

Evaluating Improvement in Situation Awareness and Decision-Making Through Automation

(Poster)

Erin Zaroukian

*Computational and Information
Sciences Directorate*

Army Research Laboratory
Aberdeen Proving Ground, USA
erin.g.zaroukian.civ@mail.mil

Justine Caylor

*Computational and Information
Sciences Directorate*

Army Research Laboratory
Aberdeen Proving Ground, USA
justine.p.caylor.ctr@mail.mil

Michelle Vanni

*Computational and Information
Sciences Directorate*

Army Research Laboratory
Adelphi, MD, USA
michelle.t.vanni.civ@mail.mil

Sue Kase

*Computational and Information
Sciences Directorate*

Army Research Laboratory
Aberdeen Proving Ground, USA
sue.e.kase.civ@mail.mil

Abstract— Automated systems such as information extraction (IE) pipelines are designed to facilitate situation awareness by providing human decision makers with relevant information, but beyond the validity of the pipeline itself, designing the output of the pipeline for optimal human understanding should be a goal. This paper presents results comparing comprehension of text documents with and without markup from a (simulated) IE pipeline in a simulated intelligence task. While previous work suggests that markup hurts both objective and subjective measures of performance and preference, this paper uses hand-generated markup designed to be maximally accurate and task relevant, finding more favorable results. These results, however, still point toward the limitations of markup and the importance of the task it is intended to facilitate.

Keywords—information extraction, situation awareness, deductive reasoning, visual search, workload, usability, automation, decision making

I. INTRODUCTION

Good intelligence is critical for maintaining situation awareness (SA), but finding resources to tackle large amounts of intelligence documents is a widespread bottleneck. To address big data situations such as this, information extraction (IE) pipelines have been developed to automatically extract relevant information from large collections of resources, and the output of these pipelines, especially during development, is often presented to human analysts as markup on text input. Research on these pipelines typically focuses on precision and recall of the outputs, and while these measures are important for extracting correct information, the actual usefulness of this information to the end (human) user is often overlooked. Work that has focused on human users typically looks at user experience with highly specific, low-level features such as font size, color, and serifs, e.g., [1-4], or high-level metrics such as text coherence [5]. Little research is available on complex markup and its effect, both in terms of the information chosen for markup and the way in which it is marked up, on reading comprehension. In particular, while many different markup schemes are in use, there is a paucity of work comparing text with markup to text without markup to assess the value added

by markup.

This paper asks whether markup actually improves human comprehension of text documents. Participant performance is measured objectively as the accuracy and speed with which participants answer comprehension questions about the text, and it is measured subjectively through ratings of workload. Subjective ratings of preference are also collected. In a previous paper [6], we asked this question using a pre-existing IE pipeline [7-8] and found, somewhat surprisingly, that marked-up text leads to worse performance (lower accuracy, slower response times, and higher workload ratings, as well as lower preference ratings) than comparable text without markup. This paper presents a follow up experiment in which we attempt to stack the cards in favor of markup, using hand-generated markup that aims to be as accurate and task relevant as possible. Nonetheless, find little evidence that participants perform better with markup. We do, however, find suggestions that participants generally preferred to have this markup present. Further investigation may shed more light on why this pattern emerged and what it means for creating useful markup for improved SA.

II. PREVIOUS EXPERIMENT

A previous experiment [6] compared human comprehension of simple text documents with and without markup from an existing IE pipeline to determine whether this markup improves human comprehension of these text documents. Performance was measured objectively as the accuracy and speed with which participants answer comprehension questions about text scenarios describing a hypothetical adversary attack, and it was measured subjectively through participant ratings of workload. A subjective measure of preference was also collected.

A. Participants, Materials, and Procedure

This experiment was identical to the experiment that will be presented in Section III below with the following notable exceptions:

This research was supported in part by an appointment to the Student Research Participation Program at the Army Research Laboratory administered by the Oak Ridge Institute for Science and Education through an interagency agreement between U.S. Department of Energy and ARL.

