



RISIS

Research infrastructure for research
and innovation policy studies



The nano S&T dataset

Philippe Larédo, Lionel Villard, Michel Revollo

Marne la Vallée, 8/10/2015

UNIVERSITÉ —
— PARIS-EST



The presentation



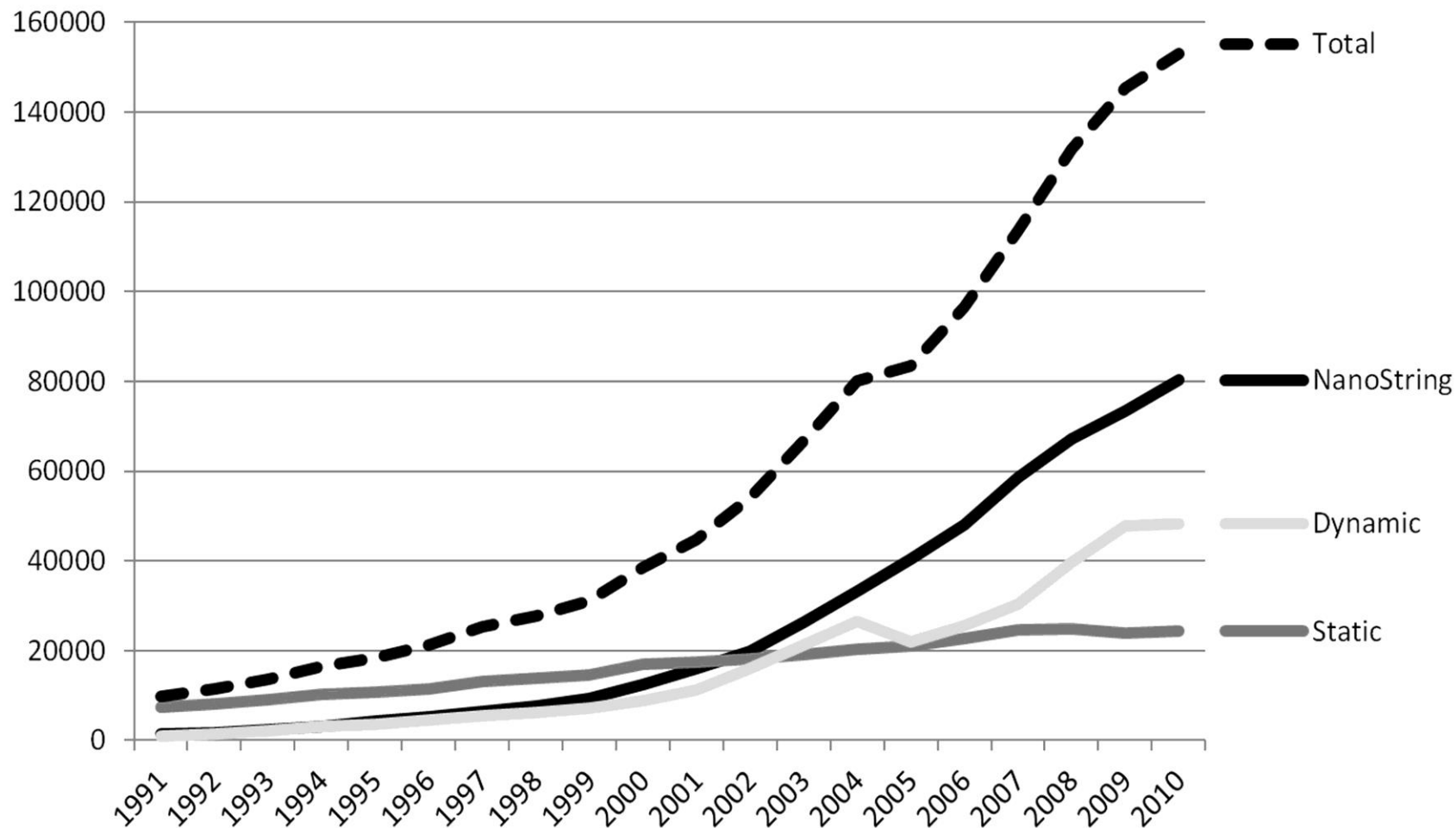
- Why this dataset
- 3 major developments for building the DB
- The DB structure
- Some perspectives

Why nano S&T DB (2)



- But multiple questions about what do we speak about
- Some articles (even recent ones) speak of some thousands patents while others speak of some hundreds thousands
- Are we speaking of a new science per se (emerging from interdisciplinary research) or of simply a new reduced scale of analysis for existing disciplines
- Some analysts consider that a new industry will emerge (like for computer science) or that an industry will be radically transformed (like biotechnology), others have different views (based on a renewed view of general purpose technologies)
- Where does it take place? In the same countries and the same places as usual? Or do we witness transformations?

An overview of results



Delineating 'nano' papers - 2



1- Building the core dataset

- The choice for a simple query for initial downloading
 - far faster to download
 - using then the different papers produced by colleagues for further targeted exclusions (e.g. plankton, flagel, or nanomolar)

The query for the nanostring

```
TI=((NANO* OR A*NANO* OR B*NANO* OR C*NANO* OR D*NANO* OR E*NANO* OR F*NANO* OR G*NANO*  
OR H*NANO* OR I*NANO* OR J*NANO* OR K*NANO* OR L*NANO* OR M*NANO* OR N*NANO* OR O*NANO*  
OR P*NANO* OR Q*NANO* OR R*NANO* OR S*NANO* OR T*NANO* OR U*NANO* OR V*NANO* OR  
W*NANO* OR X*NANO* OR Y*NANO* OR Z*NANO*) NOT (NANO2 OR NANO3 OR NANO4 OR NANO5 OR  
NANOSECOND* OR NANOLITER*)) OR TS=((NANO*) NOT (NANO2 OR NANO3 OR NANO4 OR NANO5 OR  
NANOSECOND* OR NANOLITER*))
```

Delineating 'nano' papers - 3



- Important results:
 - 517 000 articles from 1991 to 2010
 - a yearly growth of 20% during 15 years (40000 articles in 2005), doubling in 2010 (80000 articles)
 - overall the nanostring represents 14% of papers in 1991, 30% in 1999, 48% in 2005, stable since around 50%

Delineating 'nano' papers - 6



4- the static extension

- Building the external specificity of multi-terms
 - ratio presence in nanostring/total presence in WoS
 - done year by year, provides yearly sequence, individual year & average ratios
- Selecting multi-terms
 - select terms present over half of the period and part of the 250 highest termhoods (1105 terms identified),
 - rank them in descending order of external specificity
 - stop at the theoretical addition of articles nearest to the core set (517000 articles).
- Results:
 - we find the same **external specificity threshold** as in 1st dataset: 26%
 - static extension based on 114 multi-terms present 18 years out of 20 (with a skewed distribution: 13 terms bring 66%)

