

# Research Intelligence

# Новый подход к полномасштабному картированию науки. Выдающиеся научные направления России

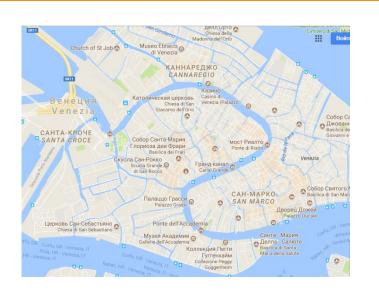
По материалам презентации Кэвина Бояка (Kevin W. Boyack,) SciTech Strategies, Inc. (август 2017)

Галина П. Якшонок, консультант по аналитическим решениям Elsevier

**5-я международная конференция НЭИКОН «Электронные научные и образовательные ресурсы: создание, продвижение и использование»** 27 сентября 2017



**ELSEVIER** Research Intelligence



#### ●Impact OF Air Travel ON Global Spread OF Infectious Diseases ●



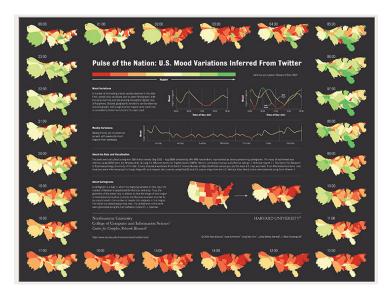
Epidemic spreading pattern changed dramatically after the development of modern transportation systems.



#### Forecasts OF THE Next Pandemic Influenza

April 5 \$ April 5 \$ April 5 \$ April 5 \$ \$





#### Happiness Depends on Various **Factors** Social scientists are starting to

include relative happiness with hard data on economic status, health, and other factors as they assess quality of life. They rely on surveys of "subjective well-being"-how good people feel about their lives. A world map of one "happiness index" shows many, but not all, wealthy northern countries faring well. Residents of sub-Saharan Africa and the former Soviet Union, meanwhile, report particularly low levels of

Any attempt to measure happiness will fall short—each life is a series of joys, struggles, and sorrows, and satisfaction can depend as much on outlook as on circumstances. Averages obscure the happy moments in struggling nations, as well as people who suffer from poor health, poverty, or discrimination in countries that rank high. Still, happiness indices can help researchers move beyond simple economics as they track progress—or backeliding—over time.

contentment.

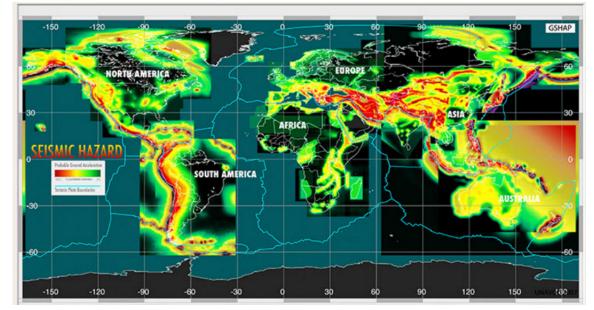
MEASURING THE INTANGIBLE The map is derived from the New Economics Foundation's 2006 "Happy Planet Index," which drew on over 100 surveys of subpethic well-being, its "satisfac-tion with life scale" — a happiness index—narks the relative happi-ness of nations, from a high of 273 (Denmark and Switzerland) to a low of 100 (Burund).

■ Unimpov

#### "It's time we admitted there's more to life than money."

RANKING THE WORLD'S HAPPIEST PLACES and several wealthy cox make the lat, but so do prosperous island natio 1 DENMARK SWITZERLAND

- 2 AUSTRIA ICELAND
- 3 BAHAMAS FINLAND
- 4 BHUTAN
- BRUNES CANADA PELAND LUXEMBOURG
- 5 COSTARICA MALTA NETHERLANDS 6 ANTIGUA AND BAR
- MALAYSIA NEW ZEALAND NORWAY SEYCHELLES ST. KITTS AND NEVI UNITED ARAB EMA UNITED STATES VENEZUELA



#### Анализ научного ландшафта

Если исследование это продукт, то анализ элементов множества продуктов поможет оптимальному распределению ресурсов:

- для грантодающих/финансирующих организации какие научные темы спонсировать
- для администраторов какие научные темы поддерживать и кого нанимать
- для исследователей над какими темами работать и, соответственно, подавать заявку на грант

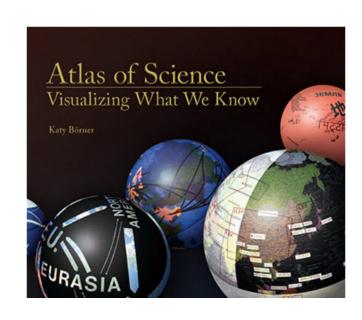
В терминах «спрос - предложение»:

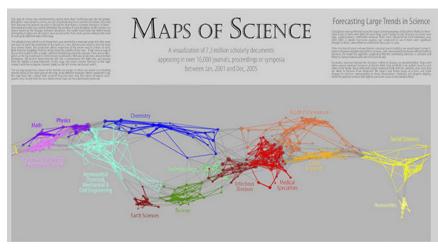
- Темы (группа публикаций) представляют предложение
- Грант (покупка исследования) представляет спрос

Выиграет тот, у кого есть информация и о спросе и о предложениях

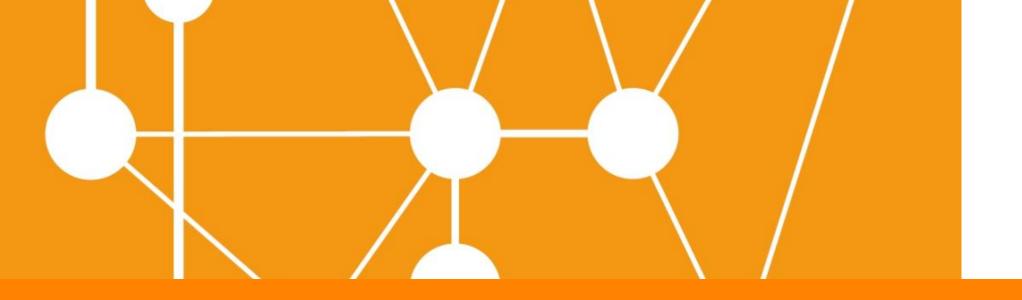
#### Пробелы и подходы

- ПРОБЕЛ: в настоящий момент нет единого решения по комплексной модели науки или списку научных тем (и их относительной ценности)
- Подход в заполнении этого пробела путем:
  - о Создания детальной модели научных тем
  - Создания индикатора спроса (ценности) для тем, который коррелирует с финансированием