All the [ORG **military**] [FAC **bases**] in [CPE **Perchland**] are heavily protected.
 There is no new information about Raven [CPE **Perchland**].
 Perchland is land locked.
 [PER **Locals**] in [CPE **Sharkland**] are being
 The [PER **Turtle**] <lost> [PER **his**] right eye
 The Bronco [ORG **group**] does not <attack>
 [PER **Members**] of the [ORG **Charger**], [ORG
 Bronco] [ORG **groups**] have experience with
 The shopping [FAC **malls**] in the [CPE **coalition**] [LOC **area**] are not well defended.
 Charger and Titan [ORG **group**] [PER **members**] have <entered> Perchland and [LOC
 Salmonland].
 The [PER **Panther**], [PER **Charger**], [ORG **Titan**], and Raven [ORG **groups**] prefer to
 <attack> in daylight.

Fig. 1. Excerpt from an ELICIT scenario showing markup with mouse-over information for “entered”.

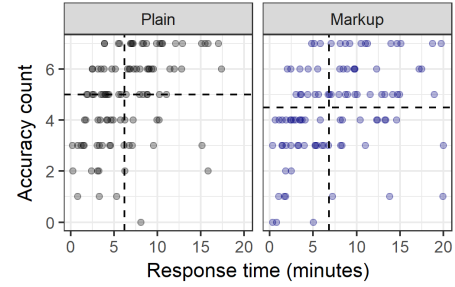


Fig. 2. Accuracy count (number of correctly answered questions, 0-7) versus response time in minutes for each participant in each condition. Medians are shown as dotted lines.

1. This first experiment treated the presence of text markup as a within-subjects manipulation for 100 participants. The second experiment, presented in this paper, treats the presence of text markup as a between-subjects manipulation with 100 participants per condition. This change was made to avoid asymmetric transfer seen between conditions in the first experiment [9].

2. The version of the workload survey used in the first experiment was modified to directly compare both version of the task (with/without markup). The second experiment, with a between-subject design, used the original workload survey as written.

3. The first experiment did not include a trust in automation survey. The second experiment does.

4. The markup in the first experiment was generated using an IE pipeline developed at Rensselaer Polytechnic Institute [7-8] which highlights a variety of entities (e.g., person, vehicle, geo-political entity) and where mouse-over reveals additional information (e.g., a relation’s arguments, the class an entity belongs to). See Fig. 1 for an example of text marked up through this IE pipeline. The markup in the second experiment was created by hand to be as accurate and relevant as possible.

5. The first experiment drew from four separate scenarios. In the second experiment, only two separate scenarios were used. This decision was made for convenience.

For more detail on the materials and procedure for this experiment, see [6] or Section III below.

B. Results

Participants’ accuracy and response times are shown for plain and markup trials separately in Fig. 2. While this markup was intended to improve performance, these results point to a small advantage for text without markup over text with markup.

Participants answered significantly more questions correctly in the Plain condition than in the Markup condition, and they responded significantly faster in the Plain condition than in the Markup condition. Additionally, participants overall associated higher workload with the markup trials and showed a preference for plain trials.

The poorer scores, both objective and subjective, seen on markup trials in this experiment suggest that markup is detrimental to performance in this task. The markup used here, however, is far from perfect and should by no means signal an

indictment of all forms of markup for all tasks. To better understand how markup might be of use, the markup in the next experiment was designed to be optimal, i.e., as accurate and as relevant to the task as possible.¹ If participants perform better with such markup than without, further testing can be conducted to estimate thresholds in accuracy and relevance required for markup to be considered a helpful addition, as well as to explore the role or markup density, the number of distinct categories, the way categories are indicated within the text, etc. If participants do not perform better with this markup than without, more serious thought should be given to the use of such markup.

III. CURRENT EXPERIMENT

The current experiment tests whether participants show improvements in objective and subjective measures using maximally-accurate and maximally-relevant text markup over non-marked-up text in uncovering hypothetical adversarial attacks.

A. Participants

Two hundred participants were recruited through Amazon Mechanical Turk to take part in this experiment. Each participant was compensated \$2.00.

