

Introduction

- There were over 2.07 billion monthly active FB users in 2017 (Q3), and 1.37 billion daily active users. In every minute 510 000 comments are posted and around 300k status are updated
- In Hungary **90 percent of the active internet users** have a profile on Facebook
- Hungarians mainly use FB to **contact** friends and relatives, but half of them use it as a **news feed**, so it plays an extremely important role in their everyday life, according to their level of knowledge and way of thinking.
- Facebook is **not the main target of researchers**: a vast number of researches deal with Twitter data (Tinati et al 2014), but relatively few with FB data
 - it is easier to structure Twitter data: tweets are short and the number of possible actions a user can do is rather limited
 - it is much easier to get data from Twitter, through its API-s
- The case of FB, contents of the public sites (eg. restaurants, institutions, telco providers, universities, public figures, etc.) are possible to collect (through its Graph API), but **contents of pages of users are prohibited to gather**, even if these contents and activities are public and not private

Possible ways of (legal) data collection

- 1. Public FB sites data through Graph API**
 - This data collection strategy is good for some specific research questions, for example to study the FB activity of politicians
 - Public sites contain only part of the FB traffic, and not a random part!
 - FB **restricts** now this access
- 2. Through Facebook application**
 - Fill out a test and in exchange, let the researchers download the FB profiles of respondents
 - myPersonality Project as an example
 - more than 6 million test results, and more than 4 million FB profiles
 - data collection stopped at 2012
 - After Cambridge Analytica scandal, this data collection method is quasi **forbidden**
- 3. Facebook gives access to data for researchers**
 - Previously only few chosen researchers have the privilege for that
 - Social Science One is a new project, with a partnership of academy and Facebook
 - Launched in 2018 April
 - Research teams have to go through a grant process
 - If their project is supported, they can get access to anonymized Facebook data
 - Crowdtangle API
 - Ad Library API
 - URL shares dataset
- 4. Through web-browser plugin**
 - Installed web-browser plugins could monitor the users FB activities if consent given
 - LMU Munich plugin - <https://www.fbforschung.de>
 - Firefox and Chrome
 - collects public posts

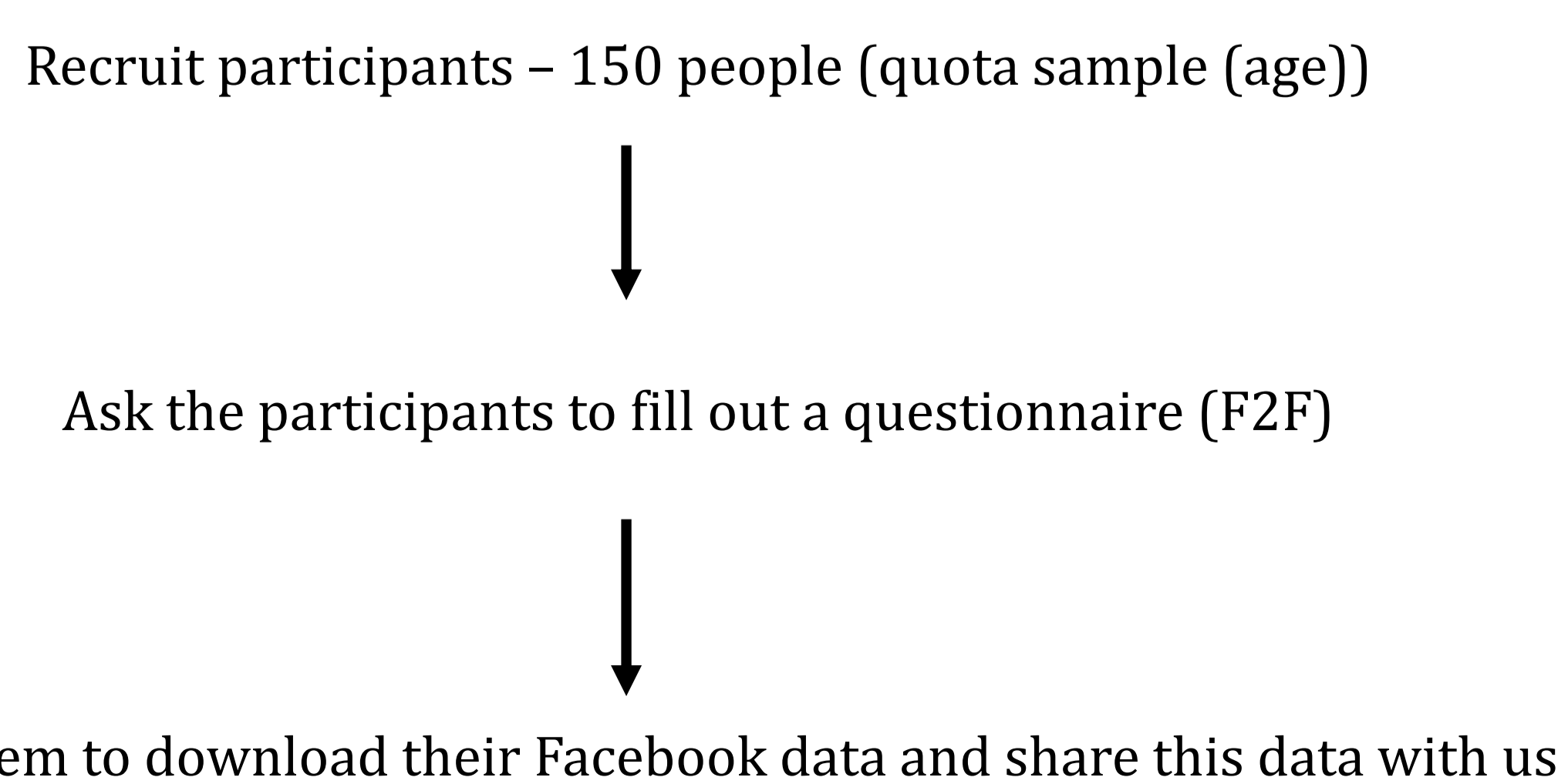
- 5. Users share their downloaded FB data with researchers**
Our project follow this strategy

