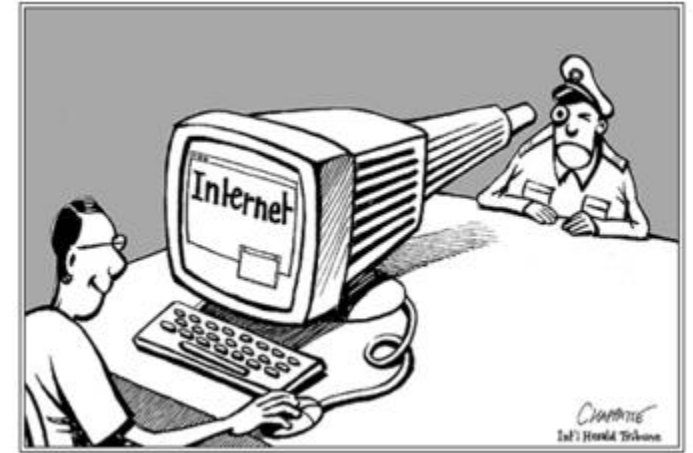


Insights from analysis of users' web browsing behavior

Yuliia Lut

Joint work with Rachel Cummings, Elizabeth Krizay, and Elissa Redmiles

Introduction



- People are constantly online. According to Pew Research Center
 - 28% of American adults report that they go online “almost constantly”,
 - roughly 8 out of 10 U.S. adults go online at least daily
- In 2017, Internet Service Providers in the US can collect, share, and sell sensitive consumer data without the user's consent.
- Analysis and modeling online browsing behavior play a key role in understanding users and technology interactions.

Introduction

- It is important to understand structural properties of the data
 - Amount of time people spend on browsing
 - Categories of websites that people browse more frequently
- Can users correctly perceive their browsing behavior?
- Can user's self reported data ensure an accurate analysis? E.g., can we use survey responses for analysis and modeling if we are unable to collect real data?

Literature review

There is correlation between user's browsing behavior and the user's type:

- Online behavior depends on demographics (gender, age, ethnicity) [KT'10]
- Educational background [GH'12]

The fact that the participants were observed changes their browsing behavior [LB'15], broad literature in behavioral economics.

People tend to overestimate the time they spend on a website and underestimate the number of visits on the website [KBLE'20]

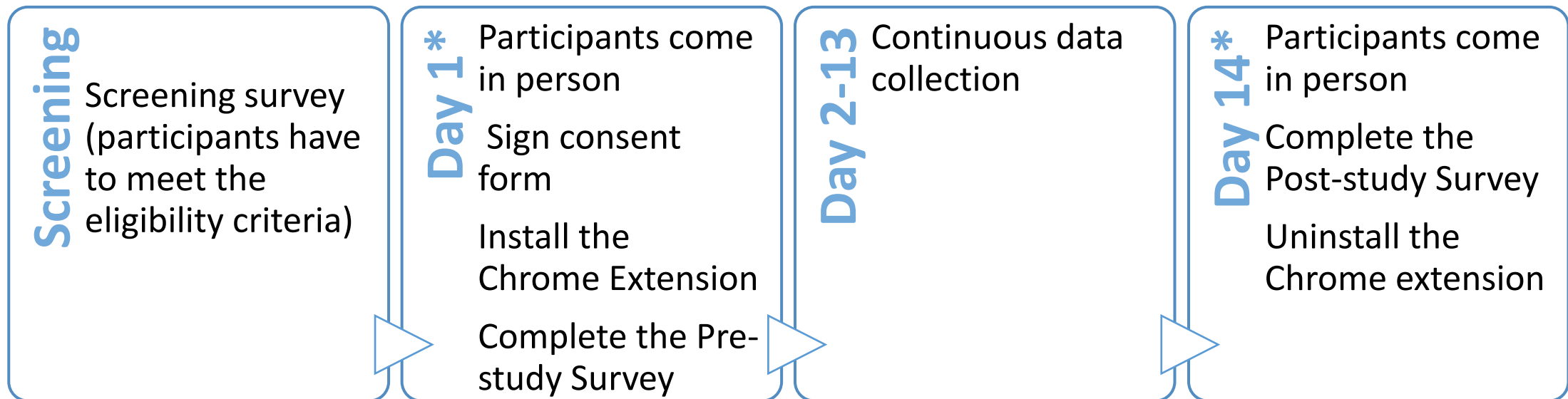
Main questions

1. Do people have correct perceptions of their behavior online?
2. Do people change browsing behavior if they are aware of being observed?
3. Does browsing behavior change in different settings (demographic groups, website categories)?
4. Can we learn structural properties of browsing patterns (e.g., that will enable realistic-looking synthetic data generation)?

Experiment design: Demographics

- 32 participants - Georgia Tech students
- Data was continuously collected for 14 days.
- We built a Chrome Extension (thanks to Michael Wang, GT'19) that enabled continuous collection of users' browsing data

Experiment outline



Experiment design: Demographics

- **Screening criteria:**

- Native English speakers
- More than 5 hours of browsing per week
- Older than 18

Metric (%)		
Gender	Female	40.6%
	Male	59.4%
Age	18-24	43.7%
	25-34	46.9%
	35-44	9.4%
Race	Asian	25%
	Black	31.2%
	White	28.1%

Experiment design: Surveys

Pre-study survey:

1. Demographics
2. Nationality
3. Hobbies
4. Average number of browsing hours per day
5. Favorite/most frequent categories of websites to visit.

Post-study survey:

1. How many hours per day browsing
2. Favorite/most frequent categories of websites to visit.
3. Questions about online privacy awareness.
4. During the study, did you change your browsing behavior to prevent information from being learned about you?

Experiment design: Privacy protection

- Participants were assigned random user ID and their browsing data is linked to ID not their name
- We truncated URLs that might be disclosive
 - facebook.com/groups/very_private_group → facebook.com
- Participants had an option to pause data collection any time
 - turn off/on the browsing extension
- Participants could opt out from the experiment at any time

Experiment design: website categorization

- Interested whether browsing behavior differs among this categories.
- Websites are categorized according to **Symantec WebPulse Site Review**
- The categorization includes categories and groups of categories.

Experiment design: website categorization

Category	Subcategory (examples)
Adult Related	Adult/Mature Content, Gore/Extreme
Liability Concerns	Piracy/Copyright Concerns, Violence/Intolerance
Security Threats	Malicious Outbound Data/Botnets, Phishing
Security Concerns	Compromised Sites, Hacking, Spam
File Transfer	File Storage/Sharing, Peer-to-Peer (P2P)
Society/Government	Charitable/Non-Profit, Government/Legal
Social Interaction	Personal Sites, Social Networking
Multimedia	Audio/Video Clips, Media Sharing
Communication	Email, Internet Telephony, Online Meetings
Health Related	Health, Restaurants/Food, Tobacco
Leisure	Art/Culture, Entertainment, Games
Commerce	Cryptocurrency, Job Search/Careers, Shopping
Technology	Cloud Infrastructure, Computer/Information Security
Information Related	Education, News, Reference, Search Engines/Portals

