



# **NAMSOR DEEP LEARNING FOR MIGRATION ANALYTICS : DECRYPTING IDENTITY IN SPACE AND TIME THROUGH PERSONAL NAMES, GEOGRAPHIC, SEMANTIC, SOCIAL GRAPH**

2018-01

Elian CARSENAT, NamSor

# Founder Bio



2

**Elian CARSENAT**, a computer scientist trained at ENSIIE/INRIA, started his career at JP Morgan in Paris in 1997. He later worked as consultant and managed business & IT projects in London, Paris, Moscow and Shanghai.

In 2012, Elian created **NamSor**, a piece of sociolinguistics software to mine the 'Big Data' and better understand international flows of money, ideas and people.

<http://fr.linkedin.com/in/eliancarsenat/en>

# NamSor sorts Names



3

- Classification with various taxonomies
  - Gender (female/ male / unknown)
  - Script (LATIN, ARABIC, GUJARATI,...)
  - Origine (Country ex. France vs. Inde)
  - Region (ex. Gujarat vs. Andhra Pradesh)
  - Diaspora (ex. Indian Diaspora in US vs Indian Diaspora in Mauricius)
- Sorting according to a numerical score, allowing combining NamSor with other algorithm (graph, semantics, predictive ...)
- Flexibility to learn new taxonomies (machine learn.)
- Ease of integration (NamSor API, Java/Python SDK, ESRI, RapidMiner, NationBuilder ...)



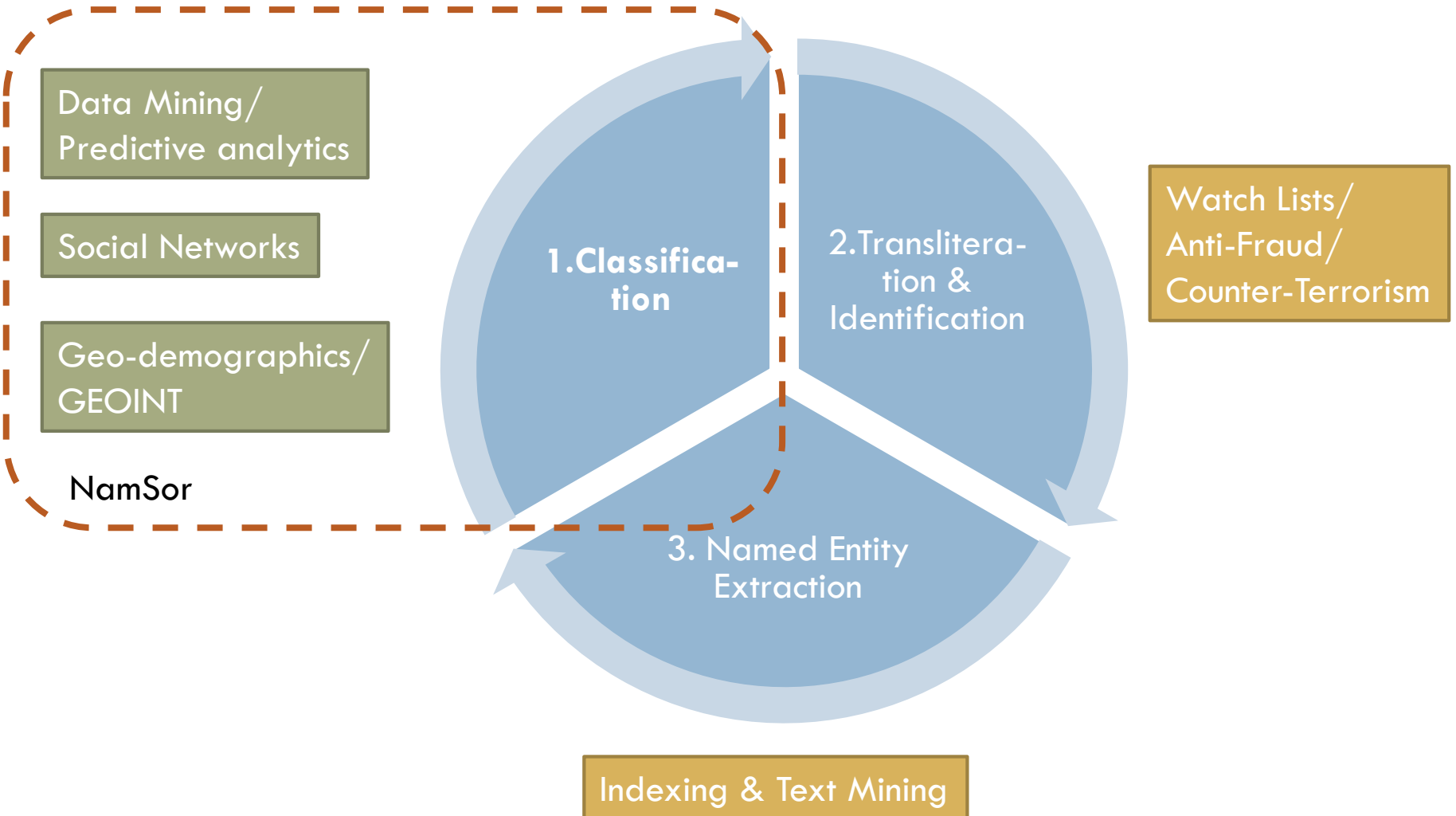
# A global coverage - 142+ countries

4

DIMENSION	CURRENT COVERAGE
SCRIPT (22)	LATIN, ARABIC, CYRILLIC, ARMENIAN, BENGALI, DEVANAGARI, GEORGIAN, GREEK, GUJARATI, GURMUKHI, HAN, HANGUL, HEBREW, HIRAGANA, KANNADA, KATAKANA, MALAYALAM, MYANMAR, ORIYA, TAMIL, TELUGU, THAI
COUNTRY (142+)	AE, AF, AL, AM, AO, AR, AT, AZ, BA, BD, BE, BF, BG, BH, BI, BJ, BN, BR, BT, BW, BY, CA, CD, CF, CG, CH, CI, CL, CM, CN, CO, CR, CV, CY, CZ, DE, DK, DZ, EE, EG, ER, ES, ET, FI, FJ, FR, GA, GB, GE, GH, GM, GN, GR, HK, HR, HT, HU, ID, IE, IL, IN, IQ, IR, IS, IT, JO, JP, KE, KG, KH, KM, KP, KR, KW, KZ, LA, LB, LK, LR, LS, LT, LU, LV, LY, MA, MD, ME, MG, MK, ML, MM, MN, MR, MU, MV, MW, MX, MY, MZ, NA, NE, NG, NL, NO, NP, OM, PE, PH, PK, PL, PS, PT, QA, RO, RS, RU, RW, SA, SD, SE, SI, SK, SN, SO, SR, SY, TD, TG, TH, TJ, TM, TN, TO, TR, TT, TW, TZ, UA, UG, US, UZ, VE, VN, YE, ZA, ZM, ZW
COUNTRY/ REGION (15)	RU (80), IN (~30), FR (22), IT (17), LB (14), BF (13), CD (8), TR (7), ID (7), GB (4), ES (17), ML (50), GN (8), CI (34), AF(16)
COUNTRY/ DIASPORA	US, CA, SG, GB, (EU)

# NamSor can enrich any nominative data

5



# Two complementary approaches

6

## NamSor CORE (Origin, Diaspora)

- Optimized for global coverage : coding names to a large multi-class taxonomy (all countries / regions / ethnicities)
- The **only input is NAMES** : not other information is required

## NamSor ML

- Deep-learning capability to re-train models towards a focused research or a customized taxonomy (binary classifier, or just a few classes)
- **Name information is combined with other data** (geographic, behavioural, semantic ...)



