

From raw survey data to a structured register : evidence from the Census



Yves-Laurent Bénichou
Frédéric Comte
Julie Djiriguan
Benjamin Sakarovitch
Eric Sigaud

census is a boring party



*identifying (in the register)
the declared employer (in
the census)*

1. the problem and the current state of affairs
2. the proposed solution
3. in the context of innovation in the NSI

employer in the census

the *who do you work for* question

- what is it ?
- why asking it ?
- how is it processed ?

declaring the employer in the census

- what is the **name of the company** that employs you ?
- what is the **address** of the place you mostly go to work ?
- what is the **activity** of the company that employs you ?
- + what is your job ?

⇒ then **INSEE** tries to match it to the exact company in its register

The image shows a page from the 2012 French population census (Recensement de la population - 2012) titled 'Bulletin individuel'. The form is divided into several numbered sections:

- 1 Sexe:** Masculin (1) / Féminin (2)
- 2 Date et lieu de naissance:** Né(e) le... à... (commune de naissance)
- 3 Quelle est votre nationalité ?** (Française, Étrangère)
- 4 Êtes-vous inscrit(e) dans un établissement d'enseignement pour l'année scolaire en cours ?** (Oui/Non)
- 5 Où habitez-vous le 1^{er} janvier 2011 ?** (Même logement, autre logement, autre commune)
- 6 La suite du questionnaire s'adresse aux personnes de 14 ans ou plus.**
- 7 Visez-vous en couple ?** (Oui/Non)
- 8 Quel est votre état matrimonial légal ?** (Célibataire, Marié, Divorcé, etc.)
- 9 Quels diplômes avez-vous ?** (Diplôme de 1^{er} cycle universitaire, etc.)
- 10 Quelle est votre situation principale ?** (Employé, Étudiant, Chômeur, etc.)
- 11 Travaillez-vous actuellement ?** (Oui/Non)

The form includes a 'Cache à temps par l'agent recenseur' box and a 'Continuer page suivante et n'oubliez pas de signer' instruction at the bottom.

what indicators does it fuel ?

1/ Coding the socio-professional status

exact company in the register ⇒ activity code at the finest level ⇒ useful to code the socio-pro status

2/ Home-to-work commuting distance

currently it is computed as the distance between the 2 town halls...

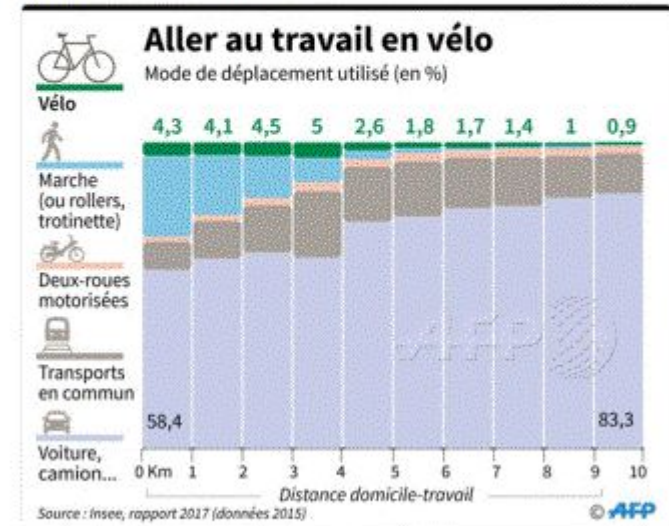


Gary Dagorn

@garydagorn

Suivre

Quasi 60% des gens prennent leur voiture pour faire moins d'un kilomètre jusqu'à leur travail. Moins de 1000 mètres en voiture.