Experiment design: website categorization

Category	Subcategory (examples)
Adult Related	Adult/Mature Content, Gore/Extreme
Liability Concerns	Piracy/Copyright Concerns, Violence/Intolerance
Security Threats	Malicious Outbound Data/Botnets, Phishing
Security Concerns	Compromised Sites, Hacking, Spam
File Transfer	File Storage/Sharing, Peer-to-Peer (P2P)
Society/Government	Charitable/Non-Profit, Government/Legal
Social Interaction	Personal Sites, Social Networking
Multimedia	Audio/Video Clips, Media Sharing
Communication	Email, Internet Telephony, Online Meetings
Health Related	Health, Restaurants/Food, Tobacco
Leisure	Art/Culture, Entertainment, Games
Commerce	Cryptocurrency, Job Search/Careers, Shopping
Technology	Cloud Infrastructure, Computer/Information Security
Information Related	Education, News, Reference, Search Engines/Portals

facebook.com



Experiment design: website categorization

Category	Subcategory (examples)
Adult Related	Adult/Mature Content, Gore/Extreme
Liability Concerns	Piracy/Copyright Concerns, Violence/Intolerance
Security Threats	Malicious Outbound Data/Botnets, Phishing
Security Concerns	Compromised Sites, Hacking, Spam
File Transfer	File Storage/Sharing, Peer-to-Peer (P2P)
Society/Government	Charitable/Non-Profit, Government/Legal
Social Interaction	Personal Sites, Social Networking
Multimedia	Audio/Video Clips , Media Sharing
Communication	Email, Internet Telephony, Online Meetings
Health Related	Health, Restaurants/Food, Tobacco
Leisure	Art/Culture, Entertainment, Games
Commerce	Cryptocurrency, Job Search/Careers, Shopping
Technology	Cloud Infrastructure, Computer/Information Security
Information Related	Education, News, Reference, Search Engines/Portals

facebook.com

youtube.com

Experiment design: website categorization

Category	Subcategory (examples)
Adult Related	Adult/Mature Content, Gore/Extreme
Liability Concerns	Piracy/Copyright Concerns, Violence/Intolerance
Security Threats	Malicious Outbound Data/Botnets, Phishing
Security Concerns	Compromised Sites, Hacking, Spam
File Transfer	File Storage/Sharing, Peer-to-Peer (P2P)
Society/Government	Charitable/Non-Profit, Government/Legal
Social Interaction	Personal Sites, Social Networking
Multimedia	Audio/Video Clips , Media Sharing
Communication	Email , Internet Telephony, Online Meetings
Health Related	Health, Restaurants/Food, Tobacco
Leisure	Art/Culture, Entertainment, Games
Commerce	Cryptocurrency, Job Search/Careers, Shopping
Technology	Cloud Infrastructure, Computer/Information Security
Information Related	Education, News, Reference, Search Engines/Portals

facebook.com

youtube.com

Outlook.office365.com

Experiment design: website categorization

Category	Subcategory (examples)
Adult Related	Adult/Mature Content, Gore/Extreme
Liability Concerns	Piracy/Copyright Concerns, Violence/Intolerance
Security Threats	Malicious Outbound Data/Botnets, Phishing
Security Concerns	Compromised Sites, Hacking, Spam
File Transfer	File Storage/Sharing, Peer-to-Peer (P2P)
Society/Government	Charitable/Non-Profit, Government/Legal
Social Interaction	Personal Sites, Social Networking
Multimedia	Audio/Video Clips , Media Sharing
Communication	Email , Internet Telephony, Online Meetings
Health Related	Health, Restaurants/Food, Tobacco
Leisure	Art/Culture, Entertainment, Games
Commerce	Cryptocurrency, Job Search/Careers, Shopping
Technology	Cloud Infrastructure, Computer/Information Security
Information Related	Education, News, Reference, Search Engines/Portals

facebook.com

youtube.com

Outlook.office365.com

amazon.com

Experiment design: website categorization

Category	Subcategory (examples)
Adult Related	Adult/Mature Content, Gore/Extreme
Liability Concerns	Piracy/Copyright Concerns, Violence/Intolerance
Security Threats	Malicious Outbound Data/Botnets, Phishing
Security Concerns	Compromised Sites, Hacking, Spam
File Transfer	File Storage/Sharing, Peer-to-Peer (P2P)
Society/Government	Charitable/Non-Profit, Government/Legal
Social Interaction	Personal Sites, Social Networking
Multimedia	Audio/Video Clips , Media Sharing
Communication	Email , Internet Telephony, Online Meetings
Health Related	Health, Restaurants/Food, Tobacco
Leisure	Art/Culture, Entertainment, Games
Commerce	Cryptocurrency, Job Search/Careers, Shopping
Technology	Cloud Infrastructure, Computer/Information Security
Information Related	Education, News, Reference, Search Engines/Portals

facebook.com

youtube.com

Outlook.office365.com

amazon.com

google.com

Experiment design: Types of actions

awake	This action indicates that a user is online but they do not perform any action <ul style="list-style-type: none">• Appears if it is 5 min since the last action• Can happen because a user is not actively browsing (e.g., watching a movie or reading).
backButton	Clicking on the back button
click	Click that does not cause URL change <ul style="list-style-type: none">• 'Like' on social media
newTab	Opening new tab (manually or by opening a link in a new tab)
omnibox	Typing in omnibox – an address bar that can be also used as a search engine
tabChange	Altering between existing tabs
type	Typing
urlChange	Click that causes URL change <ul style="list-style-type: none">• Clicking on a product on shopping websites

Experiment design: User experiment and collected data

Event ID	Action type	Tab ID	User ID	Time Start	Time End	URL	subcategory	category	group
1	click	1	0034	9:00:00.00	9:00:00.00	amazon.com	Shopping	Commerce	Business related
2	click	1	0034	9:00:00.10	9:00:00.10	amazon.com	Shopping	Commerce	Business Related
3	tabChange	2	0034	9:00:01.00	9:00:01.00				
4	urlChange	2	0034	9:00:01.10	9:00:01.10	mail.google.com	Email	Communication	Non-productive
5	type	2	0034	9:00:05.00	9:00:20.00	mail.google.com	Email	Communication	Non-productive

- Action types: *awake, backButton, click, newTab, omnibox, tabChange, type, urlChange*
- Tab ID (for generating realistic looking browsing behavior)
- Timestamps

Q1. Do people have correct perceptions of their behavior online?

Q1: self perceptions of browsing (time)

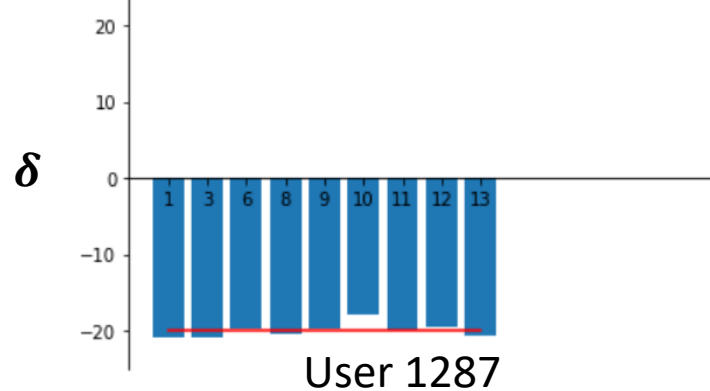
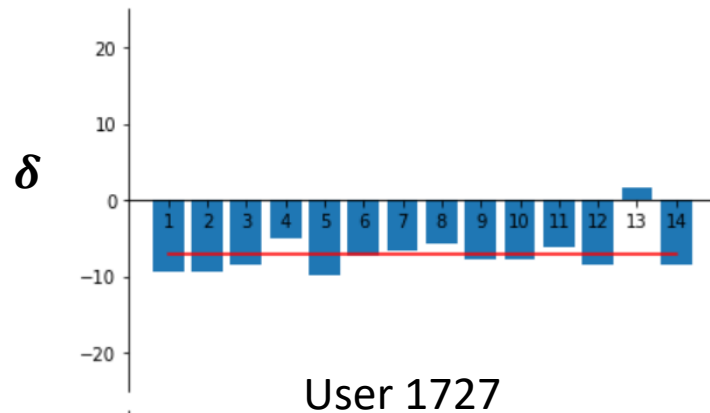
- In [KBLE'20] it is shown that most people overestimate how much time they spend online and underestimate the number of visits.
- We asked in the pre-survey how much time participants think they spend online. We want to compare these responses with real data.
- Most actions are instantaneous therefore we measure duration of browsing sessions instead of the duration of each action.
- A session ends if there are no actions for 5 min.
 - Previous work: 30 sec [KBLE'20], 30 min [AMGS'13]

