

**Abstract**—In this report we investigate ways to automatically analyze ICLS and CSCL papers to gain an understanding of the community formed by conference participants and of the research being conducted within the community. After describing how to extract data from the base dataset we derive insights into the community using network analysis, natural language processing and text mining. We discovered patterns of collaboration between participants, their migration and the global spread of the community. We extracted keywords and found clusters of referenced papers and documents within the conferences.

## I. INTRODUCTION

The starting point of the project are the proceedings of the International Conference of the Learning Sciences (ICLS), respectively the International Conference on Computer-Supported Collaborative Learning (CSCL). The two conferences are organized by the International Society of the Learning Sciences (ISLS) and held in biennial alternation with each other. At the time of writing the dataset contains 4 years of the proceedings, from 2015 to 2018. While ICLS covers the entire field of the learning sciences, CSCL focuses on learning through collaboration with the help of communication technologies [1]. To understand the community formed by the participants of the conferences as well as to understand the research being conducted within the community, we explored methods to bring insight into the following aspects:

- The institutions and countries dominant in the conferences
- The differences between ICLS and CSCL with respect to citations, research conducted and contributors
- Collaboration network patterns with respect to countries, institutions and authors
- First insight into changes over the years with respect to contributors and collaboration patterns
- The migration of participants across institutions
- Popularity of keywords over the years
- Trends in regards to cited papers and the evolution of their prevalence from 1977 until 2018
- Clusters of the cited documents
- Clusters of the papers written within the conferences

To this means we first present an extensive pipeline to process the textual data and extract features relevant to bibliographic analysis from the textual data. After this we perform a first bibliographic analysis based on the extracted data.

## II. DATA PROCESSING

The provided raw dataset consists of papers in *pdf* format, as well as associated metadata in *xml* format. To extract information from the dataset we first converted the data to a format that is better adapted to the task of data analysis and more compatible with a variety of external libraries. We stored all the resulting data in *csv* format, as *csv* is a human readable format and can easily be processed by all the data analysis tools we used.

### A. Processing the Metadata

We extracted the metadata from the *xml* files using the *lxmltree* python library. We associated the contents of each *xml* tag to a column in a *csv* file. As the *xml* files do not exactly follow the *xml* standard, we handled errors by iteratively parsing the *xml* tree. For some files, the section containing the keywords is malformed and can not be extracted by the parser. This was remedied using regular expressions to match the section containing the keywords. To each file we associated the name of the source file and the order in which an author name occurs in the *xml* file. Each author is associated to one line in the *xml* file. From the associated citation string we added the first part of the citation containing the author names, year, and title. This is later used to identify the paper when it is referenced in another paper. Then we added a shortened version of author names (i.e. West, R.) by considering the author order. Note that it happens that the short name extracted from the citation and the corresponding long name (Robert West) in the metadata do not match up. Thus the author order listed in the citation string does not always correspond to the author order listed in the metadata. After extraction, we cleaned the data by unifying names, as slight spelling differences exist within the same author name (i.e. Robert West, Robert A. West). When unifying the names we needed to make sure not to overmerge by labeling two different authors as having the same name. First, we used a strict condition to catch small differences in names, such as an extra middle name. This is done by splitting up names into substrings based on whitespaces and commas. For two names in the dataset we checked whether the intersection of their sets of substrings contained at least 2 names of length at least 3. If the condition was met we considered these two names to be identical, and added the mapping from one name to the other to a dictionary. Once all tuples of names (s.t. the first name appears before the second one in order) were checked and added to the dictionary we unified the

names. Using this method we reduced the number of authors in the dataset from 1951 down to 1879. Further, to merge names with misspellings or people who use nicknames (Christian vs Chris) and to avoid overmerging we considered only names of authors that have similar collaboration patterns. To do this we built a graph based on co-authorship relations and checked names of authors that were in the same neighborhood of a node on the graph (see methodology for further details). Then, we used the *difflib* library to check the similarity of two strings and considered a name to belong to the same person if the similarity measure reached the threshold of 0.8. Again, names considered to belong to the same person were added to the mapping dictionary which was then applied to all names. This method again reduced the number of authors in the dataset from 1879 to 1851.

### B. Handling of pdf data

To extract information from the papers themselves we needed first to extract the text from the pdfs. Various parsers performing this task exist. However, parsers such as pdfminer or PyPDF2 that claim to be able to extract the text as well as metadata from papers have not been kept up to date, which leads to issues with dependencies. We found that the most reliable parser is *poppler* [2] which can be installed and then used via the command line to extract the text contained in the pdfs to *txt* format. All methods describing how to extract data from the papers worked on the extracted *txt* files. The characters contained in the *txt* files were not all encoded uniformly, which means that characters that look the same to humans may be encoded in different ways. Hence, for many applications we converted all characters to ‘Normal Form Composed’ form.

A convenient property of papers is that they are structured fairly regularly. The ISLS Author Guidelines [3] gave us guidance on how the papers in our dataset should be structured. We used this to our advantage when extracting information from the paper text. While the vast majority of papers in the dataset follow these guidelines, exceptions do exist. We thus made sure to handle deviations from the structure.

### C. Extracting the papers cited by a paper

Each valid paper contains a reference section at the end listing all the references cited in the paper. While tools are being developed to extract the reference section directly from pdfs, notably Scholarcy [4], they did not produce satisfactory results. That is, they did not split the references correctly and categorized parts of the reference inaccurately. This may occur because they are trained to extract any reference format, and thus can not make many assumptions on the structure. But in

our dataset papers following the guidelines should list their references in APA format. Hence using regular expressions and background on APA referencing gave us better results than out of the box methods.

**Regular expressions in python:** Python has two modules that implement regular expressions, *re* and *regex*. We used *regex* as –unlike the *re* module– this module allows to handle strings with non-ascii characters by allowing the use of unicode regular expressions. As many authors have names containing non-ascii characters this is an important capability.

**To get the reference section for of a paper** we used a regular expressions to find the reference header and determine the end of the section by checking for acknowledgements or appendix sections in the text file previously extracted using *poppler*. Then, to extract individual references we split the section into a list of substrings at each new line character and concatenated substrings back together if they satisfied a set of conditions. We found this method to be robust as it makes sure that all properly formatted references are separated.

**To unify substrings** we exploited the fact that APA references all have the same underlying structure: *Lastname F.M. (when published) Title which may contain colons and other non-alphabetic chars. Where published. Who published* [5]. We thus checked that the beginning of each reference contains “Lastname F.M. (when published)”. It is relatively easy to check for this using regular expressions as author names and year have a given structure. If a string does not contain such a section, then it must be part of an other reference, which lies before in the list of reference substrings. Note that this beginning is also what is added to the metadata as identifying string. We found that in practical application it is important to first merge lines if they were clearly cut off, as papers with many authors may have the beginning on multiple lines. We could find split up lines by checking for lines starting with non-alphabetic characters or starting in a sentence. We merged these with the line above. We could identify lines by their ending as being split up, such as if they end in a comma, colon or only contain author names. We merged these lines with the lines below them. Only then could we check whether a line contained a valid citation start. If it did not, we merged it with the line above it. After this merging process we validated the returned reference split by looking at the overall length and the initials, as well as by checking that a random sample returned satisfactory splitting accuracy. We found that the above method of extracting the references produced good results for proper APA citations and on APA citations indeed outperformed Scholarcy when it came to splitting up the lines.

