# Mass digitization and digitization projects at National library of Florence

Giovanni Bergamin
Biblioteca Nazionale
Centrale Firenze

IRCDL 2016

*12th Italian Research Conference on Digital Libraries*

# Some definitions …

*Mass digitization of books (**MDB)** =*
conversion of materials (books) *on an
industrial scale* (not just on a large-scale);
*conversion* of whole libraries *without making
a selection* of individual materials

# Two main MDB projects:

Google Books



Internet Archive (Open Content Alliance)





Building a digital archive of global content
for universal access

# Some notes on Google Books_1

Started in 2004

Planned end of the project 2020

The Google Books aim is the Google aim:

"organize the world's information and make it universally accessible and usable"  ---  so the ***content of all published books*** has to be searchable together with the ***content of all web*** pages

# Some notes on Google Books_2

"Just how many books are out there?"

**Books of the world, stand up and be counted! All 129,864,880 of you.**
Thursday, August 05, 2010 at 8:26 AM
Posted by Leonid Taycher, software engineer

How many books have already been digitized by Google Books?
25-30M (non ci sono statistiche ufficiali)

# Numbers ...

# A famous debate on GB in 2005_1

# A famous debate on GB in 2005_2

Jean-Noël Jeanneney, historian and  former President of National Library of France  wrote in 2005 that:

> *The promise of Google is enchanting [...]: everyone with access to the Internet can soon view the recorded memory of the ages in the palm of their hand and search this universe in a fraction of a second"* …   however ...

# A famous debate on GB in 2005_3

*"We are faced with several possible dangers with respect to: works of various cultural heritages that have fallen into the public domain, the **list of priorities will likely weigh in favor of Anglo-Saxon culture**; works still under copyright, of which only excerpts, or "snippets," will be offered for the time being, the **weight of American publishers may be overwhelming**; journals and books disseminating ongoing research, the **dominance of work from the United States** may become even greater than it is today"*

