# The English Wikipedia Database
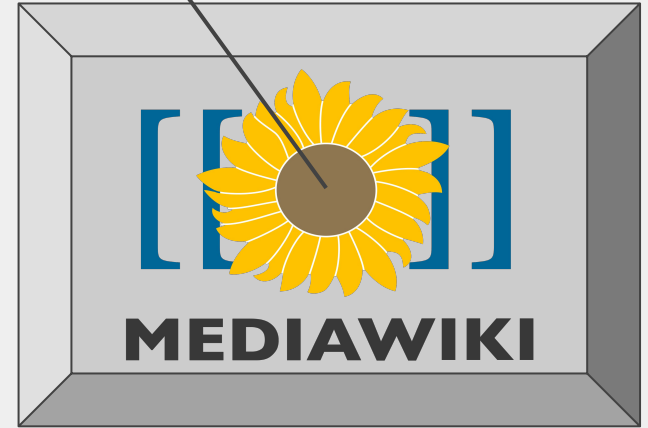
Paolo Tamagnini
Benedetta Checcarelli

WIKIPEDIA
The Free Encyclopedia

Local database of some tables of the whole schema,
containing most of the data of the today Wikipedia

Data content from
english Wikipedia

Database MySQL
infrastructure

**user** (partial)
- user_email TINYTEXT
- user_touched BINARY(14)
- user_token BINARY(32)
- user_email_authenticated BINARY(14)
- user_email_token BINARY(32)
- user_email_token_expires BINARY(14)
- user_registration BINARY(14)
- user_editcount INT
- user_password_expires VARBINARY(14)

**user_former_groups**
- ufg_user INT
- ufg_group VARBINARY(255)

**user_newtalk**
- user_id INT
- user_ip VARBINARY(40)
- user_last_timestamp VARBINARY(14)

**ipblocks** (partial)
- ipb_user INT
- ug_user INT
- ug_group VARBINARY(255)
- ipb_timestamp BINARY(14)
- ipb_auto
- ipb_anon_only
- ipb_create_account
- ipb_enable_autoblock
- ipb_expiry VARBINARY(14)
- ipb_range_start TINYBLOB
- ipb_range_end TINYBLOB
- ipb_deleted
- ipb_block_email
- ipb_allow_usertalk
- ipb_parent_block_id INT

**logging** (partial)
- log_namespace INT
- log_title VARCHAR(255)
- log_page INT
- log_comment VARCHAR(767)
- log_params BLOB
- log_deleted TINYINT

**change_tag**
- ct_id INT
- ct_rc_id INT
- ct_log_id INT
- ct_rev_id INT
- ct_tag VARCHAR(255)
- ct_params BLOB

**recentchanges** (partial)
- rc_comment VARCHAR(767)
- rc_minor TINYINT
- rc_bot TINYINT
- rc_new TINYINT
- rc_cur_id INT
- rc_this_oldid INT
- rc_last_oldid INT
- rc_type TINYINT
- rc_source VARCHAR(16)
- rc_patrolled TINYINT
- rc_ip VARBINARY(40)
- rc_old_len INT
- rc_new_len INT
- rc_deleted TINYINT
- rc_logid INT
- rc_log_type VARBINARY(255)
- rc_log_action VARBINARY(255)
- rc_params BLOB

## Pages

**archive**
- ar_id INT
- ar_namespace INT
- ar_title VARCHAR(255)
- ar_text MEDIUMBLOB
- ar_comment VARBINARY(767)
- ar_user INT
- ar_user_text VARCHAR(255)
- ar_timestamp BINARY(14)
- ar_minor_edit TINYINT
- ar_flags TINYBLOB
- ar_rev_id INT
- ar_text_id INT
- ar_deleted TINYINT
- ar_len INT
- ar_page_id INT
- ar_parent_id INT
- ar_sha1 VARBINARY(32)
- ar_content_model VARBINARY(32)
- ar_content_format VARBINARY(64)