# NamSor Core : Mapping Tunisian Diaspora

## La BIAT lance un road show pour les tunisiens en île de France

Le 17 mars 2016, BIAT France lance le « BIAT France Tour » et part à la rencontre des Tunisiens résidents à Paris et en région parisienne afin de leur présenter les produits et services qui leur sont destinés.



C'est aux 720 000 Tunisiens vivant en France que la BIAT s'adresse avec sa filiale BIAT France, sous la signature « Ici pour vous », un service de transfert d'argent leur est proposé à des prix très compétitifs et dans des conditions de rapidité et de sécurité exemplaires.

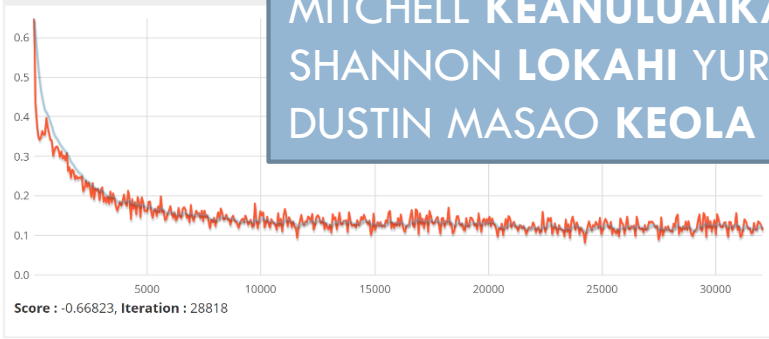


# NamSor ML : Native Hawaiian Names - A Binary Classifier Example

8

- Overview
- Model
- System
- Language

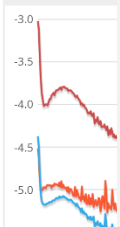
Model Score vs. Iteration



Multi-cultural names : NHawaiian, Filipino, Japanese, Anglo-Saxon ...  
**MITCHELL KEANULUAIKALANI KALUHIWA**  
**SHANNON LOKAHI YURONG ESTOCADO**  
**DUSTIN MASAO KEOLA IWABUCHI ...**

Last Update	2017-11-20 09:46:31
Total Parameter Updates	32071
Updates/sec	.82
Examples/sec	824.20

Update:Paran



Nov 20, 2017 9:49:39 AM com.namsor.namsorml.us\_hi.MLPNamSorHawaiiTrainer main

INFO:

- Examples labeled as 0 classified by model as 0: 21857 times
- Examples labeled as 0 classified by model as 1: 7731 times
- Examples labeled as 1 classified by model as 0: 2194 times
- Examples labeled as 1 classified by model as 1: 221318 times

