novelty

- ❑ Table object model assigning roles to data items, not cell
- ❑ CRL, domain-specific language to express user-defined rules for table analysis and interpretation
- ❑ CRL-to-Java translator to synthesize executable programs for spreadsheet data transformation



E. Cherkashin, et al.

Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

Table Object Model

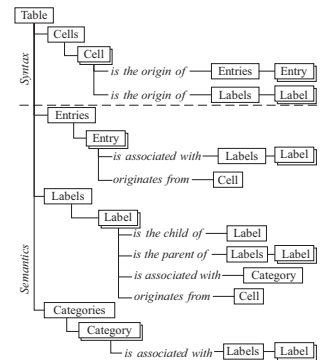
Physical Layer

Cells characterized by layout, style, and content features

Logical Layer

Functional data items and their relationships:

- ❑ entries (values)
- ❑ labels (keys)
- ❑ categories (concepts)
- ❑ entry-label pairs
- ❑ label-label pairs
- ❑ label-category pairs



E. Cherkashin, et al.

Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

CRL Grammar

```
rule = 'rule' <a Java integer literal> 'when' condition
      'then' action 'end' <EOL> (rule) <EOF>
condition = query identifier ['<' constraint '>' constraint]
           ['<' assignment ['<' assignment '>'] <EOL> (condition)]
constraint = <a Java boolean expr>
assignment = identifier ':' <a valid Java expr>
query = 'cell' | 'entry' | 'label' | 'category' | 'no cells' |
        'no entries' | 'no labels' | 'no categories'
action = merge | split | set text | set indent | set mark |
         new entry | new label | add label | set parent |
         set category | group <EOL> (action)
merge = 'merge' identifier 'with' identifier
split = 'split' identifier
set text = 'set text' <a Java string expr> 'to' identifier
set indent = 'set indent' <a Java integer expr> 'to' identifier
set mark = 'set mark' <a Java string expr> 'to' identifier
new entry = 'new entry' identifier ['<' <a Java string expr> '>']
new label = 'new label' identifier ['<' <a Java string expr> '>']
add label = 'add label' identifier ['<' <a Java string expr> '>']
           'of' identifier | <a Java string expr>
           'to' identifier
set parent = 'set parent' identifier 'to' identifier
set category = 'set category' identifier | <a Java string expr>
              'to' identifier
group = 'group' identifier 'with' identifier
identifier = <a Java identifier>
```

E. Cherkashin, et al.

Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

细胞清洁

这些操作可纠正手工编码表格中不准确的布局和内容

- ❑ **<merge>** 当相邻两个单元格共享一个边界时, 将其合并
- ❑ **<split>** 将跨 n 格 (行列交叉点) 的合并单元格划分为 n 单元格
- ❑ **<set text>** 修改单元格的文本内容
- ❑ **<set indent>** 修改单元格的文本缩进

Example

```
when
cell corner: cl == 1, rt == 1, blank
cell c: cl > corner.cr, rt > corner.rb
then
split c
```

E. Cherkashin, et al.

Logical Approach in Software and Data Design软件和数据设计中的逻辑方法

Role Analysis

The actions recover entries and labels as functional data items presented in a table

- ❑ **<set mark>** annotates a cell with a user-defined tag that can be used in subsequent table analysis
- ❑ **<new entry>** (**<new label>**) creates an entry (label) from a cell content with the use of an optional string processing

Example

```
when
cell corner: cl == 1, rt == 1, blank
cell c: cl > corner.cr, rt > corner.rb
then
new entry c
```

E. Cherkashin, et al.

Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

Structural Analysis

The actions recover pairs of two kinds: entry-label and label-label

- ❑ **<add label>** associates an entry with a label
- ❑ **<set parent>** binds two labels as a parent and its child

Example

```
when
cell c1: cl == 1
cell c2: cl == 1, rt > c1.rt, indent == c1.indent + 2
no cells: cl == 1, rt > $c1.rt, rt < $c2.rt, indent == $c1.indent
then
set parent c1.label to c2.label
```

