



Kahoona



Y-DATA

BEHAVIORAL CROSS-SESSION USER IDENTIFICATION

Evgenia Amineva
Rima Rohana



THE TEAM

Y-Data Students

Evgenia Amineva, background in Project Management

Rima Rohana, background in Software Engineer

Kahoona

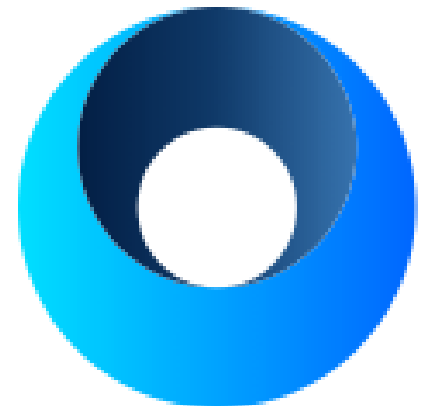
Noy Nissim, Data scientist, Y-Data graduated

Nadav Kallenberg, Data Science Team lead

The Mentor

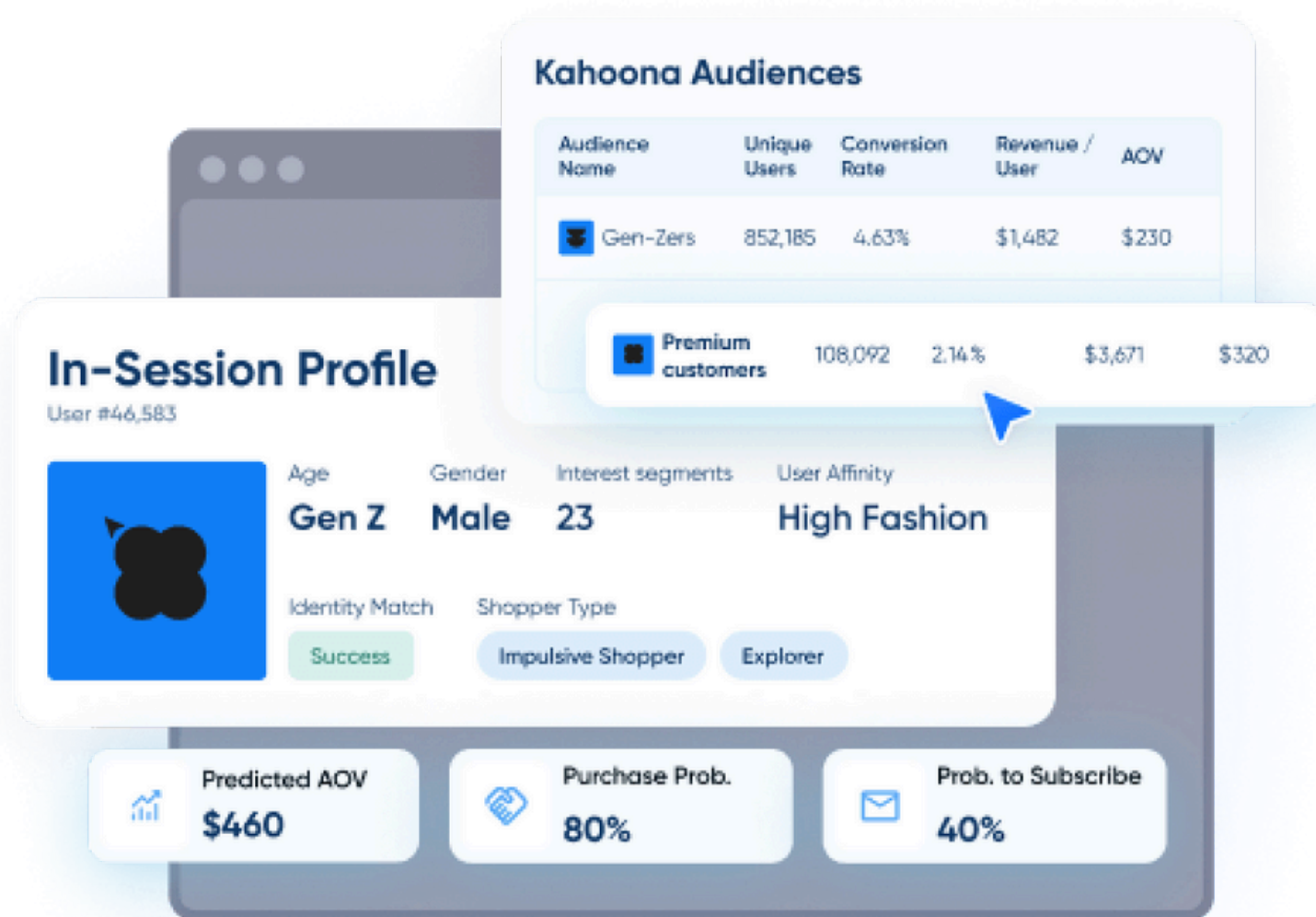
Oren Elisha, Data Science Manager, Forter

