

REUTERS/Yuya Shino

## People Disambiguation Engine and Analysis of Russian Science in Web of Science (2007-2011)

---

**Ilya Ponomarev**, Victor Shyu, Pawel Sulima, Eva Darian, Ronan Sorensen, Brian Lawton, Ciaran Bolger, Etienne Godard, and Joshua Schnell

[ilya.ponomarev@thomsonreuters.com](mailto:ilya.ponomarev@thomsonreuters.com)

Custom Analytics, Rockville, MD



THOMSON REUTERS

# Outline

---

- 1. People disambiguation for research management: why it is important?**
- 2. Custom Analytics at Thomson Reuters: PDE solution**
- 3. Example: Russian Map of Science Project (Challenges, Solutions, Validation)**
- 4. Enabling new analytical capabilities with new bibliometric indicators**

# Name Disambiguation

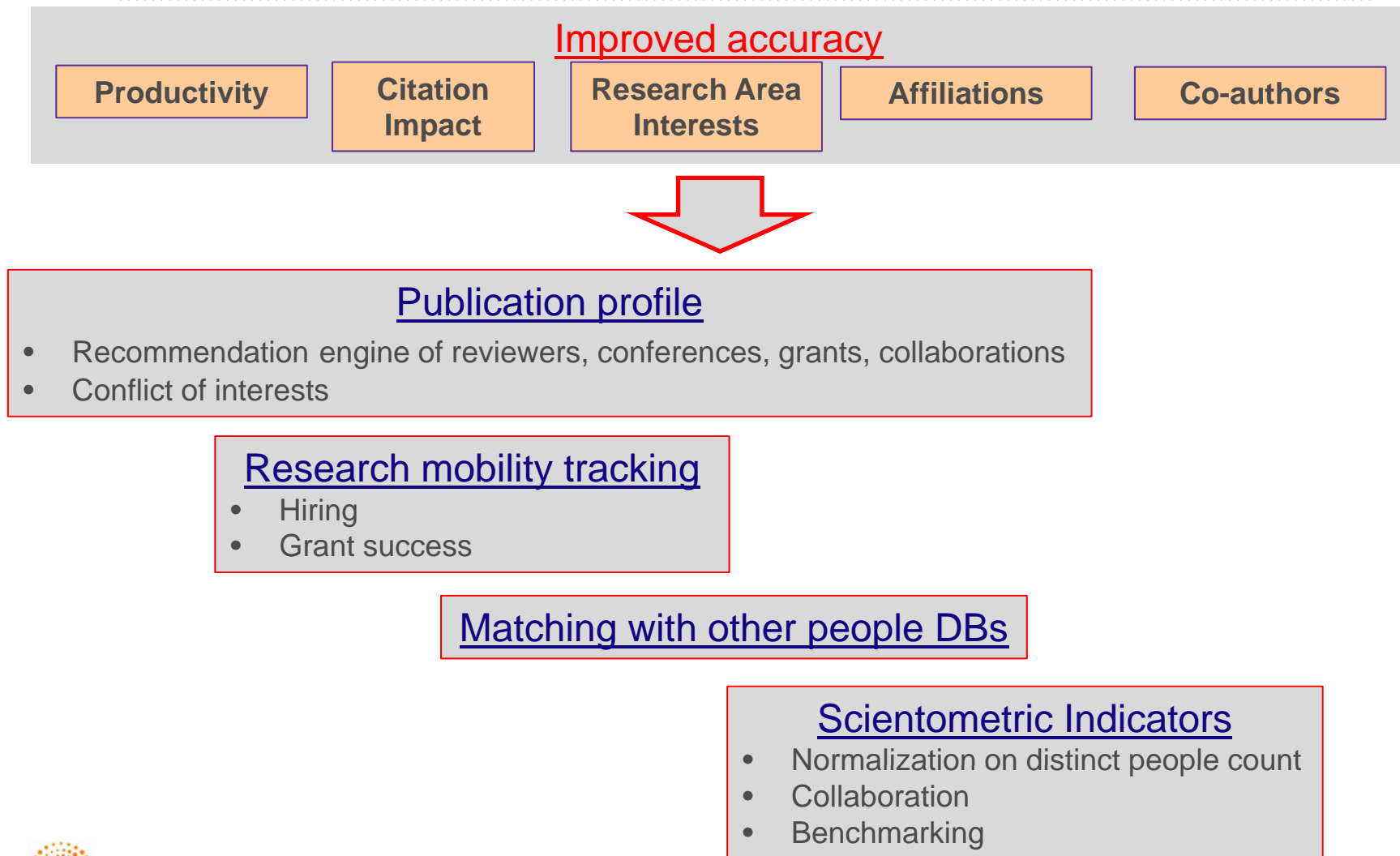
---

**Research Management DBs about events: Publications, Patents, Grants, Documents**

**Name Disambiguation:** the process of detecting and correcting ambiguous *named entities* that represent the same real-world *object*  
**Synonyms:** name unifications, matching, entity resolution, authorship

Approaches: supervised, unsupervised, semi-supervised, ORCHID  
How to build scalable approach?