E. Cherkashin, et al.

Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

Interpretation

The actions serve to recover label-category pairs

- ❑ **<set category>** associates a label with a category
- ❑ **<group>** places two labels to one group that can be considered as an undefined category

Example

```
when
label l1: cell.mark == " stub"
label l2: cell.mark == " stub" , cell.rt == l1.cell.rt
then
group l1 with l2
```

E. Cherkashin, et al.

Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

Illustrative Example

The transformation of arbitrary tables with the same layout features (a and c) to their canonical versions (b and d)

a1	a2		DATA	A	B	
b1	1	b4	4	1	a1	b1
b2	2	b5	NA	2	a1	b2
b3		b6	6	4	a2	b4
				6	a2	b6

a

a1	a2	a3	DATA	A	B		
b1	b3	1	b5	1	2	a1	b2
b2	2	b4	NA	b6	4	a2	b4
				6	a2	b6	

c

DATA	A	B
1	a1	b1
2	a1	b2
4	a2	b4
6	a2	b6

b

DATA	A	B
2	a1	b2
5	a2	b5
6	a2	b6

d

The ruleset for the cell cleansing (a), role analysis (b, c), structural analysis (d, e), and interpretation (f, g)

```
a when cell c: c.text.matches("NA")
then set text "" to c
c when cell c: (cl % 2) == 1
then new label c
when
entry e
label l1: cell.cr == e.cell.cr
then add label l to e
e label l1: cell.rt == e.cell.rt, cell.cl == e.cell.cl - 1
then add label l to e
f when label l1: cell.rt == 1
then set category "A" to 1
g when label l1: cell.rt > 1
then set category "B" to 1
b when cell c: (cl % 2) == 0, !blank
then new entry c
when
entry e
label l1: cell.cr == e.cell.cr
then add label l to e
d when label l1: cell.cr == e.cell.cr
then add label l to e
e label l1: cell.rt == e.cell.rt, cell.cl == e.cell.cl - 1
then add label l to e
f when label l1: cell.rt == 1
then set category "A" to 1
g when label l1: cell.rt > 1
then set category "B" to 1
```

This example is reproducible at <https://codeocean.com/capsule/5326436>

E. Cherkashin, et al.

Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

角色分析

操作恢复作为功能数据项的条目和标签，以表格形式呈现

- ❑ **<set mark>**为单元格标注用户定义的标签，该标签可用于后续表格分析
- ❑ **<new entry>** (**<new label>**)使用可选的字符串处理，从单元格内容创建条目（标签）。

Example

```
when
cell corner: cl == 1, rt == 1, blank
cell c: cl > corner.cr, rt > corner.rb
then
new entry c
```

E. Cherkashin, et al.

Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

结构分析

恢复的操作对有两种：条目-标签和标签-标签

- ❑ **<add label>** 将条目与标签关联
- ❑ **<set parent>** 将两个标签绑定为父标签和子标签

Example

```
when
cell c1: cl == 1
cell c2: cl == 1, rt > c1.rt, indent == c1.indent + 2
no cells: cl == 1, rt > $c1.rt, rt < $c2.rt, indent == $c1.indent
then
set parent c1.label to c2.label
```

E. Cherkashin, et al.

Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

口译

这些操作有助于恢复标签-类别对

- ❑ **<set category>** 将标签与类别关联
- ❑ **<group>** 将两个标签贴在一个可视为未定义类别的组上

Example

```
when
label l1: cell.mark == " stub"
label l2: cell.mark == " stub" , cell.rt == l1.cell.rt
then
group l1 with l2
```