====Scores=====

# of classes: 2  
 Accuracy: 0.9608  
 Precision: 0.9375  
 Recall: 0.8644  
 F1 Score: 0.9781

Nov 20, 2017 9:49:39 AM com.namsor.namsorml.us\_hi.MLPNamSorHawaiiTrainer main

INFO: \*\*\*\*\*Example finished\*\*\*\*\*

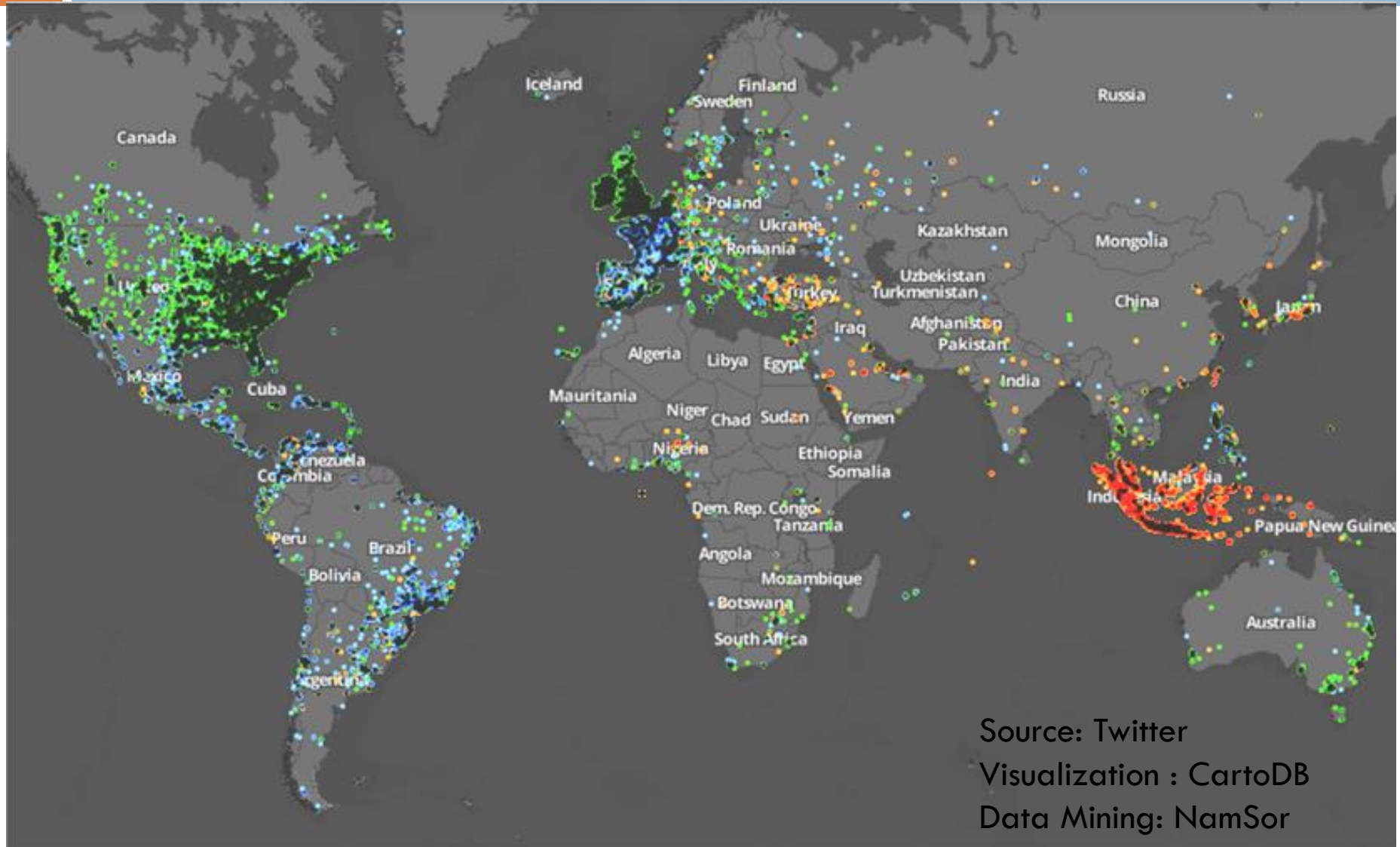


- NamSor use cases

# Mining 3M twitter names to map *Diasporas*

*Who are they, where are they and what are they doing?*

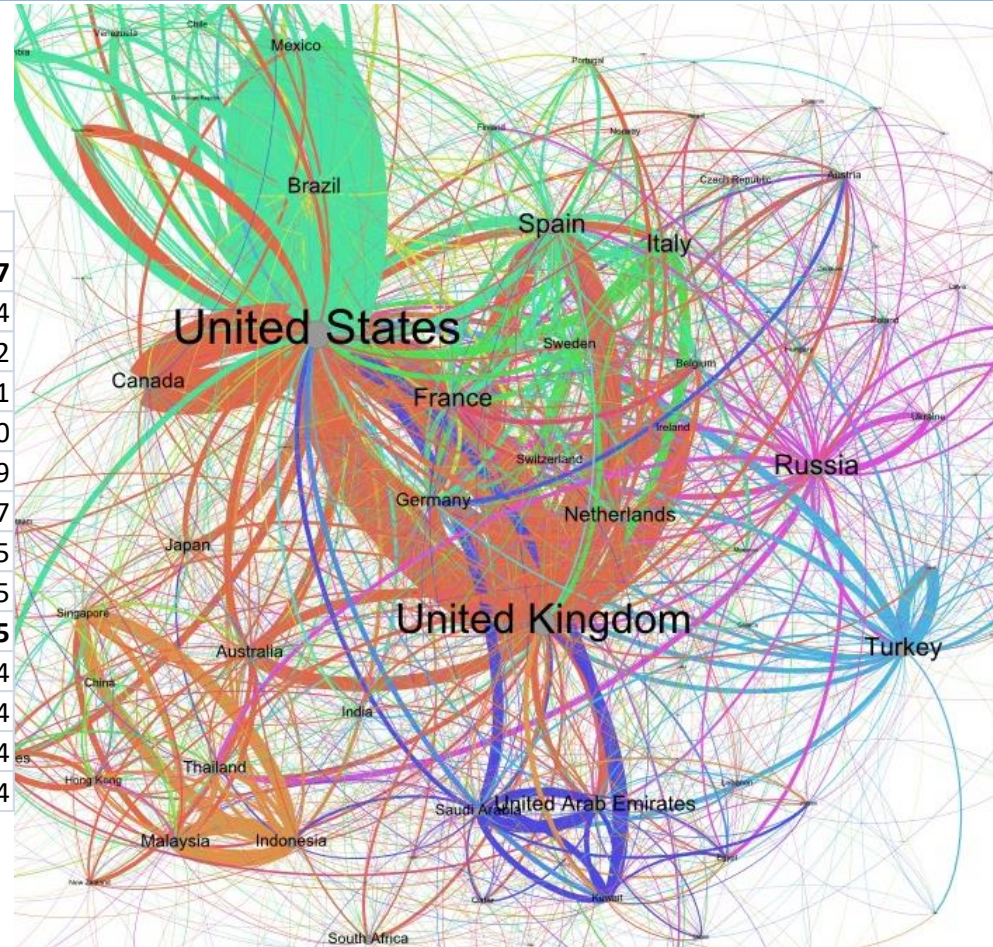
10



# Flow view – who travels where?

11

Source	Target	Type	Id	Onoma	Weight
<b>United Kingdom</b>	<b>France</b>	<b>Directed</b>	<b>16</b>	<b>Great Britain</b>	<b>37</b>
Spain	France	Directed	55	Spain	14
United States	France	Directed	75	Great Britain	12
Turkey	France	Directed	79	Turkey	11
Brazil	France	Directed	87	Portugal	10
United Kingdom	France	Directed	112	Ireland	9
Italy	France	Directed	152	Italy	7
Switzerland	France	Directed	226	France	5
Belgium	France	Directed	247	France	5
<b>United Kingdom</b>	<b>France</b>	<b>Directed</b>	<b>258</b>	<b>France</b>	<b>5</b>
Mexico	France	Directed	287	Spain	4
Ireland	France	Directed	317	Great Britain	4
United Kingdom	France	Directed	333	Italy	4
United States	France	Directed	375	France	4



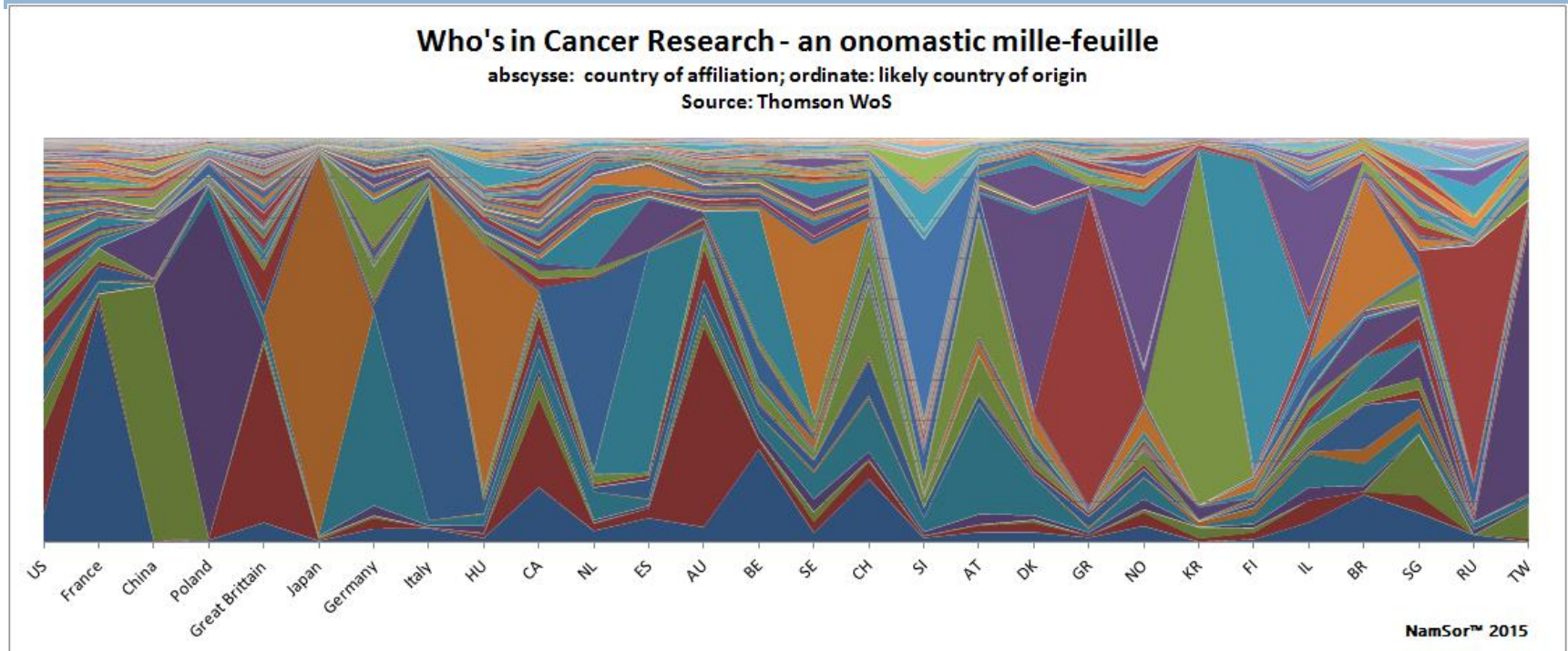
Source: Twitter  
 Visualization : Gephi  
 Data Mining: NamSor



# Mapping Talents in Cancer Research

(in collaboration with French INSERM)

12



## Thomson Reuters WebOfScience (6 countries, 250k scientists, 50k papers)

“Analysts uncovered amazing patterns in the way scientists’ names correlate with whom they publish, and who they cite in their papers - not just in case of a particular country, but globally. Tania Vichnevskaja of the French National Institute for Health (INSERM) presented the paper ‘Applying onomastics to scientometrics’ at IREG International symposium 2015 organised by University of Maribor and Shanghai Jiao Tong University. The paper was prepared jointly with NamSor, a private start-up company specialized in mapping international Diasporas.”

