

Semistructured Data: The TSIMMIS Experience

Joachim Hammer, Jason McHugh, and Hector Garcia-Molina

Department of Computer Science
Stanford University
Stanford, CA 94305-9040
U.S.A.

{joachim,mchughj,hector}@db.stanford.edu
<http://www-db.stanford.edu>

Abstract

In this paper we discuss the management of semi-structured data, i.e., data that has irregular or dynamically changing structure. We describe components of the Stanford TSIMMIS Project that help extract semi-structured data from Web pages, that allow the storage and querying of semi-structured data, and that allow its browsing through the World Wide Web. A prototype implementation of the TSIMMIS system as described here is currently installed and running in the database group testbed.

1 Introduction

At a recent workshop on management of semistructured data [15], the workshop attendants defined *semistructured data* as data that does not have a regular and static structure like data found in a relational database but whose schema is dynamic and may contain missing data or types. For example, if we look at weather forecasts on the Web, the "fields" and their structure may differ across sites. Even at a single site, some forecasts may be missing information, or may have extra information depending on the geographical location of the affected region (e.g., cities in the Rocky Mountains usually include ski reports in the winter months whereas forecast for tropical resorts do not). However, semistructured data is not just limited to the World Wide Web (WWW), but is also found in many other interesting sources including file systems, news wires, electronic mail systems, etc. just to name a few. In addition to occurring natively in the above-mentioned classes of sources, semistructured data is often a "by-product" of the integration process when multiple heterogeneous schemas are involved. In those cases, semistructured (rather than "fully structured") data arises because the integrated objects may be based on complementary, sometimes conflicting, and often dynamic information from multiple sources, forcing the integrator to filter, merge, or omit certain fields when performing the integration.

The goal of the TSIMMIS project at Stanford [4, 7, 11, 13] is to provide integrated access to a wide variety of heterogeneous data sources (e.g., databases, object stores, knowl-

edge bases, digital libraries) including sources containing semistructured data (e.g., WWW, file system). In this paper, we present the TSIMMIS approach to managing semistructured data. In particular, we discuss three critical aspects semistructured data management: (1) extracting the intended content from its native source (how to get it?), (2) reading the extracted data (how to query it?), and (3), exploring the result in a graphical, easy-to-understand manner (how to browse it?). In TSIMMIS we have developed components that address all of the above issues and together provide an integrated solution to the problem of managing semistructured data. Several other recent projects have similar goals (e.g., LORE [10], Garlic [3], Information Manifold [9], Rufus [14]), but we do not survey them here.

2 Representing Semistructured Data in TSIMMIS

For the TSIMMIS project we have adopted a simple *self-describing* (or *tagged*) object model. Similar models have been in use for years; we call our version the *Object Exchange Model*, or OEM [4]. OEM is a flexible model that is particularly well suited for representing semistructured data. Data represented in OEM constitutes a graph, with a unique root object at the top and zero or more nested subobjects. The fundamental idea is that all objects, and their subobjects, have *labels* that describe their meaning. For example, the following object represents a Fahrenheit temperature of 80 degrees:

```
temp-in-Fahrenheit, int, 80
```

Here, the string "temp-in-Fahrenheit" is a human-readable label, "int" indicates an integer value, and "80" is the value itself. If we wish to represent a *complex* object, then each component of the object has its own label. For example, an object representing a set of two temperatures may look like:

```
set-of-temps, complex, {  
  temp-in-Fahrenheit, int, 80  
  temp-in-Celsius, int, 20 }
```