**To then get the subcomponents of a reference** was

Susan A. Yoon, Jessica Koehler-Yom, Emma Anderson, and Chad Evans  
 yoonsa@upenn.edu, jkoehl@gse.upenn.edu, emmaa@gse.upenn.edu, echad@sas.upenn.edu  
 University of Pennsylvania

Figure 1: First format used to reference institution

Susan A. Yoon, University of Pennsylvania, yoonsa@upenn.edu  
 Cindy Hmelo-Silver, Indiana University, chmelosi@indiana.edu

Figure 2: Second format used to reference institution

rather trivial, as we could again use the structure of APA citations. The string before the first parenthesis contains the authors names, inside the parentheses we find the year or some other reference to publication time, after which the title follows until the first dot. Then a reference to the publishing venue follows.

#### D. Extracting Institution and Location

Now we consider the section before the abstract, which contains the authors names, their email addresses and the institutions they are affiliated with. Affiliation with institution can give us a lot of insight where people are working and gives us a way estimate where they are from. To extract the exact affiliation of an author with an institution two approaches were used.

1) Approach 1: The most robust way to get a person's institution affiliation is to use their email address. Email addresses can easily be extracted from the header section using a regular expression. Then, the part of the email address containing the university domain can be extracted using another regular expression. After this we used a domain-university mapping [6] to map domains to institutions. The additional advantage of this method is that using the additional metadata associated to each university in the mapping, we also access information about the country of the university.

2) Approach 2: The header section itself contains information on the institution, but different ways to list authors are accepted, notably two ways as can be seen in Figure 1 and Figure 2. Hence, we first detected the structure of the header section before extracting the association.

We tagged the substring present in this section using a combination of two methods. Firstly, *named entity detection*, which automatically classifies strings as belonging to some pre-defined category, using the implementation provided by *spacy*[7]. This tagger categorized strings as organization or a person. Secondly, we used simple string comparison matching of the names of authors and institutions to determine if a string is an author name (person) or an institution (organization). We extracted the

data needed for this from the metadata respectively from Approach 1. To detect sections that are email addresses, we used regular expression. We assigned these strings the category "Email". Then, taking the tagged section, we detected which subsections corresponded to which author-affiliation listing type, and extracted the authors email address - institution mapping accordingly.

Approach 2 allows to match many more institutions to authors than Approach 1 as not every institution is present in the mapping and many participants use gmail and similar providers as their official email. However, Approach 2 may contain differences in spelling in the institution name, which needed to be unified. Thus we combined both approaches, and used the corresponding institutions in Approach 1 to unify the institutions obtained in Approach 2. Furthermore, we later merged names by considering the location of the institutions.

But first, as with Approach 2 institution affiliation did not allow us to automatically associate a country to the institution, we needed to get this information through other means. A first step was to use the top level domains present in the emails to associate a country of origin to the participants and their associated institutions. Then, to plot the community we augmented the dataset with geolocation information. As a further benefit, having this information we could then easily associate a country to the institution for the emails that do not have a top level domain indicating location and further merge the names of institutions based on location.

As there are a variety of ways in which the name of an institution can be written – consider EPFL, EPF Lausanne– crawling the English, French and German version of Wikipedia (not using the API, or the API wrapper, as this does not give access to the desired information) was a quick and accurate way to get the location. The advantage of sending requests over using the API is due to the website redirecting to the relevant article on the institution even when misspellings are present or multiple ways to spell the university exist. Requests were sent using *urllib3*, and *beautifulsoup* was used to parse them and get the coordinates present in the article. Regrettably not all institutions, even well known ones such as Tokyo University, have their coordinates listed on Wikipedia. To get the remaining locations we used *geopy*, which is a Python client for geocoding web services. These results however are more error prone as *geopy* will always return a result. It is thus recommended to verify the correctness of the results. To get the country we then again used *geopy*. The resulting data containing paper identifier, the authors email, its associated domain, order in which the author was cited, the institution, country, longitude and latitude were saved to a CSV file. We used the paper identifier and author order to associate this data with the metadata to link information about

authors with the papers they authored.

#### E. Keyword extraction

1) *Text preprocessing*: The pdf parser returns txt files that are encoded using a variety of conventions, thus when extracting the strings from the text files, we were met with non-ascii characters, which had to be dealt with. We preprocessed the text by applying the following steps to it:

- We used *beautiful soup* and *string* libraries to only keep sections of the text that we needed and remove special characters.
- We applied *regular expression* to remove numbers because numbers were not important for our analysis.
- We put all words to lower case
- We removed words that were connecting parts of a sentence rather than showing subjects, objects or intent (*stopwords*). For the purpose of this task these words were not important and were only going to increase the bag of words length
- In some methods we applied *lemmatization* and *stemming* to the words. In other methods we did not because we implement a mapping with *GloVe* vocabulary instead, which we will explain later on.

**Lemmatization** is a word variations to the root of the word (e.g. working, works, worked changed to work). **Stemming** is the process of reducing inflected and derived words to their word stem, base or root form—generally a written word form. As the use of the stop words to reduce the BOW (bag of words) size we also used lemmatization and stemming to be able to reduce the number of words without taking off any valuable information.

**Lemmatization** and different forms of **Stemming** were only used with **TF-IDF** while doing the Keywords extraction, otherwise we prefered keeping the words as they were because we were doing a mapping with a pretrained Stanford Word Embedding [8] and sometimes root forms were not in the dictionary of words of the pre-trained words-vectors.

For the Keywords extraction we used three methods : The first one is a trivial approach that consists of searching in the pdf itself for keywords and extract them, this method is the most accurate but we found out that for 492 papers we had no keywords included. So we deciced to use other two methods for keywords extraction, namely one based on **TF-IDF** and a second one based on **Rake**. **TF-IDF** is one of the most popular term-weighting schemes today, in information retrieval, tf-idf or TFIDF, short for term frequency-inverse document frequency, is a numerical statistic that is intended to reflect how important a word is to a document in a

collection or corpus. It is often used as a weighting factor in searches of information retrieval, text mining, and user modeling. The tf-idf value increases proportionally to the number of times a word appears in the document and is offset by the number of documents in the corpus that contain the word, which helps to adjust for the fact that some words appear more frequently in general [9]. With TF-IDF we sorted for each document the words with the highest score and then selected the 7 first words of each document.

**Rake**: Rake stands for Rapid Automatic Keyword Extraction. It is an existing python implementation that uses the NLTK toolkit for the score calculations.

#### How the method works :

- It splits the text into sentences and generates some words as candidates.
- Various punctuation signs will be treated as sentence boundaries.
- All words listed in the stopwords file will be treated as phrase boundaries. This helps generate candidates that consist of one or more non-stopwords. However, it won't work in cases where the stopword is part of the phrase. For example, 'new' is listed in RAKE's stopword list. This means that neither 'New York' nor 'New Zealand' can be ever a keyword.
- For each keyword candidate generated, the algorithm computes the score of the pretended keyword, which is the sum of the scores for each of its words (if it's a composed word). The words are scored according to their frequency and the typical length of a candidate phrase in which they appear.
- The last step is ranking the scores and selecting the ones with the highest score

**Keywords in text**: Select from the text the keywords already given by the authors of the paper.

Once we have all the keywords from the three methods, we tried to select the most accurate keywords. So we created a method that selected the intersection of TF-IDF keywords and the Rake keywords with the keywords of the authors (if there are any keywords in the text), if the intersection is empty and we have the keywords in the text we only take the keywords of the text. If the paper does not contain any keywords, we take the intersection of the TF-IDF and the Rake outputs. If again, the intersection is empty then we only take the TF-IDF keywords because they were closer to the text keywords than the Rake Keywords.

Once we have our final keywords, we create a method that goes into the metadata and assign to every document a publication year.