12:40 - 15 sept. 2018

542 Retweets 474 J'aime



the current (costly) procedure

2 **coding machines** (from the 80s) with decision rules

- from company name and address : identification in the register
- in case of doubt : coding the activity to decide
- if unsuccessful : goes to employees doing it manually

44% of declared employers are coded automatically (from a total of 1.7M)

70 persons do the 56% left for 5 months

even by hand they find the exact company in only 70% of the cases (rest is activity-coded or blanked)

why is automatisisation a challenge ?

imprecisions in the declaration :

- incoherent address (street or municipality)
- declared activity does not correspond to the one registered
- vague description of the employer : eg ministry of education without further indication (middle school...) and/or address

registration of companies :

name not corresponding to the usual denomination

- management contracts
- acronyms
- persons names



how to assess performance ?

No ground truth

Special process to assess accuracy of handmade identification :

- every five years or so - last from 2014
- on 32000 census sheets (30000 manual / 2000 auto)
- double additional manual check + third for decision in case of disagreement

⇒ on the 2000 automatically coded : final coding coherent in 86% of cases

⇒ on the 30000 manually coded : final coding coherent in 66% of cases

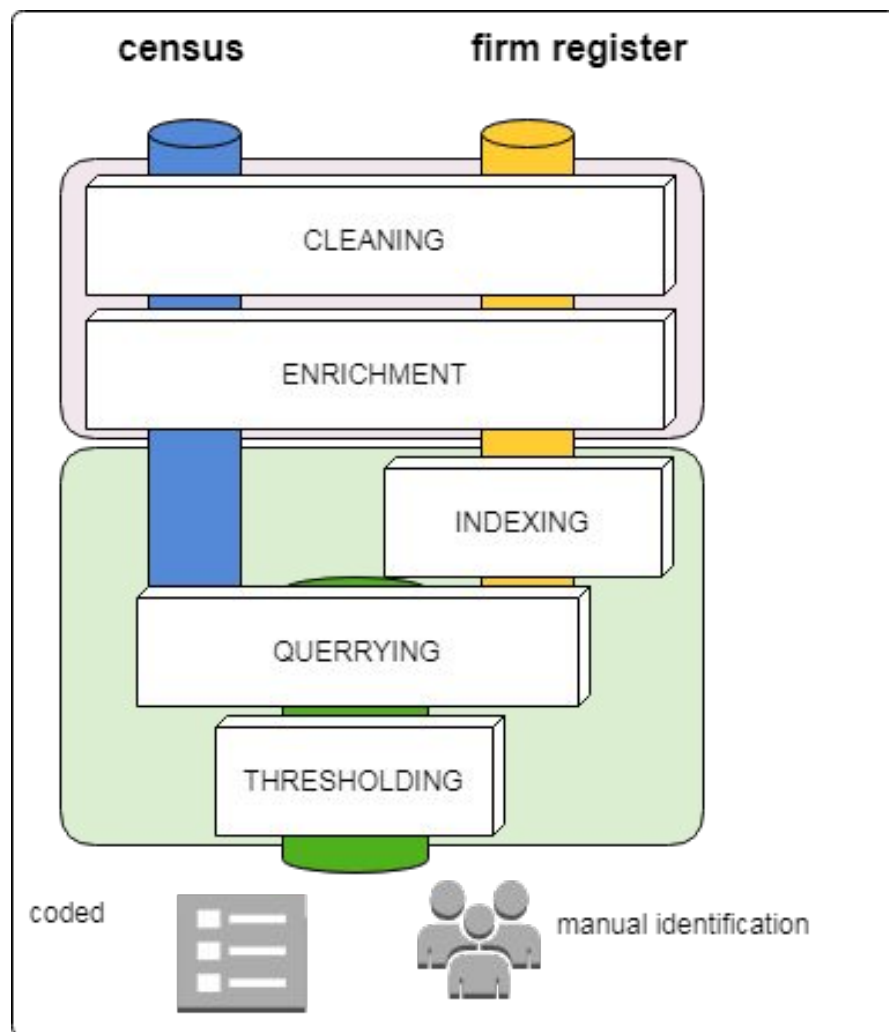
*only metrics are : **coverage** (number of employers automatically coded) and **consistency** (with current process of coding)*

