

# Developing standards and systems for MOOC data science

Kalyan Veeramachaneni, Franck Dernoncourt  
Colin Taylor, Zach Pardos, Una-May O'Reilly

Any Scale learning for All  
CSAIL, MIT



# Overview

- Background and motivation
  - What did the data looked like?
  - Where are the bottlenecks?
- Data science @ scale
  - Organize- MOOCdb
  - Create multiple views- MOOC En Images
  - Provide APIs- MOOCdb Access
- Why standardize?

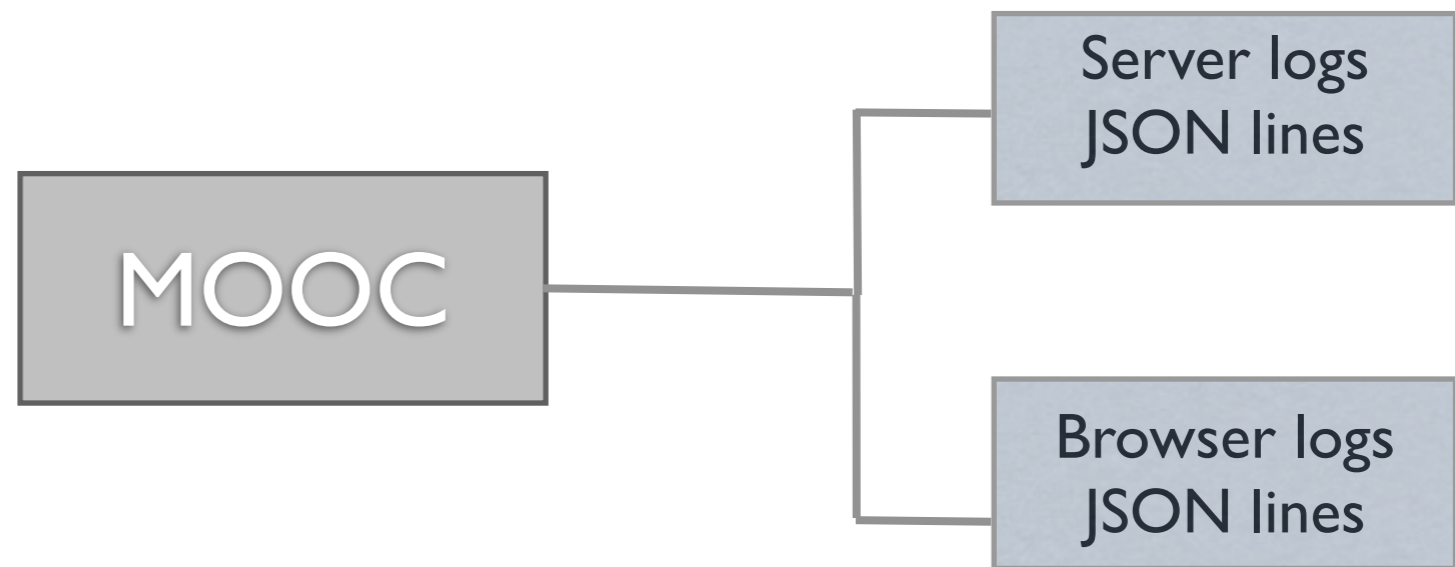


# What did the data look like?

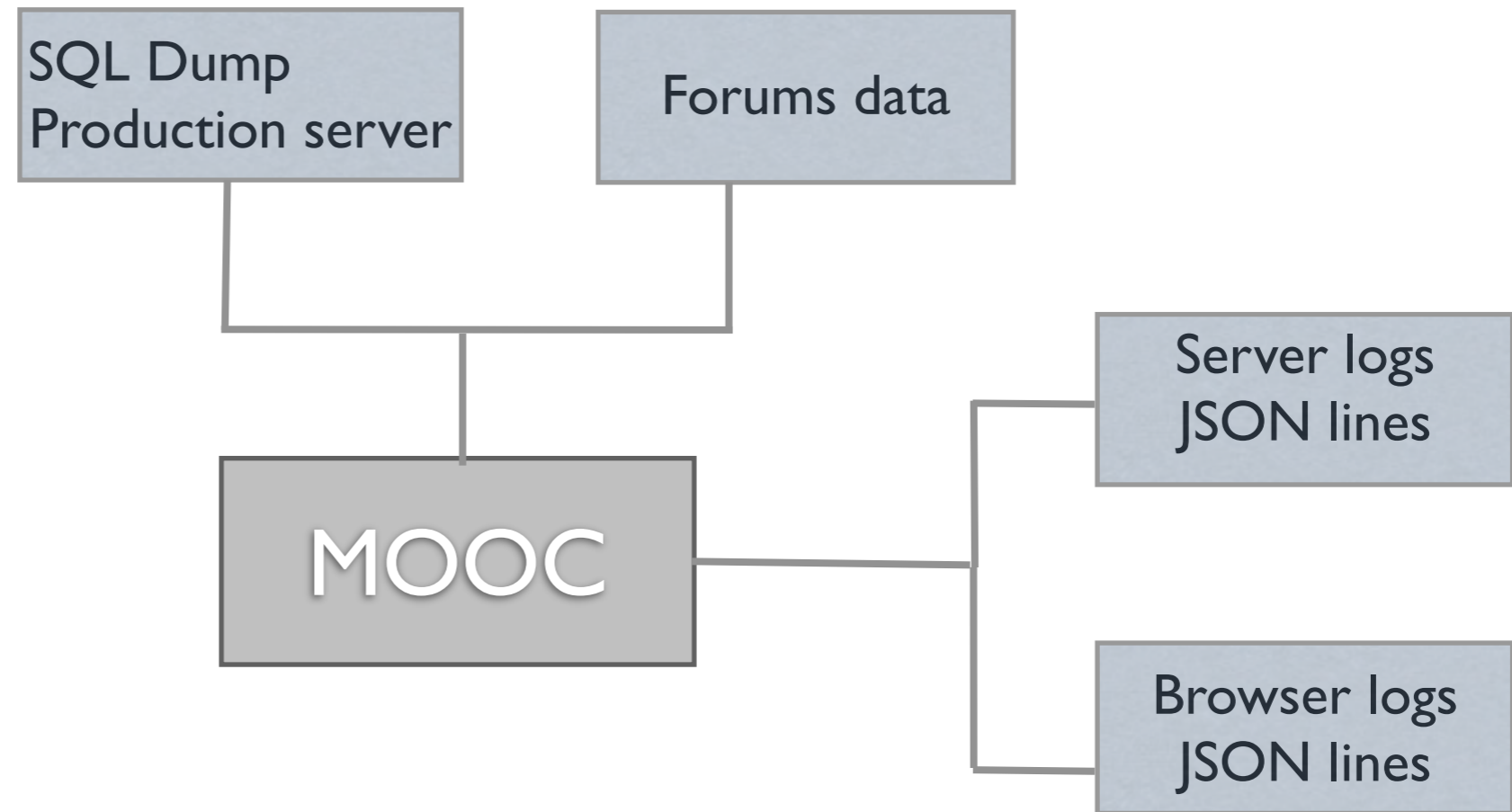
MOOC



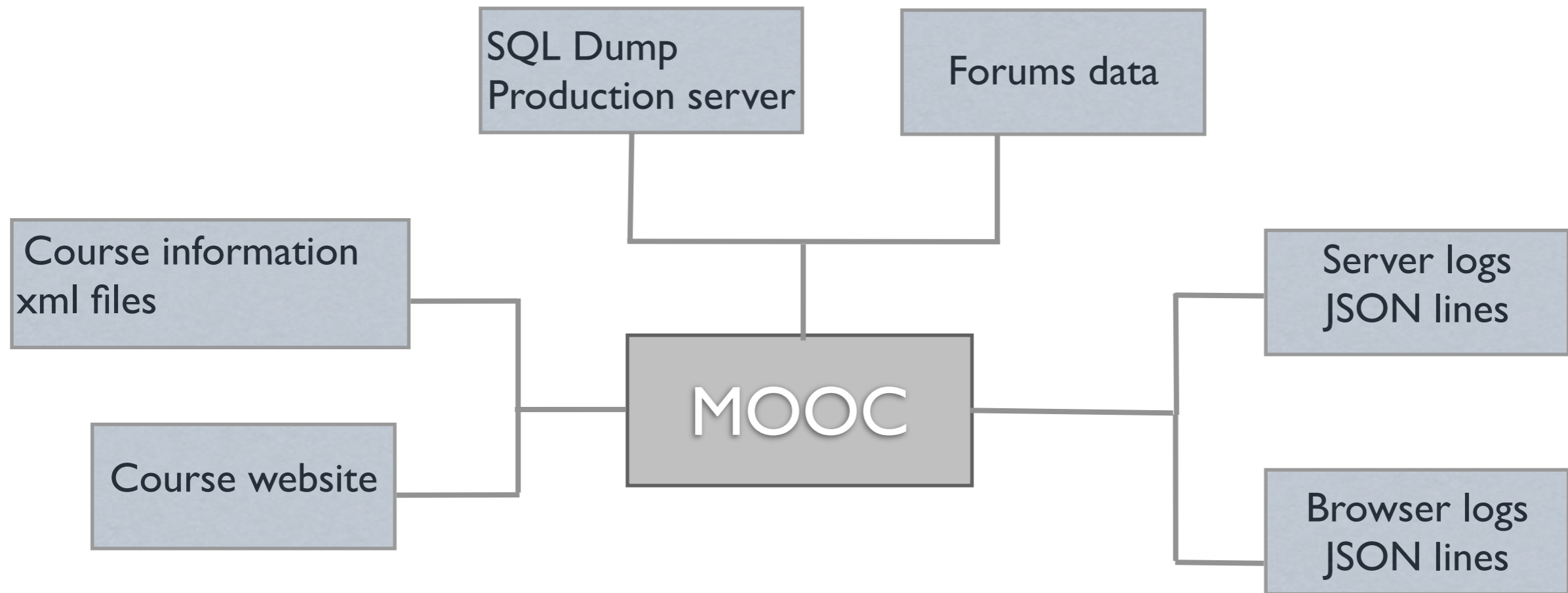
# What did the data look like?