Since we don't have the same number of documents for each year we needed to scale over the years and

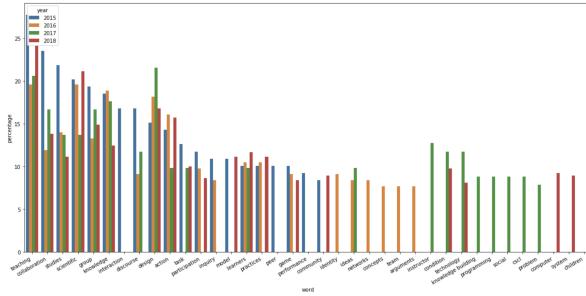


Figure 3: Trends of keywords over the years

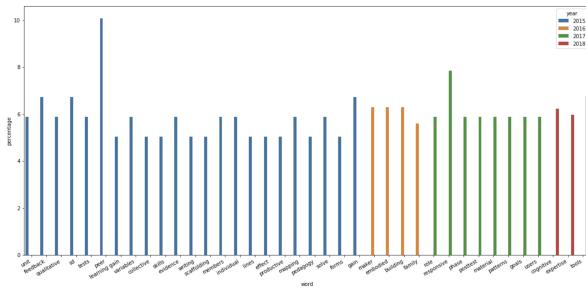


Figure 4: Distinct keywords over the years

use the percentage of appearance of a word. So after grouping the words by year and scaling them we found the graph below 3.

We wanted then to know what are the trends of words for each specific year, so we selected distinct words for each year. In other words it means that we selected for each year the words that were only present for this specific year and never appeared in the others. We get the figure below. As we can see from the graph 4, the most used word in the papers, more than 10% of 2015 papers used the word ‘peer’, more than 6% of the papers used the words ‘building’, ‘embodied’ in 2016. In 2017 the most frequent keyword is ‘responsive’ and in 2018 ‘tools’ is present with more than 6%.

### III. METHODS FOR DATA ANALYSIS

#### A. Document Clustering

To cluster the documents we decided to use two different approaches. The first one is based on Natural Language Processing (NLP) and the second one is rooted in network analysis. We find that network clustering using the references graph is most accurate.

##### 1) Approach 1 : Natural Language Processing (NLP):

###### a) Baseline - TF-IDF combined with K-means:

With the natural language processing, we started with a naive method which consisted of clustering the documents using only TF-IDF. Each document has a vector

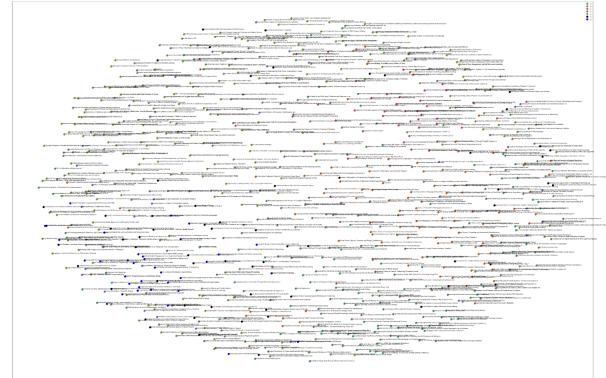


Figure 5: Baseline model document clusters

representation with values for each word present in the document depending on the frequency of the word in this document regarding its frequency in all the other documents. So if a word is a keyword to the document and is not a common word than its value will be more important than the other words in the vector of the document. Once we obtained our vector representation, we selected the max features parameter as 2500, it means that we only considered the top max\_features ordered by term frequency across the corpus. On this word representation we run K-means.

**K-means** clustering is a method used for clustering analysis, especially in data mining and statistics. It aims to partition a set of observations into a number of clusters ( $k$ ), resulting in the partitioning of the data into Voronoi cells. It can be considered as a method of finding out which group a certain object really belongs to.

This method did not give good results as we can see it in the graph below 5. All the vectors are together there is no clear separation between the documents. So we decide to use other methods using words embeddings combined with properties of TF-IDF.

b) *Second method : Average of GloVe vectors combined with K-means:* The second method was to use pre-trained words embeddings. We decided to use *GloVe* library pre-trained on Wikipedia dataset.

*GloVe*, coined from Global Vectors, is a model for distributed word representation. The model is an unsupervised learning algorithm for obtaining vector representations for words. Training is performed on aggregated global word-word co-occurrence statistics from a corpus (Wikipedia), and the resulting representations showcase interesting linear substructures of the word vector space. It is developed as an open-source project at Stanford.[10] The goal of *GloVe* is to group words with similar meaning together in a new vector space.

We constructed a dictionary from the words embeddings mapping each word of a document to its vector

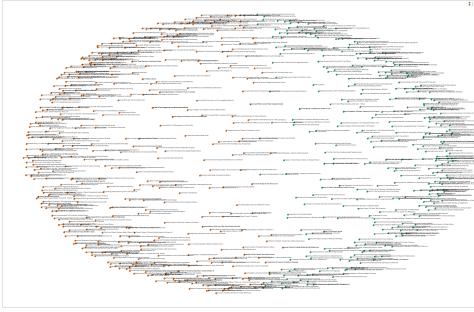


Figure 6: Cluster using Average of GloVe vectors combined with K-means

representation, then we averaged all the vectors of each document. Since each words has its own representation in space, the average of all the words for each document will lead to a vector representing a document.

Once we obtained a document vector representation, we applied K-means to it. The results were slightly better, but we had only two clear separated clusters. If we increase the number of clusters, then there is no clear separation between clusters visible anymore. The plot in figure 6 represents the documents belonging to each cluster.

*c) Third method : Weighted average of TF-IDF + GloVe combined with K-means:* With this method, we did not use the parameter max features of the TF-IDF vector, we took all the vocabulary generated by all the documents and we created a method that sorts the document frequency vector. We selected the 500 first words and computed a weighted average of the every word embedding vector with its respectively TF-IDF score.

Once we get the weighted average vector, since we were using a words embedding of 50 dimensions, we decided to use feature reduction technique to reduce the number of dimensions and represent our vectors in another vector space where all the features are orthogonal to each others.

To apply feature reduction to our vectors we used **Principal Component Analysis (PCA)** algorithm.

PCA is a statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables (entities each of which takes on various numerical values) into a set of values of linearly uncorrelated variables called principal components. If there are  $n$  observations with  $p$  variables, then the number of distinct principal components is  $\min(n-1, p)$ . This transformation is defined in such a way that the first

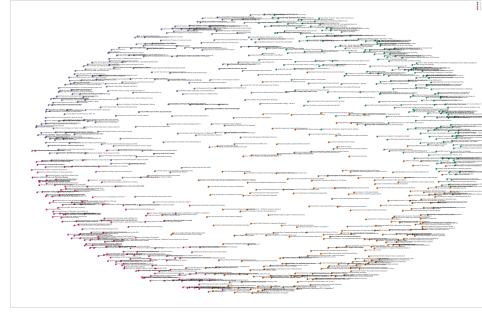


Figure 7: Cluster using Weighted average of TF-IDF + GloVe combined with K-means

principal component has the largest possible variance (that is, accounts for as much of the variability in the data as possible), and each succeeding component in turn has the highest variance possible under the constraint that it is orthogonal to the preceding components. The resulting vectors (each being a linear combination of the variables and containing  $n$  observations) are an uncorrelated orthogonal basis set. [11].

Once we get our new vector representation, we applied K-means and we obtained the best separation, with a clear separation of 4 clusters as we can see it in figure 7.

*2) Approach 2 : References Network Analysis:* For the references graph we chose to construct 4 different graphs to be able to extract as much information as possible. Firstly, to be able to cluster the documents that are about the same subject and secondly to gain more insights about the most influencial papers and trends of citations.

