

# Supporting Collaborative Information Analysis: A Classroom Study

**1st Author Name**

Affiliation  
City, Country  
e-mail address

**2nd Author Name**

Affiliation  
City, Country  
e-mail address

**3rd Author Name**

Affiliation  
City, Country  
e-mail address

## ABSTRACT

Collaborative information analysis involves representation and strategic analysis of complex information space over extended time periods. Current tools either do not fully support such workflows, or not support collaboration at all. We developed *CAnalytics* to address two critical challenges: to enable data modeling and data analysis within an integrated workspace, and to maintain activity awareness in a complex interface. We conducted a classroom study of student analysts using the tool over multiple usage sessions. The results suggest positive evidence of awareness support and an all-in-one environment. We also observe the distinction between accretion and filtering in data modeling, different approaches to information uncertainty in artifact analysis, as well as different collaboration strategies. Given these results, we suggest several design implications for collaborative information analysis tools.

## Author Keywords

Collaborative information analysis; visualization; classroom study;

## ACM Classification Keywords

H.5.3. Information Interfaces and Presentation (e.g. HCI): Computer-supported cooperative work

## INTRODUCTION

Collaborative information analysis is a form of sensemaking wherein a team analyzes a complex information space of facts and relationships to identify and evaluate causal hypotheses. A common example of collaborative information analysis is crime investigation; a variety of putative facts are assembled, including financial records, witness observations and interviews, and social connections of various sorts among persons of interest, from which investigators collaboratively assess means, motives, and opportunities, articulate and investigate further hypotheses and deductions, and develop one or more theories of the crime. Other examples include intelligence analysis, business intelligence, scientific research, and social constructivist learning.

A critical challenge for information analysts is building an adequate preliminary data model from textual documents, and insuring that the data model is employed effectively in hypothesis development and evaluation. This is an open challenge [1]. Standard methods often do not support it at all. For example, Analysis of Competing Hypotheses (ACH) assumes that data has been modeled, and that relevant evidence can be adduced appropriately to various hypotheses, but provides no structured support for either. Other techniques, such as Information Extraction and Weighting (IEW), help structure modeling of evidence, but do not extend utilization of evidence to hypothesis generation. These methods seem to assume information analysis as a “waterfall” model, in which data modeling, data analysis, and hypothesis development are separate staged activities. Analysis is thus often broken into places in separate workspaces, and requires considerable efforts to duplicate and coordinate low-level data structures in different software systems. We therefore are motivated to develop an integrated workspace in which analysts can model and analyze data in one place, and we utilize the system as an experimental instrument to investigate analytic behavior afforded by such integration.

Any work of information analysis at a non-trivial scale is fundamentally collaborative. Yet computer tools to support this work (e.g. Analyst’s Notebook and PARC ACH) do not support collaboration. Thus, analysts must coordinate their work outside of their tool support, through frequent face-to-face meetings, or other messaging exchanges. A key enabler for effective collaboration is *activity awareness*, defined as team’s awareness of its own sustained collaborative activity [6]. Derived from Activity Theory [16], activity awareness transcends synchronous awareness of who collaborators are, where a collaborator is looking, etc. It encompasses issues of many different kinds of information covering all aspects of an activity, such as events, tasks, goals, mediating artifacts, social interactions, and group values and norms, which becomes higher demanding as the activity becomes more complicated. Awareness support in such a complex activity of information analysis is perhaps also more challenging than many other situations (e.g. collaborative writing). Teammates could be working with much more complex data structure (e.g. spatial data, temporal data, and relational data, as opposed to only text), coordinating through multiple analytic artifacts (e.g. map, timeline, network, as opposed to only a document), and making sense of different levels of analysis, assumptions, and hypotheses, both synchronously and asynchronously throughout a long-term course of collaborative

Paste the appropriate copyright statement here. ACM now supports three different copyright statements:

- ACM copyright: ACM holds the copyright on the work. This is the historical approach.
- License: The author(s) retain copyright, but ACM receives an exclusive publication license.
- Open Access: The author(s) wish to pay for the work to be open access. The additional fee must be paid to ACM.

This text field is large enough to hold the appropriate release statement assuming it is single spaced.

Every submission will be assigned their own unique DOI string to be included here.

interaction. Hence we will investigate how technology can mediate team collaboration in a sufficiently complex analytic task over multiple usage sessions.

We situate our study in classroom learning of information analysis. Classroom study provides a natural environment in which participants engage in multi-sessions, sufficiently complex class projects. Due to difficulty in accessing professional analysts or having them in long term design loops, analysts in training who are learning to be information analysts become a good compromise. They already have some knowledge and experience with state-of-the-art analytic techniques and tools. In class projects students are graded on their ability to understand and enact professional practices of information analysis. This strong normative emphasis on problem solving practices is a great evaluation context for new interactive tools: Tools are only valuable to the students insofar as they actually support better practices and better outcomes.

We thus are motivated to investigate the feasibility, effectiveness and consequence of supporting integrated data modeling and analysis, as well as supporting activity awareness in complex information analysis, in the context of classroom study. We have developed a tool that includes annotation for data modeling, interactive visualization for data analysis, and a set of awareness features. For the balance of this paper, we describe the tool we have developed, the classroom settings, and our observations. We conclude with design implications derived from the study as well as future work.

## RELATED WORK

Many studies have been reported to investigate specific design features to support collaborative information analysis. For example, Goyal and Fussell [10] studied the effect of hypotheses sharing on sensemaking. Mahyar and Tory [17] designed a visualization to connect collaborators' common findings and evaluated its support for team performance. Hajizadeh et al. [12] explored how sharing teammate's interactions affects awareness. These studies report interesting results of controlled lab studies to validate hypotheses of specific design features. However, they do not provide insight on how teams would collaborate with a complete tool as nexus of features in the real world over extended period of time.

Field studies were conducted aimed to understand design requirements of collaborative information analysis in more realistic settings. Chin et al. [7] observed and analyzed the analytic strategies, work practices, tools and collaboration norms of professional intelligence analysts. Kang and Stasko [15] studied how student analysts, as in our study, completed in-class intelligence projects. Carroll et al. [5] attempted to model a complex analytic task scenario in a lab setting, and examined the development of team awareness in a four-hour-long task. These studies helped improve understanding of current work practice with state-of-the-art tools or no tools at all. We built our tool based upon their study findings, and pursue to further explore design implications by investigating tool usage in a similar naturalist environment.