B. Materials

The experiment was prepared using the Ibex tool for running behavioral psycholinguistic experiments (<https://code.google.com/archive/p/webspr/>) and run online through Amazon Mechanical Turk. The markup used in this experiment was generated by hand by the first author and checked by the other authors, and it separately highlights phrases relevant to four types of answers (*Who*, *What*, *Where*, and *When*) participants are required to provide, as described in the next paragraph. While there are many ways to judge relevance, the researchers believe that highlighting all and only potential answers made the markup much more relevant than the markup in the first experiment without making it too computationally unrealistic or causing it to directly give away any answers. See Fig. 3 for an example of marked-up text. The markup in this experiment drops the bracketing and labeling used in the first experiment, as participants often commented that they found this distracting. Furthermore, the markup here is expressed through background color instead of font color, as

¹ See work such as [10-11] where participants perform better with sparser, higher quality markup.

The **Turquoise group** focuses on destroying **energy infrastructure**.

No attacks are being planned on **religious organizations** in **Sigmaland**.

The **target** is in a coalition country (**Muland**, **Xiland**, **Omicronland**, **Piland**, or **Sigmaland**).

An attack is being planned for **the first month of the year**.

The Rose group may be involved.

There is a lot of activity involving **the Rose group**.

Embassies in **Piland** were recently attacked and evidence of more attacks has been found.

There are reports that spent nuclear fuel is missing in **Muland**.

The largest **bank** in **Piland** has 4 machine gun emplacements on its roof.

No traces of members from **the Blue group** have been found in **Sigmaland**.

Fig. 3. Excerpt from an ELICIT scenario showing hand-generated markup designed to be as accurate and relevant as possible.

we felt this better allowed us to maintain four visually distinct categories (*Who*, *What*, *Where*, and *When*) without sacrificing the contrast between text and background color [3].

The text used in this experiment was drawn from ELICIT, the Experimental Laboratory for the Investigation of Collaboration, Information Sharing, and Trust [12]. ELICIT is a simulated intelligence task containing a number of hypothetical adversary attack scenarios, each a list of 68 simple sentences that together allow a reader to deduce the Who, What, When, and Where of an anticipated adversary attack. These questions are answered in this experiment through seven drop-down menus (When is broken down into separate menus for month, date, time of day, and am/pm), allowing possible accuracy score to range from 0 to 7. See Figs. 1 and 3 for example sentences from ELICIT scenarios.

This experiment included three questionnaires: a demographic questionnaire, a trust in automation questionnaire, and the NASA Task Load Index (NASA-TLX) [13] with an added preference question. The trust in automation questionnaire, proposed by the United States Air Force Research Laboratory [11,14-15], consisted of 12 questions that capture how the participants feel about automation. Each question is rated from 1-7 on level of agreement (1 representing “disagree” and 7 representing “agree”). These questions can be seen in Table II. The NASA-TLX asks participants to rate the task on a variety of workload measures, and a question was added asking participants, if they were to do the task again, whether they would prefer to do the task with or without markup. These questions can be seen in Table I. Participants responded to each question by choosing a point on a 21-point scale where 1 represents “very low” or “perfect”, and 21 represents “very high” or “failure”. The additional preference question asks participants to report which version of the task they prefer, where 1 indicates a strong preference for the plain version, and 21 indicates a strong preference for the markup version.

C. Procedure

At the beginning of the experiment, participants completed a demographic questionnaire and read a page of instructions explaining the experiment. Participants then completed an abbreviated practice scenario in each condition (Plain and Markup) in order to familiarize them with the scenario presentation and the method for answering questions, as well

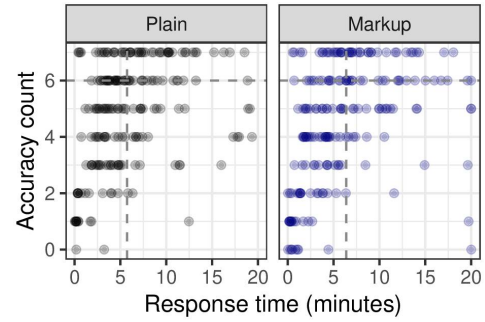


Fig. 4. Accuracy count (number of correctly answered questions, 0-7) versus response time in minutes for all participants, separated by each condition. Medians (with filtering criterion applied) are shown as dotted

as to ensure that each participant had some exposure to the markup (important for the trust in automation questionnaire). Participants then completed a trust in automation survey to gauge their baseline trust in “the system that automatically generated the text highlighting” that they saw (though, recall that this highlighting was not actually automatically generated in this experiment). At the end of the experiment, participants again completed the trust in automation questionnaire, as well as the workload and preference questionnaire. In addition, feedback was obtained from each participant on the strategies that they used to help them complete the scenarios, and they were given an opportunity to provide comments.