*a) First reference graph : Intersection between referenced documents:* Our first graph is constructed using the references of each document. It is constructed as follows: for each two documents, if there is an intersection between the references of both documents, then there is a link connecting the two documents. The weights of the edges are the length of the intersection set. The intuition behind it is that the higher the weight, the more the documents are likely to be about the same subject, since they are citing the same documents. From this graph, we create a subgraph of it. We filter the edges keeping only those who have a weight strictly higher than 3. This means that, each node of the subgraph, is related to at least one other node with 4 same documents referenced.

The resulting graph, plotted in figure 8, is an undirected graph with 805 nodes and 6699 edges. On this

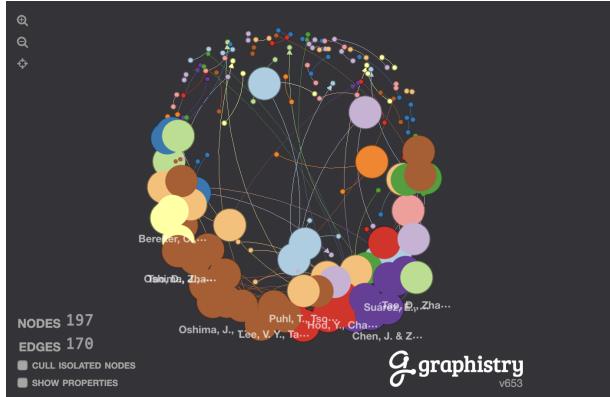


Figure 8: Intersection between references graph

graph we extracted two main features:

- Strongly connected components to be able to detect clusters and related documents.
- The in-degree in order to check if there are documents that are related to a lot of documents and if those documents are more general than others, and will not necessarily belong to one specific cluster.

*b) Second graph: Relations between references:*

The second graph relates more to the referenced documents and is not about the documents themselves. The aim of this graph is to see if there are documents that are more likely to be cited together. If this is the case, then these documents are likely to be about the same subject. This weighted graph is constructed as follows: for each document, we create a link between the references of the document. We create a dictionary of links, for each new document, if two references of the latter have already a link in the dictionary then we increase the weight of the link. For instance, if two references are present together in the 5 distinct documents, then the weight of the edge would be of 5. We obtained a graph of 10845 nodes and 148594 edges but we filtered the edges that have a weight higher than 4 and we get the graph below represented in this figure 9

We can directly select from the interactive graph, which can be found in the jupyter notebook corresponding to this report, the nodes with degree higher than 100. These nodes represent the most cited documents, but we will construct a subgraph to specifically tackle the most important referenced documents. On this graph, to be able to detect similar documents we also run the Strongly connected components algorithm. On top of this, we inspected the documents with higher weights. We will explain more about what we extracted from this graph in the result section.

*c) Third graph: Exploratory references graph:*

With this graph we aimed to find the most cited doc-

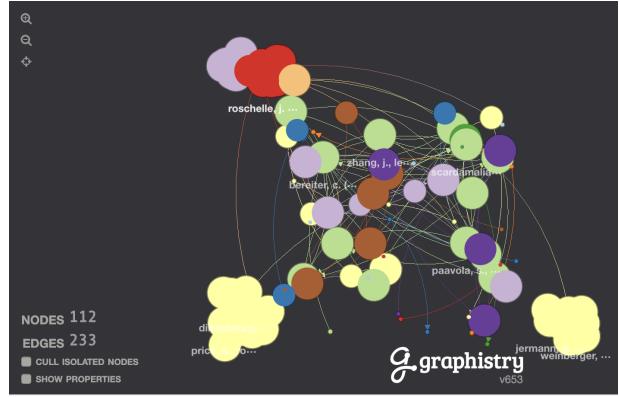


Figure 9: Relations between references filtered

uments, the documents that are the most influential and how citations evolve over time. This is a directed graph constructed in the following way: for each document, if the document cites another document, then there is a link between the document and the cited one. We get a graph of 11650 nodes and 13243 edges. We also constructed a subgraph of this graph, which is a co-citations in conferences, but due to a lack of data, we found only three documents that were cited by other documents.

On this graph, we run **PageRank** algorithm to be able to extract the most influential documents.

**PageRank (PR)\*** is an algorithm originally used by Google Search to rank web pages in their search engine results. In our case, PageRank is a way of measuring the importance of referenced papers. Compared to the in-degree analysis we get almost the same results. This is because almost all nodes have a zero out degree.

## B. Co-authorship

To get insight into collaborations within the community we built graphs based on co-authorship of papers. To build the graph structure and visualize it the python package *networkx* was used. When constructing the graph, we associated qualitative features –such as the name and institution of the author that is represented by a node– to the nodes and edges of the graph. We also associated the number of collaborations between different parties to the graphs as the weight of the edge.

*1) Graph Construction:* First, we built a co-authorship graph based on individual authors. Each node represented an author, and two nodes were connected if the authors collaborated on a paper. Other works [12], [13] have used a directed graph that only adds a link from the first author of the paper to all other authors. However, as we were more interested in collaboration patterns and less in finding the most influential first author, we built an undirected graph with edges between all co-authors. On

Pädagogische Hochschule Freiburg  
Universität Duisburg-Essen  
Carnegie Mellon University  
Ruhr-Universität Bochum  
Stanford University  
EPFL Lausanne  
Norwegian University of Science and Technology  
ETHZ - ETH Zürich  
Northwestern University

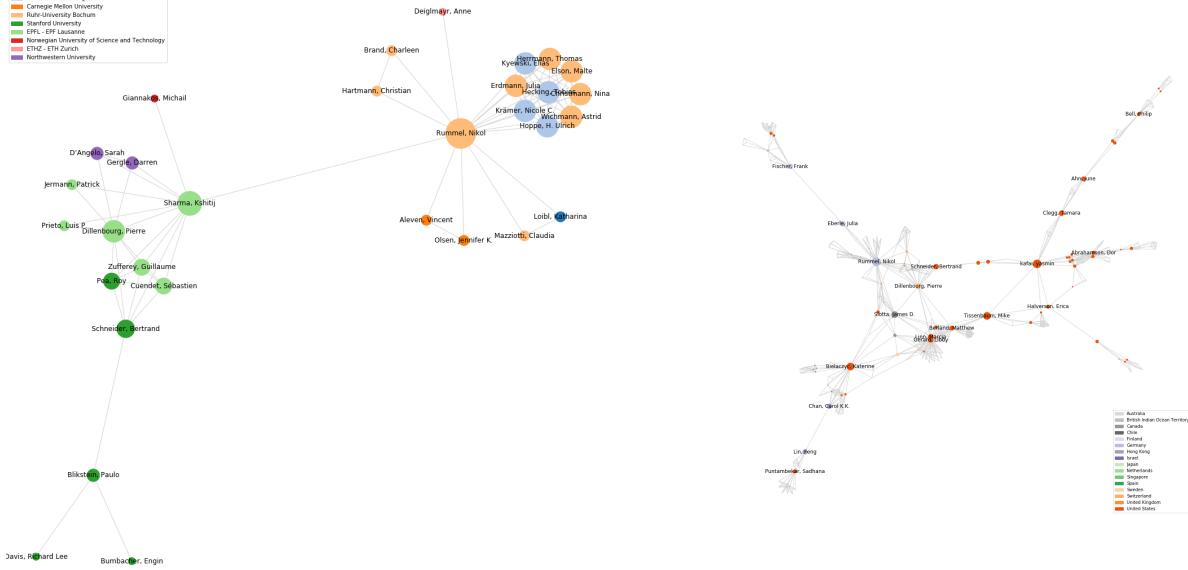


Figure 10: One of the connected components of the multi-year co-authorship graph

the downside, it yielded a very dense graph. Hence, we pruned the initial graph by removing nodes that did not satisfy additional conditions. Intuitively it is clear that a singular collaboration on a paper does not necessarily indicate that two people even know each other. A first condition was to link two nodes together only if the authors corresponding to each node appeared as co-authors in more than one conference. The graph produced by this approach was much more sparse and contained much fewer authors. Additionally, all connected components contained fewer than 35 authors, making it possible to consider each author individually. Such a component is shown in figure 10.