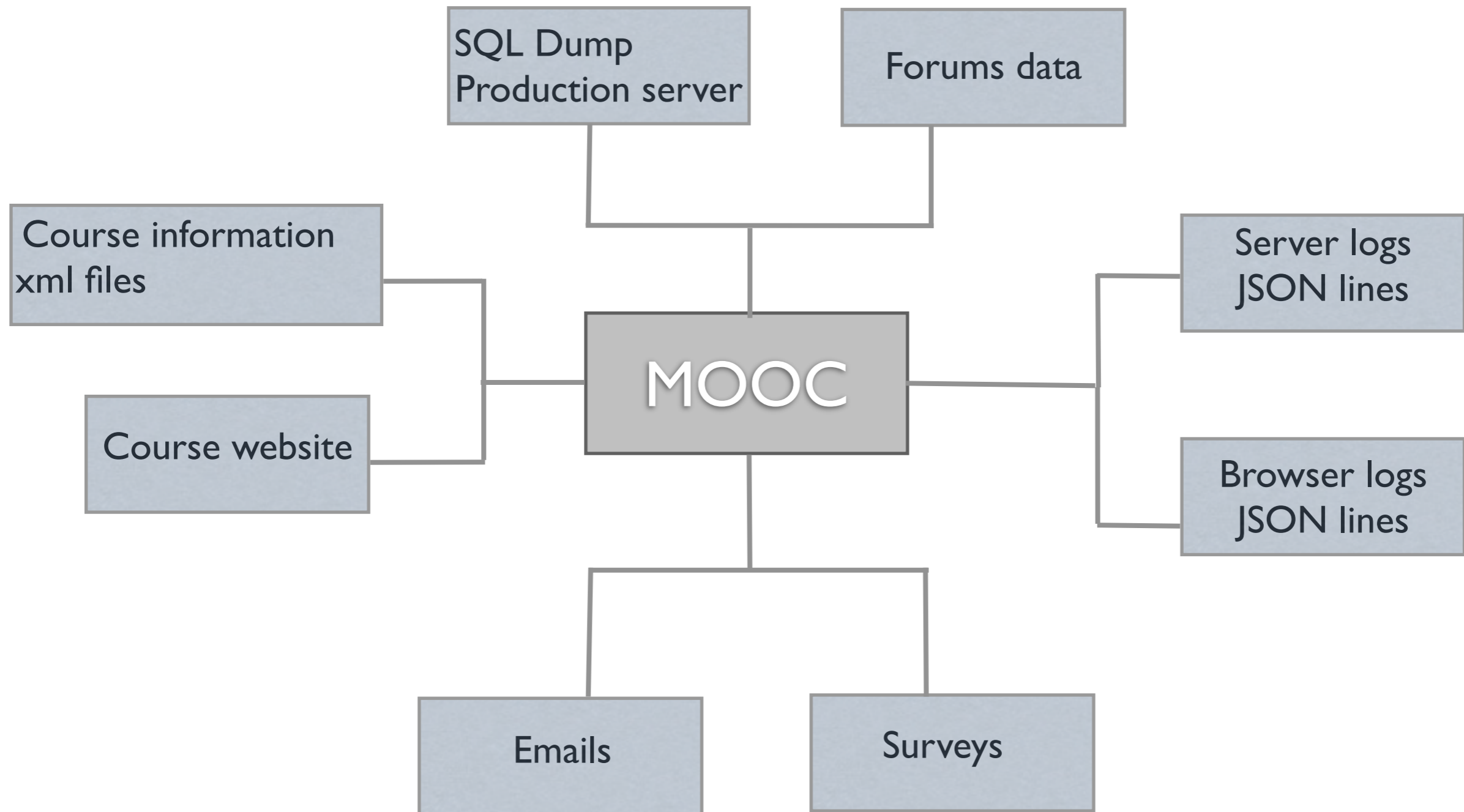
# What did the data look like?



# What did the data look like?

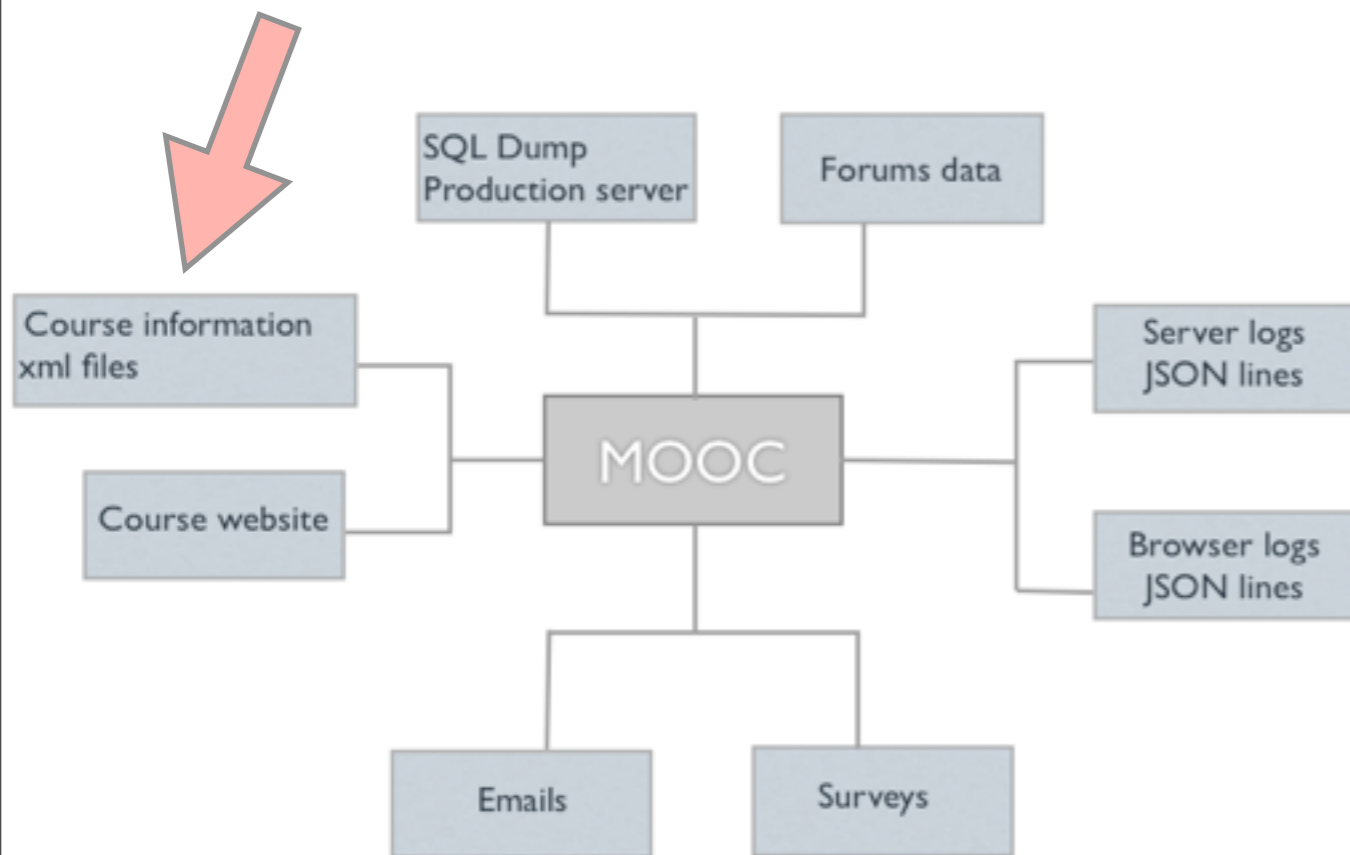


# What did the data look like?



# Example Research Question

RQ1: How does amount of time spent on the video correlate to performance on the homework?



What is the time period between two consecutive homeworks?

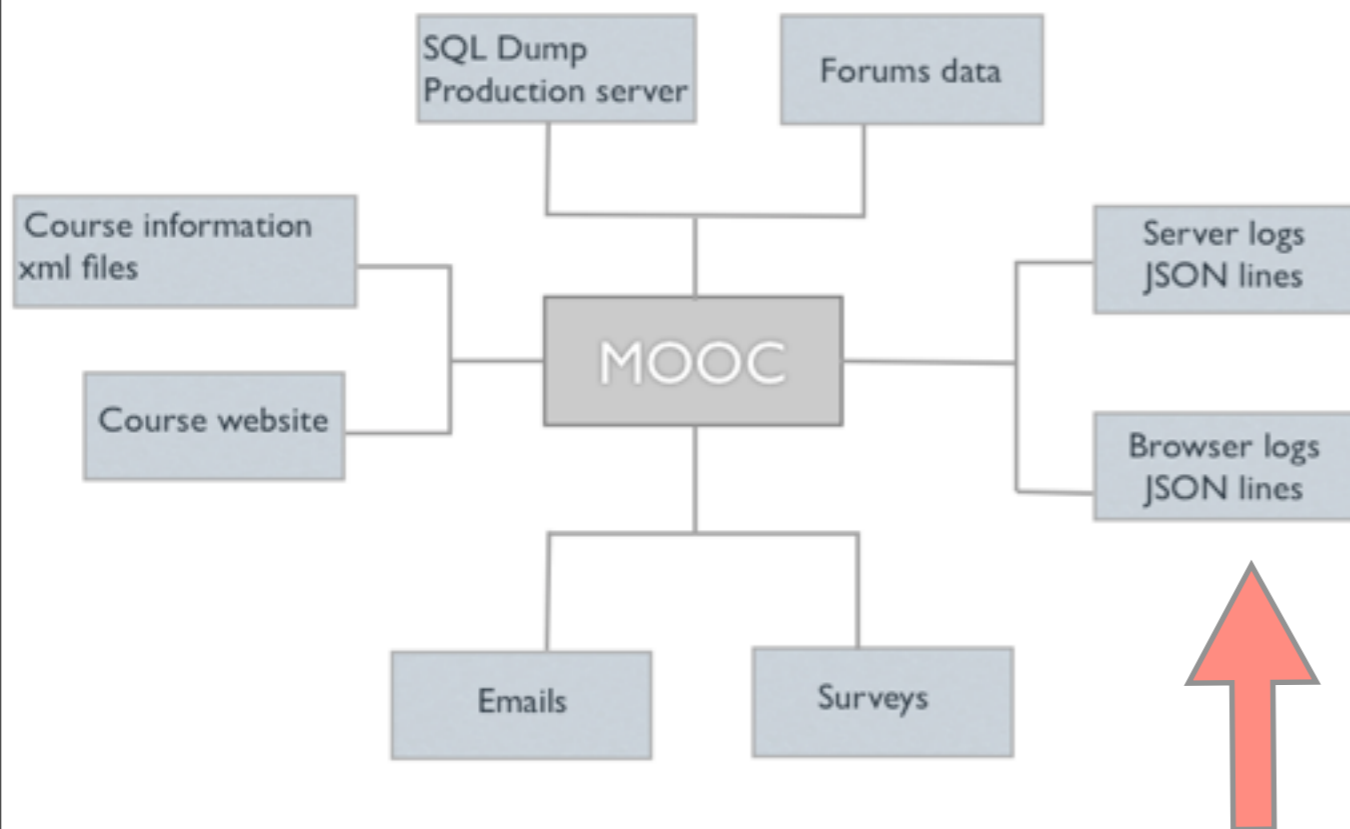
Identify what are the video urls that correspond to these homeworks?





