

# Applying mix de-anonymisation techniques for good

Steven Murdoch  
University College London

# Wi-Fi data collection

We collect Wi-Fi connection data at this station to better understand journey patterns and improve your services.

We will not identify individuals.

You can opt out by turning off your device's Wi-Fi.



# TfL monitor wifi MAC addresses to track mobility



 **Steven Murdoch** @sjmurdoch · Jul 15

My issues are that the poster implies there's an opt-out when it isn't really and also there would be privacy-preserving ways to achieve the same goal

1 1 1

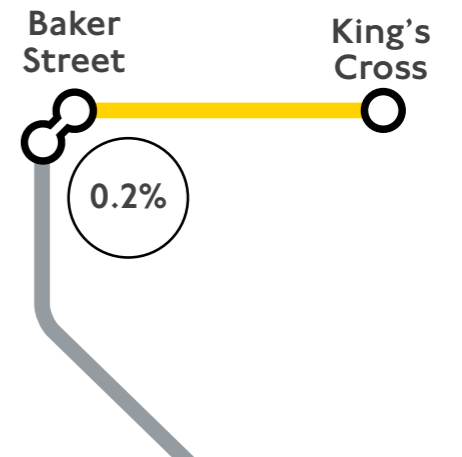
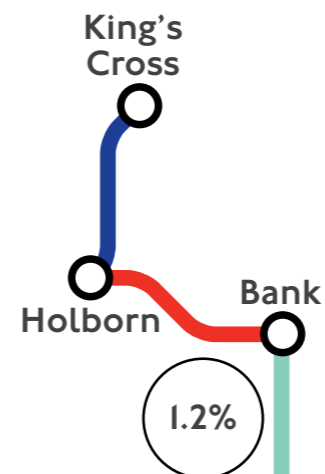
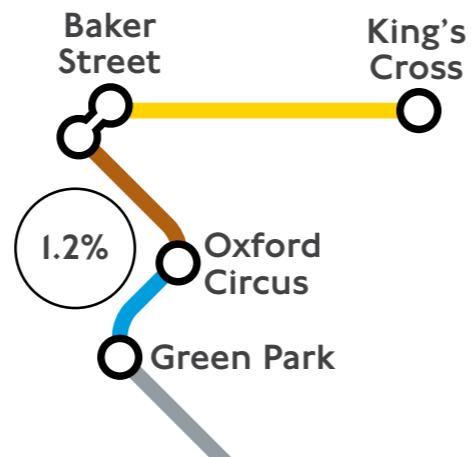
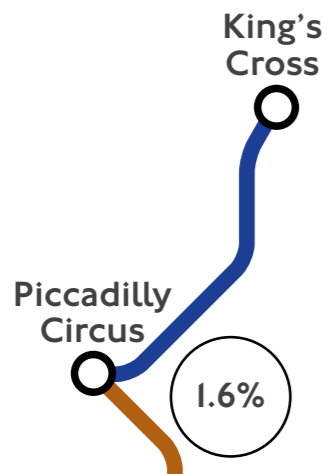
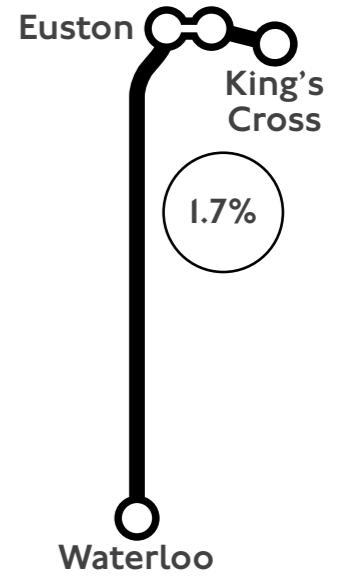
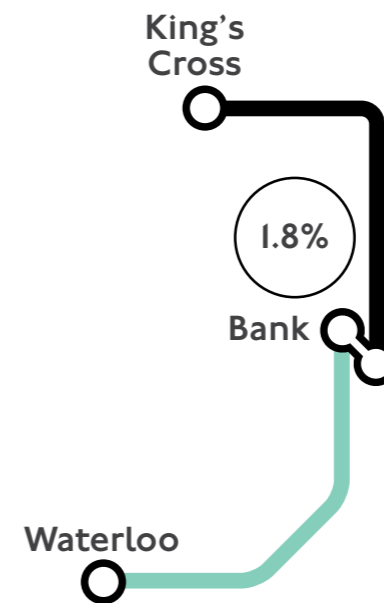
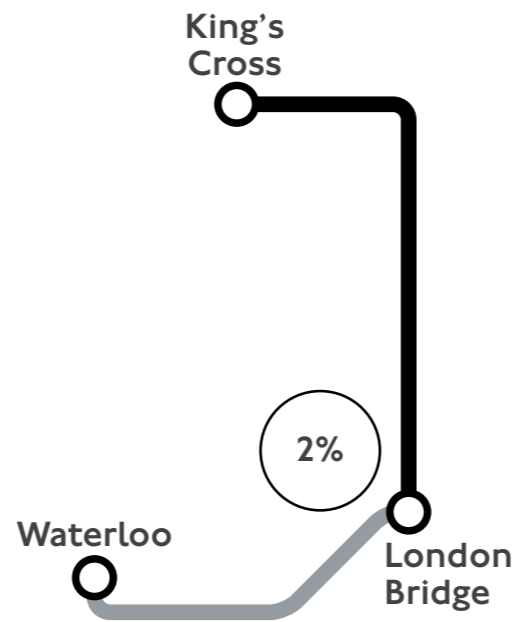
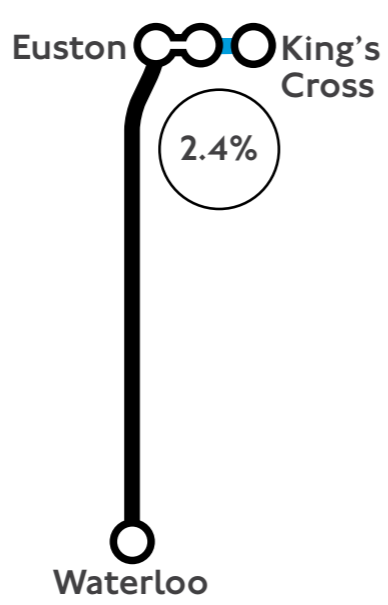
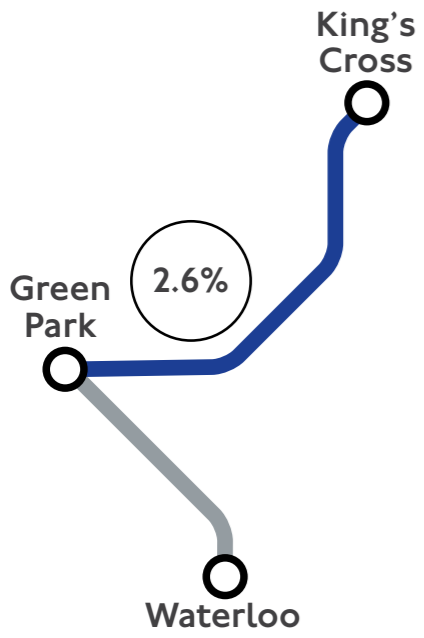
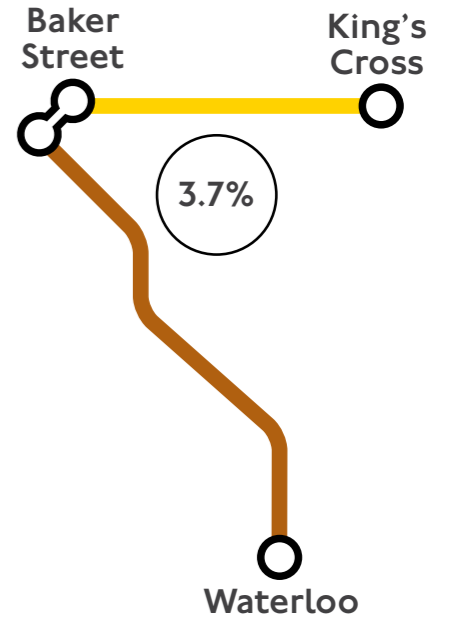
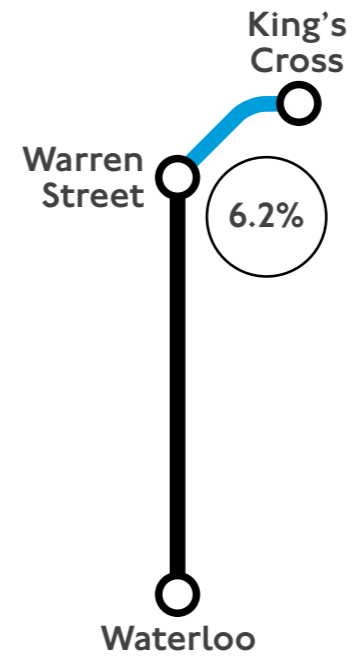
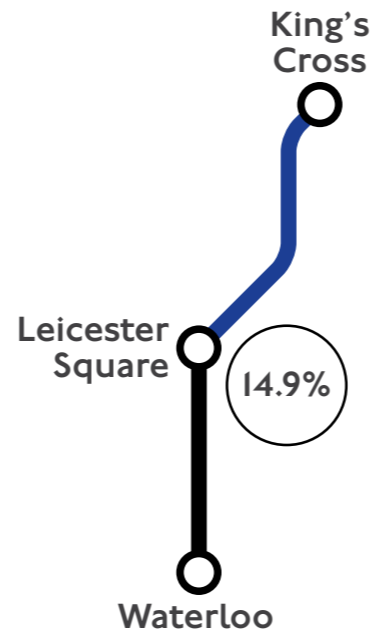
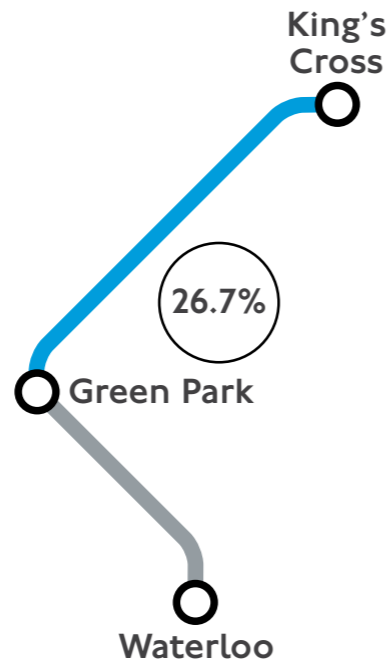
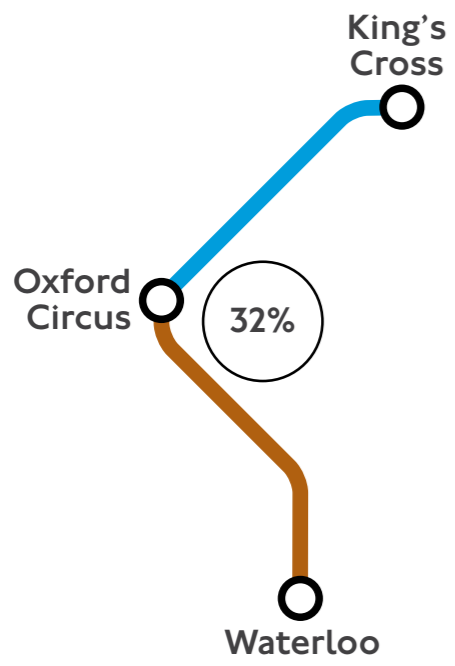
 **Kate Bevan** 🙄🙄👍  
@katebevan [Follow](#)