# 11 years later ...

according to "reliable sources" the highest percentage of the digitized books is in English (close to 50% out of 450 languages of books in GB)
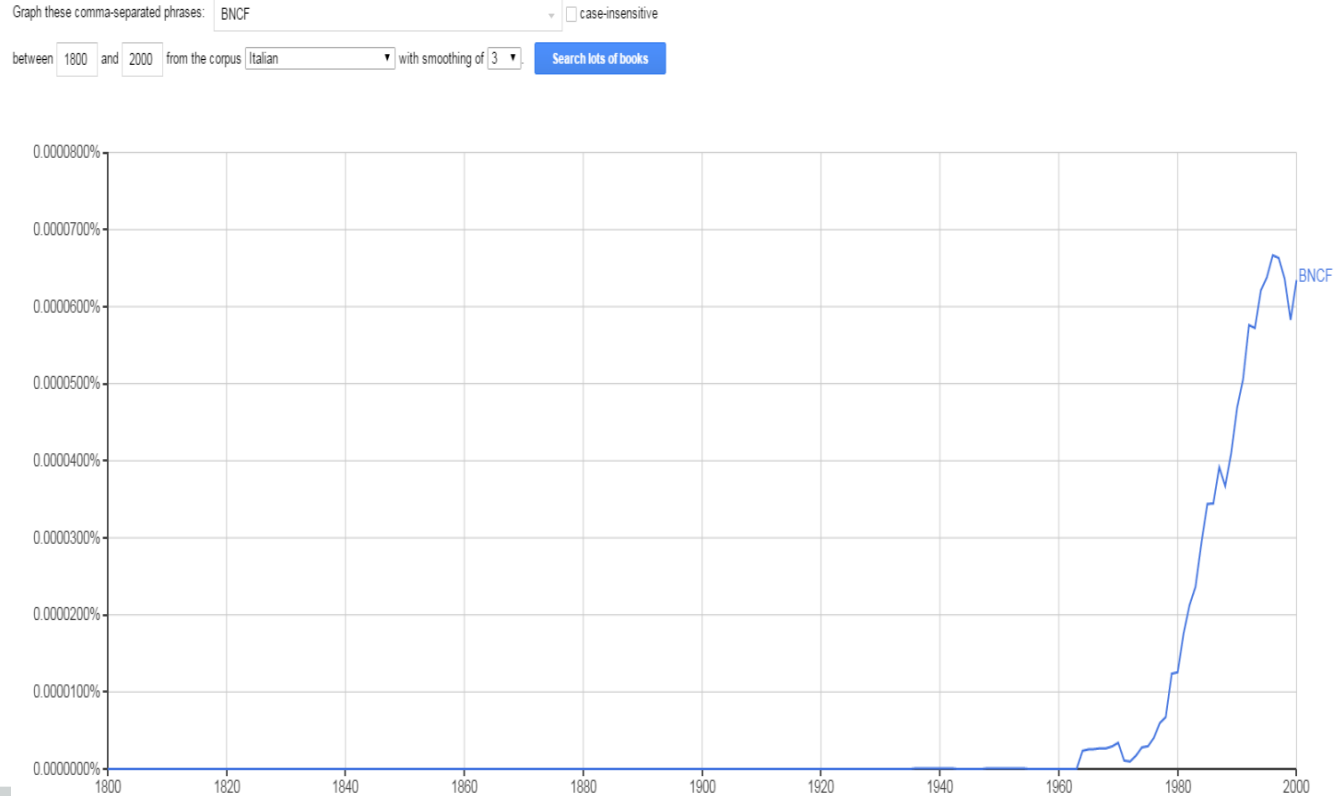
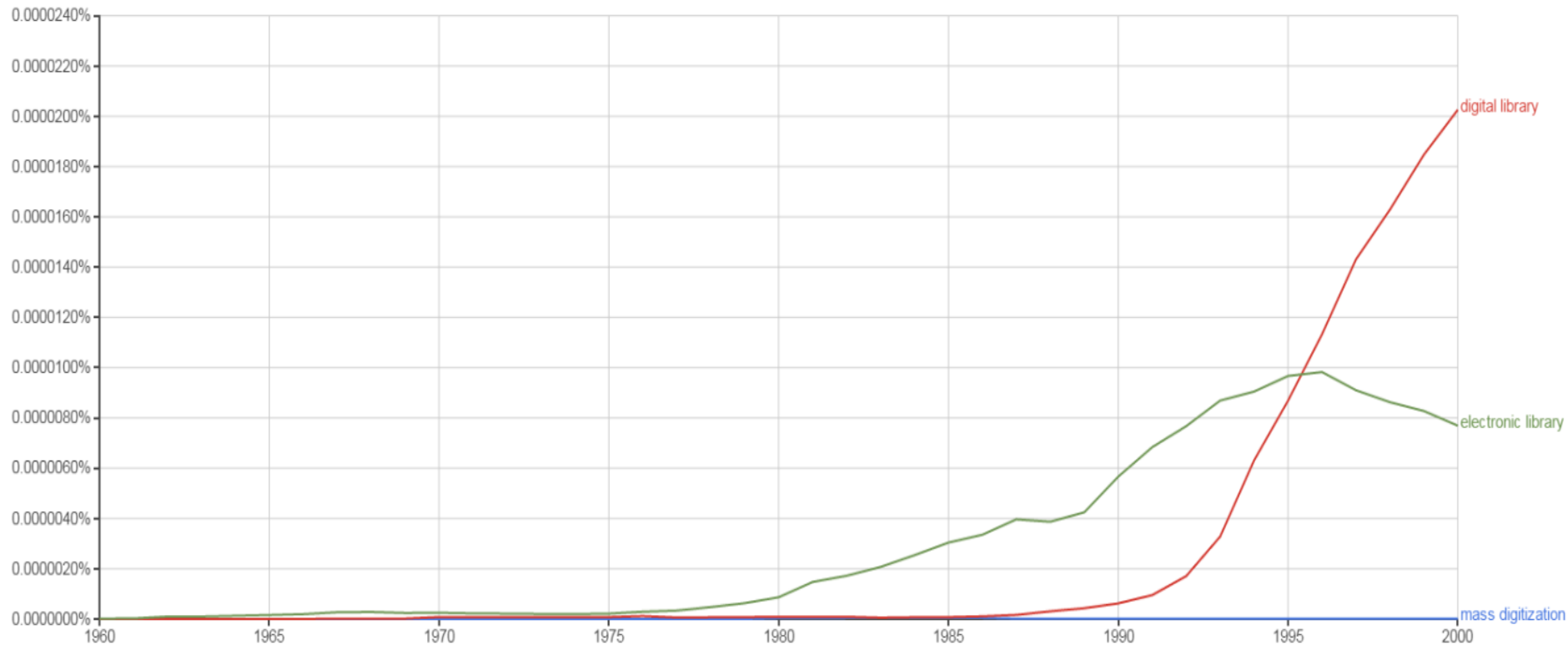The interest of Google for non English languages is growing