# Example Research Question

RQ1: How does amount of time spent on the video correlate to performance on the homework?



What is the time period between two consecutive homeworks?

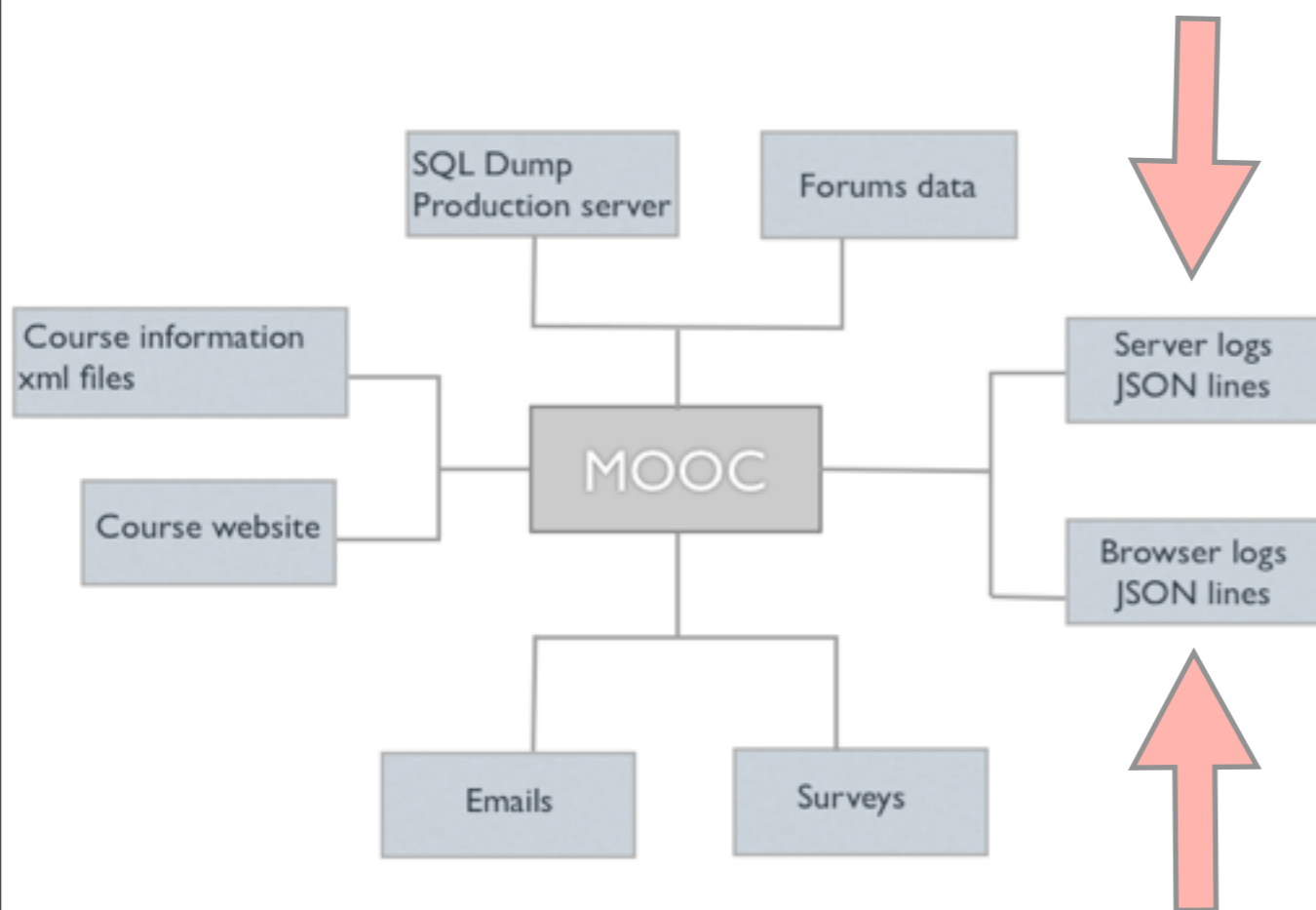
Identify what are the video urls that correspond to these homeworks?

Calculate the time spent on these urls by the user?



# Example Research Question

RQ1: How does amount of time spent on the video correlate to performance on the homework?



What is the time period between two consecutive homeworks?

Identify what are the video urls that correspond to these homeworks?

Calculate the time spent on these urls by the user?

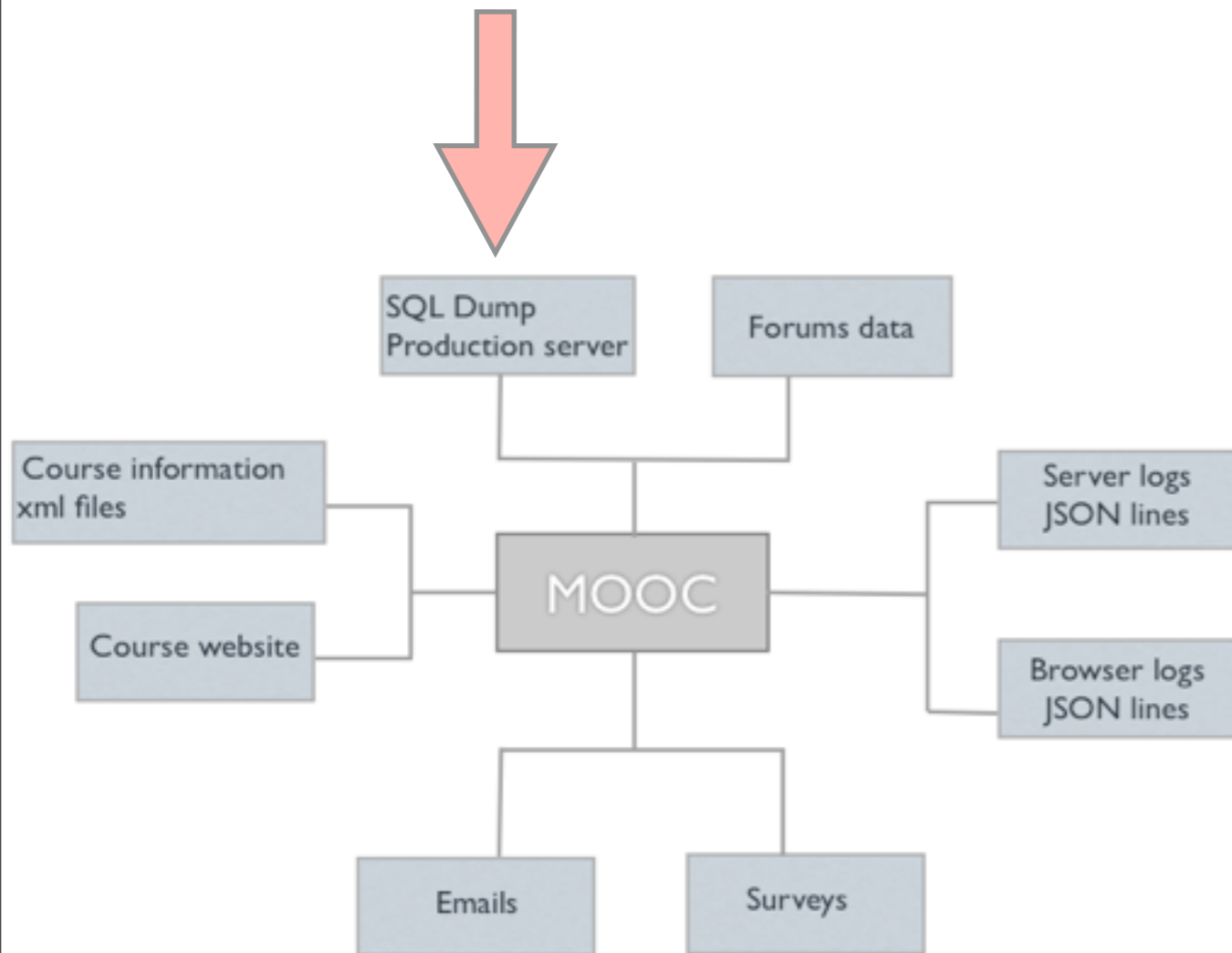
What answers did the user submit during this period?

What are the correct answers for those problems?



# Example Research Question

RQ1: How does amount of time spent on the video correlate to performance on the homework?



What is the time period between two consecutive homeworks?

Identify what are the video urls that correspond to these homeworks?

Calculate the time spent on these urls by the user?

What answers did the user submit during this period?