THE COMPANY



Kahoon

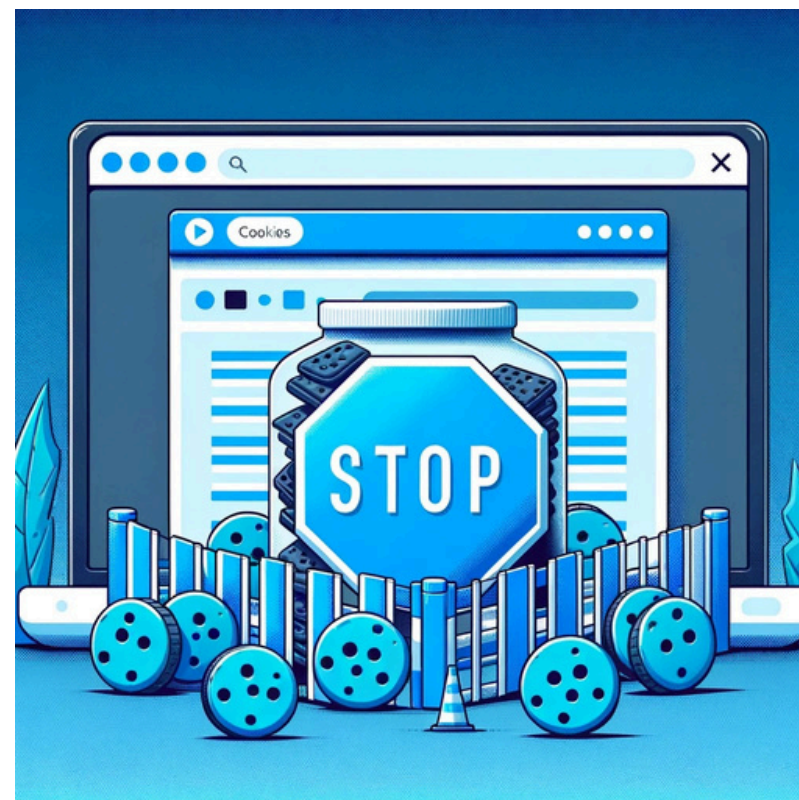
A first-party data activation platform that provides an essential solution for a cookie-free and identity-free digital ecosystem



