

Helena Kazenski - CS333 Final Project

Project Overview:

This hanzi database organizes foundational components of the Chinese language, such as: words, characters (hanzi), radicals, example sentences, semantic fields, and word-to-word relationships. The goal of this project is to model how the Chinese language is structured from the ground up— from characters and the radicals they are built from, to words and the characters they contain, and then connecting those words with meaningful example sentences and semantic fields.

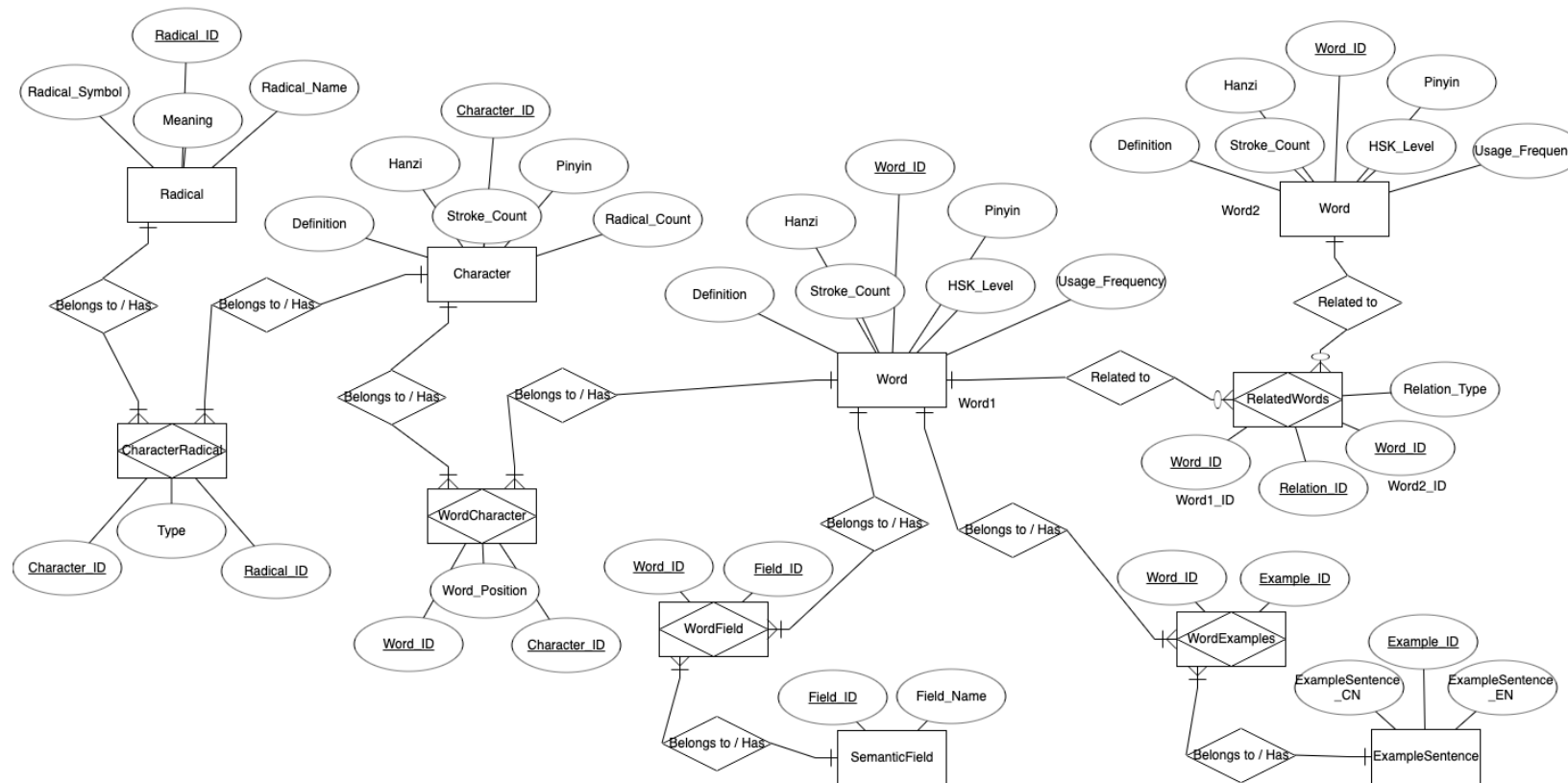
This database has the potential to be useful for language learners, educators, and anyone studying the morphology and semantics of Chinese. It builds the foundations of what could eventually be expanded into a larger-scope project — one that could grow beyond dictionary functions. It could even be a foundation for an educational app that molds itself to best fit the users' needs, or a deeper linguistics analysis tool, or an NLP project for machine translation and/or speech recognition.

