

DSA5104 Assignment1

Wang Jiangyi, A0236307J
National University of Singapore

1 Task 2: Schema Diagram

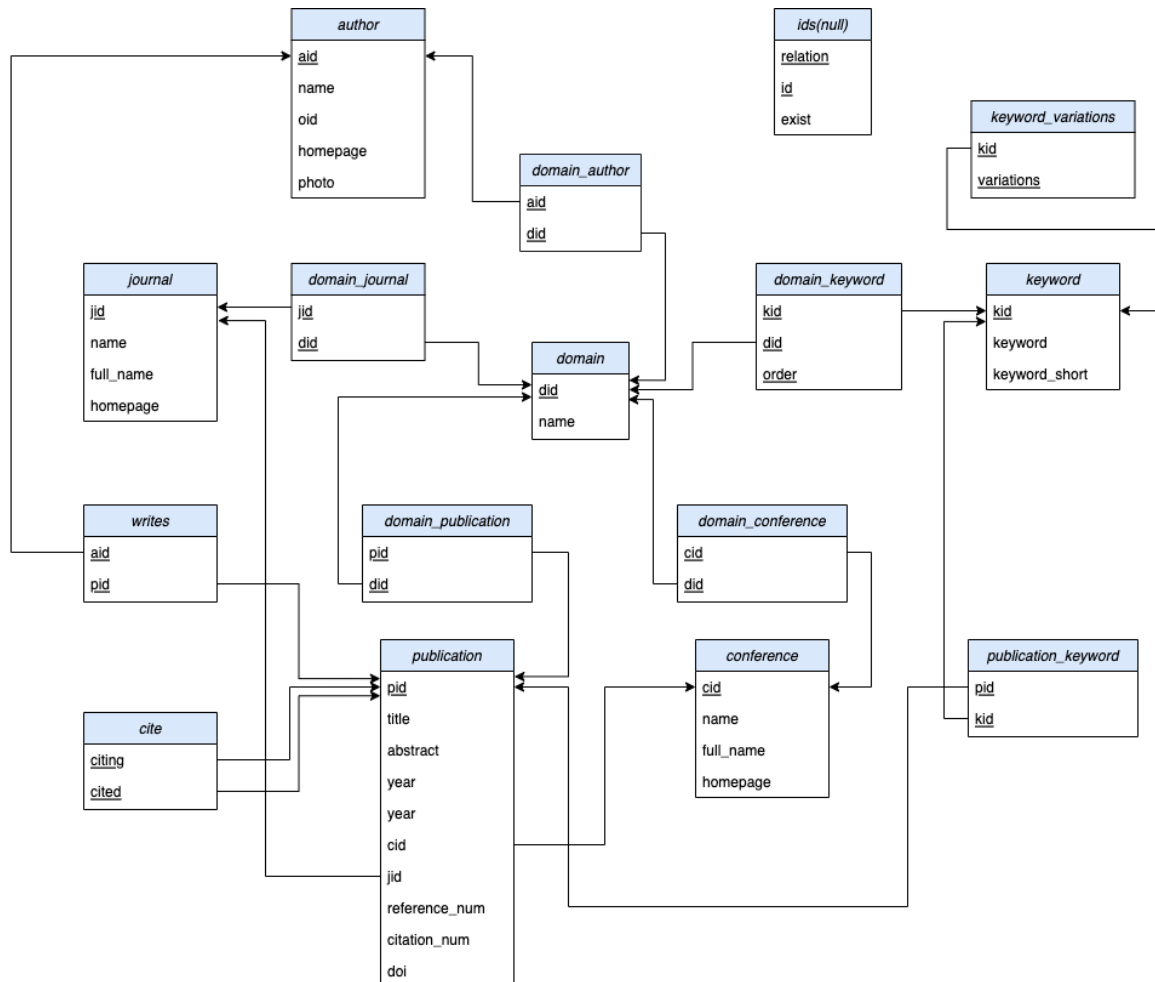


Figure 1: Relational Diagram for MAS

2 Task 3: NLQueries

In the beginning, I pointed out that, for Query 7 and Query 13, I have different understanding of NL Query. Therefore, I give all the answers corresponding to different understanding.

2.1 Query 1

NL Query: Return me the authors who have papers in PVLDB.