OEM is very simple, while providing the expressive power and flexibility needed for representing semistructured data from a wide range of heterogeneous sources. Our primary reason for choosing a simple model is to be able to accommodate a wide variety of external data models and to facilitate integration. As pointed out in [2], a simple model such as OEM has an advantage over complex models when used for representing and integrating heterogeneous data, since the

```

1 <HTML>
2 <HEAD>
3 <TITLE>INTELLICAST: europe weather</TITLE>
4 <BASE href="http://"></BASE>
5 <TABLE border="0" cellpadding="0" cellspacing="0" width="800">
6 <TR>
7 <TD colspan="4">Click on a city for forecasts</TD></TR>
8 </TD>
9 </TR>
10 <TR>
11 <TD colspan="4">temperatures in degrees celsius</TD></TR>
12 </TD>
13 </TR>
14 <TR>
15 <TD colspan="4">Helsinki</TD></TR>
16 </TD>
17 </TR>
18 <TR>
19 <TD colspan="4">
20 <TABLE cellpadding="0" cellspacing="0" width="610">
21 <TR>
22 <TH align="left">
23 <TR>
24 <TR>
25 <TR>
26 <TR>
27 <TR>
28 <TR>
29 <TR>
30 <TR>
31 <TR>
32 <TR>
33 <TR>
34 <TR>
35 <TR>
36 <TR>
37 <TR>
38 <TR>
39 <TR>
40 <TR>
41 <TR>
42 <TR>
43 <TR>
44 <TR>
45 <TR>
46 <TR>
47 <TR>
48 <TR>
49 <TR>
50 <TR>
51 <TR>
52 <TR>
53 <TR>
54 <TR>
55 <TR>
56 <TR>
57 <TR>
58 <TR>
59 <TR>
60 <TR>
61 <TR>
62 <TR>
63 <TR>
64 <TR>
65 <TR>
66 <TR>
67 <TR>
68 <TR>
69 <TR>
70 <TR>
71 <TR>
72 <TR>
73 <TR>
74 <TR>
75 <TR>
76 <TR>
77 <TR>
78 <TR>
79 <TR>
80 <TR>
81 <TR>
82 <TR>
83 <TR>
84 <TR>
85 <TR>
86 <TR>
87 <TR>
88 <TR>
89 <TR>
90 <TR>
91 <TR>
92 <TR>
93 <TR>
94 <TR>
95 <TR>
96 <TR>
97 <TR>
98 <TR>
99 <TR>
100 <TR>
101 <TR>
102 <TR>
103 <TR>
104 <TR>
105 <TR>
106 <TR>
107 <TR>
108 <TR>
109 <TR>
110 <TR>
111 <TR>
112 <TR>
113 <TR>
114 <TR>
115 <TR>
116 <TR>
117 <TR>
118 <TR>
119 <TR>
120 <TR>
121 <TR>
122 <TR>
123 <TR>
124 <TR>
125 <TR>
126 <TR>
127 <TR>
128 <TR>
129 <TR>
130 <TR>
131 <TR>
132 <TR>
133 <TR>
134 <TR>
135 <TR>
136 <TR>
137 <TR>
138 <TR>
139 <TR>
140 <TR>
141 <TR>
142 <TR>
143 <TR>
144 <TR>
145 <TR>
146 <TR>
147 <TR>
148 <TR>
149 <TR>
150 <TR>
151 <TR>
152 <TR>
153 <TR>
154 <TR>
155 <TR>
156 <TR>
157 <TR>
158 <TR>
159 <TR>
160 <TR>
161 <TR>
162 <TR>
163 <TR>
164 <TR>
165 <TR>
166 <TR>
167 <TR>
168 <TR>
169 <TR>
170 <TR>
171 <TR>
172 <TR>
173 <TR>
174 <TR>
175 <TR>
176 <TR>
177 <TR>
178 <TR>
179 <TR>
180 <TR>
181 <TR>
182 <TR>
183 <TR>
184 <TR>
185 <TR>
186 <TR>
187 <TR>
188 <TR>
189 <TR>
190 <TR>
191 <TR>
192 <TR>
193 <TR>
194 <TR>
195 <TR>
196 <TR>
197 <TR>
198 <TR>
199 <TR>
200 <TR>
201 <TR>
202 <TR>
203 <TR>
204 <TR>
205 <TR>
206 <TR>
207 <TR>
208 <TR>
209 <TR>
210 <TR>
211 <TR>
212 <TR>
213 <TR>
214 <TR>
215 <TR>
216 <TR>
217 <TR>
218 <TR>
219 <TR>
220 <TR>
221 <TR>
222 <TR>
223 <TR>
224 <TR>
225 <TR>
226 <TR>
227 <TR>
228 <TR>
229 <TR>
230 <TR>
231 <TR>
232 <TR>
233 <TR>
234 <TR>
235 <TR>
236 <TR>
237 <TR>
238 <TR>
239 <TR>
240 <TR>
241 <TR>
242 <TR>
243 <TR>
244 <TR>
245 <TR>
246 <TR>
247 <TR>
248 <TR>
249 <TR>
250 <TR>
251 <TR>
252 <TR>
253 <TR>
254 <TR>
255 <TR>
256 <TR>
257 <TR>
258 <TR>
259 <TR>
260 <TR>
261 <TR>
262 <TR>
263 <TR>
264 <TR>
265 <TR>
266 <TR>
267 <TR>
268 <TR>
269 <TR>
270 <TR>
271 <TR>
272 <TR>
273 <TR>
274 <TR>
275 <TR>
276 <TR>
277 <TR>
278 <TR>
279 <TR>
280 <TR>
281 <TR>
282 <TR>
283 <TR>
284 <TR>
285 <TR>
286 <TR>
287 <TR>
288 <TR>
289 <TR>
290 <TR>
291 <TR>
292 <TR>
293 <TR>
294 <TR>
295 <TR>
296 <TR>
297 <TR>
298 <TR>
299 <TR>
300 <TR>
301 <TR>
302 <TR>
303 <TR>
304 <TR>
305 <TR>
306 <TR>
307 <TR>
308 <TR>
309 <TR>
310 <TR>
311 <TR>
312 <TR>
313 <TR>
314 <TR>
315 <TR>
316 <TR>
317 <TR>
318 <TR>
319 <TR>
320 <TR>
321 <TR>
322 <TR>
323 <TR>
324 <TR>
325 <TR>
326 <TR>
327 <TR>
328 <TR>
329 <TR>
330 <TR>
331 <TR>
332 <TR>
333 <TR>
334 <TR>
335 <TR>
336 <TR>
337 <TR>
338 <TR>
339 <TR>
340 <TR>
341 <TR>
342 <TR>
343 <TR>
344 <TR>
345 <TR>
346 <TR>
347 <TR>
348 <TR>
349 <TR>
350 <TR>
351 <TR>
352 <TR>
353 <TR>
354 <TR>
355 <TR>
356 <TR>
357 <TR>
358 <TR>
359 <TR>
360 <TR>
361 <TR>
362 <TR>
363 <TR>
364 <TR>
365 <TR>
366 <TR>
367 <TR>
368 <TR>
369 <TR>
370 <TR>
371 <TR>
372 <TR>
373 <TR>
374 <TR>
375 <TR>
376 <TR>
377 <TR>
378 <TR>
379 <TR>
380 <TR>
381 <TR>
382 <TR>
383 <TR>
384 <TR>
385 <TR>
386 <TR>
387 <TR>
388 <TR>
389 <TR>
390 <TR>
391 <TR>
392 <TR>
393 <TR>
394 <TR>
395 <TR>
396 <TR>
397 <TR>
398 <TR>
399 <TR>
400 <TR>
401 <TR>
402 <TR>
403 <TR>
404 <TR>
405 <TR>
406 <TR>
407 <TR>
408 <TR>
409 <TR>
410 <TR>
411 <TR>
412 <TR>
413 <TR>
414 <TR>
415 <TR>
416 <TR>
417 <TR>
418 <TR>
419 <TR>
420 <TR>
421 <TR>
422 <TR>
423 <TR>
424 <TR>
425 <TR>
426 <TR>
427 <TR>
428 <TR>
429 <TR>
430 <TR>
431 <TR>
432 <TR>
433 <TR>
434 <TR>
435 <TR>
436 <TR>
437 <TR>
438 <TR>
439 <TR>
440 <TR>
441 <TR>
442 <TR>
443 <TR>
444 <TR>
445 <TR>
446 <TR>
447 <TR>
448 <TR>
449 <TR>
450 <TR>
451 <TR>
452 <TR>
453 <TR>
454 <TR>
455 <TR>
456 <TR>
457 <TR>
458 <TR>
459 <TR>
460 <TR>
461 <TR>