Complexity:

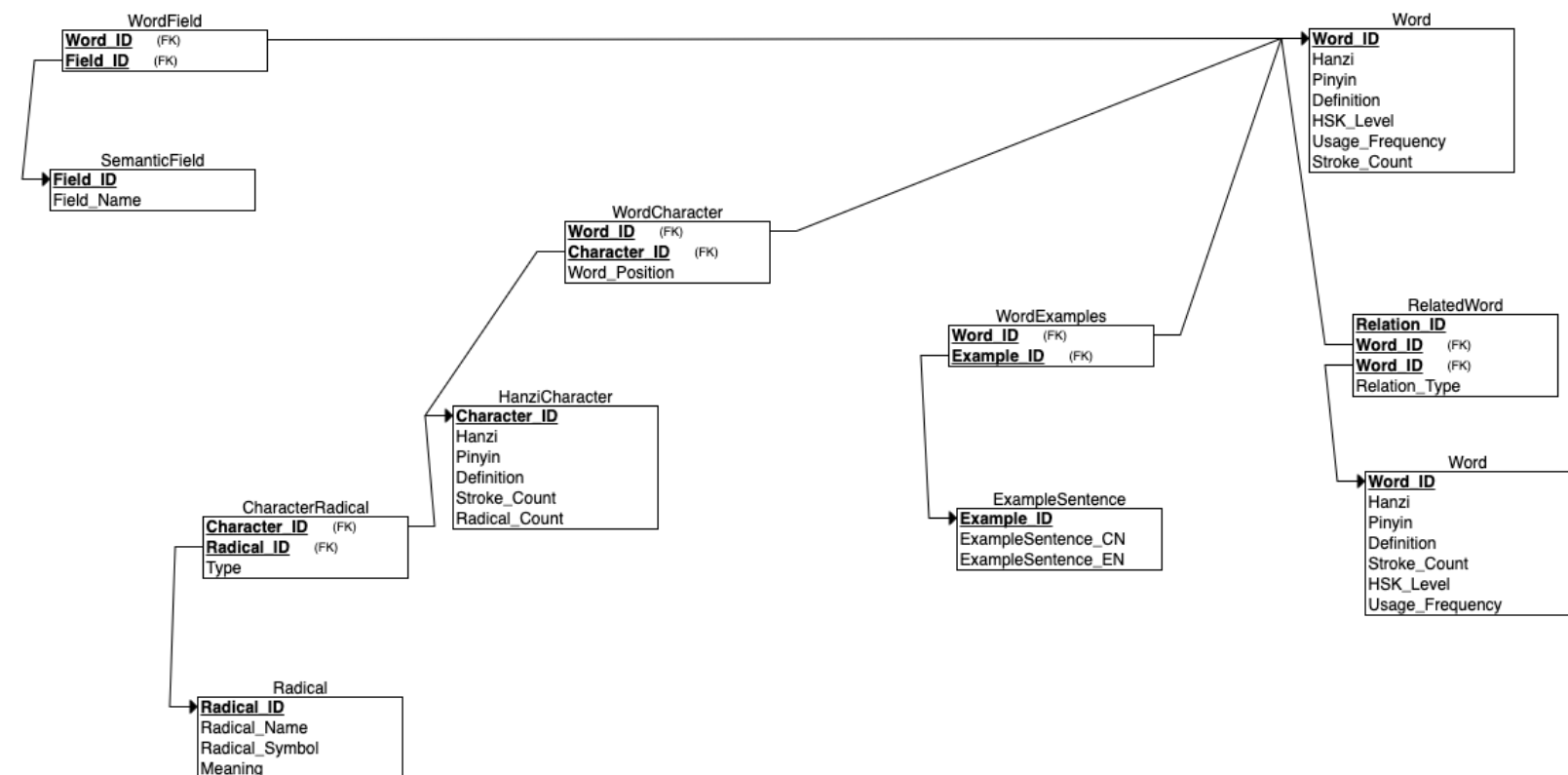
This project has the appropriate technical depth for a 300-level course because it includes multiple connected tables that model real relationships that occur in the Chinese language. Setting up tables and the connecting tables that tackle many-to-many relationships required careful planning to ensure that the data stayed organized and avoided duplication.