Analysis of multi-terms of the static extension



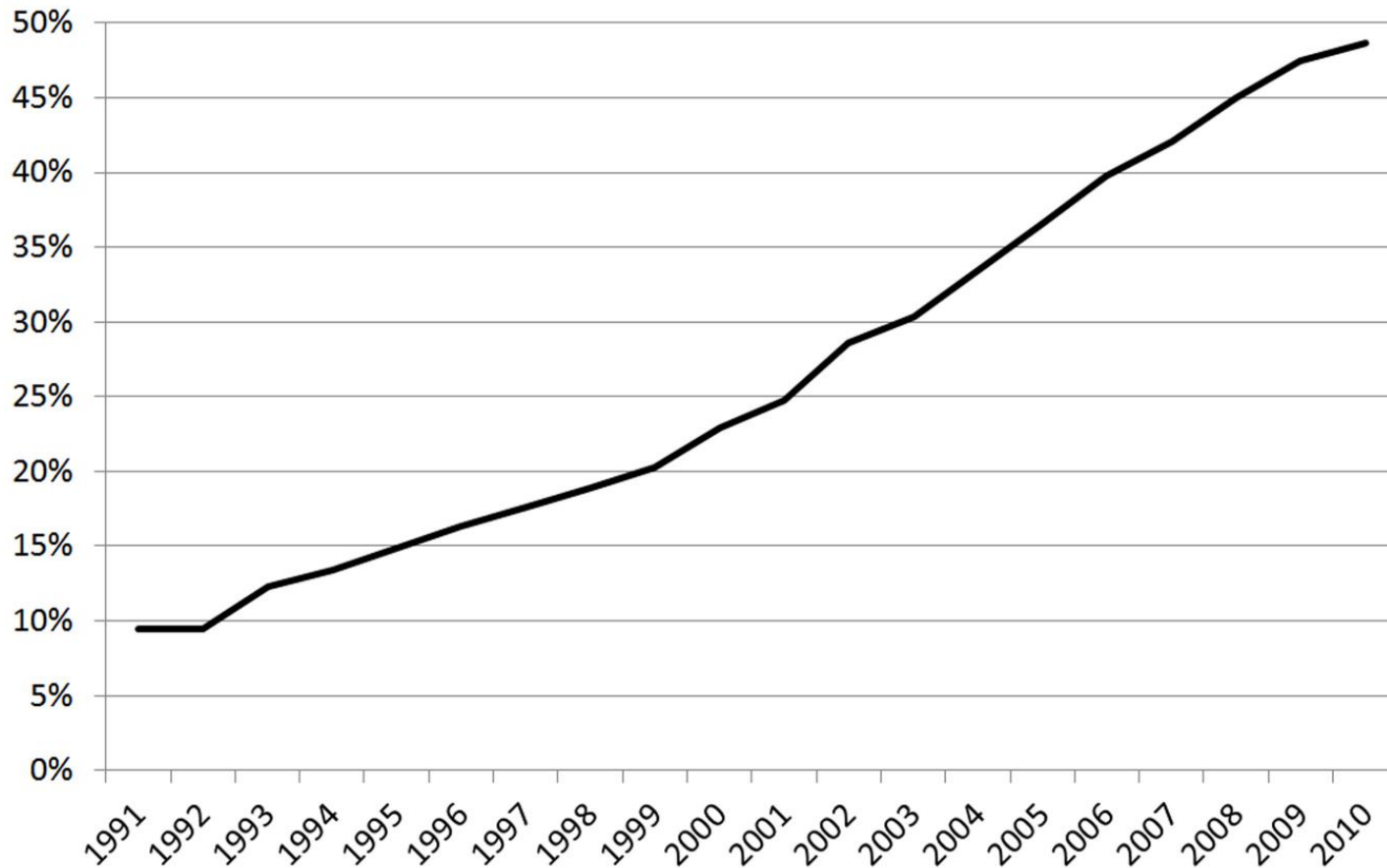
Type of multi-term	Nber of multi-terms	% of theoretical addition
Observation, manipulation & control (TEM, AFM, STM...)	30	57%
Materials (TiO ₂ , graphene,..) & CNT, nanowires	37	23%
Characteristics & properties of materials, molecules, genes	36	12%
Fabrication / expression techniques	11	8%
	114	100%