Replying to @sjmurdoch @futureidentity and 5 others

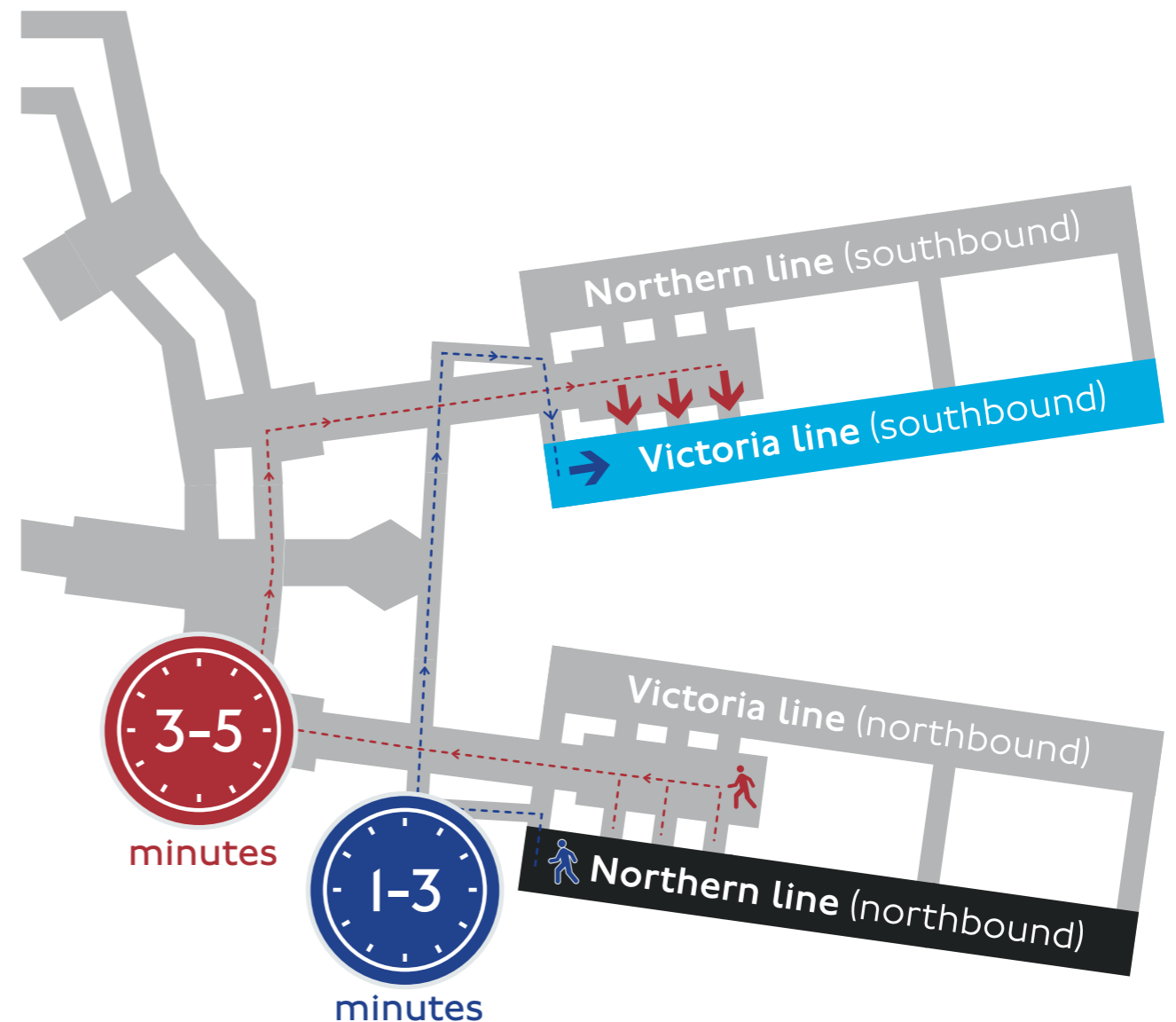
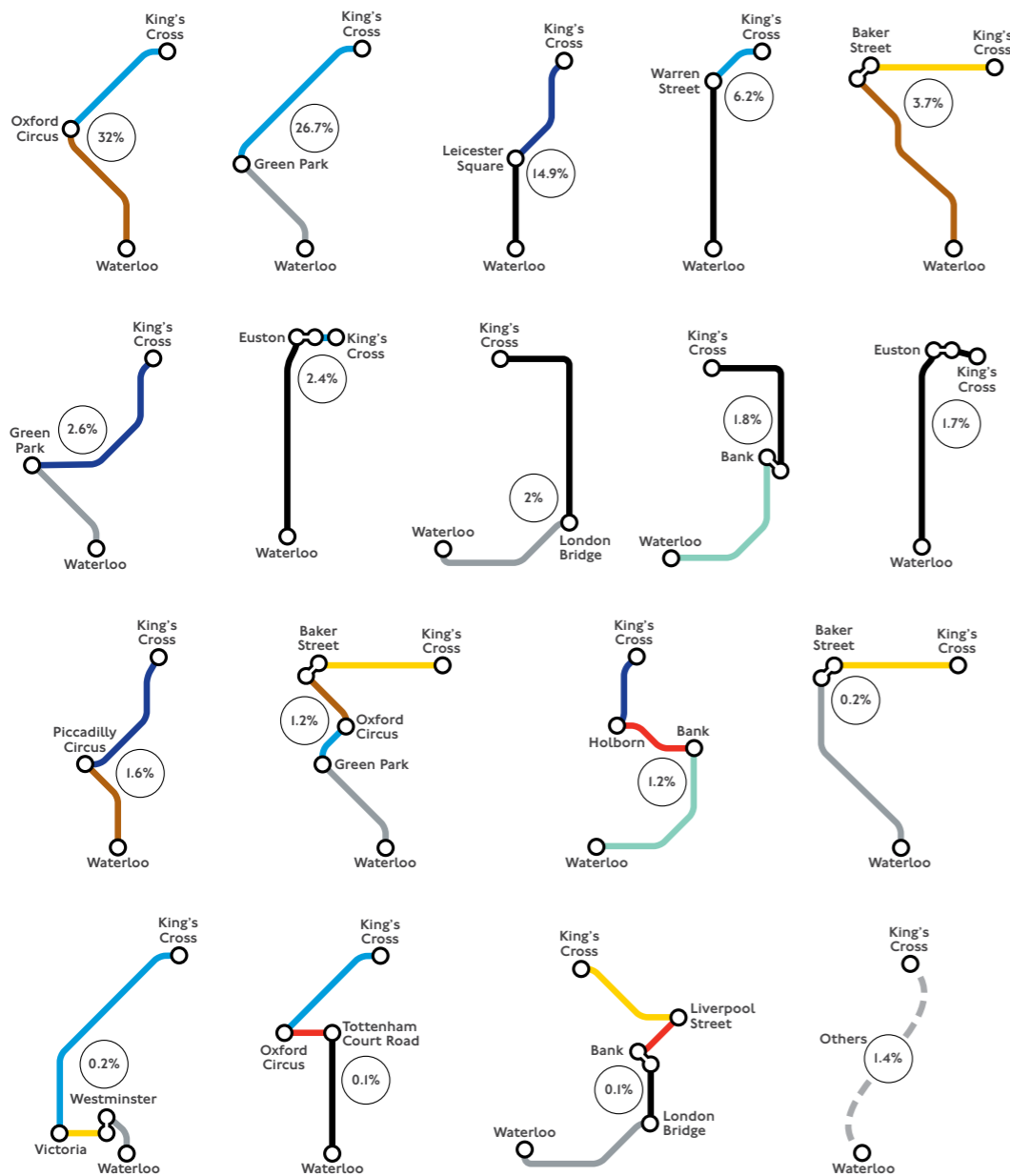
You can turn off your WiFi, though as Michael points out, it's less quick if you've got a laptop in your bag. But seriously, this is a good project. It's less perfectly implemented than some might like, but it's data for good. Not all data collection is terrible.