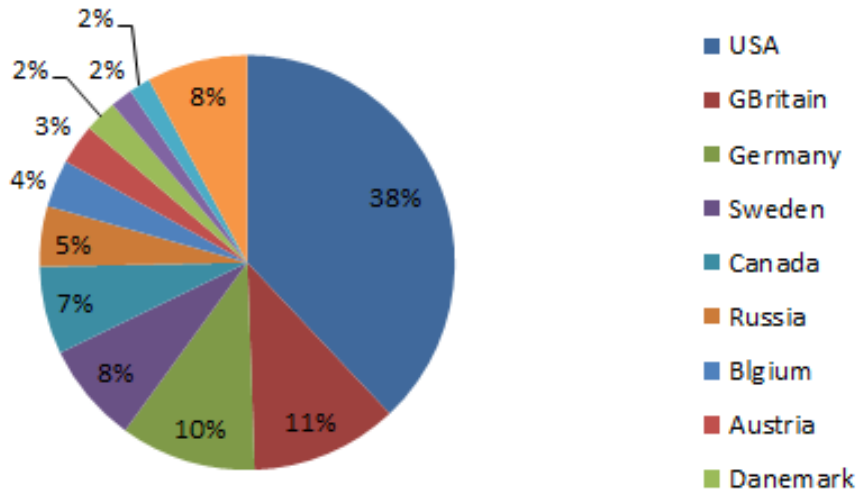
Source: WoS; Data Mining: INSERM with NamSor

# Cancer Research in Poland and Slovenia

## Examining the 'brain drain'

13

### The Polish "brain drain"

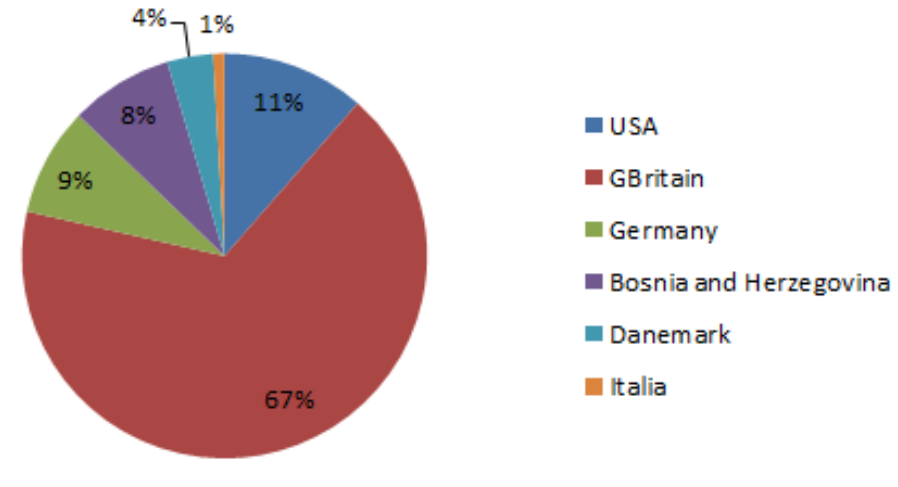


In the Polish Corpus, we look at co-authors with Polish names, affiliated abroad.

Top countries:

1. **US,**
2. **Great-Britain,**
3. **Germany.**

### The Slovenien "brain drain"



In the Slovenian Corpus, we look at co-authors with Slovenian names, affiliated abroad.

Top countries:

1. **Great-Britain,**
2. **US,**
3. **Germany.**

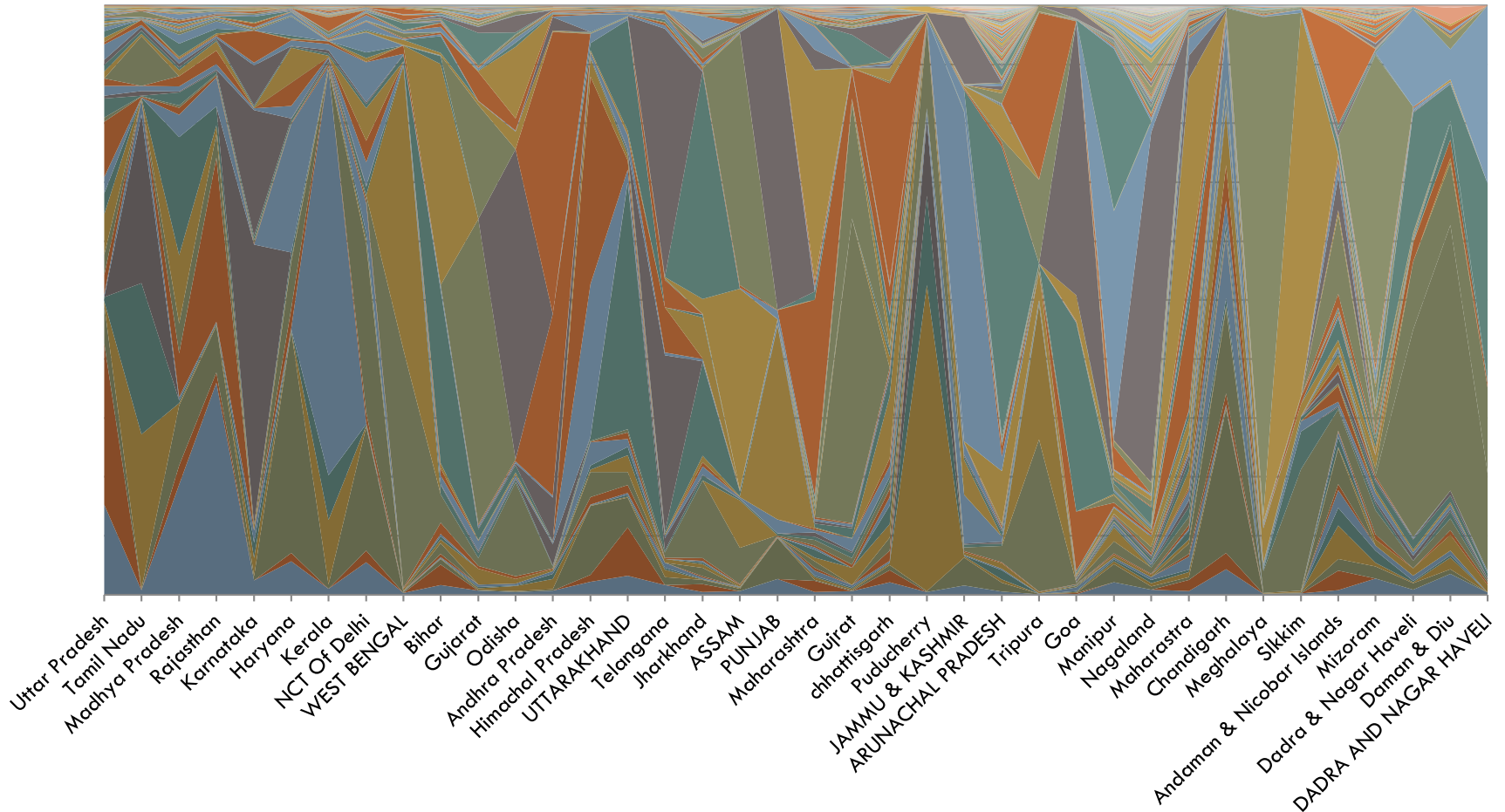


# “Incredible India” – 1.2 BN People

## Indian onomastics by State/Union Territory

14

Names in LATIN, BENGALI, DEVANAGARI, GUJARATI, GURMUKHI, KANNADA, MALAYALAM, ORIYA, TAMIL, TELUGU, ARABIC

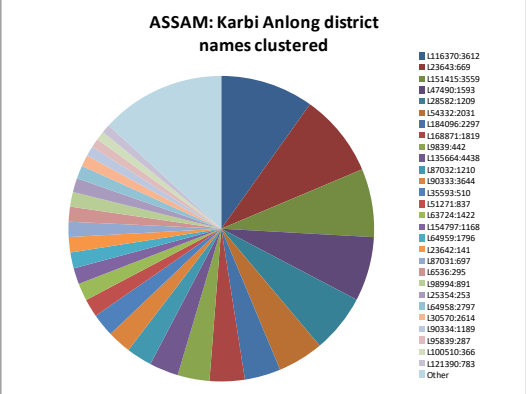


# ASSAM: Karbi Anglong, within district Inter-caste marriages ?



15

output	Input	Input
clusterId clusterParentId	Firstname LastName parent is	FirstParen LastParen
L25354:25! L64958:2797	Aṁ 1 1}ššā husband	ṁṁṁṁ ṁṁṁṁ
L47490:15! L64958:2797	ṁṁ 1 ṁṁṁṁ father	ṁṁṁṁ ṁṁṁṁ
L28582:12! L47490:1593	ṁṁ >à ṁṁṁṁ ššā husband	ṁṁṁṁ 1 ṁṁṁṁ
L23643:66! L35593:510	ṁṁṁṁ à ṁṁṁṁ ššā father	ṁṁṁṁṁṁ ṁṁṁṁ
L23643:66! L35593:510	3 à à ṁṁṁṁ ṁṁṁṁ ššā father	ṁṁṁṁṁṁ ṁṁṁṁ
L47490:15! L35593:510	ṁṁṁṁ ṁṁṁṁ ṁṁṁṁ father	ṁṁṁṁ ṁṁṁṁ
L23643:66! L35593:510	Aṁ 1 t ṁṁ 1 ššā husband	ṁṁṁṁ ṁṁṁṁ
L35593:51! L47490:1593	ṁṁṁṁ š ṁṁṁṁ father	ṁṁṁṁ ṁṁṁṁ
L23643:66! L47490:1593	ṁṁ >à ṁṁṁṁ ššā father	ṁṁṁṁ ṁṁṁṁ



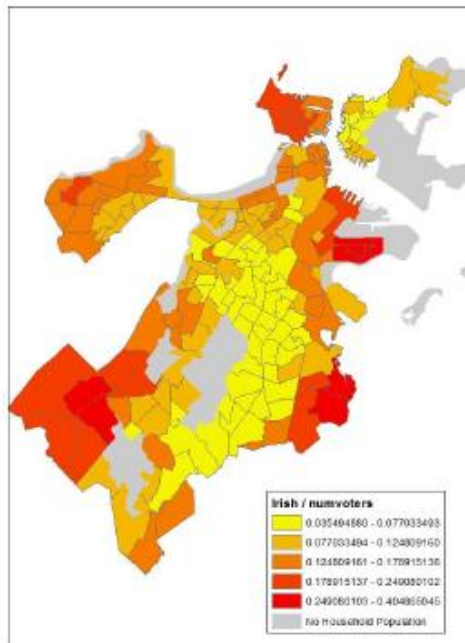
parent is	husband								
Count of serial Column Labels									
Row Labels	L47490:1593	L116370:3612	L54332:2031	L184096:2297	L35593:510	L168871:1819	L135664:4438	L51271:837	
L23643:669	6931	84	5099	15	2069	28	791	1924	
L151415:3559	18	212	11	6446	19	1217	55	6	
L28582:1209	5132	68	3565	10	1494	17	592	1323	
L116370:3612	66	10283	38	72	40	321	137	29	
L9839:442	2491	60	1851	9	774	11	321	660	
L168871:1819	7	263	6	361	8	2730	24	4	
L23642:141	1198	8	822	2	375	4	156	332	
L25354:253	1181	12	932		375	7	100	323	
L135664:4438	20	154	5	22	19	44	2212	3	
L87032:1210	11	315	13	51	14	141	37	9	
L90333:3644	3	204	2	31		190	5		
L184096:2297		13		1735	3	84	11	1	
L87031:697	4	136	4	12	3	137	4	5	
L14495:131	614	10	432		167	4	68	163	
L63724:1422	17	83	10	34	34	28	96	6	
L98994:891	31	161	46	21	19	59	21	5	