# People Disambiguation for Research Management – Why It Is Important?



# People Disambiguation Engine

**Custom Analytics**  
**Rockville, MD**

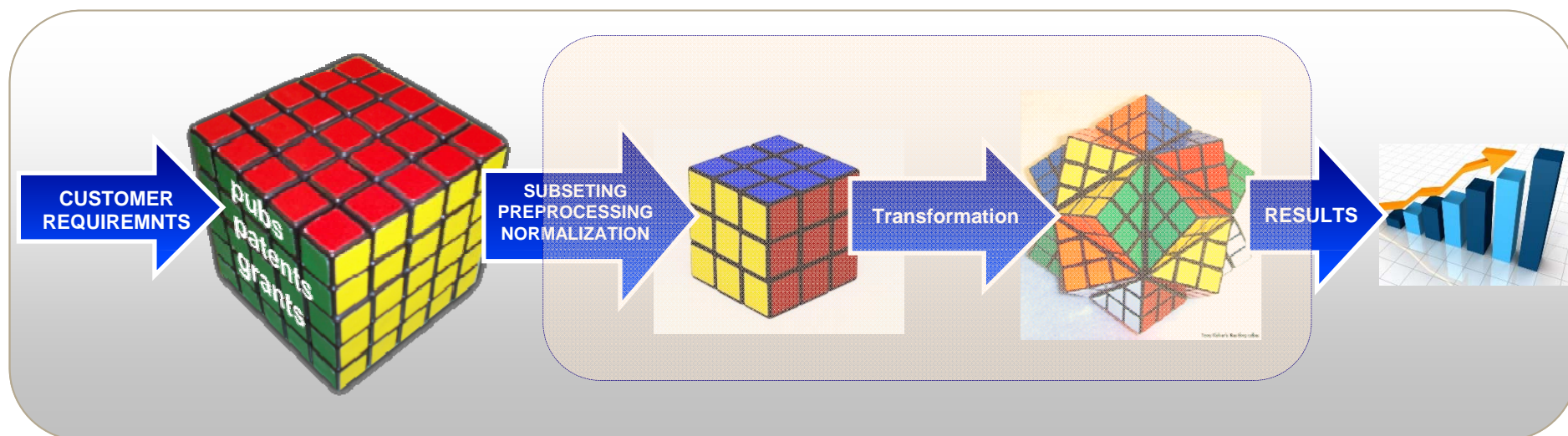
# Projects Workflow: Custom Analytics Group

1. Project oriented
2. Scalability
3. Semantic and context-aware
4. DB-centric system

ScienceWire®

PDE Engine

New Knowledge



# PDE: Vocabulary

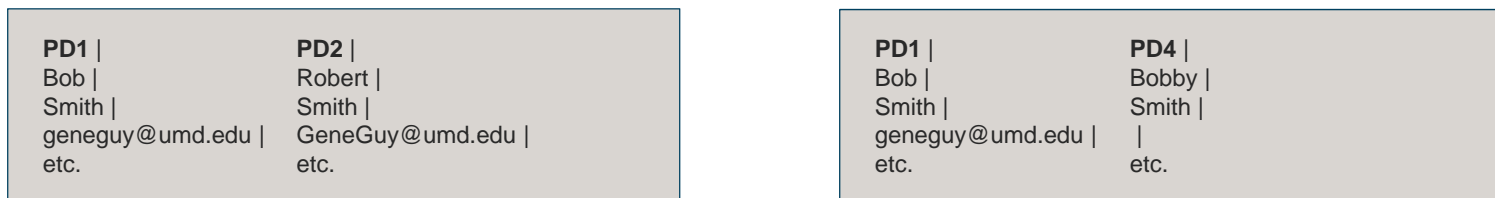
**Descriptor** - a group of attributes describing a person

