

A Large-Scale Characterization of How Readers Browse Wikipedia

Tiziano Piccardi, Martin Gerlach, Akhil Arora, Robert West



Summary - TL;TR

- 1 The **majority** of the incoming **traffic** originates from **search engines**.
- 2 The **majority of navigation** on Wikipedia is **on the surface**, with short sessions.
- 3 Access from **desktop** and **mobile** devices have a **different temporal pattern**.
- 4 Different **topics** are associated with **different navigation patterns**.
- 5 Navigation is **interrupted** on **low-quality** pages.

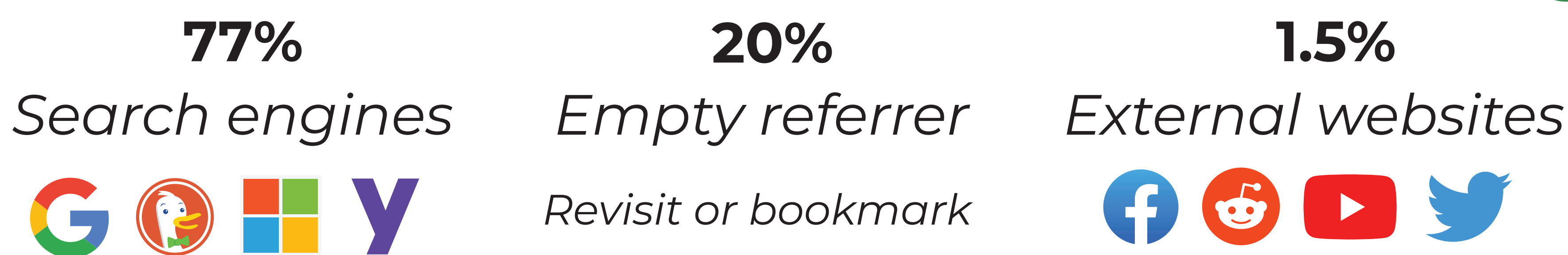
Server logs from English Wikipedia

One month 6.52B load events

User id	Timestamp	Device	Referrer	Article
d6n119fjl	2021-04-12 11:29:51	Desktop	bina.com	A
d6n119fjl	2021-04-12 11:31:26	Desktop	WP: A	B
d6n119fjl	2021-04-12 11:31:33	Desktop	WP: A	C
d6n119fjl	2021-04-12 11:36:16	Desktop	WP: C	D
d6n119fjl	2021-04-12 11:37:50	Desktop	facebook.com	E
...

