

The University Database Integration: An AutoMed Example

AutoMed Technical Report 12: Draft Version A

Peter McBrien

Dept. of Computing, Imperial College, pjm@doc.ic.ac.uk

Abstract

This reports describes a contrived set of database schemas, and how they are integrated using the AutoMed approach. It serves to illustrate several common techniques used in a AutoMed integration of databases.

1 Introduction

AutoMed [2] supports many methodologies for performing data integration and hence forming a **network** of pathways joining schemas together. Here we describe a simple methodology based on forming union-compatible schemas, the general structure of which is illustrated in Figure 1. Each of the n local schemas LS_i is first transformed into a “union” schema US_i . These n union schemas are syntactically identical, and this is asserted by creating a sequence of *id* transformation steps between each pair of union schemas US_i and US_{i+1} , of the form $id(US_i : c, US_{i+1} : c)$ for each schema construct c . These *id* transformations between pairs of union schemas are generated automatically by the AutoMed software. An arbitrary one of the union schemas can then be selected for further transformation into the global schema GS . This two stage process reflects first schema conformance, followed by schema integration and restructuring.

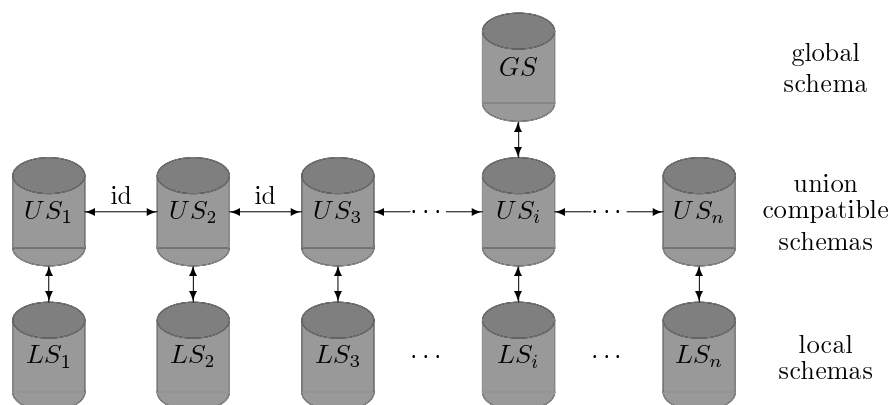


Figure 1: A general AutoMed Integration

| | | | |
|-----------------|-------------------------------------------------------------------------------------------------------------------------------------------|-----------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| LS ₁ | dept(<u>dname</u>) staff(<u>id</u> , name, sex, dname) | LS ₅ | university(uname) college(<u>cname</u> , uname) dept(<u>dname</u> , street, cname) staff(<u>id</u> , name, sex, dname) |
| LS ₂ | staff(<u>id</u> , name, dname) male(<i>id</i>) female(<i>id</i>) | US | college(<u>cname</u>) dept(<u>dname</u> , street, cname) degree(<u>dcode</u> , title, dname) staff(<u>id</u> , name, sex, dname) student(<u>id</u> , sex, dname) |
| LS ₃ | dept(<u>deptname</u>) degree(<u>dcode</u> , title, dname) person(<u>id</u> , dname) male(<i>id</i>) female(<i>id</i>) | GS | college(<u>cname</u>) dept(<u>dname</u> , street, cname) degree(<u>dcode</u> , title, dname) person(<u>id</u> , name#, sex, dname) |
| LS ₄ | dept(<u>dname</u>) student(<u>id</u> , sex, dname) degree(<u>dcode</u> , dname) | | |

Figure 2: Example schemas

2 The University Global Database Example

Figure 2 gives some specific schemas that illustrate the integration approach of Figure 1. Primary key attributes are underlined, foreign key attributes are in italics and nullable attributes are suffixed by #.