SQL Query: (2 methods)

```
1 # Q1: Return me the authors who have papers in PVLDB.
2
3 ## Method 1:
4
5 WITH pid_in_PVLDB AS (
6     SELECT pid
7     FROM publication
8     WHERE jid = (
9         SELECT jid
10        FROM journal
11        WHERE name = 'PVLDB'
12    )
13 )
14 SELECT a.name AS author_name
15 FROM author a
16 WHERE EXISTS (
17     SELECT pid
18     FROM pid_in_PVLDB
19     WHERE pid IN (
20         SELECT b.pid
21         FROM writes b
22         WHERE b.aid = a.aid
23     )
24 );
25
26 ## Method 2:
27
28 WITH author_flag AS (
29     SELECT a.aid, b.name AS author_name
30     , CASE
31         WHEN a.pid IN (
32             SELECT pid
33             FROM publication
34             WHERE jid = (
35                 SELECT jid
36                 FROM journal
37                 WHERE name = 'PVLDB'
38             )
39         )
40         THEN 1
41         ELSE 0
42     END AS author_flag
43 FROM writes a
44     LEFT JOIN author b ON a.aid = b.aid
45 )
46 SELECT aid, author_name
47 FROM author_flag
48 GROUP BY aid, author_name
49 HAVING sum(author_flag) > 0;
```

Query Result: (1320 rows in total)

aid	author_name
1241	Chun Chen
1329	Philip S. Yu
4766	Yannis Sismanis
5061	Kyuseok Shim
8789	Jens Teubner
9074	Fernando C. Pereira
9728	Stefan Manegold
11196	Juliana Freire
11281	Minghua Chen
11314	Christian S. Jensen
13495	Anthony K. H. Tung
13599	Vivek Narasayya
16868	Phokion G. Kolaitis
20614	Stavros Harizopoulos
24944	Mark D. Hill
28766	Jian-Hua Feng
32465	Sabrina De Capitani d...
35950	Daniel J. Abadi

Figure 2: Result for query 1

2.2 Query 2

NL Query: Return me the organization H. V. Jagadish is in.

SQL Query:

```
1 # Q2: Return me the organization H. V. Jagadish is in.
2
3 SELECT a.name AS author_name, b.name AS organization_name
4 FROM (
5     SELECT oid, name
6     FROM author
7     WHERE name = 'H. V. Jagadish'
8 ) a
9     LEFT JOIN organization b ON a.oid = b.oid;
```

Query Result: (1 row in total)

author_name	organization_name
H. V. Jagadish	University of Michigan

Figure 3: Result for query 2

2.3 Query 3

NL Query: Return me the authors who have papers in VLDB conference before 2002 after 1995.

SQL Query:

```
1 # Q3: Return me the authors who have papers in VLDB conference before 2002 after 1995.
2
3 WITH pid_VLDB_1995_2002 AS (
4     SELECT pid
5     FROM publication
6     WHERE cid = (
7         SELECT cid
8         FROM conference
9         WHERE name = 'VLDB'
10     )
11     AND year > 1995
12     AND year < 2002
13 )
```

```

14 SELECT b.aid AS author_id, b.name AS author_name
15 FROM (
16     SELECT DISTINCT aid
17     FROM writes
18     WHERE pid IN (
19         SELECT pid
20         FROM pid_VLDB_1995_2002
21     )
22 ) a
23 LEFT JOIN author b ON a.aid = b.aid;

```

Query Result: (984 rows in total)

author_id	author_name
636544	Praveen Seshadri
1037685	Miron Livny
1480663	Raghu Ramakrishnan
61208	Ashish Kumar Gupta
850509	Sunita Sarawagi
939673	Jeffrey Naughton
2016347	Prasad M. Deshpande
2209729	Rakesh Agrawal
56151118	Sameet Agarwal

Figure 4: Result for query 3

2.4 Query 4

NL Query: Return me the authors who have cooperated both with "H. V. Jagadish" and "Divesh Srivastava".

SQL Query:

```

1 # Q4: Return me the authors who have cooperated both with "H. V. Jagadish" and "Divesh
2   Srivastava".
3 WITH coop_HVJ AS (
4     SELECT DISTINCT aid AS coop_HVJ_aid
5     FROM writes
6     WHERE pid IN (
7         SELECT pid AS pid_HVJ
8         FROM writes
9         WHERE aid = (
10             SELECT aid
11             FROM author
12             WHERE name = 'H. V. Jagadish'
13         )
14     )
15     AND aid <> (
16         SELECT aid
17         FROM author
18         WHERE name = 'H. V. Jagadish'
19     )
20 ),
21 coop_DS AS (
22     SELECT DISTINCT aid AS coop_DS_aid
23     FROM writes
24     WHERE pid IN (
25         SELECT pid AS pid_DS
26         FROM writes
27         WHERE aid = (
28             SELECT aid
29             FROM author
30             WHERE name = 'Divesh Srivastava'
31         )
32     )