**revision**
- rev_id INT
- rev_page INT
- rev_text_id INT
- rev_comment VARCHAR(767)
- rev_user INT
- rev_user_text VARCHAR(255)
- rev_timestamp BINARY(14)
- rev_minor_edit TINYINT
- rev_deleted TINYINT
- rev_len INT
- rev_parent_id INT
- rev_sha1 VARBINARY(32)
- rev_content_model VARBINARY(32)
- rev_content_format VARBINARY(64)

**text**
- old_id INT
- old_text MEDIUMBLOB
- old_flags TINYBLOB

**page**
- page_id INT
- page_namespace INT
- page_title VARCHAR(255)
- page_restrictions TINYBLOB
- page_is_redirect TINYINT
- page_is_new TINYINT
- page_random DOUBLE
- page_touched BINARY(14)
- page_links_updated VARBINARY(14)
- page_latest INT
- page_len INT
- page_content_model VARBINARY(32)
- page_lang VARBINARY(35)

**page_props**
- pp_page INT
- pp_propname VARBINARY(60)
- pp_value BLOB
- pp_sortkey FLOAT

**page_restrictions**
- pr_id INT
- pr_page INT
- pr_type VARBINARY(60)
- pr_level VARBINARY(60)
- pr_cascade TINYINT
- pr_user INT
- pr_expiry VARBINARY(14)

**redirect**
- rd_from INT
- rd_namespace INT
- rd_title VARCHAR(255)
- rd_interwiki VARCHAR(32)
- rd_fragment VARCHAR(255)

**category**
- cat_id INT
- cat_title VARCHAR(255)
- cat_pages INT
- cat_subcats INT
- cat_files INT

**protected_titles**
- pt_namespace INT
- pt_title VARCHAR(255)
- pt_user INT
- pt_reason VARCHAR(767)
- pt_timestamp BINARY(14)
- pt_expiry VARBINARY(14)
- pt_create_perm VARBINARY(60)

## Link tables

**pagelinks**
- pl_from INT
- pl_from_namespace INT
- pl_namespace INT
- pl_title VARCHAR(255)

**imagelinks**
- il_from INT
- il_from_namespace INT
- il_to VARCHAR(255)

**templatelinks**
- tl_from INT
- tl_from_namespace INT
- tl_namespace INT
- tl_title VARCHAR(255)

**iwlinks**
- iwl_from INT
- iwl_prefix VARBINARY(20)
- iwl_title VARBINARY(255)

**externallinks**
- el_id INT
- el_from INT
- el_to BLOB
- el_index BLOB

**langlinks**
- ll_from INT
- ll_lang VARBINARY(20)
- ll_title VARCHAR(255)

**categorylinks**
- cl_from INT
- cl_to VARCHAR(255)
- cl_sortkey VARBINARY(230)
- cl_sortkey_prefix VARCHAR(255)
- cl_timestamp TIMESTAMP
- cl_collation VARBINARY(32)
- cl_type ENUM(...)

## Multimedia

**image**
- img_name VARCHAR(255)
- img_size INT
- img_width INT
- img_height INT
- img_metadata MEDIUMBLOB
- img_bits INT
- img_media_type ENUM(...)
- img_major_mime ENUM(...)
- img_minor_mime VARBINARY(100)
- img_description VARCHAR(767)
- img_user INT
- img_user_text VARCHAR(255)
- img_timestamp VARBINARY(14)
- img_sha1 VARBINARY(32)

**filearchive**
- fa_id INT
- fa_name VARCHAR(255)
- fa_archive_name VARCHAR(255)
- fa_storage_group VARBINARY(16)
- fa_storage_key VARBINARY(64)
- fa_deleted_user INT
- fa_deleted_timestamp BINARY(14)
- fa_deleted_reason VARCHAR(767)
- fa_size INT
- fa_width INT
- fa_height INT
- fa_metadata MEDIUMBLOB
- fa_bits INT
- fa_media_type ENUM(...)
- fa_major_mime ENUM(...)
- fa_minor_mime VARBINARY(100)
- fa_description VARCHAR(767)
- fa_user INT
- fa_user_text VARCHAR(255)
- fa_timestamp BINARY(14)
- fa_deleted TINYINT
- fa_sha1 VARBINARY(32)