462 <TR>
463 <TR>
464 <TR>
465 <TR>
466 <TR>
467 <TR>
468 <TR>
469 <TR>
470 <TR>
471 <TR>
472 <TR>
473 <TR>
474 <TR>
475 <TR>
476 <TR>
477 <TR>
478 <TR>
479 <TR>
480 <TR>
481 <TR>
482 <TR>
483 <TR>
484 <TR>
485 <TR>
486 <TR>
487 <TR>
488 <TR>
489 <TR>
490 <TR>
491 <TR>
492 <TR>
493 <TR>
494 <TR>
495 <TR>
496 <TR>
497 <TR>
498 <TR>
499 <TR>
500 <TR>
501 <TR>
502 <TR>
503 <TR>
504 <TR>
505 <TR>
506 <TR>
507 <TR>
508 <TR>
509 <TR>
510 <TR>
511 <TR>
512 <TR>
513 <TR>
514 <TR>
515 <TR>
516 <TR>
517 <TR>
518 <TR>
519 <TR>
520 <TR>
521 <TR>
522 <TR>
523 <TR>
524 <TR>
525 <TR>
526 <TR>
527 <TR>
528 <TR>
529 <TR>
530 <TR>
531 <TR>
532 <TR>
533 <TR>
534 <TR>
535 <TR>
536 <TR>
537 <TR>
538 <TR>
539 <TR>
540 <TR>
541 <TR>
542 <TR>
543 <TR>
544 <TR>
545 <TR>
546 <TR>
547 <TR>
548 <TR>
549 <TR>
550 <TR>
551 <TR>
552 <TR>
553 <TR>
554 <TR>
555 <TR>
556 <TR>
557 <TR>
558 <TR>
559 <TR>
560 <TR>
561 <TR>
562 <TR>
563 <TR>
564 <TR>
565 <TR>
566 <TR>
567 <TR>
568 <TR>
569 <TR>
570 <TR>
571 <TR>
572 <TR>
573 <TR>
574 <TR>
575 <TR>
576 <TR>
577 <TR>
578 <TR>
579 <TR>
580 <TR>
581 <TR>
582 <TR>
583 <TR>
584 <TR>
585 <TR>
586 <TR>
587 <TR>
588 <TR>
589 <TR>
590 <TR>
591 <TR>
592 <TR>
593 <TR>
594 <TR>
595 <TR>
596 <TR>
597 <TR>
598 <TR>
599 <TR>
600 <TR>
601 <TR>
602 <TR>
603 <TR>
604 <TR>
605 <TR>
606 <TR>
607 <TR>
608 <TR>
609 <TR>
610 <TR>
611 <TR>
612 <TR>
613 <TR>
614 <TR>
615 <TR>
616 <TR>
617 <TR>
618 <TR>
619 <TR>
620 <TR>
621 <TR>
622 <TR>
623 <TR>
624 <TR>
625 <TR>
626 <TR>
627 <TR>
628 <TR>
629 <TR>
630 <TR>
631 <TR>
632 <TR>
633 <TR>
634 <TR>
635 <TR>
636 <TR>
637 <TR>
638 <TR>
639 <TR>
640 <TR>
641 <TR>
642 <TR>
643 <TR>
644 <TR>
645 <TR>
646 <TR>
647 <TR>
648 <TR>
649 <TR>
650 <TR>
651 <TR>
652 <TR>
653 <TR>
654 <TR>
655 <TR>
656 <TR>
657 <TR>
658 <TR>
659 <TR>
660 <TR>
661 <TR>
662 <TR>
663 <TR>
664 <TR>
665 <TR>
666 <TR>
667 <TR>
668 <TR>
669 <TR>
670 <TR>
671 <TR>
672 <TR>
673 <TR>
674 <TR>
675 <TR>
676 <TR>
677 <TR>
678 <TR>
679 <TR>
680 <TR>
681 <TR>
682 <TR>
683 <TR>
684 <TR>
685 <TR>
686 <TR>
687 <TR>
688 <TR>
689 <TR>
690 <TR>
691 <TR>
692 <TR>
693 <TR>
694 <TR>
695 <TR>
696 <TR>
697 <TR>
698 <TR>
699 <TR>
700 <TR>
701 <TR>
702 <TR>
703 <TR>
704 <TR>
705 <TR>
706 <TR>
707 <TR>
708 <TR>
709 <TR>
710 <TR>
711 <TR>
712 <TR>
713 <TR>
714 <TR>
715 <TR>
716 <TR>
717 <TR>
718 <TR>
719 <TR>
720 <TR>
721 <TR>
722 <TR>
723 <TR>
724 <TR>
725 <TR>
726 <TR>
727 <TR>
728 <TR>
729 <TR>
730 <TR>
731 <TR>
732 <TR>
733 <TR>
734 <TR>
735 <TR>
736 <TR>
737 <TR>
738 <TR>
739 <TR>
740 <TR>
741 <TR>
742 <TR>
743 <TR>
744 <TR>
745 <TR>
746 <TR>
747 <TR>
748 <TR>
749 <TR>
750 <TR>
751 <TR>
752 <TR>
753 <TR>
754 <TR>
755 <TR>
756 <TR>
757 <TR>
758 <TR>
759 <TR>
760 <TR>
761 <TR>
762 <TR>
763 <TR>
764 <TR>
765 <TR>
766 <TR>
767 <TR>
768 <TR>
769 <TR>
770 <TR>
771 <TR>
772 <TR>
773 <TR>
774 <TR>
775 <TR>
776 <TR>
777 <TR>
778 <TR>
779 <TR>
780 <TR>
781 <TR>
782 <TR>
783 <TR>
784 <TR>
785 <TR>
786 <TR>
787 <TR>
788 <TR>
789 <TR>
790 <TR>
791 <TR>
792 <TR>
793 <TR>
794 <TR>
795 <TR>
796 <TR>
797 <TR>
798 <TR>
799 <TR>
800 <TR>
801 <TR>
802 <TR>
803 <TR>
804 <TR>
805 <TR>
806 <TR>
807 <TR>
808 <TR>
809 <TR>
810 <TR>
811 <TR>
812 <TR>
813 <TR>
814 <TR>
815 <TR>
816 <TR>
817 <TR>
818 <TR>
819 <TR>
820 <TR>
821 <TR>
822 <TR>
823 <TR>
824 <TR>
825 <TR>
826 <TR>
827 <TR>
828 <TR>
829 <TR>
830 <TR>
831 <TR>
832 <TR>
833 <TR>
834 <TR>
835 <TR>
836 <TR>
837 <TR>
838 <TR>
839 <TR>
840 <TR>
841 <TR>
842 <TR>
843 <TR>
844 <TR>
845 <TR>
846 <TR>
847 <TR>
848 <TR>
849 <TR>
850 <TR>
851 <TR>
852 <TR>
853 <TR>
854 <TR>
855 <TR>
856 <TR>
857 <TR>
858 <TR>
859 <TR>
860 <TR>
861 <TR>
862 <TR>
863 <TR>
864 <TR>
865 <TR>
866 <TR>
867 <TR>
868 <TR>
869 <TR>
870 <TR>
871 <TR>
872 <TR>
873 <TR>
874 <TR>
875 <TR>
876 <TR>
877 <TR>
878 <TR>
879 <TR>
880 <TR>
881 <TR>
882 <TR>
883 <TR>
884 <TR>
885 <TR>
886 <TR>
887 <TR>
888 <TR>
889 <TR>
890 <TR>
891 <TR>
892 <TR>
893 <TR>
894 <TR>
895 <TR>
896 <TR>
897 <TR>
898 <TR>
899 <TR>
900 <TR>
901 <TR>
902 <TR>
903 <TR>
904 <TR>
905 <TR>
906 <TR>
907 <TR>
908 <TR>
909 <TR>
910 <TR>
911 <TR>
912 <TR>
913 <TR>
914 <TR>
915 <TR>
916 <TR>
917 <TR>
918 <TR>
919 <TR>
920 <TR>
921 <TR>
922 <TR>
923 <TR>
924 <TR>
925 <TR>
926 <TR>
927 <TR>
928 <TR>
929 <TR>
930 <TR>
931 <TR>
932 <TR>
933 <TR>
934 <TR>
935 <TR>
936 <TR>
937 <TR>
938 <TR>
939 <TR>
940 <TR>
941 <TR>
942 <TR>
943 <TR>
944 <TR>
945 <TR>
946 <TR>
947 <TR>
948 <TR>
949 <TR>
950 <TR>
951 <TR>
952 <TR>
953 <TR>
954 <TR>
955 <TR>
956 <TR>
957 <TR>
958 <TR>
959 <TR>
960 <TR>
961 <TR>
962 <TR>
963 <TR>
964 <TR>
965 <TR>
966 <TR>
967 <TR>
968 <TR>
969 <TR>
970 <TR>
971 <TR>
972 <TR>
973 <TR>
974 <TR>
975 <TR>
976 <TR>
977 <TR>
978 <TR>
979 <TR>
980 <TR>
981 <TR>
982 <TR>
983 <TR>
984 <TR>
985 <TR>
986 <TR>
987 <TR>
988 <TR>
989 <TR>
990 <TR>
991 <TR>
992 <TR>
993 <TR>
994 <TR>
995 <TR>
996 <TR>
997 <TR>
998 <TR>
999 <TR>
1000 <TR>