The schema US_1 is the simplest to integrate into US since the only constructs it contains are exactly the same in extent as the same constructs in US . Hence the pathway $LS_1 \rightarrow US$ in Example 1 involves only *extend* transformations to add each of the tables *student*, *college* and *degree*, plus the two attributes of *dept* which are present in US but not LS_1 , and *add* transformations to assert the missing primary and foreign key constraints:

Pathway 1 $LS_1 \rightarrow US$

- ① extendTable(⟨⟨student⟩⟩)
- ② extendField(⟨⟨student, id⟩⟩)
- ③ extendField(⟨⟨student, sex⟩⟩)
- ④ extendField(⟨⟨student, dname⟩⟩)
- ⑤ addPK(⟨⟨student, ⟨⟨student, id⟩⟩⟩⟩)
- ⑥ addFK(⟨⟨student, ⟨⟨student, dname⟩⟩, dept, ⟨⟨dept, dname⟩⟩⟩⟩)
- ⑦ extendTable(⟨⟨college⟩⟩)
- ⑧ extendField(⟨⟨college, cname⟩⟩)
- ⑨ addPK(⟨⟨college, ⟨⟨college, cname⟩⟩⟩⟩)
- ⑩ extendTable(⟨⟨degree⟩⟩)
- ⑪ extendField(⟨⟨degree, dcode⟩⟩)
- ⑫ extendField(⟨⟨degree, title⟩⟩)
- ⑬ extendField(⟨⟨degree, dname⟩⟩)
- ⑭ addPK(⟨⟨degree, ⟨⟨degree, dcode⟩⟩⟩⟩)
- ⑮ addFK(⟨⟨degree, ⟨⟨degree, dname⟩⟩, dept, ⟨⟨dept, dname⟩⟩⟩⟩)
- ⑯ extendField(⟨⟨dept, street⟩⟩)
- ⑰ extendField(⟨⟨dept, cname⟩⟩)
- ⑱ addFK(⟨⟨dept, ⟨⟨dept, cname⟩⟩, college, ⟨⟨college, cname⟩⟩⟩⟩)

In Example 2, transformations ①9–③0 use *extend* transformations to state that the tables `student`, `college` and `degree` in `US` cannot be derived from `LS2`. Then ③3–③6 use the `dname` attribute of `person` to derive the `dept` table in `US`, and use *extend* transformations for the two attributes `street` and `cname` that cannot be derived from `LS2`. Finally, in ④1–④9 the `male` and `female` relations of `LS2` are restructured into the single `sex` attribute of `staff`.