**uploadstash**
- us_id INT
- us_user INT
- us_key VARCHAR(255)
- us_orig_path VARCHAR(255)
- us_path VARCHAR(255)
- us_source_type VARCHAR(50)
- us_timestamp VARBINARY(14)
- us_status VARCHAR(50)
- us_chunk_inx INT
- us_props BLOB
- us_size INT
- us_sha1 VARCHAR(31)
- us_mime VARCHAR(255)
- us_media_type ENUM(...)
- us_image_width INT
- us_image_height INT
- us_image_bits SMALLINT

**oldimage**
- oi_name VARCHAR(255)
- oi_archive_name VARCHAR(255)
- oi_size INT
- oi_width INT
- oi_height INT
- oi_bits INT
- oi_description VARCHAR(767)
- oi_user INT
- oi_user_text VARCHAR(255)
- oi_timestamp BINARY(14)
- oi_metadata MEDIUMBLOB
- oi_media_type ENUM(...)
- oi_major_mime ENUM(...)

## Interwiki

**sites**
- site_id INT
- site_global_key VARBINARY(32)
- site_type VARBINARY(32)
- site_group VARBINARY(32)
- site_source VARBINARY(32)
- site_language VARBINARY(32)

**site_identifiers**
- si_site INT
- si_type VARBINARY(32)
- si_key VARBINARY(32)

## Caching tables

**objectcache**
- keyname VARBINARY(255)
- value MEDIUMBLOB
- exptime DATETIME

**l10n_cache**
- lc_lang VARBINARY(32)
- lc_key VARCHAR(255)
- lc_value MEDIUMBLOB

**transcache**
- tc_url VARCHAR(255)
- tc_contents TEXT
- tc_time BINARY(14)

## Maintenance

**updatelog**
- ul_key VARCHAR(255)
- ul_value BLOB

**job**
- job_id INT
- job_cmd VARBINARY(60)
- job_namespace INT
- job_title VARCHAR(255)
- job_timestamp VARBINARY(14)
- job_params BLOB

## page

- 🔑 page_id INT
- 🔷 page_namespace INT
- 🔷 page_title VARCHAR(255)
- 🔷 page_restrictions TINYBLOB
- 🔷 page_is_redirect TINYINT
- 🔷 page_is_new TINYINT
- 🔷 page_random DOUBLE
- 🔷 page_touched BINARY(14)
- 🔷 page_links_updated VARBINARY(14)
- 🔷 page_latest INT
- 🔷 page_len INT
- 🔷 page_content_model VARBINARY(32)
- 🔷 page_lang VARBINARY(35)

Indexes ▶

## page_props

- 🔷 pp_page INT
- 🔷 pp_propname VARBINARY(60)
- 🔷 pp_value BLOB
- 🔷 pp_sortkey FLOAT

Indexes ▶

## page_restrictions

- 🔑 pr_id INT
- 🔷 pr_page INT
- 🔷 pr_type VARBINARY(60)
- 🔷 pr_level VARBINARY(60)
- 🔷 pr_cascade TINYINT
- 🔷 pr_user INT
- 🔷 pr_expiry VARBINARY(14)

Indexes ▶

## protected_titles

- 🔷 pt_namespace INT
- 🔷 pt_title VARCHAR(255)
- 🔷 pt_user INT
- 🔷 pt_reason VARCHAR(767)
- 🔷 pt_timestamp BINARY(14)
- 🔷 pt_expiry VARBINARY(14)
- 🔷 pt_create_perm VARBINARY(60)