Our research design

Objectives

- Analyze **private FB data**, not only posts shared on public sites
- Link survey and FB data**
Analyze which kind of **research questions** can be answered based on individuals FB data

Strategy



Our research design

Questionnaire cover different topics

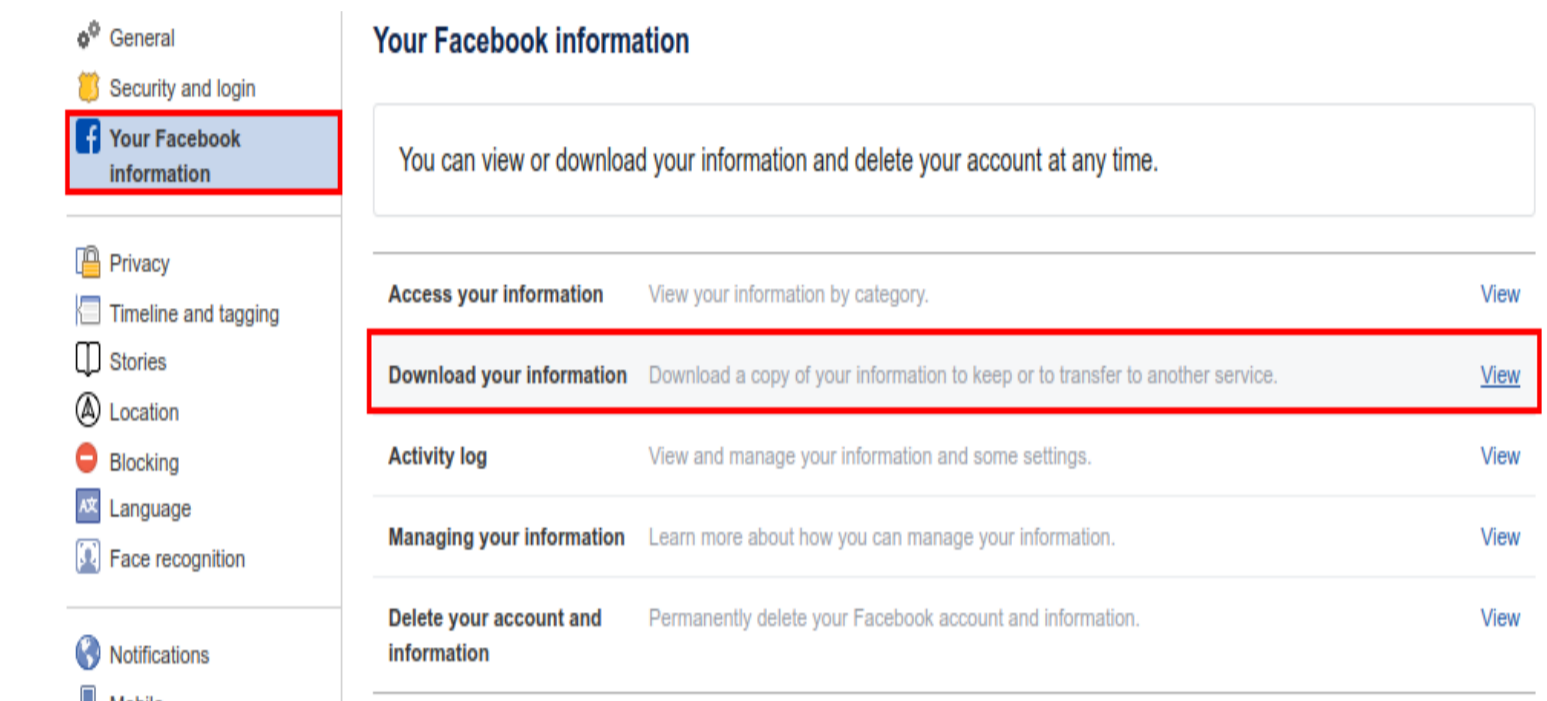
- General FB usage
- Self-representation in FB
- News consumption in FB
- Politics (political network homophily, party preferences)
- Leisure time activities
- Music preferences
- Depressions scale
- Demographic