Landing on Wikipedia

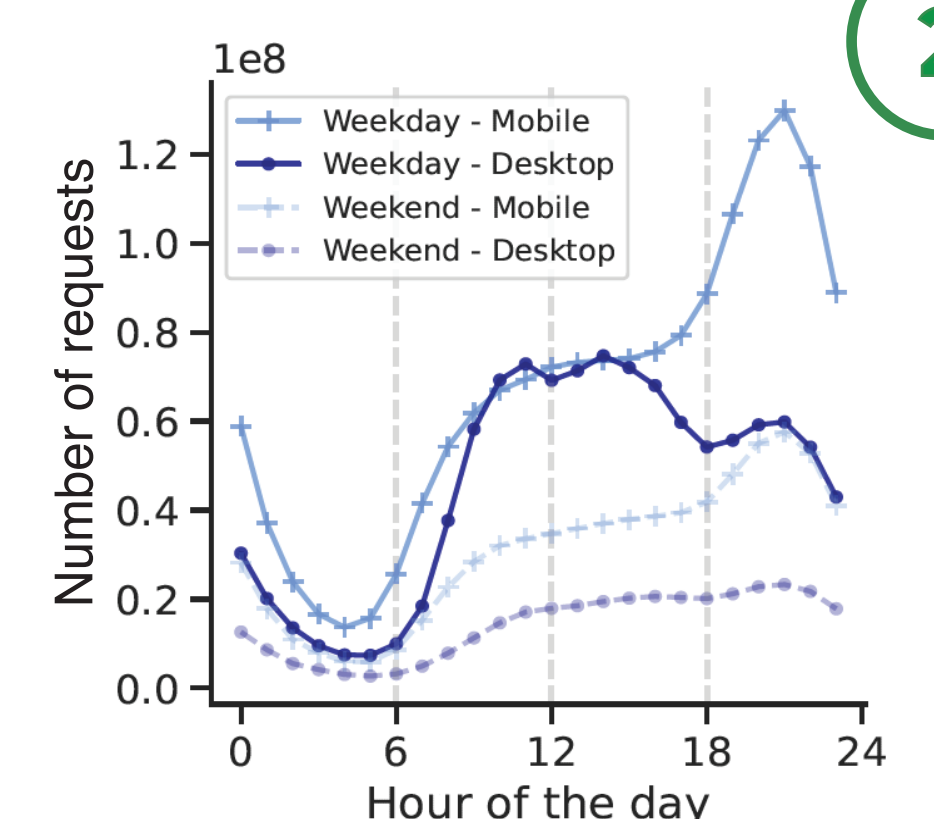
Incoming traffic origin



1

Daily pattern