FOUNDATIONS

95%
of E-commerce
websites visitors are
unidentified

**Third-party
cookie
restriction**

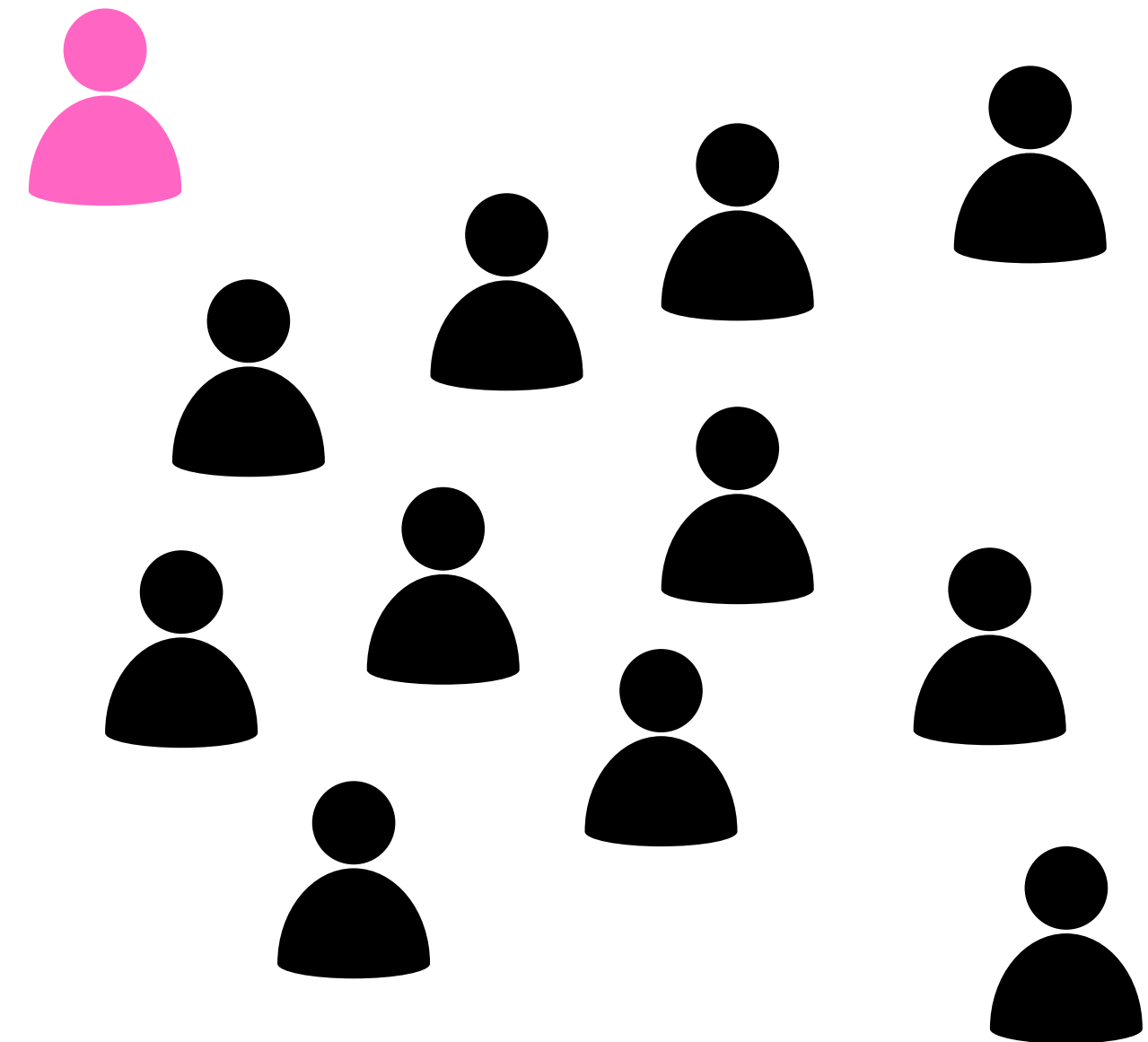
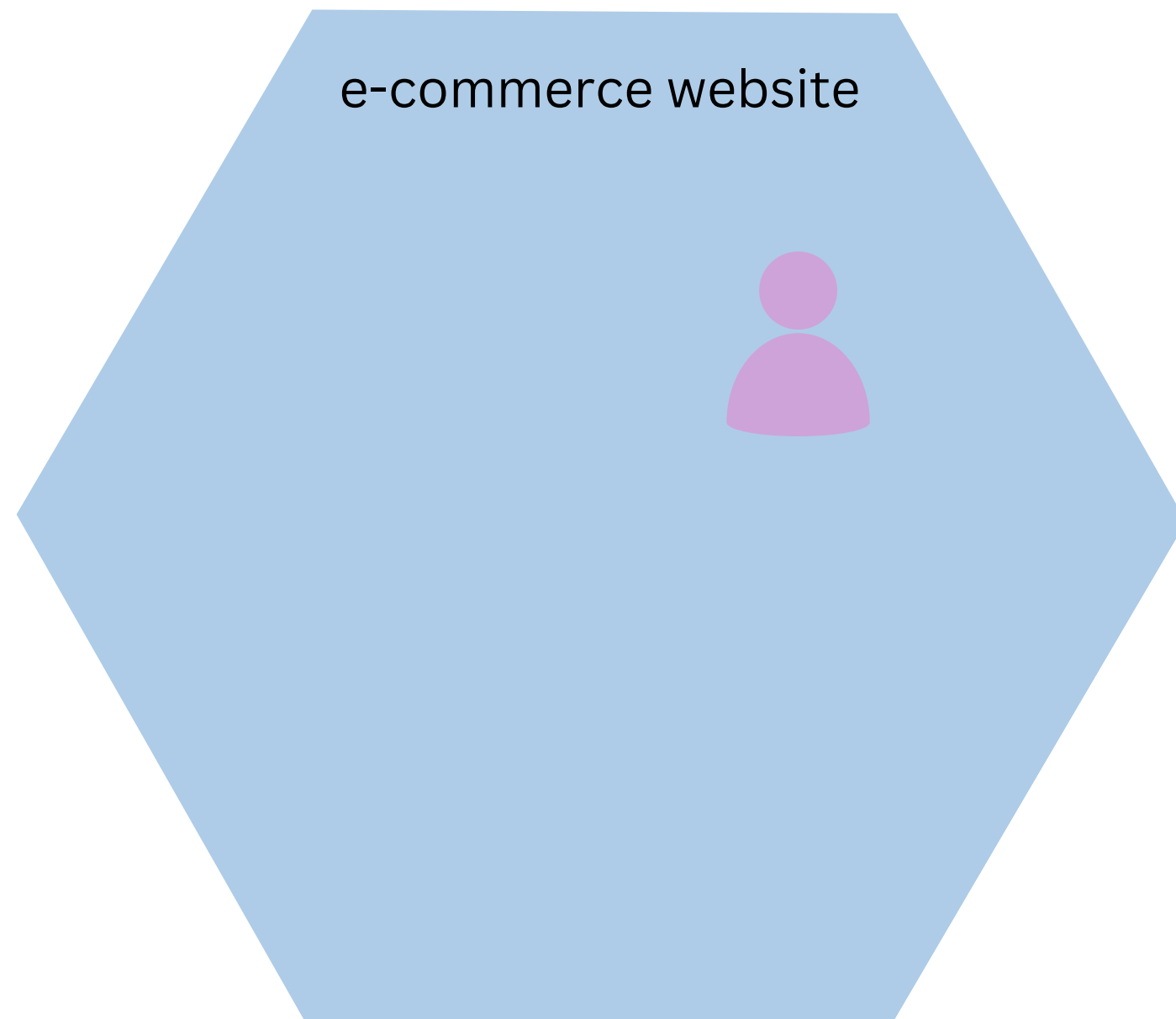