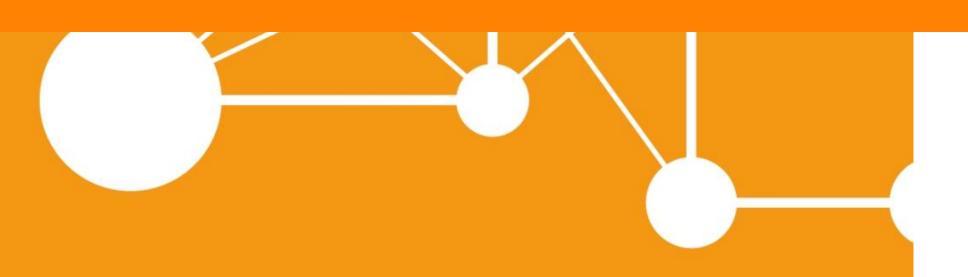


http://scimaps.org/





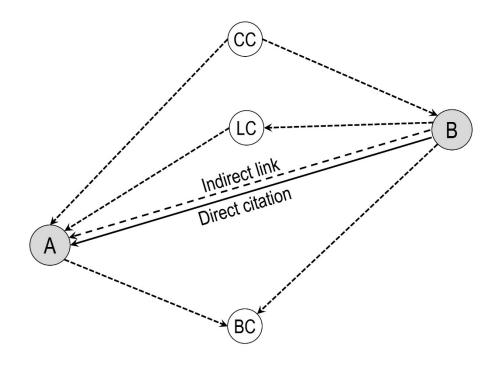
# Решение





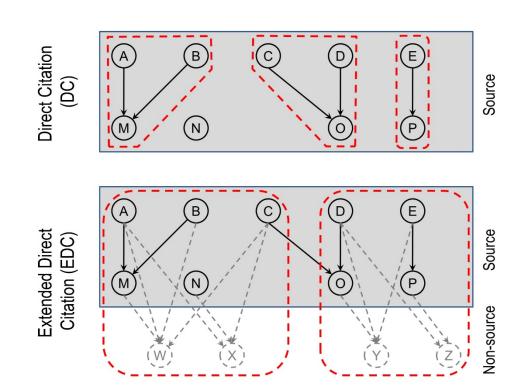
#### Технические аспекты

- Выделение направлений/тем требует определения взаимосвязей между парами статей
- Существует множество подходов
  - Схожесть ключевых слов (keyword similarity)
    - о Не очень точный
  - Текстовая схожесть (textual similarity)
    - Чрезвычайно трудоёмкий
  - Связи по ссылкам
    - о Прямое цитирование (Direct citation (DC))
      - Вычислительно просто
    - Библиографическое сочетание (Bibliographic coupling (BC))
    - Колитирование (Со citation (СС))
    - Опосредованное цитирование (Longitudinal coupling (LC))
      - Все − вычислительно трудоёмки



#### Технические аспекты

- Прямое цитирование имеет варианты
- Только для индексируемых документов
  - 32 млн статей
  - 439 млн связей
  - Но ... только 29 млн связано
- Включая цитируемые, но не индексируемые документы
  - 32 млн статей
  - 31 млн не индексируемых документов
  - 830 млн связей
  - Намного больший сигнал для кластеризации 32 млн статей



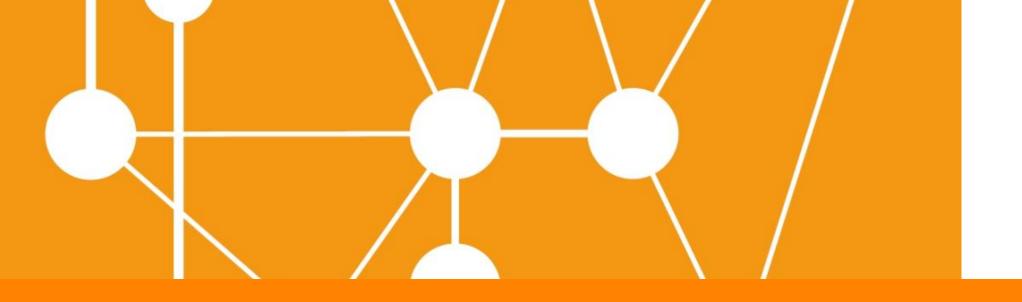
# Требования и решения

#### ТРЕБОВАНИЯ

- Охват: Полный список тем в науке
- Гранулярность: Темы соответствующего размера и их число
- о Точность: Точные темы, которые содержат относящиеся к теме статьи
- Стабильность: Темы с реалистичной динамикой

#### РЕШЕНИЕ

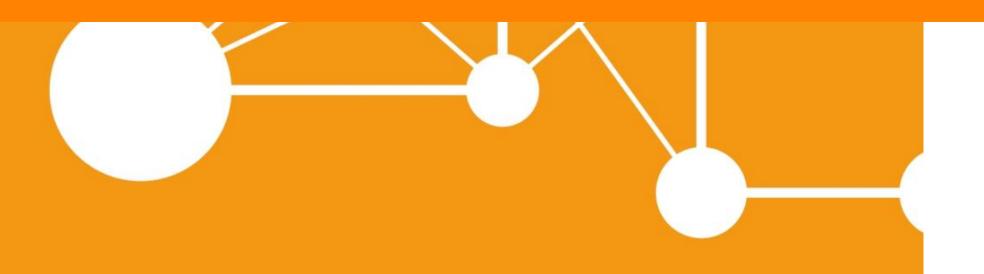
Определение около 100 000 тем в науке, используя прямое цитирование по связям цитирования (в том числе цитируемых неиндексированных документов) по полной базе данных Scopus\*





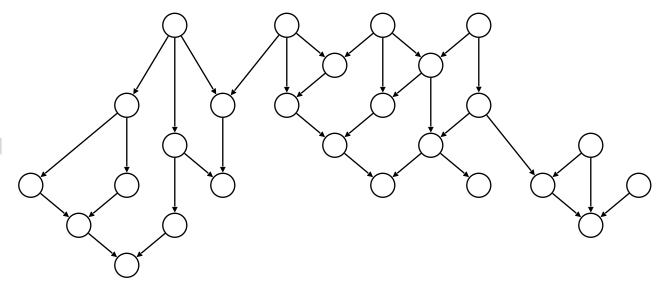
# Создание модели





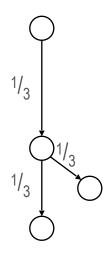
#### Моделирование тем – процесс

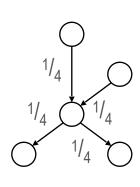
- Создание списка цитирующихцитируемых (статья-ссылка) пар (EIDs)
- Вычисление значения связей для каждой пары, на основе количества ссылок/связей
- Используя весь список ссылок и значений связей, проведение кластеризации документов
- В результате список EIDкластерID определяет набор тем



## Моделирование тем – процесс

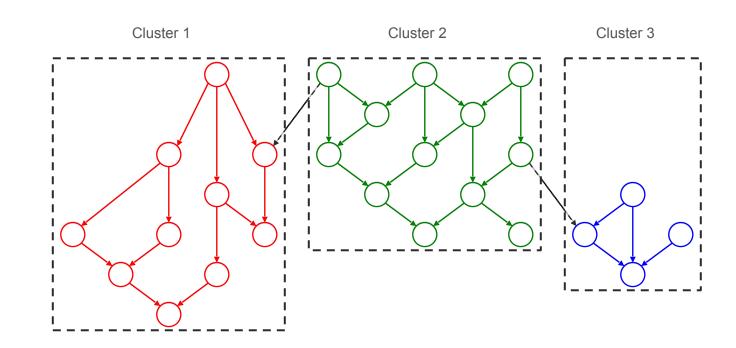
- Создание списка цитирующихцитируемых (статья-ссылка) пар (EIDs)
- Вычисление значения связей для каждой пары, на основе количества ссылок/связей
- Используя весь список ссылок и значений связей, проведение кластеризации документов
- В результате список EIDкластерID определяет набор тем





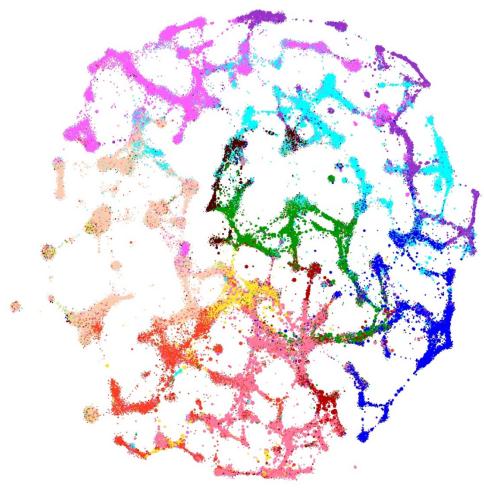
#### Моделирование тем – процесс

- Создание списка цитирующих-цитируемых (статья-ссылка) пар (EIDs)
- Вычисление значения связей для каждой пары, на основе количества ссылок/связей
- Используя весь список ссылок и значений связей, проведение кластеризации документов
- В результате список EIDкластерID определяет набор тем