Our study took place during the 10th week of the course. Before that students learned several analytic techniques, in-

cluding IEW (a technique to extract and assess values of evidence), ACH (a technique to evaluate multiple hypotheses against evidence), timeline analysis and network analysis, as well as state-of-the-art tools to implement these techniques. In the eighth weeks students practiced applying these techniques in a hands-on project. A typical workflow started with IEW to extract and model evidence from documents, followed by building analytic artifacts such as an ACH Matrix in PARC ACH, a timeline and a network graph in Analyst's Notebook. Since data from IEW table could not be shared or extended to other tools, students had to replicate data for each different tool. Most tools they used lacked serious collaboration support (except that some teams used Google Doc to construct an IEW table). Analysts were unable to contribute simultaneously (an issue known as production blocking [9]). The analysis work was often divided by tools: each person created and analyzed an artifact with a tool on their own. This had the consequence that findings and hypotheses be made without integrating collective efforts and diverse knowledge. Analysts must coordinate work by manually sharing notebooks or graphs through email or cloud storage service (e.g. Dropbox), resulting in a scattered placement of results, requiring repeated manual resynchronizing to identify redundant or missing pieces of information, analysis of information, and analytic hypotheses. The instructor and students in our study were all aware of the shortcomings of available tools with respect to support of collaboration.

## CANALYTICS FEATURES

We developed a collaborative information analysis tool, CANalytics, to support teams of analysts in identifying, visualizing, integrating and assessing facts from multiple sources. The design is informed by our paper prototype studies, where we examined communication patterns and team's spontaneously created artifacts when teams were engaged in a complex crime scenario. We also take into account findings from empirical studies conducted by Chin et al. [7] and Kang and Stasko [15] when making design decisions.

CAnalytics supports evidence modeling through annotation. In the document view users can select and highlight a snippet of text and annotate it as a type of entity such as a person, location, event, etc., or as a relationship between entities. Unlike in other entity-based systems such as [3, 21], we use annotation to allow analysts to manually create evidence objects of interest. Manual annotation allows for greater user control, allows more integrated source data objects to be identified, and avoids the user problems associated with automatic identification of disaggregated people, locations and times [2]. Users can decide their own information of interest and granularity that best suits their ad-hoc analytic needs.

Users can add attributes to the annotated object, e.g. adding time attribute in an event, and placing coordinate in a location. Users can also make reference to other objects in the attribute; for example, users can add people objects to an event indicating that these people were involved in the event. In this case, a relationship between the people and the event is automatically created. Users can also explicitly create a rela-

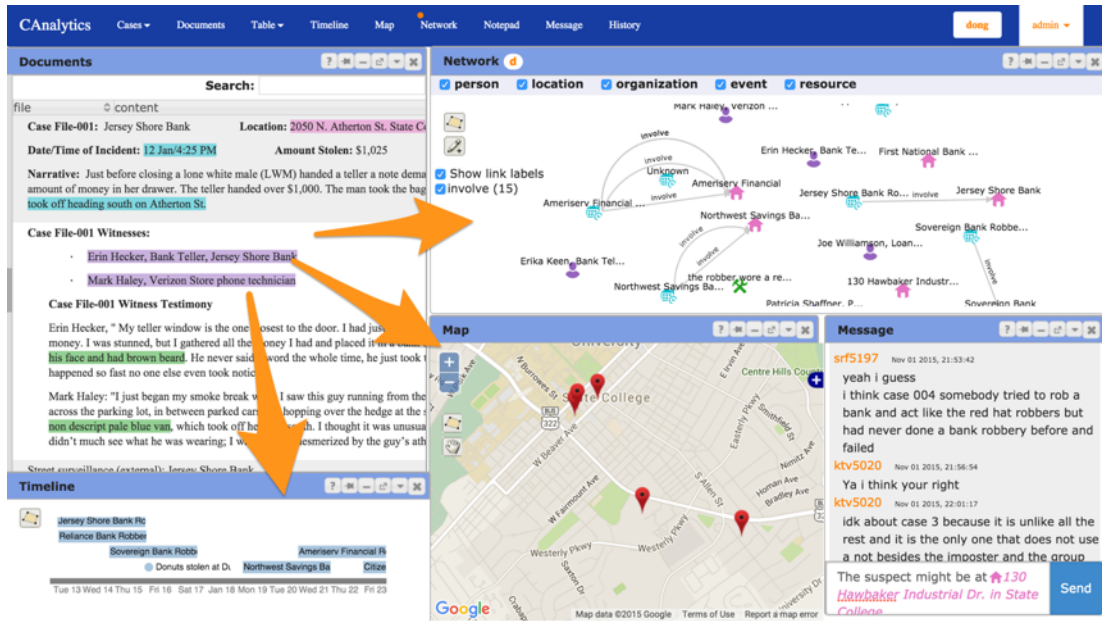


Figure 1. CAnalytics user interface

tionship between two entities by selecting a source entity and a target entity and labeling the relation name.

When users are explicitly creating an annotation, they are also implicitly creating a provenance of the modeled entity—annotation records the source where a data object was created. As observed in Carroll et al.’s [5] empirical study, while integrating information in a visual artifact helps sharing and pooling information with teammates, the action also removes problem information from its original context. Participants often forget what an entity refers to and why it is such positioned in an artifact. Entities in CAnalytics are linked back to the location of documents where they were created through annotation. Users can always re-access the data objects for provenance, a critical requirement emphasized by Chin et al. [7].

The data modeled through annotation is then displayed in multiple coordinated views in the same workspace, including table, timeline, map, and network graph—artifacts frequently constructed to hold attribute data, temporal data, spatial data and relational data respectively [5]. Figure 2 shows an example of the tool interface: when an annotation is created in the document view with information about time, location, participants, and their relationships, a new event is created in the timeline view, a new location is created in the map view, and new people are added to the network graph with a labeled edge representing the relationship (or new edges are added to existing nodes). Hovering mouse over an entity will activate an entity detail window that displays attributes in detail, and analysts can modify, or re-model the entity in situ.