Analysis of multi-terms of the dynamic extension

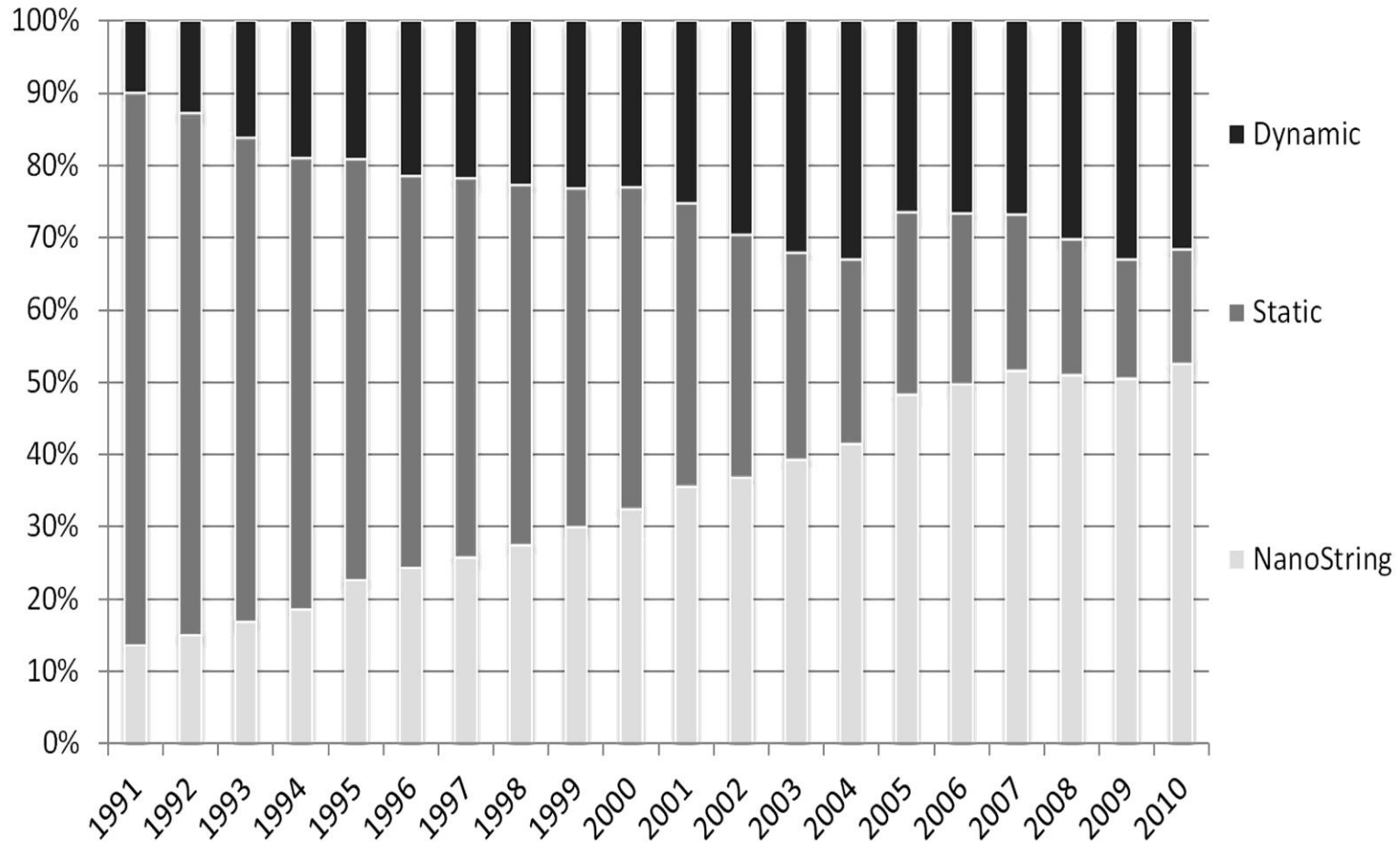


Type of multi-term	Nber of multi-terms	% of theoretical addition
Observation, manipulation & control (TEM, AFM, STM...)	32	19%
Materials (TIO2, graphene,..) & CNT, nanowires	38	28%
Characteristics & properties of materials, molecules, genes	28	14%
Measures	22	8%
Fabrication / Production	30	28%
Applications	12	3%
	162	100%

Evolution of the external specificity over time for the dynamic extension



The overall publication dataset



Moving from publications to patents 2

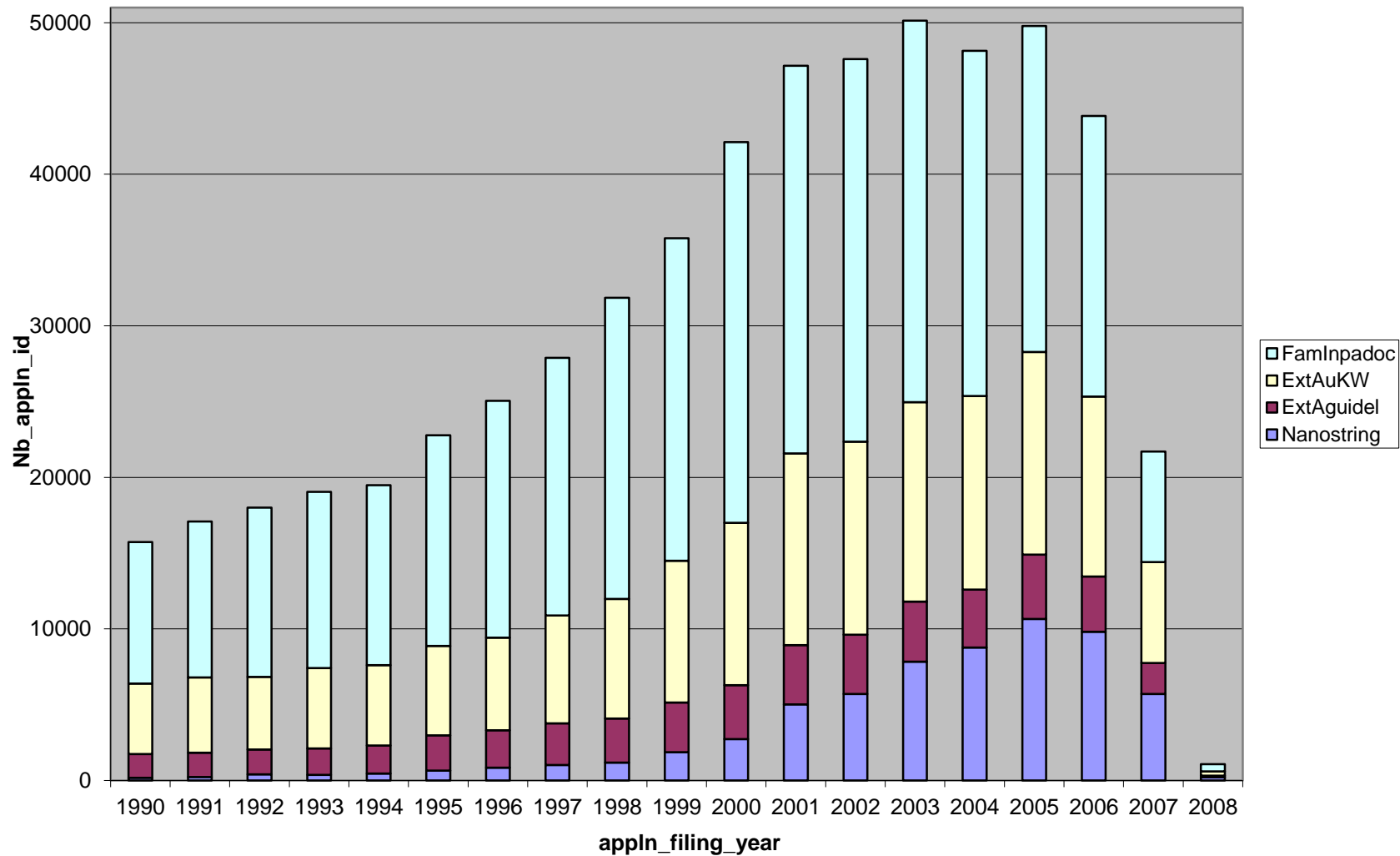


- NanoDB V1
 - we use the same vocabulary as for publications for selecting applications (on top of nanostring)
 - argument: Bonaccorsi (2007): 70% of nanoinventors are also academic authors.
- NanoDB V2 (under construction)
 - testing the method on 1 year: 80% of the vocabulary is different but 80% of the patents are the same...
 - testing the full method: works with the static extension (5 times more multi-terms), but not with the dynamic one (words too specific driving to very low external specificity thresholds)
 - developing a combined approach: using publication multi-terms not present in the patent static extension, selecting those above the average presence (5 years) and the average external specificity (25%).