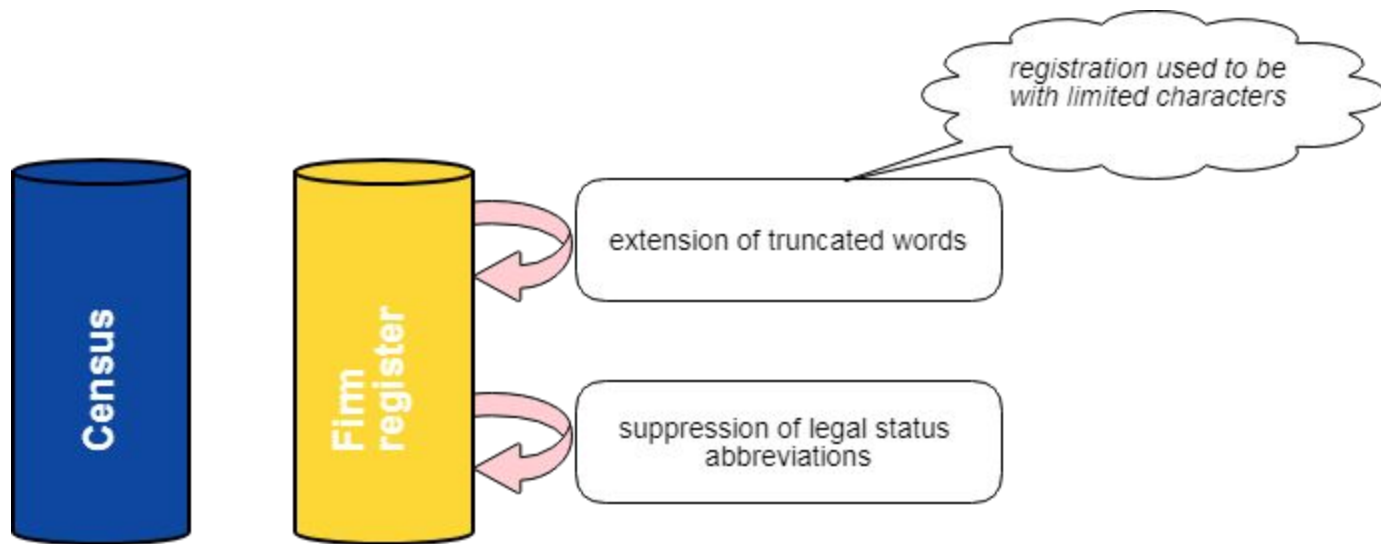
Developing a prototype

a search engine to
solve the scalability
challenge

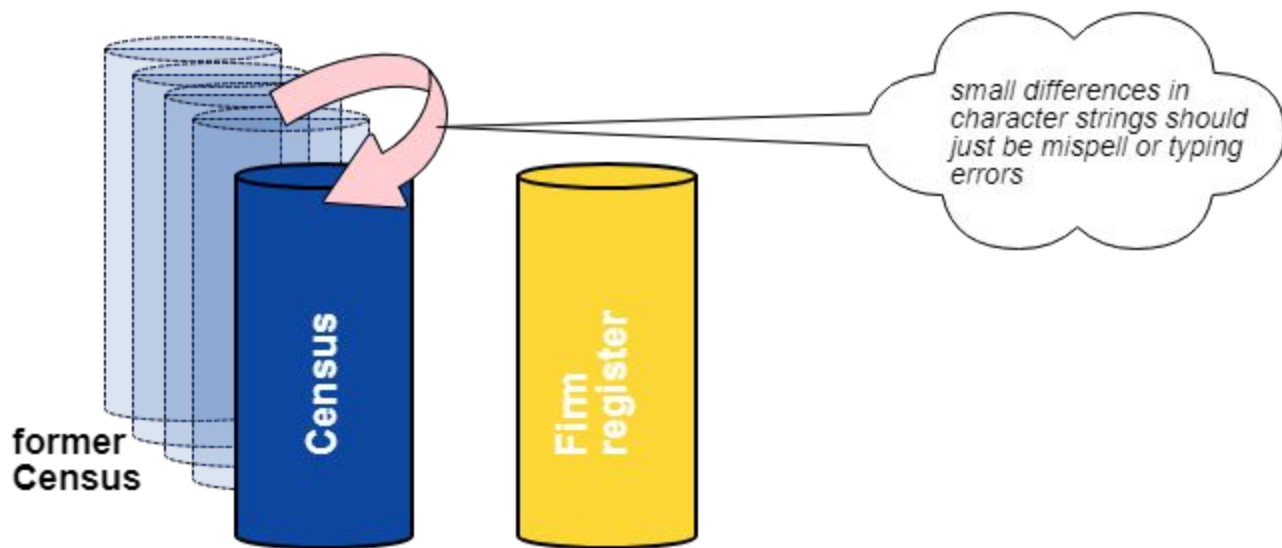


elasticsearch

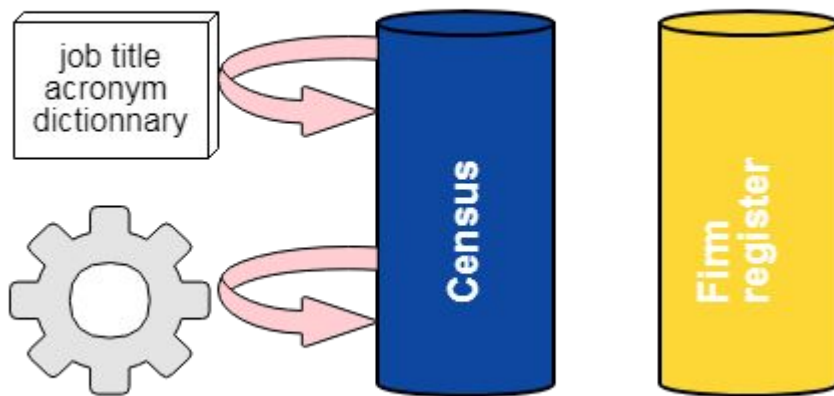




clustering similar observation
(per field) : uniformisation to
most common declaration