The views are coordinated and afford brushing and linking interaction; that is, when users apply a graphic filter in one view, related information is displayed in other views. Thus the analyst can retrieve entities within a time range using

timeline filter, make a spatial query with map filter, or select a cluster of entities by drawing a filter area in the network view.

CAnalytics supports real-time collaborative editing, similar to the Google Tools. Users can open several concurrent editors to collaboratively edit multiple annotations. Annotations are immediately shared within a team and are automatically added to teammates’ document views and other visualizations. More than Google Docs, CAnalytics supports the editing of more complex data structures than text, including spatial data, temporal data, and object attribute data.

In addition to real-time data sharing, CAnalytics is built in with several other awareness features, including a notification system, a feature we named “tool coordinator”, a message tool, and a history tool, and a collaborative editor. A notification system sends individual’s actions to the team, in the form of a text box in the top right corner of the workspace. The tool coordinator is an iconic indicator on top of a tool window, suggesting who is working on the tool. The message tool is a real time chat window that enables team communication with persistent message history. The system also maintains a traceable log of time-stamped individual activities in a history tool. Users can learn team activity about who did what to which object at when. Entities and relationships mentioned in the message tool and history tool are hyperlinks that will trigger pop-up detail window when being moused over. With these awareness features, users who work synchronously can be informed of others’ activity continually; users who work asynchronously will be able to use the history to reconstruct their work status and become aware of changes beyond the point of their last interaction.

We also included a simple notepad to support collaborative hypothesis development. We integrated Etherpad <sup>1</sup>, an open-

<sup>1</sup><http://etherpad.org>

source collaborative editor similar to Google Doc, for teams to compose their hypotheses. Users can insert tables (e.g. an ACH matrix) and images (e.g. screenshots of the tool views).

### CLASSROOM STUDY SETTINGS

The context for this study was an undergraduate course in an intelligence training program in a US university. The program was designed to train students to become professional intelligence analysts. A key requirement of the course is to emphasize hands-on practice on team-based intelligence analysis. During the first nine weeks, students learned strategic knowledge (e.g. bottom-up analysis and top-down analysis) and structured analytic techniques, such as IEW, ACH and network analysis, and practiced to apply these techniques to solving two small projects with state-of-the-art tools including PARC ACH and IBM Analyst's Notebook.

Our study began from the 10th week of the course and lasted for one week. The task was to investigate a series of bank robberies fabricated by the course instructor. Teams were provided with a set of documents pertaining to seven robberies, including police reports, witnesses reports, video records, and news media. The task was designed open-ended, which meant that there was no single answer to the task. The instructor explained that the task was to simulate real world scenarios, in which analysts always reasoned in the circumstances of uncertainty, ambiguity, and complexity. The instructor told the students that 6 hours was expected to complete the project, including in-class and outside-class work. In the end of the project students were required to submit a team report, describing their hypotheses, assumptions, conclusions and supporting evidence.

Students were given a tutorial on CAnalytics a week before the project began. One of the authors walked through features of CAnalytics and then let students accomplish a small case analysis on their own pace. During the study week, one author was always available to help with any technical issues. Although students were encouraged to make full use of CAnalytics, to ensure a naturalist environment students were always free to employ any other tools that they believed useful.

Of the 98 students enrolled in the course (from two sections), 73 consented to participate in the study. Students were randomly assigned into 25 teams (23 three-person teams and 2 two-person teams). Research suggested that group size be an important factor in group collaboration, thus two-person teams might behave differently from three-person teams. We thus excluded data from the two-person teams in our analysis in this paper. Also, from the log (and confirmed by their questionnaire), we found that one team made little use of CAnalytics and opted for other tools (Google Doc). Hence their data was also excluded. Thus in this paper we reported the result from 22 teams.

All the students held major in the program of Security and Risk Analysis. Most (75%) of them were in the third academic year (3.05 years in average), indicating that participants in our study had relatively advanced experience and knowledge in intelligence analysis. Participants' age ranged



Figure 2. Classroom setting

from 19 to 28 (20.3 in average). 77% of the participants were male.

We employed several data collection approaches. We administered a post-study questionnaire, which included several 7-index questions measuring individual's self-reported awareness (adapted from [8]), team communication (adapted from [8]), collective efficacy (adapted from [8]), perceived performance (adapted from [11]), and cognitive load [13]. The questionnaire also included open-ended questions asking how the tool helped or impeded their work. We captured user interactions with system logs. Instead of simply logging low-level events like mouse click and keyboard strokes, we recorded actions such as creating an annotation and deleting an entity. Finally, we reviewed team reports and graded them as an indicator of team performance. Since the task was open-ended, there was no single right answer. We constructed an assessment rubric together with the course instructor by listing all possible hypotheses and evidence from the documents, with a full score of 16, but we also credit teams if they provide other reasonable arguments. The first author and a research assistant graded the reports independently. If the grades differ by less than 2, an average is set as the final grade (14 out of 22 reports). Otherwise (the rest 8 reports), the two graders review the reports together and make an agreement.

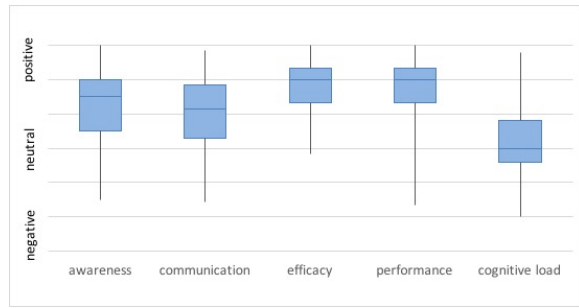
### RESULT

Over the week, teams created 1805 entities and 1529 relationships in total. The number of entities teams created ranged from 24 to 223 ( $M=82$ ,  $SD=39.9$ ), and the number of relationships ranged from 7 to 237 ( $M=69.5$ ,  $SD=51.0$ ). The large variety was related to team data modeling strategy, which will be detailed later.

Overview of the survey items indicates that students rated positive on CAnalytics overall, as shown in Figure [fig:survey]. CAnalytics were ranked favorably in all aspects except cognitive load, towards which they had a close to neutral feeling.