$$\textit{Browsing time} = \textit{time}(S_1) + \textit{time}(S_2) + \dots + \textit{time}(S_N)$$

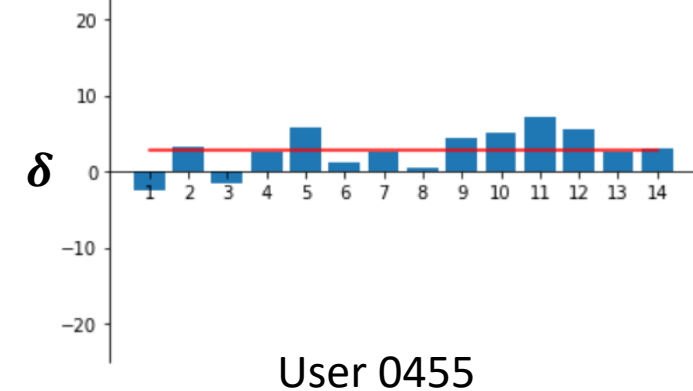
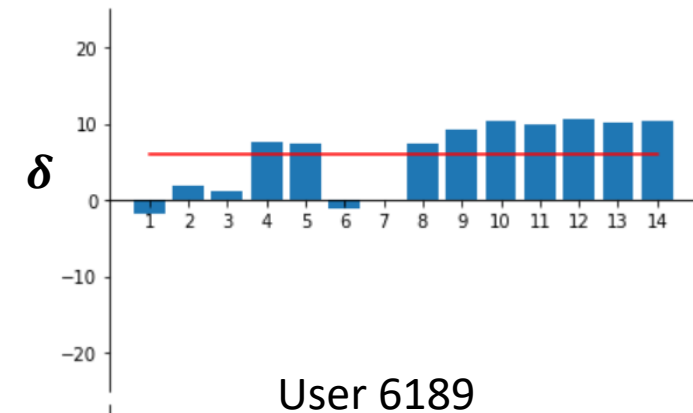
Q1: self perceptions of browsing (time)

δ = real time of browsing – pre-estimation

Overestimated time

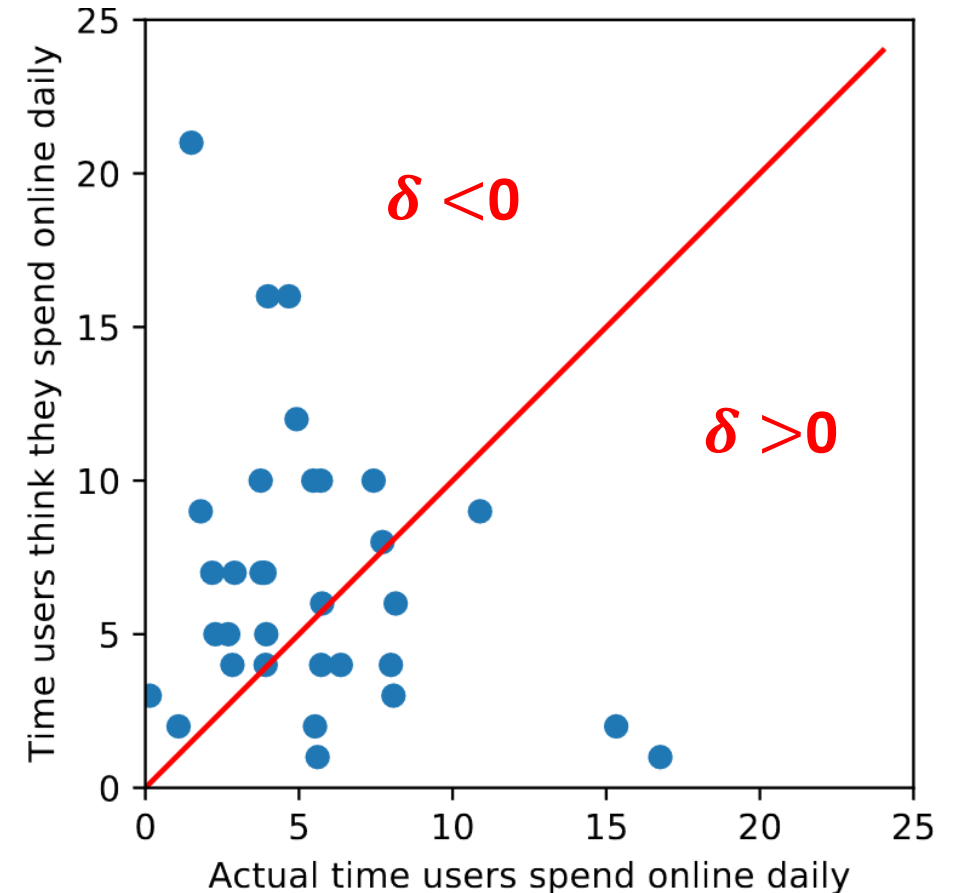


Underestimated time



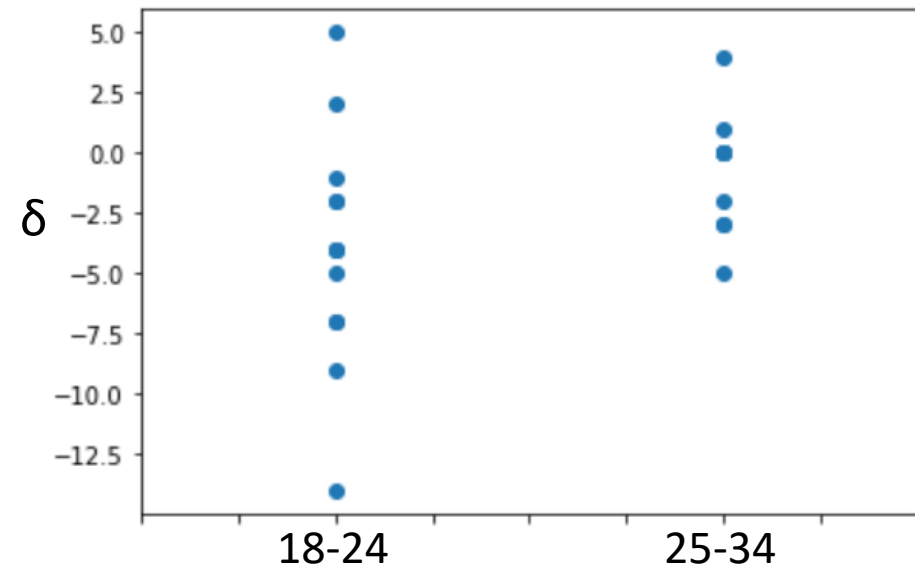
Q1: self perception on browsing (time)