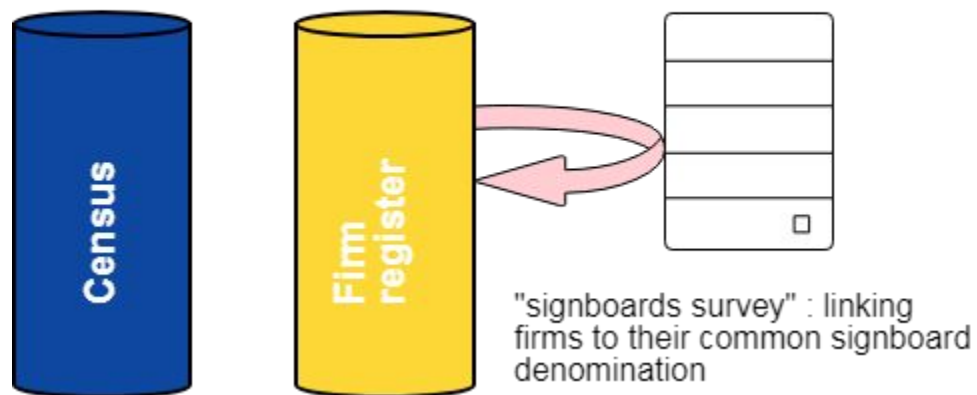


enriching the declared activity of the employer

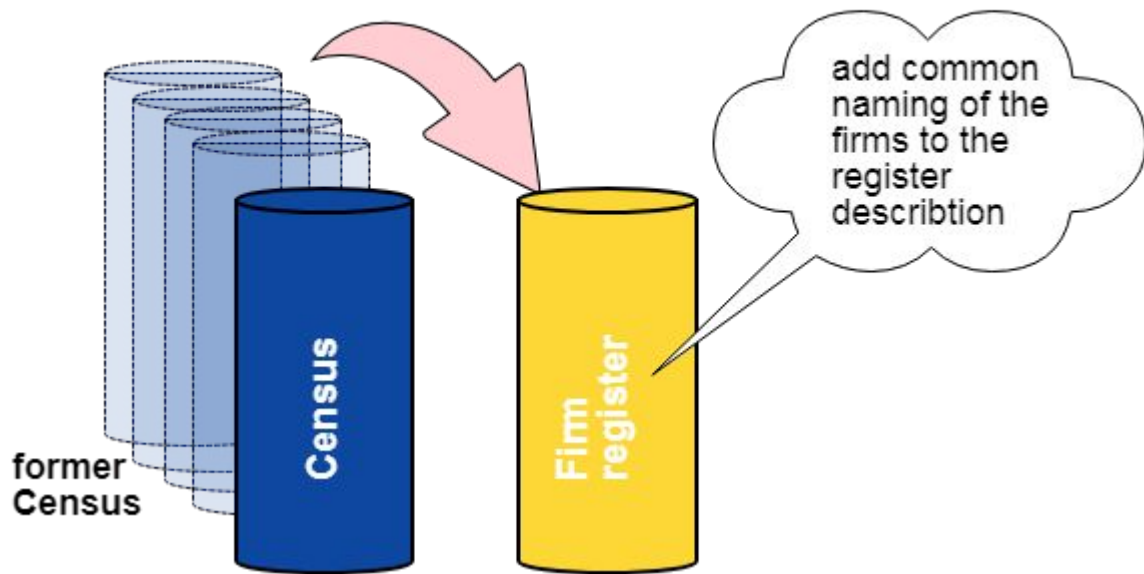


coding declared activity to "NACE"