The queries show sufficient complexity and showcase skill in using SQL. They use multi-table joins, grouping, aggregates like AVG and COUNT, and string functions like GROUP_CONCAT to analyze patterns in the data, such as shared radicals, average stroke counts, and how words are used in sentences. Because the database supports linguistic breakdowns and semantic analysis, it demonstrates the level of design thinking and SQL skills expected for a 300-level course.

ERD:

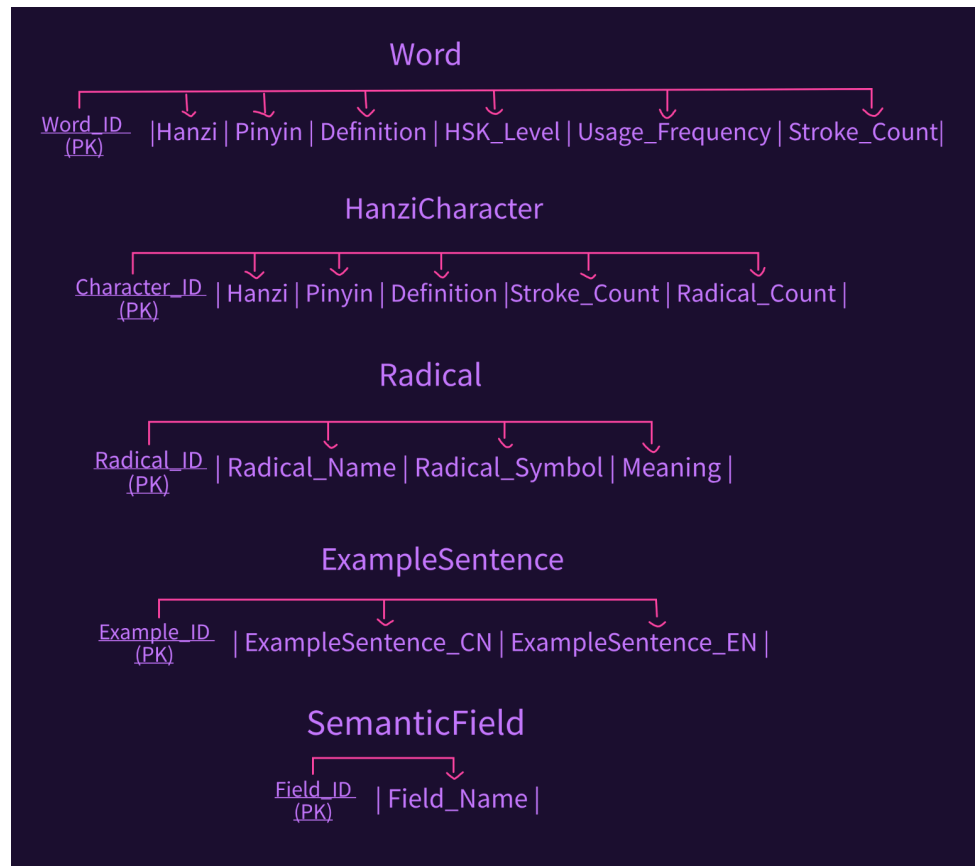


RS:

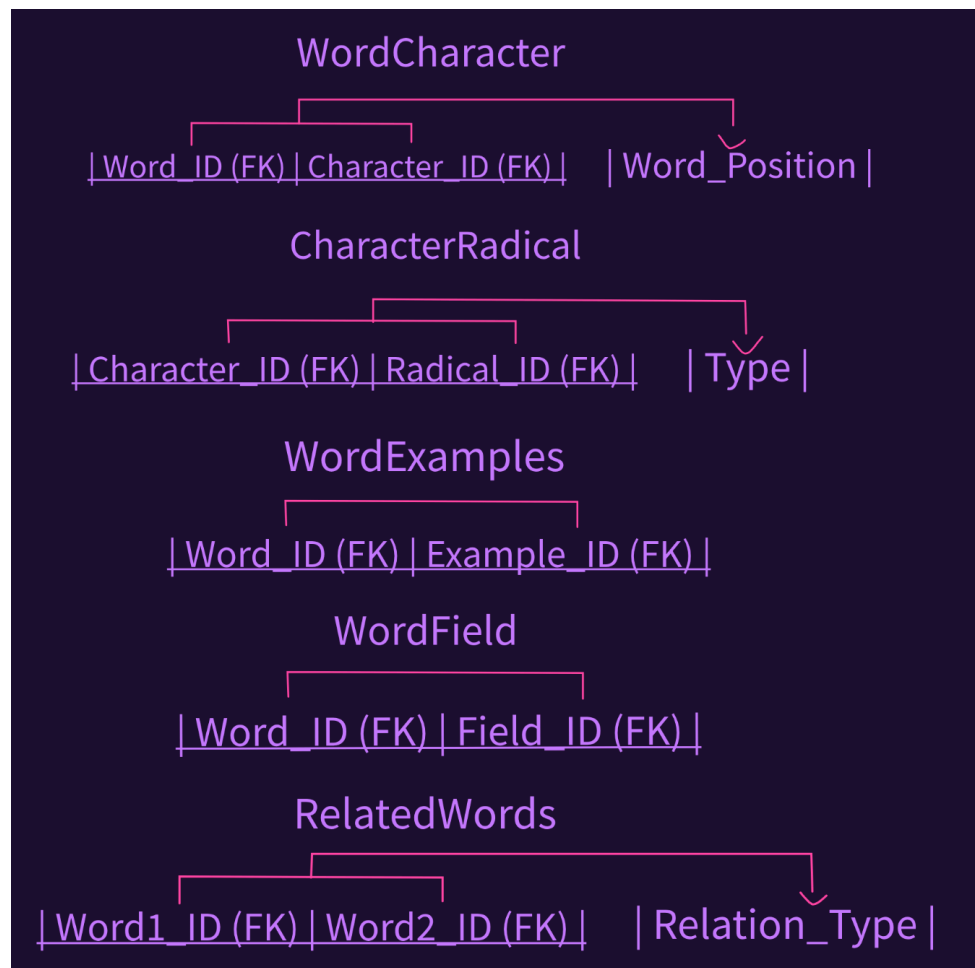


Dependency Diagrams:

Entities:



Associative Tables:



Dependency:

Yes, my database is 3NF. It doesn't have any repeating groups, so it is 1NF. It doesn't have any partial dependencies either — even in the associative tables, all non-key attributes depend on the full composite key — so it is 2NF. Finally, it doesn't have any transitive dependencies because every attribute in every table depends directly on the table's primary key (or composite key), so it is 3NF.

Queries:

```
27  -- Query 1
28  -- Select all characters with pinyin ending in 'i'
29  SELECT *
30  FROM HanziCharacter
31  WHERE Pinyin LIKE '%i'
32  ;
```

Character_ID	Hanzi	Pinyin	Definition	Stroke_Count	Radical_Count
1	你	nǐ	you	7	2
4	是	shì	to be	9	3
13	谁	shéi	who	10	2
18	几	jǐ	how many	2	1
21	一	yī	one	1	1
24	四	sì	four	5	2
27	七	qī	seven	2	1
30	十	shí	ten	2	1
32	岁	suì	years of age	6	2
35	块	kuài	piece; lump; money	7	3
37	没	méi	not; don't have	7	2
39	太	tài	too; very	4	2
42	在	zài	at; in; exist	6	2
47	喂	wèi	hello (on the phone)	12	3
54	医	yī	medicine	7	2
58	北	běi	north	5	2
65	里	lǐ	inside	7	2
75	期	qī	period; week	12	3
79	时	shí	time	10	3
83	子	zǐ	child	3	1
85	师	shī	teacher	10	3
90	衣	yī	clothes	6	2
92	米	mǐ	rice	6	2
96	杯	bēi	cup	8	2
98	飞	fēi	to fly	3	1
99	机	jī	machine; airplane	6	2
104	视	shì	vision	11	2
107	气	qì	air; gas	4	2
112	西	xī	west; thing	6	3
114	字	zì	character; word	6	2
119	椅	yǐ	chair	12	4
122	再	zài	again	6	2
125	对	duì	correct; toward	5	2
126	起	qǐ	to rise	10	3
134	来	lái	to come	7	1
135	回	huí	to return	6	2
137	吃	chī	to eat	6	2
139	睡	shuì	to sleep	13	3
142	买	mǎi	to buy	6	2
143	开	kāi	to open; drive	4	1
146	爱	ài	love	10	3
147	喜	xǐ	happy	12	3
151	识	shí	knowledge	7	2
152	会	huì	can; able to	6	2
162	丽	lì	beautiful	7	2