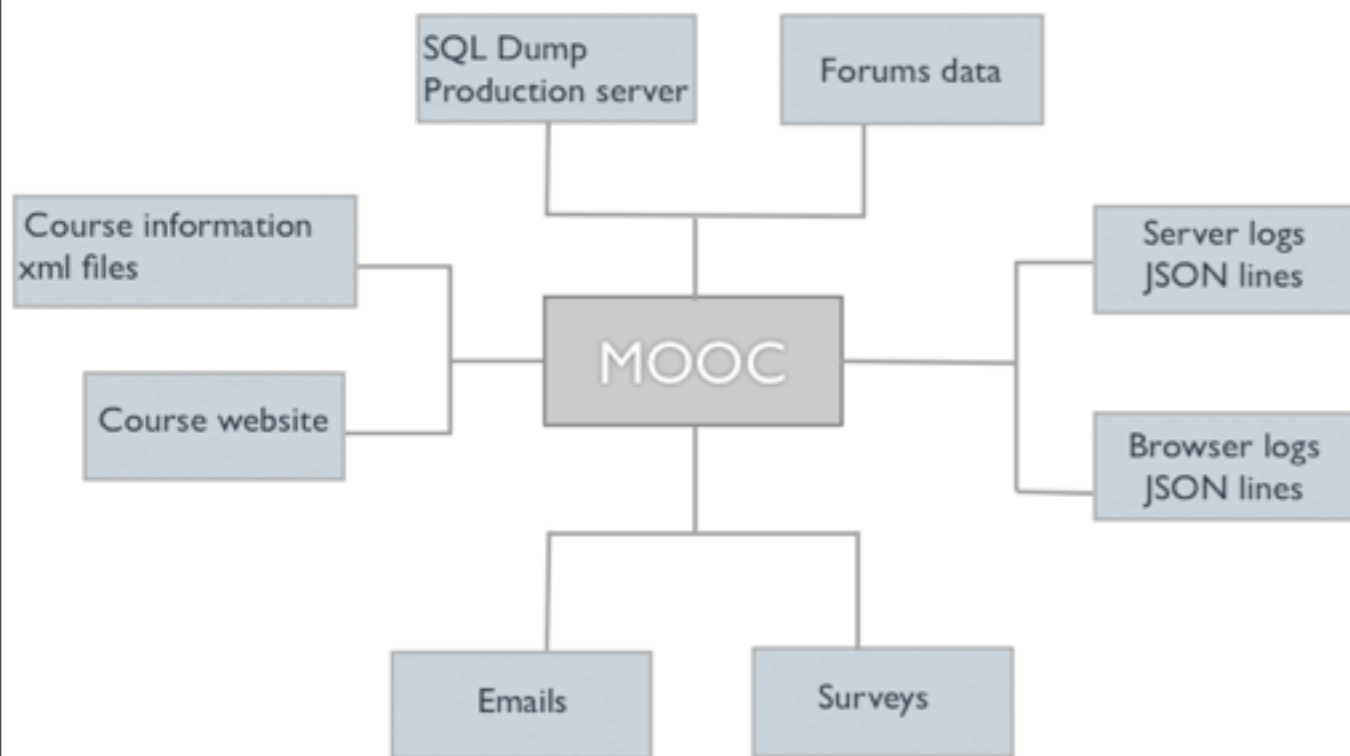
What are the correct answers for those problems?

How many did the student get right?



# Example Research Question

RQ1: How does amount of time spent on the video correlate to performance on the homework?



What is the time period between two consecutive homeworks?

Identify what are the video urls that correspond to these homeworks?

Calculate the time spent on these urls by the user?

What answers did the user submit during this period?

What are the correct answers for those problems?

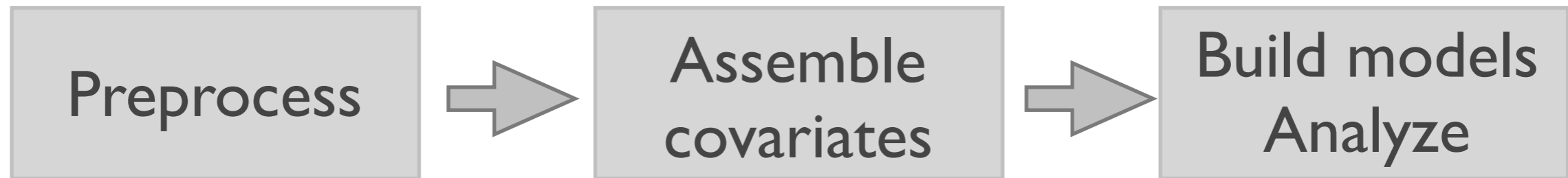
How many did the student get right?

Run `corr` function in MATLAB



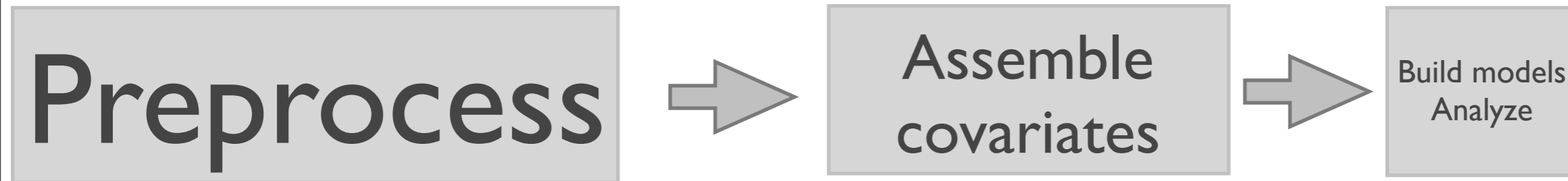
# Where are the bottlenecks?

A typical process

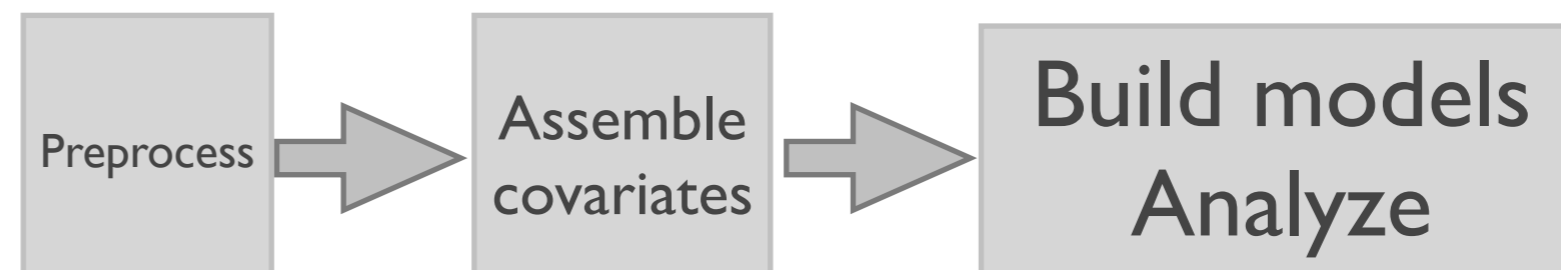


# Where are the bottlenecks?

What is now ?



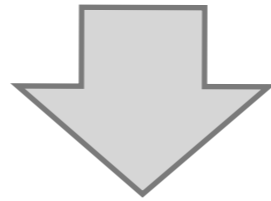
What we desire ?



sizes represent the time spent in doing the activity



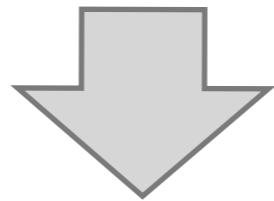
How about we organize the data using a concise schema that we will always adhere to ?



Then all the scripts we write can be **reused**



How about we organize the data using a concise schema that we will always adhere to ?



Then all the scripts we write can be **reused**

**Challenge:** Generalizable, loss-less





# Steps to Data science @ scale

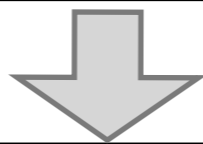
Step 1: Develop a generalizable, loss-less schema for the data



Step 2: Define and design APIs to extract multiple views of data



Step 3: Design user friendly APIs



Step 4: Design a platform where users can contribute and share the scripts



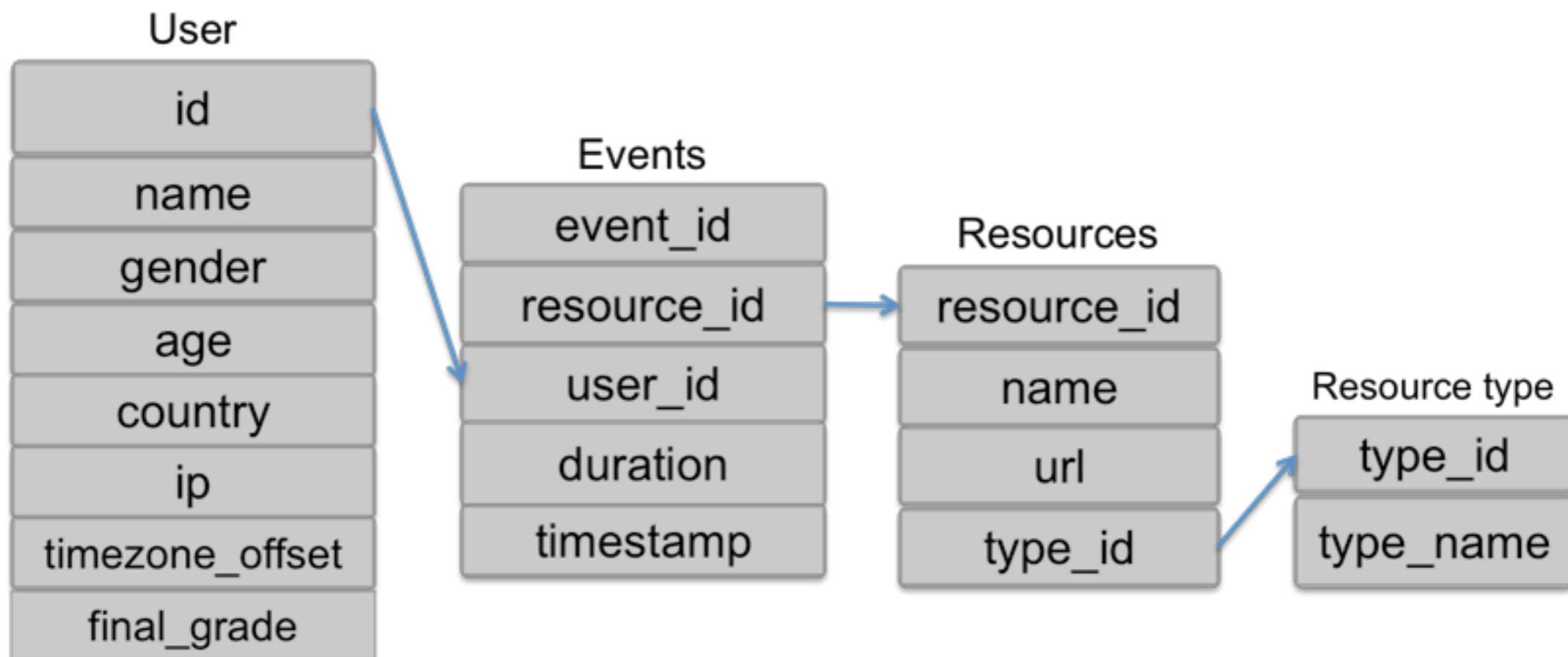
## Step 1: Develop a generalizable, loss-less schema