**enriching firms description with survey
from the NSI**

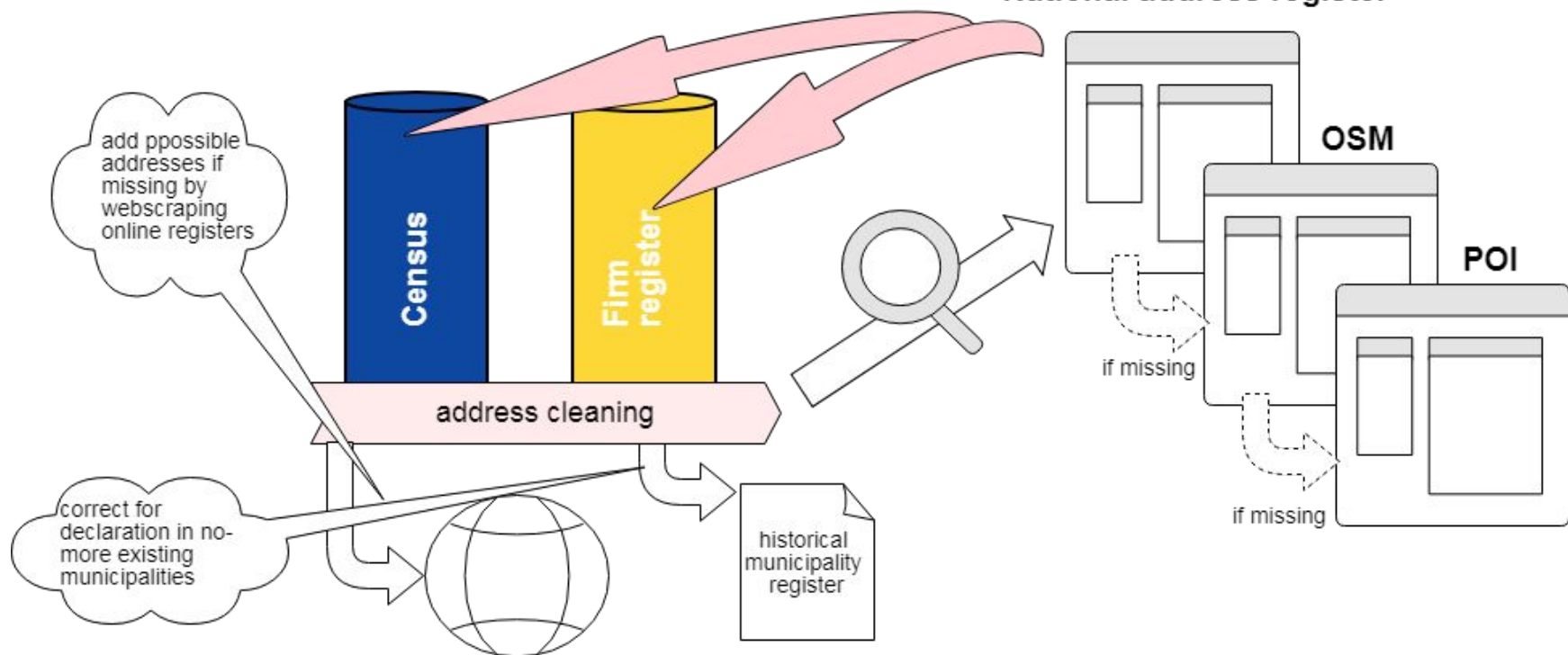