9:02 AM - 15 Jul 2019

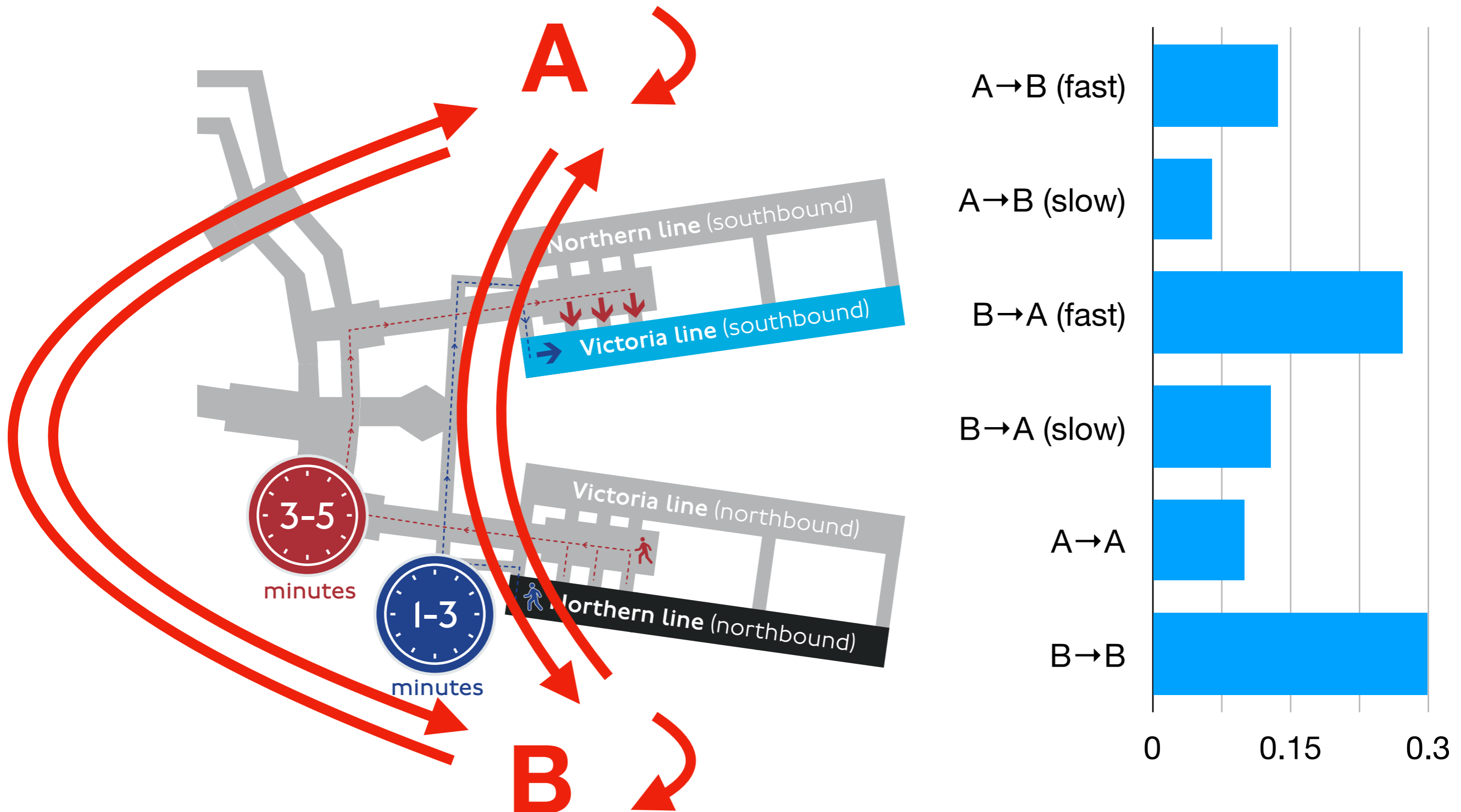
4 1 1



# Ticketing data doesn't explain movements in stations

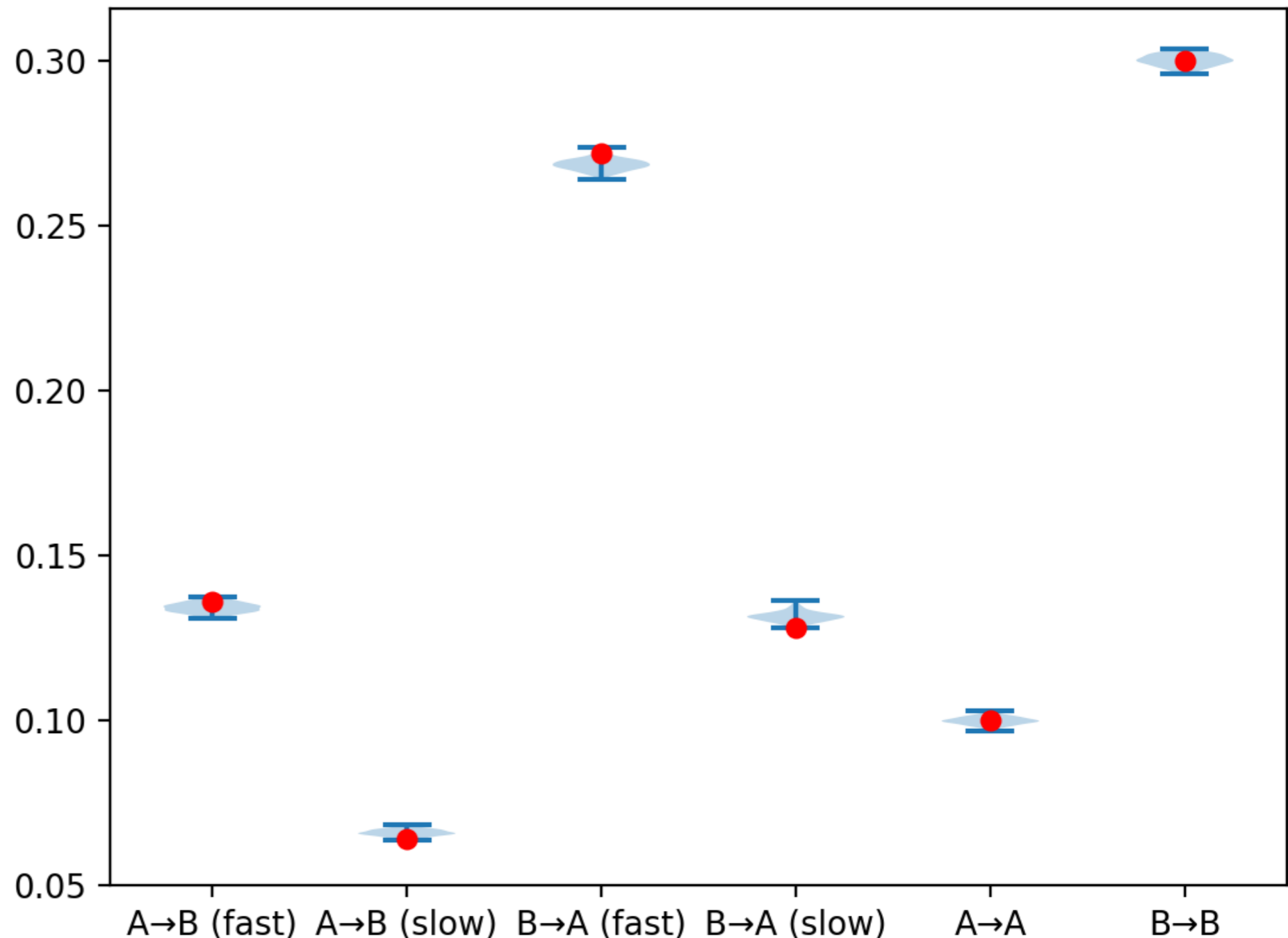


# We can simulate wifi observations in a station based on user