#### Пример модели и карта

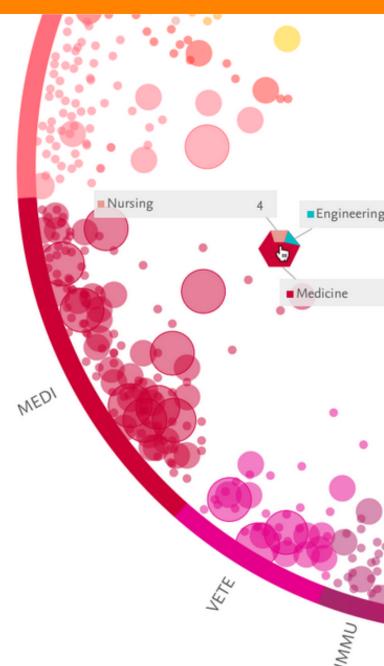
- Данные Scopus 1996-2012
- 582 млн цитирующих-цитируемых пар, 24.6 млн источников EID, 23.8 млн цитируемых не индексируемых EID
- Расчет значения связей для 582 млн пар
- Использование SLM (smart local moving algorithm)
- Несколько кластеров с <50 единицами влились в более крупные кластеры
- Результат 91,726 кластеров (научных тем)

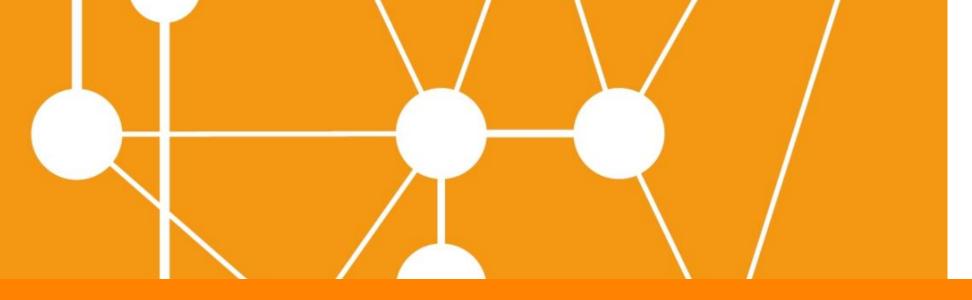


Klavans, R. and K.W. Boyack, Research portfolio analysis and topic prominence. Journal of Informetrics, 2017 (under review).

#### Дополнение модели

- Работы 2013-2015 были добавлены к существующей модели с их ссылками (90% точности)
- Нет необходимости в построении новой модели каждый год
- Такая стабильность приветствуется пользователями

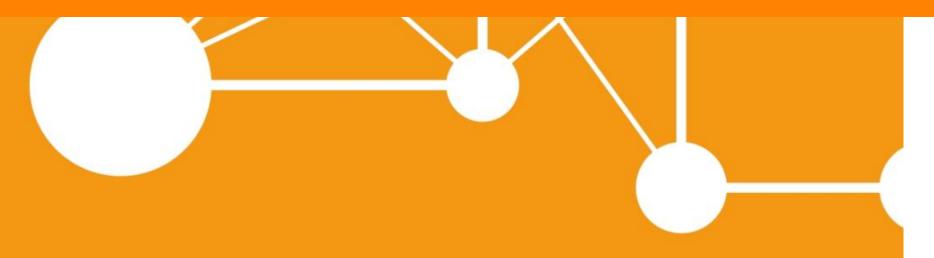






# Характеристики модели – prominence





## Выдающиеся направления (topic prominence)

- Составной показатель
- Рассматриваемые параметры
  - Количество ссылок в году n на статью опубликованную в году n и n-1
  - Scopus Просмотры (Views Count) в году n на статью опубликованную в году n и n-1
  - Средний CiteScore для года n
  - Среднее число авторов на статью для года n
  - Витальность противоположность возрасту ссылок, подобие "state-of-the-art" в Competencies (SciVal)

- Изначально также рассматривались ссылки в патентах и авторство в отрасли, но они были отброшены, чтобы избежать особенностей, связанных с экономическими мотивами

# Выдающиеся направления (topic prominence)

- Параметры Количество ссылок, Просмотры, CiteScore и Авторы были преобразованы для уменьшения искажений
- Высокая корреляция между Количеством ссылок, Просмотрами и Citescore
- Низкая для Авторов и Витальности

Table 1. Correlation matrix for candidate variables. "L:" denotes log transform.

	L:Citations	L:Views	L:CiteScore	L:Authors	Vitality
L:Citations	1.000				
L:Views	0.810	1.000			
L:CiteScore	0.533	0.483	1.000		
L:Authors	0.395	0.395	0.509	1.000	
Vitality	0.313	0.288	0.290	0.425	1.000

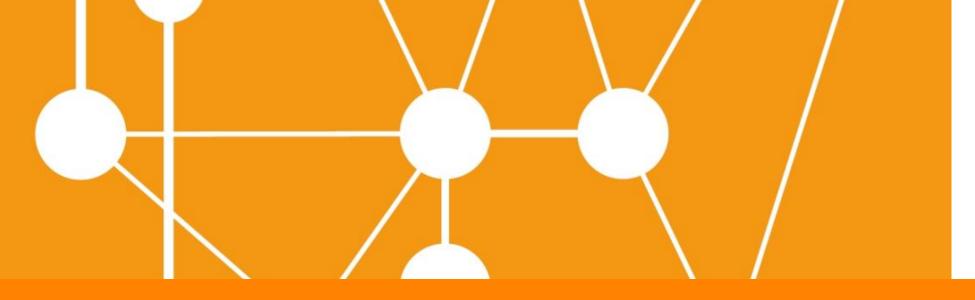
- Решение о включении в индикатор (prominence) Числа ссылок, Просмотров и CiteScore
- Формула Prominence:  $P_j = 0.495 (C_j mean(C_j))/stdev(C_j) + 0.391 (V_j mean(V_j))/stdev(V_j) + 0.114 (CS_j mean(CS_j))/stdev(CS_j),$

Table 2. Factor loadings and scoring coefficients used to calculate topic prominence.

1 metor roughly and seering eventerents used to emediate to pre-pro-									
	Factor 1	Factor 2	Normalized Score						
L:Citations	0.837	- 0.244	0.495						
L:Views	0.812	- 0.262	0.391						
L:CiteScore	0.653	0.154	0.114						
L:Authors	0.593	0.334	(not used)						
Vitality	0.441	0.269	(not used)						

## Выдающиеся направления (topic prominence)

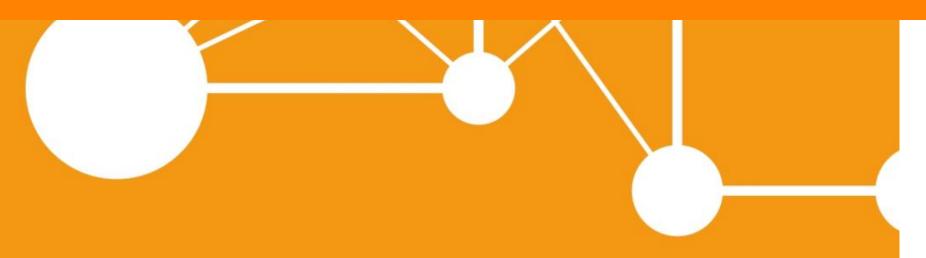
- Почему "Prominence"
- Prominence ≠ Importance (тема может быть важной, но не выдающейся)
- Prominence ~ Visibility или Momentum





**Характеристики модели – корреляция с грантами** 





#### Корреляция с грантами

- Распределение грантов по темам
- Корреляция грантов с выдающимися направлениями

#### Распределение грантов по темам

- U.S. Star Metrics data (2008-2014)
  - В основном из NIH, NSF
  - 364,000 грантов , \$253 млрд



- Наличие описания (word profiles) для каждой из 91,726 тем модели
- Заглавия и аннотации также доступны для большинства грантов
- В результате: 314,000 грантов, \$203 млрд были соотнесены по темам созданной модели
  - Некоторые гранты имели недостаточно текста

# Корреляция грантов с выдающимися направлениями

- Гранты были разделены на два периода (2008-10, 2011-13)
- Показатель Prominence был рассчитан для 2010

Table 4. Co	orrelation r	natrix for	variables	considered	in the	funding	prediction	analysis.
I WOIC II C	JII CIMUIOII I	HUULIA IOI	' til ittbies	combiactea	III CIIC	Iununi	prediction	ttitti y 515.