### Intertwined data modeling and data analysis





**Figure 3. Survey responses (box shows Q1-Q3 and median; whiskers show maximum and minimum)**

We examined the pattern of data modeling and data analysis first by qualitatively looking at a visualization of the entire interaction log (e.g. Figure [fig:intertwined]a shows one team's interaction). All teams worked intensively on data modeling as they started off the project. This was the phase when teams were getting themselves familiar with the documents and made initial input into CAnalytics. Starting from certain point (e.g. indicated by the blue dash line in the figure), however, teams started to work on visualizations. They did not wait to start analysis till they finished data modeling, as they returned to making annotations again later. Indeed, the activity of data modeling and data analysis were highly intertwined since then. Participants switched from one activity to the other activity frequently. The state transition diagram [fig:intertwined]b demonstrates the interweaving in a quantitative way, in which we encode the number of transitions as width of the link. This result confirmed our design expectation that data modeling and data analysis should not be separate staged activities, and that an integrated environment streamlines the workflow.

### Data modeling: accretion vs. filtering

From granularity of entities, We noted a distinction between accretion and filtering strategies in data modeling, similar to what was reported in the paper prototype study. Filtering is selectively modeling of data and adding to an artifact. Users must decide what information is relevant, and thus what is to be excluded, as well as what granularity of information is to model. Filtering requires more team coordination, because teammates must reach a common ground of the current problem as well as information needed to answer the problem. Figure [fig:network'accretion]a is an example of filtering, representation only the key information of robberies.

Accretion is an attempt to comprehensively represent the problem by adding all information to an artifact. Users extract every fact from the document, regardless of its immediate relevance to the problem. Accretion costs less coordination as it is relatively mechanical note taking. A disadvantage of accretion is that it could be time consuming to model all details and the produced artifact could be fairly complex. The strategy causes a huge variety of number of created entities and relationships mentioned earlier. An example is Team 108, who modeled every step the suspects took, which resulted in far more entities than the average and more cluttered network view (Figure [fig:network'accretion]b). Users reflected that

they spent too much time in details that they lost the bigger picture:

We would find ourselves glued to our computer screens, and spent too much time on intelligence gathering rather than analysis (P135)

While similar data modeling strategy was reported in the paper prototype study [5], users with CAnalytics seemed far more tempted to accretively add information, with far more entities and cluttered views. Students realized in the end that many annotations did not help them solve the problem at all because many entities were unrelated to their problem.

I felt that after we were done annotating, we hadn't really accomplished anything and that we were no closer to solving the case than when we had started. In the end it didn't really help that we had annotated the data. (P86)

### Artifact construction: fact vs. inference

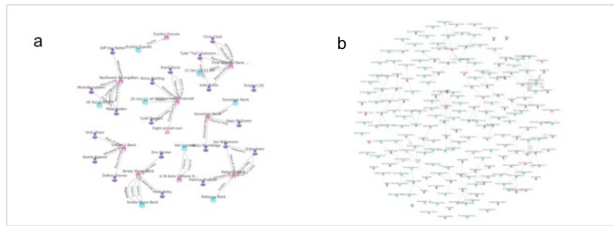
We examined the analytic artifacts teams created, the network graph in particular because social relationships played the most critical role in this specific scenario and teams spent most time on network analysis (as reflected from the log). We found the network artifacts fell into one of two categories: the networks consisted of 1) separate clusters, or 2) connected clusters. For example, networks from 8 teams (36%, mean performance=7.8) consist of separate clusters (Figure [fig:network'cluster]a). Nodes within a cluster are connected, representing information space of a robbery case; Nodes between clusters are nonetheless not connected, indicating each robbery is a self-contained case. However, these teams still claimed connections between robberies in their report. Where did they externalize these connections? Or did the teams simply share orally and held them in mind? It turns out that teams documented possible relationships between robberies in the notepad tool.

In contrast, 6 other teams (mean performance=8.3) created networks composed of connected clusters. While a cluster is still a representation of a robbery, some of them are connected through an evidence node. An example is Figure [fig:network'cluster]b, in which we mark four *connectors* that link the clusters. These connectors were key evidence that led the teams to hypothesize that those robberies were related and might be committed by the same criminal group.

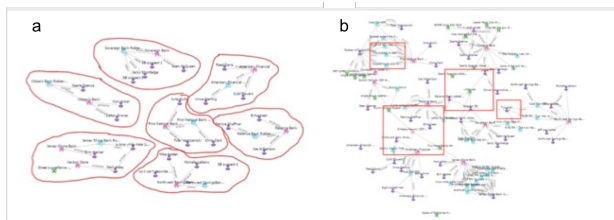
While many causes might account for the different network views, we attempt to interpret the difference from a perspective of *uncertainty*. For instance, links within a cluster are often factual relationships modeled from raw documents (e.g. a white van was witnessed at a location), but links between clusters are often inferences beyond literally documented (e.g. a white van at location A is the same van witnessed at location B). Teams creating separate clusters only represented facts in the network and held evidence with uncertainty in a separate artifact. One advantage of distinguishing facts and inferences is that teams can be aware of assumptions made when making a conclusion. And since all inferences are held in one place, teams are forced to confront them and review their uncertainty iteratively in the process. However, the strategy also adds difficulty to analysis as analysts



**Figure 4.** (a) Visualization of interaction logs of Team 107. Each row of colored marks indicates the sequence of top level activities a participant performed. (b) State transition diagram of interaction logs of Team 107. Each node is an activity, whose size represents the time spent on the it (M: data modeling; A: data analysis; H: hypothesis development; C: coordination); a link represents a switch from one activity to another, whose width encodes the number of switches. We see highly frequent transitions between data modeling and data analysis



**Figure 5.** Network artifact comparison: filtering (a) vs. accretion



**Figure 6.** Network artifact comparison: separate clusters (a) vs. connected clusters (b). The parts highlighted in red squares in (b) are key evidence that connects clusters

may overlook or fail to combine evidence scattered in different artifacts.