45 rows in set (0.00 sec)

```

34  -- Query 2
35  -- Count the number of characters with more than two radicals
36  SELECT Count(Character_ID) AS Num_Over_2_Radicals
37  FROM HanziCharacter
38  WHERE Radical_Count > 2
39  ;

```

```

+-----+
| Num_Over_2_Radicals |
+-----+
|                      |
+-----+
1 row in set (0.00 sec)

```

```

41  -- Query 3
42  -- Select the first 20 words and all the radicals used in each character (w/ no duplicate radicals in a word)
43  SELECT Word.Word_ID, Word.Hanzi AS Word_Hanzi,
44         GROUP_CONCAT(DISTINCT Radical_Symbol ORDER BY Radical_Symbol) AS Radicals_Used
45  FROM Word
46  JOIN WordCharacter ON Word.Word_ID = WordCharacter.Word_ID
47  JOIN CharacterRadical ON WordCharacter.Character_ID = CharacterRadical.Character_ID
48  JOIN Radical ON CharacterRadical.Radical_ID = Radical.Radical_ID
49  GROUP BY Word.Word_ID, Word.Hanzi
50  ORDER BY Word.Word_ID
51  LIMIT 20
52  ;

```

```

+-----+-----+-----+
| Word_ID | Word_Hanzi | Radicals_Used |
+-----+-----+-----+
| 1 | 你 | イ,刀,小 |
| 2 | 我 | 戈,扌 |
| 3 | 你们 | イ,刀,小,冂 |
| 4 | 他 | 亻,匕,イ |
| 5 | 她 | 亻,匕,女 |
| 6 | 它 | 匕,宀 |
| 7 | 我们 | イ,戈,扌,冂 |
| 8 | 他们 | 亻,匕,イ,冂 |
| 9 | 她们 | 亻,匕,イ,女,冂 |
| 10 | 什么 | ノ,二,イ,厶 |
| 11 | 谁 | 讠,隹 |
| 12 | 哪儿 | 亻,二,儿,口,阝 |
| 13 | 怎么 | 一,丨,ノ,二,厶,心 |
| 14 | 多少 | ノ,夕,小 |
| 15 | 几 | 几 |
| 16 | 这 | 文,辶 |
| 17 | 那 | 亻,二,阝 |
| 18 | 一 | 一 |
| 19 | 二 | 二 |
| 20 | 三 | 一,二 |
+-----+-----+-----+
20 rows in set (0.01 sec)

```



```

54 -- Query 4
55 -- Count the number of example sentences that 20 words appear in
56 SELECT Word.Word_ID, Word.Hanzi, COUNT(WordExamples.Example_ID) AS Sentence_Count
57 FROM Word
58 JOIN WordExamples ON Word.Word_ID = WordExamples.Word_ID
59 GROUP BY Word.Word_ID, Word.Hanzi
60 ORDER BY Sentence_Count DESC
61 LIMIT 20
62 ;

```

Word_ID	Hanzi	Sentence_Count
2	我	9
30	是	6
4	他	5
1	你	5
28	个	4
5	她	4
7	我们	4
8	他们	3
29	岁	2
20	三	2
16	这	2
11	谁	1
17	那	1
6	它	1
9	她们	1
22	五	1
12	哪儿	1
14	多少	1
10	什么	1
25	八	1

20 rows in set (0.00 sec)

```

64 -- Query 5
65 -- Select all characters containing the selected 20 radicals
66 SELECT Radical.Radical_Symbol,
67        GROUP_CONCAT(HanziCharacter.Hanzi ORDER BY HanziCharacter.Hanzi) AS Characters
68 FROM Radical
69 JOIN CharacterRadical ON Radical.Radical_ID = CharacterRadical.Radical_ID
70 JOIN HanziCharacter ON CharacterRadical.Character_ID = HanziCharacter.Character_ID
71 GROUP BY Radical.Radical_Symbol
72 ORDER BY Radical.Radical_Symbol
73 LIMIT 20
74 ;

```

Radical_Symbol	Characters
一	一,七,三,不,喂,怎,是,本
丨	个,他,哪,块,她,怎,那
丶	太,的
㇀	商
ノ	不,么,九,在,少,怎
し	七,他,喂,她
丿	了
二	三,二,五,些,什,哪,怎,那,院
亠	六,商,校
人	个,块
イ	什,他,们,你,在
儿	儿,四,院
八	八,六,商,校
冂	商
冂	学
几	几
刀	你
力	五
勹	的
匕	些,呢,它

