

MoocDB

Taming MOOC Big Data while Fostering Collaboration in Online Education Research

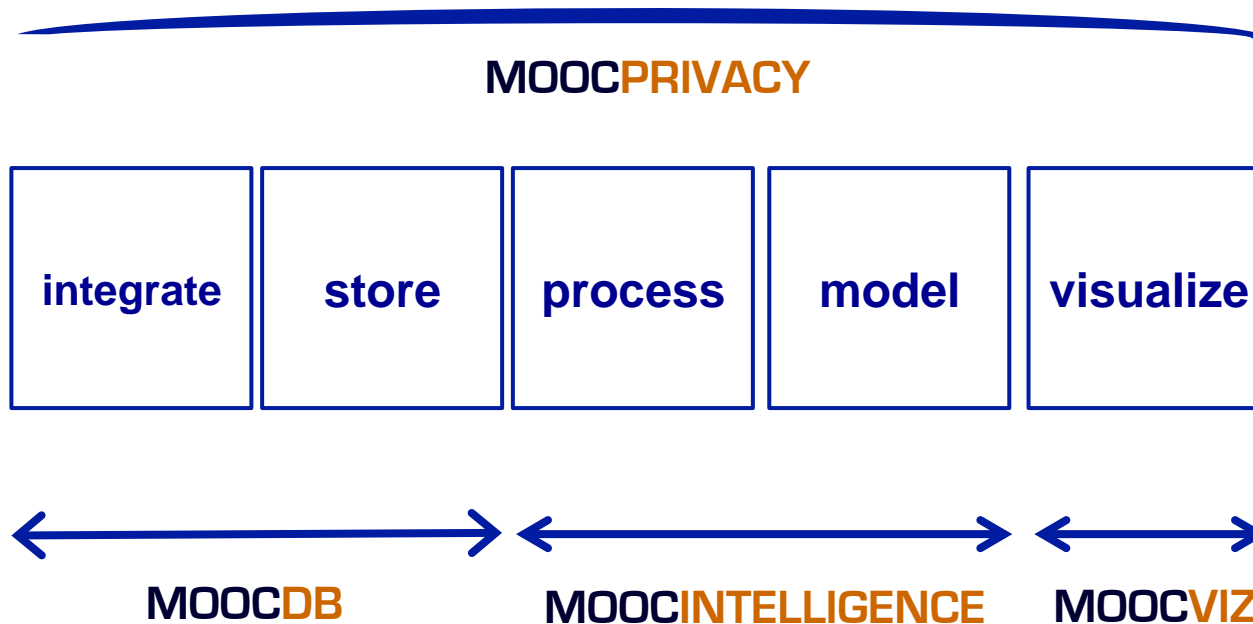
Una-May O'Reilly

**AnyScale Learning for All Group: ALFA
Computer Science and Artificial Intelligence Lab
MIT**

<http://groups.csail.mit.edu/ALFA/groupWebSite/index.php?n=Site.AlfaX>



ALFA MOOC Data Science



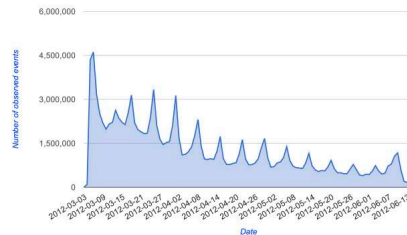
integration... to... insight

ALFA MOOC Research



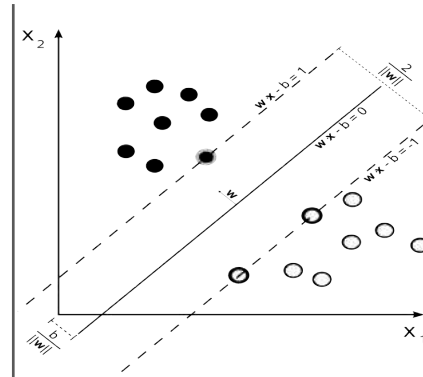
Shared data model

MOOCDB



Shared analytics

MOOCVIZ



Who is likely to stopout?
Community detection
Weekly topic analysis
Crowd Sourcing

Machine Learning

MOOCINTELLIGENCE



Access
Privacy Protection
Policy
Differential privacy

Privacy as a service

MOOCPRIVACY

Massive Online Open Courses



Circuits & Electronics
6.002x

Enroll in 6.002x Circuits & Electronics

6.002x (Circuits and Electronics) is an experimental on-line adaptation of MIT's first undergraduate analog design course: 6.002. This course will run, free of charge, for students worldwide from March 5, 2012 through June 8, 2012.

6.002x on MITx

If you successfully complete the course, you will receive an electronic certificate of accomplishment from MITx. This certificate will indicate that you earned it from MITx's pilot course. In this prototype version, MITx will not require that you be tested in a testing center or otherwise have your identity certified in order to receive this certificate.

ABOUT THE COURSE STAFF



Anant Agarwal

Director of MIT's Computer Science and Artificial Intelligence Laboratory (CSAIL) and a professor of the Electrical Engineering and Computer Science department at MIT. His research focus is in parallel computer architectures



MOOC Stakeholders

- ...Instructors
- ...Students
- ...Data custodians/guardians
- ...Course designers
 - ...Education technology specialists
- ...Education/Learning researchers

MOOC Research Questions

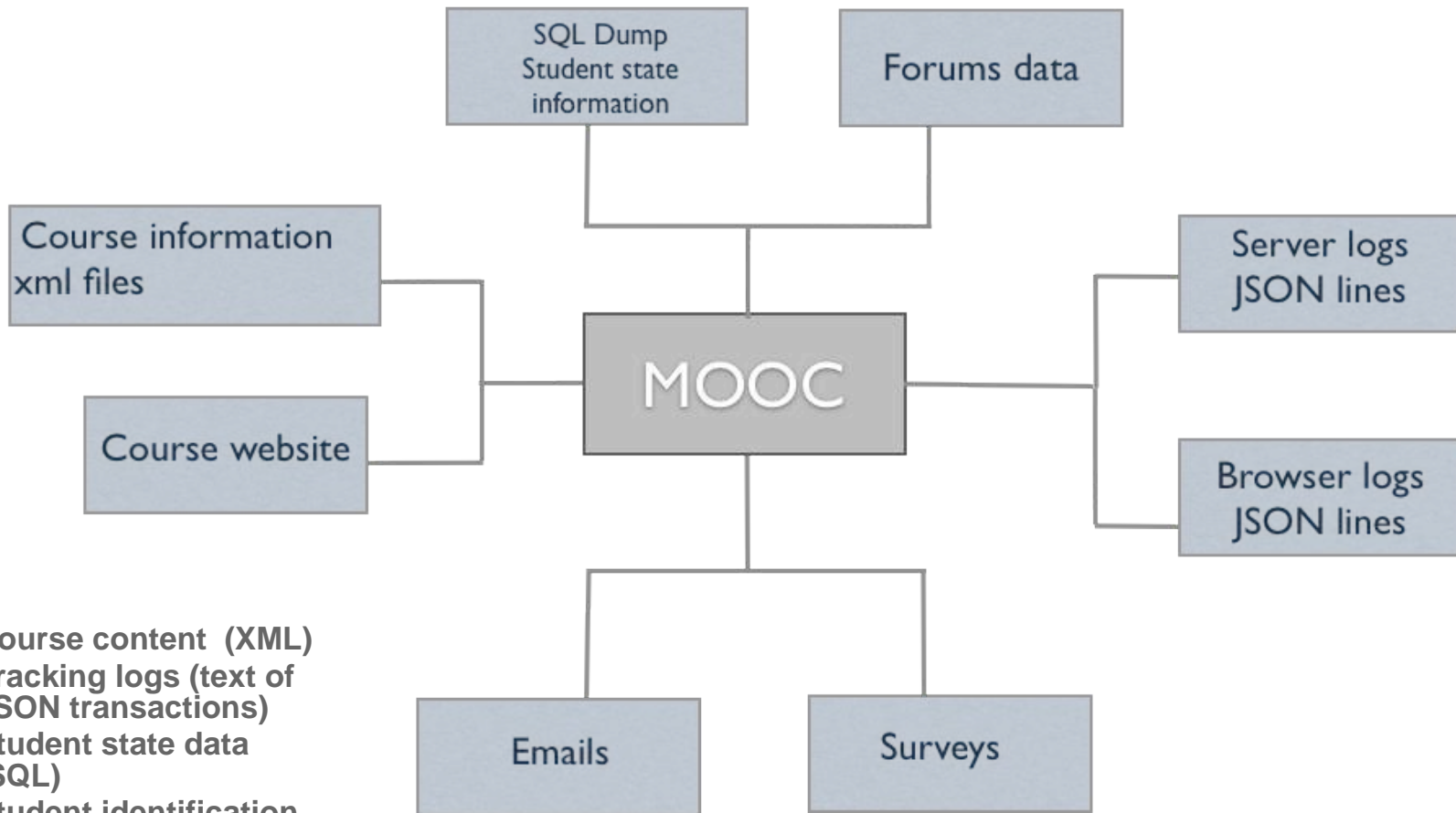
- **Descriptive information**
 - Who? When?
 - » Demographics and grades, statistical correlations
- **MOOC specific**
 - Trajectory related
 - Resource related
 - Using the crowd
 - Response related
- **General questions about learning and education**
 - Learning styles?
 - Knowledge acquisition
 - Flipped classrooms, blended learning

