

Brief Article

The Author

August 13, 2021

1 Methodology

1.1 Study Setup

We conducted a remotely moderated, within subjects user study with X number of participants who had an average of X number of years of programming experience and regularly used Kubernetes in their work. Participants completed a Kubernetes task using the kubectl CLI and the OpenShift console. The order in which participants used the tools was counter balanced to mitigate learning bias between sessions, with half of the participants using the kubectl CLI first and the other half using the OpenShift console first [1]. We employed a classic "Think Aloud" user evaluation method during the sessions to avoid the effects of reactivity [2]. Video and audio from the sessions were recorded for later transcription and analysis.

1.2 Pilot Testing

We conducted pilot testing with two participants to identify any issues in our study design. One participant identified a misleading error in one of the task instructions so we rewrote it. Initially, we had only provided task hints for the kubectl CLI, and one participant realized that they could not apply these hints to the OpenShift console, so we also included task hints for the OpenShift console. We found that both participants were able to complete the task using each tool in under 30 minutes.

1.3 Study Task

The goal of the Kubernetes task is to create an application in a namespace, find the deployment pods that have the string "magic key" in their logs, delete the application, and then delete the namespace. This task was chosen for our study because it is simple enough to be accomplished by a Kubernetes novice in multiple steps or by an expert in fewer steps and it requires participants to perform search and navigation using each tool. Participants were given a GitHub repo with the application files and task instructions. If participants

got stuck during the session, hints specific to each tool interface were provided in links below the task instructions. Participants were told to think aloud during the session and given a prompt to continue thinking aloud if they fell silent for more than 15 seconds.

1.4 Data Collection

Before the study, participants filled out a background survey to provide information about their programming experience, use of Kubernetes, and familiarity with the two tools. After completing the study, participants filled out a follow up survey and ranked the tools based on their preference. In the follow up survey, participants were also asked to rate the tools based on their perceived cognitive load using the NASA TLX scale [1]. To give us a better understanding of preference, participants also rated the quality of each tool on a five-point Likert scale and were asked to give reasons for why they either liked or didn't like using each tool. Participants were also asked to share their recommendations for how each tool could be improved.

1.5 Data Analysis

The video recordings were analyzed to manually count how frequently participants used the keyboard to enter a command (k), how frequently participants clicked (c), and how frequently participants moved their mouse and/or scrolled to click (m). Participant input type frequencies were stored in tuples as (k, c, m). Rather than count characters, we chose to count the number of total commands entered to normalize the input frequencies as we are neither counting the duration of the clicks nor the length of the mouse movements. In the OpenShift console, we counted commands as the number of times participants typed something followed by an enter or click to enter. Mouse movements were only noted when they were followed by a purposeful click (for example, moving the mouse to click on the browser back button), so random mouse movements were not counted in our analysis. Specific input combinations were also noted such as how frequently participants copied and pasted text. Table X provides a list of all the input measures we collected for analysis and their exact definitions. Regular expressions were also used to create shortened descriptions of the patterns of input types used to complete the task with each tool. Table X provides a list of the regular expressions participants used to complete the task with each tool.

2 Results

We compared the tuples of total input type frequencies (k, c, m) to participants' perceived cognitive load ratings and tool preferences.

We compared the regular expressions of input types to participants' perceived cognitive load ratings to determine whether there was a correlation with patterns of task switching.