- δ is significantly different from 0 (t-test, $p=0.00032$)
 - 22 out of 32 participants overestimated their browsing time



Q1: self perception on browsing (time)

- Dependence of δ on demographic features
 - Linear regression with gender and race as independent variables does not have significant coefficients ($p > 0.1$, $R^2 = 0.17$)
 - Distribution of δ does not differ among demographic features :
 - Gender: Mann-Whitney, $p = 0.29$
 - Age: Mann-Whitney, $p = 0.055$
 - Race: Mann-Whitney, $p > 0.2$



Q1: self perception on browsing (category)

- In Pre- and Post-study survey we asked participants what is ***their favorite/most frequent categories of websites to visit.***
- We want to check if their responses are supported by the collected data.

Mostly Underestimated

Shopping			
Actual data			
		YES	NO
Survey data	YES	16%	38%
	NO	15%	31%

Reference			
Actual data			
		YES	NO
Survey data	YES	9%	41%
	NO	3%	47%

Social Network			
Actual data			
		YES	NO
Survey data	YES	16%	56%
	NO	9%	19%

Entertainment			
Actual data			
		YES	NO
Survey data	YES	25%	59%
	NO	3%	13%

Mostly Overestimated

Business			
Actual data			
		YES	NO
Survey data	YES	34%	6%
	NO	22%	38%

Search			
Actual data			
		YES	NO
Survey data	YES	56%	28%
	NO	9%	9%

Q1: self perception on browsing

- We show that people mostly do not have a correct perception of their browsing.
 - Overestimate how much time they spend on browsing.
 - Wrong perception on what are their most frequently visited websites.

Q2. Do people change browsing behavior if they are aware of being observed?

Q2: Do people change browsing behavior if they are aware of being observed?

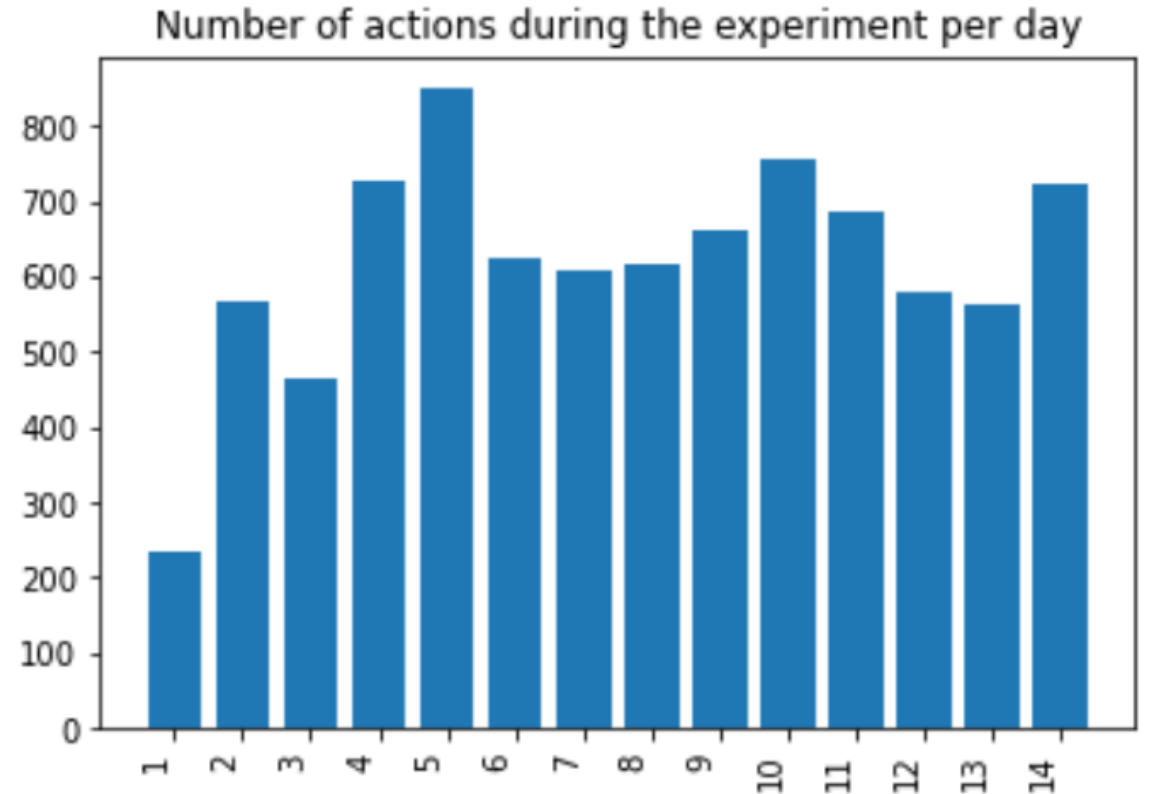
- One of the questions in Post-study survey:

“During the study, did you change your browsing behavior to prevent information from being learned about you?”

- Do people change browsing behavior over course of study because they aware they’re being observed?
- In literature, there are evidence that people’s behavior can depend on whether they are being tracked or not [LB15]

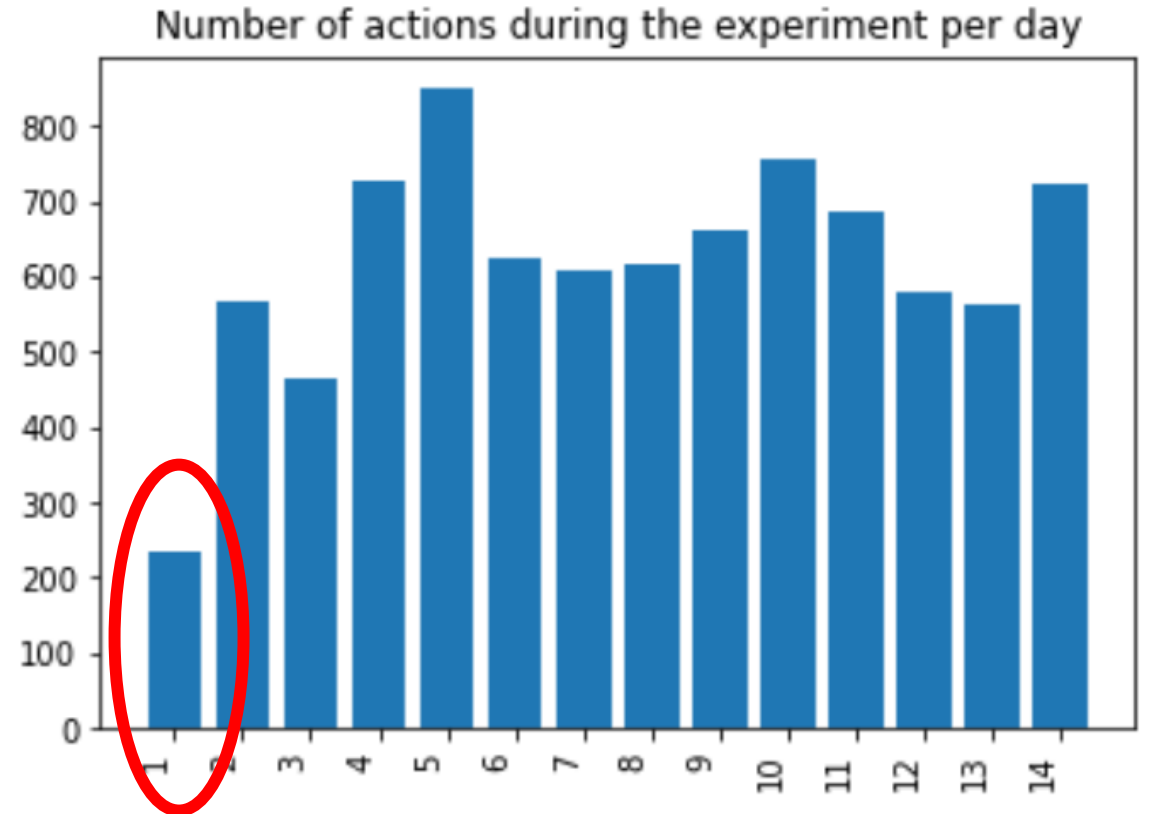
Q2: Do people change browsing behavior if they are aware of being observed?