- Can be complicated (e.g. – a list of co-authors)
- Includes a link to source



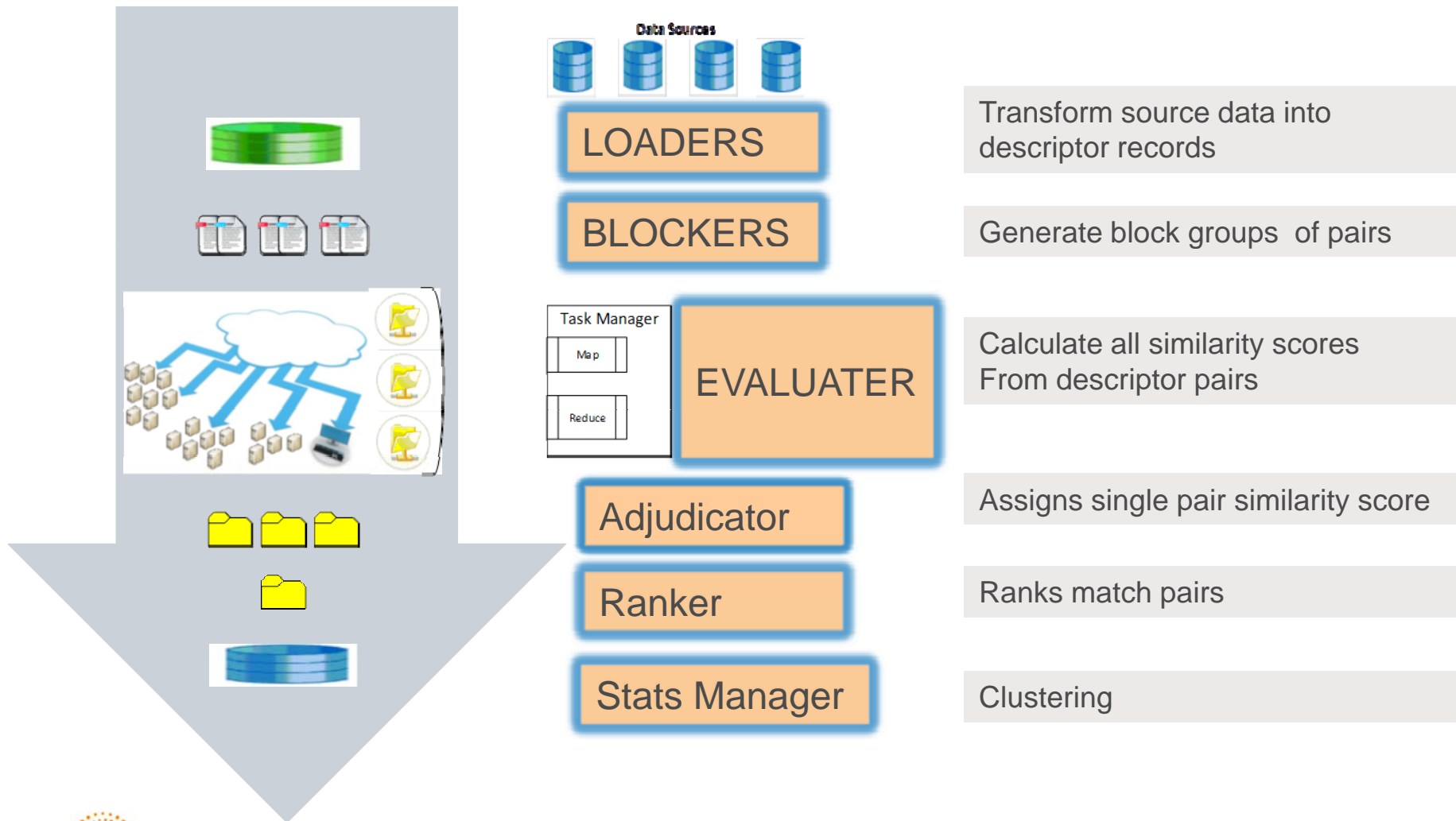
Pub1

**Descriptor pair** - two descriptors linked together





# OVERALL STRUCTURE





# EVALUATORS: List of Comparisons Performed

ALL DATABASES

<< Back to results list Record 10 of 2,680 Record from Web of Science®

## Stokes' dream: Measurement of fluid viscosity from the attenuation of capillary waves

Full Text Print E-mail Add to Marked List Save to EndNote Web Save to EndNote, RefMan, ProCite Save to RefWorks more options

Author(s): Behroozi F (Behroozi, F.)<sup>1</sup>, Smith J (Smith, J.)<sup>1</sup>, Even W (Even, W.)<sup>1</sup>

Source: AMERICAN JOURNAL OF PHYSICS Volume: 78 Issue: 11 Pages: 1165-1169 Published: NOV 2010