Behavioral Analysis

- ...Hypothesis
- ...Assemble data and features
- ...Statistical model
 - ...Validate, inspect, interpret
 - ...visualize



A MOOC Data Management Problem



- ... Course content (XML)
- ... Tracking logs (text of JSON transactions)
- ... Student state data (SQL)
- ... Student identification data (SQL)
- ... Forum data (NOSQL)
- ... Wiki data (SQL)

Pain Points and Bottlenecks

- ...Heterogeneous data formats
- ...Bloated raw data storage
- ...Lack of a comprehensive view of the data
 - ...Needs to be organized according to use!
- ...Un-identified cross-platform compatibilities
- ...Wasted effort replicating efforts of others

What about ...

Multiple courses?

Multiple platforms?

How can we bring many eyes to the data?

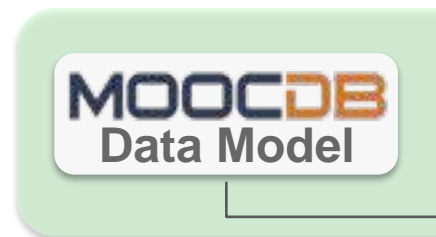
Enable and encourage community reflection and intellectual engagement around it



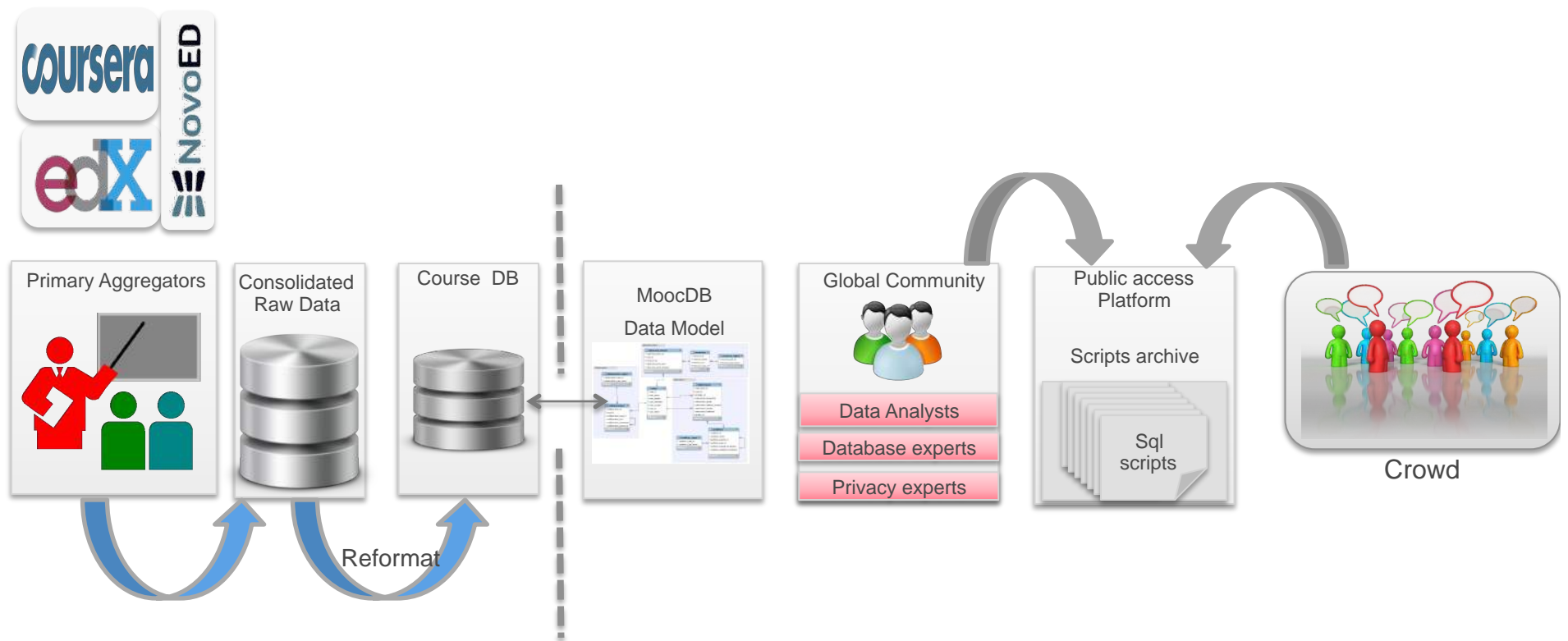
<http://www.flickr.com/photos/hagdorned/7434861784/>