- We want to check whether users' browsing changes with time.
- We show that a change is not significant:
 - T-test of equal mean, $p=0.37$
 - Levene test for equal variance, $p=0.18$



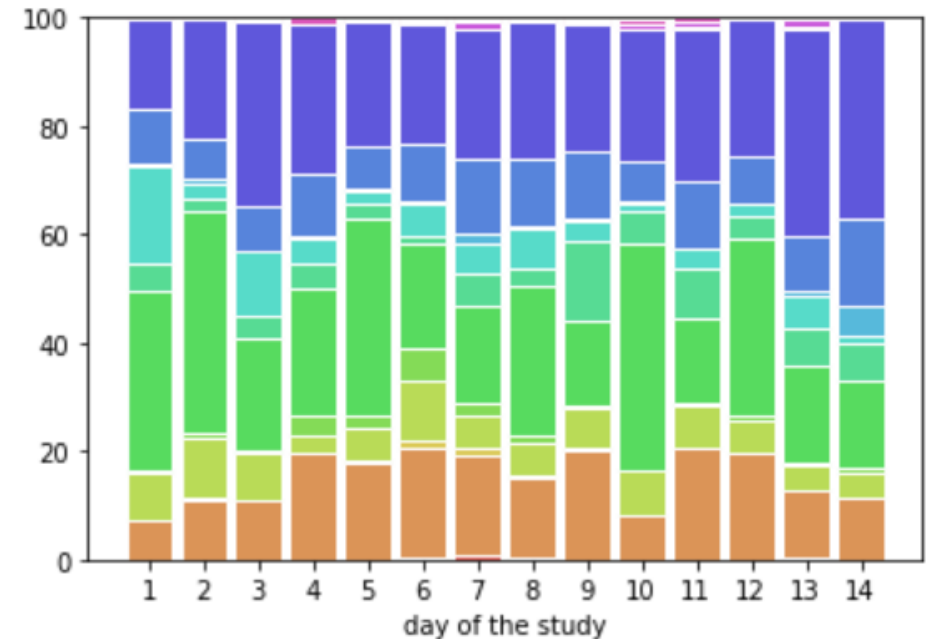
Q2: Do people change browsing behavior if they are aware of being observed?

- We want to check whether users' browsing changes with time.
- We show that a change is not significant:
 - T-test of equal mean, $p=0.37$
 - Levene test for equal variance, $p=0.18$



Q2: Do people change browsing behavior if they are aware of being observed?

- Consider a distribution of actions among categories of websites
- Does this distribution change over course of study?
 - No significant change($p>0.2$)



Q2: Do people change browsing behavior if they are aware of being observed?

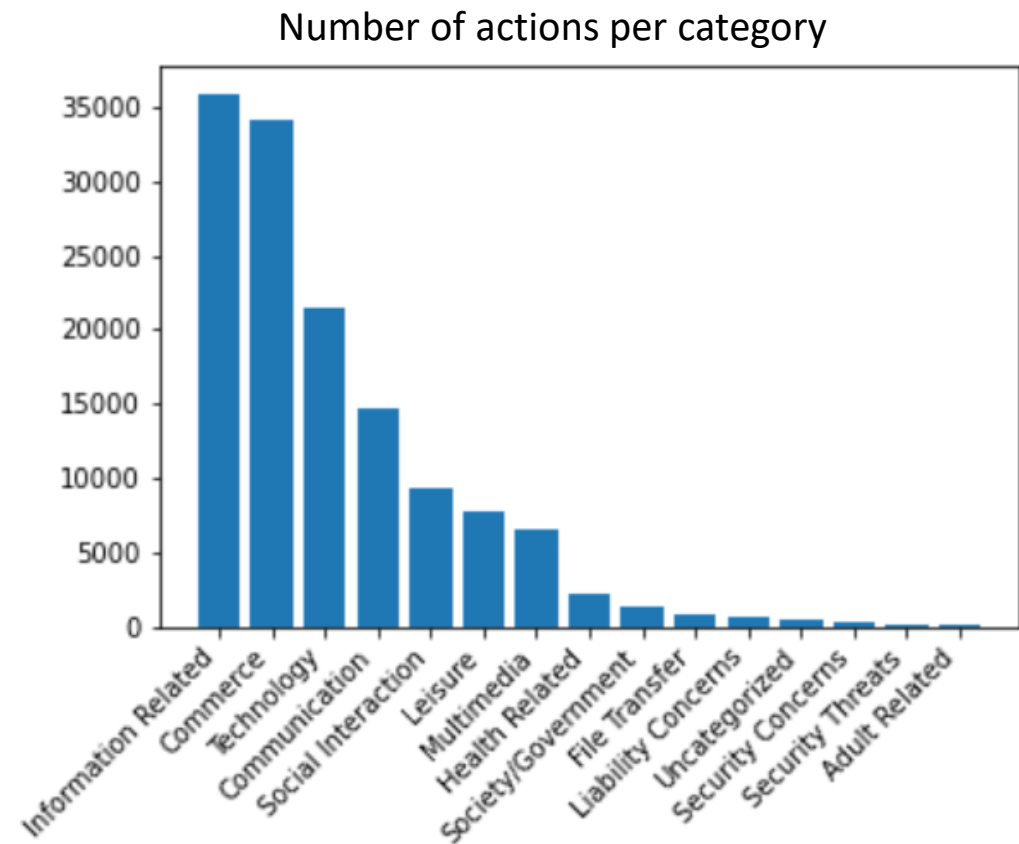
- We do not observe the change:
 - Participants were informed from the beginning that their browsing will be recorded, 14 days is not enough time to 'forget' about being observed.
 - Motivates for longer term of study.
- This result coincides with the Post-study survey: only 2 out of 32 participants claimed that they altered their behavior during the study.

Q3. Does browsing behavior change for different settings (demographic groups, website categories)?

Q3: Does browsing behavior change for different settings?