Each participant completed two different test scenarios, both either with markup (Markup condition) or without (Plain condition), determined randomly. Accuracy and response time were collected for each test scenario. Each test scenario was presented along with a countdown timer, and participants were cautioned that their responses would be submitted automatically if 20 minutes elapsed during that scenario.

D. Results

After concerns about quality of responses in the previous experiment, several filtering criteria were established based on response time and response quality (e.g., the coherence and relevance of their free-text strategy description). Unfortunately, this led to a loss of nearly half of all participants. In the results presented below, only the criteria based on response time, requiring participants to have spent at least two minutes on each test scenario, was maintained, as it was the least subjective.² This caused 50 participants to be removed from analysis. One additional participant was removed due to a technical failure, leaving 80 participants in the Plain condition and 69 participants in the Markup condition. All participants’ accuracy counts and responses times are shown in Fig. 4. All other measures include only the filtered 149 participants.

1) *Accuracy*: A Wilcoxon rank sum test found no significant difference in the number of correctly answered questions between conditions (Plain median = Markup median = 6, $W=10976$, $p=0.93$, $r=0.005$).

² The authors considered 2 minutes to be the minimum amount of time in which a participant could read and roughly comprehend a scenario, with upwards of 20 minutes needed to deduce the answers from the text.

2) *Response time*: An additional Wilcoxon rank sum test found no significant difference in response time between conditions (Plain median = 5.73, Markup median = 6.38, $W=12005$, $p=0.19$, $r=0.08$). While the effect size here is quite small, it suggests that with more power significantly faster response times may emerge in the Plain condition, as was seen in the first experiment (though, again, the effect in this experiment is quite small).

3) *Workload and preference*: Responses to the NASA-TLX and preference questions were compared using bootstrapped Kolmogorov-Smirnov tests, where the distributions of responses from participants in the Plain condition to were compared to the distributions of responses from participants in the Markup condition. Results, including median responses, are shown in Table I. Recall that 1 indicates “very low” (questions 1-3,5-6), “perfect” (question 4), or “definitely without highlighting” (question 7). After correcting for multiple comparisons³, only preference emerged as differing significantly between conditions, with responses skewing higher in the Markup condition, as shown in Fig. 6.⁴

TABLE I. WORKLOAD AND PREFERENCE

Question	Median response with Kolmogorov-Smirnov test	
	Plain	Markup
1. How mentally demanding was the task?	18	18
	$D=0.08$, $p=0.772$	
2. How physically demanding was the task?	2	2
	$D=0.07$, $p=0.819$	
3. How hurried or rushed was the pace of the task?	11	11
	$D=0.09$, $p=0.712$	
4. How successful were you in accomplishing what you were asked to do?	9	11
	$D=0.20$, $p=0.047$	
5. How hard did you have to work to accomplish your level of performance?	18	17
	$D=0.10$, $p=0.603$	
6. How insecure, discouraged, irritated, stressed, and annoyed were you?	9	8
	$D=0.15$, $p=0.232$	
7. If given the choice, which version of the task would you prefer to work with?	11.5	20
	$D=0.32$, $p<0.001$	

While there are no clear differences among objective measures of performance, the subjective measures indicate an overall preference for markup. This contrasts with the first experiment, where all advantages were in favor of plain trials.

4) *Trust in automation*: Table II shows the median responses to the trust in automation questionnaire that was given before and after the test scenarios. Recall that in this questionnaire the markup seen in training (and for some participants in test as well) was described as having been

³ The Holm-Bonferroni method was used to correct for multiple comparisons, yielding the following significance levels for questions 1-7 in Table 1: 1) 0.05/2, 2) 0.05/3, 3) 0.05/3, 4) 0.05/6, 5) 0.05/4, 6) 0.05/5, 7) 0.05/7.