- Students interact with the system in the following ways
  - Observe
  - Submit
  - Collaborate
  - Give feedback



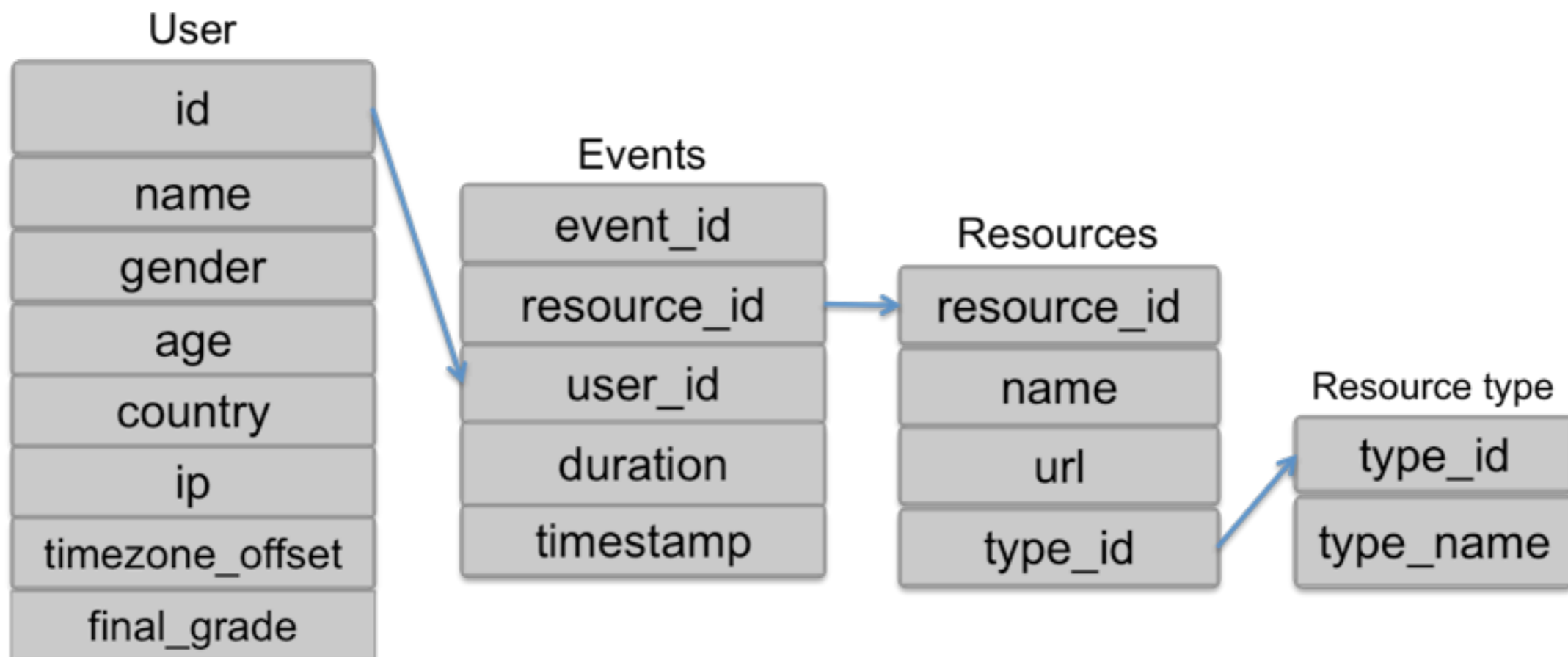
# Step 1: Develop a generalizable, loss-less schema

## The observing mode



# Step 1: Develop a generalizable, loss-less schema

## The observing mode

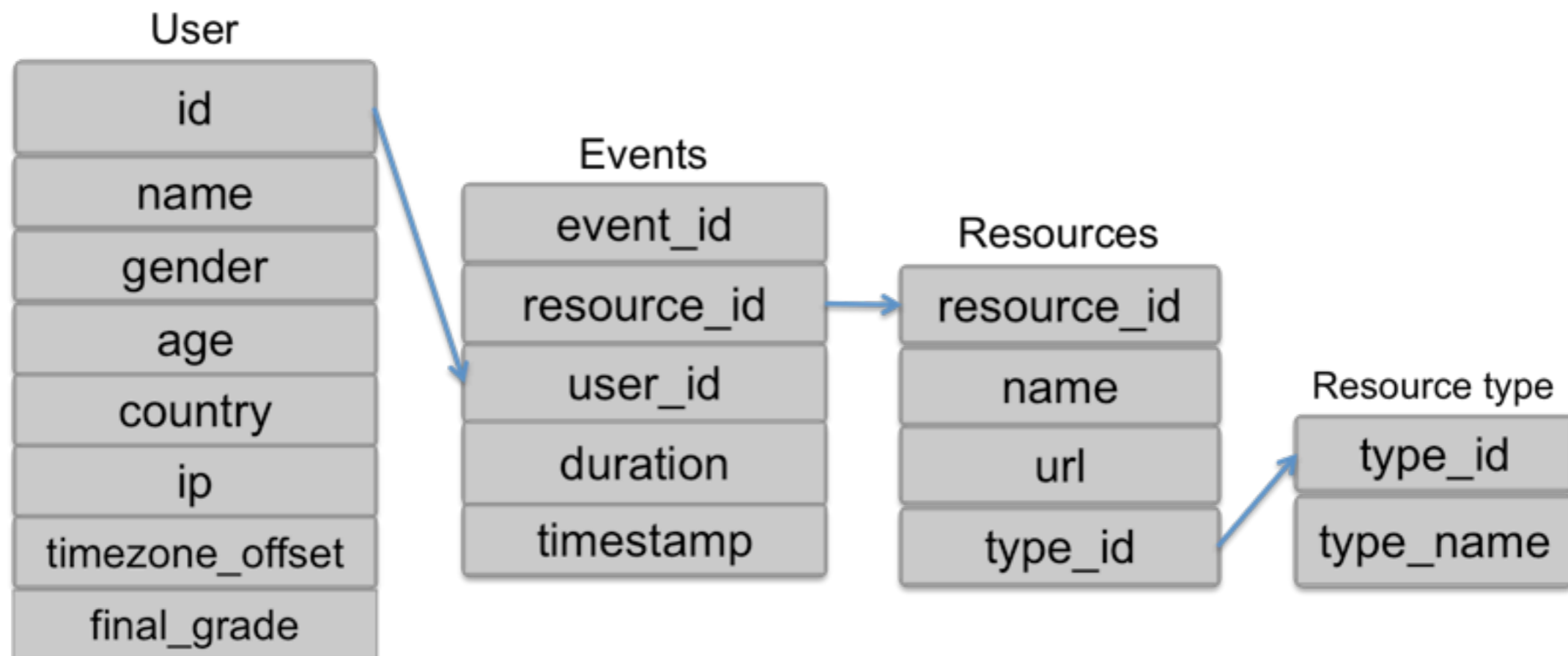


Is this generalizable?