On the contrary, some other teams overlaid facts and inferences in the same artifact. Both facts and inferences drove the layout of the network, thus influencing team’s framing of the problem. Most teams made evaluation of the uncertainty of inferences when adding them to the network. This strategy was relatively more interactive among teammates: they needed to negotiate, evaluate, and reach consensus on the value and validity of every inference. To some extent teams might forget whether a relationship is factual or inferred, and ask whether conclusion derived from the visualization can be trusted under uncertainty.

### Collaboration and awareness

One recurring theme in the subject feedback we collected was that the collaboration features were helpful for solving the problem on a team basis. In the survey 88% of the students rated positively on their group awareness. Participants appreciated that the tool complemented traditional analytic tools such as Analyst’s Notebook with real time synchronization features like Google Doc, and added visual analytics to Google Doc which only supports text writing. To quote one participant, “*CAnalytics is like an analysts notebook that multiple people could work on at once [...and] an analysts version of a Google Doc.*” (P65). For example, participants

reflected that they could now contribute simultaneously without concerns of interference and could have everything in one place instead of manually sharing documents using cloud service.

It was much easier to coordinate as a team with CAnalytics because we could all work on the same system at the same time. Without CAnalytics, we were forced to do the work separately and compile all the work onto one system after we had finished. (P156)

Students also reported that being able to see teammate’s status made the task more motivating and engaging:

During class I wasn’t sure if my teammates were doing work for that class or another thing but then seeing their dot [tool indicator] switch between applications on the software and updates pop up on my screen I knew they were doing work for 231. (P141)

The fact that you can see what other teammates are doing and they can see what you are doing creates a sense of accountability in terms of separating the work load. (P51)

Another repeated theme was awareness that all teammates were executing the team plan. Participants reflected on their experience that a common team breakdown was misunderstanding of the team plan. Teammates thought they had reached agreement to a plan till in the end to find no collaborators were doing the things as expected. CAnalytics made it easier because they could always see if teammates were doing the expected work; and if not, they could communicate immediately rather than in the end of the project.

Participants reported many other instances of awareness they realized using CAnalytics. We categorized them based on the element of awareness, or the essential problem of awareness of *what* [19], into social awareness, information awareness, action awareness, history awareness, and intention awareness, as shown in Table [tab:awareness].

When asked what features helped them stay aware of team activities, 28 participants mentioned the tool coordinator, 24 mentioned the notification system, 19 mentioned the history tool, 14 mentioned the real-time update of user-generated data, 12 mentioned the collaborative editor, and 7 mentioned the message tool. While the number of mentions does not simply indicate tool usefulness, it suggests users appropriate these features and were explicitly aware of their support.

Students’ positive feedback on awareness was further corroborated by interaction logs. For example, we measured the

Table 1. Subject feedback of awareness aspects

Element	Example
Social awareness <i>who is present?</i>	CAAnalytics helped me stay aware,of my teammates activities because I could see who was logged on in the top,right corner (P123)
History awareness <i>Who has done what?</i>	The way you are able to view when and where your teammate made or updated annotations/information was the key to staying aware of what your team has done. It is a great tool in respects to that. For example, I was able to view the changes my team made while I was not using the CAAnalytics tool at the same time they were using the history tab. (P171)
Information awareness <i>What is being changed?</i>	CAAnalytics was very helpful in keeping us updated on what was being changed/noted/amended by whom and when. This was very beneficial for staying on the same page and knowing what changes were being made so no one individual was out of the loop. (P157)
Action awareness <i>Who is doing what?</i>	I liked how you could always see what your teammate were viewing on the website. For example I was working on the bluf when my teammates were working on the network part of the program. If I were to come across a piece of information that I thought might be helpful to them I would just tell them. My teammates did the same thing in return. (P51)
Intention awareness <i>Who is going to do what?</i>	CAAnalytics showed what tab [tool] my teammates were working on which helped me be aware of what they were working on. For example, if I saw that one of my teammates was on the network tab, I knew that they were attempting to connect the information that was relevant to one another.,I would then be able to mention any new findings I had that could influence their work (P160)

number of entities accessed by collaborators versus by the author only. While data generated by users is automatically shared, it is up to collaborators to choose to read the shared information or ignore information altogether. A high awareness team would keep updated with collaborators' generated information and read information soon after it is shared; whereas a low awareness team might experience a significant delay or even never access it. We found that most teams shared a high proportion of entities ( $M=77.6\%$ ). We found that in average, 77.6% of the created entities were accessed by at least one teammate.

One major critique is the lack of sharing of intermediate analytic insights in close collaboration. To gain insights individuals are constantly exploring and rearranging visualizations. For example, analysts can apply a filter to have a reduced data view of interest, highlight an area to sharpen analytic focus, and re-layout the node-link graph to make their own cluster of relevant entities. While the shared data pool represents the information the team have in common, views of data reflect and contextualize individual's *interpretations* toward the data. The view together with its interpretation represents user's intermediate analytic status. While it may not be a mature hypothesis or conclusion, the intermediate insight is still worth sharing because it may inspire collaborators. With CAAnalytics participants complaint that they could not easily communicate insight together with the associate views to the team. The team could "*be looking at the same information but arranged in completely different ways*" (P131).

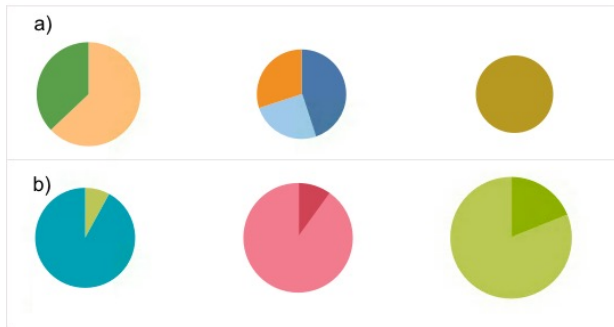
### Collaboration strategies

We noted different collaboration strategies from interaction logs. Seven teams followed *document-based collaboration*: they divided their work by evenly distributing the provided documents among team members (as shown in Figure [fig:labor' division]a). Each member read through and made annotations on their own set of documents. An advantage of this strategy is that individuals get less workload and thus have more time to think deep into their own documents. Individuals also implicitly take over the responsibility of their assigned documents to gain insights and share them when the team synthesize findings in later analysis. When the team

needs information from one document, they rely on the "document owner" to share his/her finding. The fallback of this strategy is thus in case an individual fails to identify or convey evidence in the document, the team may overlook the information altogether [4].