profiles



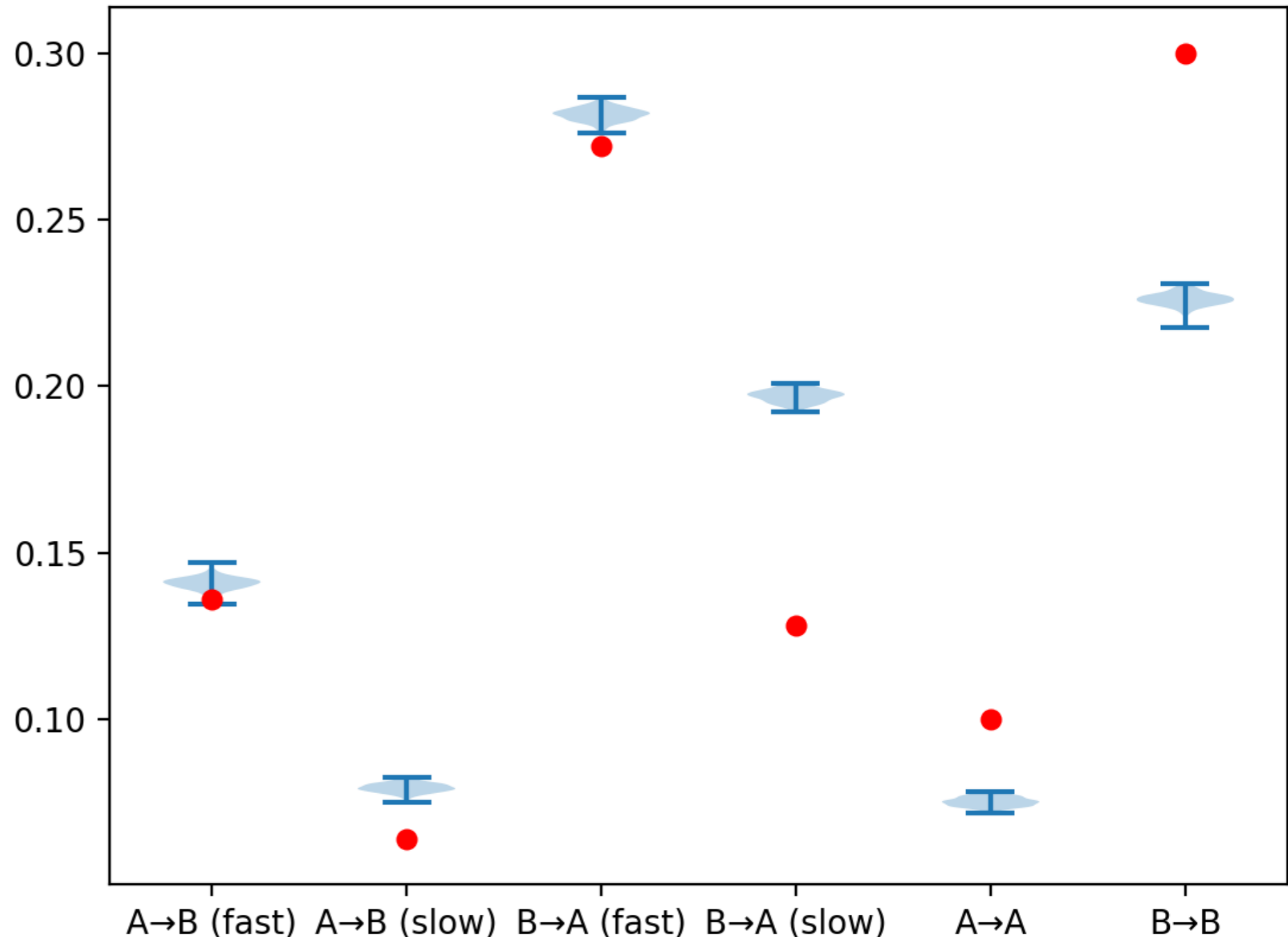
# Analysis of 64 bit MAC addresses gives good results

- Poisson arrival
  - $\lambda = 1$  per s
- Normal distribution for walking speed
  - $\sigma = 30$  s
- 100 simulations
  - Each 1 day