Times Cited: 0 References: 33 Citation Map

**Abstract:** The determination of viscosity from the attenuation of capillary waves was first suggested by Stokes more than a century ago. At the time, it was not practical to measure the attenuation of surface waves with the requisite precision to render the method useful. We describe a noncontact method for measuring the wavelength and amplitude of single-frequency capillary waves to obtain reliable values of the surface tension and viscosity. The attenuation data for several glycerin-water mixtures are used to obtain the viscosity as a function of glycerin concentration. For a wide range of viscosity, the method yields results that are in good agreement with the most reliable published data. (C) 2010 American Association of Physics Teachers. [DOI: 10.1119/1.3467887]

**Document Type:** Article

**Language:** English

**Keywords Plus:** SURFACE LIGHT-SCATTERING; DISPERSION-RELATION; LASER INTERFEROMETRY; LIQUID INTERFACES; GRAVITY WAVES; WATER; ENERGY; CONSERVATION; MIXTURES

**Reprint Address:** Behroozi, F (reprint author), Univ No Iowa, Dept Phys, Cedar Falls, IA 50614 USA

**Addresses:** 1: Univ No Iowa, Dept Phys, Cedar Falls, IA 50614 USA

**E-mail Addresses:** fbehroozi@uni.edu

**Funding Acknowledgement:**

Funding Agency	Grant Number
Carver Trust	
UNI Applied Technology Fund	
Iowa Space Grant Consortium	

[Show funding text]

**Publisher:** AMER ASSOC PHYSICS TEACHERS AMER INST PHYSICS, STE 1 NO 1, 2 HUNTINGTON QUADRANGLE, MELVILLE, NY 11747-4502 USA

**IDS Number:** 667LZ

**ISSN:** 0002-9505

**DOI:** 10.1119/1.3467887

**Cited by: 0**  
This article has been cited 0 times (from Web of Science).  
[Create Citation Alert]

**Related Records:**  
Find similar records based on shared references (from Web of Science).  
[view related records]

**References: 33**  
View the bibliography of this record (from Web of Science).

**Additional information**

- View the journal's impact factor (in Journal Citation Reports)
- View the journal's Table of Contents (in Current Contents Connect)

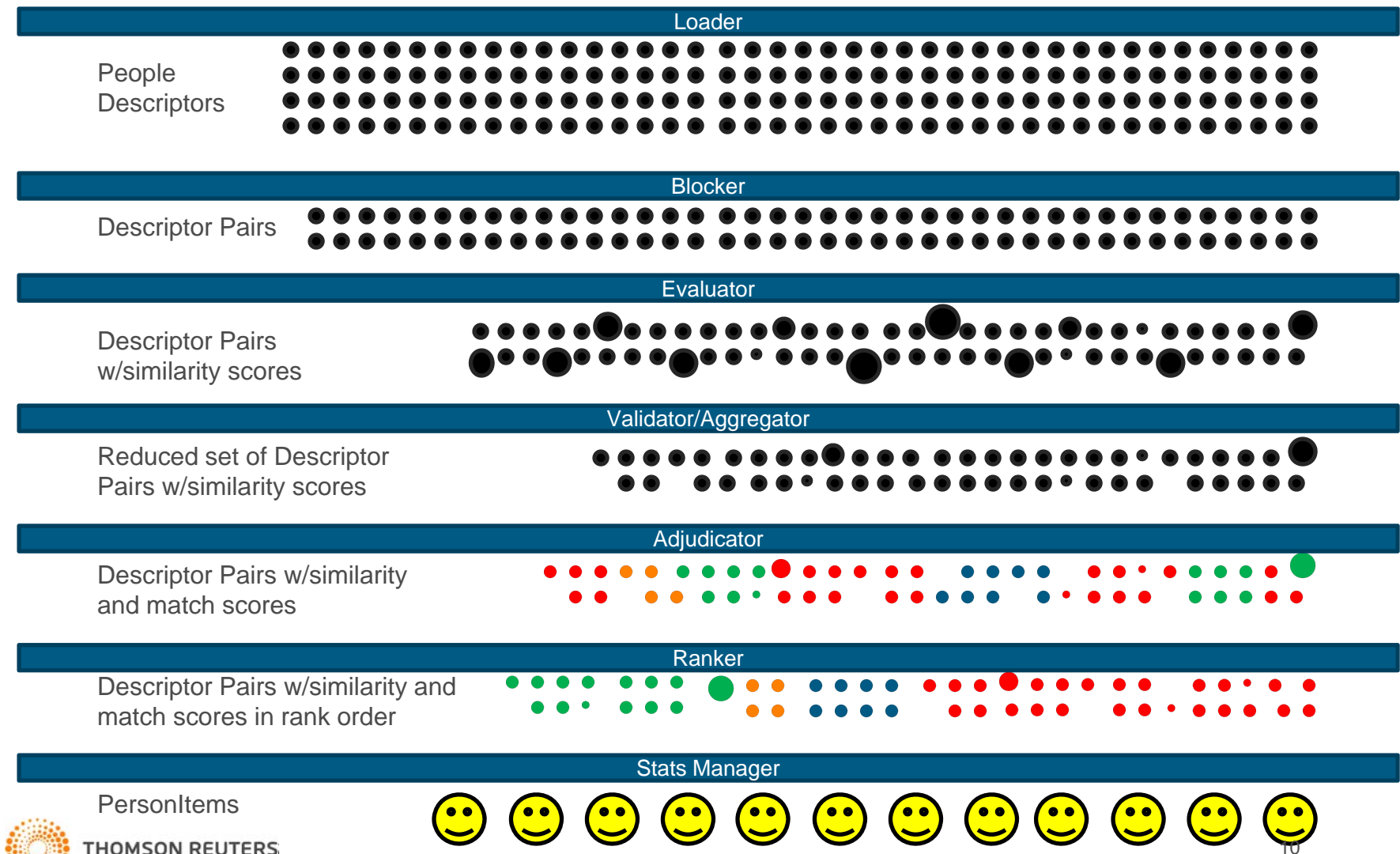
**View this record in other databases:**