E. Cherkashin, et al.

Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

示例

将具有相同布局特征的任意表格 ((a) 和 (c)) 转换为其规范版本 ((b) 和 (d))

a1	a2	DATA	A	B
b1	1	b4	4	1
b2	2	b5	NA	2
b3		b6	6	4

a

a1	a2	a3	DATA	A	B
b1	b3	1	b5	1	2
b2	2	b4	NA	b6	4
				6	a2

c

DATA	A	B
1	a1	b1
2	a1	b2
4	a2	b4
6	a2	b6

b

DATA	A	B
2	a1	b2
3	a2	b3
5	a3	b5
6	a3	b6

d

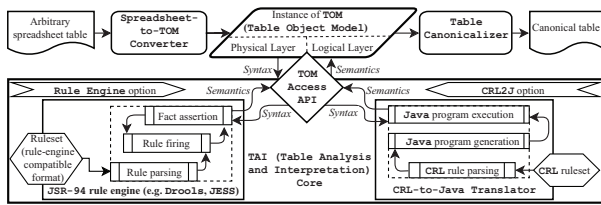
细胞清洗的规则集 (a), 角色分析 (b, c), 结构分析 (d, e), 和解释 (f, g)

```
a when cell c: c.text.matches("NA")
then set text "" to c
c when cell c: (cl % 2) == 1
then new label c
when
entry e
label l1: cell.cr == e.cell.cr
then add label l to e
e label l1: cell.rt == e.cell.rt, cell.cl == e.cell.cl - 1
then add label l to e
f when label l1: cell.rt == 1
then set category "A" to 1
g when label l1: cell.rt > 1
then set category "B" to 1
b when cell c: (cl % 2) == 0, !blank
then new entry c
when
entry e
label l1: cell.cr == e.cell.cr
then add label l to e
d when label l1: cell.cr == e.cell.cr
then add label l to e
e label l1: cell.rt == e.cell.rt, cell.cl == e.cell.cl - 1
then add label l to e
f when label l1: cell.rt == 1
then set category "A" to 1
g when label l1: cell.rt > 1
then set category "B" to 1
```

该示例可在以下网址复制 <https://codeocean.com/capsule/5326436>

E. Cherkashin, et al.

Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法



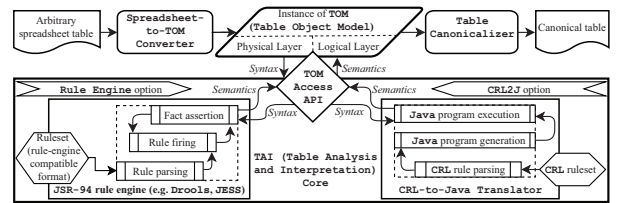
Two options are provided

Rule Engine option

Executing a ruleset in an appropriate format with a JSR-94 compatible rule engine (e.g. Drools, Jess)

CRL2J option

Translating a ruleset expressed in CRL to an executable Java program



提供两种选择

规则引擎选项

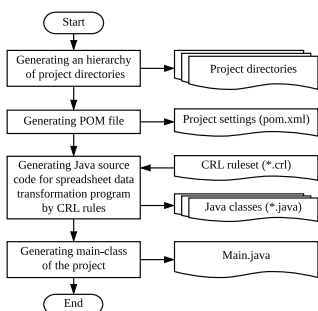
使用与 JSR-94 兼容的规则引擎（如 Drools、Jess）以适当格式执行规则集

CRL2J 选择权

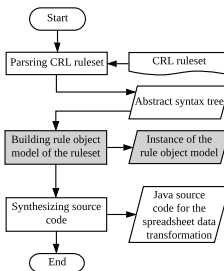
将 CRL 表达的规则集转换为可执行 Java 程序

CRL2J Translation

Workflow for generating a Maven-project of a spreadsheet data transformation program