```

Figure 1: A section of the HTML source file

operations to transform and merge data will be correspondingly simpler. Meanwhile a simple model can still be very powerful: advanced features can be “emulated” when they are necessary (e.g., subclass/superclass relationships, inheritance, etc.). For additional information on OEM, please refer to [12].

3 Extracting Data

Continuing with our weather example, let us assume that we have an application that needs to process weather data, such as temperature and forecast, for a given city. As one of its information sources, we want to use a Web site called Intellicast, which reports daily weather data for most major cities across the world. Since this site cannot be queried directly from within another application (e.g., “What is the forecast for Helsinki for May 7, 1997?”) we first have to extract the contents of the weather table from the underlying HTML page¹ which is displayed in Figure 1.

3.1 The Extraction Process

Our *configurable extraction* program parses this HTML page based on the specification file shown in Figure 2. The specification file consists of a sequence of *commands*, each defining one extraction step. Each command is of the form

```
[ variables, source, pattern ]
```

where *source* specifies the input text to be considered, *pattern* tells us how to find the text of interest within the source, and *variables* are one or more extractor variables that will hold the extracted results. The text in variables can be used as input for subsequent commands. (If a variable contains an extracted URL, we can also specify that the URL be followed, and that the linked page be used as further input.)

¹The line numbers shown on the left-hand side of this and the next figures are not part of the content but have been added to simplify the following discussion.

After the last command is executed, some subset of the variables will hold the data of interest. Later we describe how the contents of these variables are packaged into an OEM object.

Looking at Figure 2, we see that the list of commands is placed within the outermost brackets ‘[’ and ‘]’, and each command is also delimited by brackets. The extraction process in this example is performed by five commands. The initial command (lines 1-4) fetches the contents of the source file whose URL is given in line 2 into the variable called *root*. The ‘#’ character in line 3 means that everything (in this case the contents of the entire file) is to be extracted. After the file has been fetched and its contents are read into *root*, the extractor will filter out unwanted data such as the HTML markup commands and extra text with the remaining four commands.

The second command (lines 5-8) specifies that the result of applying the pattern in line 7 to the source variable *root* is to be stored in a new variable called *_temperature*. The pattern can be interpreted as follows: discard everything until the first occurrence of the token *</TR>* (‘*’ means discard) in the second table definition and save the data that is stored between *</TR>* and *</TABLE>* (‘#’ means save). The two *<TABLE>* tokens between the ‘*’ are used as navigational help to identify the correct *</TR>* token since there is no way of specifying a numbered occurrence of a token (i.e., discard everything until the third occurrence of *</TR>*). After this step, the variable *_temperature* contains the information that is stored in lines 22 and higher in the source file in Figure 1 (up to but not including the subsequent *</TABLE>* token which indicates the end of the temperature table). The underscore at the beginning of the name *_temperature* indicates that this is a temporary variable; its contents will not be included in the resulting OEM object.

The third command (lines 9-12) instructs the extractor to split the contents of the temperature variable into “chunks” of text, using the string *(TRALIGN = left)* (lines 22, 30, 38, etc. in Figure 1) as the “chunk” delimiter. Note, each “chunk” represents one row in the temperature table. The result of each split is stored in a temporary variable called *_citytemp*. The split operator can only be applied if the input is made up of equally structured pieces with a clearly defined delimiter separating the individual pieces. If one thinks of extractor variables as lists (up until now each list had only one member) then the result of the split operator can be viewed as a new list with as many members as there are rows in the temperature table. Thus from now on, when we apply a pattern to a variable, we really mean applying the pattern to *each* member of the variable, much like the apply operator in Lisp.