# ex. g. Ngram Viewer

Ngram service now available also for texts in:
German,French, Italian, Spanish, Russian, Hebrew, Chinese
From 2009 year after year they are adding new languages

# Somes notes on Internet Archive (OCA)_1

The Open Content Alliance (OCA) is a consortium of organizations contributing to a permanent, publicly accessible archive of digitized texts.
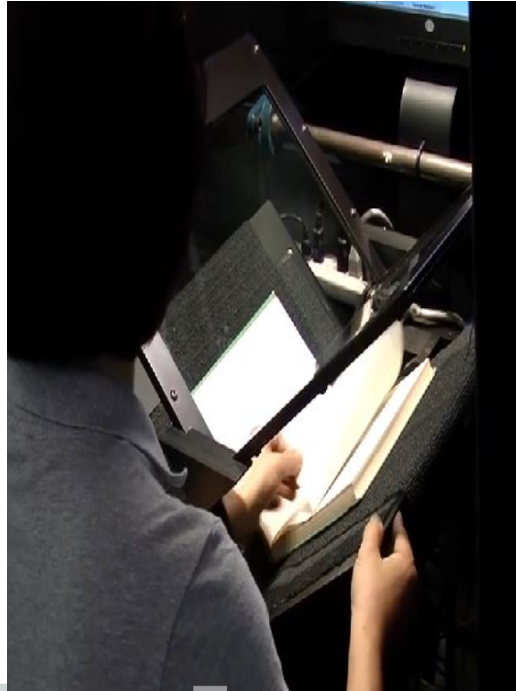
Its creation was announced in October 2005 by Yahoo!, the Internet Archive, the University of California, the University of Toronto and others.

Scanning for the OCA is administered by the Internet Archive, which also provides permanent storage and access through its website

# Somes notes on Internet Archive (OCA)_2

More than 8,7 million of texts available up to now
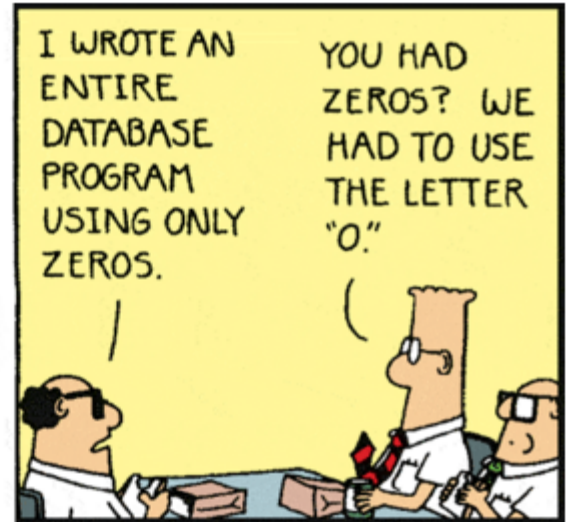
# Some differences





private owned (1 company)
huge amount of
resources available for
digitization
research and development

consortium between
companies and
nonprofit institutions
depends on donations
and self-financing
(limited resources
available)

# Digitization projects at BNCF (DPB)

started in the early 1990s  when the size of HD was 40 Megabytes (and there was no WWW)

