# _Italian Dataset_

602,380    546,313    1,320,892

No. of missing articles(un-scraped) = 56,067
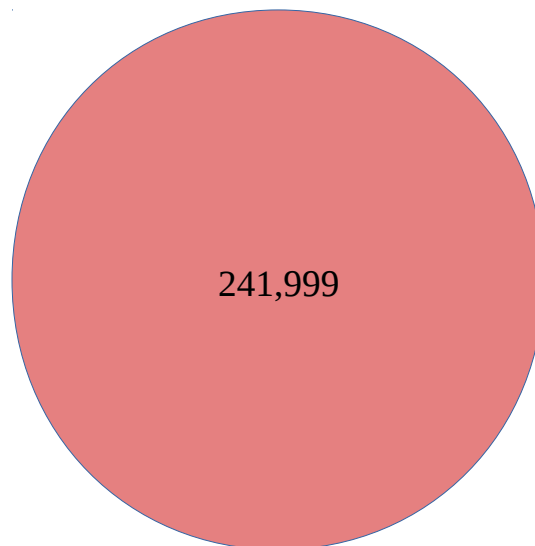
## PDF 2015

435,252    418,268    544,557

No. of missing articles(un-scraped) = 16,984
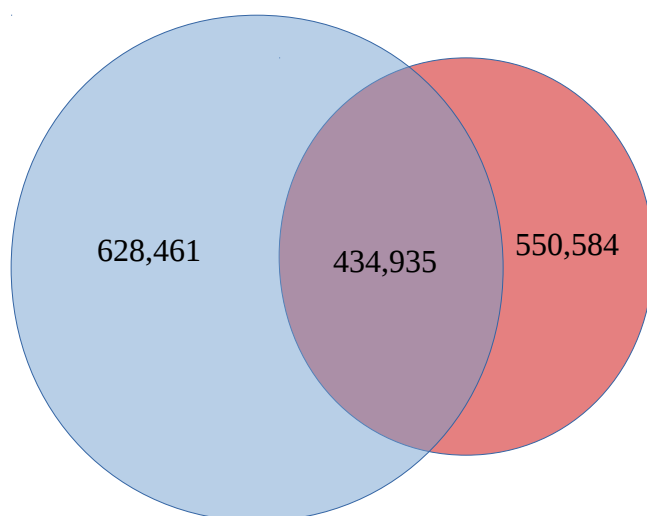
## PDF 2016

No. of missing urls(un-scraped) =  6271

**POLI 2015**



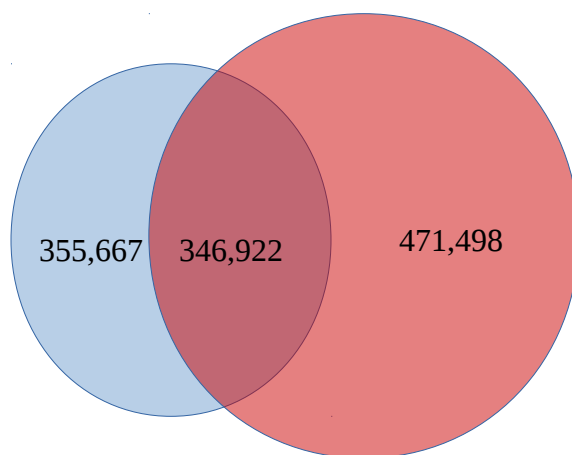No. of missing urls(un-scraped) = 241,999
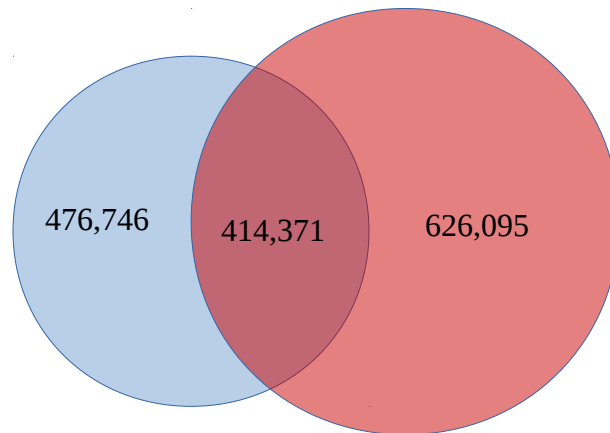
**POLI 2016**

*Extra urls scraped =  193526
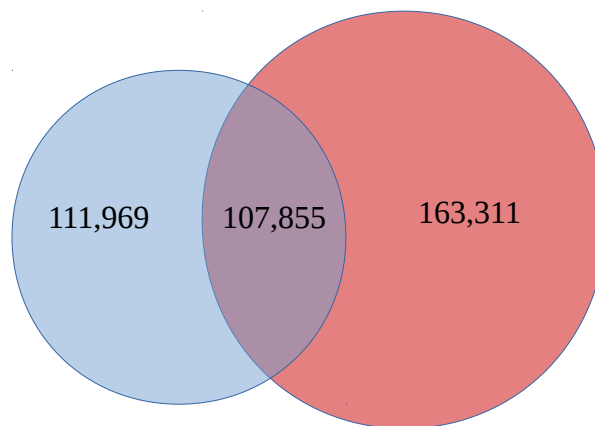
**PUL 2015**



No. of missing urls(un-scraped) = 8745

**PUL 2016**

No. of missing urls(un-scraped) = 62,375

**UMB 2015**



No. of missing urls(un-scraped) = 4114

**UMB 2016**

# Data comparison statistics:



Bar chart showing data comparison for PDF, POLI, PUL, and UMB categories across '15 (blue) and '16 (red).

Vertical axis labels (months): Jan, Feb, Mar, Apr, May, Jun, Jul, Aug, Sep, Oct, Nov, Dec

| Category | '15 | '16 |
|----------|-----|-----|
| PDF | Dec | May |
| POLI | Mar | Jan |
| PUL | Dec | Apr |
| UMB | Feb | May |