Indexes ▶

## redirect

- 🔑 rd_from INT
- 🔷 rd_namespace INT
- 🔷 rd_title VARCHAR(255)
- 🔷 rd_interwiki VARCHAR(32)
- 🔷 rd_fragment VARCHAR(255)

Indexes ▶

## category

- 🔑 cat_id INT
- 🔷 cat_title VARCHAR(255)
- 🔷 cat_pages INT
- 🔷 cat_subcats INT
- 🔷 cat_files INT

Indexes ▶

## pagelinks

- pl_from INT
- pl_from_namespace INT
- pl_namespace INT
- pl_title VARCHAR(255)

Indexes

## imagelinks

- il_from INT
- il_from_namespace INT
- il_to VARCHAR(255)

Indexes

## categorylinks

- cl_from INT
- cl_to VARCHAR(255)
- cl_sortkey VARBINARY(230)
- cl_sortkey_prefix VARCHAR(255)
- cl_timestamp TIMESTAMP
- cl_collation VARBINARY(32)
- cl_type ENUM(...)

Indexes

## image

- img_name VARCHAR(255)
- img_size INT
- img_width INT
- img_height INT
- img_metadata MEDIUMBLOB
- img_bits INT
- img_media_type ENUM(...)
- img_major_mime ENUM(...)
- img_minor_mime VARBINARY(100)
- img_description VARCHAR(767)
- img_user INT
- img_user_text VARCHAR(255)
- img_timestamp VARBINARY(14)
- img_sha1 VARBINARY(32)

Indexes

## English Wikipedia statistics

| Number of user accounts | Number of articles | Number of files | Number of administrators |
|---|---|---|---|
| 30,534,803 | 5,367,222 | 846,990 | 1,268 |

Uncompressed approximated data sizes:
- Article content ~ 50 GB
- User talk and data ~ 50 GB
- Full history of changes ~ 10 TB
- All files ~ 30 TB

## Our local Wikipedia database

| # of articles (without text) | # of files (names and infos) | # of categories |
|---|---|---|
| 3,096,190 | 372,825 | 1,569,810 |

Our database size is just ~ 5 GB

① 

```
mysql> select distinct pp_propname
    -> from page_props;
+------------------------------+
| pp_propname                  |
+------------------------------+
| defaultsort                  |
| disambiguation               |
| displaytitle                 |
| forcetoc                     |
| graph_specs                  |
| hiddencat                    |
              [...]
| templatedata                 |
| wikibase-badge-Q17437796     |
| wikibase-badge-Q17437798     |
| wikibase-badge-Q20748092     |
| wikibase_item                |
+------------------------------+
32 rows in set (0.47 sec)
```

②

```
mysql> select page_title, page_len
    -> from page
    -> order by page_len DESC
    -> limit 10;

+------------------------------------------------+----------+
| page_title                                     | page_len |
+------------------------------------------------+----------+
| RichardWeiss/Archivehistory                    |  1751049 |
| Johnfreez                                      |  1457543 |
| Reference_desk_archive/Science/January_2006    |  1356342 |
| Reference_desk_archive/Science/April_2006      |  1220726 |
| Upload_log_archive/September_2004_(1)          |  1219495 |
| Tarawneh/Archive_2016                          |  1200041 |
| Reference_desk_archive/Science/October_2005    |  1175930 |
| Euchiasmus                                     |  1172395 |
| MRacer                                         |  1161703 |
| Upload_log_archive/May_2004_(1)                |  1137338 |
+------------------------------------------------+----------+
10 rows in set (0.00 sec)
```

