

BACKGROUND INFORMATION

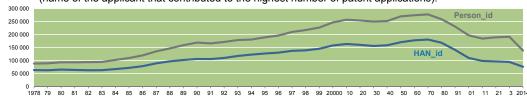
The OECD HAN database provides a grouping of patent applicant's names which has been elaborated with business register data. The names of patent applicants were originally extracted from *European Patent Office's (EPO) Worldwide Statistical Patent Database (PATSTAT, October 2011 edition).*

The grouping has been propagated to applicant identifiers up to *Spring 2016 edition* of PATSTAT. The database also includes the list of patent documents filed to the EPO, the US Patent and Trademarks Office (USPTO) or through the Patent Co-operation Treaty (PCT).

METHODOLOGY

The grouping of patent applicant names (PATSTAT's PERSON table) has been performed as follows:

- Cleaning and harmonising: names are corrected from punctuation, accents, abbreviations and legal
 information, using dictionaries developed on a country basis. A preliminary grouping is generated upon
 the harmonised name.
- Consolidating: cleaned/harmonised names were matched against company names from business register data (as provided in the ORBIS© database from Bureau van Dijk Electronic Publishing, June 2011). The matching was performed using series of algorithms (approximate string matching; weighted token-based comparisons; distance measures) within the IMALINKER system developed for the OECD by IDENER, Seville 2013. Each algorithm computes a matching score per pair of names, assessing therefore for the likelihood of names similarity. The matched pairs of names are selected according to high thresholds of matching scores in order to maximise the precision of the match. Finally, names are further grouped together according to either the matched ORBIS© company name or the cleaned/harmonised names resulting from the algorithms. The harmonisation of names was propagated to new applicant names from the latest edition of PATSTAT's PERSON table.
- Grouping: A unique identifier HAN_id is automatically generated for each grouping of patent applicants.
 A common name is then attributed to each HAN_id group according to the first applicant of the grouping (name of the applicant that contributed to the highest number of patent applications).

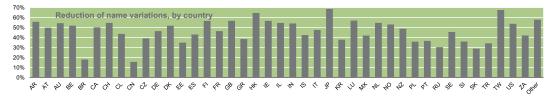


Due to the large volume of data processed, it was not possible to control each names grouping. Therefore, errors may be encountered: any feedback on incorrect harmonisation would be highly appreciated.

RAGE

The OECD HAN database, Spetember provides groupings of patent applicant's names for most OECD countries and countries in the BRIICS. The list of patents filed to the EPO, the USPTO and through the PCT is

made available for each grouping of applicants. Further improvements are expected in future versions, notably on the countries coverage.



BASE
The OECD HAN database, September 2016, consists of 3 distinct tables presented in flat files (UTF-8 format, columns separated using the pipe "|" character). Applicant's idenfiers from the last 4 editions of PATSTAT are

linked to 2,882,783 unique HAN identifiers (HAN_id). Changes in the HAN_Id may occur from one version to the next.

HAN_PERSON Correspondance table between HAN_id and Person_id 5,449,554 rows		
HAN_id	Unique identifier - surrogate key Modified at each data release	
Person_id	Surrogate key - applicant identifier in PATSTAT, October 2013 to Sp	oring 2016 editions
Apr13_Person_id	Surrogate key - applicant identifier in PATSTAT, April 2013	
Oct12_Person_id	Surrogate key - applicant identifier in PATSTAT, October 2012	
Apr12_Person_id	Surrogate key - applicant identifier in PATSTAT, April 2012	
Oct11_Person_id	Surrogate key - applicant identifier in PATSTAT, October 2011	
Person_name_clean	Harmonised applicant name	

HAN_NAME Names associated to each HAN_id 2,882,783 rows			rows
HAN_id	Unique identifier - surrogate key Modified at each data release		
Clean_name	Proposed harmonised name (top applicant in the proposed grouping)		
Person_ctry_code	Applicant's country		
Matched	Indicator of sucessful match to ORBIS© (=1 if matched)		

HAN_PATENTS Patents filed by each HAN_id (for EPO, USPTO, PCT only) 12,620,410 rows			
HAN_id	Unique identifier - surrogate key Modified at each data release		
Appln_id	Surrogate key - patent application identifier in PATSTAT		
Publn_auth	Publication authority		
Patent_number	Patent publication number - normalised format EPXXXXXXX (patent published by to the EPO) USXXXXXXX (patent granted by USPTO) USYYYYXXXXXX (patent published by USPTO) WOYYYYXXXXXX (publication of patent application filed through the PCT) where YYYY represents the filing year and X in {0-9}		

RESTRICTIONS SOURCE& CONTACT

Please note that the OECD HAN database is provided for research and analytical work. When publishing the results of your analysis, make sure it is quoted as: "OECD, HAN database, September 2016".

For further information about OECD patent related work, methodologies and access to patent indicators, please visit our web page at: oe.cd/ipstats. Comments and questions about this dataset should be sent to STI.Microdata@oecd.org. For further information on EPO's PATSTAT, please contact patstat@epo.org.