	L:Fund1113	L:Fund0810	Prominence	Vitality	L:Authors
L:Fund1113	1.000				
L:Fund0810	0.837	1.000			
Prominence	0.606	0.616	1.000		
Vitality	0.166	0.162	0.314	1.000	
L:Authors	0.160	0.171	0.242	0.202	1.000

- Корреляция грантов обоих периодов между собой
- Prominence коррелируется с грантами по обоим временным периодам

#### Корреляция грантов с выдающимися направлениями

- Прогнозирование
  - $\circ$  Fund1113 predicted by Fund0810 R<sup>2</sup> = 0.707
  - Fund1113 predicted by Prominence  $R^2 = 0.367$
  - $\circ$  Fund1113 predicted by Fund0810 and Prominence  $R^2 = 0.713$
- Будущие гранты могут быть спрогнозированы прошлыми грантами, но добавление показателя prominence повышает прогнозирование
- Тем не менее, мы не можем использовать прошлые гранты для прогнозирования в большинстве случаев из-за отсутствия информации
  - Star Metrics уникальная база, но с ограниченным объемом данных
  - Практически не существует данных подобного уровня, с суммами, публично доступные
- Поэтому, prominence чрезвычайно ценный показатель, поскольку может выступать индикатором для финансирования

#### Примеры с выдающимися направлениями и финансированием

- High prominence high funding (HH)
- High prominence no funding (HN)
- Low prominence high funding (LH)

## **High prominence – high funding**

Topic	<b>#U.S.</b>	Prom	Funding	#Pub	#Pub	Description	Discipline
	Author	2010	2010 (\$	(2008-	(2011-		
	(2010)	(pctl)	million)	2010)	2013)		
				High I	Promine	nce – High Funding (HH)	
2538	674.2	99.8	420.6	807	2379	next-generation DNA sequencing	Cell Biology
73	364.5	99.3	305.2	1597	1960	T-lymphocytes	Immunology
1544	181.7	98.3	189.4	747	1297	orbitofrontal cortex and reward	Neurodeg Disease
1493	223.7	99.0	180.1	742	2070	default mode network (brain)	Brain, Vision, Hea
2771	246.1	98.4	150.2	753	1204	inflammation and obesity	Diabetes
5042	209.1	95.1	143.4	396	518	autism phenotype	Psychiatry
236	338.7	99.1	80.1	1215	1675	peptide identification in proteomics	Analytical Chemi
205	216.5	98.2	41.6	1053	1523	amyloid function in Alzheimer's	Neurodeg Disease
2646	215.1	98.1	37.3	677	945	solid-state nanopores	Nanochemistry
2877	128.3	96.4	33.0	472	912	BPA and endocrine disruption	Environ Chemistr

• Очевидно: финансирование высоко значимых тем

### **High prominence – no funding**

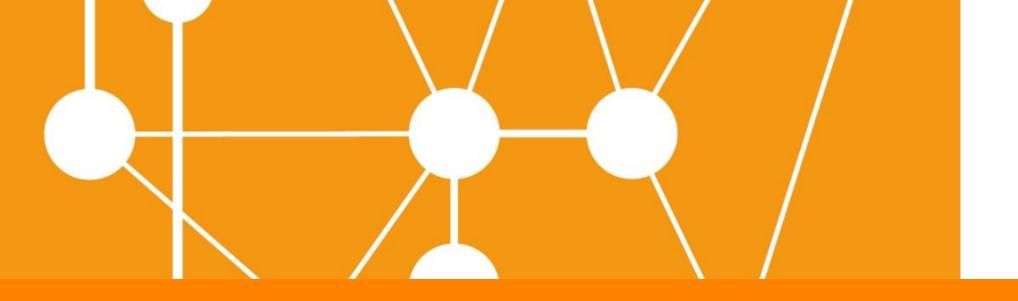
Topic	#U.S. Author (2010)	Prom 2010 (pctl)	Funding 2010 (\$ million)	#Pub (2008- 2010)	#Pub (2011- 2013)	Description	Discipline
		`•		High	Promine	ence – No Funding (HN)	
2187	0.5	95.8	0	591	849	electrochem degradation in wastewater	Electrochemistry
25	1.6	96.5	0	785	1532	corrosion inhibitors (steel)	Materials
135	2.7	98.2	0	1118	2136	dye remediation in effluents	Electrochemistry
15	4.2	98.7	0	1263	1961	biosorption of heavy metals	Electrochemistry
4003	7.8	98.8	0	547	878	dispersive liquid-liquid micro-extraction	Environ Chemistr
566	9.2	94.7	0	608	753	properties of olive extracts	Animal Science
4594	14.4	96.9	0	540	841	hollow nanoparticles	Electrochemistry
580	21.1	96.8	0	1010	1847	phosphors for LEDs	Optical Materials
644	21.1	95.4	0	725	777	hydrogen energy storage	Electrochemistry
7	74.3	99.3	0	1963	2407	Zn0 nanostructures	Semicond Physics

- Исследования «за пределами» США отсутствие грантов США
  - Вопросы окружающей среды для крупных производств, которых нет в U.S.
  - Оливки не растут в США
  - Исследования материалов, вопросы энергии в которых доминируют китайские университеты

### Low prominence – high funding

Topic	#U.S. Author (2010)	Prom 2010 (pctl)	Funding 2010 (\$ million)	#Pub (2008- 2010)	#Pub (2011- 2013)	Description	Discipline
				Low P	rominer	nce – High Funding (LH)	
83249	0.1	0.3	11.9	6	10	loose topic - clinical investigation centers	Patient Care
25667	2.2	2.2	24.8	6	13	academic medical centers, enthusiasm	Patient Care
54378	2.3	12.8	89.3	20	20	forestry education	Agricultural Polic
38569	2.7	14.2	13.7	27	29	agroecology, sustainability	Agricultural Polic
54158	5.0	4.6	12.0	18	60	loose topic – DMZ, networks, protocols	Computing
33105	10.9	6.1	29.6	42	60	role of nursing in clinical trials	Patient Care
18741	30.1	3.5	23.0	71	137	capstone projects, engineering	Learning
33702	30.4	11.2	18.3	51	77	extension programs and learning	Agricultural Polic
38645	31.1	8.6	10.9	72	92	systems engineering competency training	Management
26483	39.8	9.5	44.2	99	176	civil engineering program criteria	Learning