https://en.wikipedia.org/wiki/User_talk:RichardWeiss/Archivehistory

③

```
mysql> select pt_title
    -> from protected_titles
    -> where instr(pt_reason,"deprecating") and
    -> (instr(pt_title,"wiki") or instr(pt_title,"Wiki"))
    -> limit 100;
+------------------------------------------------------------------------+
| pt_title                                                               |
+------------------------------------------------------------------------+
| All_wiki_admin_are_bastards                                            |
| Fuck_Wikipedia                                                         |
| I_hate_wikipedia                                                       |
| NawlinWiki                                                             |
| Wikiabuse.com                                                          |
| Wikimocracy                                                            |
| Wikinazi                                                               |
| Wikipedia_and_OCD                                                      |
| Wikipedia_is_gay                                                       |
| Wikipedia_nerds                                                        |
| Wikipedia_sucks                                                        |
| Wikipedo                                                               |
| Wikislavia                                                             |
| Economic_Left_and_Social_Libertarian_Wikipedians                      |
| Fascist_Wikipedians                                                    |
| Neutral_Good_Wikipedians                                               |
| Stressed_Wikipedians                                                   |
| Suspected_Wikipedia_sockpuppets_of_Tjstrf                             |
| Wikipedians_born_in_1999                                              |
| Wikipedians_born_in_2000                                              |
| Wikipedians_born_in_2001                                              |
                               [...]
| Wikipedians_born_in_2006                                              |
| Wikipedians_born_in_2007                                              |
| Wikipedians_born_in_the_2000s                                         |
| Wikipedians_favoring_BJAODN                                           |
| Wikipedians_who_believe_West_Virginia_is_in_the_Southern_U.S.         |
| Wikipedians_who_dislike_George_W._Bush                                |
+------------------------------------------------------------------------+
58 rows in set (0.34 sec)
```

④

```
mysql> select rd_title, count(*) as numRed
    -> from redirect
    -> where rd_from < 10000
    -> group by rd_title
    -> order by numRed DESC
    -> limit 1;
+---------------------------+--------+
| rd_title                  | numRed |
+---------------------------+--------+
| List_of_sovereign_states  |     28 |
+---------------------------+--------+
1 row in set (0.06 sec)
```

https://en.wikipedia.org/wiki/List_of_sovereign_states

⑤

```
mysql> select page_title
    -> from page
    -> where page_id in
    -> (select rd_from
    -> from redirect
    -> where rd_title = "List_of_sovereign_states");
```

```
+----------------------------------+
| page_title                       |
+----------------------------------+
| CountriesOfTheWorld              |
| Countries_of_the_World           |
| List_of_independent_countries    |
| List_of_National_Titles          |
| List_of_countries                |
| List_of_countries_by_name        |
| Independent_States               |
| List_of_countries_of_the_world   |
               [...]
+----------------------------------+
36 rows in set (9.37 sec)
```