User behavior

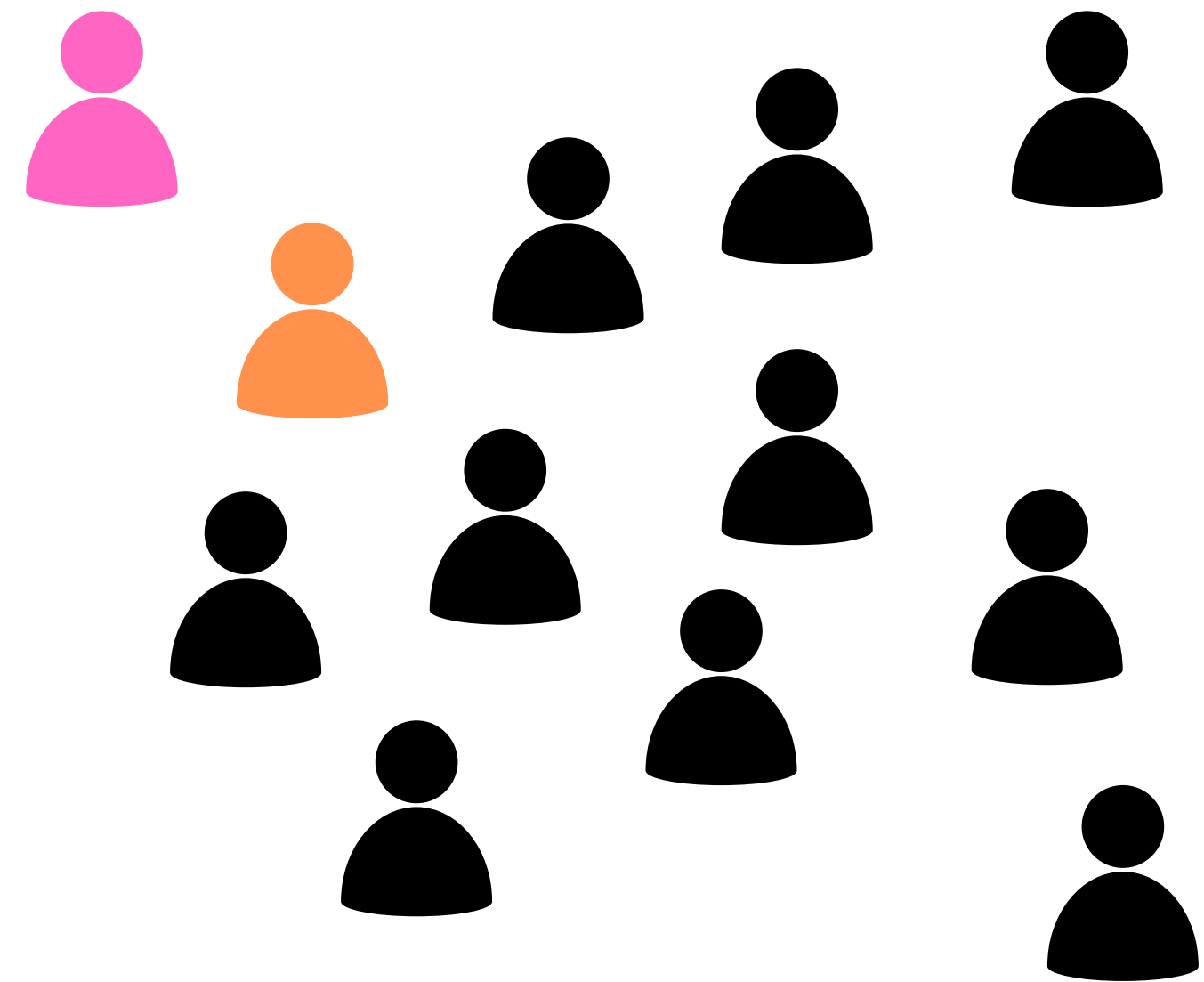
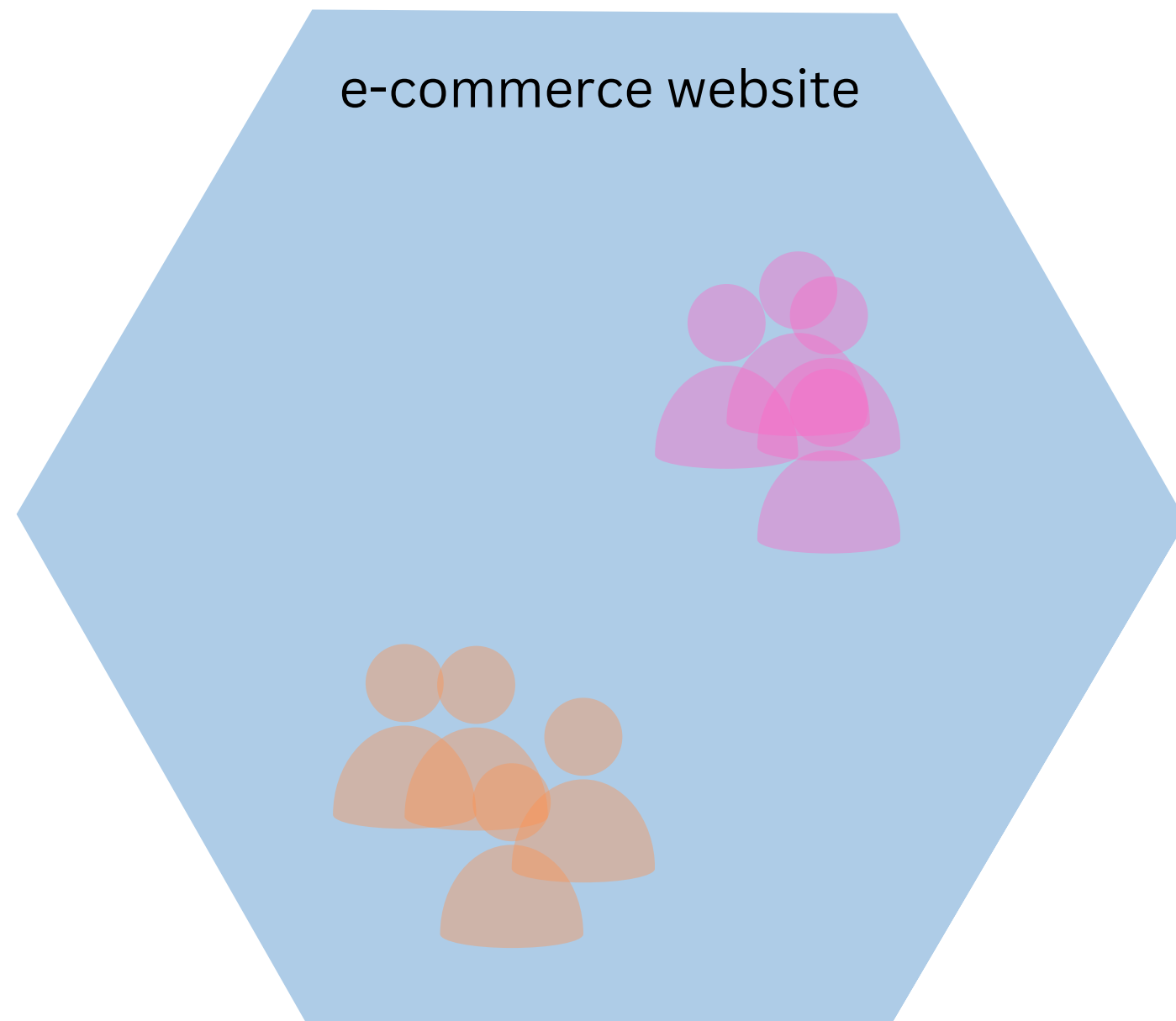
New perspective on user
behavior in E-commerce
websites that focuses on
the 'how' of user
interaction rather than
'what'.