```

```

33         AND aid <> (
34             SELECT aid
35             FROM author
36             WHERE name = 'Divesh Srivastava'
37         )
38     )
39 SELECT b.name AS coop_author_name
40 FROM (
41     # coop_HVJ_aid intersects coop_DS_aid
42
43     SELECT coop_HVJ_aid AS coop_aid
44     FROM coop_HVJ
45     WHERE coop_HVJ_aid IN (
46         SELECT coop_DS_aid
47         FROM coop_DS
48     )
49 ) a
50 LEFT JOIN author b ON a.coop_aid = b.aid;

```

Query Result: (48 rows in total)

coop_author_name
S Muthukrishnan
Panagiotis G. Ipeirotis
Luis Gravano
Nick Koudas
Lauri Pietarinen
Shurug Al-khalifa
Cong Yu
Andrew Nierman
Jignesh M. Patel
Laks Lakshmanan

Figure 5: Result for query 4

2.5 Query 5

NL Query: Return me the authors who have more papers on VLDB than ICDE.

SQL Query:

```

1 # Q5: Return me the authors who have more papers on VLDB than ICDE.
2
3 WITH VLDB AS (
4     SELECT a.aid AS author_id, count(publication_VLDB.pid) AS count_VLDB
5     FROM writes a
6         LEFT JOIN (
7             SELECT pid
8             FROM publication
9             WHERE cid = (
10                 SELECT cid
11                 FROM conference
12                 WHERE name = 'VLDB'
13             )
14         ) publication_VLDB
15     ON a.pid = publication_VLDB.pid
16     GROUP BY a.aid
17 ),
18 ICDE AS (
19     SELECT a.aid AS author_id, count(publication_ICDE.pid) AS count_ICDE
20     FROM writes a
21         LEFT JOIN (
22             SELECT pid
23             FROM publication

```

```

24         WHERE cid = (
25             SELECT cid
26             FROM conference
27             WHERE name = 'ICDE'
28         )
29         ) publication_ICDE
30         ON a.pid = publication_ICDE.pid
31     GROUP BY a.aid
32 )
33 SELECT b.name AS author_name
34 FROM (
35     SELECT a.author_id AS author_id, a.count_VLDB AS count_VLDB, b.count_ICDE AS
36         count_ICDE
37     FROM VLDB a
38     LEFT JOIN ICDE b ON a.author_id = b.author_id
39     WHERE a.count_VLDB > b.count_ICDE
40 ) a
41 LEFT JOIN author b ON a.author_id = b.aid;

```

Query Result: (2627 rows in total)

author_name
Debra E. Vandermeer
Yannis Sismanis
Michael J. Lopez
Kyuseok Shim
Hans-dieter Ehrich
David Botzer
Jurgen Wasch
Ricardo A. Baeza-yates
Shlomo Moran
Peter M. Schwarz

Figure 6: Result for query 5

2.6 Query 6

NL Query: Return me the authors who have cited papers of H. V. Jagadish.

SQL Query: (2 methods)

```

1 # Q6: Return me the authors who have cited papers of H. V. Jagadish
2
3 WITH citing_HVJ AS (
4     SELECT citing
5     FROM cite
6     WHERE cited IN (
7         SELECT pid
8         FROM writes
9         WHERE aid = (
10             SELECT aid
11             FROM author
12             WHERE name = 'H. V. Jagadish'
13         )
14     )
15 )
16 SELECT DISTINCT b.aid AS author_id, b.name AS author_name
17 FROM citing_HVJ a
18     LEFT JOIN (
19         SELECT a1.pid, a2.aid, a2.name
20         FROM writes a1
21         LEFT JOIN author a2 ON a1.aid = a2.aid
22     ) b
23     ON a.citing = b.pid
24 WHERE b.aid IS NOT NULL;

```

```

25
26 # method 2
27
28 WITH citing_HVJ AS (
29     SELECT citing
30     FROM cite
31     WHERE cited IN (
32         SELECT pid
33         FROM writes
34         WHERE aid = (
35             SELECT aid
36             FROM author
37             WHERE name = 'H. V. Jagadish'
38         )
39     )
40 )
41 SELECT aid AS author_id, name AS author_name
42 FROM author
43 WHERE aid IN (
44     SELECT aid
45     FROM writes
46     WHERE pid IN (
47         SELECT citing
48         FROM citing_HVJ
49     )
50 );

```

Query Result: (7245 rows in total)

author_id	author_name
1241	Chun Chen
1242	Chun-Yi Shi
1329	Philip S. Yu
1486	Lukas Relly
2087	Rossana Maria de Ca...
2945	Clara Pizzuti
3253	Philip W. Trinder
3503	Young-whun Lee
3662	Debra E. Vandermeer
3715	Christian Capelle

Figure 7: Result for query 6

2.7 Query 7

NL Query: Return me all the papers, which contain the keyword "Natural Language".