MoocDB

Shared data model

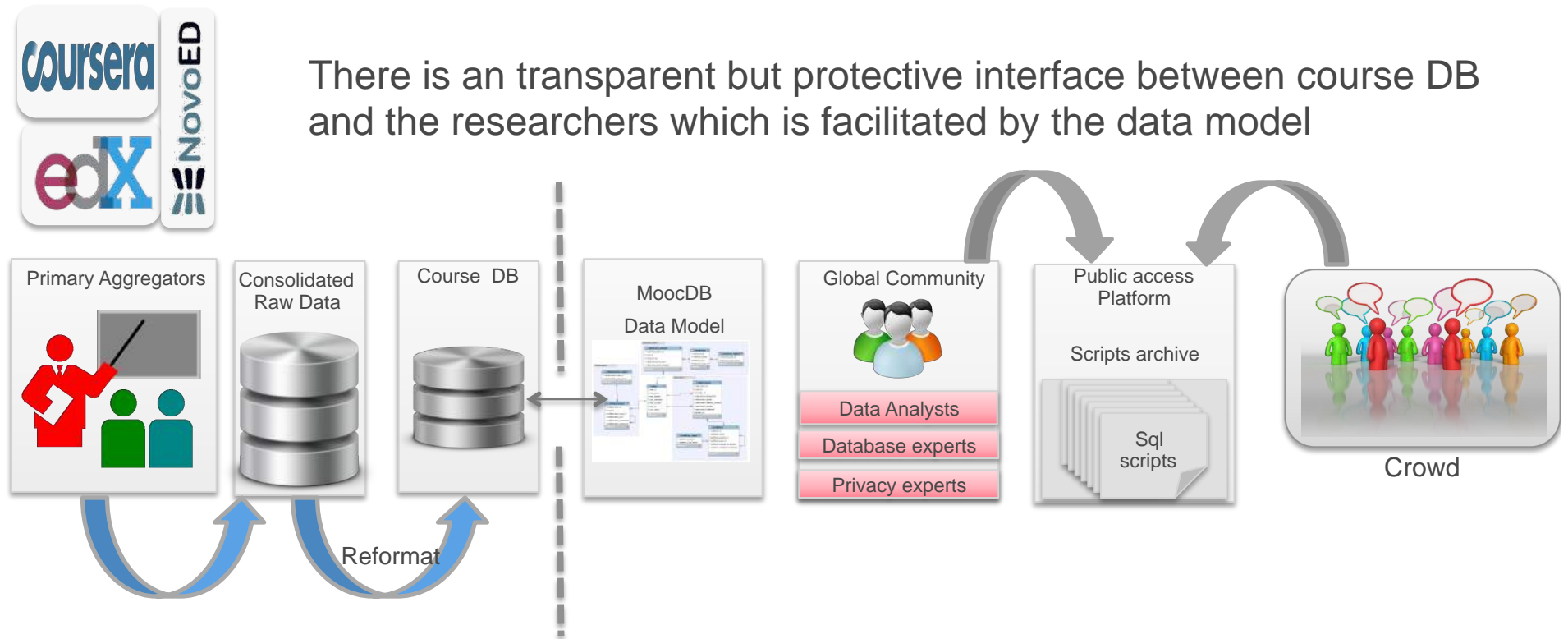


MoocDB



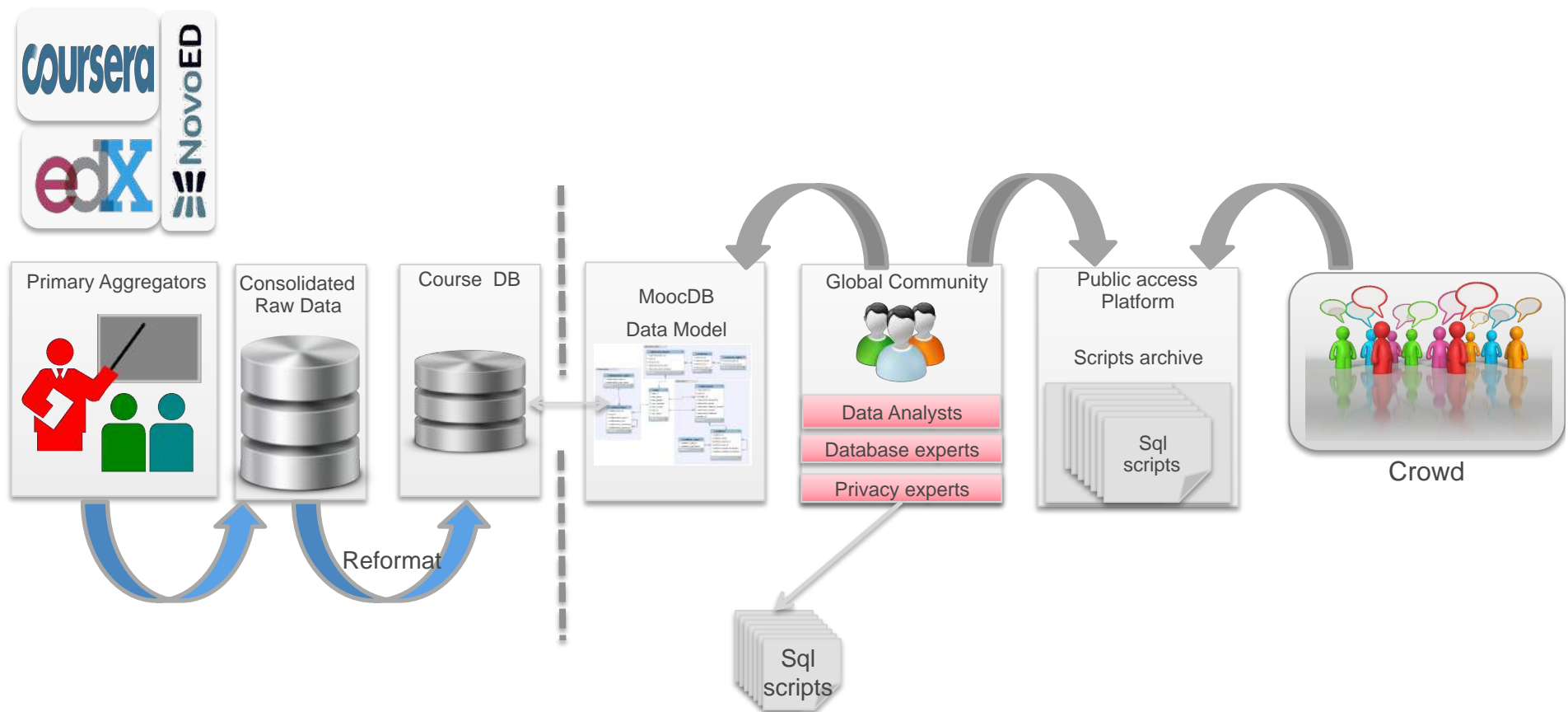
MoocDB: Data organization to support many eyes on the data