- View citation data (in Web of Science)
- View the most recent data (in Current Contents Connect)
- View engineering data (in Inspec)

## Attributes:

1. Names
2. Emails
3. Co-authors
4. Self-citations
5. Citations
6. References
7. Common words in titles
8. Subject categories (similar journals)
9. Affiliations

# Summary of Process



# RMOS PROJECT

## RMOS Project Goals

---

- Create Map of the Most Productive Research Organizations in Russia  
(Russian Ministry of Edu and Sci, PWC, NICON)
  1. Collect all papers with at least 1 Russian affiliation in 2007-2011
  2. Disambiguate organizations in these papers
  3. Disambiguate people on these pubs

Sounds simple

# Challenges

## **Author disambiguation in Publication DBs**

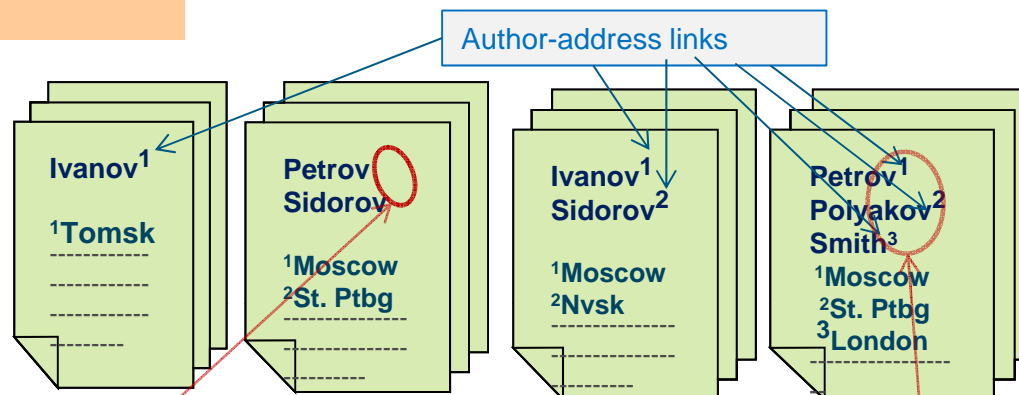
- 1. Complexity: there are 180,000,000 descriptions of people in WoS (51000 years of calculations by brut force algorithm)**
- 2. Orthographic and spelling name variants (author, publisher, OCR-DB soft)**
- 3. Real name changes over time (marriage)**
- 4. Incomplete data (First Name vs First Initial vs NULL)**
- 5. Parsing errors (First name in MiddleName field)**
- 6. Incorrect or ambiguous links to other attributes, incorrect attributes**
- 7. Name commonality (Li, Wang, Ivanov)**