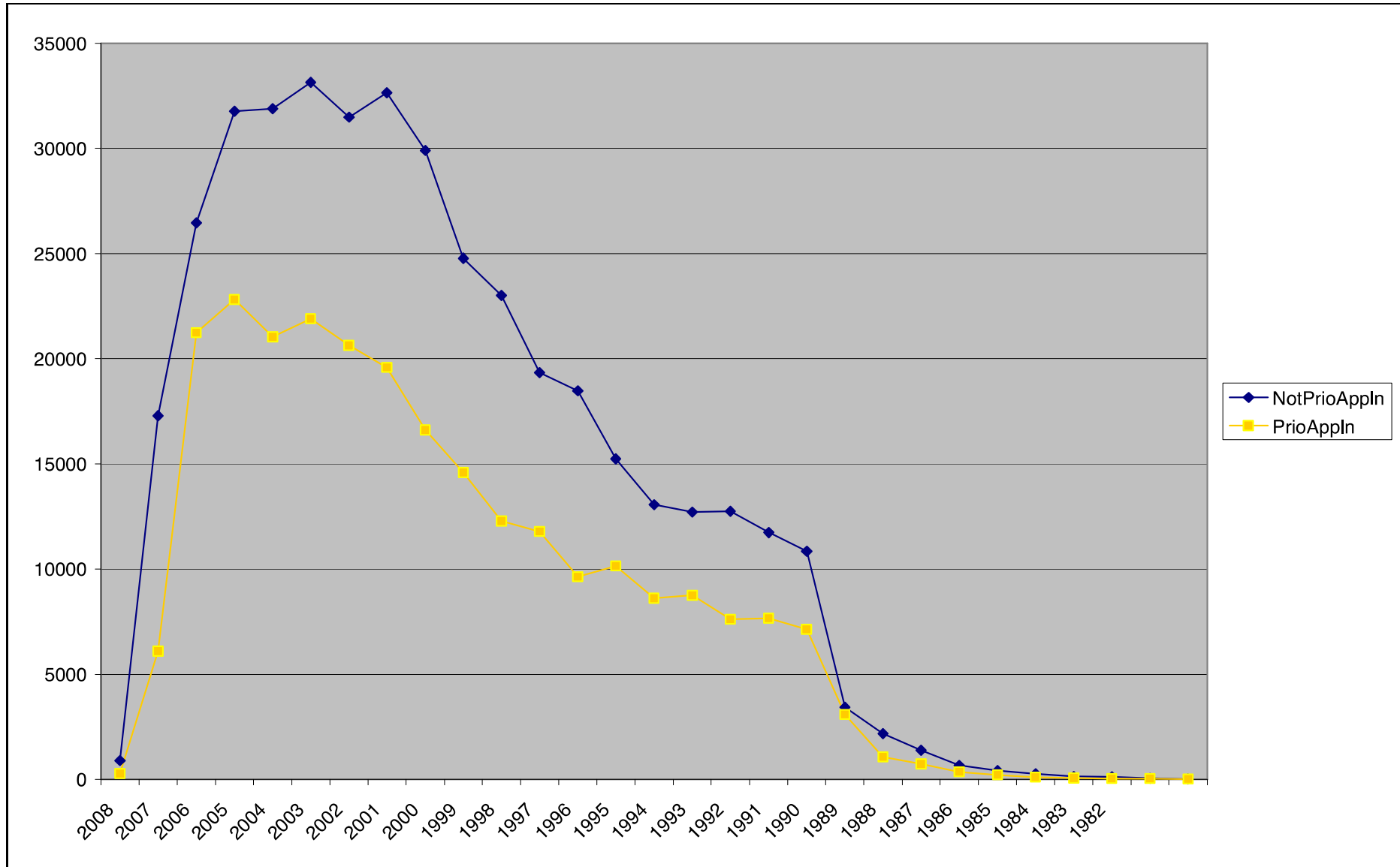
Nano patstat V1



- Construction in 3 steps
 - using patstat 2009 (ie until 2006)
 - step 1: extracting the patents using the nanostring → gives 63580 patents
 - step 2: extracting the extension following the 2007 approach (Mogoutov & Kahane) → gives 207000 patents
 - step 3: incorporating patents from INPADOC families → gives 327000 patents
- Results
 - overall 660000 patents but only **254000 priority patents, 190000 INPADOC families**



A fast growth followed by a plateau since 2000



The importance of 'one-off' patents



Nbr applications per INPADOC family	Nber of inpadoc families	Nber of priority patents	% priority patents
1 only (singletons)	97800	97800	38%
2 to 5	61216	74862	29%
6 to 10	21610	35800	14%
11 to 50	8540	29205	11%
51 to 100	314	4577	2%
101 to 500	164	6631	3%
501 to 1000	10	2476	1%
More than 1000	4	2733	1%
total	189658	254084	100%

Few countries concentrate priority patents ...

office	nber
US	93759
CA	798
AU	1929
RU	2827
IN	330
ZA	103
BR	Below 100

office	nber
JP	68937
KR	27060
CN	24798
TW	5076
SG	165

office	nber
DE	12568
GB	3863
EPO	3088
FR	2893
SE	697
IT	661
CH	414
NL	378
IL	337

...BUT

- Are all patents equal
- Can we trust offices as a source for identifying where knowledge is produced?

Strong country differences
 if we consider the 'value' of patents
 (i.e. those valorised are those 'extended' (not singletons))

office	total	% single-ton	Net total
US	93759	11,5%	82966
JP	68937	47,5%	36189
KR	27060	65,5%	9335
CN	24798	95,1%	1218
EU	25240	21,4%	19829

region	% crude total	% net total
US	38,3%	54,7%
JP	28,2%	23,9%
Other Asia	23,3%	8,3%
EU	10,3%	13,1%

D2-Localising inventive activities



- Principle: Moving from offices to the addresses of actors
- Choice: not stay at organisation's address and consider the location of inventors
- Limitation: many addresses not filled (wrongly located by Patstat, artificial patents, simply missing) → drove to important developments (OECD, Munich, Leuven, and Paris); see their integration in Patstat IFRIS.
- Hypotheses:
 - very strong agglomeration process in metropolitan areas
 - important collaborations mostly between metropolitan areas
 - countries thus not a good entry point to understand dynamics; need for developing clustering methods

Localising inventive activities - 2



- Three steps:
 - (i) reshaping addresses → see presentation on Patstat IFRIS
 - (ii) geolocalising addresses
 - (iii) clustering addresses
 - see overall chart next slide
- Here only key elements presented
- Reference for full presentation: Villard-Revollo, 2015
'geographical concentration of S&T activities', RISIS website