A second alternative approach was to consider the strength of bounds. In this graph authors are only linked if they have written at least 2 papers together. This method yielded a much larger component –containing 441 authors– as this second condition was weaker. To get an understanding of author to author relations in this component we needed to further prune away parts of the graph that were not insightful. When we looked at individual author within the community, we were interested in collaborative individuals. Hence we removed nodes that had fewer than 3 connections to other nodes. The resulting graph can be seen in figure 11. Details creating such a plot can be found in the section dedicated to this below.

Focusing in on the authors is not the only way to interpret the co-authorship graph. One should also consider the overall structure of the graph. As the co-



Figure 11: Largest Component of the multiple collaboration co-authorship graph

authorship graph is not connected, and contains one large main component with 1141 authors and 208 very small components with less than 25 authors in them we focused on analyzing the structure of the main component of the graph. First we aimed to detect sub-communities within the graph. Modularity is one measure used to find such communities. It measures the connection strength within a community relative to the strength of outgoing connections. Networks with high modularity have dense connections between the nodes within the subcommunities but sparse connections between nodes in different subcommunities. A non-parametric algorithm for community detection is the Louvain Method [14]. We used the implementation of Thomas Aynaud [15]. This gives a good partitioning of the graph into 26 communities of median size 42.84. This is quite a large number of communities, so we considered an alternative way to yield bigger partitions. For this we looked at normalized cuts, which partition the nodes into groups with many within-group connections and relatively few between-group connections. The advantage of this method is that we can choose how many cuts to make, and we are thus able to define the size of subcommunities. Finding normalized cuts in a graph can be done using *sklearns* spectral clustering algorithm by passing it the adjacency matrix of the graph. With this method we were able to neatly separate the main component graph in three large communities, as can be seen in figure 16

To further analyze these subcommunities we looked at the distribution of nationalities and institutions over

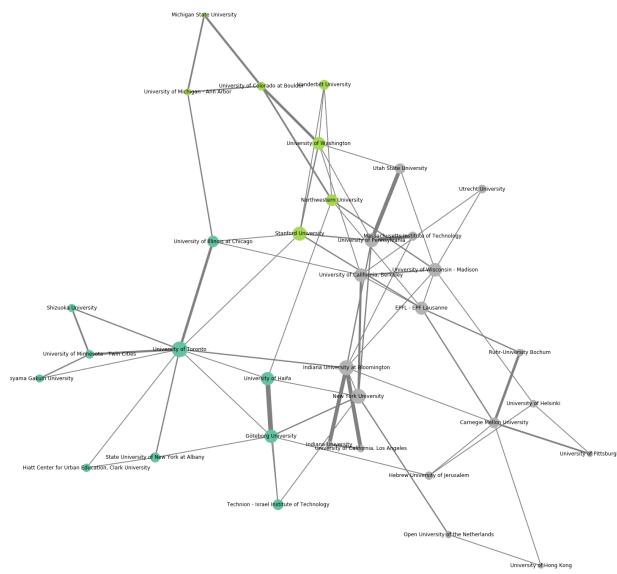


Figure 12: Different formats used to reference institution

the years. Additionally, to understand differences in the research being conducted, we analyzed the frequency of keywords used for the papers that were writing within that.

To get a better overview of how co-authorship happens across countries and institutions, we built graphs where every node represented a country, respectively an institution. Two countries or institutions were linked if two authors from these two institutions collaborated on a paper together. The country graph is simple enough that it does not have to be pruned. The institution graph on the other hand is again not comprehensible without further processing. Thus, similar to the first approach on pruning presented, we only considered institutions who collaborated in at least two conferences. The main component of this graph is shown in figure 12.

2) *Numerical analysis of graphs:* To analyze the structure of the graphs themselves we calculated the following measures:

- Diameter: longest shortest path in the graph.
  - Average clustering coefficient: measure of how much the nodes in a graph tend to cluster together. Average over all clustering coefficients. The clustering coefficient of a node is the fraction of triangles passing through the node over all possible triangles through that node [16].
  - Density: the number of edges in the graph divided by the maximum number of possible edges for a graph with the same number of nodes

3) *Plotting of co-authorship graphs:* Graphs were visualised with the built in plot functions and graph layouts

of networkx. The force directed layout was used as it clusters together nodes with strong connections. As this layout is produced by iterations and based on random initial conditions, we set a fix random seed to be able to reproduce the positioning. The size of each node was changed based on the position of the node in the graph. Central nodes with a lot of connections should be larger, while nodes with fewer connections should be smaller. To achieve this we used the measures of degree centrality for graphs with few nodes and edges and betweenness centrality for larger graphs. Degree centrality measures the fraction of degrees (number of connections) a node has with respect to the total number of connections in the graph. Betweenness centrality of a node is the sum of the fraction of all-pairs shortest paths that pass through it [17]. We found that degree centrality is more appropriate for smaller graphs as their structure is still easy to interpret. As by definition betweenness centrality takes into consideration the degrees of a nodes neighboring nodes, it better structures the nodes large graphs into central and less central nodes. To decide which nodes to label, we use the same mesures by selecting the nodes with the largest centrality, only plotting nodes who's cenrality is in the 90th percentile of the graph. The width of the edge is based on the weight of the edge, which in turn reflects the number of collaboration between two nodes.

### C. Creating Interactive Visualizations

To better understand and interact with the data we created interactive graphs on top of the static graphs. Many libraries exist to allow for interaction. While working on the project, we worked with *bokeh*, *holoviews*, *mpld3* and *plotly*. For some time implementing a custom visualization using *D3* was also contemplated. However, functionality provided by *plotly* is most general and allows to easily build interactive figures, while not being as heavy as *bokeh* and *holoviews*. Additionally, it is the most well maintained. *D3* requires a running server to be viewed and has a hard time dealing with large graphs. To plot figures we created them in *networkx* to grab the positioning of the nodes. Then we passed this positioning to *plotly* to plot.

## IV. RESULTS

In this section we present first insights gathered during the Bibliographic analysis.