# Challenges

## Real people, real full source documents



1. Emails are not always good
2. Common name
3. Poison descriptors
4. Incorrect affiliations

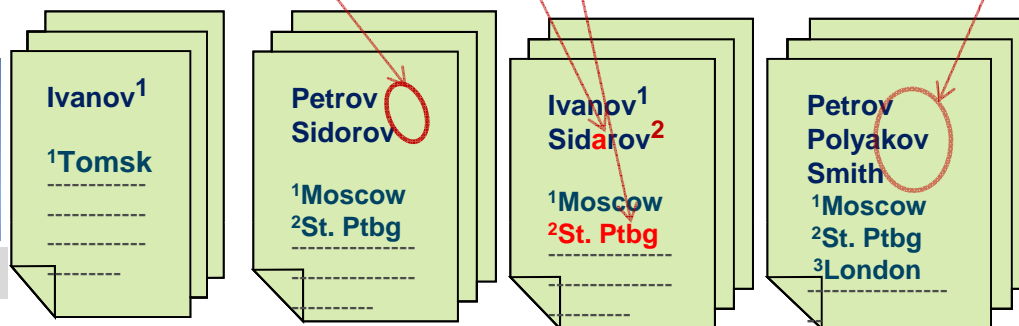


No links

Typos, wrong links

Missing links (90% before 2007)

## Web of Science



# Project Specific Challenges: Preprocessing Organizations Unification and Ambiguous links



Dusamb Org	# addresses	# WoS orgs
Moscow MV Lomonosov State Univ	12352	113
St Petersburg State Univ	4072	53

## Ambiguous links recovery

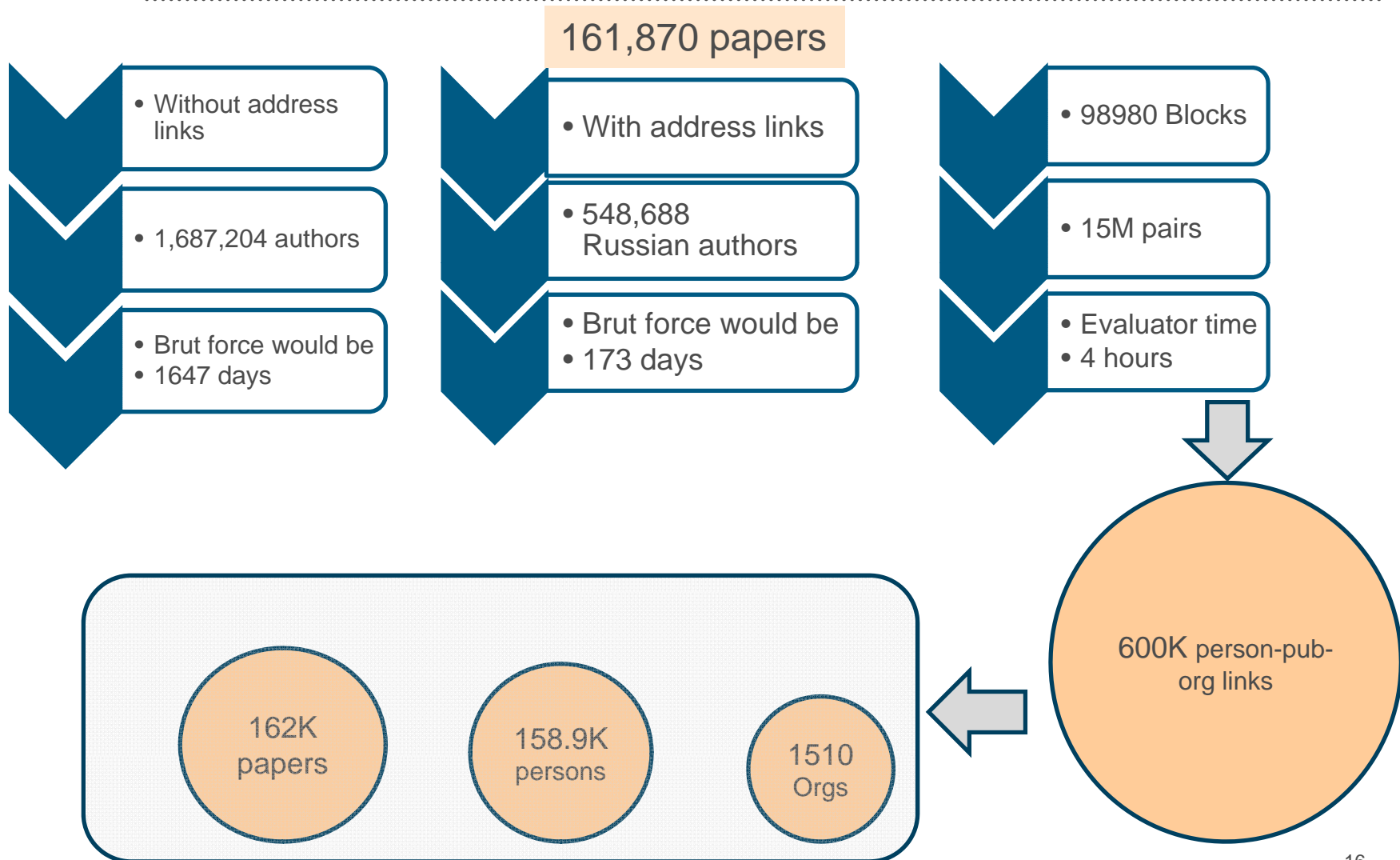
Year	Total Number of publications	Number pubs without Author-address link	% without Author-Address link (before)	Number pubs with recovered links	Recovery rate, %
2007	31.0K	24.6K	79%	12.6K	52%
2008	33.3K	5.3K	16%	5.0K	96%
2009	33.3K	5.4K	16%	5.2K	98.8%
2010	31.9K	4.8K	15%	4.7K	98.7%
2011	32.5K	4.6K	14%	4.5K	98.7%
Total	161.9K	44.6K	27%	32.0K	72%