Methodological steps

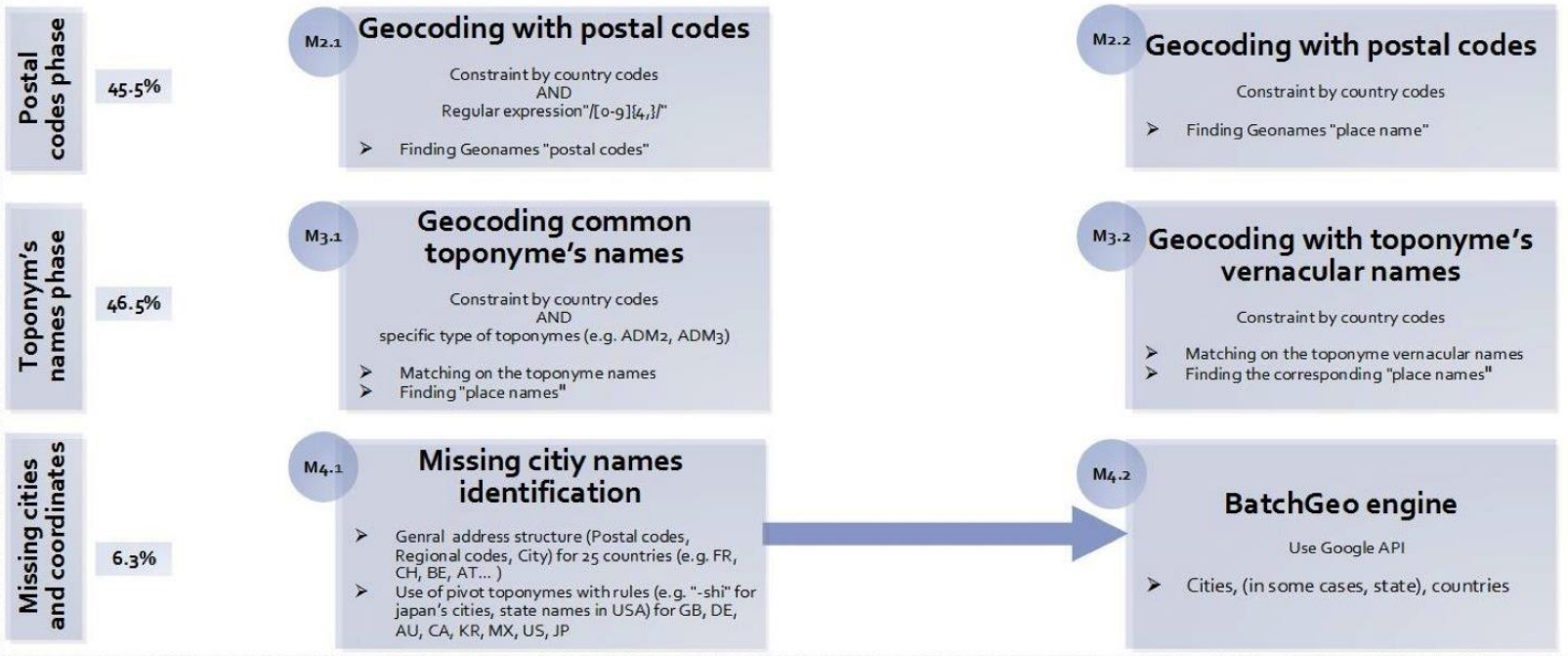
Pre-processing

Preparing addresses



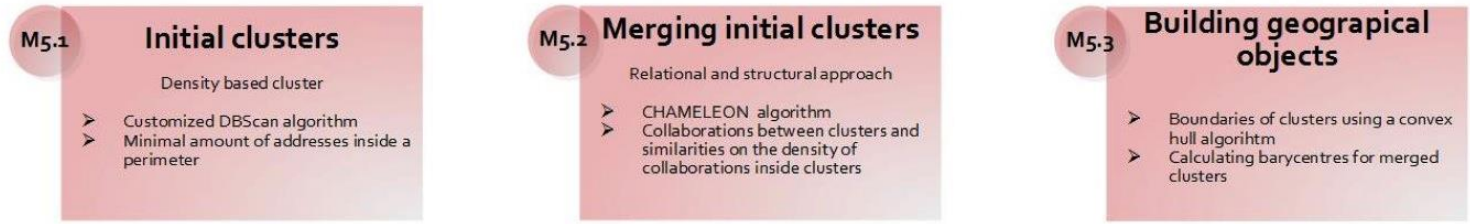
Geocoding

Identifying toponyms and attribution of the geographic coordinates



Clustering

Identification of the areas of aggregation of the S&T activities



Geolocalisation: 3 phases



- use of postal codes by comparing extracted postal code with GeoNames postal codes (over 900000 codes proposed with their place-name, latitude & longitude) → 45.5% of addresses geocoded
- use of toponymes by comparing extracted city names with corresponding toponymes of GeoNames (using country code & selected administrative & PPL features of GeoNames (for disambiguation) → 46.5% of addresses geocoded
- Dealing with non identified addresses → specific address structure identified for 35 countries (including state/regional division for 8) & submitted to BatchGeocode (using the 9 level accuracy filtering) → 6.3% of addresses geocode
- Result: Only 1.7% of total addresses not geocoded

Geolocalisation: an important comment



- As the geolocalisation process works well, the coverage of the DB depends on the existence of addresses → we decided recently to **extend coverage**
- Even with Patstat IFRIS solutions for additions (addresses in other Patstat fields, Regpat & other additions, propagation for artificial patents), we have only **50% of addresses filled**.
- We thus decided of a **logical extension** (called standard name extension) after clusterisation and organisation identification, **in 2 steps**
 - (i) allocating the institution address to authors of singleton applications: brings 7% more addresses
 - (ii) extending the geocoding within 1 patent for non covered authors when all addresses for that patent are in one cluster and that there is a number of such patents by the same organisation in that cluster: adds 62% addresses to the initial set!
- Overall this drives to **88% total coverage of addresses** without changing the **coverage of patent applications: 85%**
- **Note: all data below do not take into account this extension**

Building clusters -2



- (ii) agglomerating interconnected clusters
 - in V1 we use the overlap between clusters in term of addresses (clusters are merged when overlap $> 20\%$)
 - in V2 we use the 'chameleon approach' (Karypis, Han & Kumar, 1999) which considers 2 parameters
 - * relative interconnectivity (measured as a ratio between bilateral links and internal links) under a maximum geographical constraint
 - * relative closeness: measures the similarity of collaboration profiles of the 2 clusters