# DPB: faithful copy or searchable text? - 1

Early projects aim:
   enrichment of bibliographic records through
   the digitization of title pages, table of
   contents etc (OCR)
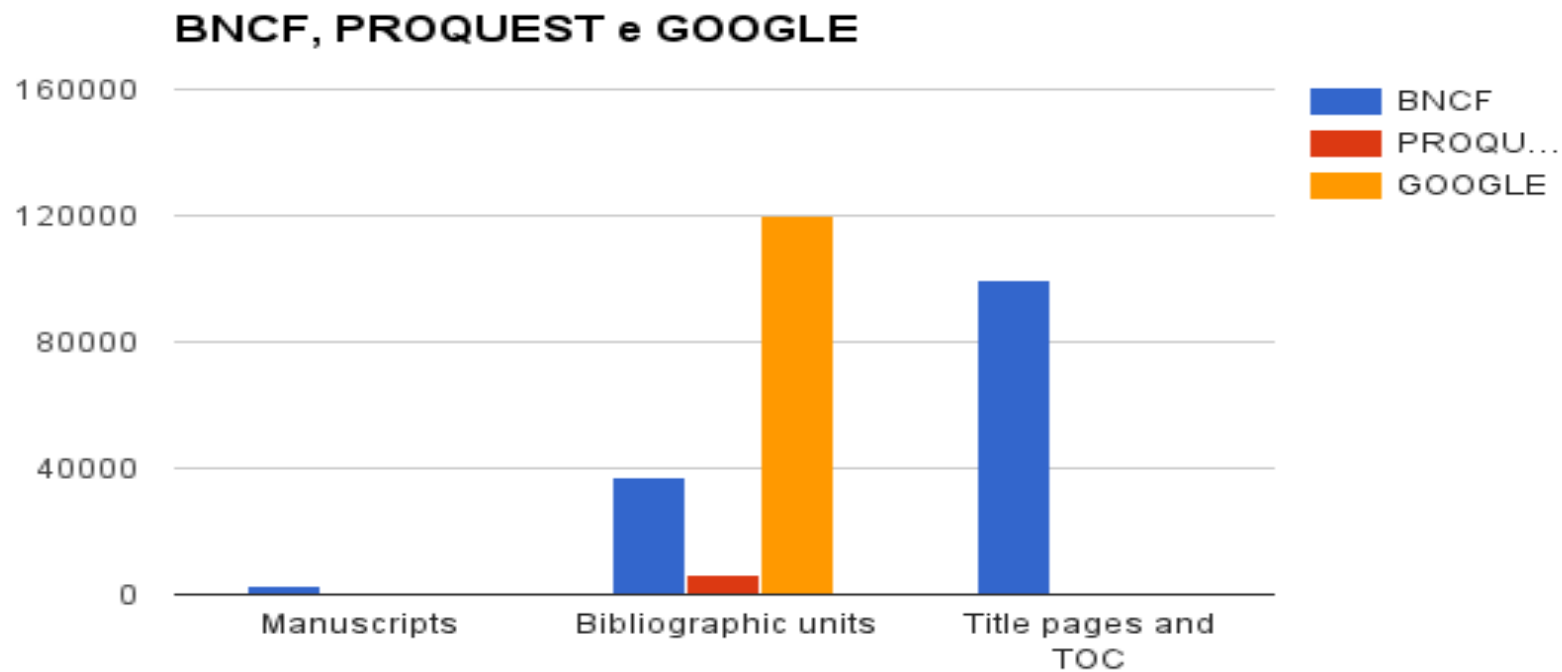
# DPB: faithful copy or searchable text? - 2

Following projects aim:
  faithful copy for manuscripts, ancient books,
    maps etc

# DPB results



BNCF, PROQUEST e GOOGLE

# Google Books and Proquest EEB at BNCF - 1

GB
- *range:* 1701-1875
- ebooks with "liquid" and searchable text
- national project
- *costs:* books circulation (inside lib.)

EEB
- *range:*      -1700
- faithful copy
- BNCF project
- *costs:* none

# Google Books and Proquest EEB at BNCF - 2

GB
  *scanning location:* outside lib. (Italy)
  *outcomes:* GRIN and free worldwide accessibility

EEB
  *scanning location:* inside lib.
  *outcomes:* "master" files and
      free access from Italian IP
      access fee outside and
        royalties for BNCF

# MD problems (Google)

limitations, ex. g.:

  size of books

  foldouts (from 2016 it will be possible)

note:

  MPOB Modified Process for older books pre 1700
    (color *and* text)

# Il gesuitismo svelato

Parigi : presso Pagnerre editore, 1846.
45 p. ; 16 cm.

Monografia - Materiale a stampa - Pubbl. in: fr - Lingua: it

75

a poca scala, e alla me-

no: nelle *uu, a tríi* o *staff,* ontate o tte.

ini: pic- *álcor).* stoffa di , talora n filo di

e nello scendere — cra- canape, cne na qualche
dino : ciascuno somiglianza col frusta-
stessi piani de gno, m più fina e
delle chiese ri più fo molto.
nobili edifi *Basgia* : quel co-
ognuno vo fa della ca-
posano sul n ine o
l'alta ana, el-
A nata : lino
aglioni mma-
*a*) co ita.
tri ana : la

distanza dalla scab, e che acceniiatio alla medesima. *Basellin.* Predellino nelle carrozze — *a duu, a trii pass* o *pastad,* o *staff,* a due, a Ire **montate** o ballenti o palelle. *Basgerilt.* Calcolini: piecole calcóle *(calcar).* Basen. Basilio: **stoffa** di filo di **colone** , **talora** mista anco con filo di canape, **che** ha qualclie col fruelapiú **fina e molto.** : qui' fa

rinchiuso in baccello *(sgorbia)* erboso , cra.s so, più facile a pulrefarsi che non a disseccarsi.