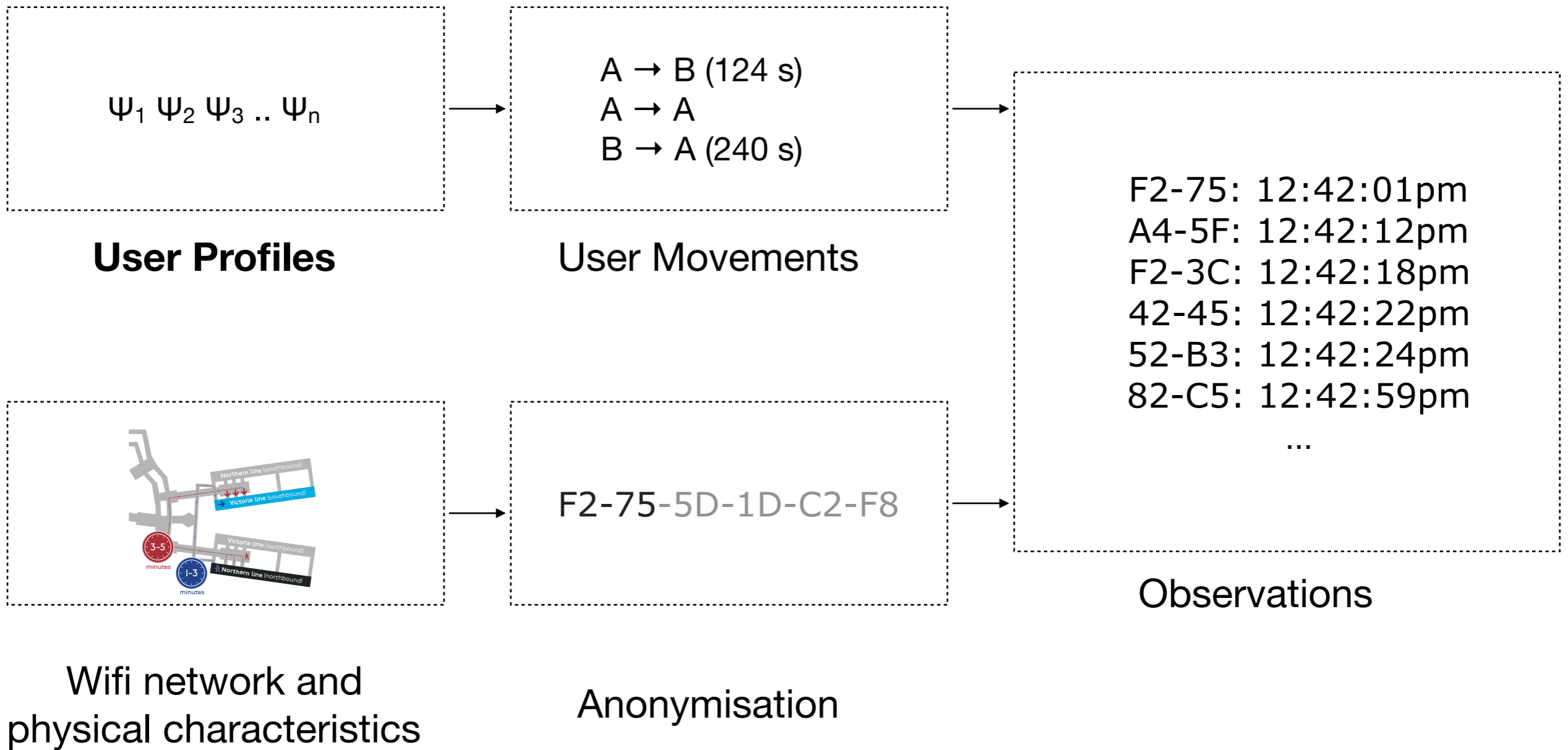


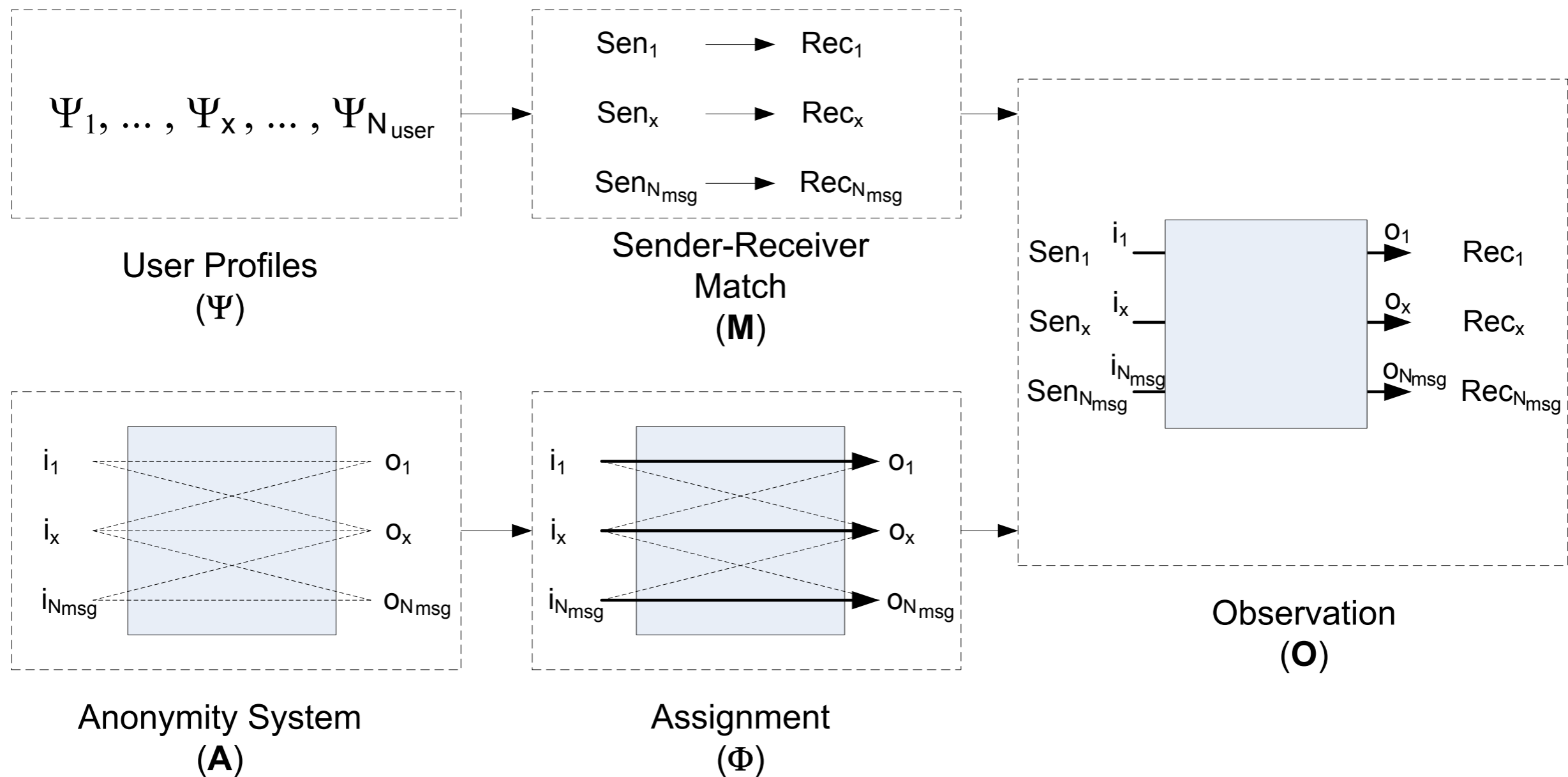
# Truncated 16 bit MACs don't work as well