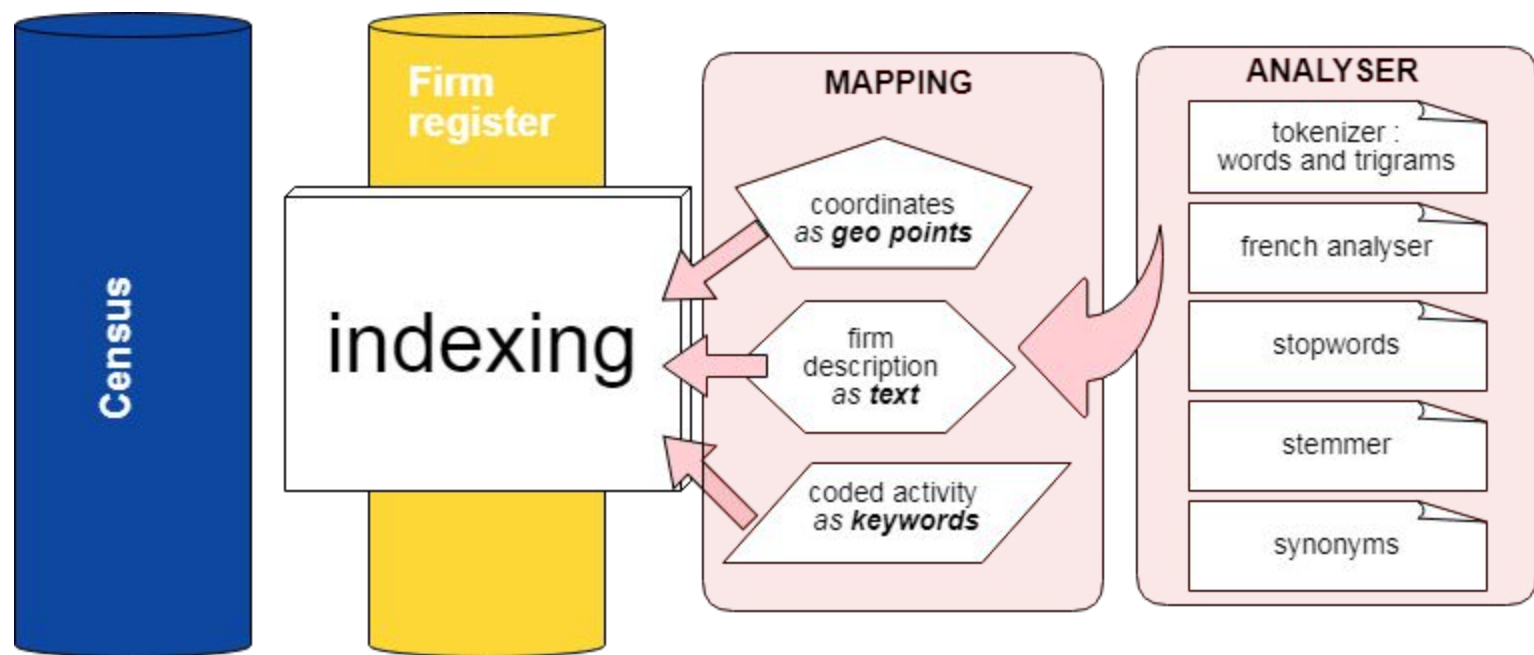
crowd sourcing from the past :
enriching the register from
passed years identifications