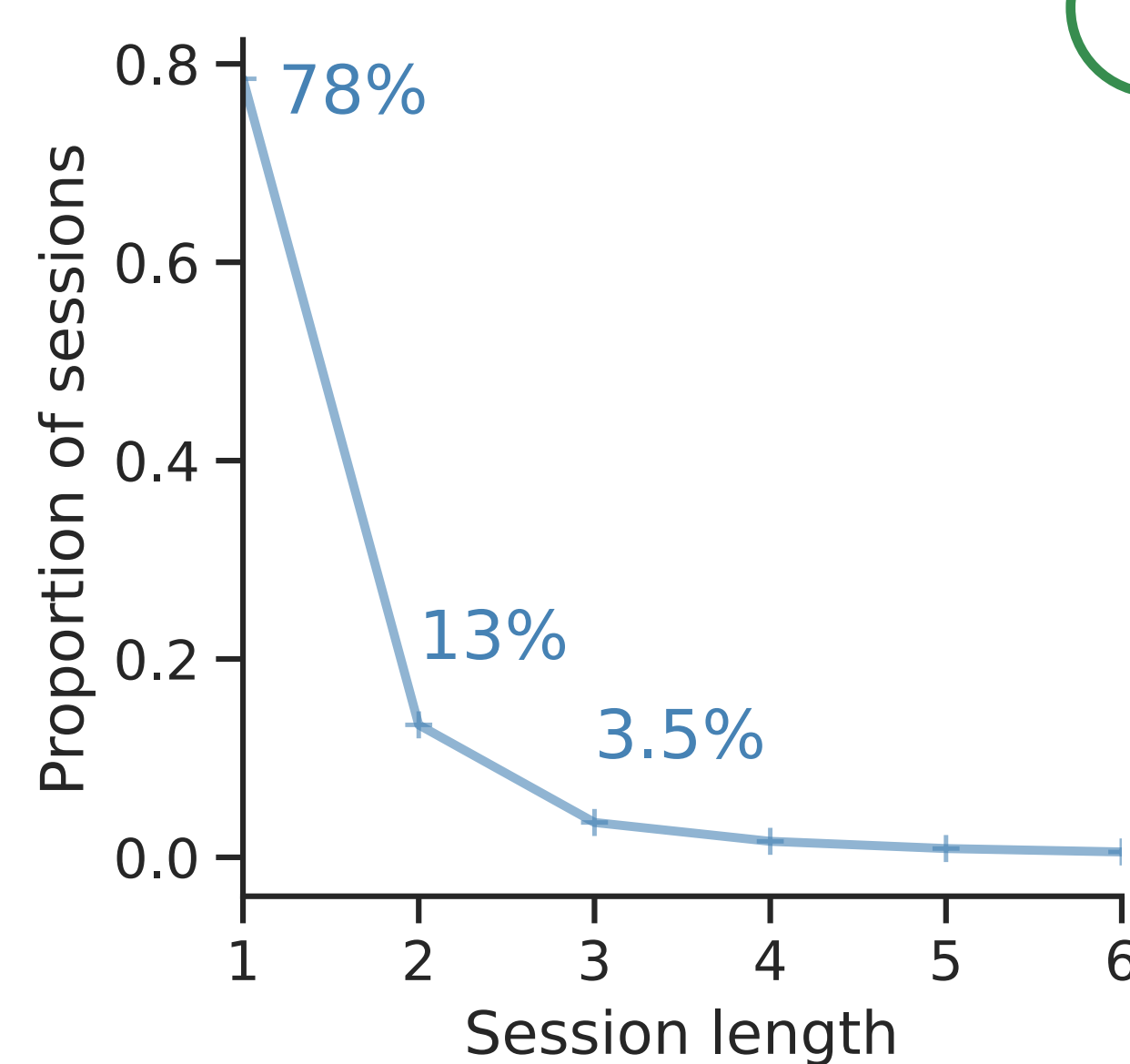
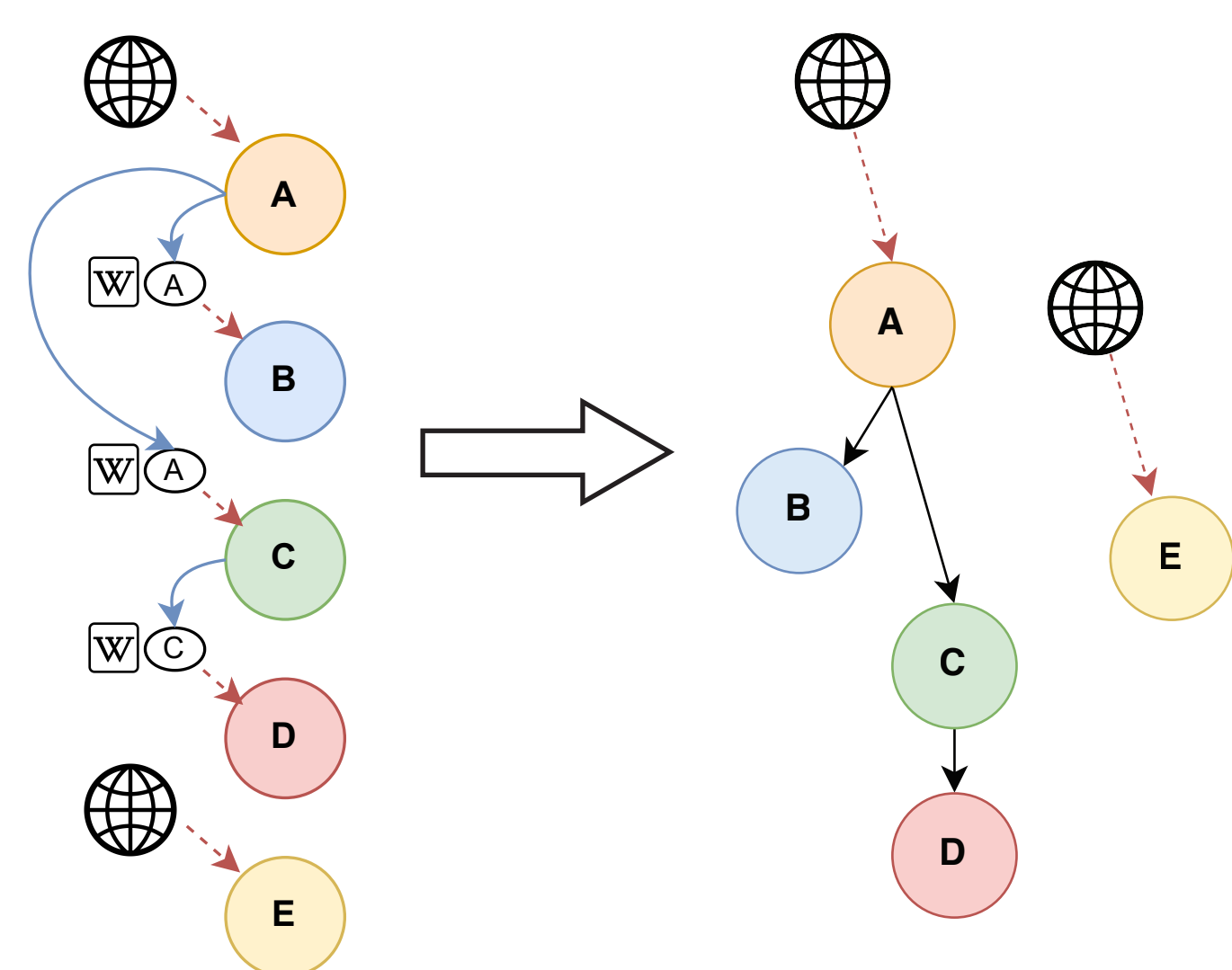
Differences between Desktop and Mobile