## A. Community Composition & Structure

Participants from 372 institutions participated over the last 4 years, with a median of 128 institutions per year. The number of institutions participating increased each year, with a noticeable peak in 2018, with participants

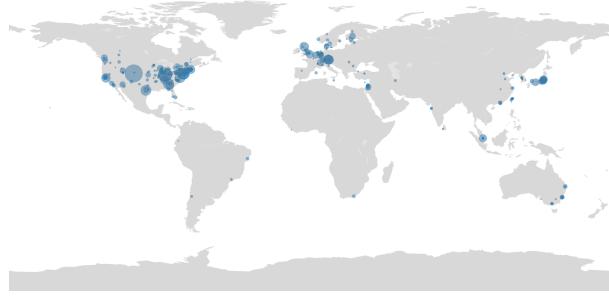


Figure 13: Global spread of community. Size of dot represent number of participants from each institution

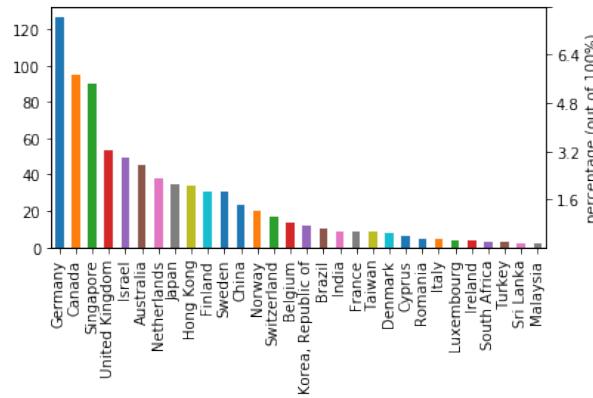


Figure 14: Number of Contributors by country, excluding the US

from 261 institutions. This is a considerable increase from the last conference ICLS in 2016 which had 135 participating institutions. Most participants are from the US, with a average of 50% of all participants being from the US. We could also observe a trend of increase in US participants with only 35% of US participants in 2015 and 51% in 2017, followed by 64% in 2018. Overall, participants from 41 nations have participated over the 4 years, with an average of 28 nations represented at each conference. After the US, participants of Germany, Canada and Singapore make up the biggest proportion of participants 14.

**1) Collaboration: International Collaboration:** The U.S. again stands out as a major player, with most collaborations happening between the U.S. and other countries. Out of 851 papers, 153 were international collaborations. Out of these 153, 108 involved at least one U.S. participant. Excluding US participation we found the collaboration graph presented in figure 15. We note that some countries, such as Japan, are not present in this graph. This means that all international collaboration of these countries, in this case japanese

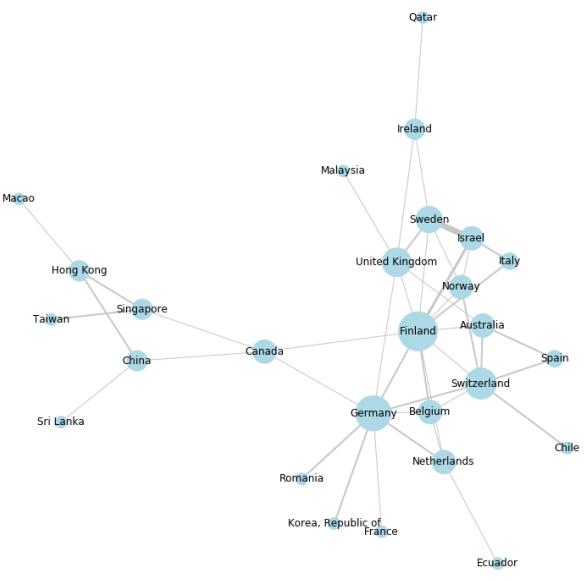


Figure 15: Collaboration between countries excluding the US

authors, also involves American authors. Additionally we observed that many Asian nations seem to collaborate almost exclusively with each other or Canada.

**Inter-institution collaboration:** In regards to collaboration between institutions, 318 out of the 815 papers were written in collaboration of authors from different institutions. The most collaborative university is the University of Indiana at Bloomington, followed by the University of Toronto. Considering the graph of university collaboration, we saw that none of the nodes is solitary - all institutions in the conference have collaborated with another institution at least once. There is a large main component and 6 other smaller clusters of universities. The cluster of main institutions collaborating in two conferences can be seen in figure 12. We see that the topology of the cluster is rather odd, and that while some strong links exist, most links are weak.

**Author collaboration:** As the graph yielded by only considering repeat collaboration during different years is sparse, containing connected components of size at most 31, we deduce that repeat collaboration over the years is rather infrequent. As the graph conditioned on multiple collaboration is more dense (having a connected component with more than 400 authors), we deduce that when multiple collaborations with the same authors happen, they yield work submitted at the same time.

The fact that normalized graph cuts will split the community into three prompted further investigation. As will be discussed in the next section, the graph of

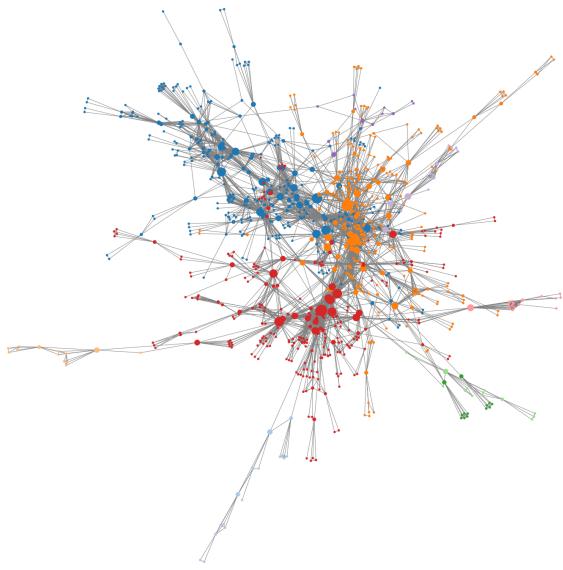


Figure 16: Subcommunities via normalized cuts in largest component of co-authorship graph

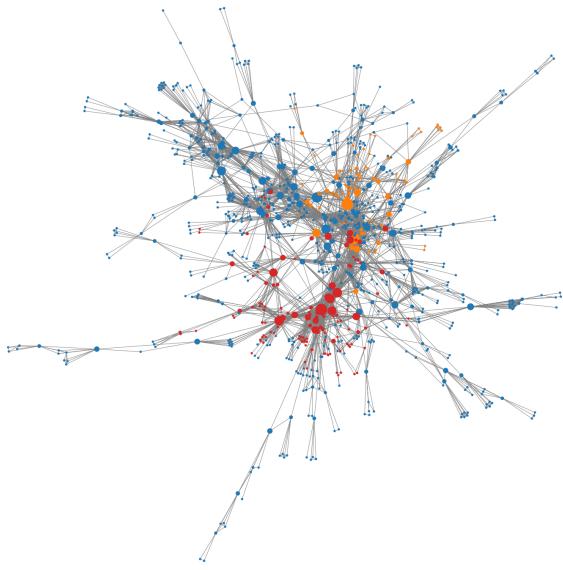


Figure 17: Subcommunities by participation in CSCL. Blue: no participation in main components of CSCL. Red and orange each mark participation in one of the two main components of CSCL



Figure 18: Keywords of the blue community

authors that participated in CSCL yields two large connected components. By coloring the authors contained in the two CSCL components with different colors, and authors in neither with a third color, we found a similar partitioning of the graph, as can be observed in figure 17. This is somewhat trivial, as the graph is built based on collaboration in these conferences. However, the analysis of nationalities present in the graph yielded an interesting results. All Asian participants from nations other than Singapore are contained in the same cluster. We also observed that one of the CSCL clusters is closer to the group of ICLS-only participants than the other. Additionally, through normalized splits, we found that the section of the normalized cut containing ICLS only authors is overwhelmingly American.

Considering the wordclouds in figure 18, 19, 20 of the most common keywords in each of the large splits, we found a first indication that the splits also reflect some thematic differences in research.

2) *ICLS and CSCL*: In our analysis, we could find many structural differences between the two conferences. Firstly, ICLS is much larger and american dominated than CSCL. We also observed differences in the most cited papers and the average number of authors of a paper. By constructing the graph based on collaboration of authors in each conference, we found that while ICLS contains one giant component, CSCL contains two large main components. All these results are however to be taken with a grain of salt, as for each conference we only have data for two years, and the year 2018 is a marked outlier. Additionally, we are unable to find a strong difference between the metrics of the graphs, that



Figure 19: Keywords of the red community



Figure 20: Keywords of the orange community

is diameter, density, clustering coefficient and average path length.