MoocDB



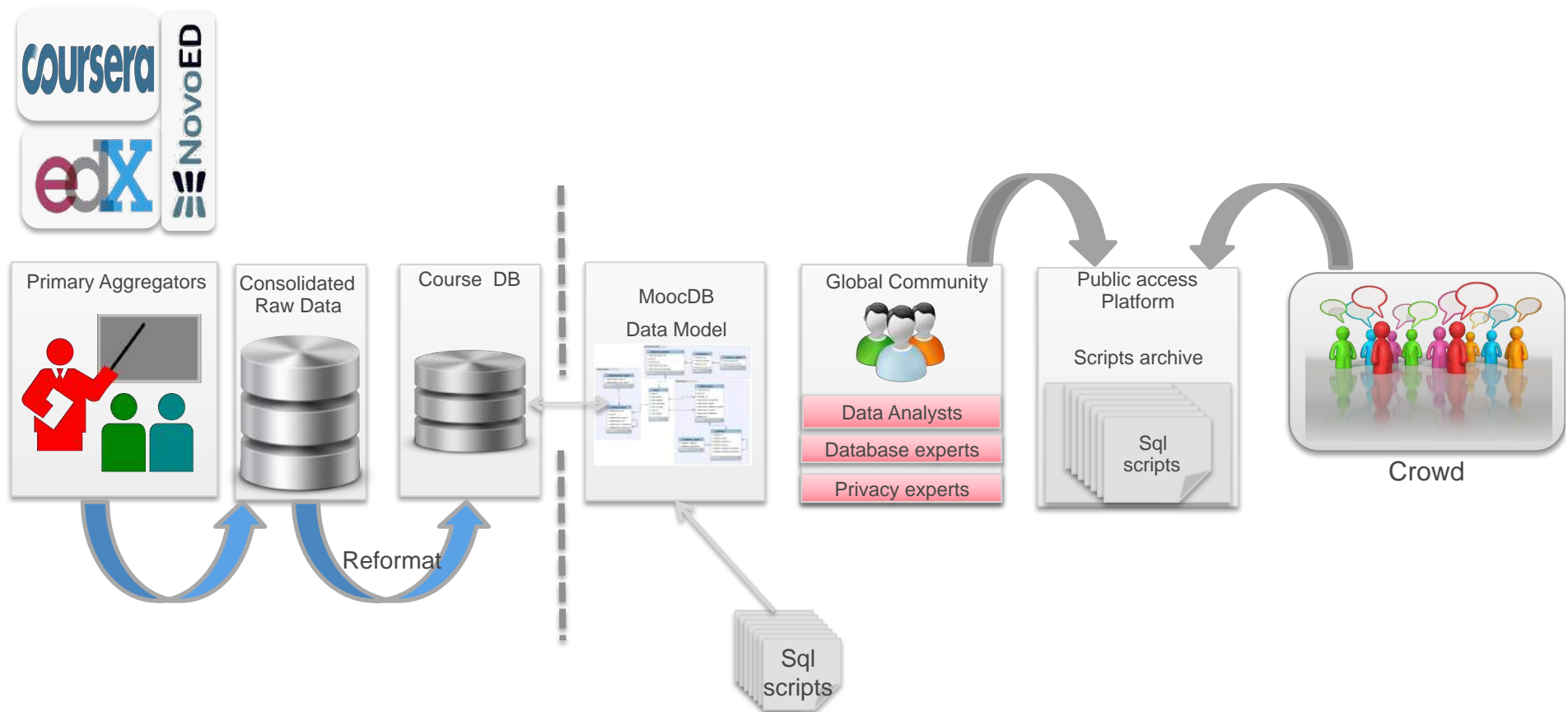
MoocDB: Data organization to support many eyes on the data

MoocDB



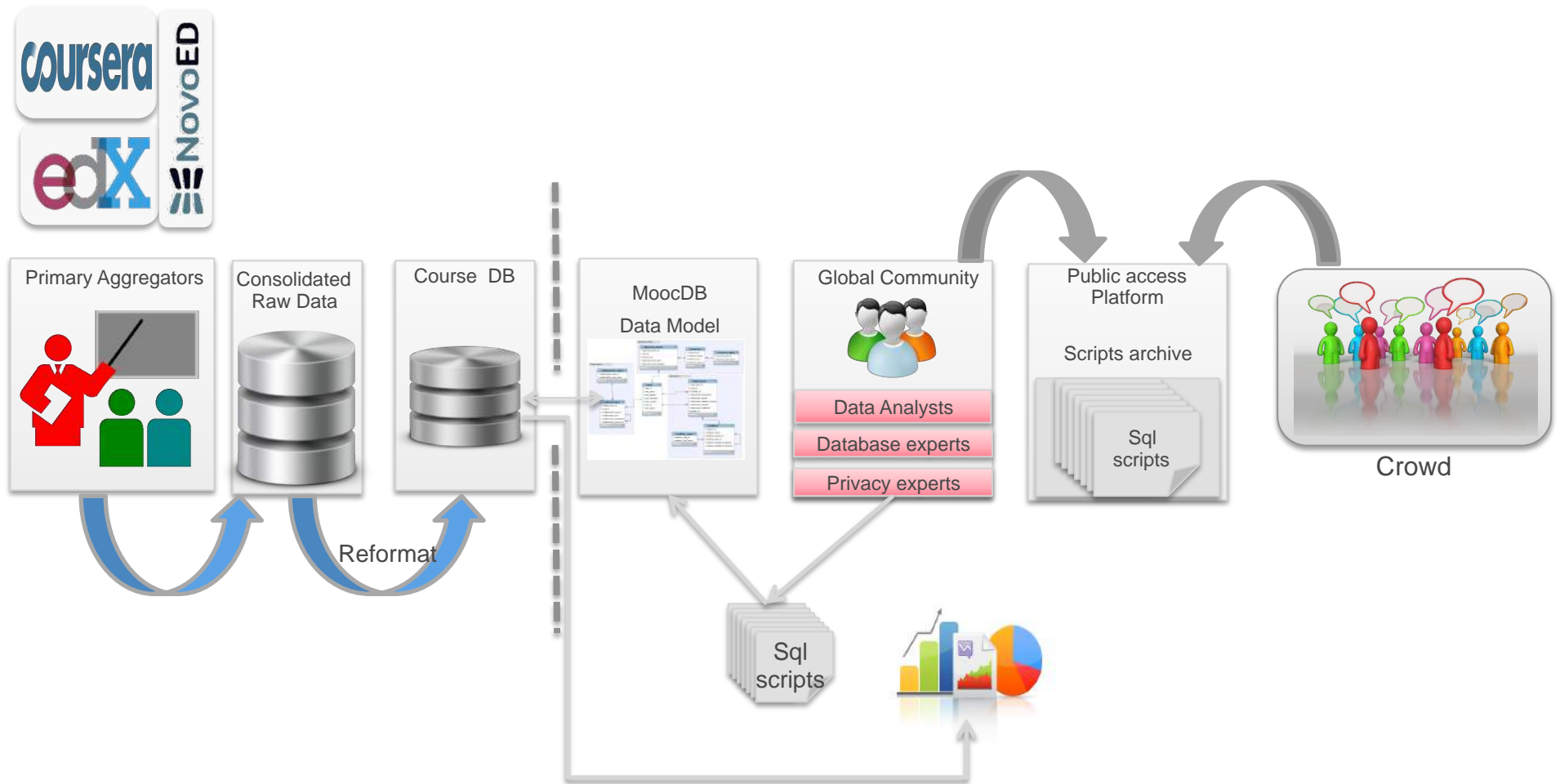
Community Visualization and Data Analysis:
Step 1: analysts write scripts by consulting the data model