In contrast, four teams followed an *entity-based collaboration* strategy. Instead of dividing by documents, they divided work by entity types: each individual went through all documents but only annotated entities of certain types, e.g. teammate A only annotated persons and teammate B annotated locations (as shown in Figure [fig:labor' division]b). This strategy saves teammate's time on data modeling. And since each person focuses only on certain entities, they are more likely to identify recurring patterns, for example, the white van used in multiple robberies. However, focusing on certain entities could lead individuals to superficial syntactic scanning of documents instead of deep reading. This could further lead to extremes of annotating all entities of the type, whether they are related to the problem at hand or not. We found from the interaction log that EBC teams created more entities ( $M=101$ ) in average than DBC teams ( $M=94.5$ ), which seems to back our guess. Moreover, with emphasis on certain entities, individuals are likely to know only partial aspects of a robbery and hence have difficulty connecting and synthesizing facts to deduce any conclusion. As one participant reflected, "*we broke up by entity type, which reduced our individual involvement in each other's entity types* (P99)." The average performance of EBC teams ( $M=7.63$ ) was lower than that of DBC teams ( $M=9.25$ ), which seems to confirm our analysis, although the difference was not statistically significant.

The rest (eleven) of the teams did not show specific labor division patterns. Indeed, teams did not necessarily have to divide their work in order to collaborate, especially when the collaborative tool provides possibility to work closely together. Teammates could read and annotate the same document because they could see new annotations by others in real time and build on other's annotation. Figure [fig:close' collaboration] shows an example where one team worked on the same document simultaneously. This has the



**Figure 7.** Pie charts showing different labor division strategies. Each pie chart corresponds to one team member. (a) Document-based collaboration. The pie chart shows the documents one team member annotated, color coded by document ID. (b) Entity-based collaboration. The pie chart shows the entity types one team member created, color coded by entity types.



**Figure 8.** Graph showing the timeline of one team creating annotations. Each row corresponds to one team member. Each bar represents an annotation, color coded by document ID. The red blocks highlights the periods when all teammates worked on the same documents simultaneously.

advantage that teammates are always on the same page and can discuss hypotheses throughout the analysis process. Participants did have concern for possible duplication. As one participant complaint, “we could not actively see the changes our teammates were making until well after they had made them.(P46)” This was because an annotation was shared only *after* it was created, yet another teammate might be drafting an annotation in the meantime.

## DISCUSSION

### Reflections on method

The goal of the study is to explore design opportunities to support collaborative information analysis by evaluating tool usage in a natural environment over multiple usage sessions. Our work builds upon prior empirical studies (e.g. [5, 4, 15, 7]) and embodies their design implications in our tool. We are also to complement research that only tests tools in short term lab studies (e.g. [8, 10, 12]). Due to the constraint of time (mostly within one hour), these studies had to employ a simplified task with reduced content and complexity. Participants would thus have no need to create complex information artifacts (e.g. with a single artifact or few items in an artifact). More complex task would push participants to create more sophisticated artifacts (e.g. multiple views or cluttered display that requires filtering) and to balance between team coordination and individual analysis, which would have provided more insights into team-based analytic process.

Our classroom study attempts to gain deeper insights on collaborative information analysis behavior by better simulating real world settings in two aspects:

First, the study spanned multiple usage sessions over one-week long. Teams are able to learn to adapt to team functions and to appropriate the tool to best serve their team purpose [20]. They can explore different strategies in seeking to solve the problem, and to make changes if they get stuck with one strategy. For example, we noticed that two teams decided to change the use of the tool halfway in their analysis. One team started with dividing work by case documents, but later decided members should annotate different entity types. Another team started with an accretion strategy by annotating all entities. Later they discovered that this strategy brought too much noise, and decided to clean out irrelevant entities (filtering strategy). Such change occurs as a consequence of increased awareness of team functions and tool capabilities, which takes time to develop.

Second, participants in this study are being trained to become professional analysts. Before our study they had already been introduced to the information analysis techniques and the state-of-the-art analytic tools. In their reflections, participants often compared CAnalytics to those tools, as well as the different teamwork experiences. Therefore their feedback is likely to provide deeper insight into strength and weakness of CAnalytics. In addition, the students are young learners that are willing to employ new work practices supported by features in tools. They are important parts of the future analytic community. In some sense, their practice can be treated as a view into the future of practice of the community [18].

Yet classroom study also has limits. For example, many factors and variables could exist that affect team performance. The fact that these factors are often impossible to model or control adds to the difficulty in data analysis (e.g. teammate absence). Also, data collection is challenging because team interactions are not always accessible. Teams can choose to work synchronously or asynchronously, and it is difficult to predict when or where the interaction of most interest is to occur. Verbal communication is not accessible, which could be useful to infer team awareness as a complement to interaction logs. Our work identifies both positive evidence and problematic situations, and propose potential solutions and possible hypotheses, yet rigorously evaluating these solutions and validating hypotheses is beyond this study. Lab studies and case studies can be conducted in the future to address these problems with greater control and deeper data access.

### Reflections on result

The study provides encouraging results on supporting collaborative information analysis with an integrated workspace and awareness design. Participants appreciated an all-in-one environment where they could share raw documents, evidence snippets, views of evidence and hypotheses in one place. They liked the fact that they could contribute simultaneously without blocking or interfering each other. Another benefit of the collaborative tool is to keep teammates aware of each other’s activities. In addition to the many aspects of awareness listed in Table [tab:awareness], which help establish a common ground for plan execution, obtaining real time feedback and staying aware that teammates are indeed following the team plan is equally important. Moreover, partici-



pants suggest the awareness features in the tool have positive social impact such as *social facilitation*, whilst CAnalytics provides an environment where individuals simultaneously engage in the same activity in awareness of each other (co-action effect) and the notification and other awareness features act as a source of arousal [23].