In command 4 (lines 13-16), the extractor copies the contents of each cell of the temporary array into the array *citytemp* starting with the second cell from the beginning. The first integer in the instruction *_citytemp[1 : 0]* indicates the beginning of the copying (since the array index starts at 0, 1 refers to the second cell), the second integer indicates the last cell to be included (counting from the end of the array). As a result, we have excluded the first row of the table which contains the individual column headings. Note, that we could have also filtered out the unwanted row in the second command by specifying an additional **</TR>* condition before the ‘#’ in line 7 of Figure 2. The final command (lines 17-20) extracts the individual values from each cell in the *citytemp* array and assigns them into the variables listed in line 17 (*country, c_url, city, etc.*).

```

1 ["root",
2  "get('http://www.intellicast.com/weather/europe/')",
3  "#",
4 ],
5 ["_temperatures",
6  "root",
7  "**<TABLE*<TABLE*</TR>#</TABLE>*"
8 ],
9 ["_citytemp",
10 "split(temperatures, '<TR ALIGN=left>')",
11 "#",
12 ],
13 ["city_temp",
14  "_citytemp[1:0]",
15  "#",
16 ],
17 ["country,c_url,city,w_tody,hgh_tody,low_today,w_tomorrow,hgh_tomorrow,low_tomorrow",
18  "city_temp",
19  "**<TD>#</TD>*HREF=#>#</A>*<TD>#</TD>*<TD>#</TD>*<TD>#</TD>*<TD>#</TD>*<TD>#</TD>*"
20 11

```

Figure 2: A sample extractor specification file

```

root    complex {
        city_temp complex {
            country string "Finland"
            city_url url http://www...
            city string "Helsinki"
            weather_today string "rain"
            high_today string "16"
            low_today string "4"
            weather_tom string "snow"
            high_tomorrow string "10"
            low_tomorrow string "4"
        }
        city_temp complex {
            country string "France"
            city_url url http://www...
            city string "Bordeaux"
            ...
        }
    }

```

Figure 3: The extracted information in OEM format

After the five commands have been executed, the variables hold the data of interest. This data is packaged into an OEM object, shown in Figure 3, with a structure that follows the extraction process. Notice that this sample object reflects the structure of our extractor specification file. That is, the root object of the OEM answer will have a label *root* because this was the first extracted variable. This object will have children objects with label *city_temp* and so on. Notice that the variables *_temperature* and *_citytemp* do not appear in the final result because they are declared as temporary variables.

3.2 Additional Capabilities

In addition to the basic capabilities described in our example, the extractor has components for automatic handling of HTML tables, for conditional parsing, and other services. The extractor can also follow URLs in the process, extracting data from multiple Web pages into a single OEM object. Overall, we believe the extractor provides natural facilities for extracting data, as well as for structuring it in different ways into OEM objects. For more details on the extractor, please refer to [8].

4 Querying Semistructured Data

In this section we introduce the LOREL query language, primarily through examples. LOREL is an extension of OQL and a full specification can be found in [1]. Here we highlight those features of the language that have an impact on the novel aspects of the system—features designed specifically for handling semistructured data. Many other useful features of LOREL (some inherited from OQL and others not) that are more standard will not be covered.

4.1 Simple LORE Examples

Our first example query introduces the basic building block of LOREL: the simple path expression, which is a name followed by a sequence of labels. For example, *Root.City.Location* is a simple path expression. Its semantics consists of the set of objects that can be reached starting with the *Root* object, following an edge to objects labeled *City*, then following an edge to objects labeled *Location*. Range variables can be assigned to path expressions, e.g., "*Root.City.LocationX*" specifies that *X* ranges over the set of locations.

Continuing with our European weather example, the following example query retrieves the locations of all cities located in England when evaluated over the sample OEM database shown in Figure 4.

Query 4.1 (LOREL)

```

SELECT Root.City.Location
WHERE Root.City.Country = "England"

```

At a high level, the query execution engine will find all objects which satisfy the path *Root.City.Location* and for each of these will check whether the where clause is satisfied. The result of Query 4.1 is shown here:

```

answer complex {
  location string "Southern"
  location complex {
    longitude float -0.167
    latitude float 51.5 }
}

```

The database over which this query is evaluated presents a number of irregularities, as discussed earlier. A guiding principle in LOREL is that, to write a query, one should not

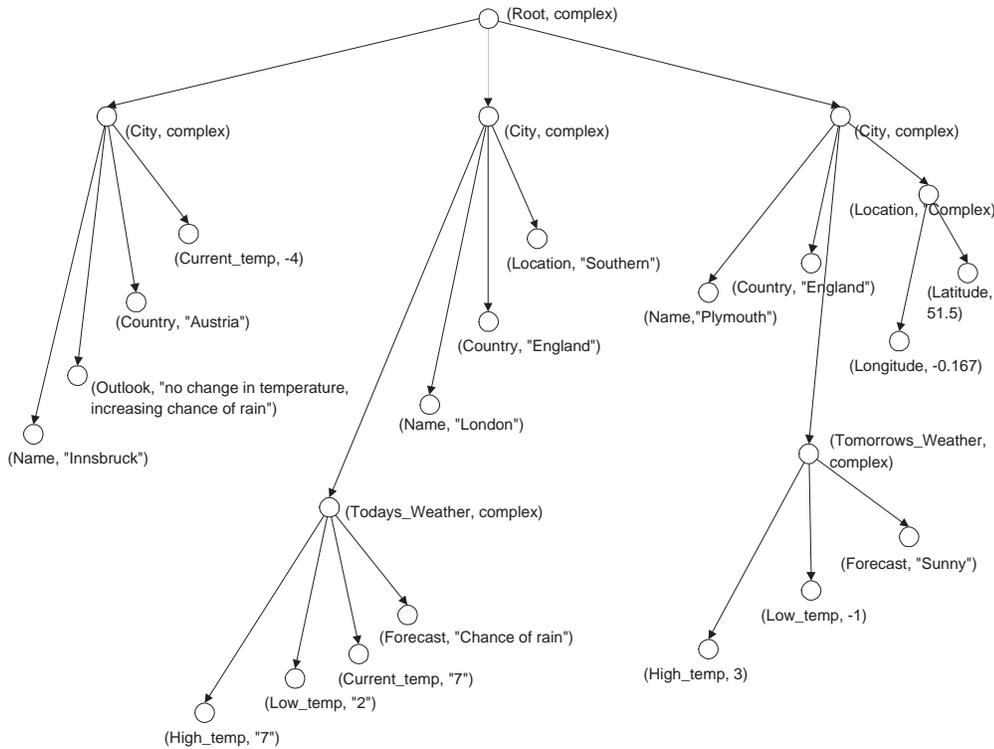


Figure 4: Sample OEM database

have to worry about such irregularities or know the *precise* structure of objects (e.g., the structure of location objects), nor should one have to bother with *precise types*. This query will not yield a run-time error if a *Location* object has a string value or is complex, or if *Country* objects are single-valued, set-valued, or even absent for some cities. Indeed, the above query will succeed no matter what the actual structure of the database is, and will return an appropriate answer. Of course, this query was written with some obvious knowledge of how the graph is laid out within our database. In Sec. 4.2 we discuss how an end user can discover the structure of the database.