3) *Movement of participants:* Knowing the affiliation of participants to an institution in a given year, we can now find members that change institutions over the years. Over the years of recording we found 14 participants that have changed countries and 47 that have moved institutions. We can see that most movement occurred within the US respectively Europe, see figure 21.

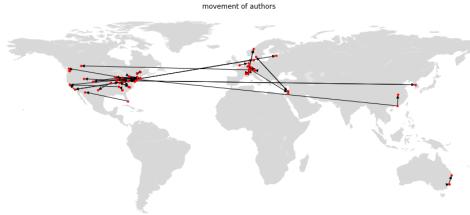


Figure 21: Movement of participants across institutions

#### B. Changes over the years

We observe a growth of the number of participants over the years, and by building co-authorship graphs for each year of the conference, we observed that the density and clustering coefficient of the graph indicate a steady decrease, while the diameter and average shortest path length tend to increase. This hints that the community is attracting new member that do not collaborate with the established group of participants. But, as the years 2018 shows a massive increase of participation, in number of participants and number of institutions represented, as well as a large decrease in density of the graph –which is reduced by more than half from the previous ICLS conference– it is impossible to draw any statistically significant conclusion.

#### C. Clustering and references graphs results

The distribution of referenced publications years are plotted in the graph 22. We can see from the graph that there is a kind of a trend, each two, three years there is a pick of documents referenced followed by a two or three years of stagnancy. We can also notice that throughout the year each pick is higher than the older one, which can be explained by the fact that the International Society of the Learning Sciences community is investigating more and more tools to facilitate the learning over the years.

Since the number of documents is increasing over the years, we also wanted to investigate more if the relevance of all the documents published is increasing as well. Using the Pagerank algorithm and the In-degree of each nodes to get the most influencial documents referenced of the exploratory citations graph that we constructed before. We filtered the first 100 documents and then contructed a histogram of the counts by year. We obtained the figure 23:

We can see that in the top 100 cited documents more than 12% are from 2006, 8% are from 2013, 6% from

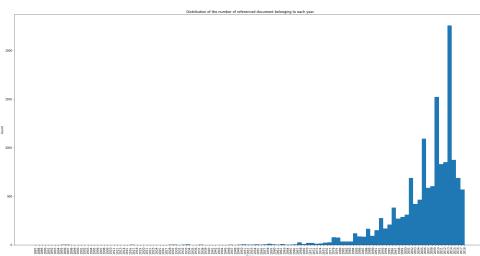
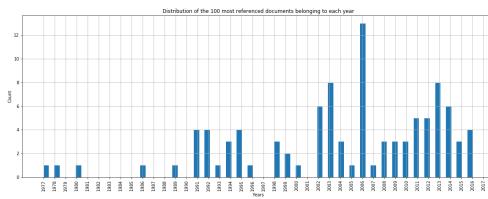


Figure 22: Distribution of the number of referenced documents belonging to each year



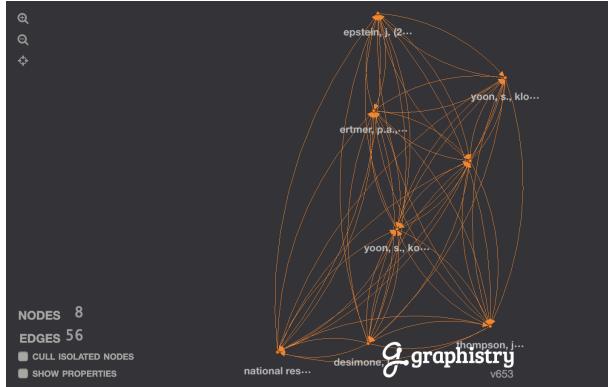


Figure 27: Fully connected references graph

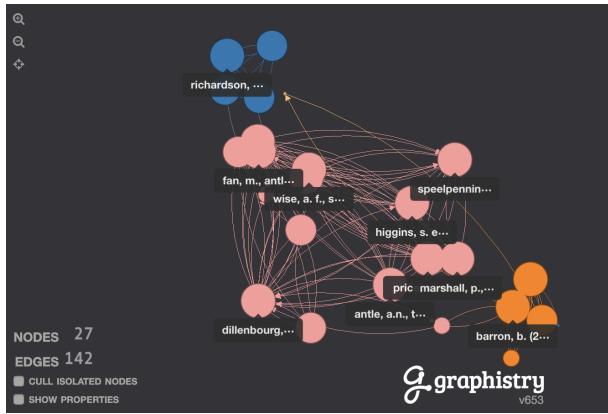


Figure 28: Giant component graph

In the table 29 below are the documents contained in the giant connected component represented in the interactive graph 28.

We wanted to apply the same logic with the papers in-conferences. So we focused on the analysis of the first

Documents	
0	white, t., & peir, r. (2011). distributed by design: on the promises and perils of collaborative learning with multiple representations
1	scheider, b., wallace, r., peir, r., & bilskiene, p. (2013). preparing for future learning with a tangible user interface: the case of neuroscience
2	dillenbourg, p., & evans, e. (2011). interactive tablets in education
3	speelpennings, t., antle, a. n., doering, t., & van den haven, e. (2011). exploring how tangible tools enable collaboration in a multi-touch tabletop game
4	price, s., regier, y., stanton, d., & smith, h. (2003). a new conceptual framework for casl
5	fan, m., antle, a. n., neustadtter, c., & wise, a. f. (2014). exploring how a co-dependent tangible tool design supports collaboration in a tabletop activity
6	jordan, b., & henderson, a. (1995). interaction analysis: foundations and practice
7	rogoff, b. (1998). apprenticeship in thinking: cognitive development in social context
8	kirschner, p. a., sweller, j., & clark, r. e. (2008). why minimal guidance during instruction does not work: an analysis of the failure of constructivist, discovery-based, experiential, and inquiry-based teaching
9	wise, a. f., saghafian, m., & padmanabhan, p. (2012). towards more precise design guidance: specifying and testing the functions of assigned student roles in online discussions
10	roschelle, j. (1995). learning by collaborating: convergent conceptual change
11	toncaso, m. (1995). joint attention as social cognition
12	jemmann, p., mullini, d., nüssl, m.-a., & dillenbourg, p. (2011). collaborative gaze footprints: correlates of interaction quality
13	higgins, a. e., mercier, e., bard, e., & hatch, a. (2011). multi-touch tables and the relationship with collaborative classroom pedagogies: a synthetic review
14	scheider, b., & peir, r. (2013). real-time mutual gaze perception enhances collaborative learning and collaboration quality
15	meier, a., spads, h., & rummel, n. (2007). a rating scheme for assessing the quality of computer-supported collaboration processes
16	barron, b. (2003). when smart groups fail
17	antle, a. n., tenenbaum, j., macarthur, a., and robinson, j. (2014). games for change: looking at models of persuasion through the lens of design
18	ulmer, b., & iahl, h. (2000). emerging frameworks for tangible user interfaces
19	shear, o., strait, m., valdes, c., feng, t., linton, m., & wang, h. (2011). enhancing genomic learning through tabletop interaction
20	feneau, y. & thöndorf, j. (2006). finding design qualities in a tangible programming space
21	richardson, d. c., & date, n. (2005). looking to understand: the coupling between speakers' and listeners' eye movements and its relationship to discourse comprehension
22	lave, j., & Wenger, e. (1991). situated learning: legitimate peripheral participation
23	scheider, b., jermann, p., zuferay, g., & dillenbourg, p. (2011). benefits of a tangible interface for collaborative learning and interaction