Note: There are two different understanding of this NL Query. One is, the keyword contains "Natural Language" and the other is, the keyword is exactly "Natural Language". The 2 different SQL Queries can be shown as follows:

SQL Query 1, exactly 'Natural Language':

```

1 # Q7: Return me all the papers, which contain the keyword "Natural Language".
2
3 SELECT b.pid AS publication_id, b.title AS paper_name
4 FROM (
5     SELECT DISTINCT pid
6     FROM publication_keyword
7     WHERE kid IN (
8         SELECT kid
9         FROM keyword
10        WHERE keyword = 'Natural Language'
11    )
12 ) a
13 LEFT JOIN publication b ON a.pid = b.pid;

```

Query Result: (11232 rows in total)

publication_id	paper_name
38	Using Natural Language for Database Design
323	Exploiting lexical regularities in designing natural language systems
474	Reasoning about Information Change
511	Typed Logics With States
902	User-Needs Analysis and Design Methodology for an Automated Documentation Generator
1426	Making Systems Sensitive to the User's Time and Working Memory Constraints
1460	MultiDimensional User Models for Multimedia I/O in the Maintenance Consultant
1474	Natural language processing using a propositional semantic network with structured variables
1922	SUPPORTING FLEXIBILITY AND TRANSMUTABILITY: MULTI-AGENT PROCESSING AND ROLE-SWITCHING I...
2152	A generic algorithm for generating spoken monologues

Figure 8: Result for query 7, version 1

SQL Query 2, contains 'Natural Language':

```

1 # Q7: Return me all the papers, which contain the keyword "Natural Language".
2
3 SELECT b.pid AS publication_id, b.title AS paper_name
4 FROM (
5     SELECT DISTINCT pid
6     FROM publication_keyword
7     WHERE kid IN (
8         SELECT kid
9         FROM keyword
10        WHERE keyword LIKE '%Natural Language%'
11    )
12 ) a
13 LEFT JOIN publication b ON a.pid = b.pid;

```

Query Result: (19208 rows in total)

publication_id	paper_name
38	Using Natural Language for Database Design
323	Exploiting lexical regularities in designing natural language systems
474	Reasoning about Information Change
511	Typed Logics With States
692	Higher-Order Logic Programming as Constraint Logic Programming
807	Concept-Based Retrieval using Controlled Natural Language
834	Empirical learning of natural language processing tasks
902	User-Needs Analysis and Design Methodology for an Automated Documentation Generator
1416	Integrated Natural Language Generation Systems
1426	Making Systems Sensitive to the User's Time and Working Memory Constraints

Figure 9: Result for query 7, version 2

2.8 Query 8

NL Query: Return me all the researchers in database area in University of Michigan.

SQL Query:

```

1 # Q8: Return me all the researchers in database area in University of Michigan.
2
3 SELECT aid AS author_id, name AS author_name
4 FROM author
5 WHERE oid = (
6     SELECT oid
7     FROM organization
8     WHERE name = 'University of Michigan'
9 )
10 AND aid IN (
11     SELECT aid
12     FROM domain_author
13     WHERE did = (
14         SELECT did
15         FROM domain

```



```

16 WHERE name = 'Databases'
17 )
18 );

```

Query Result: (146 rows in total)

author_id	author_name
73839	Andrew Niernan
75475	Shuming Bao
105297	Wee Teck Ng
141842	Mark E. Deppe
144550	T. Ceccarelli
160200	Michael J. Cafarella
244284	Alan G. Merten
326823	Toby J. Teorey
360926	Brahim Medjahed
365133	Mark S. Ackerman

Figure 10: Result for query 8

2.9 Query 9

NL Query: Return me the number of papers written by H. V. Jagadish, Yunyao Li, and Cong Yu.