MoocDB



Step 2: their scripts use the schema to reference the data in the course DB

MoocDB



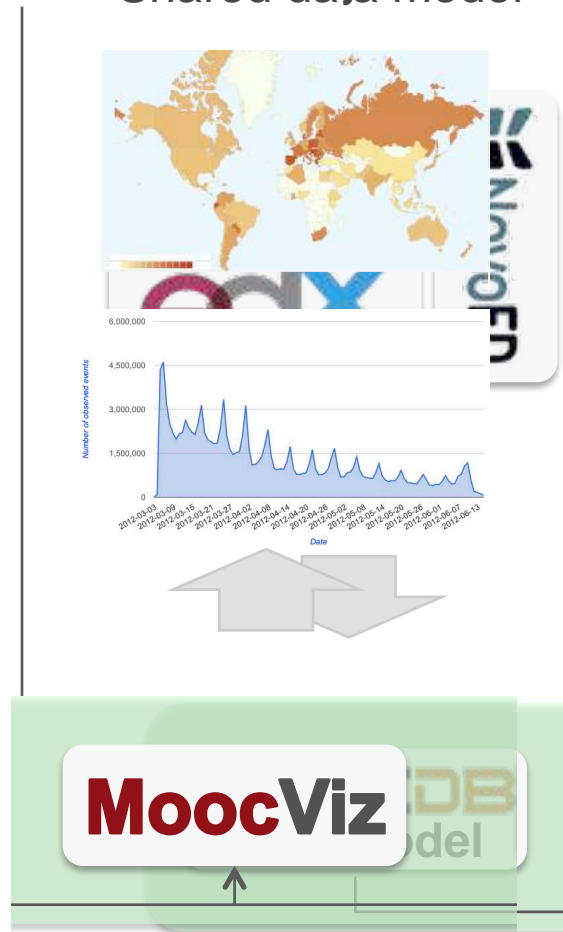
Step 3: The script executing over the data, referencing the data model, allows the insights from the course DB to be returned

MOOCDB supports multiple frameworks

- ...Our 6.002x DB using MOOCDB model
 - ...17 million submission mode events
 - ...150M observing mode events
 - ...96K collaborative events
 - ...Collapsed from 60GB to 6 GB
- ...Multiple Frameworks based on MOOCDB
 - 1....Export of data from course db
 - 2....Interoperability with programming languages
 - 3....Privacy protection via differential privacy
 - 4....Visualization and analytics -> MOOCVIZ

MoocVIZ

Shared analytics
Shared data model

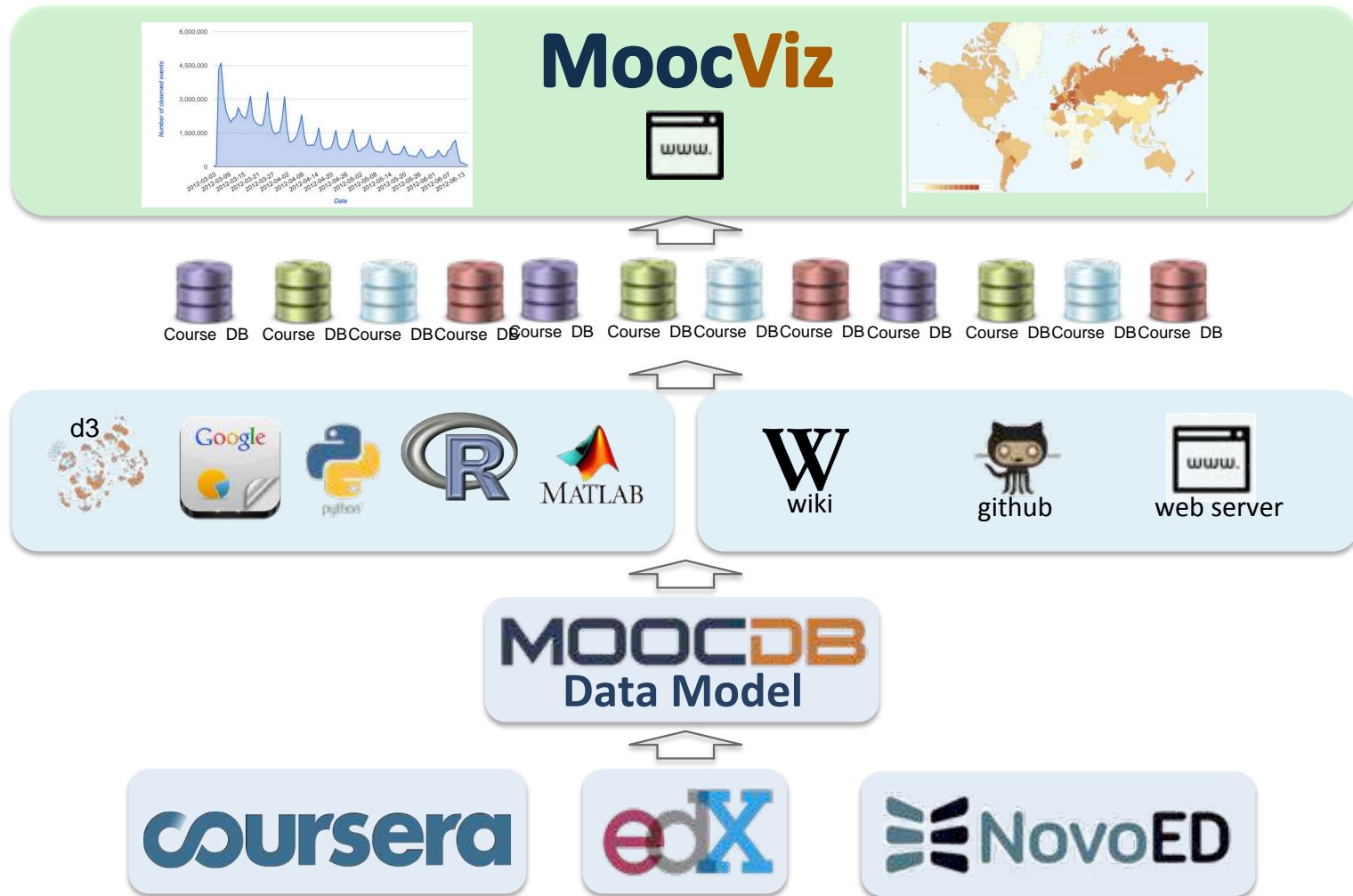


MOOCDB and MOOCVIZ

Use MOOCVIZ to demonstrate

- ...MOOCDB's collaboration support
 - ...On Stanford and MIT courses
 - ...On 2 different platforms EDX and COURSEERA
- ...New analytics around resource usage
 - ...Visualization and statistical support

MoocViz Analytics Platform

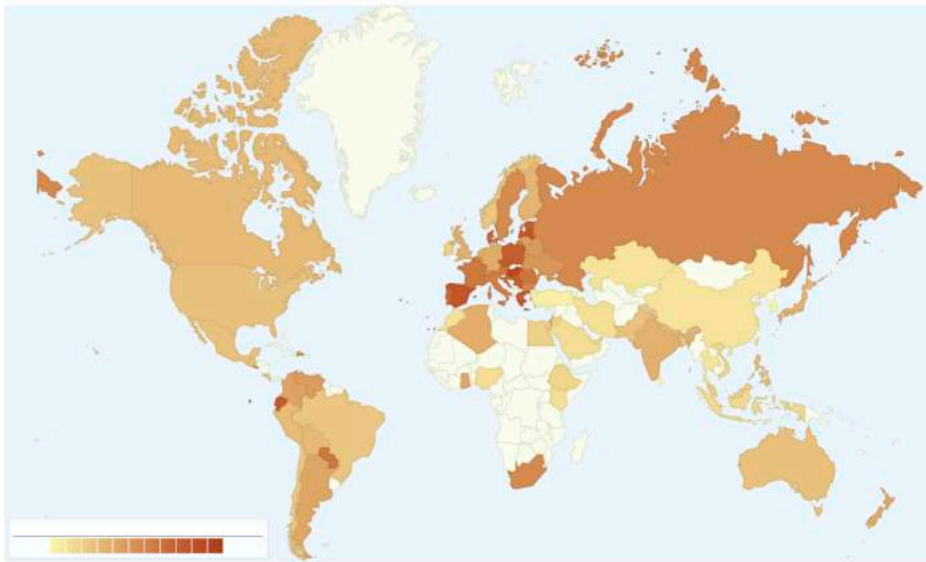


MoocVIZ User Types

- **Arms-length observers**
 - Checking the website
 - **Technology-savvy crowd**
 - Vote on utility of a visualization
 - Contribute s/w from other domains
 - **Course instructors**
 - **MOOC providers**
 - **Education researchers**
- Script developers**