Source: Voters List; Data Mining: NamSor

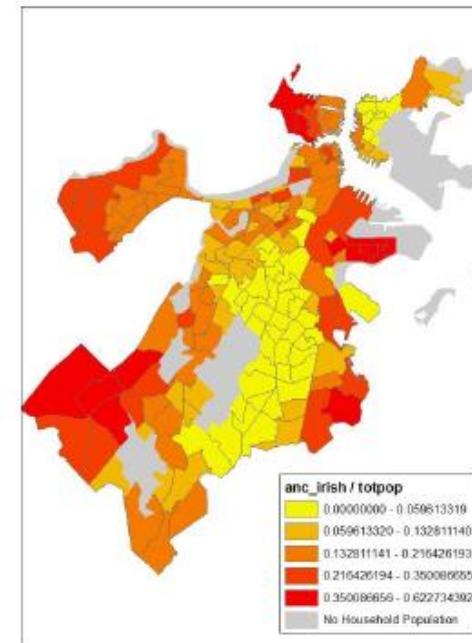
# Boston geo-demographics 1/2

16

## Irish Share, namsor



## Irish Share, 2010-2014 ACS



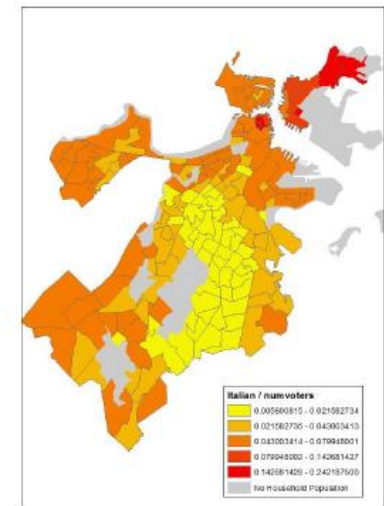
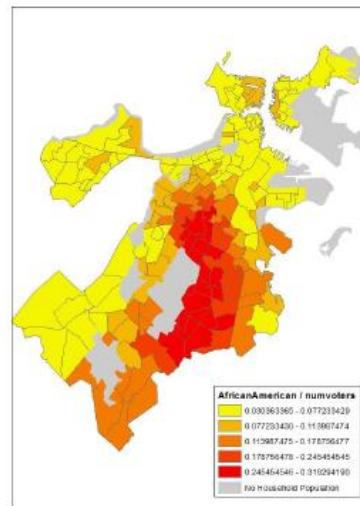
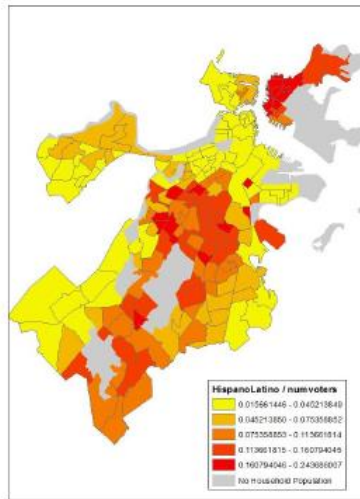
# Boston geo-demographics 2/2



Hispanic/Latino Share, namsor

Black/African-American Share, nams

Italian Share, namsor



March 7, 2016  
Presentation Title

March 7, 2016  
Presentation Title

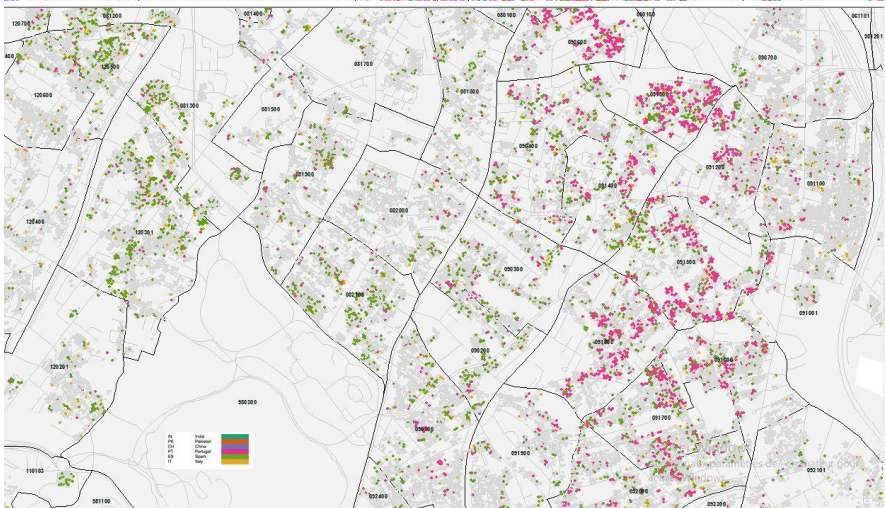
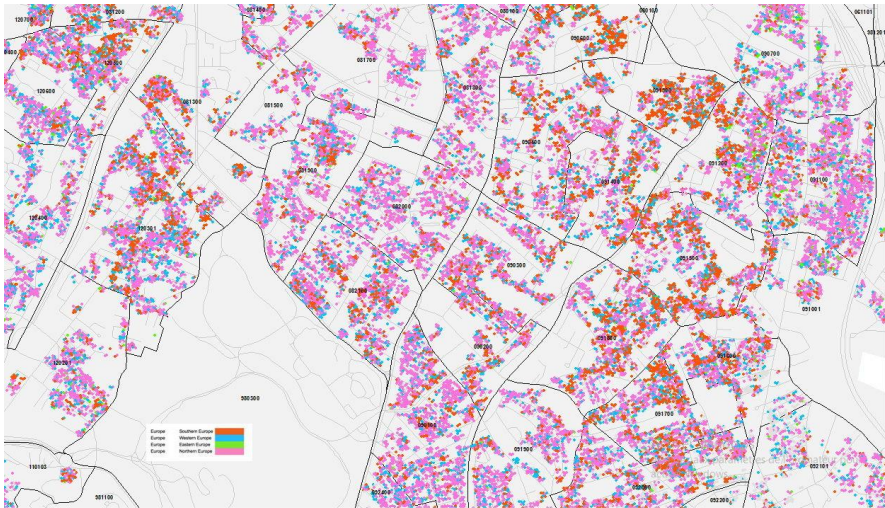
March 7, 2016  
Presentation Title

Source: Boston Voters List  
Visualization : ESRI  
Data Mining: NamSor



# Breaking down 'White' and 'Asian' into Portuguese, Spanish, Italian, India, Pakistan, China, ...

18



Source: Boston Voters List  
Visualization : ESRI  
Data Mining: NamSor



# PATENT DATABASES

19



Get our newsletter

Search

ISSUES

PUBLICATIONS

EVENTS

NEWS ROOM

MULTIMEDIA

ABOUT

## The Demographics of Innovation in the United States

[Adams Nager](#), [David M. Hart](#), [Stephen Ezell](#), and [Robert D. Atkinson](#) February 24, 2016

A groundbreaking ITIF survey shows why the country needs to broaden and deepen its pool of potential innovators with better STEM immigration and education policies.

[View Report](#)

[View Executive Summary](#)

[Event](#)



Groundbreaking @ITIFdc survey shows why US needs to broaden and deepen pool of potential innovators



.@ITIFdc releases groundbreaking survey on who innovates in the United States and where and how it occurs



# Indian diaspora names - a global airline use case

21

*'For 93% of our customers, when NamSor recognizes an Indian name, the client has travelled to India in the past.'*

At state level : ~50%

Analysis of NamSor's First Choice Country Compared to Historic Travel on [REDACTED]

NamSor's First Choice Country	Count of Unique Individuals who Have Travelled to NamSor's First Choice Country			% of Unique Individuals who Have Travelled to NamSor's First Choice Country		
	No	Yes	Grand Total	No	Yes	Grand Total
India	1,633	20,315	21,948	7%	93%	100%
Italy	281	869	1,150	24%	76%	100%
Bangladesh	524	1,456	1,980	26%	74%	100%
Ethiopia	3	8	11	27%	73%	100%
Iran	701	1,657	2,358	30%	70%	100%
Saudi Arabia	679	771	1,450	47%	53%	100%
Afghanistan	21	23	44	48%	52%	100%
Pakistan	2,171	2,309	4,480	48%	52%	100%
Jordan	148	124	272	54%	46%	100%
Kuwait	51	37	88	58%	42%	100%
Qatar	3	2	5	60%	40%	100%

Analysis of NamSor's Region Rounded Score Compared to Historic Travel on [REDACTED] for India