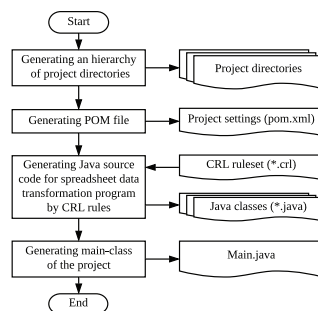


Workflow for translating a CRL ruleset to Java source code

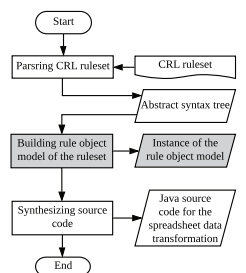


CRL2J 翻译

生成电子表格数据转换程序 Maven 项目的工作流程



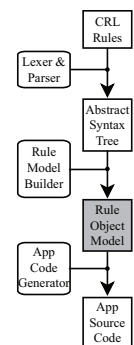
将 CRL 规则集转换为 Java 源代码的工作流程



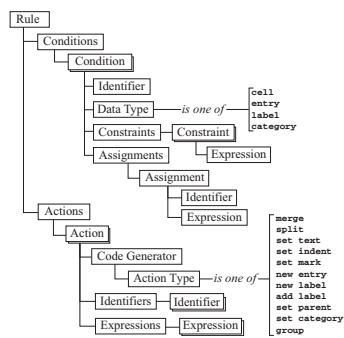
CRL2J Translation

CRL2J 翻译

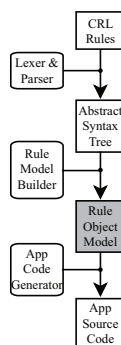
In the Workflow



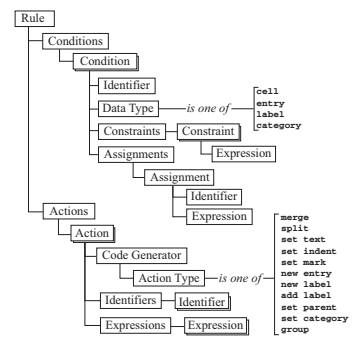
Rule Object Model



在工作流程中



规则对象模型



CRL2J Translation

CRL2J 翻译

Example (Source Rule)

```
when
  cell corner: cl == 1, rt == 1, blank
  cell c: cl > corner.cr, rt > corner.rb, ! marked
then
  set mark " @entry" to c
  new entry c
```

Example (Fragment of the Generated Java Code)

```
...
Iterator<CCell> iterator1 = getTable().getCells();
while (iterator1.hasNext()) {
  corner = iterator1.next();
  if ((corner.getCl() == 1) && (corner.getRt() == 1) && ...
    iterator<CCell> iterator2 = getTable().getCells();
    while (iterator2.hasNext()) {
  ...
```

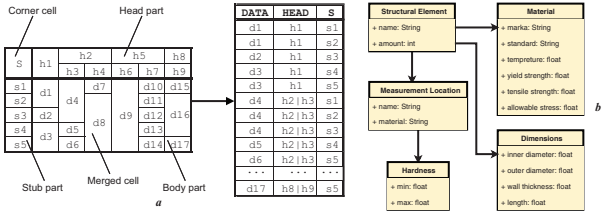
Example (来源规则)

```
when
  cell corner: cl == 1, rt == 1, blank
  cell c: cl > corner.cr, rt > corner.rb, ! marked
then
  set mark " @entry" to c
  new entry c
```

Example (生成的 Java 代码片段)

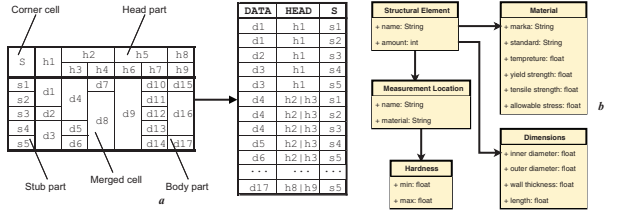
```
...
Iterator<CCell> iterator1 = getTable().getCells();
while (iterator1.hasNext()) {
  corner = iterator1.next();
  if ((corner.getCl() == 1) && (corner.getRt() == 1) && ...
    iterator<CCell> iterator2 = getTable().getCells();
    while (iterator2.hasNext()) {
  ...
```