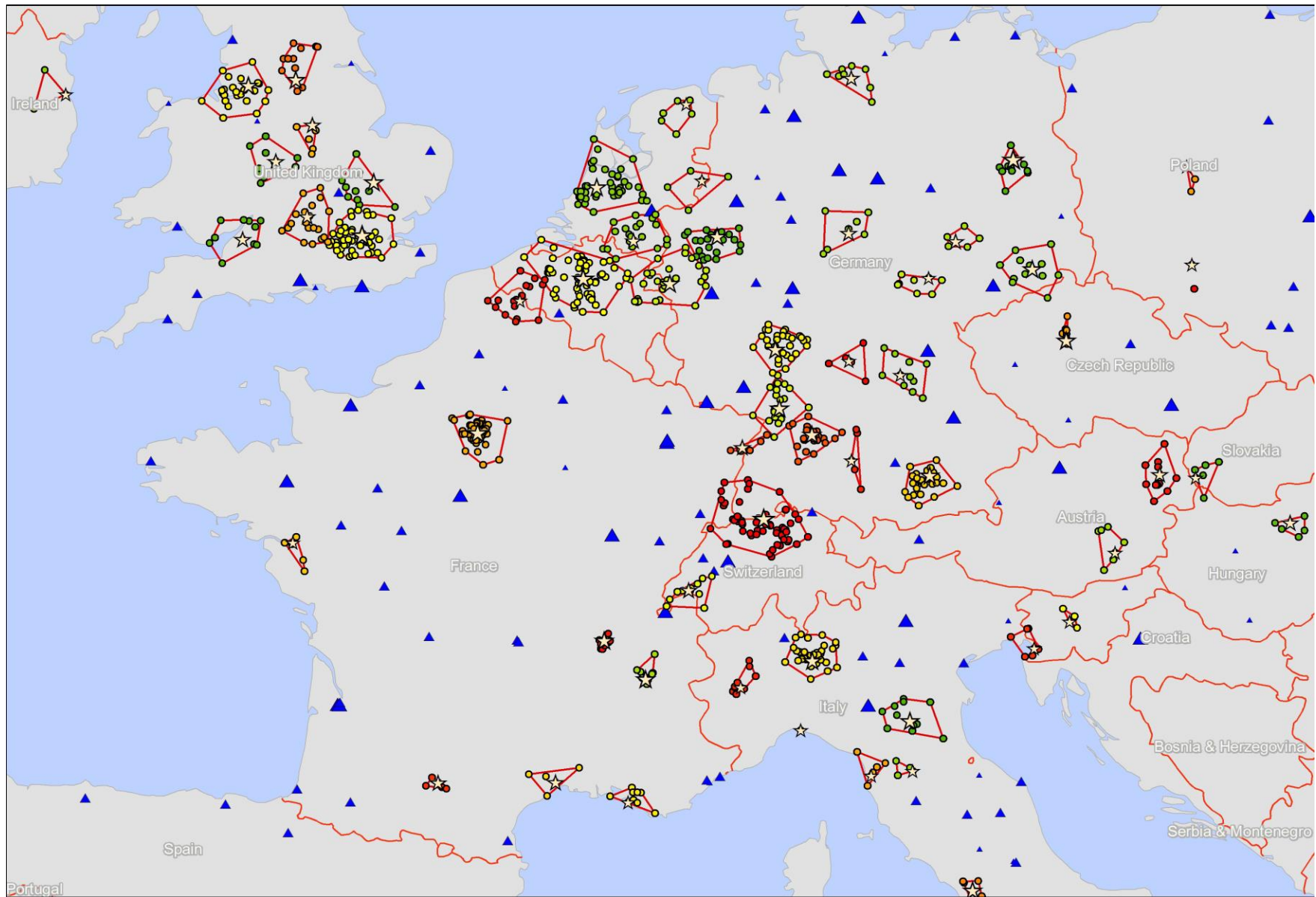
Building clusters -3



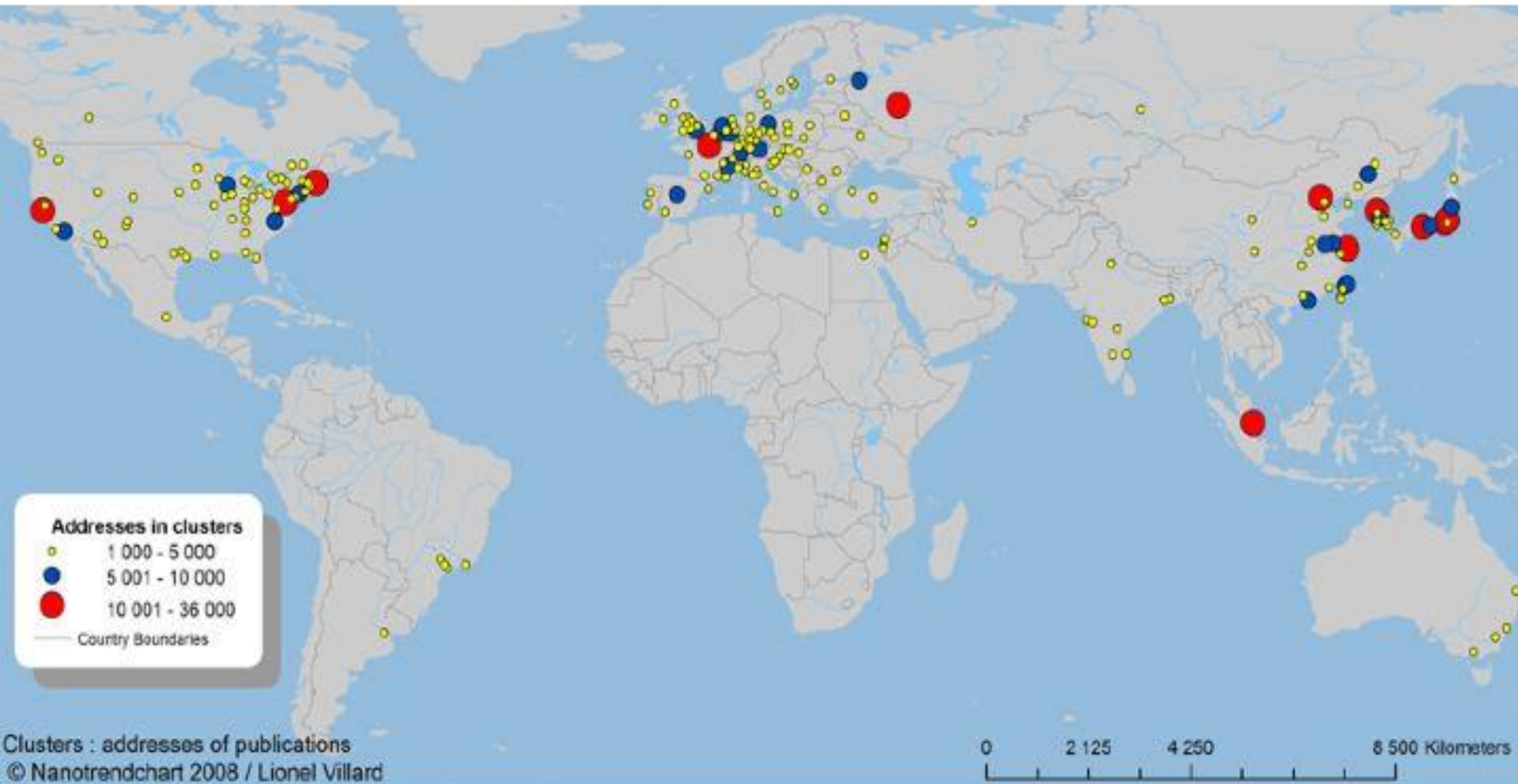
- We chose to build clusters on nano publications (1998-2006)
→ 203 clusters identified gathering 77% of addresses
- We applied the clusters to nano patents (1998-2006)*
→ the 203 clusters cover 75% of inventor addresses
- Some results:
 - A very high concentration (not considering singleton applis)
 - 10 clusters represent 53% of addresses (1998-2006) while the last 2/3rds represent only 8%
 - of the 52 top clusters (more than 1000 applis), 48% are from the 25 US clusters (San Francisco representing nearly 1/3rd alone), 40% from the 12 Asian clusters (Tokyo representing 40%) and 11% from the 14 EU clusters (8 being German)
- On-going developments: look at complementary clusters (patent specific, or by adding publications & patents)