Our study suggests that tools play an important role in shaping user's behavior towards more collaborative behavior. With traditional single-user tools, students often employ a divide-and-conquer strategy; they divide their job by tools, work individually on separate tools, and compile the results together in the end. In our study, we observed many teams spontaneously conducted closer collaboration and enjoyed being able to contribute simultaneously. A potential problem with simultaneous contribution is duplicated or conflicted efforts. To minimize the problem, further nuanced awareness features could be added. For example, [12] revealed collaborator's action even though no change is being made. We can visualize where the collaborator is making an annotation as an indicator of interest of that piece of text. This is like typing indicator bubbles in chatting tools (e.g. Facebook Messenger) which provides awareness of actions currently being performed before outcome becomes visible.

A misconception about information analysis is that data modeling and data analysis are two staged activities. This is akin to the waterfall software development model, which features a sequential process that flows downwards through the phases of requirement conception, software design, software implementation, testing and maintenance. Critics have pointed out that the staged approach may not work properly, because clients may not know exactly what their requirements are before they see the working software and designers may not be fully aware of future difficulties in a new software product. Instead, an iterative design process is often required that leads to reframe user requirements, redesign, redevelopment, and retesting.

Similarly, relying only on information that has already been modeled and delivered to analysts will probably not solve all analytical problems [14]. It will probably be necessary to look elsewhere than evidence already extracted, re-model the data, and dig for more information. Yet many analytic tools assume data has already been modeled and ready to be visualized and analyzed and affords no utility to construct or refine data models. As Ware termed as "asymmetry in data rates" [22, p. 382], analytic tools emphasized data flowing from systems to users far more than from users to systems. Functionalities are mostly designed to adjust visual representation rather than remodel data underlying the representation, which is a critical aspect in information analysis. Re-modeling of the data could lead to a different picture of the problem (e.g. adding a link between two clusters changes the layout of the network, and thus framing of relationships between two robberies), leading to completely different analytic path. The interaction log we captured demonstrated such an iterative analysis process and the positive subject feedback confirmed data modeling as an integral part of analysis. This kind of problem solving activity is better supported by an in-

tegrated workspace design approach, rather than a workflow-driven design approach.

We noted the importance of representing uncertainty. We observed teams in our study spontaneously employed two different approaches to deal with team uncertainty: either to mix them for better synthesis or to separate them for better clarification. This demonstrates both challenge and opportunity to design for uncertainty support. We propose that a richer graphic language and interaction be designed so that analysts can encode uncertainty into the network view. For example, links and entities with different uncertainty can be visualized in different transparency. Users can *filter* by uncertainty so that users can choose to have only facts to take into account or review all inferences.

We noted several cases where teams created far more entities than needed with an accretion strategy, and strikingly, also far more than that reported in paper prototype studies [5]. Why did this happen? We guess both the context of classroom study and the system design contributed. Unlike in the lab study where teams are temporarily assembled, teams in a class evaluate peers either consciously or unconsciously and value how themselves are being evaluated. Such social pressure motivate individuals to make contributions, and indeed to make *visible* contributions, more than valuable contributions. That is, participants noticed that their work activity was visible to their partners, and accordingly prioritized doing more visible work over doing less visible work. In some cases, this led to a new problem of easy and less valuable contributions that were highly visible - such as creating and therefore sharing data entities that were not particularly important, and subsequently made data models seem cluttered. For example, creating and therefore sharing an entity gets immediately notified to the team whereas weighing the importance and relevance of an entity goes silent in the system. We need to investigate approaches to making significant contributions more visible, or perhaps making it more immediately visible that less important contributions are indeed less important.

Views could still get cluttered as data volume increases and analysts dig into event details (e.g. representing suspect's all actions to identify patterns of common actions in two robberies). Indeed, analysts often engage in multi-level analysis in parallel, frequently coordinating among, say, confirmation of location of an event, to comparison of two events, to over-viewing all events as a robbery as a whole. A possible design solution is to enable collapsible data views. A collapsible view can help analysts focus attention on a certain level of details at a time while conveniently switching between levels, and when in collaboration, draw teammate's attention to a specific item.

An important improvement to CAnalytics is to support view sharing, or sharing intermediate analytic status. A common solution is to have a distinction of public view and private view [8]. We are proposing a step further that views should be treated as a team *resource*, just like data. Views as resource should be sharable, extensible, and reusable. For example, several participants reflected that there were situations when

they found a collaborator's view inspiring and wanted to continue exploring that view on their own side without manually reproducing the view. With views as resource, the owner of a view can deliberately save the view as a shared resource when they feel it appeals to collaborators. Other people can reuse the view to their need. Shared views should be interactive rather than static images, so that analysts can perform all interactions including filtering and highlighting, and are able to evolve the view with collective team efforts, a critical requirement emphasized in [5].

## CONCLUSION

In this paper, we present findings from a classroom study in which teams of information analysts in training collaboratively completed a complex intelligence project mediated by our tool. As collaborative information analysis is increasingly a typical and chronic task, it is important for research to examine, understand, and provide effective tools and environments for these long-term, real-world CSCW interactions. This requires situating research in more complex work activity contexts, and directly investigating interactions, experiences, and outcomes in those contexts. Our classroom study provides preliminary encouraging results on supporting team interactions. We also identify new problems regarding data uncertainty and contribution visibility. We need to investigate approaches to addressing these problems in our future iterative design and evaluation.