2

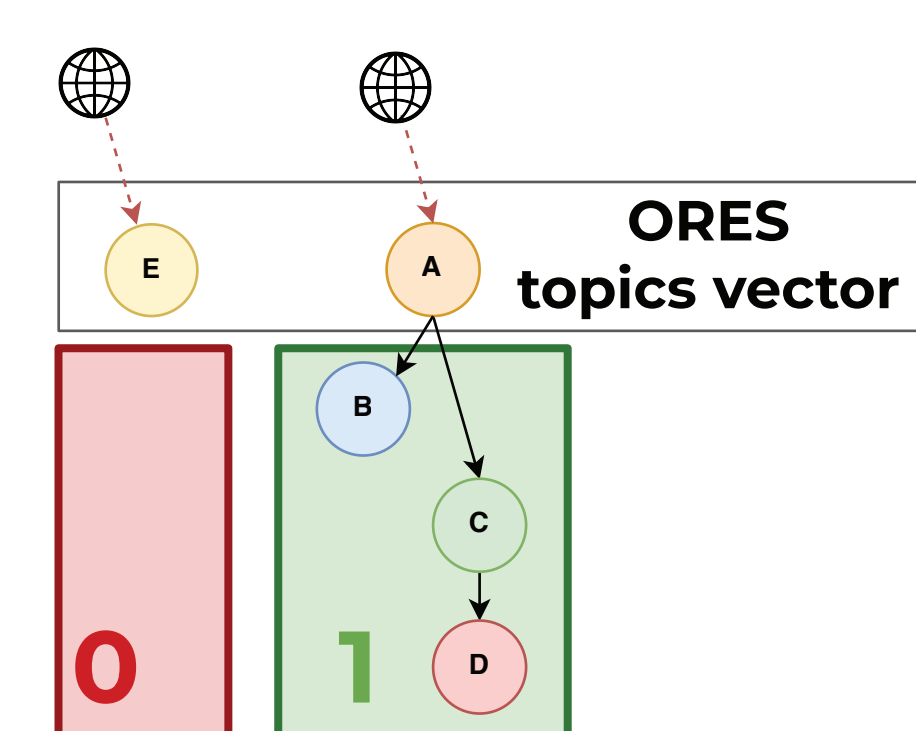
Beyond the first article