Pathway 2 `LS2 → US`

```

①9 extendTable(⟨⟨student⟩⟩)
②0 extendField(⟨⟨student,id⟩⟩)
②1 extendField(⟨⟨student,sex⟩⟩)
②2 extendField(⟨⟨student,dname⟩⟩)
②3 addPK(⟨⟨student,⟨⟨student,id⟩⟩⟩⟩)
②4 extendTable(⟨⟨college⟩⟩)
②5 extendField(⟨⟨college,cname⟩⟩)
②6 addPK(⟨⟨college,⟨⟨college,cname⟩⟩⟩⟩)
②7 extendTable(⟨⟨degree⟩⟩)
②8 extendField(⟨⟨degree,dcode⟩⟩)
②9 extendField(⟨⟨degree,title⟩⟩)
③0 extendField(⟨⟨degree,dname⟩⟩)
③1 addPK(⟨⟨degree,⟨⟨degree,dcode⟩⟩⟩⟩)
③2 addFK(⟨⟨degree,⟨⟨degree,dname⟩⟩,dept,⟨⟨dept,dname⟩⟩⟩⟩)
③3 addTable(⟨⟨dept⟩⟩, [x | (y,x) ← ⟨⟨staff,dname⟩⟩])
③4 addField(⟨⟨dept,dname⟩⟩, [(x,x) | x ← ⟨⟨dept⟩⟩])
③5 extendField(⟨⟨dept,street⟩⟩)
③6 extendField(⟨⟨dept,cname⟩⟩)
③7 addFK(⟨⟨degree,⟨⟨degree,dname⟩⟩,dept,⟨⟨dept,dname⟩⟩⟩⟩)
③8 addFK(⟨⟨staff,⟨⟨staff,dname⟩⟩,dept,⟨⟨dept,dname⟩⟩⟩⟩)
③9 addFK(⟨⟨student,⟨⟨student,dname⟩⟩,dept,⟨⟨dept,dname⟩⟩⟩⟩)
④0 addFK(⟨⟨dept,⟨⟨dept,cname⟩⟩,college,⟨⟨college,cname⟩⟩⟩⟩)
④1 addField(⟨⟨staff,sex⟩⟩, [(x,' M') | x ← ⟨⟨male⟩⟩] ++ [(x,' F') | x ← ⟨⟨female⟩⟩])
④2 deleteFK(⟨⟨male,⟨⟨male,id⟩⟩⟩, ⟨⟨person,⟨⟨person,id⟩⟩⟩⟩)
④3 deletePK(⟨⟨male,⟨⟨male,id⟩⟩⟩⟩)
④4 deleteField(⟨⟨male,id⟩⟩, [(x,x) | x ← ⟨⟨male⟩⟩])
④5 deleteTable(⟨⟨male⟩⟩, [x | (x,' M') ← ⟨⟨staff,sex⟩⟩])
④6 deleteFK(⟨⟨female,⟨⟨female,id⟩⟩⟩, ⟨⟨person,⟨⟨person,id⟩⟩⟩⟩)
④7 deletePK(⟨⟨female,⟨⟨female,id⟩⟩⟩⟩)
④8 deleteField(⟨⟨female,id⟩⟩, [(x,x) | x ← ⟨⟨female⟩⟩])
④9 deleteTable(⟨⟨female⟩⟩, [x | (x,' F') ← ⟨⟨staff,sex⟩⟩])

```

The example uses the **Intermediate Query Language (IQL)**, which is the default query language supported by the AutoMed implementation. In IQL `++` is the bag union operator and the construct `[e | Q1; ... Qn]` is a **comprehension** [1]. The expressions `Q1` to `Qn` are termed **qualifiers**, each qualifier being either a **filter** or a **generator**. A filter is a boolean-valued expression. A generator has syntax `p ← c` where `p` is a **pattern** and `c` is a bag-valued expression. In IQL, the patterns `p` are restricted to be single variables or tuples of variables.

The pathway `LS3 → US` contains *extend* steps ⑤0–⑤7 to add the missing `student` and `college` tables, which are textually the same as ①9–②5. It then renames `deptname`, adds the missing attributes of `dept`, renames `person` to `staff`, and adds the missing `name` attribute. Finally, in steps ⑥4–⑦1 it does the same restructuring as steps ④1–④9 of `LS2 → US`, converting the `male` and `female` relations into the single `sex` attribute of `staff`.