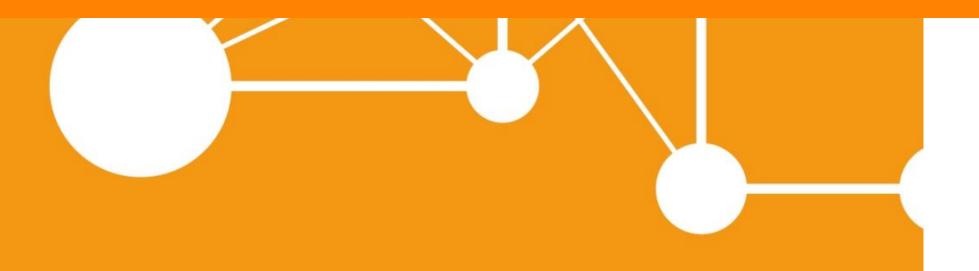
- Распределение грантов по темам из-за использования широких терминов
  - "clinical trials", "medical center"
  - Education более широкое понятие, используемое во многих грантах NIH, NSF
  - Маленькие темы с недостаточным текстом





# Обновление тем



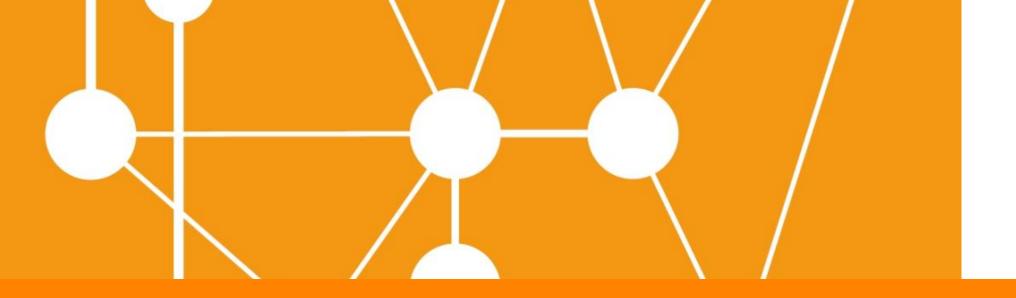


#### Обновление тем

- Для поддержания стабильности тем, нет необходимости совершенно новой кластеризации и создания модели каждый год
- Вместо этого добавляются новые работы из Scopus к существующим темам
  - Это может быть сделано очень точно с использованием ссылок из каждой работы
- Тем не менее, каждый год могут возникать новые направления и это должно быть отражено в модели (в продукте SciVal, с использованием VOS алгоритма)
- Таким образом, помимо добавления статей к уже существующим темам планируется ежегодно создавать несколько новых тем

# Создание новых направлений

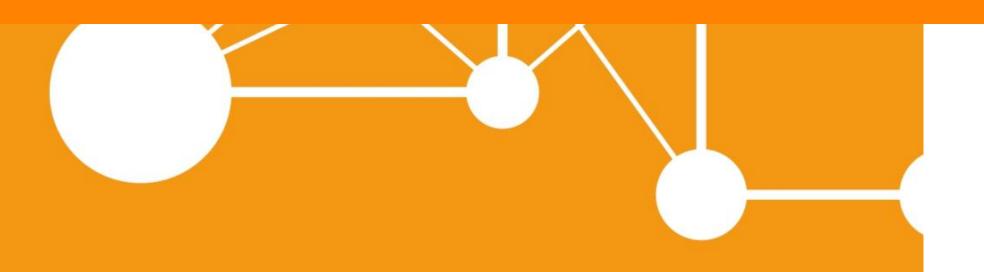
- Новые темы имеют следующие свойства
  - Относительно небольшая на момент появления с высоким темпом роста
  - Событие например, прорывная статья которое действует в качестве триггера для возникновения
  - Постоянство тема не угасает
- Новые темы, во всех случаях, связаны с существующими темами они не материализуются из воздуха
- Большинство возникающих тем встроены в более крупные существующие являются кандидатами для разделения тем
  - Основано на сравнении моделей разных лет
- По экспериментальным данным SciTech появление 30-50 новых тем ежегодно
- Алгоритмическая идентификация возникающих тем может быть дополнена ручной идентификацией новых тем на основе текущих событий в науке





# Исторический аспект





#### История

см. Atlas of Science: Visualizing What We Know, by Katy Börner

- 1985 ISI (теперь Clarivate) разработка Research Fronts
  - Библиометрический способ выявления исследовательских возможностей
- 1988 CRP (сейчас SciTech) разработка **Research Communities** 
  - Те же алгоритмы и ниже пороги для более широкого охвата
- 2007 SciTech разработка Distinctive Competencies
  - о Кластеризация научных сообществ используя научные преимущества организации
- 2015 SciTech разработка Торісs
  - о Существенно увеличивает охват и точность
- 2017 SciTech разработка индикатора Topic Prominence
  - Используя число ссылок, загрузок и метрику влияния журнала
  - Библиометрический индикатор может использоваться для прогнозирования грантовых паттернов

#### История

• Research Fronts (1985) 2% охват 10,000 кластеров

• Research Communities 4% охват 35,000 кластеров

• Distinctive Competencies 15% охват 200,000 кластеров

• **Topics** 95% охват 100,000 скластеров

- Topic Prominence (2017) прогнозирование финансирования
  - Практически полный охват и точность модели спроса и предложения для науки

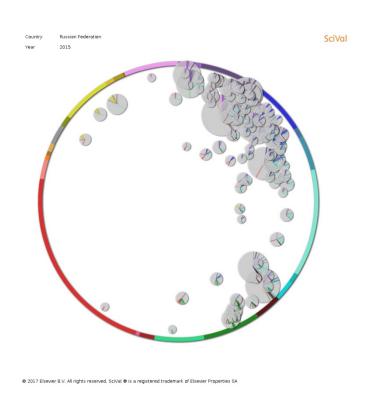
#### Заключение

- Была создана точная модель науки с около 100,000 темами, подходящая для анализа научного ландшафта
  - Методология может быть полностью воспроизведена, но требует наличия полного объема данных
- Был разработан показатель на уровне тем Prominence который коррелирует с грантовым финансированием
- Финансирование на одного автора возрастает с ростом показателя prominence
- Созданная модель науки (темы) и показатель prominence помогает в принятии решений в области науки

### Реализация проекта в SciVal

Elsevier работает совместно с SciTech Strategies с 2008 года над новыми подходами и концепциями научного планирования и прогнозирования

Релиз Topics и выдающихся направлений – <u>03 октября 2017</u> (с Днем Рождения, Clarivate Analytics!)