Customer has Flown to NamSor's Region	Count of Unique Individuals who Have Travelled to NamSor's Region			% of Unique Individuals who Have Travelled to NamSor's Region		
	No	Yes	Grand Total	No	Yes	Grand Total
5		3	3	0%	100%	100%
4	569	696	1,265	45%	55%	100%
3	2,202	2,774	4,976	44%	56%	100%
2	2,861	3,226	6,087	47%	53%	100%
1	2,442	2,523	4,965	49%	51%	100%
0	1,686	1,423	3,109	54%	46%	100%
-1	702	519	1,221	57%	43%	100%
-2	180	109	289	62%	38%	100%
-3	23	9	32	72%	28%	100%
-4	1		1	100%	0%	100%
Grand Total	10,666	11,282	21,948	49%	51%	100%

# Indian diaspora in Mauritius

22



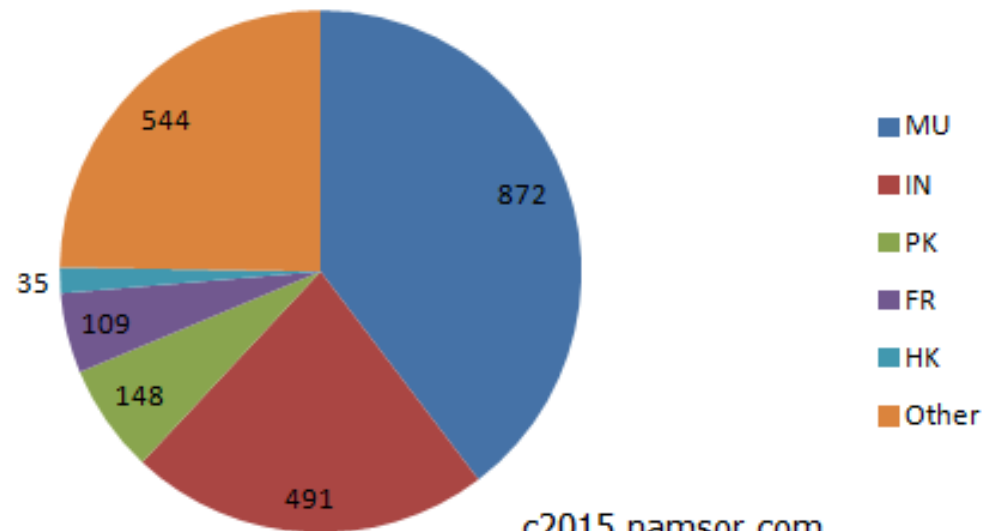
## Registered Medical Practitioners

The annual list for the reg  
category;

[Register of General Practitio](#)