ME1. Please rate the following statements on how characteristic they are for you.
1-7 scale of the answers:
1- not characteristic at all, 7- very typical of me
-1 DK/NA

I feel as if I don't know myself very well.	1	2	3	4	5	6	7	-1
I feel alienated from myself.	1	2	3	4	5	6	7	-1
I don't know how I really feel inside.	1	2	3	4	5	6	7	-1
I always stand by what I believe in.	1	2	3	4	5	6	7	-1
I am true to myself in most situations.	1	2	3	4	5	6	7	-1
I think it is better to be yourself, than to be popular.	1	2	3	4	5	6	7	-1
I live in accordance with my values and beliefs.	1	2	3	4	5	6	7	-1
I usually do what other people tell me to do.	1	2	3	4	5	6	7	-1
Other people influence me greatly.	1	2	3	4	5	6	7	-1
I am strongly influenced by the opinions of others.	1	2	3	4	5	6	7	-1
I always feel I need to do what others expect me to do.	1	2	3	4	5	6	7	-1

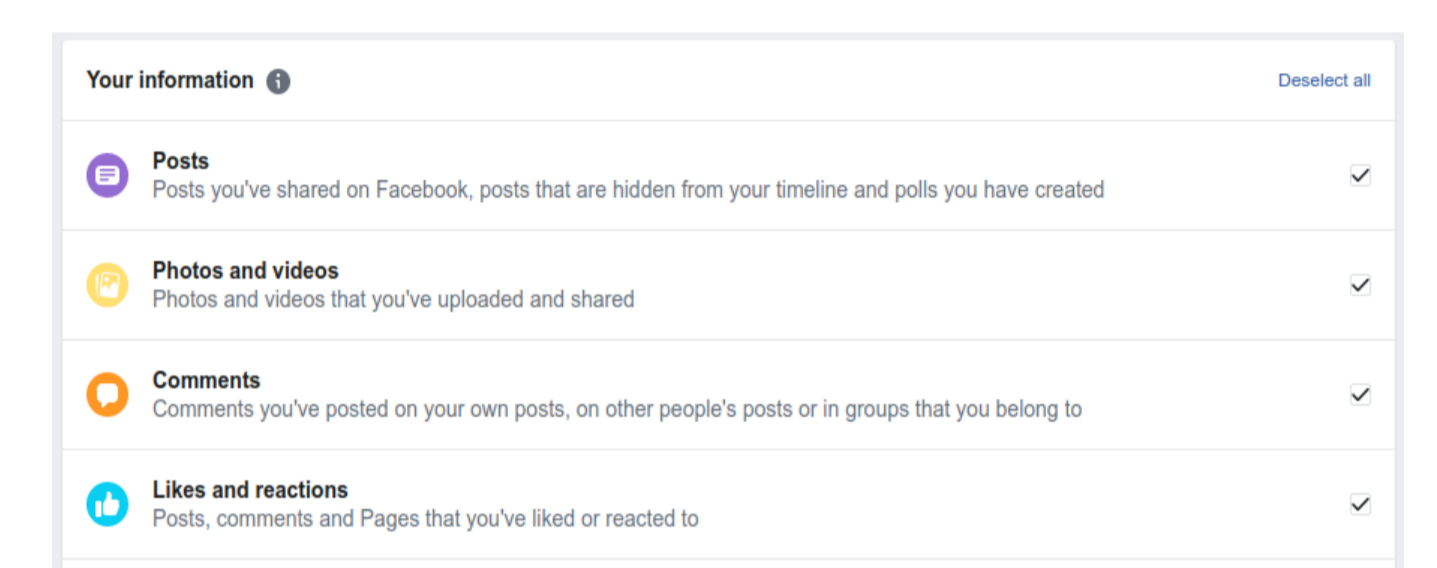
Facebook data

Ask participants to download their personal FB data, through the site



Some of the data tables contain very **sensitive information** (FB messenger), and some of them don't contain any useful information from research perspective, so we focus on specific tables:

- Posts
- Likes and reactions
- Friends
- Following and Followers
- Groups
- Profile information
- Ads



```

{
  "id": "1559492369",
  "timestamp": "1559492369",
  "comment": {
    "timestamp": "1559492369",
    "comment": {
      "author": "1559492369",
      "title": "I replied to [redacted] comment."
    }
  },
  "attachments": {
    "data": {
      "url": "https://www.facebook.com/1559492369/photos/?fbid=1559492369",
      "creation_timestamp": "1559492369",
      "media_attachments": {
        "title": "I replied to [redacted] comment."
      }
    }
  },
  "comment": {
    "timestamp": "1559492369",
    "comment": {
      "author": "1559492369",
      "title": "commented on her own photo."
    }
  }
}
    
```

Data format: Json

GDPR

Participants sign **consent form** before the interview

The research team follow **strict data privacy protocol**

We can only access and use **anonymized data**

After downloading FB data in Json form, an R script is used to pre-process the data

Data format changed to CSV (small data loss, but more handy format for social scientists)

Names are masked

MD5 method (openssl package in R)

We can only mask those names which are in the contact list