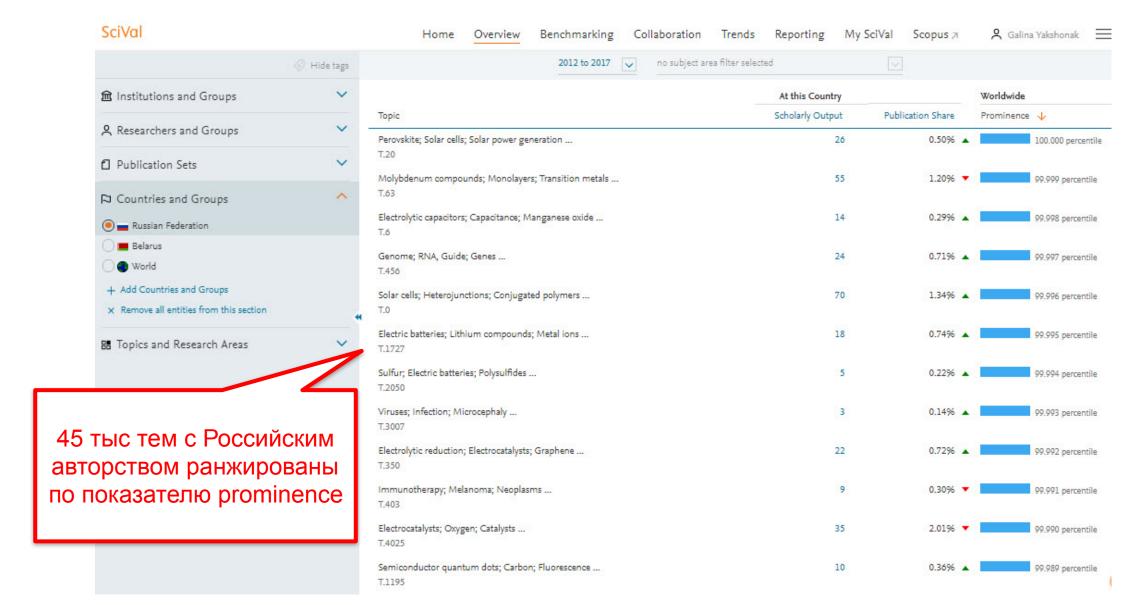


SciVal: 144 компетенции России

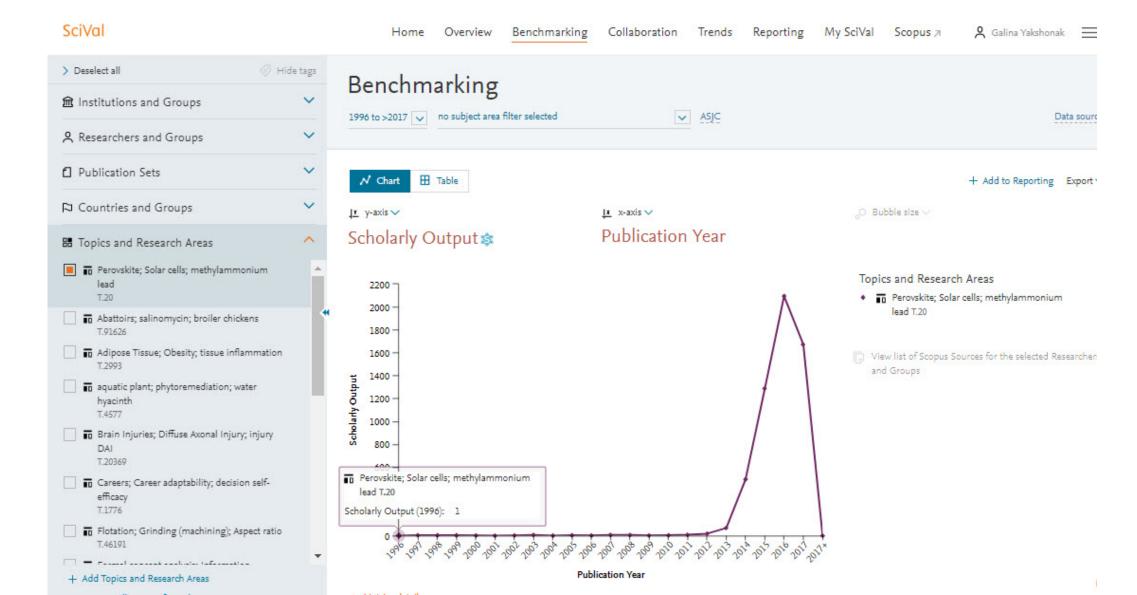


SciVal: 45090 научных тем с Российским участием

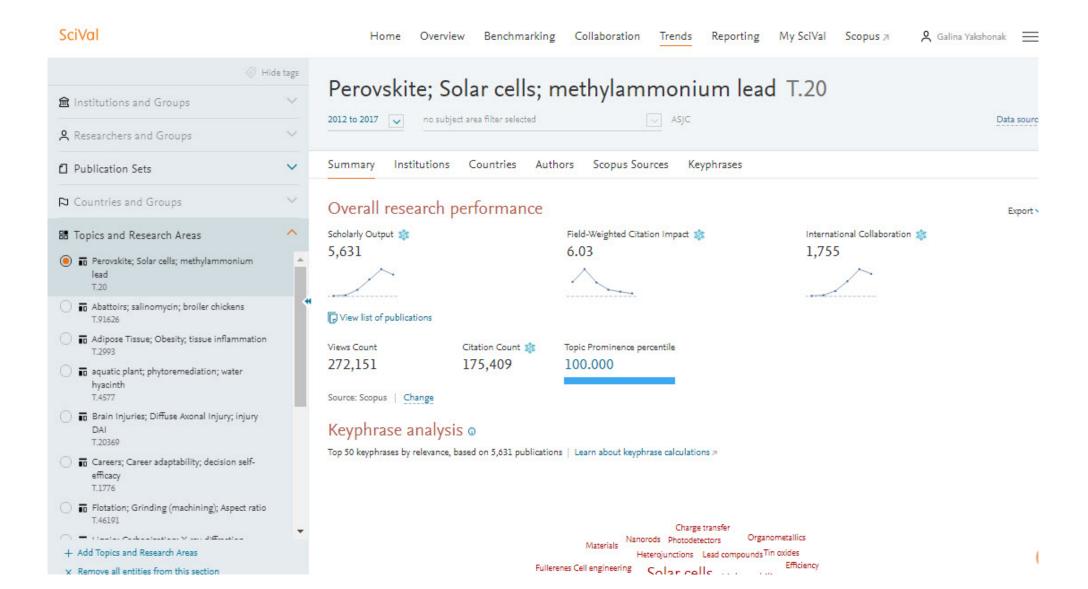
#### Топ-выдающиеся направления с участием России



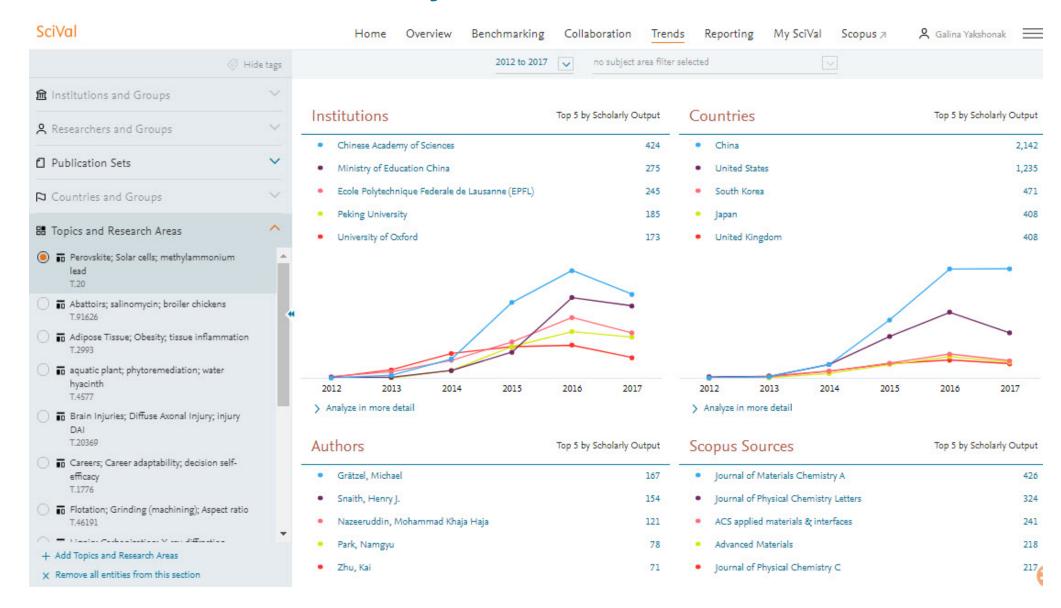
## Perovskite; Solar cells; methylammonium