# Step 1: Develop a generalizable, loss-less schema

## The observing mode

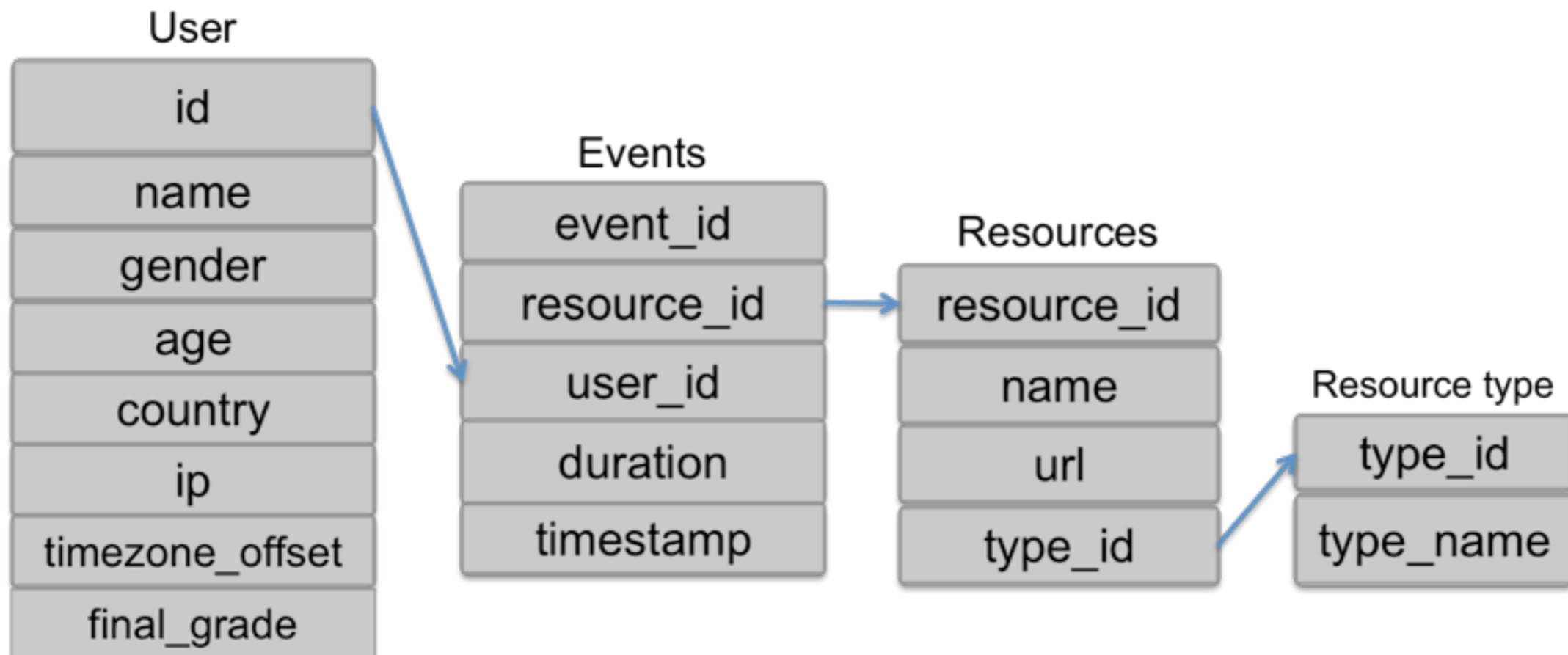


Is this generalizable? **yes**



# Step 1: Develop a generalizable, loss-less schema

## The observing mode



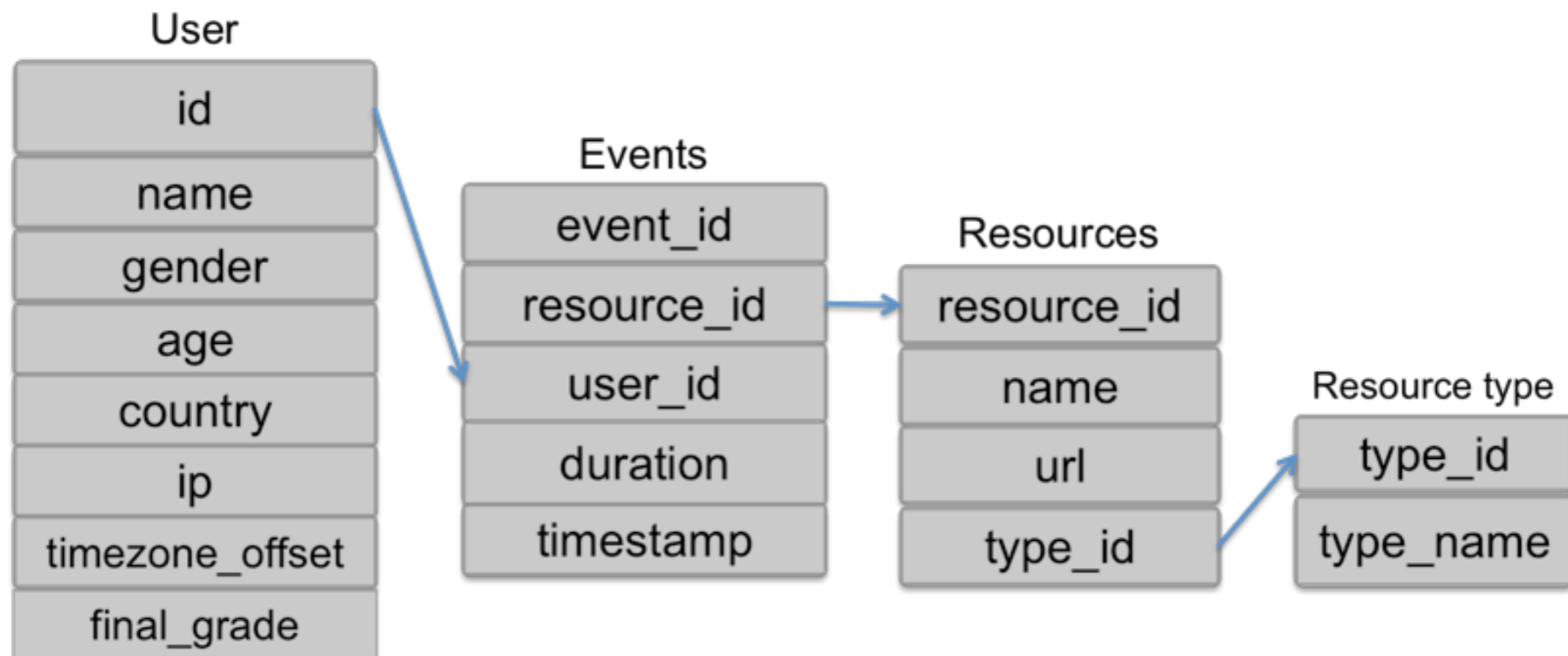
Is this generalizable? **yes**

Is this loss -less?