1. Does browsing behavior differ among demographic groups?
2. Does browsing behavior differ among website categories?

Metric (%)		
Gender	Female	40.6
	Male	59.4
Age	18-24	43.7
	25-34	46.9
	35-44	9.4
Race	Asian	25
	Black	31.2
	White	28.1



Q3: Does browsing behavior change for different settings (demographics)?

- According to χ^2 -test, distribution of activity among website categories is not significantly different among demographical features.

Category	P-value (gender)	P-value (age)	P-value (race)
Adult Related	0.06	0.18	0.71
Liability Concerns	0.62	0.09	0.74
Security Threats	0.12	0.01	0.82
Security Concerns	0.4	0.04	0.45
File Transfer	0.53	0.02	0.89
Society/Government	0.73	0.66	0.09
Social Interaction	0.45	0.35	0.86
Multimedia	0.08	0.54	0.49
Communication	N/A	N/A	N/A
Health Related	0.19	0.28	0.94
Leisure	0.23	0.3	0.77
Commerce	0.57	0.66	0.67
Technology	N/A	N/A	N/A
Information Related	0.57	0.66	0.9

Q3: Does browsing behavior change for different settings (demographics)?

- According to χ^2 -test, distribution of activity among website categories is not significantly different among demographical features.

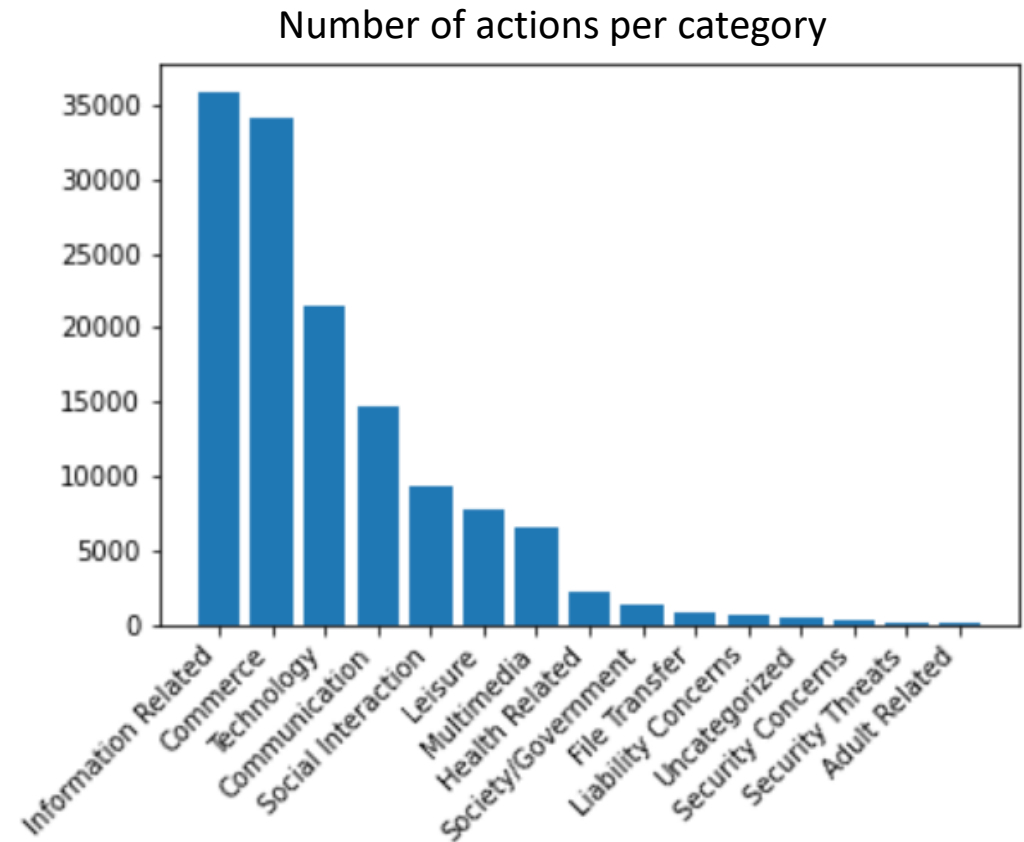
Category	P-value (gender)	P-value (age)	P-value (race)
Adult Related	0.06	0.18	0.71
Liability Concerns	0.62	0.09	0.74
Security Threats	0.12	0.01	0.82
Security Concerns	0.4	0.04	0.45
File Transfer	0.53	0.02	0.89
Society/Government	0.73	0.66	0.09
Social Interaction	0.45	0.35	0.86
Multimedia	0.08	0.54	0.49
Communication	N/A	N/A	N/A
Health Related	0.19	0.28	0.94
Leisure	0.23	0.3	0.77
Commerce	0.57	0.66	0.67
Technology	N/A	N/A	N/A
Information Related	0.57	0.66	0.9

Q3: Does browsing behavior change for different settings?

1. Does browsing behavior differ among demographic groups?

Metric (%)		
Gender	Female	40.6
	Male	59.4
Age	18-24	43.7
	25-34	46.9
	35-44	9.4
Race	Asian	25
	Black	31.2
	White	28.1

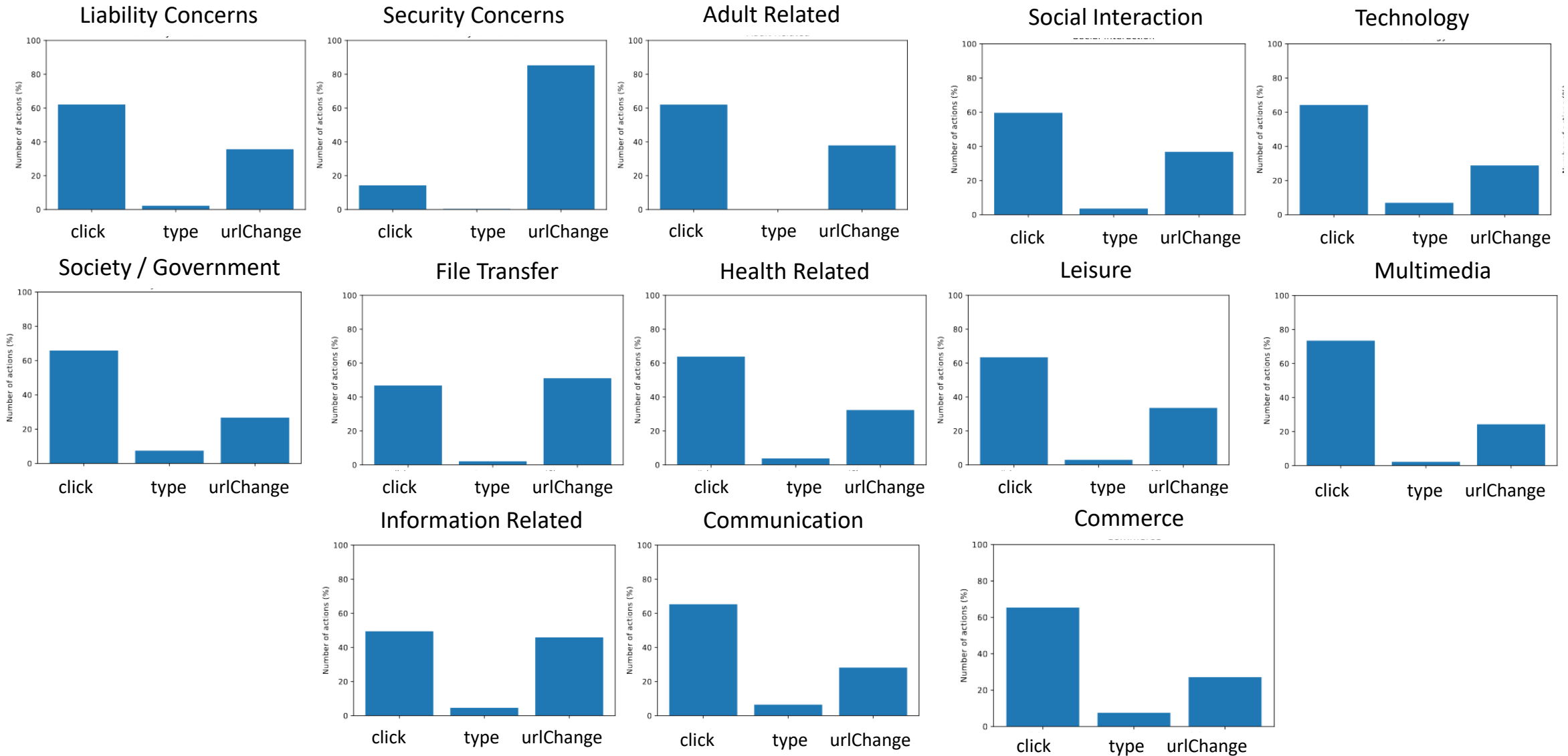
2. **Does browsing behavior differ among website categories?**