Aggregate the sessions

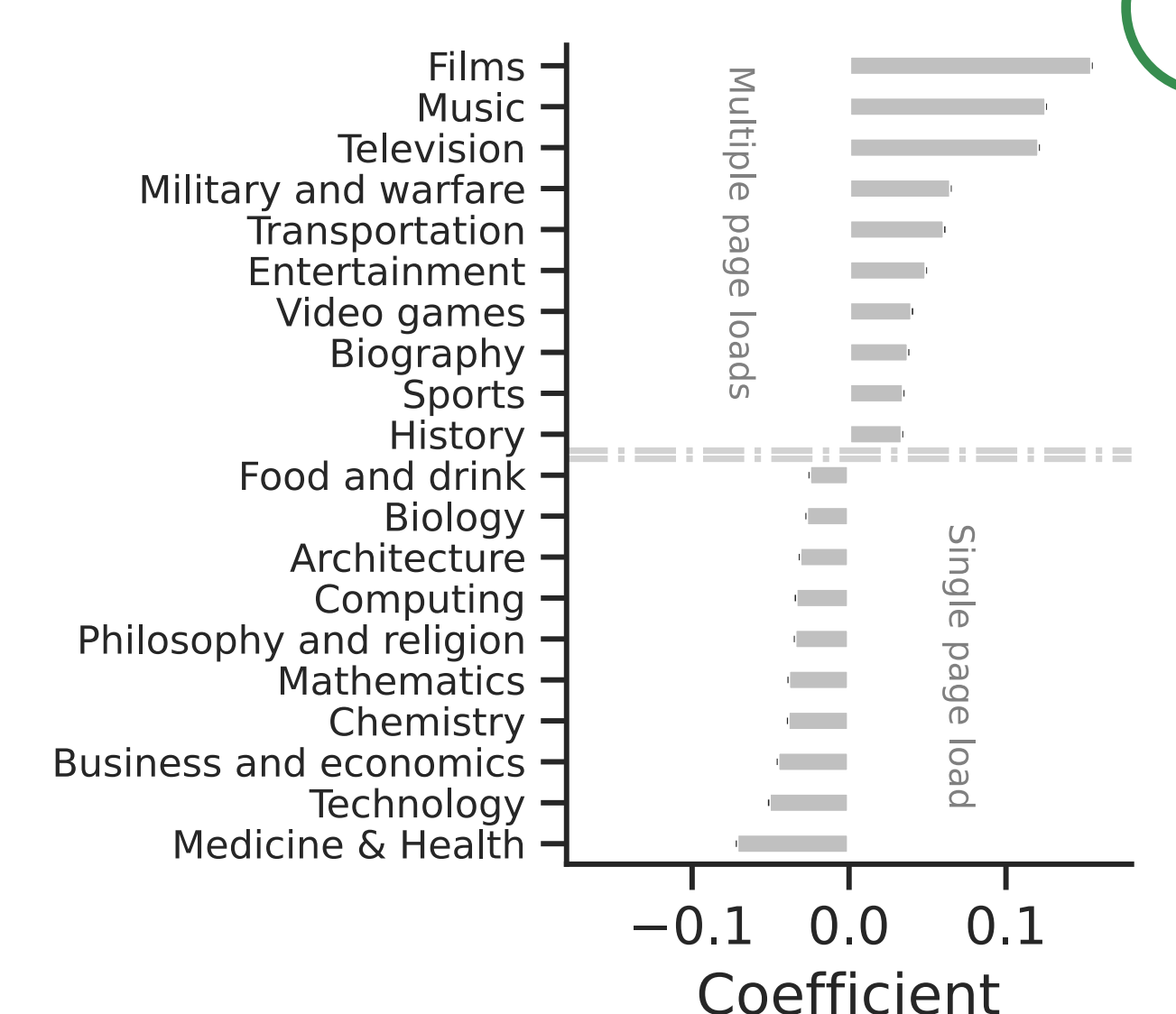


3

Topics