Json files are deleted researchers only have access to processed files

V1	V2	V3
1560340774	@70e4d846a97b62c0d63739054460 replied to [redacted] comment.	[redacted] Igy még jobbi!
1560320607	@70e4d846a97b62c0d63739054460 replied to [redacted] comment.	@1321778da160b147b35237b36c8181 de azért az meno, hogy az embereknek a Napló jut eszébe
1560319900	@70e4d846a97b62c0d63739054460 replied to [redacted] comment.	@1321778da160b147b35237b36c8181
1560194900	@70e4d846a97b62c0d63739054460 commented on @2950a650a033896970a058d37763d 's video.	En azt?
1560107107	@70e4d846a97b62c0d63739054460 commented on @65d09063001e46e480080564609e 's photo.	Nagyon boldog!!!
1560028107	@70e4d846a97b62c0d63739054460 commented on @1321778da160b147b35237b36c8181 's photo.	Point itt voltak. Daravál es Daravál par her!
1560006685	@70e4d846a97b62c0d63739054460 commented on @09e3bee67f62f9040802a7ec3008a 's post.	Boldog az újragot Lujának!
1560006090	@70e4d846a97b62c0d63739054460 commented on [redacted] post.	@1317b0c89b0bea31baeb720b06c48d5
1560002309	@70e4d846a97b62c0d63739054460 replied to @9d870e28e131b0c5d46d26f85cd0d97 's comment.	@d8d70e28e131b0c5d46d26f85cd0d97 már kint vannak, így most nem is lesztek!

Challenges

- People names who are not in the **friend list** are not masked after the first round of anonymization
- **Nicknames or not tagged names** are not anonymized by the script
- **Human processing** needed

Limitations - fragmented pieces

Although we will have access to private FB data, this data will be **fragmented** in many ways

In the case of comments by participants we don't know the **original content**. So we see a comment, but don't know what was the original post

In the case of likes and reactions we also don't know the **original content**

We can not follow a **whole thread**, just the some pieces of it

This is a serious limit, but totally understandable from GDPR point of view

Planned papers

Wide spectrum of papers are planned:

- Methodological paper (pros and cons of this research strategy)
- Self representation and posting habits
- Depression and posting habits
- Music tastes and liked pages of musicians and bands (omnivore – univore taste)
- Politics related posts in off- and on-campaign period

Collaboration is possible, we're open for joiners

We can share the survey data and aggregated FB data (not raw data)