Generating conceptual models — (b) from arbitrary tables presented in industrial safety inspection reports — (a)



The more detail can be found at <https://github.com/tabbydoc/tabbyxl/wiki/industrial-safety-inspection>

生成概念模型 — (b) 来自工业安全检查报告中的任意表格 — (a)



更多详情访问 <https://github.com/tabbydoc/tabbyxl/wiki/industrial-safety-inspection>

Conclusions & Further Work

- Impact on software development for spreadsheet data management
 - Table object model associating functional roles with data items
 - Table analysis and interpretation driven by user-defined rules
 - Formulated actions to recover missing semantics of arbitrary tables
 - Translation of rules to executable spreadsheet transformation programs
- Limitations
 - The inaccurate cell structure prevents the table analysis
 - The very limited interpretation (without external vocabularies)
- Further work
 - Rearrangement of cell structure by using visual (human-readable) cells
 - Detecting derived data by spreadsheet formulas
 - Enriching the table analysis by named entity recognition
 - Linking extracted data items with LOD cloud

结论和下一步工作

- 对电子表格数据管理软件开发的影响
 - 将功能角色与数据项关联起来的表对象模型
 - 根据用户定义的规则进行表格分析和解释
 - 为恢复任意表的缺失语义而制定的行动
 - 将规则转化为可执行的电子表格转换程序
- 局限性
 - 不准确的单元格结构妨碍了表格分析
 - 非常有限的解释（没有外部词汇表）
- 进一步的工作
 - 利用可视（人类可读）细胞重新排列细胞结构
 - 通过电子表格公式检测派生数据
 - 通过命名实体识别丰富表格分析
 - 将提取的数据项与 LOD 云连接起来

Thanks!

Read more about the project at <http://td.icc.ru>

The project source code is available at <https://github.com/tabbydoc/tabbyxl>

But it is not all ...

谢谢!

有关该项目的更多信息，请访问 <http://td.icc.ru>

项目源代码见 <https://github.com/tabbydoc/tabbyxl>

但这并不是全部 ...

Domain Knowledge Graphs Induction from Tables

Tables are the most available sources of information. They are valuable data sources for Knowledge Bases (KB)

Knowledge Base Construction Populating with document and structured table extracted data

Knowledge Base Population Populating with recognized new facts on entities from big text corpses

Knowledge base Augmentation Populating with relations with table data.

- (Ré, 2014) Ré C., et al. Feature engineering for knowledge base construction. IEEE Data Eng. Bull., 37, 26–40, (2014).
- (Balog, 2018) Balog K. Populating knowledge bases. Entity-Oriented Search. INRE, 39, 189–222, (2018).
- (Zhang & Balog, 2020) Zhang S. & Balog K. Web table extraction, retrieval, and augmentation: A survey. ACM Trans. Intell. Syst. Technol., 11, (2020).

领域知识图谱 表格归纳法

表格是最常用的信息来源。它们是知识库（KB）的重要数据源

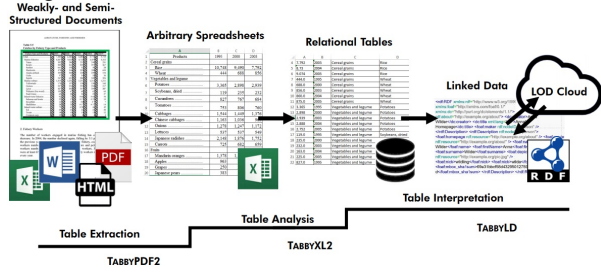
知识库建设 填充文件和结构化表格提取的数据

知识库人口 从大文本尸体中填充关于实体的公认新事实

知识库扩充 用表格数据填充关系

- (Ré, 2014) Ré C., et al. Feature engineering for knowledge base construction. IEEE Data Eng. Bull., 37, 26–40, (2014).
- (Balog, 2018) Balog K. Populating knowledge bases. Entity-Oriented Search. INRE, 39, 189–222, (2018).
- (Zhang & Balog, 2020) Zhang S. & Balog K. Web table extraction, retrieval, and augmentation: A survey. ACM Trans. Intell. Syst. Technol., 11, (2020).