* Not taking into account the standard name extension

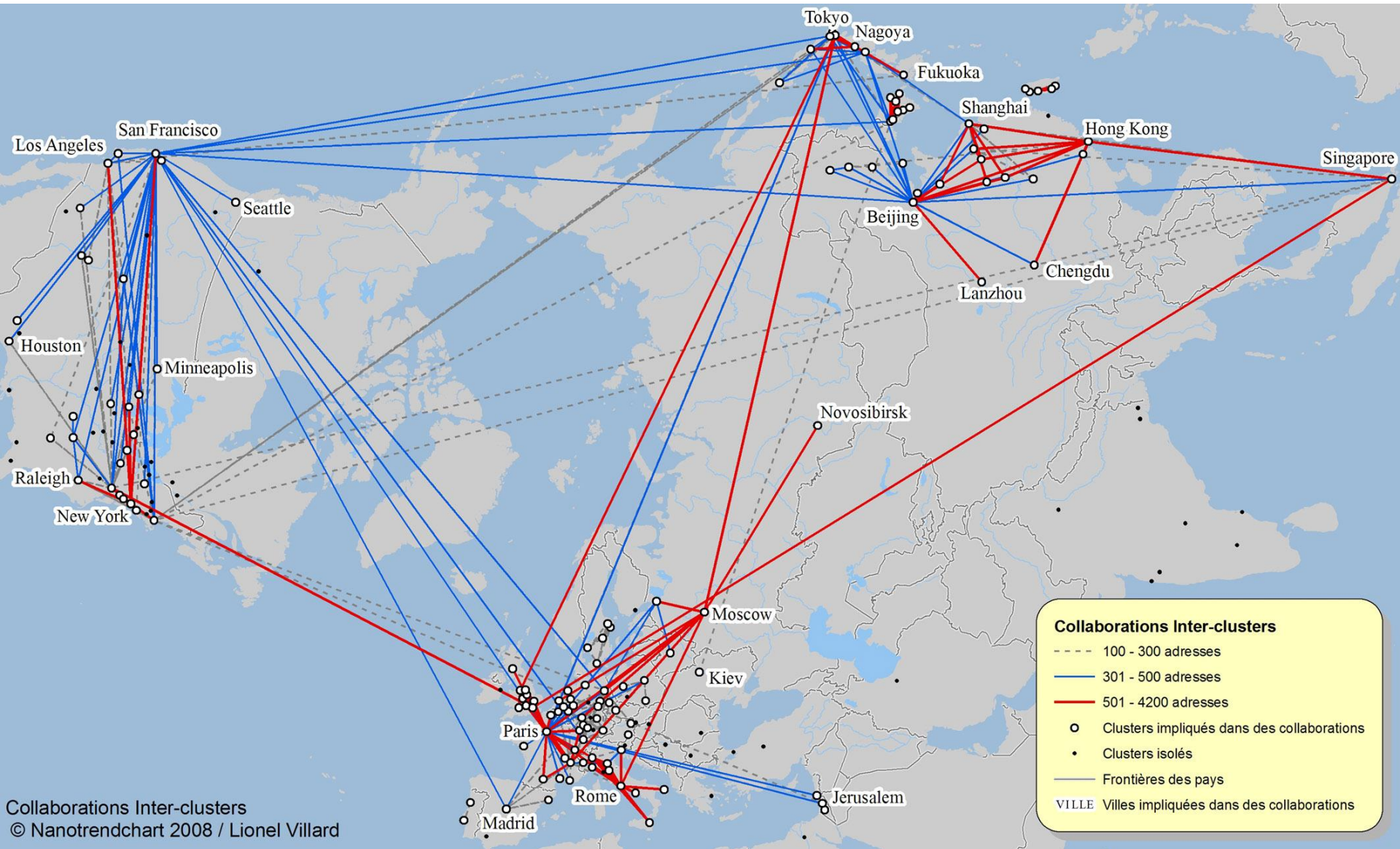
The construction of clusters: Europe as an example



The 203 clusters identified based on 1998-2006 outputs



Intercontinental collaborations concentrated in a few clusters



Characterising institutions- 1



- Based on Patstat standard names
- A lengthy manual process after using classical automated tools/
 - initial markers enabling to tag type (firm, univ, govt including PRO, other incl. hospitals)
 - addressing name similarities (manual check of proposals by automated treatment)
- An important support over time: the reference databases built in house
 - in particular CIB for large firms & their subsidiaries,
 - using other projects (on specific countries or sectors) in particular for PRO (e.g. handling Helmholtz institutes)
 - remaining problems: handling university hospitals, level of agglomeration (e.g. US DOE or individual labs from DOE)

Characterising institutions- 2



- Time consuming approach to both identifying organisations & finalizing tagging:
 - (i) work within clusters, and for organisations not in clusters, work at country level.
 - (ii) work on aggregated data per organisation for checking (insures robustness of large actors)
- Some results
 - 13300 different organisations, 61 firms with 1000+applis (37% total)
 - Patents with more than 1 assignee represent around 10% of total patents, 60% are other firms, 20% public assignees and 20% individuals
 - 87% of patents filed are by firms 7% by universities, 5% by PRO (and governments), 1% by others (including hospitals)

Top applicants (1998-2006)*



Samsung	13021	CEA	1636	Univ california	2243
Mitsubishi	7590	JST (Japan)	1571	Univ Tsinghua	825
LG	7565	KIST (Korea)	1391	Univ Texas	695
Sumitomo	5923	ITRI (Taiwan)	1341	Univ shanghai jiaotong	547
IBM	4884	Fraunhofer	1209	Caltech	512
Seiko Epson	4828	AIST (Japan)	1095		
Sony	4531	CNRS	943		
		CSIRO (Australia)	636	Harvard	332
BASF	3927			Univ oxford	131
3M	3525			Univ Cambridge	169
Gen Electric	3375				

* 2010 treatment

Some EU examples of Cluster composition*



	global	Paris	Zurich/Bale**	Leuven	Berlin	Atlanta	
Firm	87%	74%	67%	75%	63%	68%	
Univ	7%	3%	23%	24%	7%	31%	
PRO/gov	5%	22%	9%		29%	1%	
other	1%						
		77 firms			80 firms	110 firms	

* Only taking clusters around 1000 priority patent applications (1998-2006)

* Transborder cluster (PRO= Fraunhofer!)

Is nano linked with specialised industries?