## ACKNOWLEDGMENTS

To be added

## References

- [1] Richard V Badalamente and Frank L Greitzer. "Top Ten Needs for Intelligence Analysis Tool Development". In: *First Annual Conference on Intelligence Analysis Methods and Tools* May (2005).
- [2] Eric a. Bier, Stuart K. Card, and John W. Bodnar. "Entity-based collaboration tools for intelligence analysis". In: *2008 IEEE Symposium on Visual Analytics Science and Technology*. Ieee, Oct. 2008, pp. 99–106. ISBN: 978-1-4244-2935-6. DOI: [10.1109/VAST.2008.4677362](https://doi.org/10.1109/VAST.2008.4677362). URL: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4677362>.
- [3] Eric a Bier, Stuart K Card, and John W Bodnar. "Principles and tools for collaborative entity-based intelligence analysis." In: *IEEE transactions on visualization and computer graphics* 16.2 (2010), pp. 178–91. ISSN: 1077-2626. DOI: [10.1109/TVCG.2009.104](https://doi.org/10.1109/TVCG.2009.104). URL: <http://www.ncbi.nlm.nih.gov/pubmed/20075480>.
- [4] Marcela Borge et al. "Patterns of team processes and breakdowns in information analysis tasks". In: *Proceedings of the ACM 2012 conference on Computer Supported Cooperative Work - CSCW '12*. New York, New York, USA: ACM Press, 2012, pp. 1105–1114. ISBN: 9781450310864. DOI: [10.1145/2145204.2145369](https://doi.org/10.1145/2145204.2145369). URL: <http://dl.acm.org/citation.cfm?doid=2145204.2145369>.
- [5] John M. Carroll, Marcela Borge, and SI Shih. "Cognitive artifacts as a window on design". In: *Journal of Visual Languages & Computing* 24.4 (2013), pp. 248–261. URL: <http://www.sciencedirect.com/science/article/pii/S1045926X13000207>.
- [6] John M. Carroll et al. "Awareness and teamwork in computer-supported collaborations". In: *Interacting with Computers* 18.1 (Jan. 2006), pp. 21–46. ISSN: 09535438. DOI: [10.1016/j.intcom.2005.05.005](https://doi.org/10.1016/j.intcom.2005.05.005). URL: <http://iwc.oxfordjournals.org/cgi/doi/10.1016/j.intcom.2005.05.005>.
- [7] George Chin Jr, Olga A. Kuchar, and Katherine E. Wolf. "Exploring the analytical processes of intelligence analysts". In: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 2009, pp. 11–20. ISBN: 9781605582467. URL: <http://dl.acm.org/citation.cfm?id=1518704>.
- [8] Gregorio Convertino et al. "Supporting common ground and awareness in emergency management planning". In: *ACM Transactions on Computer-Human Interaction* 18.4 (Dec. 2011), pp. 1–34. ISSN: 10730516. DOI: [10.1145/2063231.2063236](https://doi.org/10.1145/2063231.2063236). URL: <http://dl.acm.org/citation.cfm?doid=2063231.2063236>.
- [9] Michael Diehl and Wolfgang Strpebe. "Productivity Loss In Brainstorming Groups : Toward the Solution of a Riddle". In: *Journal of personality and social psychology* 53.3 (1987), pp. 497–509. DOI: [10.1037/0022-3514.53.3.497](https://doi.org/10.1037/0022-3514.53.3.497).
- [10] Nitesh Goyal and Susan R Fussell. "Effects of Sense-making Translucence on Distributed Collaborative Analysis". In: *CSCW '16 Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social Computing*. 2016.
- [11] Nitesh Goyal et al. "Effects of Implicit Sharing in Collaborative Analysis". In: *Proceedings of the Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 2014, pp. 129–138. ISBN: 9781450324731. URL: [http://www.cs.cornell.edu/%7B~%7Dngoyal/1470%7B%5C\\_%7DChi2014.pdf](http://www.cs.cornell.edu/%7B~%7Dngoyal/1470%7B%5C_%7DChi2014.pdf).
- [12] Amir Hossein Hajizadeh, Melanie Tory, and Rock Leung. "Supporting awareness through collaborative brushing and linking of tabular data." In: *IEEE transactions on visualization and computer graphics* 19.12 (Dec. 2013), pp. 2189–97. ISSN: 1941-0506. DOI: [10.1109/TVCG.2013.197](https://doi.org/10.1109/TVCG.2013.197). URL: <http://www.ncbi.nlm.nih.gov/pubmed/24051785>.
- [13] Sandra G. Hart and Lowell E. Staveland. "Development of NASA-TLX (Task Load Index): Results of Empirical and Theoretical Research". In: *Advances in Psychology* 52.C (1988), pp. 139–183. ISSN: 01664115. DOI: [10.1016/S0166-4115\(08\)62386-9](https://doi.org/10.1016/S0166-4115(08)62386-9). arXiv: [arXiv:1011.1669v3](https://arxiv.org/abs/1011.1669v3).
- [14] Richards J Heuer. *Psychology of intelligence analysis*. Lulu.com, 1999.

- [15] Youn-ah Kang and John Stasko. “Characterizing the intelligence analysis process: Informing visual analytics design through a longitudinal field study”. In: *2011 IEEE Conference on Visual Analytics Science and Technology (VAST)*. Ieee, Oct. 2011, pp. 21–30. ISBN: 978-1-4673-0014-8. DOI: [10.1109/VAST.2011.6102438](https://doi.org/10.1109/VAST.2011.6102438). URL: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=6102438>.
- [16] Aleksei N Leont’ev. “The problem of activity in psychology”. In: *Soviet psychology* 13.2 (1974), pp. 4–33.
- [17] Narges Mahyar and Melanie Tory. “Supporting Communication and Coordination in Collaborative Sense-making”. In: *Visualization and Computer Graphics, IEEE Transactions on* 20.12 (2013), pp. 1633–1642.
- [18] Judith S Olson et al. “How People Write Together Now: Beginning the Investigation with Advanced Undergraduates in a Project Course”. In: *ACM Transactions on Computer-Human Interaction (TOCHI)* 24.4 (2017), pp. 1–40. DOI: [10.1145/3038919](https://doi.org/10.1145/3038919). URL: <http://dx.doi.org/10.1145/3038919>.
- [19] Kjeld Schmidt. “The Problem with ‘ Awareness ’”. In: *Computer Supported Cooperative Work (CSCW)*. 2002, pp. 285–298.
- [20] Gerry Stahl. *Group cognition: Computer support for building collaborative knowledge*. Mit Press Cambridge, MA, 2006.
- [21] John Stasko, Carsten Görg, and Robert Spence. “Jigsaw: supporting investigative analysis through interactive visualization”. In: *Information Visualization* 7.2 (2008), pp. 118–132. ISSN: 1473-8716. DOI: [10.1057/palgrave.ivs.9500180](https://doi.org/10.1057/palgrave.ivs.9500180). URL: <http://ivi.sagepub.com/lookup/doi/10.1057/palgrave.ivs.9500180>.
- [22] Colin Ware. *Information visualization: perception for design*. Elsevier, 2012.
- [23] Robert Boleslaw Zajonc. “Social facilitation”. In: *Science, New Series* 149.3681 (1965), pp. 269–274.