## Perovskite; Solar cells; methylammonium



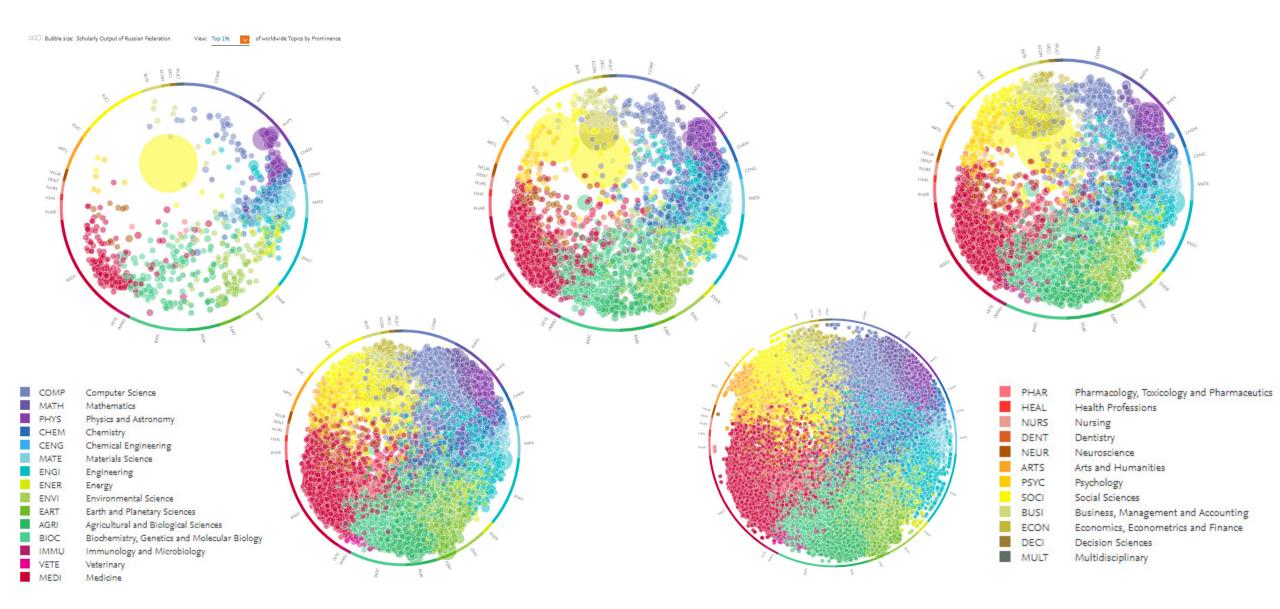
## Perovskite; Solar cells; methylammonium



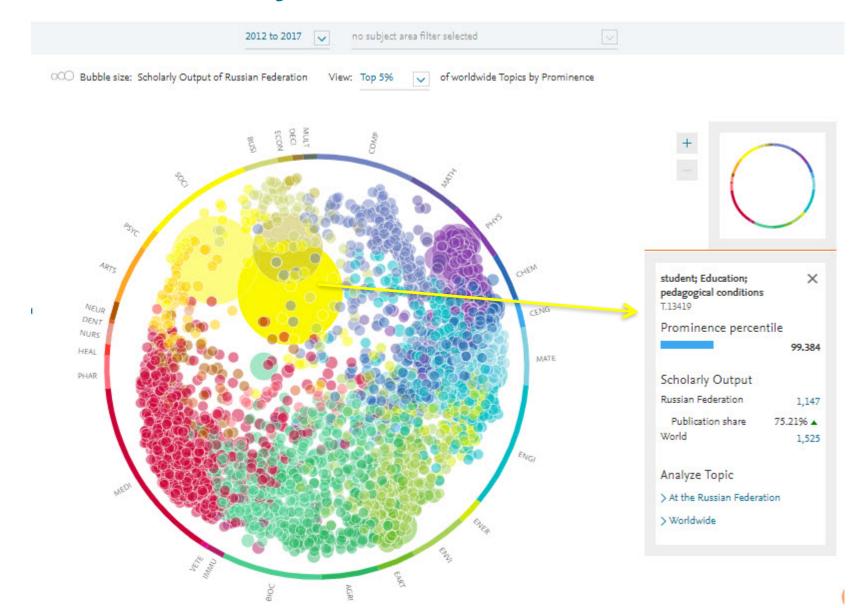
# Российские участники Perovskite; Solar cells; methylammonium

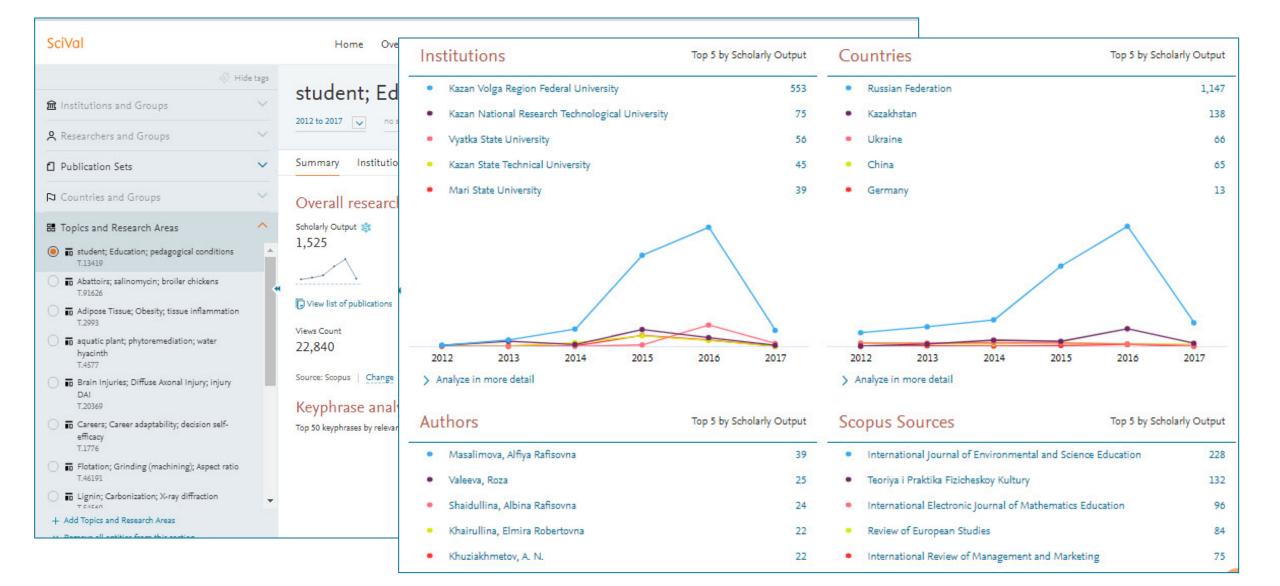


#### Russia: 1% 5 % 10 % 25% All



# Prominence ≠ важное/актуальное





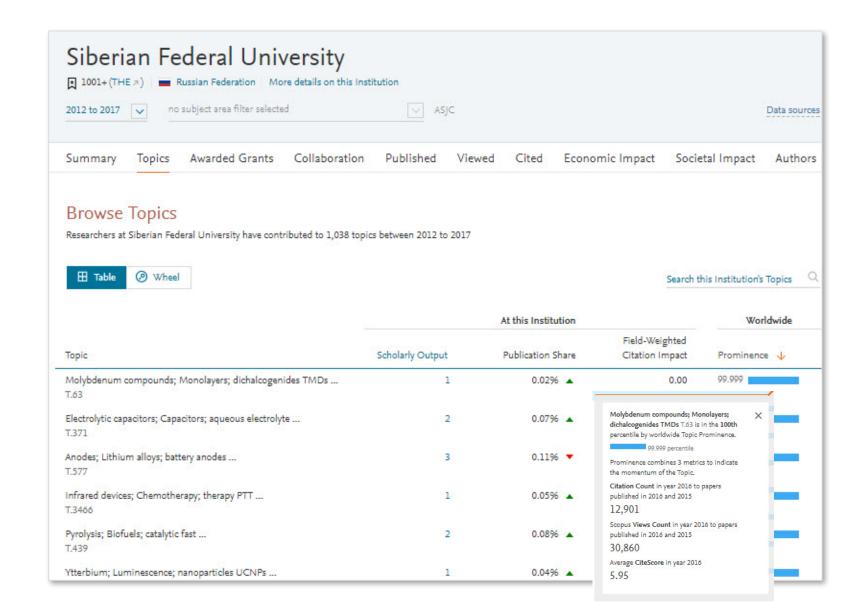
#### На практике:

"В каких выдающихся направлениях участвую я и моя организация?"