20 rows in set (0.00 sec)

```

76  -- Query 6
77  -- Select the top 10 most frequently used words
78  SELECT Word_ID, Hanzi, Definition, Usage_Frequency
79  FROM Word
80  ORDER BY Usage_Frequency DESC
81  LIMIT 10
82  ;

```

Word_ID	Hanzi	Definition	Usage_Frequency
24	七	seven	0.95
25	八	eight	0.95
21	四	four	0.95
18	一	one	0.95
23	六	six	0.95
19	二	two	0.95
26	九	nine	0.95
22	五	five	0.95
10	什么	what	0.95
20	三	three	0.95

10 rows in set (0.00 sec)


```

84  -- Query 7
85  -- Select the average stroke count of all hanzi
86  SELECT AVG(Stroke_Count) AS Avg_Stroke_Count_Characters
87  FROM HanziCharacter
88  ;
89
90  -- Query 7.5
91  -- Select the average stroke count of all radicals
92  SELECT AVG(HanziCharacter.Stroke_Count) AS Avg_Stroke_Count_Radicals
93  FROM Radical
94  JOIN CharacterRadical ON Radical.Radical_ID = CharacterRadical.Radical_ID
95  JOIN HanziCharacter ON CharacterRadical.Character_ID = HanziCharacter.Character_ID
96  ;

```

```

+-----+
| Avg_Stroke_Count_Characters |
+-----+
|                               6.9136 |
+-----+
1 row in set (0.00 sec)

+-----+
| Avg_Stroke_Count_Radicals |
+-----+
|                               6.9612 |
+-----+
1 row in set (0.00 sec)

```

```

98  -- Query 8
99  -- Select pairs of related words and their relationship type
100 SELECT w1.Hanzi AS Word1, w2.Hanzi AS Word2, RelatedWords.Relation_Type
101 FROM RelatedWords
102 JOIN Word w1 ON RelatedWords.Word1_ID = w1.Word_ID
103 JOIN Word w2 ON RelatedWords.Word2_ID = w2.Word_ID
104 ORDER BY RelatedWords.Relation_Type
105 ;

```

Word1	Word2	Relation_Type
吃	米饭	Action-Object
喝	茶	Action-Object
买	水果	Action-Object
大	小	Antonym
热	冷	Antonym
看	听	Complementary-Actions
这	那	Contrast
爸爸	妈妈	Family-Pair
儿子	女儿	Family-Pair
是	有	Functional-Pair
做	买	Functional-Pair
有	没	Negation
是	不	Negation-Pair
来	没	Negation-Pair
吃	去	Opposite-Actions
你	喝	Paired-Actions
你	你们	Plural-Form
我	我们	Plural-Form
他	他们	Plural-Form
她	她们	Plural-Form
什么	怎么	Question Words
哪儿	谁	Question Words
多少	几	Question Words
苹果	水果	Type-Of

24 rows in set (0.00 sec)

```

108 -- Query 9
109 -- Select the top 10 radicals used in the most characters
110 SELECT Radical.Radical_Symbol, COUNT(CharacterRadical.Character_ID) AS Num_Radical_Appearances
111 FROM Radical
112 JOIN CharacterRadical ON Radical.Radical_ID = CharacterRadical.Radical_ID
113 GROUP BY Radical.Radical_Symbol
114 ORDER BY Num_Radical_Appearances DESC
115 LIMIT 10
116 ;

```

Radical_Symbol	Num_Radical_Appearances
二	9
一	8
口	7
丨	7
ノ	6
イ	5
冫	4
八	4
乚	4
小	3

10 rows in set (0.00 sec)

```

119 -- Query 10
120 -- Select top 10 characters that appear in more than one word
121 SELECT HanziCharacter.Hanzi, COUNT(DISTINCT WordCharacter.Word_ID) AS Word_Count
122 FROM HanziCharacter
123 JOIN WordCharacter ON HanziCharacter.Character_ID = WordCharacter.Character_ID
124 GROUP BY HanziCharacter.Character_ID, HanziCharacter.Hanzi
125 HAVING COUNT(DISTINCT WordCharacter.Word_ID) > 1
126 ORDER BY Word_Count DESC
127 LIMIT 10
128 ;

```

Hanzi	Word_Count
们	4
儿	3
天	3
我	2
果	2
他	2
她	2
么	2
不	2
饭	2

10 rows in set (0.00 sec)