SQL Query:

```

1 # Q9: Return me the number of papers written by H. V. Jagadish, Yunyao Li, and Cong Yu.
2
3 WITH satisfied_papers AS (
4     SELECT a.pid
5     FROM (
6         SELECT pid AS pid
7              , CASE
8                  WHEN aid IN (
9                      SELECT aid
10                     FROM author
11                     WHERE name IN ('H. V. Jagadish', '
12                                Yunyao Li', 'Cong Yu')
13                  ) THEN 1
14                  ELSE 0
15              END AS flag
16         FROM writes
17     ) a
18     GROUP BY a.pid
19     HAVING sum(a.flag) = 3
20 )
21 SELECT count(pid) AS count_papers
22 FROM satisfied_papers;

```

Query Result: (1 row in total)

count_papers
3

Figure 11: Result for query 9

2.10 Query 10

NL Query: Return me the number of papers written by H. V. Jagadish in each year.

SQL Query:

```

1 # Q10: Return me the number of papers written by H. V. Jagadish in each year.
2
3 SELECT b.year AS year, count(b.pid) AS count_paper_each_year

```

```

4 FROM (
5     SELECT pid
6     FROM writes
7     WHERE aid = (
8         SELECT aid
9         FROM author
10        WHERE name = 'H. V. Jagadish'
11    )
12 ) a
13 LEFT JOIN publication b ON a.pid = b.pid
14 GROUP BY b.year
15 ORDER BY b.year ASC;

```

Query Result: (29 rows in total)

year	count_paper_each_y...
0	13
1984	1
1986	1
1987	7
1988	8
1989	14
1990	8
1991	8
1992	18
1993	8

Figure 12: Result for query 10

2.11 Query 11

NL Query: Return me the number of citations of "Making database systems usable" in each year.

SQL Query:

```

1 # Q11: Return me the number of citations of "Making database systems usable" in each year.
2
3 WITH paper_citing_MDSU AS (
4     SELECT citing AS citing_pid
5     FROM cite
6     WHERE cited = (
7         SELECT pid
8         FROM test.publication
9         WHERE title = 'Making database systems usable'
10    )
11 )
12 SELECT b.year AS year, count(b.pid) AS count_citation_each_year
13 FROM paper_citing_MDSU a
14 LEFT JOIN publication b ON a.citing_pid = b.pid
15 GROUP BY b.year
16 ORDER BY b.year ASC;

```

Query Result: (5 rows in total)

year	count_citation_each_year
0	1
2008	5
2009	11
2010	8
2011	4

Figure 13: Result for query 11

2.12 Query 12

NL Query: Return me the author who has the most number of papers in the VLDB conference.

SQL Query:

```

1 # Q12: Return me the author who has the most number of papers in the VLDB conference.
2
3 WITH count_paper_author_VLDB AS (
4     SELECT b.aid AS author_id, count(a.pid) AS count_paper_VLDB
5     FROM (
6         SELECT pid
7         FROM publication
8         WHERE cid = (
9             SELECT cid
10            FROM conference
11            WHERE name = 'VLDB'
12        )
13     ) a
14     LEFT JOIN writes b ON a.pid = b.pid
15     GROUP BY b.aid
16 )
17 SELECT b.name AS author_name, a.count_paper_VLDB
18 FROM (
19     SELECT author_id, count_paper_VLDB
20     FROM count_paper_author_VLDB
21     WHERE count_paper_VLDB = (
22         SELECT max(count_paper_VLDB)
23         FROM count_paper_author_VLDB
24     )
25 ) a
26 LEFT JOIN author b ON a.author_id = b.aid;

```

Query Result: (1 row in total)

author_name	count_paper_VLDB
H. V. Jagadish	37

Figure 14: Result for query 12

2.13 Query 13

NL Query: Return me the conferences, which have more than 60 papers containing keyword "Relational Database".

Note: There are two different understanding of this NL Query. One is, the keyword contains "Relational Database" and the other is, the keyword is exactly "Relational Database". The 2 different SQL Queries can be shown as follows:

SQL Query 1, exactly "Relational Database": (2 methods)

```

1 # Q13: Return me the conferences, which have more than 60 papers containing keyword "
2     Relational Database".
3
4 WITH paper_contain_RD AS (
5     SELECT DISTINCT pid
6     FROM publication_keyword
7     WHERE kid IN (
8         SELECT kid
9         FROM keyword
10        WHERE keyword = 'Relational Database'
11    )
12 )
13 SELECT a.name AS conference_name, count(b.pid) AS count_paper_in_conference
14 FROM conference a
15     LEFT JOIN (
16         SELECT pid, cid
17         FROM publication
18         WHERE pid IN (
19             SELECT *
20             FROM paper_contain_RD
21         )
22     ) b