Pathway 3 `LS3 → US`

```

50 extendTable(⟨⟨student⟩⟩)
51 extendField(⟨⟨student,id⟩⟩)
52 extendField(⟨⟨student,sex⟩⟩)
53 extendField(⟨⟨student,dname⟩⟩)
54 addPK(⟨⟨student,⟨⟨student,id⟩⟩⟩⟩)
55 addFK(⟨⟨student,⟨⟨student,dname⟩⟩,dept,⟨⟨dept,deptname⟩⟩⟩⟩)
56 extendTable(⟨⟨college⟩⟩)
57 extendField(⟨⟨college,cname⟩⟩)
58 addPK(⟨⟨college,⟨⟨college,cname⟩⟩⟩⟩)
59 renameField(⟨⟨dept,deptname⟩⟩,⟨⟨dept,dname⟩⟩)
60 extendField(⟨⟨dept,street⟩⟩)
61 extendField(⟨⟨dept,cname⟩⟩)
62 renameTable(⟨⟨person⟩⟩,⟨⟨staff⟩⟩)
63 extendField(⟨⟨staff,name⟩⟩)
64 addField(⟨⟨staff,sex⟩⟩,[⟨x,'M'⟩|x←⟨⟨male⟩⟩]+[⟨x,'F'⟩|x←⟨⟨female⟩⟩])
65 deleteFK(⟨⟨male,⟨⟨male,id⟩⟩⟩,⟨⟨person,⟨⟨person,id⟩⟩⟩⟩)
66 deletePK(⟨⟨male,⟨⟨male,id⟩⟩⟩⟩)
67 deleteField(⟨⟨male,id⟩⟩,[⟨x,x⟩|x←⟨⟨male⟩⟩])
68 deleteTable(⟨⟨male⟩⟩,[x|⟨x,'M'⟩←⟨⟨staff,sex⟩⟩])
69 deleteFK(⟨⟨female,⟨⟨female,id⟩⟩⟩,⟨⟨person,⟨⟨person,id⟩⟩⟩⟩)
70 deletePK(⟨⟨female,⟨⟨female,id⟩⟩⟩⟩)
71 deleteField(⟨⟨female,id⟩⟩,[⟨x,x⟩|x←⟨⟨female⟩⟩])
72 deleteTable(⟨⟨female⟩⟩,[x|⟨x,'F'⟩←⟨⟨staff,sex⟩⟩])

```

The pathway $LS_4 \rightarrow US$ contains a sequence of *extend* steps for its missing information.

Pathway 4 $LS_4 \rightarrow US$

```

73 extendTable(⟨⟨staff⟩⟩)
74 extendField(⟨⟨staff,id⟩⟩)
75 extendField(⟨⟨staff,name⟩⟩)
76 extendField(⟨⟨staff,sex⟩⟩)
77 extendField(⟨⟨staff,dname⟩⟩)
78 addPK(⟨⟨staff,⟨⟨staff,id⟩⟩⟩⟩)
79 addFK(⟨⟨staff,⟨⟨staff,dname⟩⟩,dept,⟨⟨dept,dname⟩⟩⟩⟩)
80 extendTable(⟨⟨college⟩⟩)
81 extendField(⟨⟨college,cname⟩⟩)
82 addPK(⟨⟨college,⟨⟨college,cname⟩⟩⟩⟩)
83 extendField(⟨⟨degree,title⟩⟩)
84 extendField(⟨⟨dept,street⟩⟩)
85 extendField(⟨⟨dept,cname⟩⟩)

```

The pathway $LS_5 \rightarrow US$ contains a sequence of *extend* steps for its missing information and also three *contract* steps to remove the university relation and its attributes.

Pathway 5 $LS_5 \rightarrow US$

```

86 renameTable(⟨⟨person⟩⟩,⟨⟨staff⟩⟩)
87 extendTable(⟨⟨student⟩⟩)
88 extendField(⟨⟨student,id⟩⟩)

```

```

⑧9 extendField(⟨⟨student,sex⟩⟩)
⑨0 extendField(⟨⟨student,dname⟩⟩)
⑨1 addPK(⟨⟨student,⟨⟨student,id⟩⟩⟩⟩)
⑨2 addFK(⟨⟨student,⟨⟨student,dname⟩⟩,dept,⟨⟨dept,dname⟩⟩⟩⟩)
⑨3 extendField(⟨⟨degree⟩⟩)
⑨4 extendField(⟨⟨degree,dcode⟩⟩)
⑨5 extendField(⟨⟨degree,title⟩⟩)
⑨6 extendField(⟨⟨degree,dname⟩⟩)
⑨7 addPK(⟨⟨degree,⟨⟨degree,dcode⟩⟩⟩⟩)
⑨8 addFK(⟨⟨degree,⟨⟨degree,dname⟩⟩,dept,⟨⟨dept,dname⟩⟩⟩⟩)
⑨9 deleteFK(⟨⟨college,⟨⟨college,uname⟩⟩,university,⟨⟨university,uname⟩⟩⟩⟩)
100 contractField(⟨⟨college,uname⟩⟩)
101 deletePK(⟨⟨university,⟨⟨university,uname⟩⟩⟩⟩)
102 contractField(⟨⟨university,uname⟩⟩)
103 contractTable(⟨⟨university⟩⟩)