Value comparisons are made after two objects have been coerced into comparable types. That is, if two objects do not have the same type then attempts will be made to coerce the values into comparable types before applying the comparison operator. Any types which cannot be coerced for comparison will not return type errors, but will simply evaluate to false. This reinforces our underlying principle that LOREL does not require precise knowledge of the data and is most useful when dealing with semistructured data.

The system will in fact translate all LOREL queries into OQL-like queries for evaluation. This is done for two reasons: first, LOREL is based on OQL and thus OQL gives us well defined semantics for our queries, and second it allows a user familiar with OQL to directly enter an OQL query to be evaluated over the semistructured data. In some sense, LOREL can be viewed as shorthand for OQL, however LOREL also introduces *generalized path expressions* not present within OQL. Generalized path expressions offer a richer form of “declarative navigation” in OEM databases than simple path expressions. Intuitively, the user loosely specifies a desired pattern of labels in the database: one can

specify patterns for paths (to match sequences of labels), patterns for labels (to match sequences of characters), and patterns for atomic values. A combination of these three forms of pattern matching is illustrated in the following example:

Query 4.2 (LOREL)

```
SELECT Root.City.Name
WHERE Root.City(.%weather)?(.Forecast|.Outlook)
      grep "rain"
```

Here the expression *%weather* is a label pattern that matches all labels ending with the string *weather* (e.g., *weather*, *Todays.weather*, or *Tomorrows.weather*). For path patterns, the symbol “|” indicates disjunction between two labels, and the symbol “?” is applied to the parenthesized expression to the left and indicates that the label pattern is optional. The complete syntax is based on regular expressions, along with syntactic wildcards such as “#”, which matches any path of length 0 or more. Finally, *grep* “rain” specifies that the data value should contain within it the string “rain”. The *grep* operator is similar to the Unix *grep* command. We also support *like*, based loosely on the SQL *like*, and *soundex* for phonetic matching. In English this query is asking for the names of all cities where the forecast (or outlook) of the weather contains the word “rain”. The result of Query 4.2 applied over the database in Figure 4 looks like this:

```
answer complex {
  name string "Innsbruck"
  name string "London"
}
```

During preprocessing, simple path expressions are eliminated by rewriting the query to use variables, as demon-

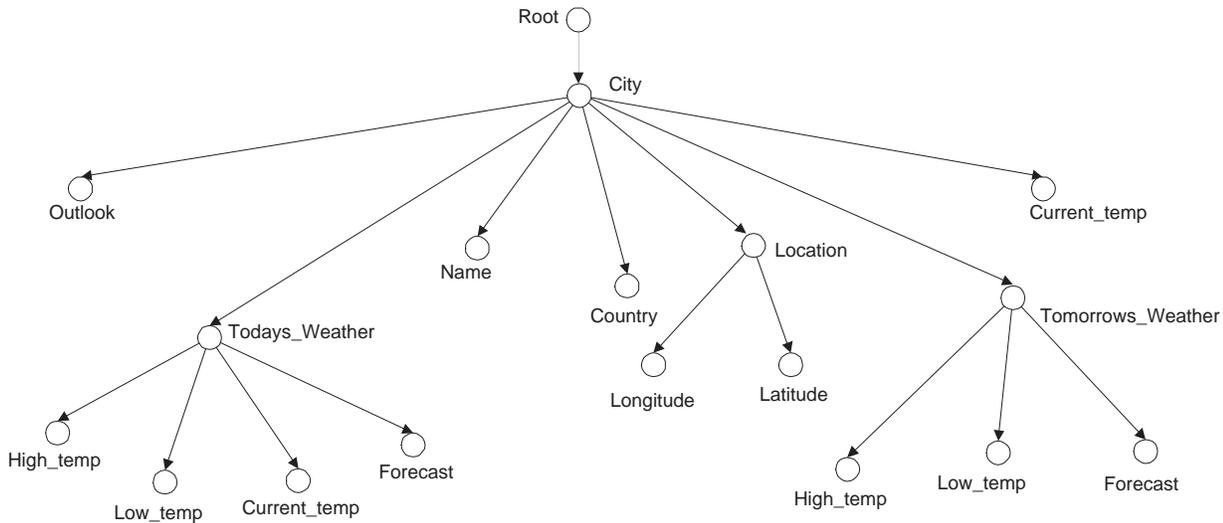


Figure 5: Sample DataGuide

strated in our first example. It is not possible to do so with general path expressions, which require a run-time mechanism. Indeed, note that if the database contains cycles, then a general path expression may match an infinite number of paths in the data. When trying to match a general path expression against the database, we match through a cycle at most once, which appears to be a reasonable simplification in practice.

We conclude with an example that illustrates advanced features of the language. The following query illustrates subqueries and constructed results. For every city in the database that satisfies the (bottom) where clause, we will select out the name of the city along with the current temperature, but only if the current temperature satisfies the WHERE clause.

Query 4.3 (LOREL)

```
SELECT C.Name,
  ( SELECT X
    FROM C.Todays_weather.current_temp X
    WHERE X < 10 )
FROM Root.City C
WHERE C.Country = "England"
```

The result is shown below. Notice that each city which provides a binding for the *C* variable and satisfies the where clause appears within the answer. Of particular interest is the fact that *Plymouth* does not have a *current_temp* field within the answer. This is filtered out as a result of the subquery appearing within the *SELECT* clause. Specifically, the *Plymouth* object does not have a subobject labeled *Todays_weather*.

```
answer complex {
  city complex {
    name string "London"
    current_temp integer 7 }
  city complex {
    Name string "Plymouth" }
}
```

4.2 Query Formulation with the DataGuide

Since our data does not have an explicit schema, query formulation and query optimization are particularly challenging. Without some knowledge of the structure of the underlying database, writing a meaningful LOREL query may be difficult, even when using general path expressions. One may manually browse a database to learn more about its structure, but this approach is unreasonable for very large databases. Further, without information about the structure of the database, the query processor may be forced to perform more work than necessary. For example, consider Query 4.1 that finds the locations of all cities whose country is England. Even if no cities have a country subobject, the execution engine would still needlessly examine every city in the database.