https://en.wikipedia.org/wiki/CountriesOfTheWorld

⑥

```
mysql> select c.cl_to as Category,
    -> p.page_title as Title,
    -> p.page_counter as Views
    -> from categorylinks c, page p
    -> where c.cl_from = p.page_id
    -> order by p.page_counter DESC
    -> limit 100;
+-------------------------------------------------------+---------------------------------+-------+
| Category                                              | Title                           | Views |
+-------------------------------------------------------+---------------------------------+-------+
| Wikipedia_fully-protected_project_pages               | Copyrights                      | 87540 |
| Wikipedia_copyright                                   | Copyrights                      | 87540 |
                            [...]
| Wikipedia_policies                                    | Copyrights                      | 87540 |
| Wikipedia_legal_policies                              | Copyrights                      | 87540 |
| 2003                                                  | Current_events/October_2003     | 53194 |
| All_BLP_articles_lacking_sources                      | List_of_French_people           | 44369 |
| BLP_articles_lacking_sources_from_July_2013           | List_of_French_people           | 44369 |
                            [...]
| Lists_of_French_people                                | List_of_French_people           | 44369 |
| Articles_including_recorded_pronunciations            | About                           | 32946 |
                            [...]
| Wikipedia_basic_information                           | About                           | 32946 |
| 2003_deaths                                           | Deaths_in_2003                  | 31832 |
                            [...]
| Wikipedia_indefinitely_move-protected_pages           | Deaths_in_2003                  | 31832 |
| Use_mdy_dates_from_October_2014                       | Mathematics                     | 30774 |
                            [...]
| Pages_using_ISBN_magic_links                          | Mathematics                     | 30774 |
| Wikipedia_features                                    | Searching                       | 24218 |
                            [...]
| Wikipedia_semi-protected_project_pages                | Searching                       | 24218 |
| Use_British_English_from_August_2016                  | World_War_II                    | 22064 |
                            [...]
| Wars_involving_Mexico                                 | World_War_II                    | 22064 |
+-------------------------------------------------------+---------------------------------+-------+
100 rows in set (1 min 22.26 sec)
```

```
⑦  mysql> create view withMoreImages(numb_images,totalSize_MB,id) as
        -> select count(il.il_to) as count,
        -> sum(i.img_size)/1000 as totalSize, il.il_from
        -> from imagelinks il, image i
        -> where il.il_to = i.img_name
        -> group by il.il_from
        -> order by count DESC
        -> limit 100;
     Query OK, 0 rows affected (0.01 sec)
```

```
mysql> select p.page_title,
    -> c.numb_images,
    -> c.totalSize_MB,
    -> c.totalSize_MB/c.numb_images
    -> as size_x_image
    -> from page p, withMoreImages c
    -> where p.page_id  = c.id;
```

```
+------------------------------------+-------------+--------------+---------------+
| page_title                         | numb_images | totalSize_MB | size_x_image  |
+------------------------------------+-------------+--------------+---------------+
| List_of_highways_in_Victoria       |          57 |     282.0410 |    4.94808772 |
| Calder_Highway                     |          52 |     272.1380 |    5.23342308 |
| Tasman_Highway                     |          51 |     293.3130 |    5.75123529 |
| List_of_highways_in_South_Australia |         46 |     242.8430 |    5.27919565 |
| Murray_Valley_Highway              |          44 |     223.7120 |    5.08436364 |
| List_of_highways_in_Queensland     |          41 |     229.4030 |    5.59519512 |
| Road_infrastructure_in_Melbourne   |          41 |     247.4980 |    6.03653659 |
| Western_Highway_(Victoria)         |          41 |     232.1900 |    5.66317073 |
| Albert_Gleizes                     |          40 |   55405.2240 | 1385.13060000 |
| Princes_Freeway                    |          34 |     169.4170 |    4.98285294 |
| National_Highway_(Australia)       |          34 |     277.5040 |    8.16188235 |
| Maroondah_Highway                  |          33 |     181.0790 |    5.48724242 |
| History_of_BBC_television_idents   |          30 |    1380.5710 |   46.01903333 |
| South_Gippsland_Highway            |          29 |     147.8680 |    5.09889655 |
| Goulburn_Valley_Highway            |          28 |     163.8150 |    5.85053571 |
| Bruce_Highway                      |          28 |     175.1960 |    6.25700000 |

                         [...]

+------------------------------------+-------------+--------------+---------------+
62 rows in set (17.66 sec)
```

⑧

```
mysql> select count(*) as num_upl_img, img_user_text
    -> from image
    -> group by img_user, img_user_text
    -> having avg(img_size) > 10
    -> order by num_upl_img desc
    -> limit 10;
+-------------+---------------------------+
| num_upl_img | img_user_text             |
+-------------+---------------------------+
|       16302 | Theo's Little Bot         |
|        5692 | DatBot                    |
|        5377 | 718 Bot                   |
|        5210 | DASHBot                   |
|        2970 | J04n                      |
|        2818 | GrahamHardy               |
|        1936 | Cavarrone                 |
|        1693 | Jasper the Friendly Punk  |
|        1662 | DISEman                   |
|        1607 | Gobonobo                  |
+-------------+---------------------------+
10 rows in set (39.89 sec)
```