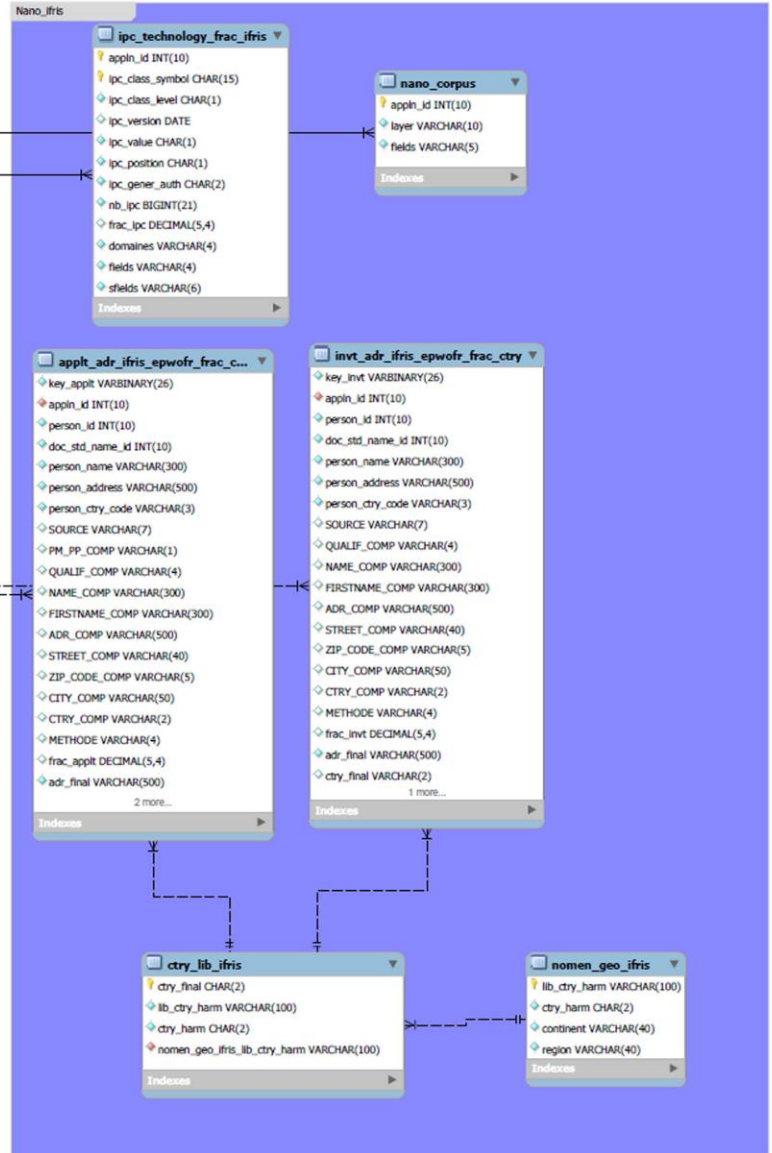
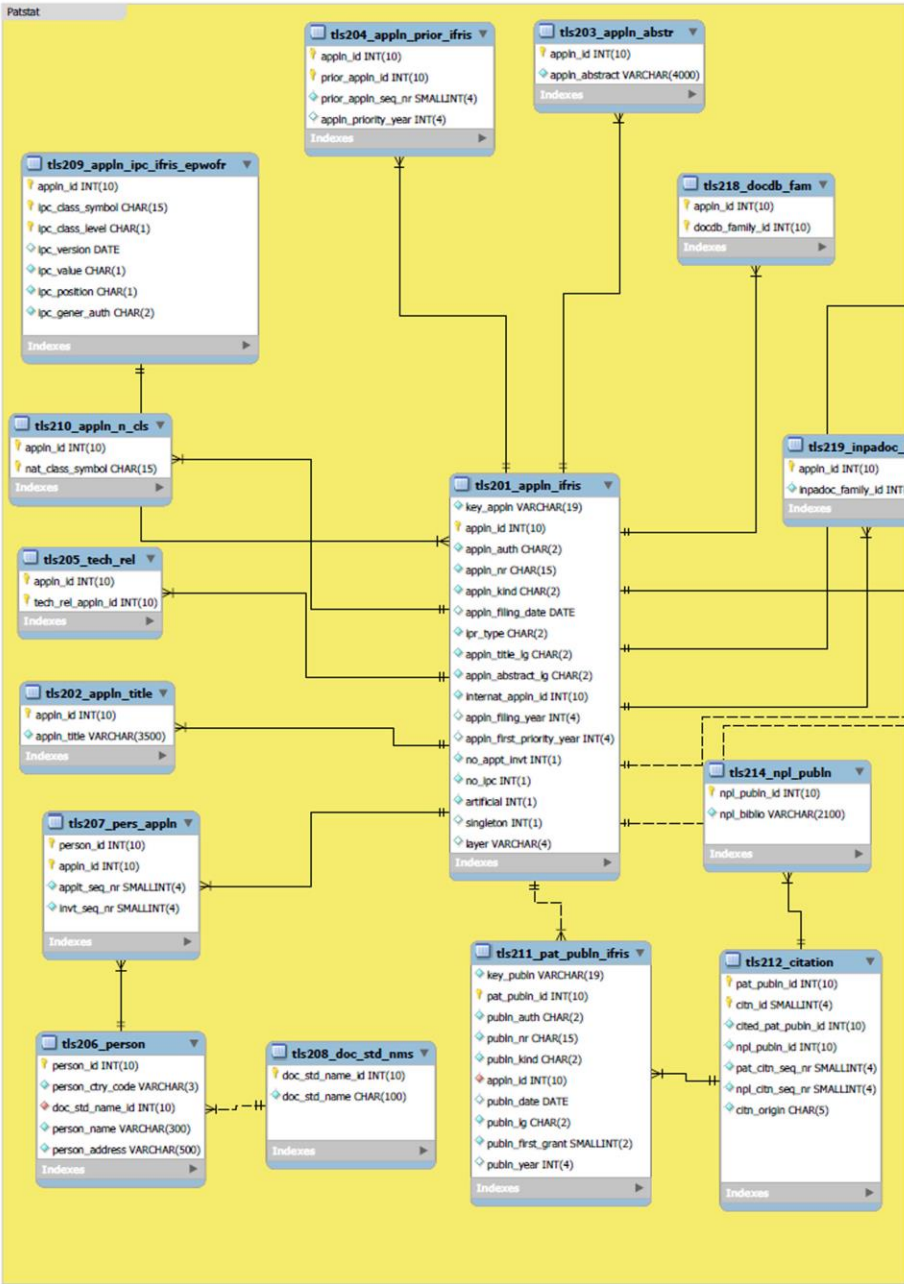
firmes du DTI scoreboard	total	nano	%
Electronic & electrical equipment	103	70	68%
Technology hardware & equipment	226	150	66%
Chemicals	96	84	88%
Pharmaceuticals & biotechnology	153	73	48%
Health care equipment & services	53	39	74%
Automobiles & transport	86	59	69%
Aerospace & defence	35	24	69%
Materials & construction	55	42	76%
Oil, Gas & Electricity	53	39	74%
Food producers inc. Beverages)	32	16	50%
General industrials	38	24	63%
Household & personal goods	40	21	53%
Industrial engineering	70	35	50%
Telecom & media	32	14	44%
Software & computer services	110	14	13%
banks, insurance, retail, leisure	49	6	12%
total	1231	710	58%

Publications: 3 main domains with a central role for Chemistry and materials

Biotechnologies
Life Science
12% total
+112% in 8 years

Electronics Physics
34% total
+104% in 8 years

Chemistry & Materials
52% total
+170% in 8 years



Next steps



- Nano Patstat V1 is available with all its improved characteristics (covering until 2006)
- In 6 months it will be complemented by Nano Patstat V2
 - based on Patstat IFRIS 2014 (until 2012)
 - with full geolocalisation process
 - completely harmonised with CIB for large firms
 - with new approach for clustering (a standard choice, but can be refined for ad-hoc questions)
- We shall progressively integrate RISIS features: identifying midsize and venture capital supported start-up firms; providing OECD standardised metropolitan areas (FUA)