A *DataGuide* is a concise and accurate summary of the structure of an OEM database, stored itself as an OEM object. Each possible path expression of a database is encoded exactly once in the DataGuide, and the DataGuide has no path expressions that do not exist in the database. As an example Figure 5 shows a DataGuide for the sample database shown in Figure 4. (Note that atomic values are usually not stored within the leaf nodes of the DataGuide since it is primarily concerned with the structure of the database.) In typical situations, the DataGuide is significantly smaller than the original database. A DataGuide plays a role similar to metadata in traditional database systems. The DataGuide may be queried or browsed, enabling user interfaces or client applications to examine the structure of the database. Assuming the role of the missing schema, the DataGuide can also guide the query processor. Of course, in relational or object-oriented systems the schema is explicitly created before any data is loaded; here, DataGuides are dynamically generated and maintained over all or part of an existing database.

In [5], formal definitions for DataGuides are provided as well as algorithms to build and incrementally maintain DataGuides that support annotations. Also given is a discussion of how DataGuides aid query formulation in practice and their use for query optimization.

5 Browsing OEM Results through MOBIE

The main idea behind our browsing tool centers around the need for displaying semistructured objects in a way that makes it easy for the user to grasp their structure and explore their contents when viewing the result of a TSIMMIS query. OEM results are typically irregular in structure and nested, containing a top-level (root) object and zero or more subobjects (sometimes referred to as children). Each subobject may itself be a nested object. In general, nested objects are structured like trees (or graphs if we allow cycles). Anybody who has worked with nested objects before can attest to the fact it becomes increasingly difficult to understand the contents of a nested object the more its structure increases in complexity (i.e., the larger the number of subobjects and the deeper the level of nesting).

For this reason, we have built a system that transforms OEM results into a “web” of hyperlinked documents that can be viewed using any WWW browser. An object that is selected for viewing is formatted as an HTML document. If the object is a complex object, the document may also include hyperlinks pointing to some or all of the object’s substructure depending on the user’s preferences. If the object is atomic, it will be displayed by itself. In addition, each document always contains a link to the parent object, unless the selected object is the root of the entire structure. The main contribution of our system is that it gives the user the option to decide which information is to be displayed, how much of the chosen information he or she wants to see, and when. Information is presented one screen at a time, allowing the user to browse complex objects, which may be too large to view all at once, in a “cafeteria-style” (pick-and-choose) fashion. This approach to browsing nested objects is analogous to how one uses the table of contents to explore the individual chapters of a book.

An important part of the functionality of our browser focuses on the layout of information on individual pages. Since this is a process that depends heavily on each user’s individual preferences as well as the data that is being displayed, we have paid careful attention to design a system that is flexible enough so that it can be tailored to satisfy many different needs. Our goal was to provide users with choices as to how information is to be displayed: from the overall layout of a screen down to the format of an individual object. Initially, the system uses default settings that maximize the amount of information that can be displayed within the given real-estate of the window. The result can then be improved upon by changing the values of *session variables*, which control the document layout, the level of nesting per screen, the number of subobjects per level, etc. By default, session variables control the formatting for the complete object hierarchy. However, by using the label names that refer to a particular object in the hierarchy, the scope of session variables can be limited: from the entire hierarchy, to a specific substructure, to one object. Although customization of the object display may be time consuming in certain cases, the state of the session variables can be saved on a per-user basis and re-used during subsequent sessions.

We have implemented a fully functional prototype system called MOBIE (Multimedia Object and Information Explorer), which currently provides the graphical interface to TSIMMIS data sources. However, MOBIE is not limited to browsing only data from TSIMMIS but can be tailored for displaying and formatting structured information from any object-based database/footnoteOne can either use a translator for converting data into OEM or modify our algorithm

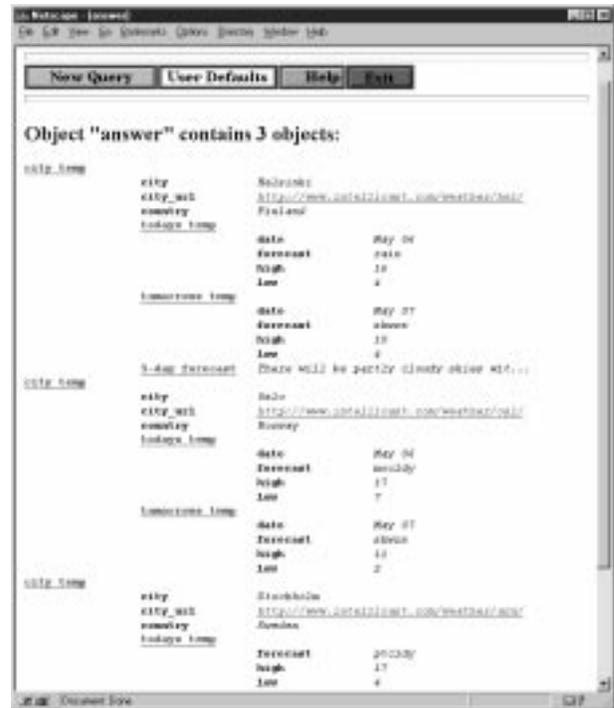


Figure 6: Result to Query 5.1

to work with other object-based data models.. Since a complete description of our browser is beyond the scope of this paper, we invite the user to obtain the details from [6]. Instead, we will briefly demonstrate some of MOBIE’s functionality using screen snapshots from a sample interaction with a TSIMMIS wrapper connected to the Intellicast weather source (via the above mentioned Web extractor). We start our description when the result is returned from the database, omitting such details as how to connect to the database server, transmission of the query and its results, etc. When displaying data, we use the following conventions. Object labels are displayed in **bold**, object values are *italicized*. Underlining indicates the existence of a hyperlink.

5.1 Sample Screen Snapshot

Let us assume that we have submitted the following LOREL query asking for all cities in Europe where tomorrow’s forecast calls for showers:

Query 5.1 (LOREL)

```
SELECT city_temp
FROM intellicast:i
WHERE i.city_temp.tomorrow's_temp.forecast = "shwrs"
```

Let us also assume that the answer to this query consists of three cities that are displayed together under one root object, labeled **answer**. Figure 6 shows the **answer** object as it is displayed in MOBIE. Each object labeled **city_temp** is a complex object exhibiting additional substructure underneath: the objects labeled **city**, **city_url**, **country**, **today's temp**, and **tomorrow's temp**. Note that the first subobject (the city of Helsinki) has one additional subobject labeled **5-day forecast** that is not present in the other results. The **city**, **city_url**, and **country** subobjects are *atomic* meaning they contain no further substructure. In those cases, the value of the object is displayed. (If there

Object Label	Object Value	
	Atomic Value	Child Object Label
city	Helsinki	
city_url	http://www.turknet.com/eng/fin/hel/	
country	Finland	
todays_temp	May 06	rate 16 f
tomorrows_temp	May 07	clear 10 f
5-day_forecast	There will be partly cloudy after wit...	