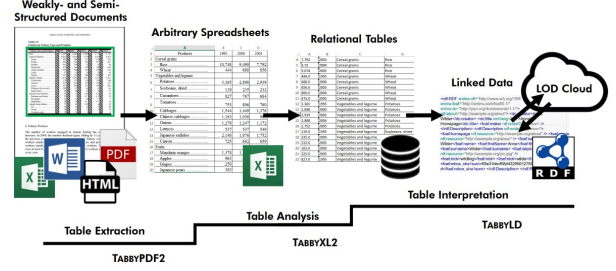
There three main stages of Automatic table interpretation (Shigarov, 2017)



- (Shigarov, 2017) Shigarov A., Mikhailov A. Rule-based spreadsheet data transformation from arbitrary to relational tables. Information Systems, 71, 123-136 (2017).

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

自动表格解释有三个主要阶段 (Shigarov, 2017)



- (Shigarov, 2017) Shigarov A., Mikhailov A. Rule-based spreadsheet data transformation from arbitrary to relational tables. Information Systems, 71, 123-136 (2017).

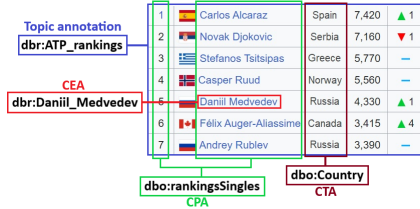
E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

Semantic Table Interpretation

Semantic Table Interpretation

Semantic interpretation (Annotation) of tables (Semantic Table Interpretation, STI) is a recognition of mutual and external relations between elements of table content. External relations relate to an enterprise KG and/or a global KG (e.g. DBpedia.org).

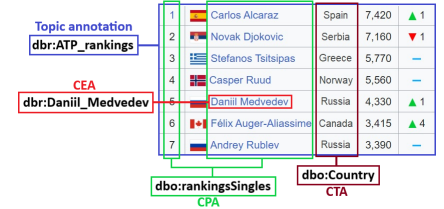
- Cell-Entity Annotation (CEA)
- Column-Type Annotation (CTA)
- Column Property Annotation (CPA)
- Topic Annotation



E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

表格的语义解释 (注释) (Semantic Table Interpretation, STI) 是对表格内容元素之间相互关系和外部关系的一种确认。外部关系涉及企业 KG 和/或全局 KG (例如 DBpedia.org)。

- 细胞实体注释 (CEA)
- 列式注释 (CTA)
- 列属性注释 (CPA)
- 主题注释



E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

Cell-Entity Annotation

细胞实体注释

CEA comprises the sequential steps as follows:

- Select a candidate entity set from DBpedia.org for each value of a cell via SPARQL endpoint and DBpedia lookup.
- Disambiguation

A SPARQL-query matching words of a phrase.

```
SELECT DISTINCT (str(?subject) as ?subject)
WHERE {
  ?subject a ?type.
  ?subject rdfs:label ?label.
  ?label <bif:contains> ".*%value1*." AND ".*%value2*." ...
  FILTER NOT EXISTS { ?subject dbo:wikiPageRedirects ?r2 }.
  FILTER (istrstarts(str(?subject), " http://dbpedia.org/resource/Category:" ) ).
  FILTER (istrstarts(str(?subject), " http://dbpedia.org/property/" ) ).
  FILTER (istrstarts(str(?subject), " http://dbpedia.org/ontology/" ) ).
  FILTER (istrstarts(str(?type), " http://dbpedia.org/ontology/" ) ).
  FILTER (lang(?label) = "en" )
}
ORDER BY ASC(strlen(?label))
LIMIT 100
```