Figure 29: Documents belonging to the giant component

Documents	
cluster 9	Mazzotti, C., Rummel, N., & Dsigmayer, A. (2016). Comparing Students' Solutions When Learning Collaboratively or Individually Within Productive Failure. Hartmann, C., Angerbach, J., & Rummel, N. (2015). Social Interaction, Constructivism and their Application within (ICSG) Theories In Lindwall, O
1	Mazzotti, C., Lohb, K., & Rummel, N. (2015). Collaborative or Individual Learning within Productive Failure: Does the Social Form of Learning Make a Difference? Hartmann, C., Rummel, N. & Lohb, K. (2016). Communication Patterns and Their Role for Conceptual Knowledge Acquisition From Productive Failure
2	McBride, E., Vitale, J., Appelbaum, L., & Linn, M. (2017). Examining the Flow of Ideas During Critique Activities in a Design Studio Environment. McBride, E., Vitale, J., & Linn, M. (2016). Learning Design Through Science vs Engineering Design Process
3	McBride, E. A., Vitale, J. M., Appelbaum, L., & CLIN, M. (2016). Use of Iterative Computer Models to Promote Integration of Science Concepts Through the Engineering Design Process
4	McBride, E., Vitale, J., & Linn, M. (2018). Middle School Student Ideas on the Relative Affordances of Physical and Virtual Models. McBride, E., Vitale, J., & Linn, M. (2018). Promoting Cognitive Processes of Knowledge Integration: Middle School Students' Ideas on the Relative Affordances of Physical and Virtual Models
cluster 7	Greenhow, C., Herathkovic, A., Baruch, A. F., Askari, E., Tsovalioti, D., Asturian, C., Puri, T., Weisberger, A., Blauter, E., & Priman, J. (2018). Teachers and Professional Development: New Contexts, Modes, and Concerns in the Age of Social Media
1	Erkens, M., Schöttborn, P., & Bodenre, D. (2016). Qualitative and Quantitative Information in Cognitive Group Awareness Tools: Impact on Collaborative Learning
2	Schneebauer, L. & Bodenre, D. (2016). How Socio-Cognitive Information Affects Individual Study Decisions
3	Asterhan, C. & Bouton, E. (2017). Secondary school peer-to-peer knowledge sharing through social network technologies in Smith, B
4	Erkens, M. & Bodenre, D. (2017). Which Visualization Guides Learners Best
5	Tsovalioti, D., Dutta, N., Puri, T., & Weisberger, A. (2017). Group and Individual Level Effects of Supporting Socio-Cognitive Conflict Awareness and its Resolution in Large SNS Discussion Groups: A Social Network Analysis in Smith, B
6	Himelrich, S. & Bodenre, D. (2016). Effects of Implicit Guidance on Contribution Quality in a Wiki-Based Learning Environment
7	Puri, T., Tsovalioti, D., & Weisberger, A. (2015). A Long-Term View on Learning to Argue in Facebook: The Effects of Group Awareness Tools and Argumentation Scripts in Lindwall, O
cluster 6	Dornfeld, C. L. & Puntambekar, S. (2016). Negotiation Towards Intersubjectivity and Impacts on Conceptual Outcomes
0	Everstone, A. L. & Puntambekar, S. (2015). Internalization of Physics Concepts and Relationships Based on Teacher Modeling of Collaborative Learning at Individual and Group Levels in Smith, B
1	Dornfeld, C., Zhao, N., & Puntambekar, S. (2017). A Mixed Methods Approach for Studying Collaborative Learning Processes at Individual and Group Levels in Lindwall, O
2	Martin, N. D., Gnedlow, D., & Puntambekar, S. (2016). Peer Scaffolding to Learn Science in Symmetrical Groups Collaborating Over Time in Lindwall, O

Figure 30: Documents in conferences belonging to the same clusters

graph that uses the intersection between the referenced documents of each two by two papers. The length of the intersection set had to be more than 4.

So if more than 4 documents are referenced by two distinct documents than the latter are more likely to be about the same subject. In the table 30 we can find some of the papers that belong to the same clusters.

When applying this kind of clustering we found that most of the clusters contains almost the same authors with different papers, which is logic because if they are specialists in one topic and write several papers on it, they might reuse some references in several documents.

## V. CONCLUSIONS

In summary, we were able to show how the dominant some institutions and countries are within the conference. Collaboration patterns between countries, institutions and authors were extracted and analyzed. Sub-communities found within the network of co-authorship reflect the division between conferences and within the conference. Using the analysis of the affiliation over the years we followed the migration of participants across institutions. Natural language preprocessing and metadata parsing enabled the extraction of trends of keywords over the years

Overall, due to lack of data, we were often unable to find non trivial results; this lack of data is aggravated by the fact that the conference of 2018 presents as an outlier in the dataset. Additionally, taking document clustering as an example, network analysis of references were more accurate than NLP regarding documents and references clustering.

We thus strongly recommend adding additional conferences to the dataset.

## REFERENCES

- [1] Isls website. <https://www.isls.org/conferences>, accessed 2019-1-10.
- [2] Poppler homepage. <https://poppler.freedesktop.org/>.
- [3] Isls author guidelines. [https://cscl2019.com/upload/ISLS\\_Author\\_Guidelines.pdf](https://cscl2019.com/upload/ISLS_Author_Guidelines.pdf), 2019.
- [4] Scholarcy reference extraction api. <https://ref.scholarcy.com/api/>, accessed 2019-1-10.
- [5] Apa formatting and style guide by purdue writing lab. [https://owl.purdue.edu/owl/research\\_and\\_citation/apa\\_style/apa\\_formatting\\_and\\_style\\_guide/reference\\_list\\_articles\\_in\\_periodicals.html](https://owl.purdue.edu/owl/research_and_citation/apa_style/apa_formatting_and_style_guide/reference_list_articles_in_periodicals.html), accessed 2019-1-10.
- [6] University domains and names data list and api. <https://github.com/Hipo/university-domains-list>, accessed 2018-12-28.
- [7] Matthew Honnibal and Ines Montani. spacy 2: Natural language understanding with bloom embeddings, convolutional neural networks and incremental parsing. *To appear*, 2017.
- [8] Glove website. <https://nlp.stanford.edu/projects/glove/>.
- [9] Tf-idf definition. <https://en.wikipedia.org/wiki/Tf%E2%80%93idf>.
- [10] Glove definition. [https://en.wikipedia.org/wiki/GloVe\\_\(machine\\_learning\)](https://en.wikipedia.org/wiki/GloVe_(machine_learning)).
- [11] Pca definition. [https://en.wikipedia.org/wiki/Principal\\_component\\_analysis](https://en.wikipedia.org/wiki/Principal_component_analysis).
- [12] France Cheong and Brian J Corbitt. A social network analysis of the co-authorship network of the pacific asia conference on information systems from 1993 to 2008. *PACIS 2009 Proceedings*, page 23, 2009.
- [13] James W Hesford, Sung-Han Sam Lee, Wim A Van der Stede, and S Mark Young. Management accounting: a bibliographic study. *Handbooks of management accounting research*, 1:3–26, 2006.
- [14] Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. Fast unfolding of communities in large networks. *Journal of statistical mechanics: theory and experiment*, 2008(10):P10008, 2008.
- [15] Thomas Aynaud. Louvain community detection. <https://github.com/taynaud/python-louvain>, 2018.
- [16] Duncan J Watts and Steven H Strogatz. Collective dynamics of ‘small-world’ networks. *nature*, 393(6684):440, 1998.
- [17] Stephen P Borgatti, Daniel J Brass, and Daniel S Halgin. Social network research: Confusions, criticisms, and controversies. In *Contemporary perspectives on organizational social networks*, pages 1–29. Emerald Group Publishing Limited, 2014.