```

Finally, we list below the pathway from the union schema US to the global schema GS:

Pathway 6 US → GS

```

104 addTable(⟨⟨person⟩⟩, ⟨⟨staff⟩⟩ ++ ⟨⟨student⟩⟩)
105 addField(⟨⟨person,id⟩⟩, ⟨⟨staff,id⟩⟩ ++ ⟨⟨student,id⟩⟩)
106 addField(⟨⟨person,name⟩⟩, ⟨⟨staff,name⟩⟩)
107 addField(⟨⟨person,sex⟩⟩, ⟨⟨staff,sex⟩⟩ ++ ⟨⟨student,sex⟩⟩)
108 addField(⟨⟨person,dname⟩⟩, ⟨⟨staff,dname⟩⟩ ++ ⟨⟨student,dname⟩⟩)
109 addPK(⟨⟨person,⟨⟨person,id⟩⟩⟩⟩)
110 addFK(⟨⟨person,⟨⟨person,dname⟩⟩,dept,⟨⟨dept,dname⟩⟩⟩⟩)
111 deleteFK(⟨⟨student,⟨⟨student,dname⟩⟩,dept,⟨⟨dept,dname⟩⟩⟩⟩)
112 deletePK(⟨⟨student,⟨⟨student,id⟩⟩⟩⟩)
113 deleteField(⟨⟨student,id⟩⟩, [(x, y) | x ← ⟨⟨student⟩⟩; (x, y) ← ⟨⟨person,id⟩⟩])
114 deleteField(⟨⟨student,sex⟩⟩, [(x, y) | x ← ⟨⟨student⟩⟩; (x, y) ← ⟨⟨person,sex⟩⟩])
115 deleteField(⟨⟨student,dname⟩⟩, [(x, y) | x ← ⟨⟨student⟩⟩; (x, y) ← ⟨⟨person,dname⟩⟩])
116 deleteTable(⟨⟨student⟩⟩, [x | x ← ⟨⟨person⟩⟩; not(member⟨⟨staff⟩⟩x)])
117 deleteFK(⟨⟨staff,⟨⟨staff,dname⟩⟩,dept,⟨⟨dept,dname⟩⟩⟩⟩)
118 deletePK(⟨⟨staff,⟨⟨staff,id⟩⟩⟩⟩)
119 deleteField(⟨⟨staff,id⟩⟩, [(x, y) | x ← ⟨⟨staff⟩⟩; (x, y) ← ⟨⟨person,id⟩⟩])
120 deleteField(⟨⟨staff,name⟩⟩, [(x, y) | x ← ⟨⟨staff⟩⟩; (x, y) ← ⟨⟨person,name⟩⟩])
121 deleteField(⟨⟨staff,sex⟩⟩, [(x, y) | x ← ⟨⟨staff⟩⟩; (x, y) ← ⟨⟨person,sex⟩⟩])
122 deleteField(⟨⟨staff,dname⟩⟩, [(x, y) | x ← ⟨⟨staff⟩⟩; (x, y) ← ⟨⟨person,dname⟩⟩])
123 deleteTable(⟨⟨staff⟩⟩, [x | x ← ⟨⟨person⟩⟩; (x, y) ← ⟨⟨person,name⟩⟩])

```

3 Conclusions

See <http://www.doc.ic.ac.uk/automed/> for more details of AutoMed.

References

- [1] P. Buneman *et al.* Comprehension syntax. *ACM SIGMOD Record*, 23(1):87–96, 1994.
- [2] P.J. McBrien and A. Poulouvasilis. Data integration by bi-directional schema transformation rules. In *Proc. ICDE'03 (to appear)*, 2003.