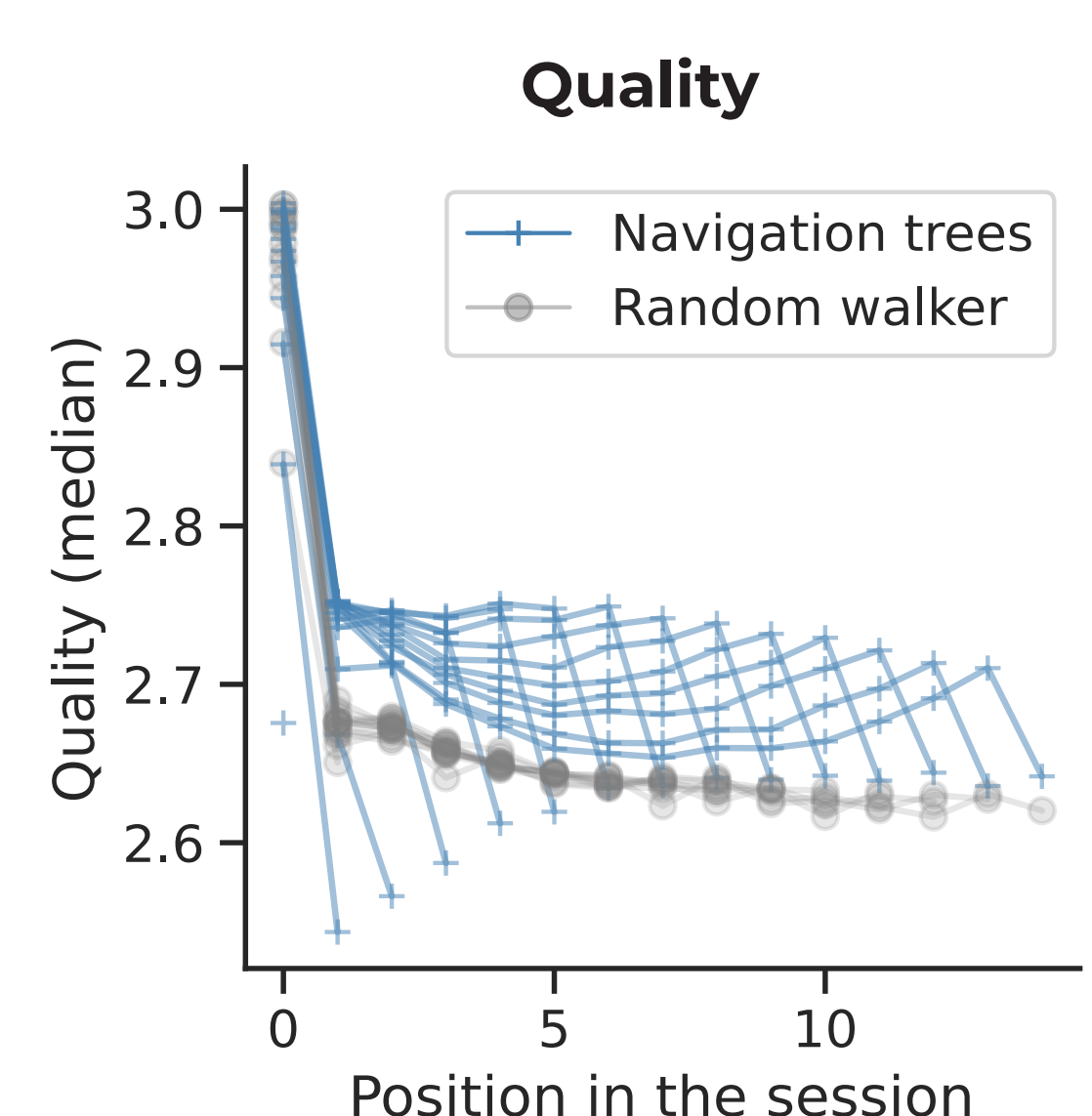
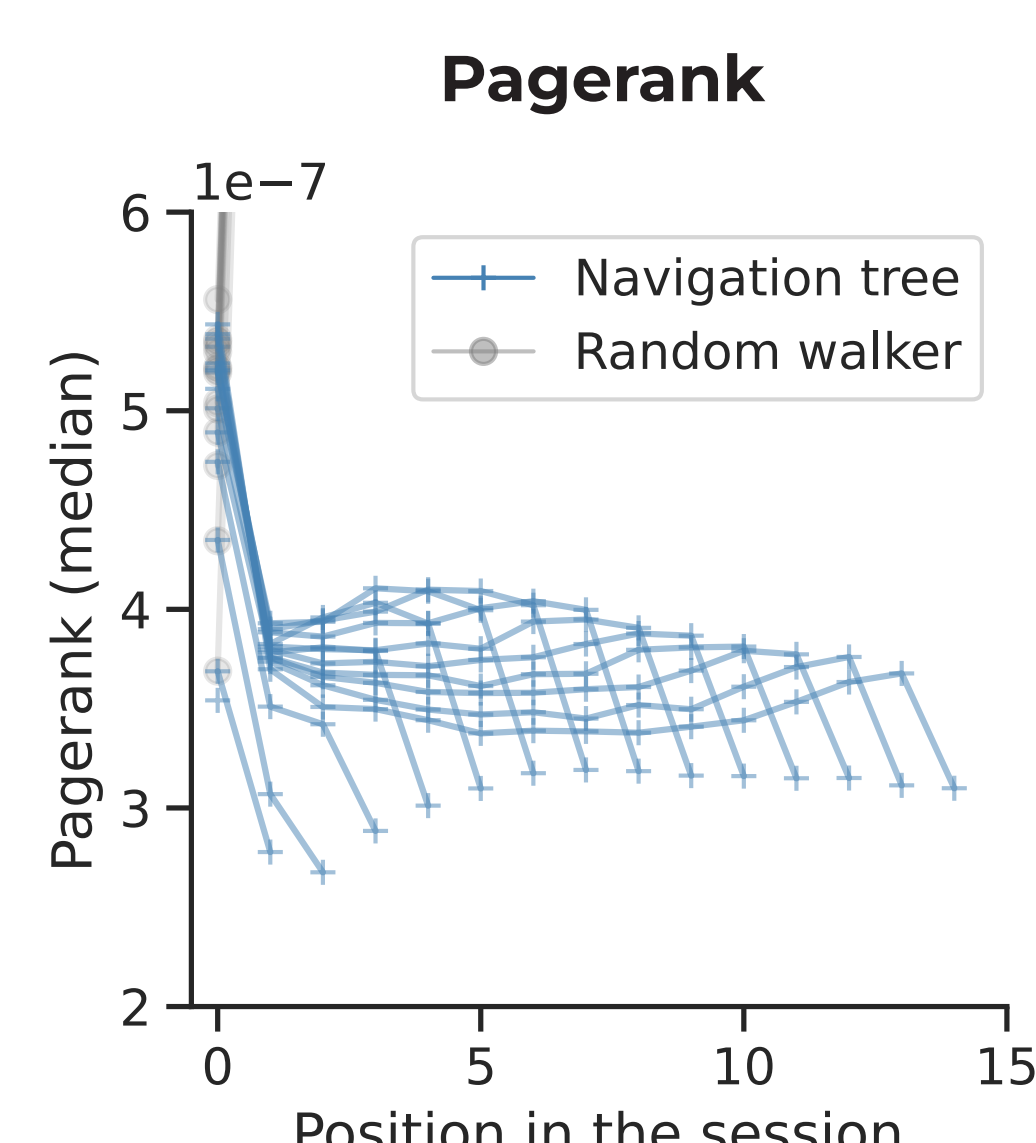