Strong recovery rate in 2008-2011



# RMOS Project Results

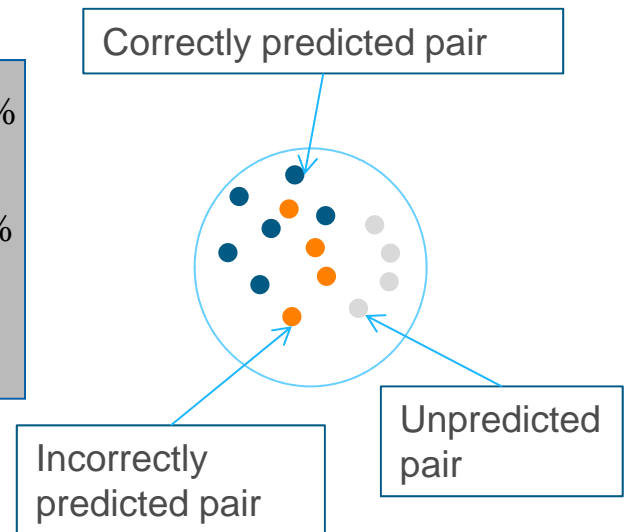


# Validation

Choosing threshold: precision/recall/accuracy tradeoff

Evaluation measures:

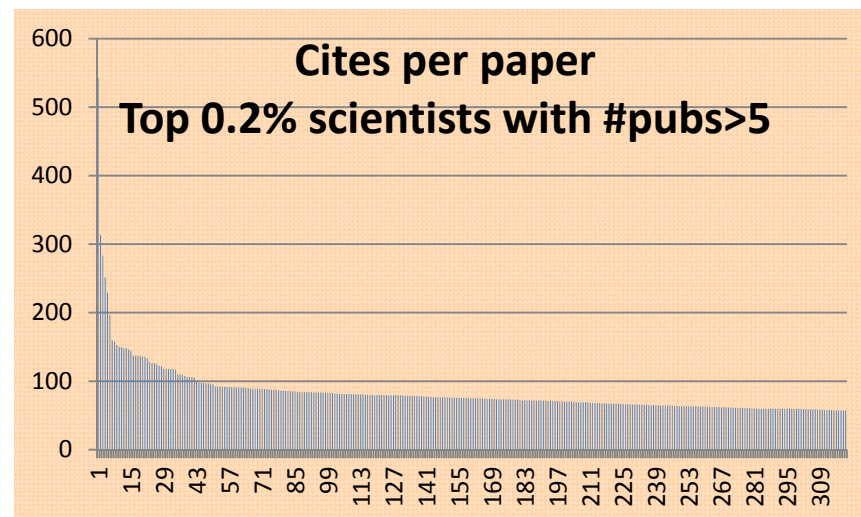
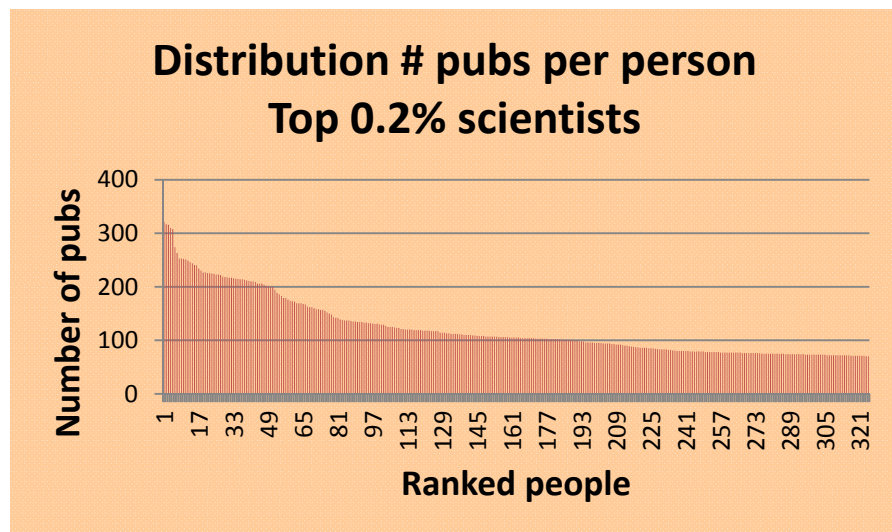
$$\begin{aligned}\text{Pairwise Precision} &= \frac{\# \text{ pairs correctly predicted to the author}}{\# \text{ total pairs predicted to the author}} = 96\% \\ \text{Pairwise recall} &= \frac{\# \text{ pairs correctly predicted to the author}}{\# \text{ total pairs to the author}} = 90\% \\ \text{Pairwise } F1 &= \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}\end{aligned}$$