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

CEA 包括以下三个连续步骤:

- 通过 SPARQL 端点和 DBpedia 查找, 为单元格的每个值从 DBpedia.org 中选择一个候选实体集。
- 消歧义

匹配短语单词的 SPARQL 查询。

```
SELECT DISTINCT (str(?subject) as ?subject)
WHERE {
  ?subject a ?type.
  ?subject rdfs:label ?label.
  ?label <bif:contains> ".*%value1*." AND ".*%value2*." ...
  FILTER NOT EXISTS { ?subject dbo:wikiPageRedirects ?r2 }.
  FILTER (istrstarts(str(?subject), " http://dbpedia.org/resource/Category:" ) ).
  FILTER (istrstarts(str(?subject), " http://dbpedia.org/property/" ) ).
  FILTER (istrstarts(str(?subject), " http://dbpedia.org/ontology/" ) ).
  FILTER (istrstarts(str(?type), " http://dbpedia.org/ontology/" ) ).
  FILTER (lang(?label) = "en" )
}
ORDER BY ASC(strlen(?label))
LIMIT 100
```

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

Evaluation on Test Table Sets

Evaluation on Test Table Sets

A well-known precision measurement (accuracy) is used for assessment

$$\text{Accuracy} = \frac{CC}{NC},$$

where CC is the number of the correctly related columns to a categorical entity, and CN is the total number of columns.

Recognition stage	T2Dv2	Tough_ Tables	Git- Tables
Stage 2, Atomic column classification	0.994	0.956	0.938
Stage 3, Column entity identification	0.924	-	-

Comparison with analogs

	TAIPAN	Table-Miner+	T2Dv2	Mantis-Table
Column entity identification	0.540	0.871	0.924	0.979

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

采用众所周知的精确测量 (精度) 进行评估

$$\text{准确性} = \frac{CC}{NC},$$

其中, CC 是与分类实体正确相关的列数, CN 是列的总数。

认可阶段	T2Dv2	Tough_ Tables	Git- Tables
第 2 阶段, 原子柱分类	0.994	0.956	0.938
第 3 阶段, 列实体识别	0.924	-	-

与类似物的比较

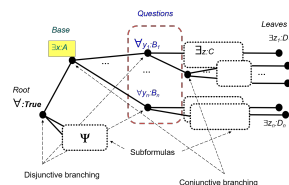
	TAIPAN	Table-Miner+	T2Dv2	Mantis-Table
列实体标识	0.540	0.871	0.924	0.979

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

Knowledge Representation and Reasoning: the PCF-Calculus

The main properties of the language of positively constructed formulas (PCF) and its calculi:

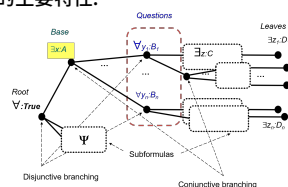
- PCFs have a large-block structure (tree-like) and consist of only positive quantifiers \exists and \forall
- the PCF-based calculus have a unique inference rule
 - the proof in the PCF-calculus is organized as a question-answering procedure
 - PCF-calculus is both machine-oriented and human-oriented; it is compatible with heuristics
 - the semantic of the PCF-calculus can be changed without modifying axioms and the inference rule



知识表示与推理：PCF 微积分

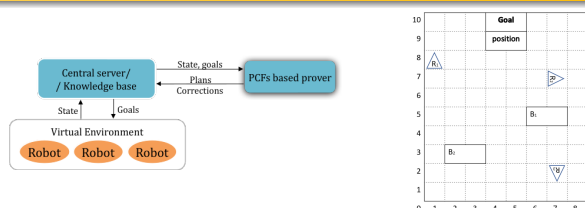
正构造公式（PCF）语言及其计算器的主要特性:

- PCF 具有大块结构（树状），仅由正量词 \exists 和 \forall 组成。
- 基于 PCF 的微积分有一个唯一的推理规则
 - PCF 微积分中的证明是以问题解答过程的形式组织的
 - PCF 微积分既面向机器，也面向人类；它与启发式方法兼容
 - 可以在不修改公理和推理规则的情况下改变 PCF 微积分的语义



E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

PCF-Based Method for Problem Solving



- The goal of the team of robot is to transport blocks to the target area
- Each block can be dragged by two or more robots
- The current state of the World and the goal of the group are formalized in PCF
- The PCF-based prover and a selection mechanism produce the optimal joint plan of actions for the team
- The current plan can be easily modified whenever the state of the World is changed

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

A Master Degree Program. Semantic Technologies and Multiagent Systems

It is a joint effort of Saint-Petersburg Electrotechnical University (LETI), Irkutsk State University, and ISDCT SB RAS. Main subjects.

- Computation Geometry, Digital Signal Processing, Internet of Things,
- Semantic web, Semantic web Information System Development,
- AI Basics, Knowledge representation, Object-oriented Logic Programming,
- Answer Set Programming (SAT), Natural Language Processing,
- Machine Learning, Neural Networks, Deep Learning,
- Multiagent Systems, Optimization with Multiagent Systems.

Started at 2022-09-01.

<https://etu.ru/sveden/education/programs/semanticheskie-tehnologii-i-mnogoagentnye-sistemy-01.04.02.html>

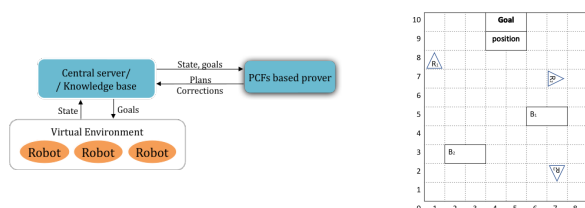
E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

Conclusion (the final one)

- Classic knowledge-based systems are powerful AI tools for solving wide class of recognition problems and synthesis of various kind: source code, data objects, control
- Contemporary means combine classic and new approaches
- Less dependent on computational resources (as compared to machine learning)
- Allow justification of the produced solutions
- Cover a larger set of tasks
- Natural for math science, and require higher level of AI education

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

基于 PCF 的问题解决方法



- 机器人团队的目标是将积木运送到目标区域
- 每个积木可由两个或多个机器人拖动
- 世界的当前状态和小组的目标在 PCF 中形式化为
- 基于 PCF 的求证器和选择机制为团队生成最优的联合行动计划
- 只要“世界”的状态发生变化，就可以轻松修改当前计划

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

A Master Degree Program. 语义技术与多代理系统

该课程由圣彼得堡电工技术大学（LETI）、伊尔库茨克国立大学（Irkutsk State University）和俄罗斯科学院空间技术研究所（ISDCT SB RAS）联合开设。

主要课题。

- 计算几何、数字信号处理、物联网、
- 语义网、语义网信息系统开发、
- 人工智能基础、知识表示、面向对象逻辑编程、
- 答案集编程 (SAT)、自然语言处理、
- 机器学习、神经网络、深度学习、
- 多代理系统，多代理系统优化。

始于 2022-09-01。

<https://etu.ru/sveden/education/programs/semanticheskie-tehnologii-i-mnogoagentnye-sistemy-01.04.02.html>

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

结论（最后）

- 经典的基于知识的系统是强大的人工智能工具，可用于解决广泛的识别问题和各种综合问题：源代码、数据对象、控制、数据处理、数据分析、数据挖掘。
- 当代手段结合了经典方法和新方法
- 较少依赖计算资源（与机器学习相比）
- 允许对所产生的解决方案进行论证
- 涵盖更多任务
- 自然适用于数学科学，需要更高水平的人工智能教育

E. Cherkashin, et al. Logical Approach in Software and Data Design 软件和数据设计中的逻辑方法

Thank You!
谢谢大家!