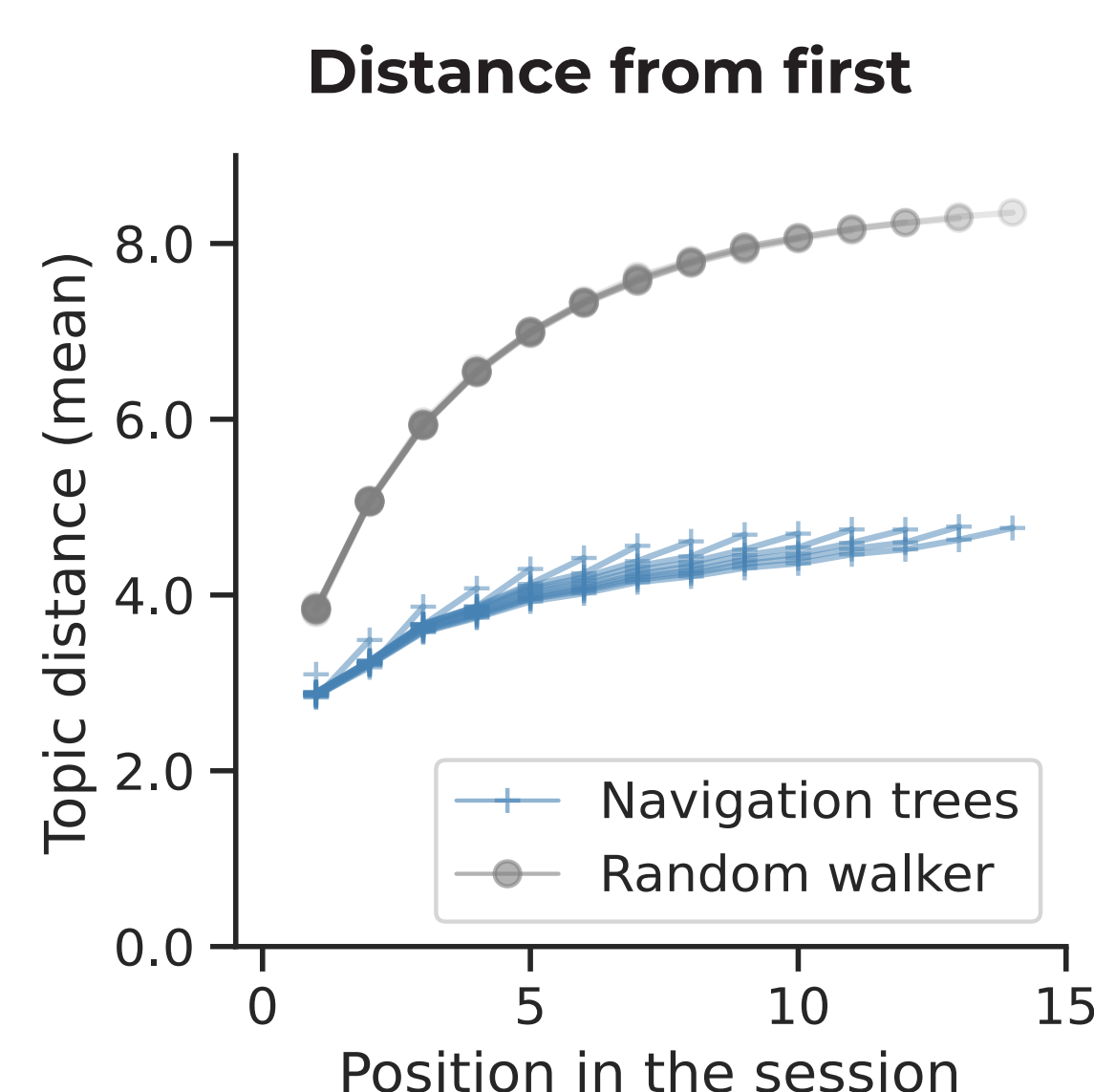
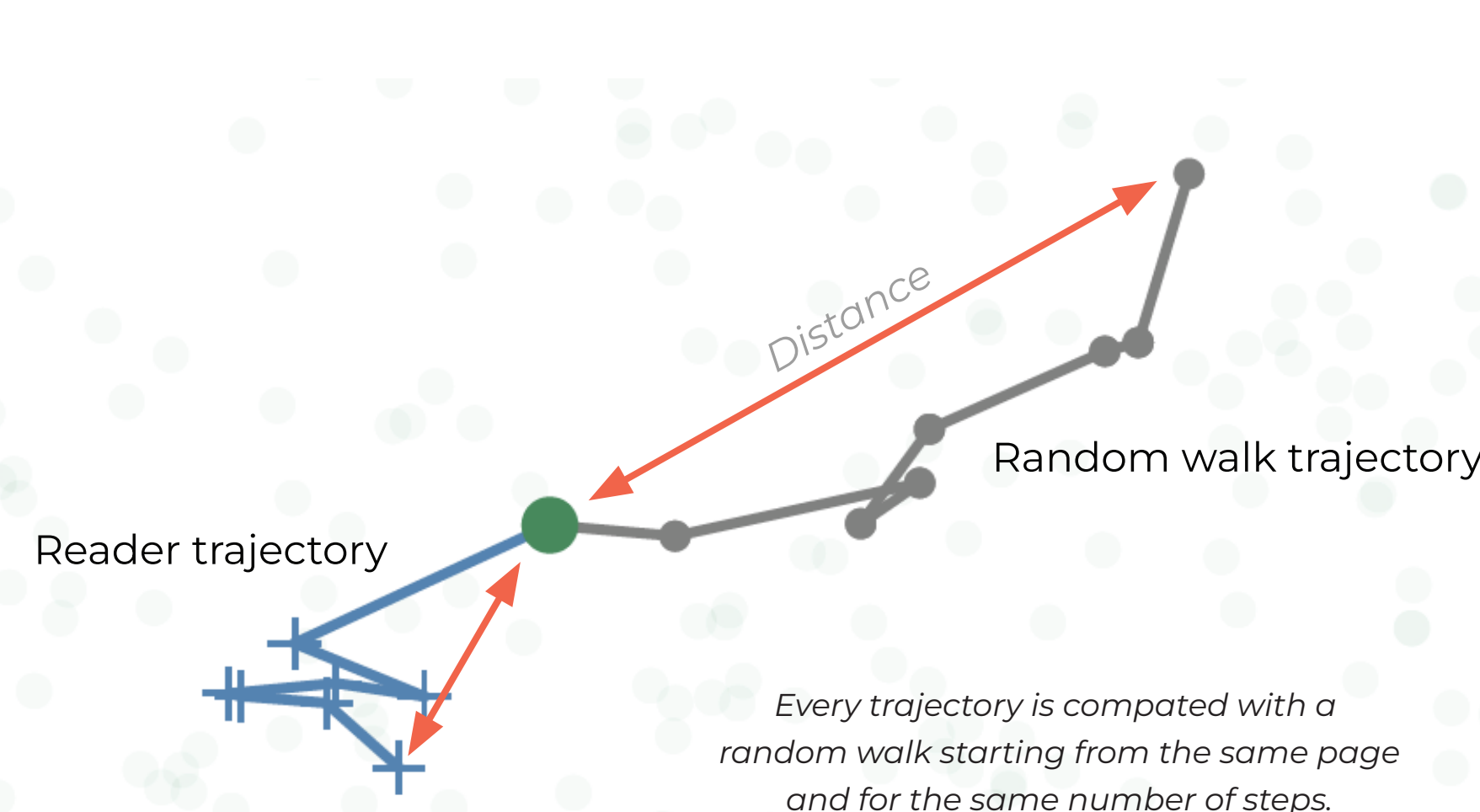


Logistic regression
Predict deep navigation



4

Evolution of the sessions



5