### Onomastics of ~2200 GPs in Mauritius

Using NamSor v0.0.27

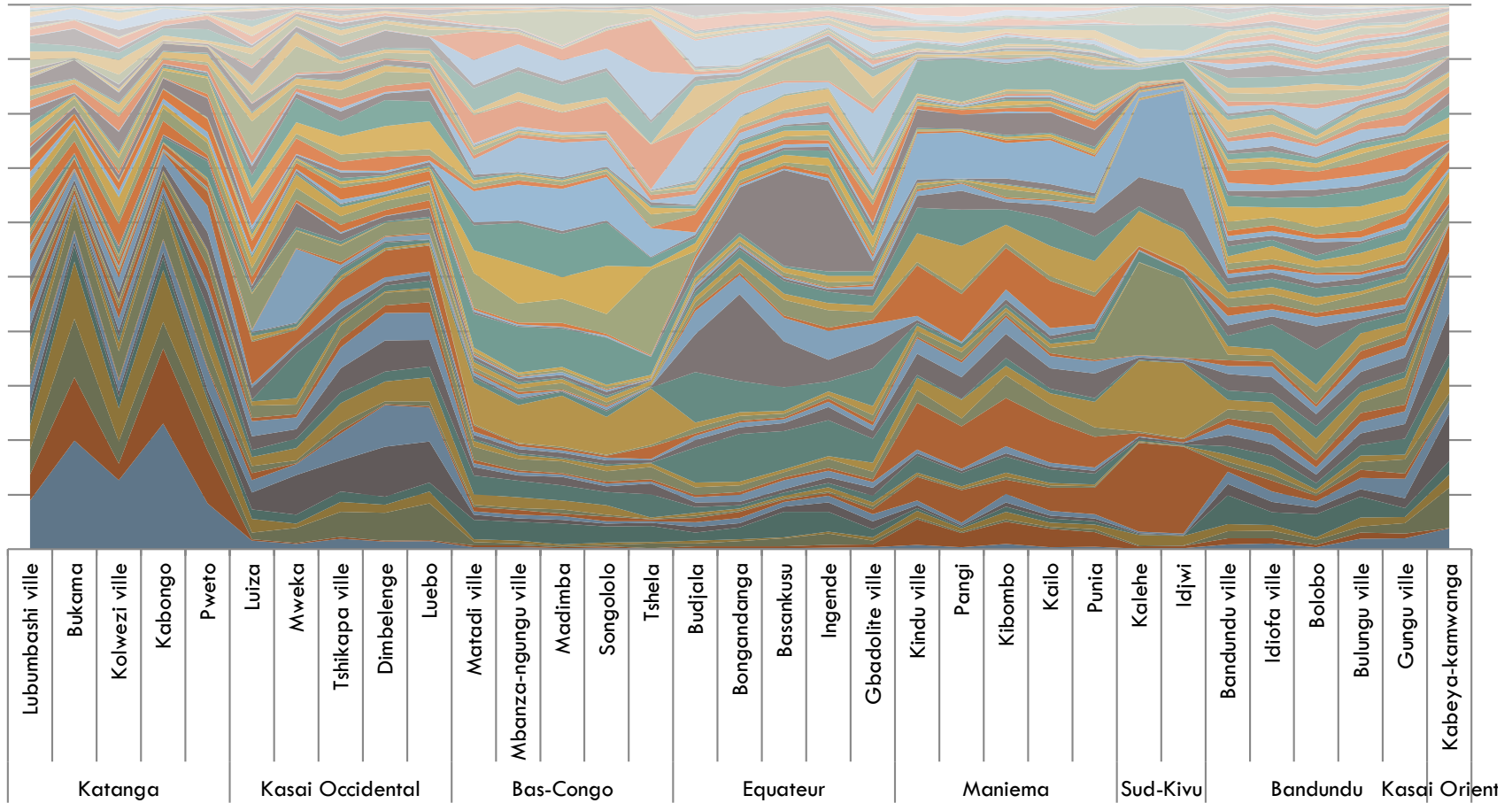


c2015 namsor.com

Source: medicalcouncilmu.org

# Africa: complex identities (Congo RDC)

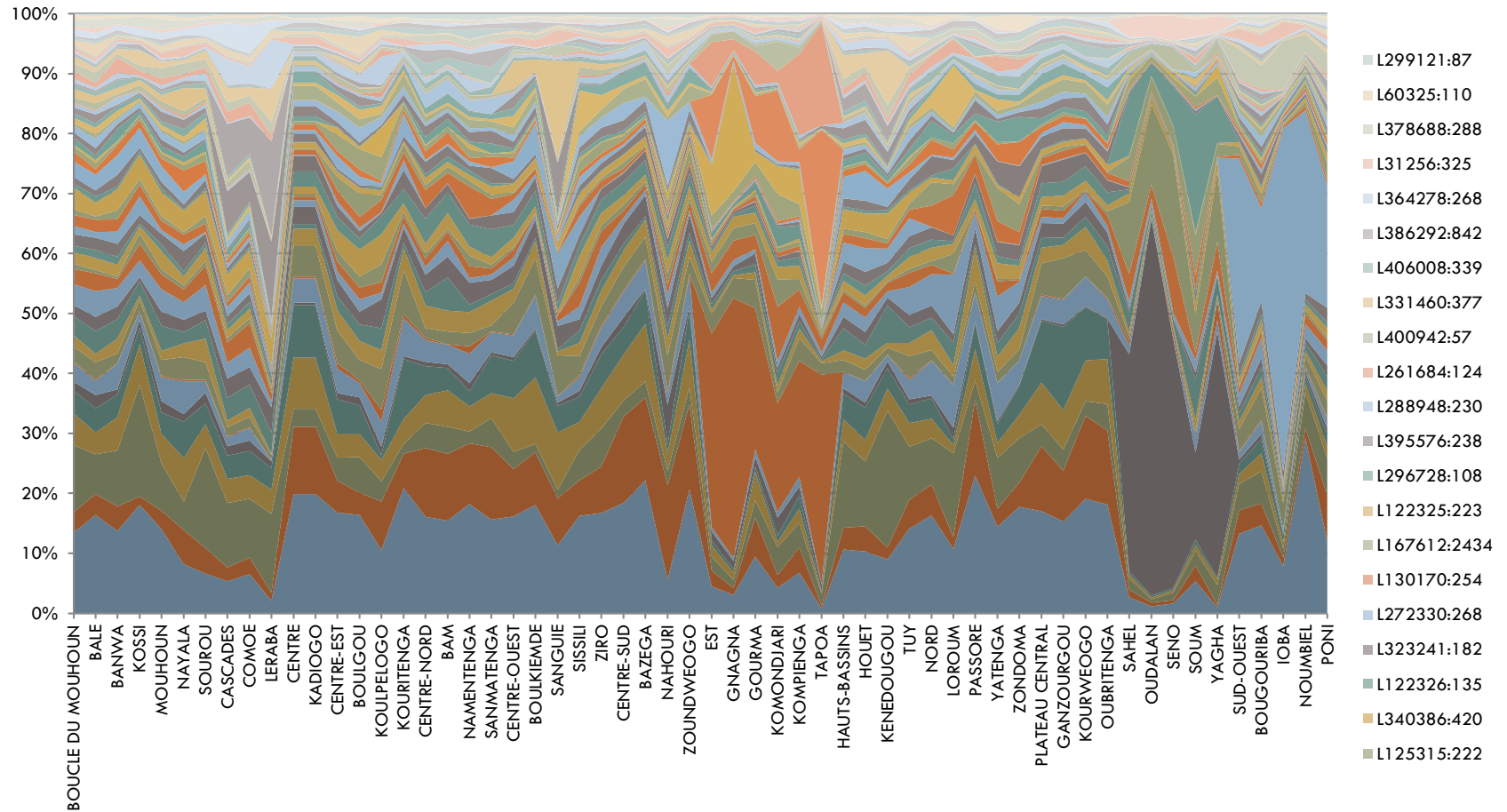
23





# Africa: complex identities (Burkina Faso)

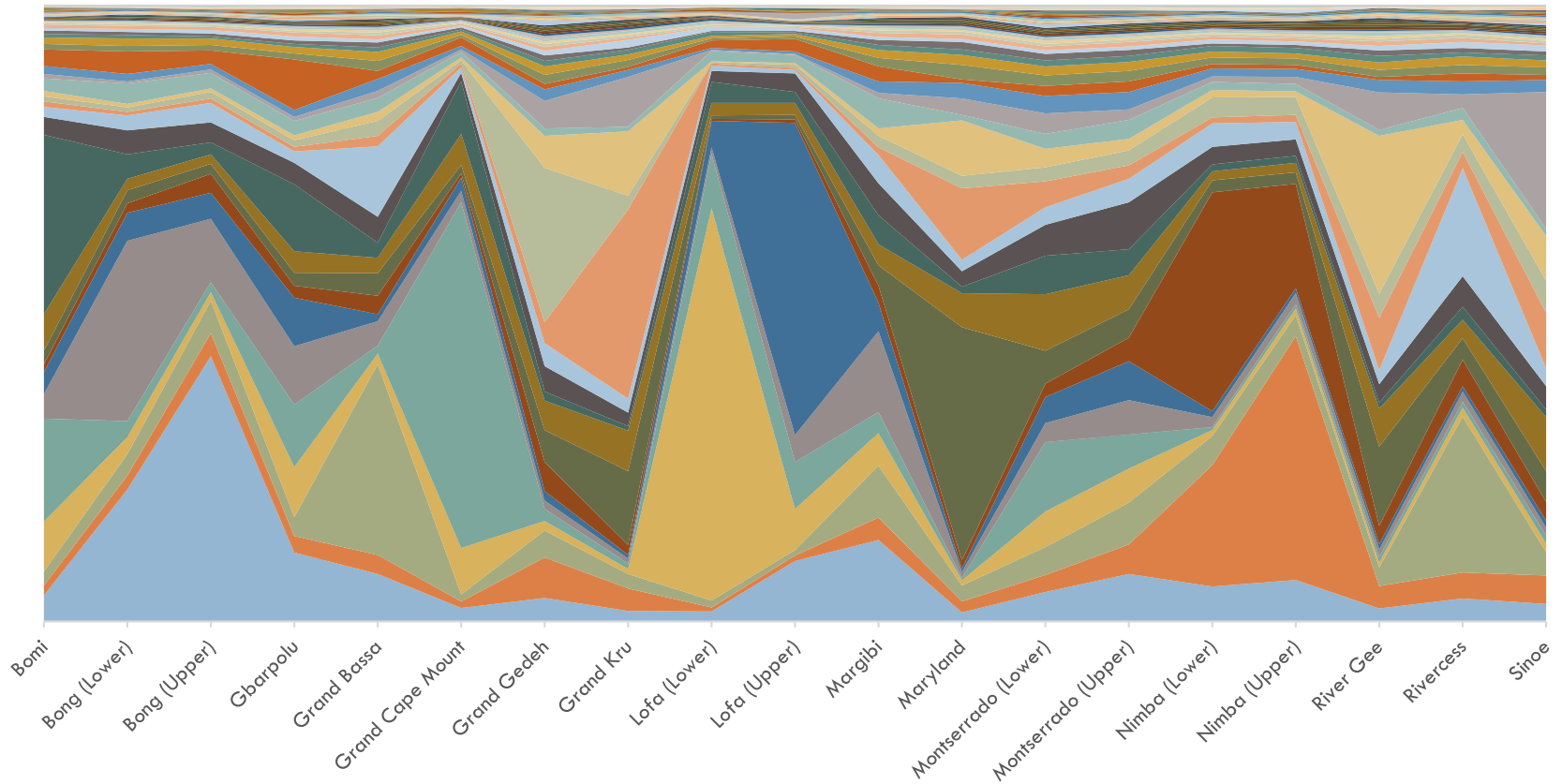
24



# Africa: complex identities (Liberia)

25

Liberia - a regional onomastics 'mille-feuille'



# Thank you !

26



Eliau CARSENAT,

[elian.carsenat@namsor.com](mailto:elian.carsenat@namsor.com)

Phone : +33 6 52 77 99 07