⁴ The same pattern of results were found using standard t-tests as well as bootstrapped Anderson-Darling tests.

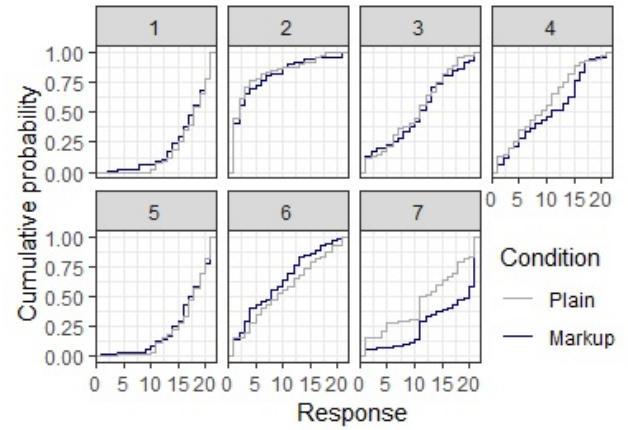


Fig. 6. Cumulative distribution of responses in the Plain vs the Markup condition for workload and preference as enumerated in Table I.

automatically computer generated, and participants were asked to rate their trust in the system that generated this markup. Overall, response scores are similar across conditions (Plain vs. Markup). A Pearson’s Chi-squared test is used to check for significance between questionnaire responses for the Before and After cases of each condition (Plain, Markup). Results indicate that for both conditions, there is no statistically significant difference between Before and After.

TABLE II. TRUST IN AUTOMATION QUESTIONNAIRE RESULTS

Trust in Automation Questionnaire	Median response			
	Plain		Markup	
	Before	After	Before	After
1. The system is deceptive	2	2	2	2
2. The system behaves in an underhanded manner	2	2	2	2
3. I am suspicious of the system’s intent, action, or outputs	2	2	2	2
4. I am wary of the system	3	3	2	3
5. The system’s actions will have a harmful or injurious outcome	2	2	2	2
6. I am confident in the system	5	4	5	5
7. The system provides security	4.5	4	5	4
8. The system has integrity	5	4	4	5
9. The system is dependable	5	4	5	5
10. The system is reliable	5	4	5	5
11. I can trust the system	5	4	5	5
12. I am familiar with the system	5	4.5	5	5

5) *Strategies*: Several strategies were utilized by the participants in each condition to help complete the experimental tasks. Participants in the Markup condition generally used the highlighted text to help them determine the relevant information, used the process of elimination, and

focused on easier questions to answer (typically *When* and *Where*). While some participants used similar strategies in the Plain condition, it appears that some faced more difficulties and had to adopt different strategies; without markup, helpful strategies included taking notes and using the Find command. While purely qualitative, this suggests that the markup assists and guides participants in an effective way that decreases the workload of completing this task, and without markup, participants will often still seek to use automation (e.g., automated search) to assist them.

IV. DISCUSSION

After participants showed better performance without markup than with markup in a previous experiment, the experiment presented here attempted to create a best-case scenario for markup and show that the current path taken by many IE researchers holds promise for improving SA. The current experiment resulted in no clear preference for text without markup over text with markup, but the only advantage for markup was seen in preference, and with more power an objective advantage for performance without markup may emerge. Still, other qualitative signs emerged that participants prefer to use markup, included reported strategies that participants actively used the markup to solve scenarios when it was available to them. While this success for markup is very modest, it suggests further changes that may lead to improved performance and SA, and even if participants only ever show a subjective preference for markup without any improvement in performance, this may be enough to justify its existence.

In this paper, a preliminary look at recently collected data was presented, and there is much further analysis that can be done. Notably, it may be informative to relate trust in automation scores to preference scores. The researchers expect to find that participants who indicated that they would prefer to perform the task in the Plain condition will show lower trust in automation than their counterparts who indicated that they would prefer to perform the task in the Markup condition. Additionally, evaluation to explore the role that demographic information (especially occupation) plays in performance and opinion will be conducted in future work. These results may point toward important points of choice and flexibility for automated systems.