```

```

22         ON a.cid = b.cid
23     GROUP BY a.name
24     HAVING count_paper_in_conference > 60
25     ORDER BY count_paper_in_conference DESC;
26
27 # method 2
28
29 WITH conference_id_count AS (
30     SELECT cid, count(pid) AS count_cid
31     FROM publication
32     WHERE pid IN (
33         SELECT pid
34         FROM publication_keyword
35         WHERE kid = (
36             SELECT kid
37             FROM keyword
38             WHERE keyword = 'Relational Database'
39         )
40     )
41     GROUP BY cid
42 )
43 SELECT name
44 FROM conference
45 WHERE cid IN (
46     SELECT cid
47     FROM conference_id_count
48     WHERE count_cid > 60
49 );

```

Query Result: (5 rows in total)

conference_name	count_paper_in_conference
ICDE	157
DEXA	122
SIGMOD	97
VLDB	93
ER(OOER)	65

Figure 15: Result for query 13, version 1

SQL Query 2, contains "Relational Database": (2 methods)

```

1 # Q13: Return me the conferences , which have more than 60 papers containing keyword "
2     Relational Database".
3
4 WITH paper_contain_RD AS (
5     SELECT DISTINCT pid
6     FROM publication_keyword
7     WHERE kid IN (
8         SELECT kid
9         FROM keyword
10        WHERE keyword LIKE '%Relational Database%'
11    )
12 )
13 SELECT a.name AS conference_name, count(b.pid) AS count_paper_in_conference
14 FROM conference a
15     LEFT JOIN (
16         SELECT pid, cid
17         FROM publication
18         WHERE pid IN (
19             SELECT *
20             FROM paper_contain_RD
21         )
22     ) b
23     ON a.cid = b.cid
24 GROUP BY a.name
25 HAVING count_paper_in_conference > 60
26 ORDER BY count_paper_in_conference DESC;
27 # method 2

```

```

28 WITH conference_id_count AS (
29     SELECT cid, count(pid) AS count_cid
30     FROM publication
31     WHERE pid IN (
32         SELECT pid
33         FROM publication_keyword
34         WHERE kid IN (
35             SELECT kid
36             FROM keyword
37             WHERE keyword LIKE '%Relational Database%'
38         )
39     )
40     GROUP BY cid
41 )
42 SELECT name
43 FROM conference
44 WHERE cid IN (
45     SELECT cid
46     FROM conference_id_count
47     WHERE count_cid > 60
48 );

```

Query Result: (7 rows in total)

conference_name	count_paper_in_conference
ICDE	212
DEXA	156
VLDB	143
SIGMOD	133
ER(OOER)	69
PODS	67
SAC	62

Figure 16: Result for query 13, version 2

2.14 Query 14

NL Query: Return me the number of papers published on PVLDB each year after 2000.

SQL Query:

```

1 # Q14: Return me the number of papers published on PVLDB each year after 2000.
2
3 SELECT year AS year, count(pid) AS count_paper_each_year
4 FROM publication
5 WHERE jid = (
6     SELECT jid
7     FROM journal
8     WHERE name = 'PVLDB'
9 )
10 AND year > 2000
11 GROUP BY year
12 ORDER BY year ASC;

```

Query Result: (5 rows in total)

year	count_paper_each_y...
2008	162
2009	160
2010	189
2011	29
2014	1

Figure 17: Result for query 14

2.15 Query 15

NL Query: Return me the paper after 2000 in VLDB conference with the most citations.

SQL Query:

```
1 # Q15: Return me the paper after 2000 in VLDB conference with the most citations
2
3 WITH VLDB_paper_after_2000_citation AS (
4     SELECT pid AS publication_id, title AS publication_title, citation_num AS
5         count_citation
6     FROM publication
7     WHERE cid = (
8         SELECT cid
9         FROM conference
10        WHERE name = 'VLDB'
11    )
12    AND year > 2000
13 )
14 SELECT *
15 FROM VLDB_paper_after_2000_citation
16 WHERE count_citation = (
17     SELECT max(count_citation)
18     FROM VLDB_paper_after_2000_citation
19 );
```

Query Result: (1 row in total)

publication_id	publication_title	count_citation
133771	Generic Schema Matching with Cupid	798

Figure 18: Result for query 15