Caveat: Error is larger in common name clusters

- Lumping – assigning articles by different people to the same person
- Splitting – assigning articles by the same person to different people

# New Knowledge – Benchmarking the Whole Country Distributions 2007-2011

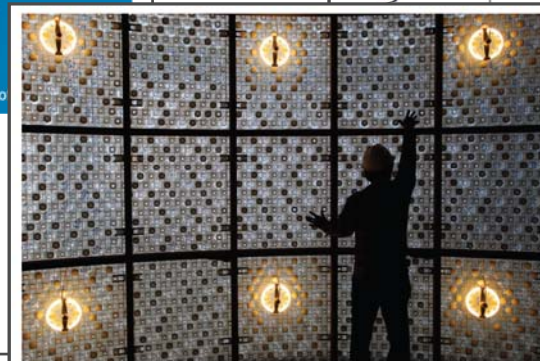


Having 1 pub – 86.5K  
 Having 2 pubs – 24.7K  
 Having 3 pubs – 12.6K  
 Having >3 pubs – 35K

Max papers – 321  
 Top 1% – 29  
 Average – 3  
 Median – 1

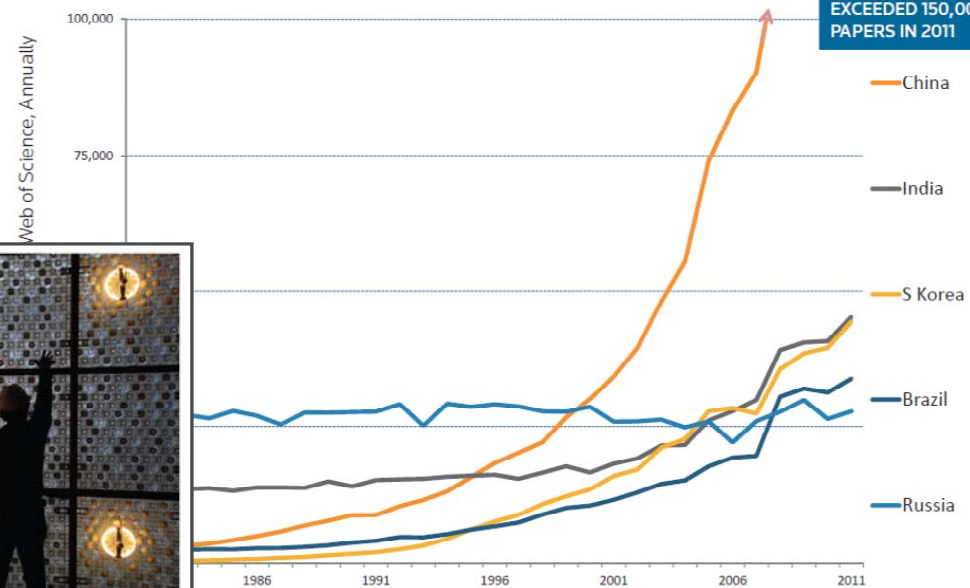
Max cites\_pub – 543  
 Top 1% – 62.  
 Average – 4.5  
 Median – 1.

# Thank You!



ANNUAL RESEARCH PUBLICATION OUTPUT OF THE FIVE BRICK COUNTRIES

FIGURE 5



Reuters Web of Knowledge. (See also Figure 7 on trajectories of patent output.)



## Contacts

---

- Joshua Schnell [joshua.schnell@thomsonreuters.com](mailto:joshua.schnell@thomsonreuters.com)
- Victor Shyu [victor.shyu@thomsonreuters.com](mailto:victor.shyu@thomsonreuters.com)
- Ilya Ponomarev [ilya.ponomarev@thomsonreuters.com](mailto:ilya.ponomarev@thomsonreuters.com)