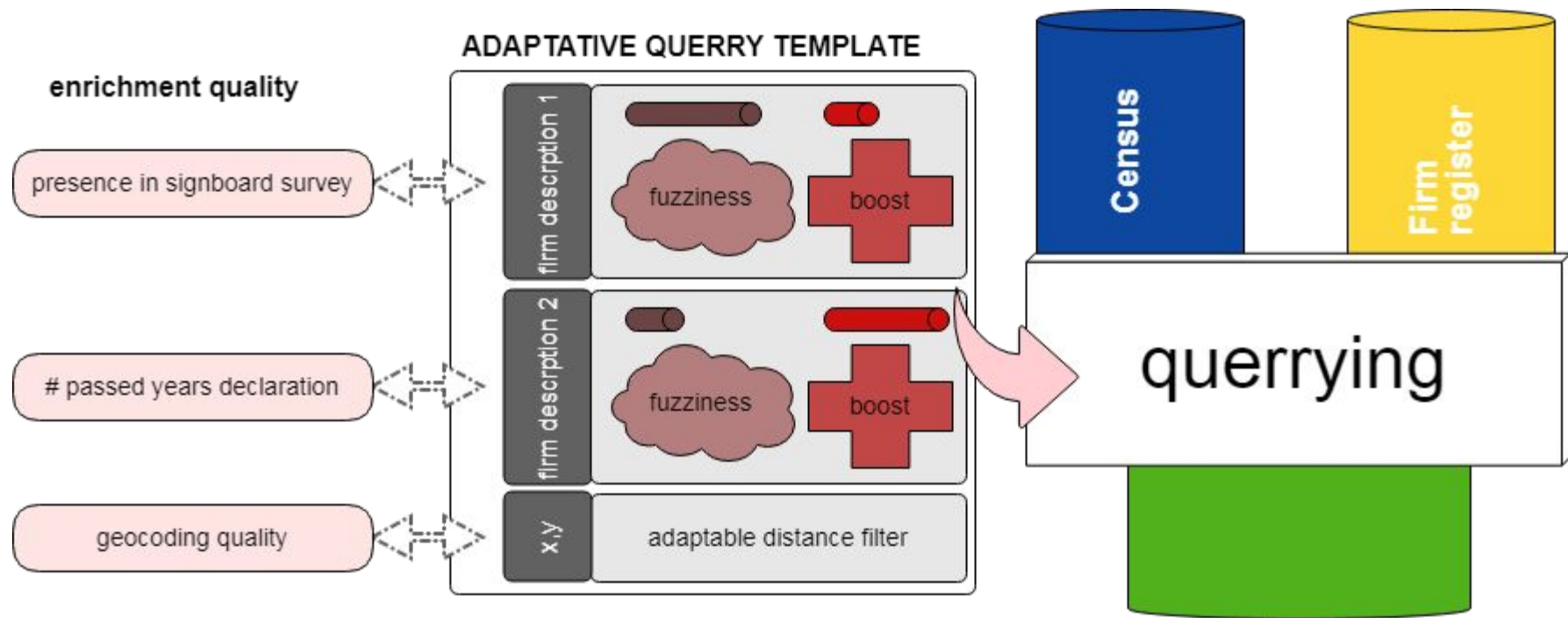


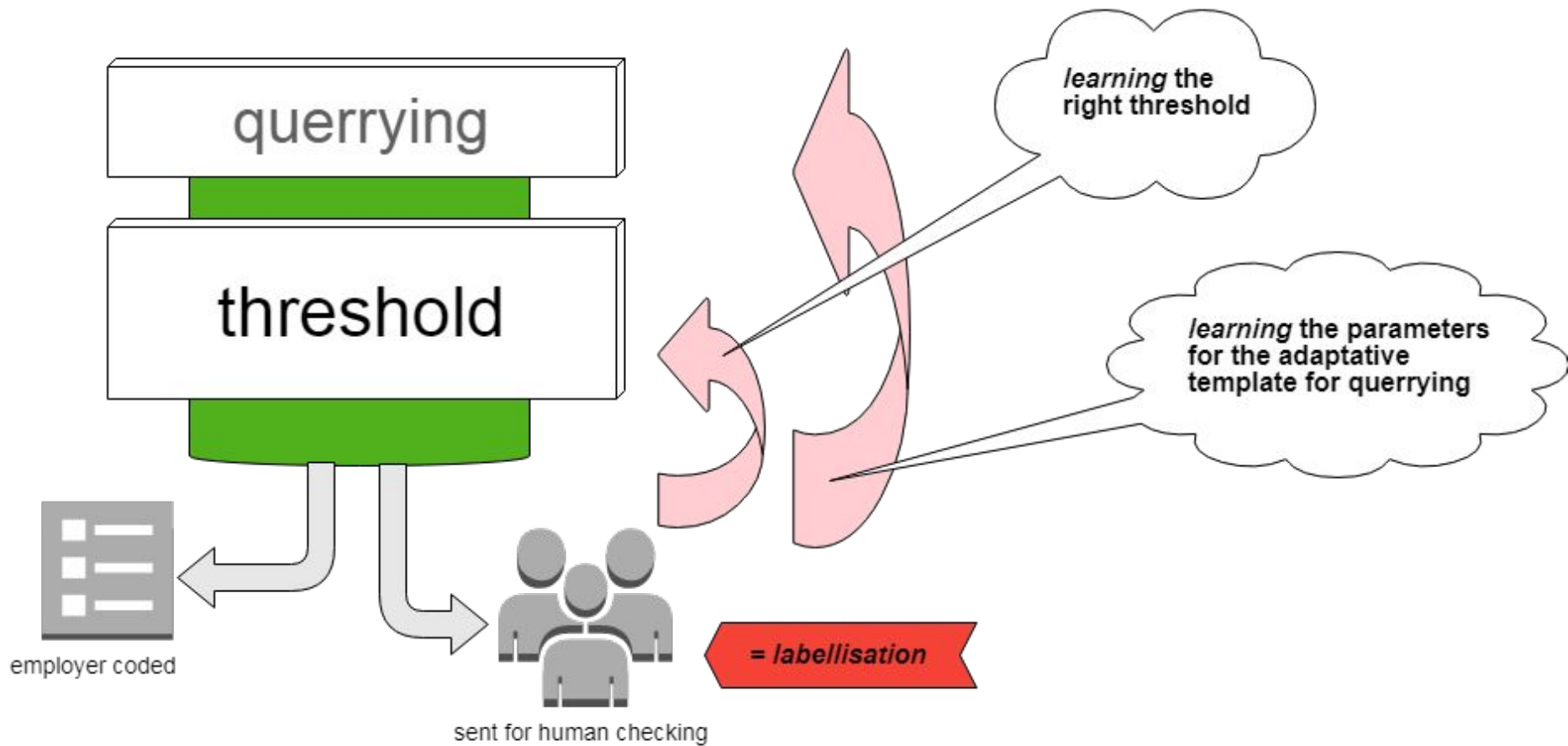
Geocoding : adding coordinates with quality of geocoding

National address register









conclusions

final word is yet to come

- actual solution leaves room for improvement
- enriching the data with many additional sources
- indexing for scalability
- identification in the firm register has many other applications
- identification in a register is a common issue

a general process of modernisation

the brand new “SSP lab” to
foster innovation for official
statistics in France

- hackathon last winter
- promoting data science in the NSI and around
- on the long run create frameworks that allow to capitalize on the work produced by gestionnaires