- Poisson arrival
  - $\lambda = 1$  per s
- Normal distribution for walking speed
  - $\sigma = 30$  s
- 100 simulations
  - Each 1 day





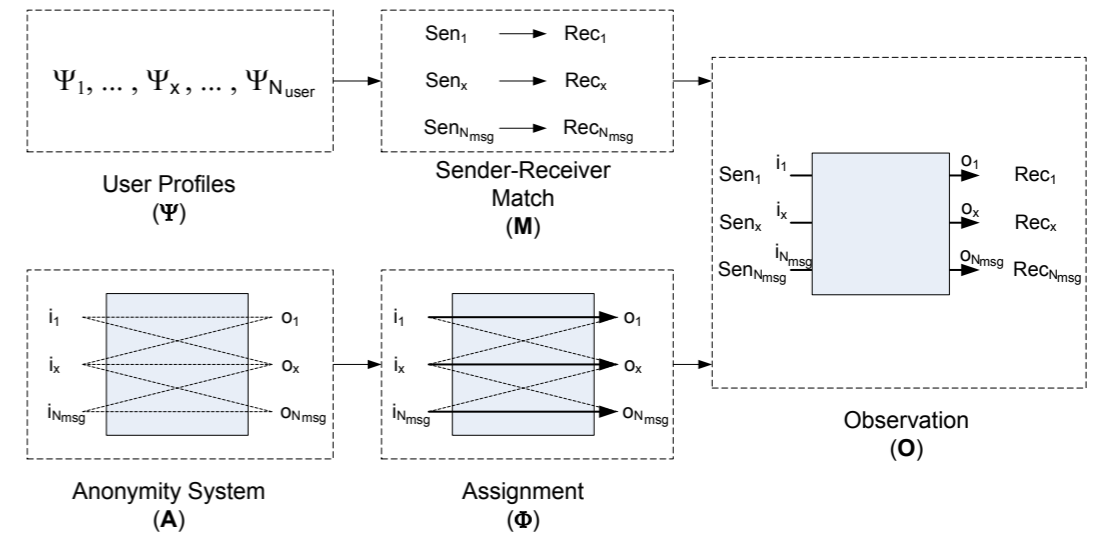




**Fig. 1.** The generative model used for Bayesian inference in anonymous communications.

We start by proposing a ‘forward’ generative model describing how messages are generated and sent through the anonymity system. We then use Bayes rule

- **De-anonymisation techniques can improve user privacy!**
- It is possible to **infer customer mobility profiles from observations of anonymised MAC addresses**
- **Model wifi network and MAC address anonymization as a mix network**
- **Take into account reasonable prior beliefs of mobility patterns**



**Fig. 1.** The generative model used for Bayesian inference in anonymous communications.

We start by proposing a ‘forward’ generative model describing how messages are generated and sent through the anonymity system. We then use Bayes rule to ‘invert’ the problem and perform inference on the unknown quantities. The broad outline of the generative model is depicted in Figure 1.

An anonymity system is abstracted as containing  $N_{\text{user}}$  users that send  $N_{\text{msg}}$  messages to each other. Each user is associated with a sending profile  $\Psi_x$  describing how they select their correspondents when sending a message. We assume, in this work, that those profiles are simple multinomial distributions, that are sampled independently when a message is to be sent to determine the receiver. We denote the collection of all sending profiles by  $\Psi = \{\Psi_x | x = 1 \dots N_{\text{user}}\}$ .

A given sequence of  $N_{\text{msg}}$  senders out of the  $N_{\text{user}}$  users of the system, denoted by  $\text{Sen}_1, \dots, \text{Sen}_{N_{\text{msg}}}$ , send a message while we observe the system. Using their sending profiles a corresponding sequence of receivers  $\text{Rec}_1, \dots, \text{Rec}_{N_{\text{msg}}}$  is selected to receive their messages. The probability of any receiver sequence is easy to compute. We denote this matching between senders and receivers as  $\mathcal{M}$ :

$$\Pr[\mathcal{M}|\Psi] = \prod_{x \in [1, N_{\text{msg}}]} \Pr[\text{Sen}_x \rightarrow \text{Rec}_x | \Psi_x].$$

In parallel with the matching process where users choose their communication partners, an anonymity system  $\mathcal{A}$  is used. This anonymity system is abstracted as a bipartite graph linking input messages  $i_x$  with potential output messages  $o_y$ , regardless of the identity of their senders and receivers. We note that completeness of the bipartite graph is not required by the model. The edges of the bipartite graph are weighted with  $w_{xy}$  that is simply the probability of the input message  $i_x$  being output as  $o_y$ :  $w_{xy} = \Pr[i_x \rightarrow o_y | \mathcal{A}]$ .

This anonymity system  $\mathcal{A}$  is used to determine a particular assignment of messages according to the weights  $w_{xy}$ . A single perfect matching on the bipartite

# UCL InfoSec are hiring!

Post-docs, PhD students, ...

New UK cybersecurity centre for doctoral training (CDT) – **55+ PhD studentships** over the next 8 years between Computer Science, Crime Science and Public Policy

[www.benthamsgaze.org](http://www.benthamsgaze.org)