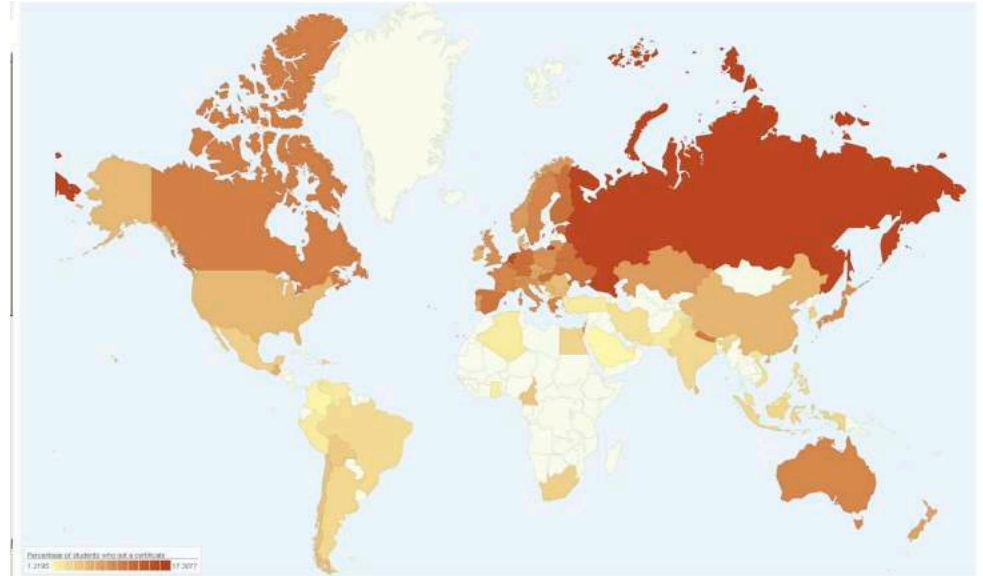
	Coursera course	edX Course
Title	Cryptography I	Circuits and Electronics (6.002x)
Instructors	Dan Boneh	Anant Agarwal, Gerald Sussman, Piotr Mitros
University	Stanford University	MIT
Length	6 weeks	14 weeks
Platform	Coursera	edX
Start date	Jan 13th, 2013	March 5, 2012
Registrants	21,744	154,763

MoocViz



6.002x
user certificates per country,
normalized, cutoff 100

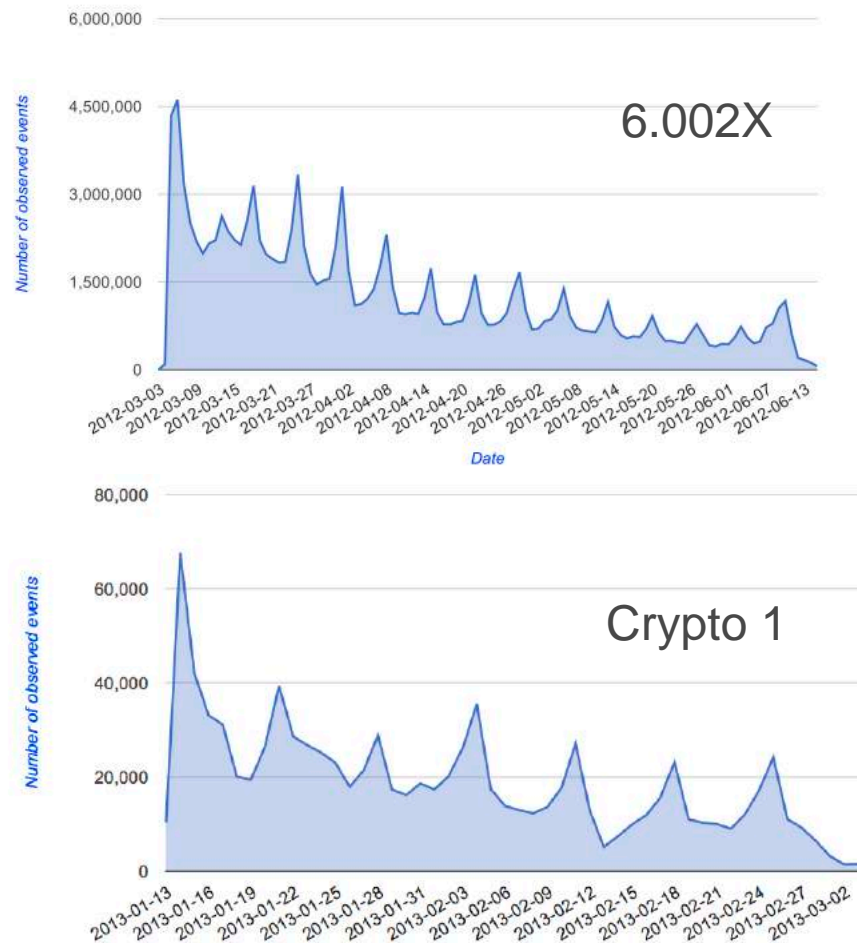
Hungary 16.2%
Spain 14.55%
Latvia 14.40%