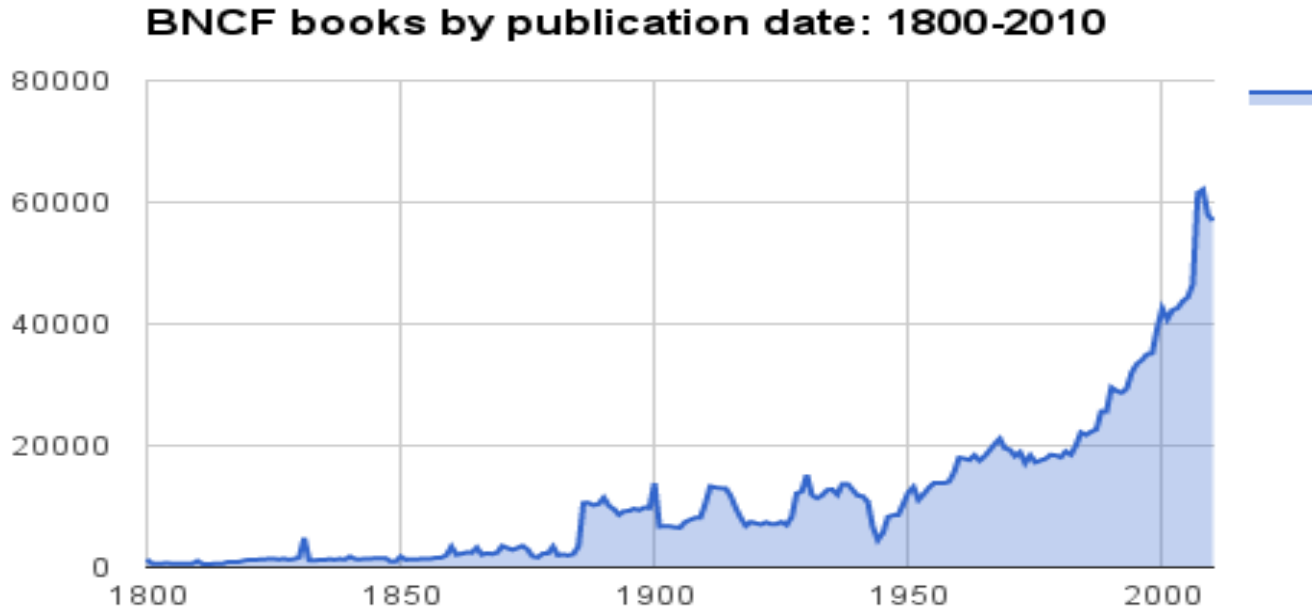*Basyinnoeura.* Fava minuta (Vicia faba minori.

, ove appender i

/» ciappa), riè a

piede,

canape, cbe ha qualche somiglianza col frusta- gno, m più fina e più fo molto. : *Basgia* : quel co- vo fa della ca- n ine el- mma- ita. ato assai *Bazzetta:* la agnellino (*be- n* nato, o da po- Fava (Vicia faba): ume *(lemm)* di for- na bislunga, scbiaccia- ta, col bellico *(oggia),* o segno del germe a una delle estremità,

# MD and copyright  (orphan works etc.)



BNCF books by publication date: 1800-2010

# BNCF and Wikisource

2014 Agreement between BNCF and Wikimedia Italia for Wikisource
- starting point: public domain book digitized by BNCF
- aim: improve access to digitized books
- results: crowdsourced text correction in Wikisource (the free library that anyone can improve)

# How it works

Manuale 150 ricette di cucina di guerra.djvu/77

Questa pagina è stata trascritta, formattata e riletta.

— 71 —

144. - **Pomodori secchi.**

Durante la stagione preparare, seccandoli al sole, dei pomodori, che prima saranno stati tagliati per metà e fatti sgocciolare, mettendoli a strati con sale in una cesta. Seccati che siano si conservano in una cesta o infilati a corona. Per servirsene metterli prima a bagno in acqua tiepida.

144. - Pomodori secchi.

Durante la stagione preparare, seccandoli al sole, dei pomodori, che prima saranno stati tagliati per metà e fatti sgocciolare, mettendoli a strati con sale in una cesta. Seccati che siano si conservano in una cesta o infilati a corona. Per servirsene metterli prima a bagno in acqua tiepida.

# Closing remarks and open questions

MDB: Is there an alternative to Google Books?

cooperation with IA and Wikisource

140 years buffer (orphan works and cooperation with publishers)

# Thank you for your patience

giovanni.bergamin@gmail.com