Q3: Does browsing behavior change for different settings?
(website category)

- Recall, we have different action types: *awake, backButton, click, newTab, omnibox, tabChange, type, urlChange*
- However, only *click, type* and *urlChange* describe actions that are performed directly on a website.

- Does the distribution of number of action per each of the three action type differ among website categories?



- Test if the distribution of activity differs among website categories.
- According to pairwise χ^2 -test, behavior on websites of some categories is **significantly different** from the rest.

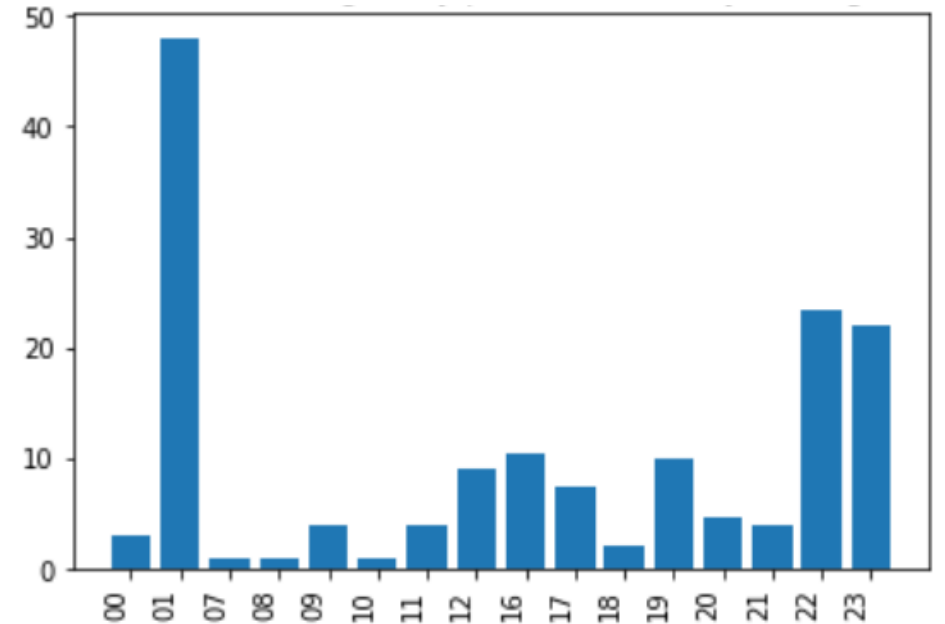
	Commerce	Technology	Information Related	Communication	Society/ Government	Social Interaction	Multi-media	Leisure	Health Related	File Transfer	Adult Related	Security Threats	Liability Concerns	Security Concerns
Commerce		0.93	0.02	0.86	0.96	0.21	0.12	0.21	0.41	0.001	0.01	0.01	0.08	$< 10^{-6}$
Technology			0.05	0.9	0.95	0.38	0.15	0.37	0.62	0.003	0.03	0.01	0.17	$< 10^{-6}$
Information Related				0.03	0.02	0.33	0.002	0.13	0.1	0.45	0.05	0.91	0.13	$< 10^{-6}$
Communication					0.95	0.4	0.2	0.45	0.72	0.003	0.04	0.01	0.22	$< 10^{-6}$
Society/ Government						0.26	0.19	0.26	0.53	0.001	0.02	0.01	0.12	$< 10^{-6}$
Social Interaction							0.1	0.84	0.78	0.11	0.25	0.16	0.7	$< 10^{-6}$
Multi-media								0.27	0.3	0.0005	0.08	0.0003	0.19	$< 10^{-6}$
Leisure									0.9	0.05	0.37	0.05	0.88	$< 10^{-6}$
Health Related										0.02	0.2	0.04	0.63	$< 10^{-6}$
File Transfer											0.09	0.51	0.1	$< 10^{-6}$
Adult Related												0.02	0.62	$< 10^{-6}$
Security Threats													0.06	$< 10^{-6}$
Liability Concerns														$< 10^{-6}$

- Test if the distribution of activity differs among website categories.
- According to pairwise χ^2 -test, behavior on websites of some categories is **significantly different** from the rest.

	Commerce	Technology	Information Related	Communication	Society/ Government	Social Interaction	Multi-media	Leisure	Health Related	File Transfer	Adult Related	Security Threats	Liability Concerns	Security Concerns
Commerce		0.93	0.02	0.86	0.96	0.21	0.12	0.21	0.41	0.001	0.01	0.01	0.08	$< 10^{-6}$
Technology			0.05	0.9	0.95	0.38	0.15	0.37	0.62	0.003	0.03	0.01	0.17	$< 10^{-6}$
Information Related				0.03	0.02	0.33	0.002	0.13	0.1	0.45	0.05	0.91	0.13	$< 10^{-6}$
Communication					0.95	0.4	0.2	0.45	0.72	0.003	0.04	0.01	0.22	$< 10^{-6}$
Society/ Government						0.26	0.19	0.26	0.53	0.001	0.02	0.01	0.12	$< 10^{-6}$
Social Interaction							0.1	0.84	0.78	0.11	0.25	0.16	0.7	$< 10^{-6}$
Multi-media								0.27	0.3	0.0005	0.08	0.0003	0.19	$< 10^{-6}$
Leisure									0.9	0.05	0.37	0.05	0.88	$< 10^{-6}$
Health Related										0.02	0.2	0.04	0.63	$< 10^{-6}$
File Transfer											0.09	0.51	0.1	$< 10^{-6}$
Adult Related												0.02	0.62	$< 10^{-6}$
Security Threats													0.06	$< 10^{-6}$
Liability Concerns														$< 10^{-6}$

Security Concerns

- In our database the following Security concerns websites were visited:
 - Suspicious (netflix.com)
 - Placeholders
 - Potentially Unwanted Software (www2.securybrowse.com)
 - Hacking (www.recoverlostpassword.com)
- Majority are Suspicious



Q3: Does browsing behavior differ among demographic groups?