Crypto 1-Stanford,
user certificates per country,
normalized, 2 columns 100

Russia 17.4%
Netherlands 16.43%
Germany 12.95%

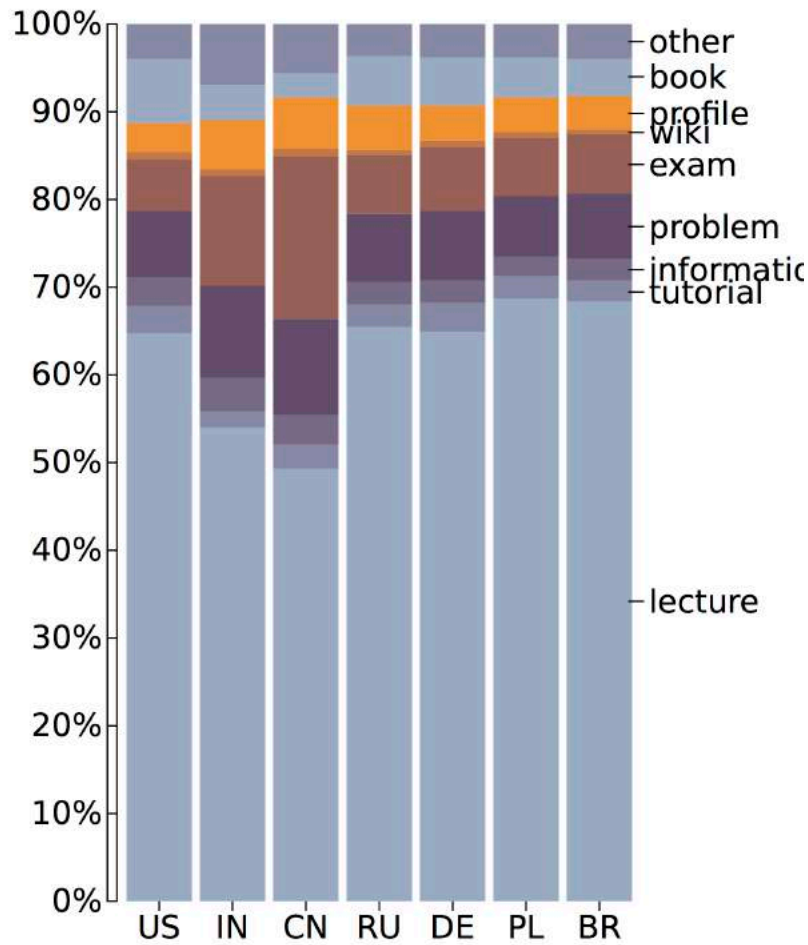
MoocViz



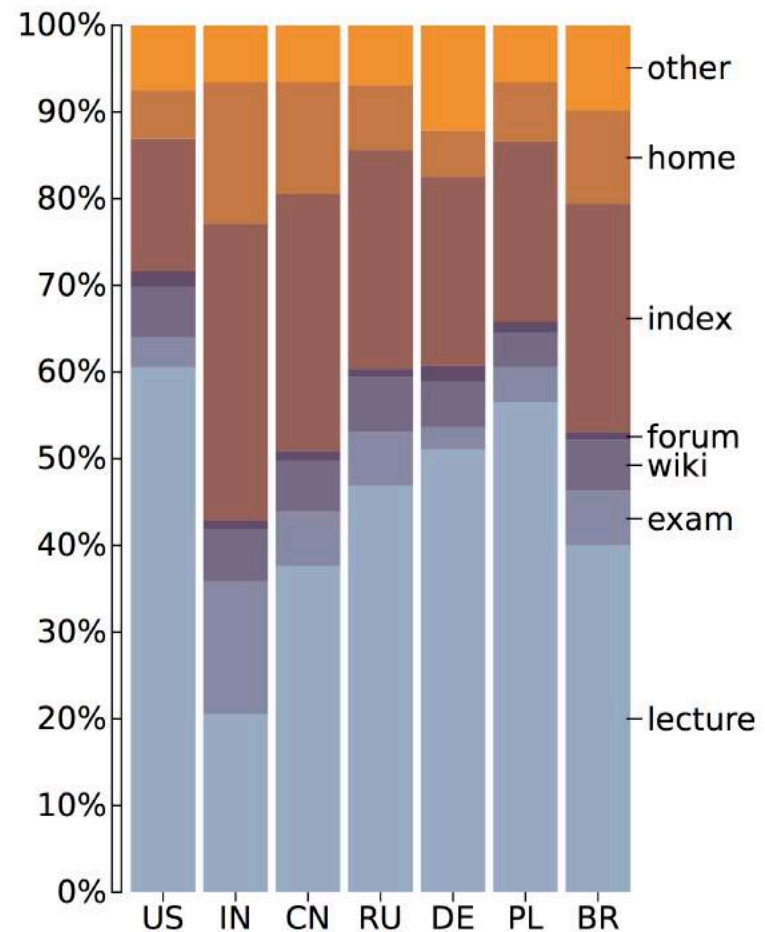
Resource Types

Resource id	Content	Medium
1	Lecture	Text
2	Lecture	Video
3	Tutorial	Text
4	Tutorial	Video
5	Informational	Any
6	Problems	Any
7	Wiki	Any
8	Forum	Any
9	Profile	Any
10	Index	Any
11	Book	Any
12	Survey	Any
13	Home	Any
14	Other	Any