"Почему это направление выдающееся?"

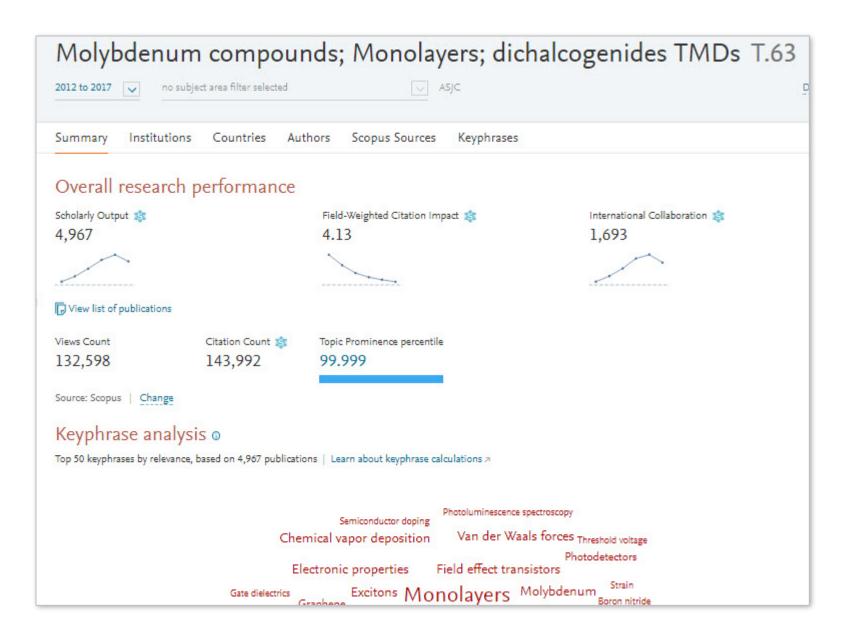




#### На практике:

"Какова динамика этой темы? Кто ее драйверы?"

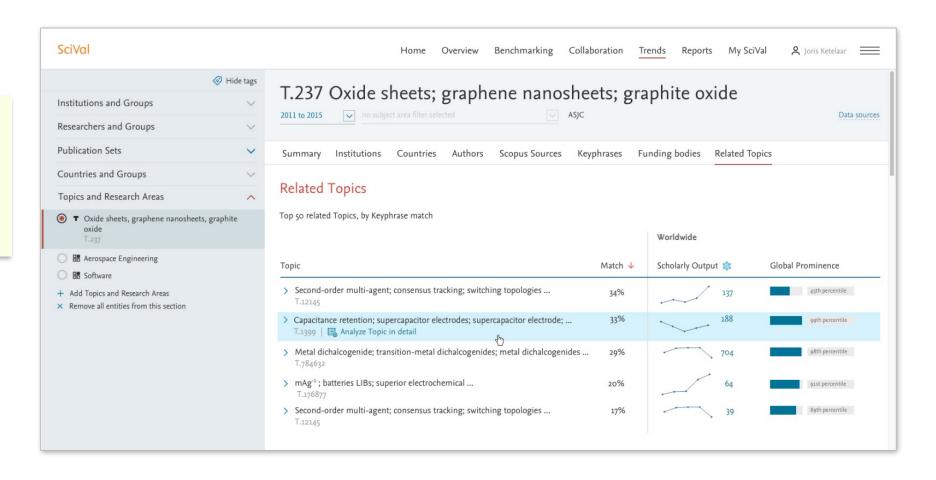


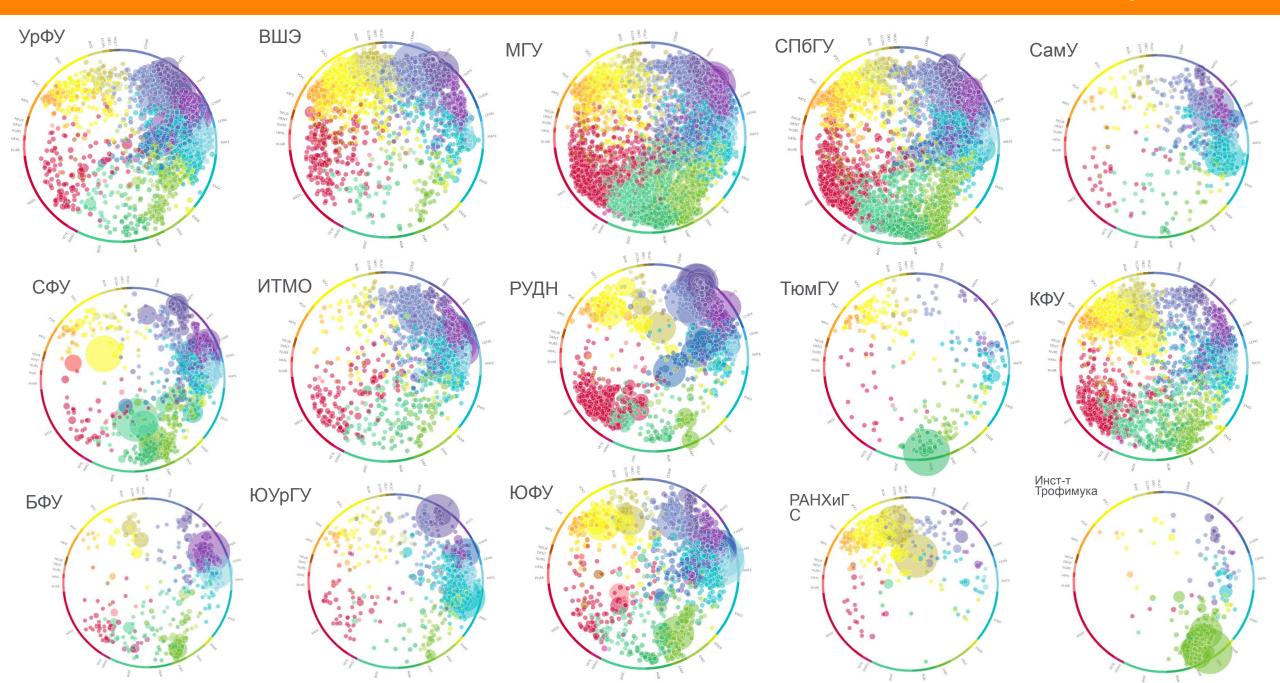


#### На практике:

"На какие схожие темы мне стоит обратить внимание?"







# Research Intelligence

# Спасибо!

www.elsevierscience.ru

#### Launching 3 October! Topic Prominence in Science

SciVal is evolving, and changing the way you work.

We are expanding SciVal from being a purely evaluative and analytical tool to being an **integral part of your** research planning process. Topic Prominence in Science will revolutionize the way in which you develop your research strategy by giving you and your colleagues unique insight into identifying new, emerging research trends.

What is Topic Using Topic Prominence Roadmap Replacing Competencies Methodology Prominence?

You are now able to run a complete portfolio analysis to see which Topics your institution is currently active in, and which Topics have high momentum, those therefore more likely to be well-funded. It will provide insight into which researchers are active in those Topics, which Topics your peers and competitors are active in and the related Topics of which you should be aware.

Topics are ranked by Prominence, an indicator of the momentum of a particular field.

The development of Topic Prominence in Science is based upon extensive research and customer feedback. Unlike other research analytics solutions, which merely scratch the surface by only analyzing top-cited articles, we take the entire world of research into account. Our