There are a range of additional issues that can be addressed in future work and that weigh on the interpretation of the results presented here. Importantly, the markup in this experiment was designed to be as accurate and relevant as possible, but that does not mean it was actually optimal. For example, while not as computationally plausible, all and only correct answers could have been highlighted within the text scenarios, providing arguably more task-relevant markup that would have led to better performance in the Markup condition. Additionally, while the researchers believe that our high contrast highlighting was an improvement over the lower contrast font color and noisy bracketing used in the first experiment (cf. Fig. 1 and 3), it cannot be asserted that there is no more advantageous way to present the text. Further steps

may include tweaking the text presentation as well as our view of task relevance (both by varying the task and the markup) to explore their impact on performance.

ACKNOWLEDGMENT

Many thanks to Stephen Tratz, Claire Bonial, Jeffrey Micher, Clare Voss, Jon Bakdash, Lucia Donatelli, and Jeff Hoyer for their assistance in designing, deploying, and interpreting this work.

REFERENCES

- [1] J. Nielsen, *Designing web usability: The practice of simplicity*, Indianapolis, IN: New Riders Publishing, 1999.
- [2] L. Rello, M. Pielot, M. Marcos, "Make it big!: The effect of font size and line spacing on online readability," *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, 2016, pp. 3637-3648.
- [3] T.R. Williams, "Guidelines for designing and evaluating the display of information on the Web," *Technical Communication*, 47(3) 2000, pp. 383-396.
- [4] L. Akhmadeeva, I. Tuhvatullin, and B. Veytsman, "Do serifs help in comprehension of printed text? An experiment with Cyrillic readers," *Vision Research*, vol 65, 2012, pp. 21-24.
- [5] P.W. Foltz, P.W. "Comprehension, Coherence and Strategies in Hypertext and Linear text," in *Hypertext and Cognition*. J.-F. Rouet, J.J. Levonen, A.P. Dillon, and R.J. Spiro, Eds. NJ: Lawrence Erlbaum Associates. 1996.
- [6] E. Zaroukian, "Information extraction for optimized human understanding and decision making," *Proceeding of the 23rd ICCRTS*, Pensacola, FL, November 2018.
- [7] Q. Li and H. Ji, "Incremental Joint Extraction of Entity Mentions and Relations," *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, June 2014, pp. 402-412.
- [8] Q. Li, H. Ji, and L. Huang, "Joint Event Extraction via Structured Prediction with Global Features," *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics*, August 2013, pp. 73-82.
- [9] E.C. Poulton and P.R. Freeman, "Unwanted asymmetrical transfer effects with balanced experimental designs," in *Psychological Bulletin*, 1966, pp. 1-8.
- [10] J. C. Jan, C. M. Chen, and P. H. Huang, "Enhancement of digital reading performance by using a novel web-based collaborative reading annotation system with two quality annotation filtering mechanisms," *International Journal of Human-Computer Studies*, vol 86, 2016, pp. 81-93.
- [11] A. Neigel, J. Caylor, S. Kase, M. Vanni, and J. Hoyer, "The Role of Trust and Automation in an Intelligence Analyst Decisional Guidance Paradigm," *Journal of Cognitive Engineering and Decision Making*, vol. 12, no. 4, 2018, pp. 239-247.
- [12] M. Ruddy, "ELICIT - The Experimental Laboratory for the Investigation of Collaboration, Information Sharing, and Trust," *Proceedings of the 12th ICCRTS*, Newport, RI, June 2007.
- [13] NASA: NASA Task Load Index (TLX), v. 1.0 Manual, 1986.
- [14] J.Y. Jian, A.M. Bisantz, and C.G. Drury, "Foundations for an empirically determined scale of trust in automated systems," *International Journal of Cognitive Ergonomics*, 4(1), 53-71 (2000).
- [15] S.E. Kase, M. Vanni, M., J. Caylor, and J. Hoyer, "Human-Assisted Machine Information Exploitation: a crowdsourced investigation of information-based problem solving," *Proceedings of SPIE Defense and Security: Next Generation Analyst V*, 2017, pp. 1-17.