Figure 7: Query result—subobject “Helsinki”

is not enough room for the value, a hyperlink is provided.) The `todays_temp` and `city_url` subobjects on the other hand, are *complex* objects that contain additional subobjects: `forecast`, `high`, `low`, and `date`. Labels belonging to complex objects are underlined meaning that a hyperlink exists that will take the user to the document containing only those subobjects. (Those subobjects are displayed in a similar fashion.) Also note, the value of the `city_url` subobjects is a standard URL that is part of the answer and has been activated by MOBIE for loading.

5.2 Formatting Options

As mentioned before, the user can control the formatting of objects through various control parameters. These parameters are called *session variables* and can be accessed from the **User Defaults** menu. Formatting options fall into two categories: “Global Settings”, which apply to the whole object structure, and “Label-Based Settings” for which the scope can be specified based on object labels. (See [6] for details and other options.)

The following parameters are available for controlling global settings:

- *Maximum levels of sub-objects* controls the number of visible levels of subobjects for each object that is displayed.
- *Sub-object indentation* controls the amount of indentation used for subobjects.

The following parameters are available for controlling label-based settings:

- *Layout* controls the overall “look-and-feel” of the output when it gets displayed in the browser window. The two options currently available are *table* and *list* layout.
- *Number of displayed sub-objects* controls the number of subobjects that are displayed on a screen.
- *Label size* and *value size* control the length of labels and values respectively.

Using these options one can format data in the way that best suits it. For example, Figure 7 shows some data formatted as a table. Labels are shown on the left side. If the subobject is an atomic object (e.g., the subobjects labeled `city`, `city_url`, `country`, and `5-day_forecast`) the first column starting from the left will contain the subobject’s value. If the subobject is a complex object, e.g., the subobjects labeled `todays_temp` and `tomorrows_temp`, the first column will be empty, and subsequent columns will contain the values of its immediate subobjects. In the latter case, the column headings are the labels of the lower-level subobjects. Note, if there are several complex subobjects with different substructure, the table will display the union of all possible headings.

As mentioned before, label-based settings apply to objects. In order to format an object, a formatting choice associated with its label must be defined. Thus it is possible, for example, to display three or more levels of nesting for the root object, and then reduce the number of visible levels to just one when viewing its subobjects. As another example, one can display the part of a result that contains numerical values as a table but leave the part that is mostly textual in list format.

6 Conclusion

In this paper we have presented an overview of the TSIMMIS approach to accessing and managing semistructured data. In particular, we have described how semistructured data can be obtained from Web pages, how it can be manipulated in a database system, and how it can be browsed. We believe that semistructured data exists in many applications, and flexible tools like the ones we have described can be very helpful for managing it.

References

- [1] S. Abiteboul, D. Quass, J. McHugh, J. Widom, and J. Wiener. The Lorel query language for semistructured data. *Journal of Digital Libraries*, 1(1), November 1996.
- [2] C. Batini, M. Lenzerini, and S. Navathe. A comparative analysis of methodologies for database schema integration. *ACM Computing Surveys*, 18(4):323–364, 1986.
- [3] M.J. Carey, L.M. Haas, P.M. Schwarz, M. Arya, W.F. Cody, R. Fagin, A. Flickner, A.W. Luniewski, W. Niblack, D. Petkovic, J. Thomas, J. H. Williams, and E.L. Wimmers. Towards heterogeneous multimedia information systems: the Garlic approach. In *In Proceedings of the Sixth International Conference on Data Engineering*, pages 123–130, Los Angeles, California, February 1995.
- [4] S. Chawathe, H. Garcia-Molina, J. Hammer, K. Ireland, Y. Papakonstantinou, J. Ullman, and J. Widom. The TSIMMIS project: Integration of heterogeneous information sources. In *Proceedings of the Tenth Anniversary Meeting*, pages 7–18. Information Processing Society of Japan, Tokyo, Japan, October 1994.
- [5] R. Goldman and J. Widom. DataGuides: Enabling query formulation and optimization in semistructured databases. In *Proceedings of the Twenty-Third International Conference on Very Large Database*, Athens, Greece, September 1997.

- [6] J. Hammer, R. Aranha, and K. Ireland. Browsing object databases through the Web. Technical report, Department of Computer Science, Stanford, California, February 1997.
- [7] J. Hammer, M. Breunig, H. Garcia-Molina, S. Nestorov, V. Vassalos, and R. Yerneni. Template-based wrappers in the TSIMMIS system. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, page 532, Tucson, Arizona, May 1997. Association of Computing Machinery.
- [8] J. Hammer, H. Garcia-Molina, Y. Cho, R. Aranha, and A. Crespo. Extracting semistructured information from the Web. In *Proceedings of the First Workshop on Management of Semistructured Data*, pages 18–25, Tucson, Arizona, May 1997.
- [9] T. Kirk, A. Levy, J. Sagiv, and D. Srivastava. The information manifold. Technical report, AT&T Bell Laboratories, 1995.
- [10] J. McHugh, S. Abiteboul, R. Goldman, D. Quass, and J. Widom. LORE: A database management system for semistructured data. *SIGMOD Record*, 26(3):50–61, 1997.
- [11] Y. Papakonstantinou, H. Garcia-Molina, and S. Abiteboul. Object fusion in mediator systems. In *Proceedings of the International Conference on Very Large Databases*, pages 234–245, Bombay, India, September 1996.
- [12] Y. Papakonstantinou, H. Garcia-Molina, and J. Widom. Object exchange across heterogeneous information sources. In *Proceedings of the Eleventh International Conference on Data Engineering*, pages 251–260. Computer Society of the IEEE, Taipei, Taiwan, March 1995.
- [13] Y. Papakonstantinou, A. Gupta, H. Garcia-Molina, and J. Ullman. A query translation scheme for rapid implementation of wrappers. In *International Conference on Deductive and Object-Oriented Databases*, pages 97–107, August 1995.
- [14] K. Shoens, A. Luniewski, P. Schwarz, J. Stamos, and J. Thomas. The RUFUS system: Information organization for semi-structured data. In *Proceedings of the International Conference on Very Large Databases*, pages 97–107, Dublin, Ireland, August 1993.
- [15] D. Suciu. Proceedings of the workshop on management of semistructured data. Tucson, Arizona, May 1997. Los Angeles. (Workshop papers are available electronically at <http://www.research.att.com/~suciu/workshop-papers.html>.)