# Step 1: Develop a generalizable, loss-less schema

## The observing mode



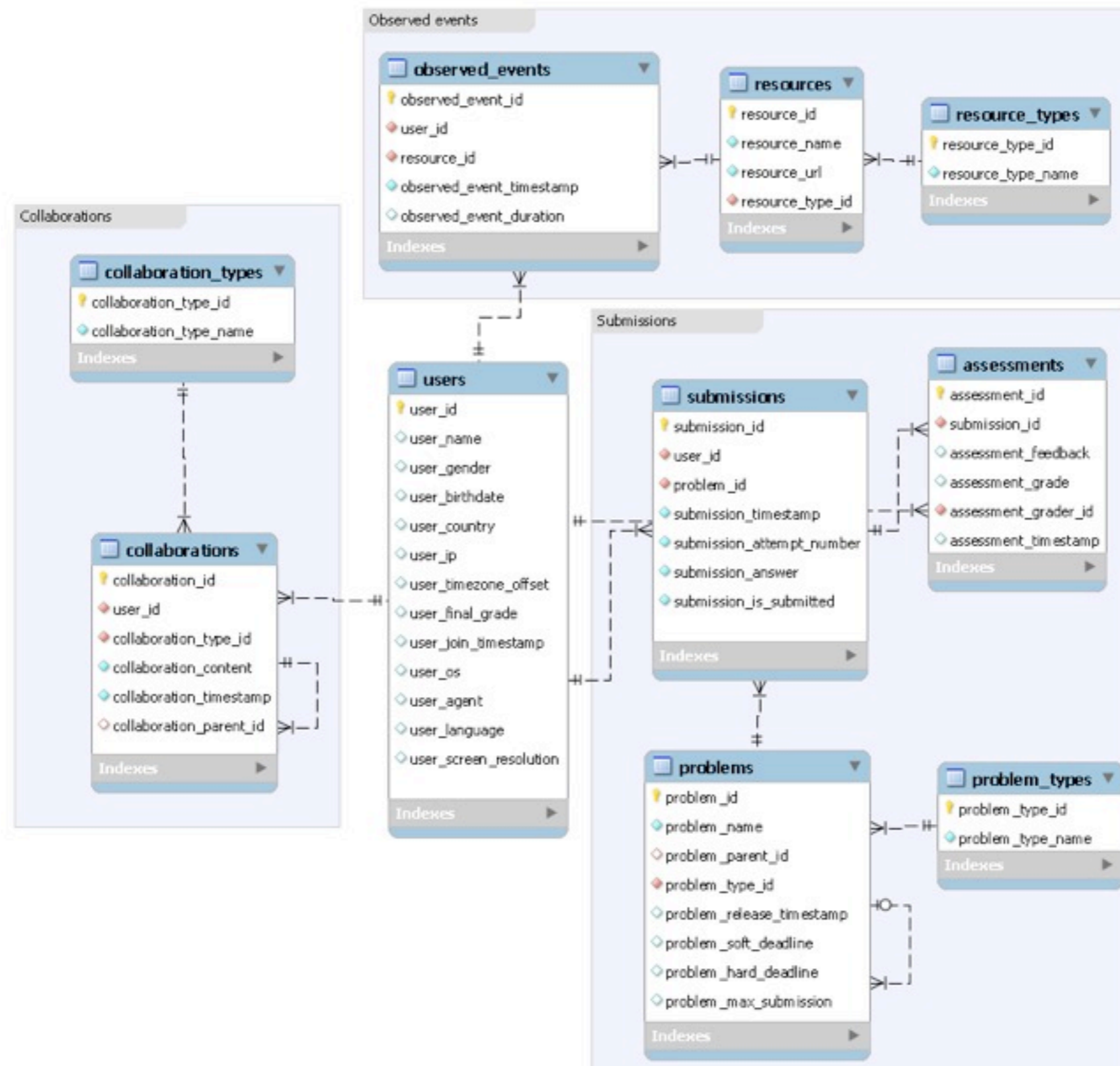
Is this generalizable? **yes**

Is this loss -less? **yes**



# Step I: Develop a generalizable, loss-less schema

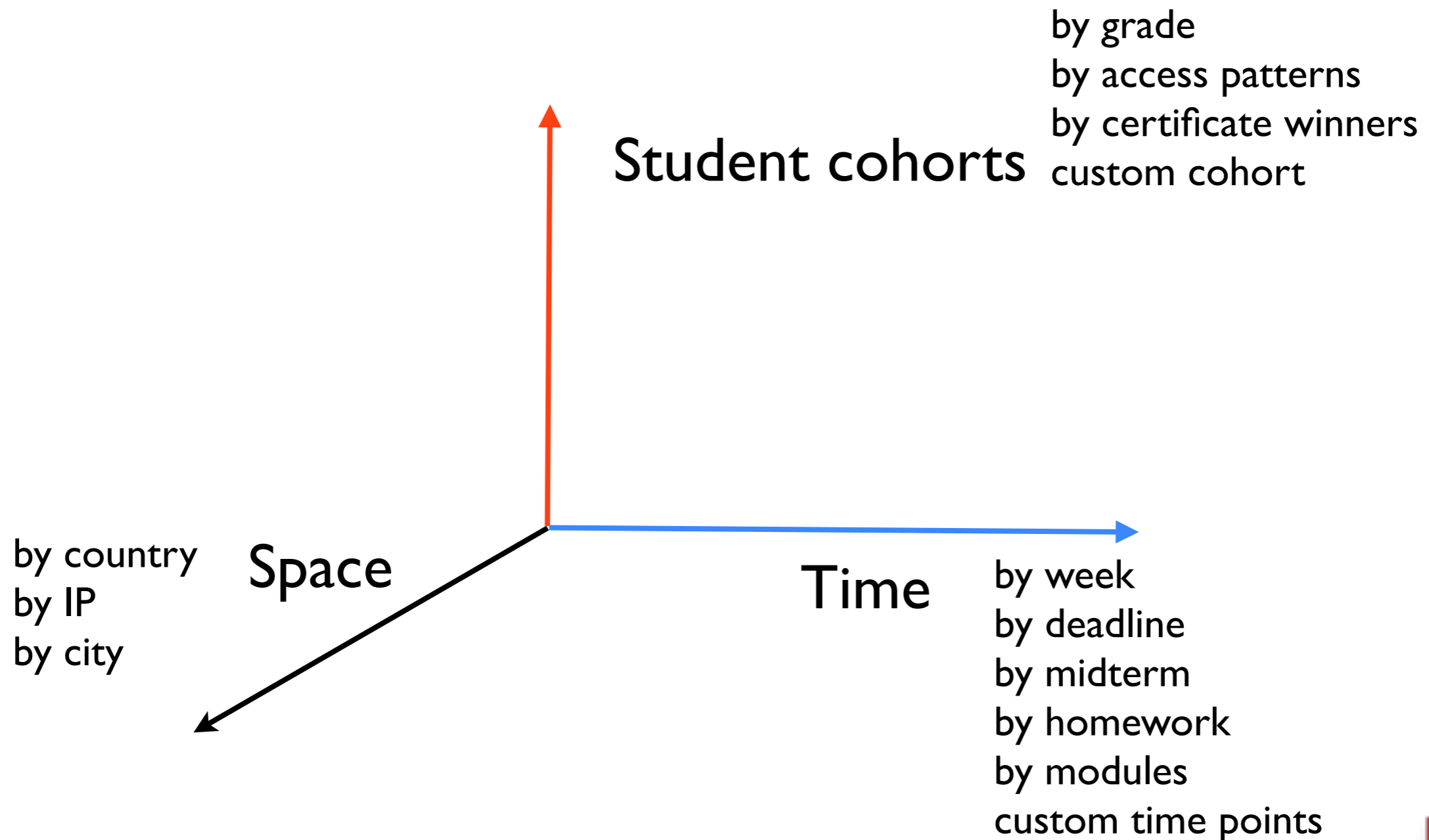
## MOOCdb





# Step 2: Define and design APIs to extract multiple views of data

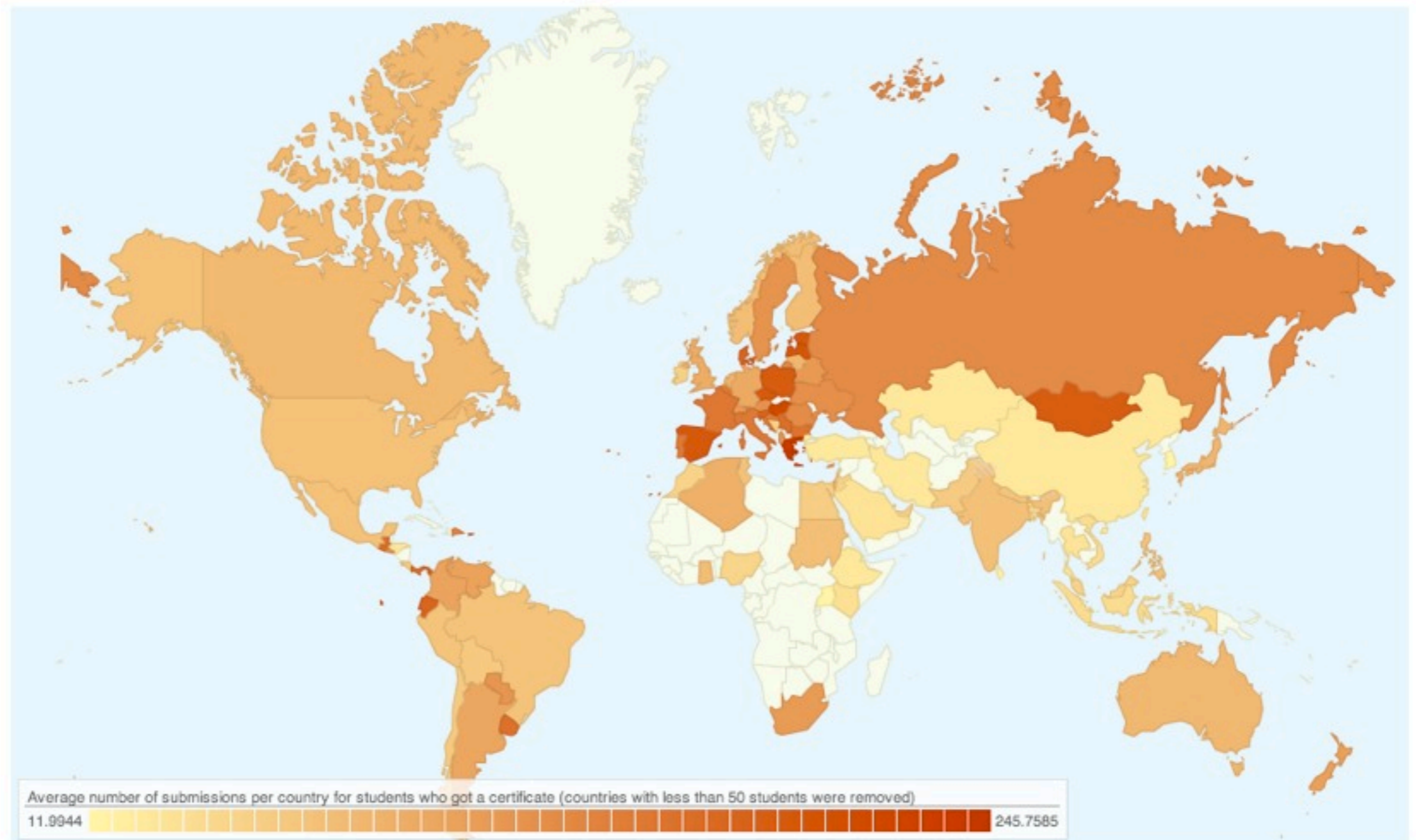
## MOOC En Images



# Step 2: Define and design APIs to extract multiple views of data

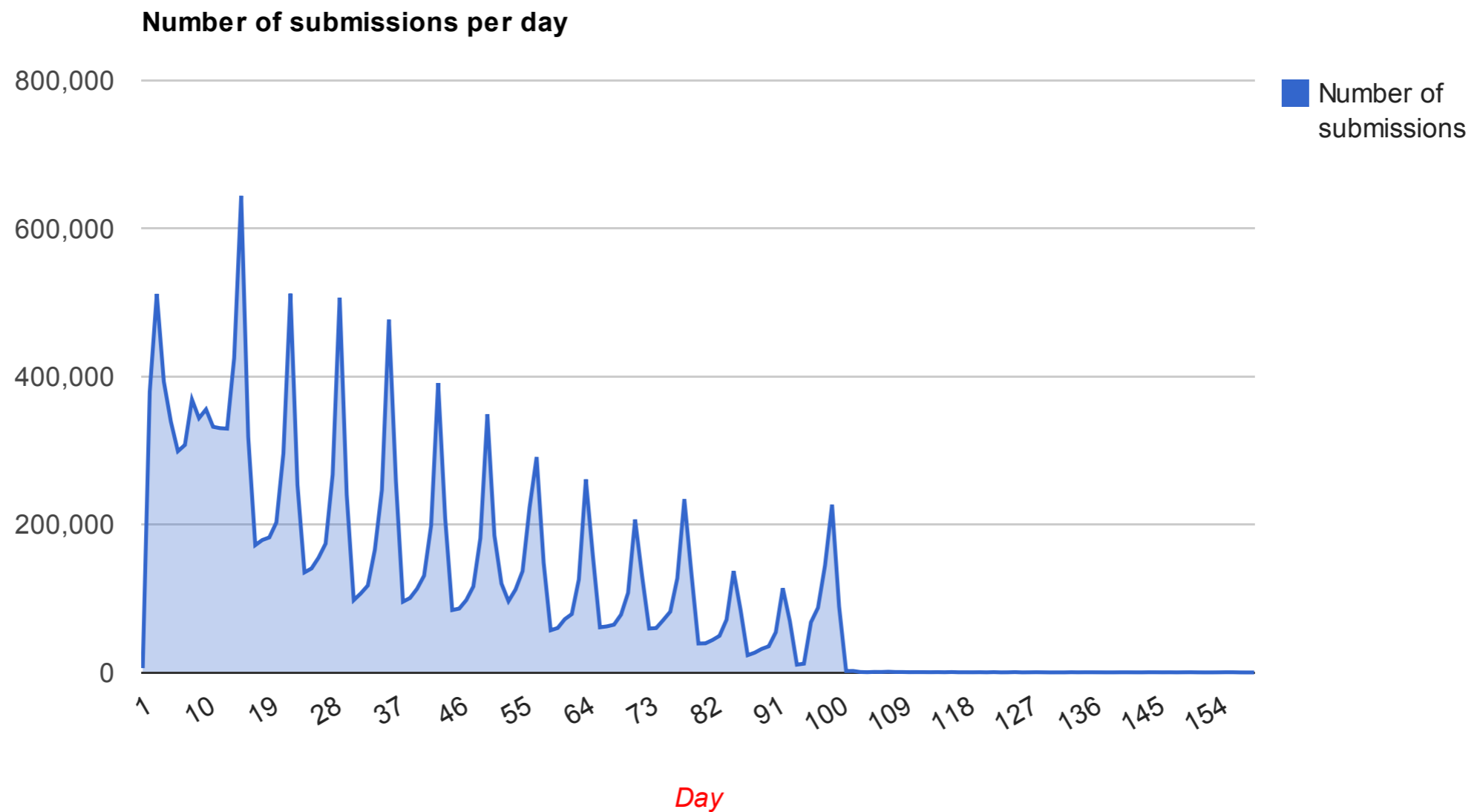
## MOOC En Images

Average number of submissions per country for students who got a certificate (countries with less than 50 students were removed)