Studying Resource Use

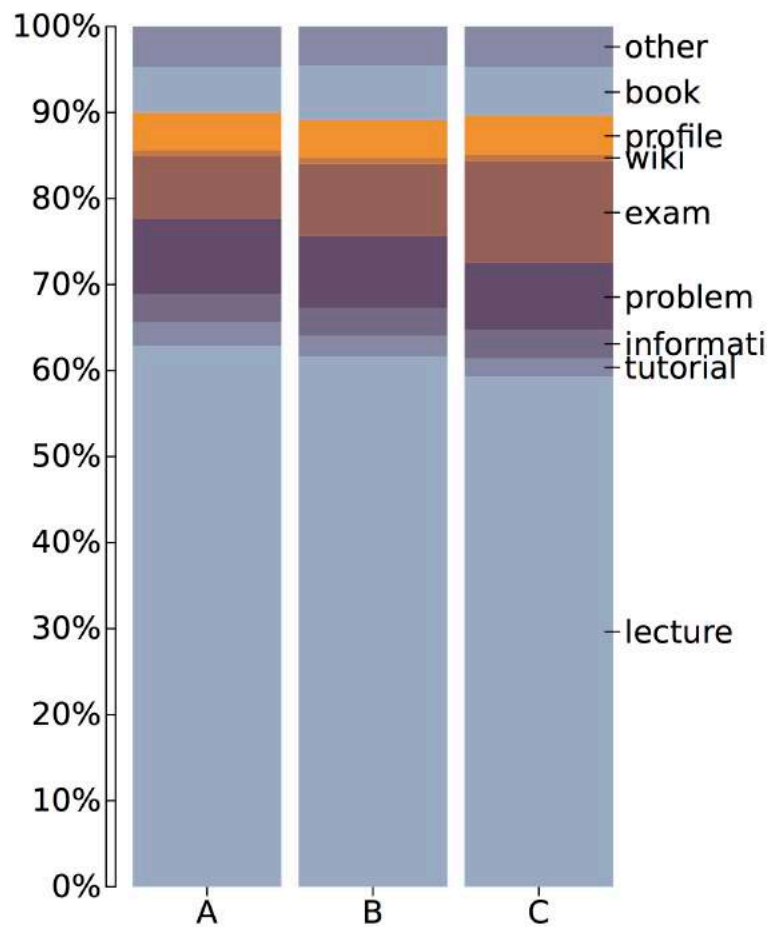


6.002X

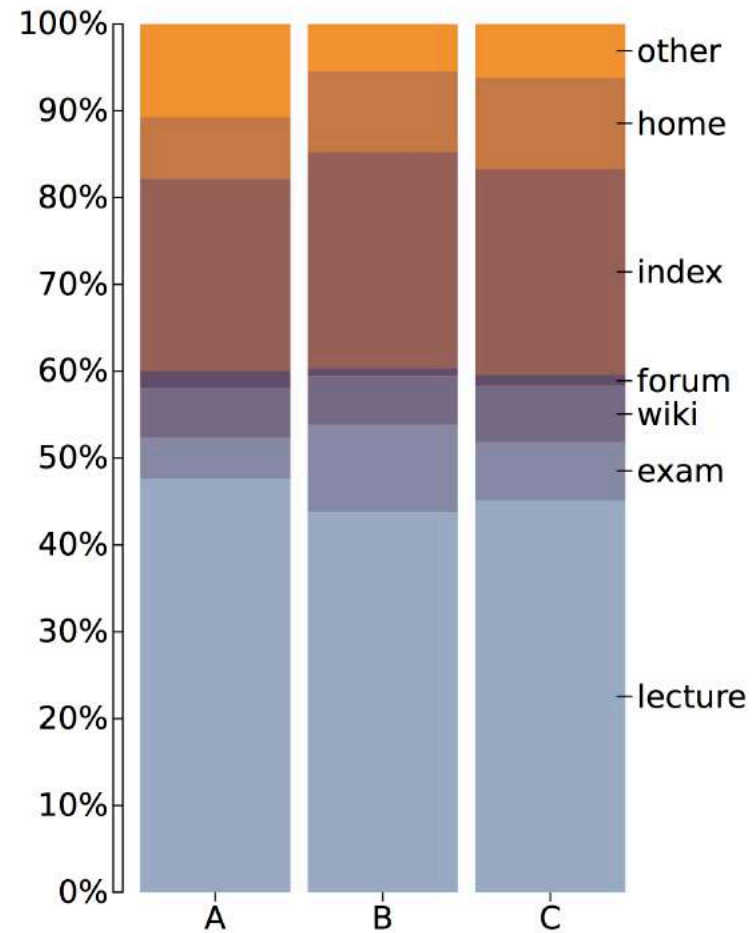


Crypto 1

Studying Resource Use



6.002x



Crypto 1

Analytics: Statistical Comparisons

	<i>Lecture</i>		<i>Exam</i>		<i>Problems</i>		<i>Book</i>	
	<i>l</i>	<i>u</i>	<i>l</i>	<i>u</i>	<i>l</i>	<i>u</i>	<i>l</i>	<i>u</i>
BR <i>vs.</i> DE	-35405	138970	-2563.1	9993.1	-7198.3	12247	-37926	13293
BR <i>vs.</i> IN	165720	287580	-1668.5	7106.2	2319.2	15908	-8910.5	26882
BR <i>vs.</i> RU	-60298	94243	-1072.3	10056	-5549.5	11684	-29058	16336
BR <i>vs.</i> US	-113120	4533	-5271.2	3201	-16114	-2993.5	-53044	-18485
DE <i>vs.</i> IN	102530	24721	-6205	4212.7	-1477.3	14656	55.123	42550
DE <i>vs.</i> RU	-121360	51742	-5455	7009.4	-9108.7	10195	-19468	31378
DE <i>vs.</i> US	-176650	-35499	-9832.1	332.03	-19948	-4207.7	-44179	-2717.7
IN <i>vs.</i> RU	-269690	-149660	-2548.6	6094.8	-12739	646.57	-32976	2281.4
IN <i>vs.</i> US	-313990	-247900	-6133.3	-1374.5	-22352	-14982	-54456	-35045
RU <i>vs.</i> US	-129150	-13383	-9695.1	-1358.9	-19076	-6166.2	-46405	-12401

Table 2: Analysis of the duration spent on resource type by country-based student cohorts for the edX course. The value of each cell is the 95% confidence interval given by the Tukey-Kramer method for the true difference of the means of the duration for the two cohorts indicated in the first column, e.g. students BR and DE in the first row. Cells with colored background indicate that the two cohorts have a significant difference in terms of mean of duration, which corresponds to the interval encompassing 0.

Analytics – Statistical Comparisons

	<i>Lecture</i>		<i>Exam</i>	
	<i>l</i>	<i>u</i>	<i>l</i>	<i>u</i>
BR <i>vs.</i> DE	-28253	-1117.2	-1024.1	592.15
BR <i>vs.</i> IN	-461.58	21961	-1335.2	0.2361
BR <i>vs.</i> RU	-20329	3236.1	-758.69	644.88
BR <i>vs.</i> US	-24271	-2475.1	-685.78	612.37]
DE <i>vs.</i> IN	15939	34930	-1017.1	114
DE <i>vs.</i> RU	-4025.8	16302	-446.33	764.43
DE <i>vs.</i> US	-7811.2	10435	-364.14	722.64
IN <i>vs.</i> RU	-25998	-12595	211.44	1009.8
IN <i>vs.</i> US	-29106	-19138	333.96	927.65
RU <i>vs.</i> US	-10989	1337	-346.88	387.29

Crypto 1: Comparison of country based cohorts

Analytics – Statistical Comparisons

	<i>Lecture</i>		Exam		Problems		Book	
	<i>l</i>	<i>u</i>	<i>l</i>	<i>u</i>	<i>l</i>	<i>u</i>	<i>l</i>	<i>u</i>
A vs. B	-16144	22245	-1346	1117.5	-3689.6	315.26	-13797	-3978.9
A vs. C	64383	119560	3935.6	7476.6	8216	13972	-2728.8	11383
B vs. C	59184	118660	3912	7728.6	9679	15884	5609.6	20820

Table 4: Analysis of the duration spent on resource type by grade-based student cohorts for the edX course. See Table 2's caption for the explanation on how to read the table.

MoocDB Resources

Publications: groups.csail.mit.edu/EVO-DesignOpt/groupWebSite/index.php?n=Site.Publications

- ...MOOCdb: Developing Data Standards for MOOC Data Science

Kalyan Veeramachaneni, Franck Deroncourt, Colin Taylor, Zachary A. Pardos, Una-May O'Reilly,
MOOCShop at Artificial Intelligence in Education, 2013

Other Resources

- ...Wiki site documenting data model
 - ...will be perpetually updated
 - ...<http://moocdb.csail.mit.edu/wiki>.
- ...Web-based software repository (not yet public)
 - ...[https://github.com/ organizations/MOOCdb](https://github.com/organizations/MOOCdb)

MoocViz Resources

- **MoocVIZ: A Large Scale, Open Access, Collaborative Analytics Platform for MOOCs**
 - DERNONCOURT, O'REILLY, VEERAMACHANENI, S. WU, C. DO, S. HALAWA
 - NIPS 2013 Workshop on Data Directed Education
- **Wiki site documenting MoocDB data model**
 - will be perpetually updated
 - <http://moocdb.csail.mit.edu/wiki>.
- **Web-based software repository (not yet public)**
 - [https://github.com/ organizations/MOOCdb](https://github.com/organizations/MOOCdb)
 - R, Python, Matlab
- **Web server**
 - Local and community versions
 - Visualizations are described in html

Future MoocDB and MOOCViz work

- **MoocDB Scaling Up**
 - Move from grass roots, bottom one step up to institutions
 - » Legacy course data
 - » EDX, Coursera and Kahn Academy joining
- **MoocVIZ: Visualization building**
- **Leveraging MoocDB for ALFA research**
 - Crowd sourcing
 - Tiger team research with fielding
 - Problem response behavior
 - Understanding MOOC attrition
 - Studying social interactions

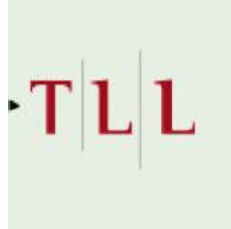
Acknowledgments

ALFA Mooc Data Science Team

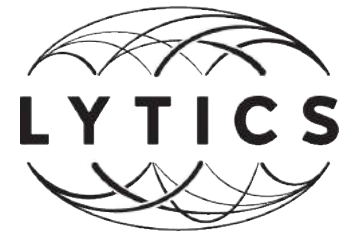
- ...Kalyan Veeramachaneni (Lead)
- ...Franck Dernoncourt
- ...Elaine Han
- ...Colin Taylor
- ...Sherwin Wu
- ...Kristin Asmus
- ...John O'Sullivan
- ...Will Grathwohl
- ...Josep Mingot

Partners

Lori Breslow
Jennifer Deboer
Glenda Stump



Sherif Halawa
Andreas Paepcke
Rene Kizilcec
Emily Schneider



Piotr Mitros
James Tauber



Chuong Do



Sponsors

- ...Mooc Research Initiative
- ...Quanta Research