- In our case, browsing behavior does not depend on demographic features. Possibly, a larger dataset needed.
- We show that browsing behavior depends on the type of the website.
 - This can help for an accurate modeling of the browsing behavior.
 - Can be used to identify malicious websites and security threats.

Q4. Can we learn structural properties of browsing patterns (e.g., that will enable realistic-looking synthetic data generation)?

Q4: Can we learn structural properties of browsing patterns?

- We are interested in capturing realistic properties of browsing behavior.
- This can be used for modeling browsing behavior such that it is indistinguishable from real data.
- Browsing data is of the form of the categorical sequences:

Session 1: *tabChange, tabChange, type, urlChange, newTab*

Session 2: *click, urlChange, newTab, click urlChange
tabChange*

Session 3: *tabChange tabChange tabChange newTab*

Types of actions:

- *awake*
- *backButton,*
- *click*
- *newTab,*
- *omnibox*
- *tabChange*
- *type*
- *urlChange*

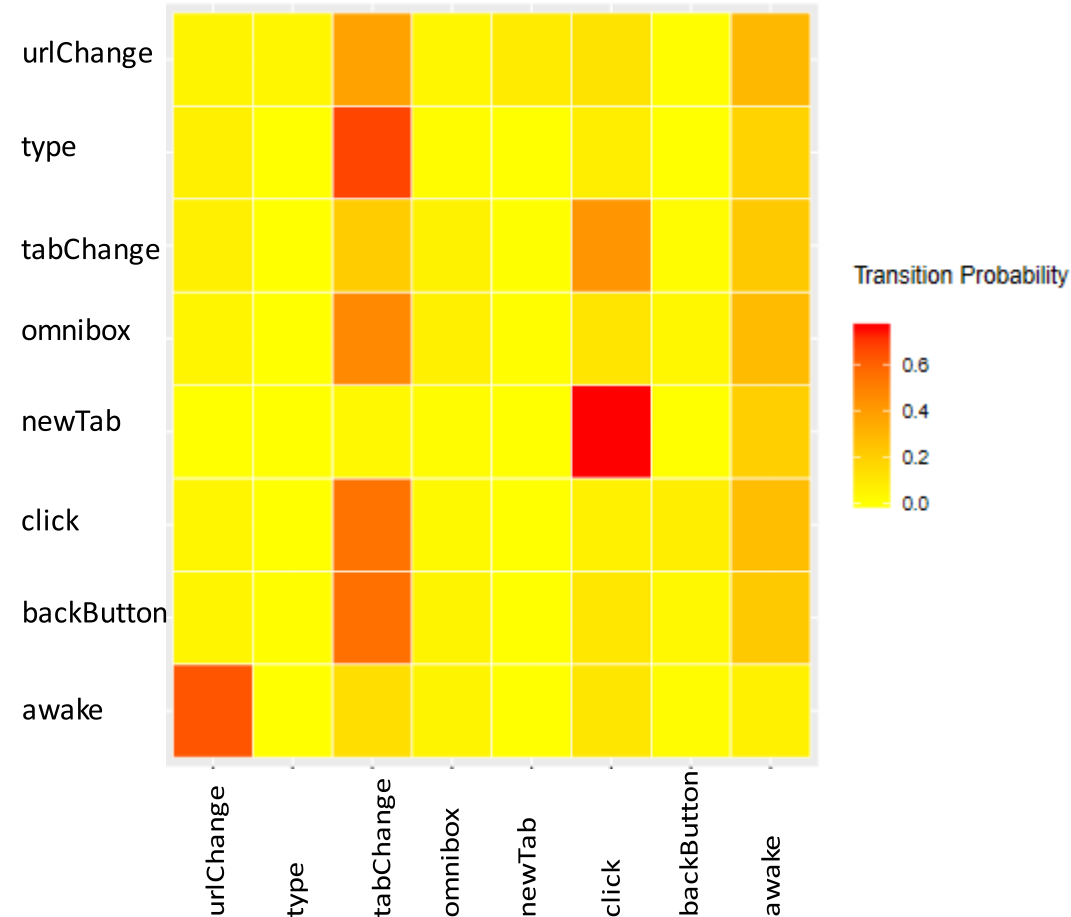
Q4: Can we learn structural properties of browsing patterns?

- Intuitively, the process has Markov property

Markov chain is a sequence $\{X_n\}_{n=1}^{\infty}$ such that

$$P(X_{n+1} = x | X_1 = x_1, \dots, X_n = x_n) \\ = P(X_{n+1} = x | X_n = x_n)$$

- For example:
 - *newTab* \rightarrow *tabChange*
 - *type* \rightarrow *click*



Q4: Can we learn structural properties of browsing patterns?

- Given all the findings, we fit different Markov models for:
 - Separate users
 - Demographic groups
 - Website categories
- Similarly to Q3, we want to compare Markov chains fitted using users' self reported data and using real data
- What features significantly impact accuracy of Markov models for browsing behavior?
- Future work: use obtained results for generating realistic synthetic data.

Conclusion

- We designed an experiment where browsing data of 32 students was collected continuously for 14 days.
- We performed a statistical analysis of the database that shows:
 - **Q1:** Mostly people do not have a correct perception of their browsing behavior (time, category)
 - **Q2:** Data does not show any evidence of the participants changing their browsing behavior.
 - **Q3:** Browsing behavior does not show significant difference among demographic groups, however we observe a significant difference in behavior for some website categories.
 - **Q4:** We suggested a future direction for applying our findings for synthetic data generation.

Insights from analysis of users' web browsing behavior

Yuliia Lut

Joint work with Rachel Cummings, Elizabeth Krizay, and Elissa Redmiles

References

- Cummings, R., Krehbiel, S., Mei, Y., Tuo, R., & Zhang, W. (2018). Differentially private change-point detection. In *Advances in Neural Information Processing Systems*(pp. 10825-10834).
- Dwork, C., McSherry, F., Nissim, K., & Smith, A. (2006, March). Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference* (pp. 265-284). Springer, Berlin, Heidelberg.
- Cynthia Dwork, Moni Naor, Toniann Pitassi, and Guy N. Rothblum. Differential privacy under continual observation. In 42nd ACM Symposium on Theory of Computing, STOC '10, 2010.
- Dwork, C., Roth, A. (2014). The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4), 211-407.
- BS Darkhovsky. A nonparametric method for the a posteriori detection of the disorder time of a sequence of independent random variables. *Theory of Probability & Its Applications*, 21(1):178{183, 1976.
- Henry B Mann and Donald R Whitney. On a test of whether one of two random variables is stochastically larger than the other. *The annals of mathematical statistics*, pages 50{60, 1947.
- Constantinos Daskalakis, Themis Gouleakis, Chistos Tzamos, and Manolis Zampetakis. Efficient statistics, in high dimensions, from truncated samples. In 2018 IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS), pages 639{649. IEEE, 2018.
- Sindhu Kiranmai Ernala, Moira Burke, Alex Leavitt, and Nicole B. Ellison. 2020.
- How Well Do People Report Time Spent on Facebook? An Evaluation of Established Survey Questions with Recommendations. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems (CHI '20)*.
- Abramson, Myriam and Gore, Shantanu. Associative Patterns of Web Browsing Behavior. AAI Fall Symposia, 2013
- Images: www.flaticon.com (articts: Freepic, monkik)