⑨

```
mysql> select pl_title
    -> from pagelinks
    -> where pl_from in (
    -> select page_id
    -> from page
    -> where page_title in (
        -> select pl_title
        -> from pagelinks
        -> where pl_from in (
            -> select page_id
            -> from page
            -> where page_title = "List_of_Mars-crossing_minor_planets")));
+---------------------+
| pl_title            |
+---------------------+
| (163364)_2002_OD20  |
| (277475)_2005_WK4   |
| (285263)_1998_QE2   |
| (29074)_5160_T-3    |
| (29076)_1972_TR8    |
| (35395)_1997_XM10   |
| (35397)_1997_YJ     |
| (410777)_2009_FD    |
| (52759)_1998_MW13   |
| (52761)_1998_MN14   |
| (53318)_1999_JV7    |
| (53320)_1999_JW8    |
| (6177)_1986_CE2     |
| (68949)_2002_QN6    |
| (68951)_2002_QH27   |
| (7888)_1993_UC      |
| (89958)_2002_LY45   |
| (89960)_2002_ND35   |
| 101955_Bennu        |
| 10301_Kataoka       |
| 10303_Fr├®ret       |
| 1035_Amata          |
| 1037_Davidweilla    |
| 1220_Crocus         |
| 1222_Tina           |
| 131_Vala            |
| 133_Cyrene          |
+---------------------+
27 rows in set (5 min 7.61 sec)
```

# List of Mars-crossing minor planets

From Wikipedia, the free encyclopedia

## Inner grazers [ edit ]

- 1951 Lick
- 4947 Ninkasi
- (10302) 1989 ML
- 15817 Lucianotesi
- (21374) 1997 WS$_{22}$

- (39774) 1997 GO$_{27}$
- (52381) 1993 HA
- (52387) 1993 OM$_7$
- (54071) 2000 GQ$_{146}$
- (67367) 2000 LY$_{27}$

- (68031) 2000 YK$_{29}$
- (68278) 2001 FC$_7$
- (85184) 1991 JG$_1$
- (85236) 1993 KH
- (85989) 1999 JD$_6$

- (87024) 2
- (87684) 2
- (88213) 2
- (90367) 2
- (90403) 2

## Inner grazers that are also Earth-crossers or grazers [ edit ]

- 1620 Geographos
- 1865 Cerberus
- 2063 Bacchus
- 3361 Orpheus

- 3362 Khufu
- 3753 Cruithne
- 4034 Vishnu
- 4581 Asclepius

- 4769 Castalia
- 6239 Minos
- (10115) 1992 SK
- 11500 Tomaiyowit

- 12711 Tu
- (17511) 1
- (22099) 2
- (68950) 2

## Mars-crossers that are also Earth-crossers or grazers [ edit ]

These objects are not catalogued as Mars-crossers in databases such as the Jet Propulsion Laboratory's online Small-body Database
are categorized as Near Earth Objects (NEOs).

- 1566 Icarus

- 5786 Talos

- (13651) 1997 BR

- (42286) 2

```
+---------------------+
| pl_title            |
+---------------------+

        [...]

| (21374)_1997_WS22   |
| (39774)_1997_GO27   |
| (52381)_1993_HA     |
| (52387)_1993_OM7    |
| (54071)_2000_GQ146  |
| (67367)_2000_LY27v  |
| (68031)_2000YK29    |

        [...]

+---------------------+
745 rows in set
(3 min 14.69 sec)
```

```
mysql> select pl.pl_title
    -> from pagelinks pl
    -> where pl.pl_from in (
        -> select p.page_id
        -> from page p
        -> where instr(p.page_title,"List_of_Mars-crossing_minor_planets"))
    -> and pl.pl_title not in (
        -> select page_title pp
        -> from page pp );
```

The existence of an internal link from a page to an existing
or non-existing page is recorded in the pagelinks table.
(https://meta.wikimedia.org/wiki/Help:Link)