# Step 2: Define and design APIs to extract multiple views of data

## MOOC En Images



# Step 3: Design user friendly APIs

## MOOCdb - Access

### MATLAB

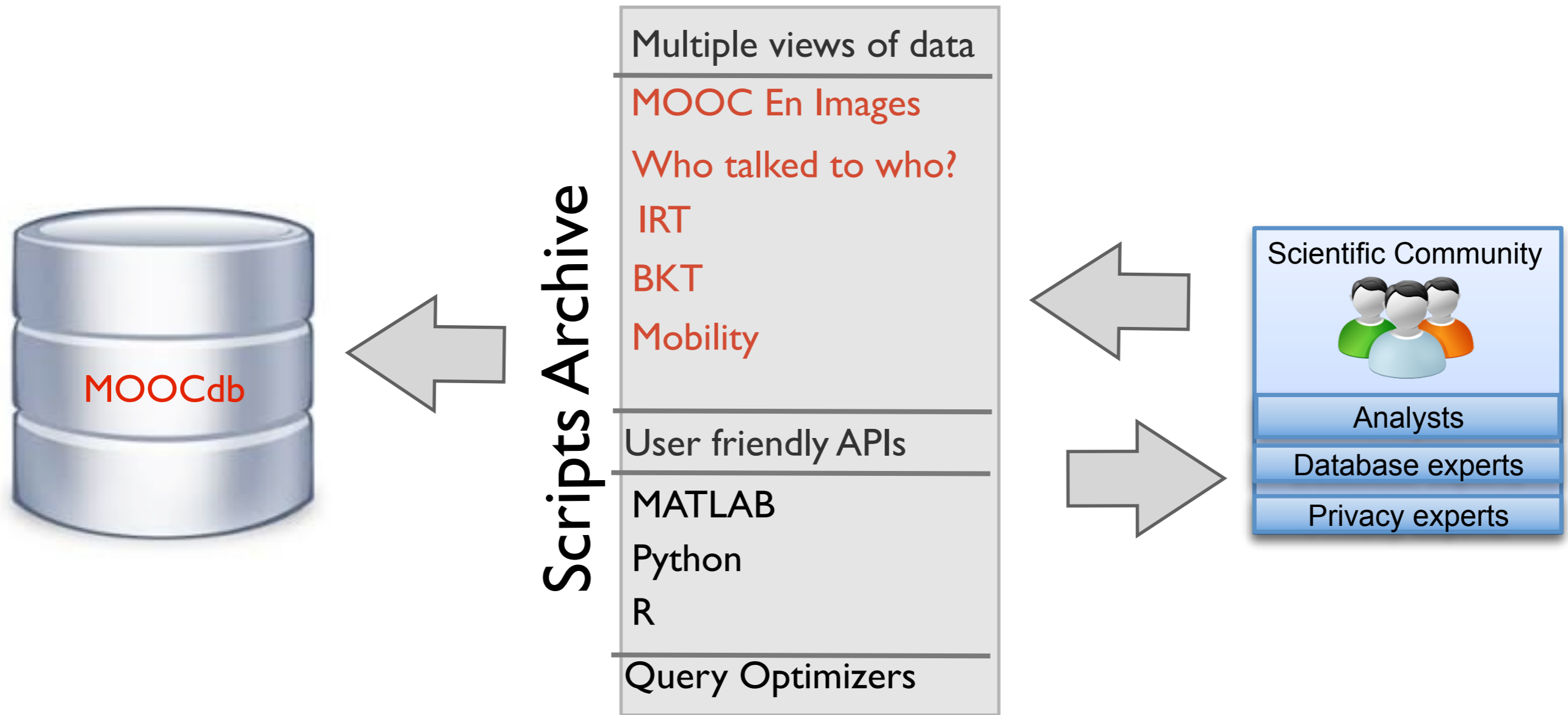
```
42 %% Running queries
43 % http://www.mathworks.com/help/database/run-sql-query.html
44
45 sql = [' SELECT observed_events.observed_event_duration ' ...
46        ' FROM moocdb.observed_events AS observed_events ' ...
47        ' WHERE observed_events.observed_event_duration < 200' ...
48        ' LIMIT 100000; '];
49
50 cursor = exec(connection,sql);
51 a = fetch(cursor);
52 data = cell2mat(a.Data);
53
54 boxplot(data)
55 title('Distribution of the duration of observed events')
56 ylabel('Duration (in seconds)')
57 print('-dpng','-r300', ['observed_events_duration_boxplot'])
58 saveas(figure(1), 'observed_events_duration_boxplot', 'fig')
59
60 plot(sort(data))
61 title('Distribution of the duration of observed events')
62 ylabel('Duration (in seconds)')
63 print('-dpng','-r300', ['observed_events_duration_plot'])
64 saveas(figure(1), 'observed_events_duration_plot', 'fig')
65
```

### Python

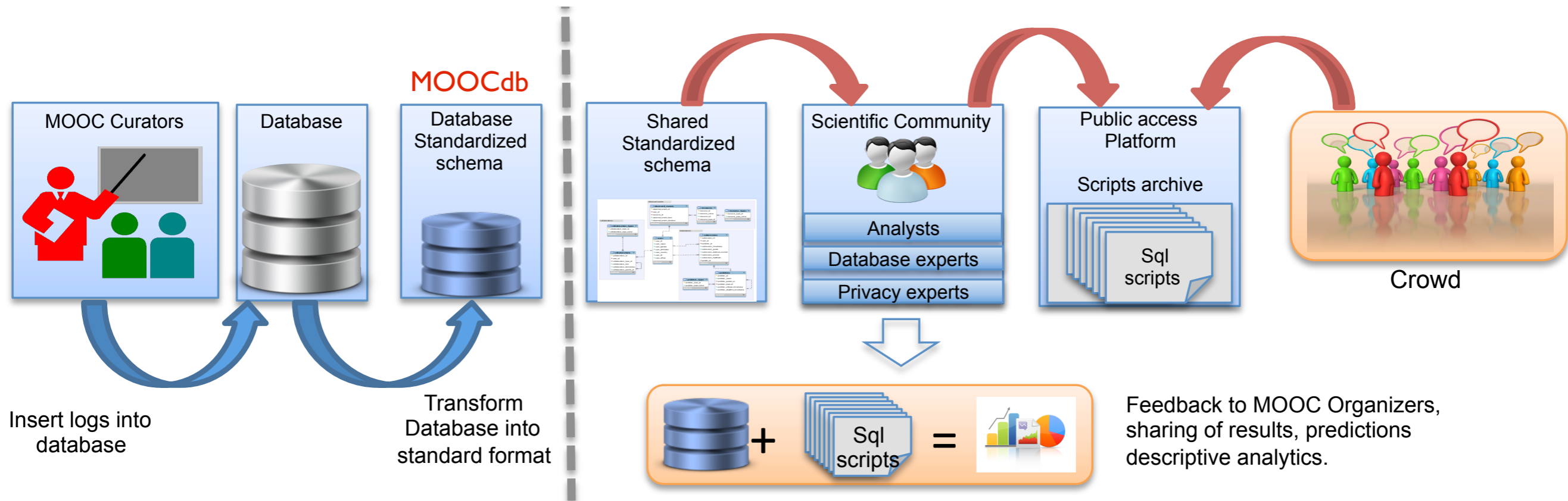
```
9 import google_charts_wrapper
10
11 def main():
12     """
13     This function plots the number of new registered students every day from 2012-02-13
14     """
15
16     sql = """
17     -- Takes 1 second to execute
18     SELECT DATEDIFF(users.user_joined_timestamp, '2012-02-13 00:00:01') AS 'Day',
19            COUNT(*) AS 'Number of new registered students'
20     FROM moocdb.users AS users
21     GROUP BY 'Day'
22     HAVING 'Day' >= 0
23     ORDER BY 'Day' ASC
24     ;
25     """
26
27     options = google_charts_wrapper.options()
28     options.set_data(google_charts_wrapper.get_data(sql))
29     options.set_chart_type("area_chart")
30     options.set_chart_title("New registered students every day from February 13, 2012")
31     options.set_height(500)
32     options.set_width(900)
33     options.set_page_title("New registered students every day from February 13, 2012")
34     options.set_h_axis("{title: 'Day #', titleTextStyle: {color: 'blue'}}")
35     options.set_v_axis("{title: 'New registered students', titleTextStyle: {color: 'blue'}}")
36     options.set_output_file("./output/users_join_date.html")
37     print options.get_data()
38     google_charts_wrapper.generate_html(options)
39
40 if __name__ == "__main__":
41     main()
42
```



# Step 4: Design a platform where community can share scripts



# Leading to standardization



# Benefits of standardization

- Reduction in time to analyze (TTA)
- Publicly available scripts to create data views
- Unified representation to DB/Privacy experts
- Crowd sourcing of feature engineering
- Crowd sourcing of data analytics
- Elimination of entry barriers



# Thank you!

Look for **MOOC En Images**, updates on **MOOCdb** and  
**open release** of all the tools.