PROBLEM

Identification of not logged-in returned users



SIMILARITY PROBLEM





INITIAL GOALS

Main goal:

Find 3 most likely matching users and beat the company's accuracy 90%

Advances goal:

Develop a mechanism that picks the most likely matched user

THE DATA

Hashed IDs: session, pageview ...

1db1bb7c-242e-
4bb8-bb63-
c3b56e4459a3

2f274e9a-5e37-
45fa-aec4-
a71d5d190150

Timestamps: session, pageview ...

1677352958273

1677353614099

1677352958249

Environmental data: user_agent, referrer

Mozilla/5.0 (iPhone; CPU
iPhone OS 15_4_1 like...

External IDs:

None

None

abFBKzTsO0FBSP8LhM2ACMzweI

Behavioral data: business_data ...

[{'reloads': 0, 'eventtype': 'khn_loaded', 'ev...

[{"events":[{"et":"load","ets":16773536...

Label:

7eb1033072cc795e9340dc7f09e57a29524dfc89652c70...

user_id

THE DATA



33 features
2 339 008 entries



Parquet
file format
8.32 GB



Period

January 21, 2023 – March 12, 2024

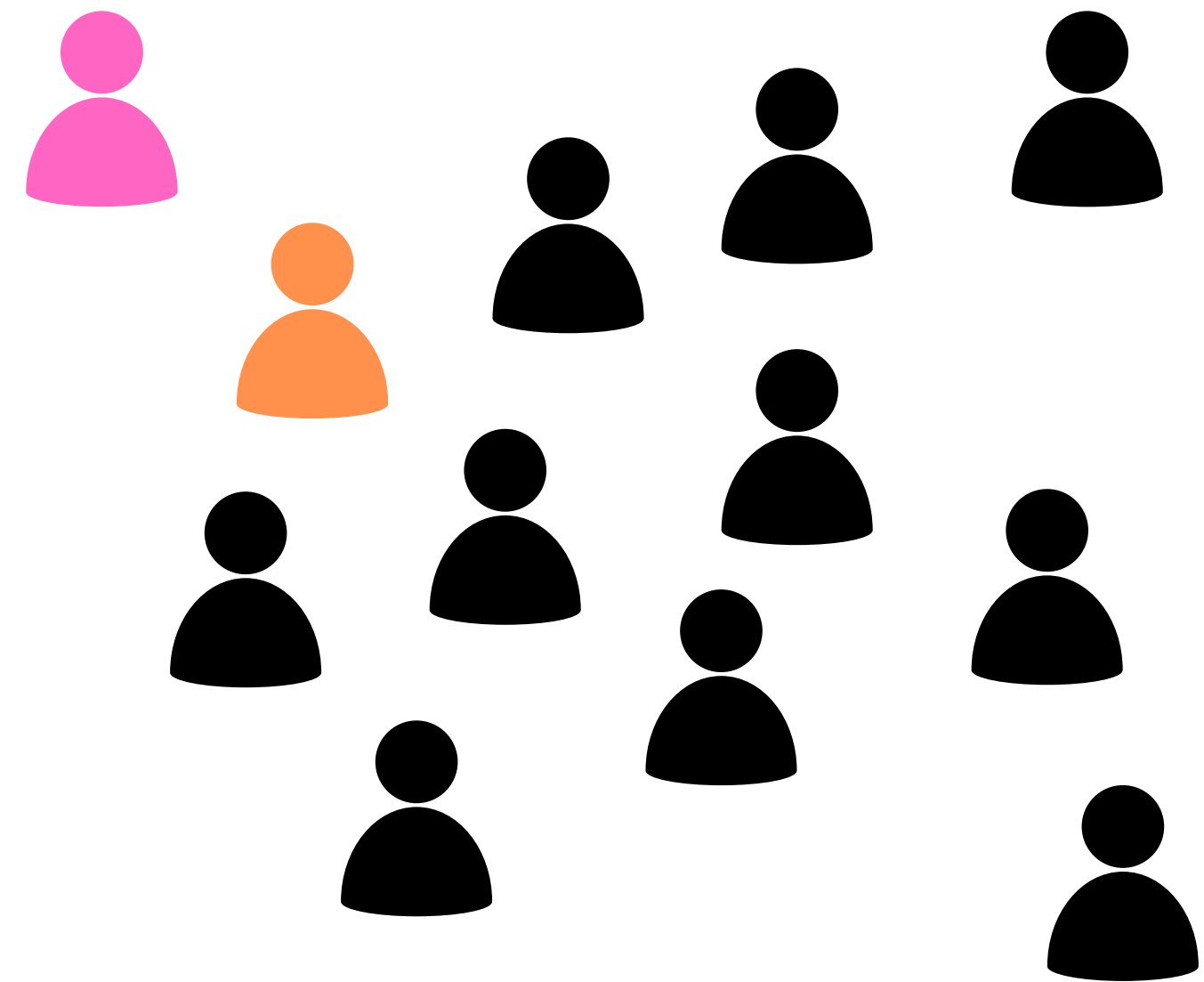
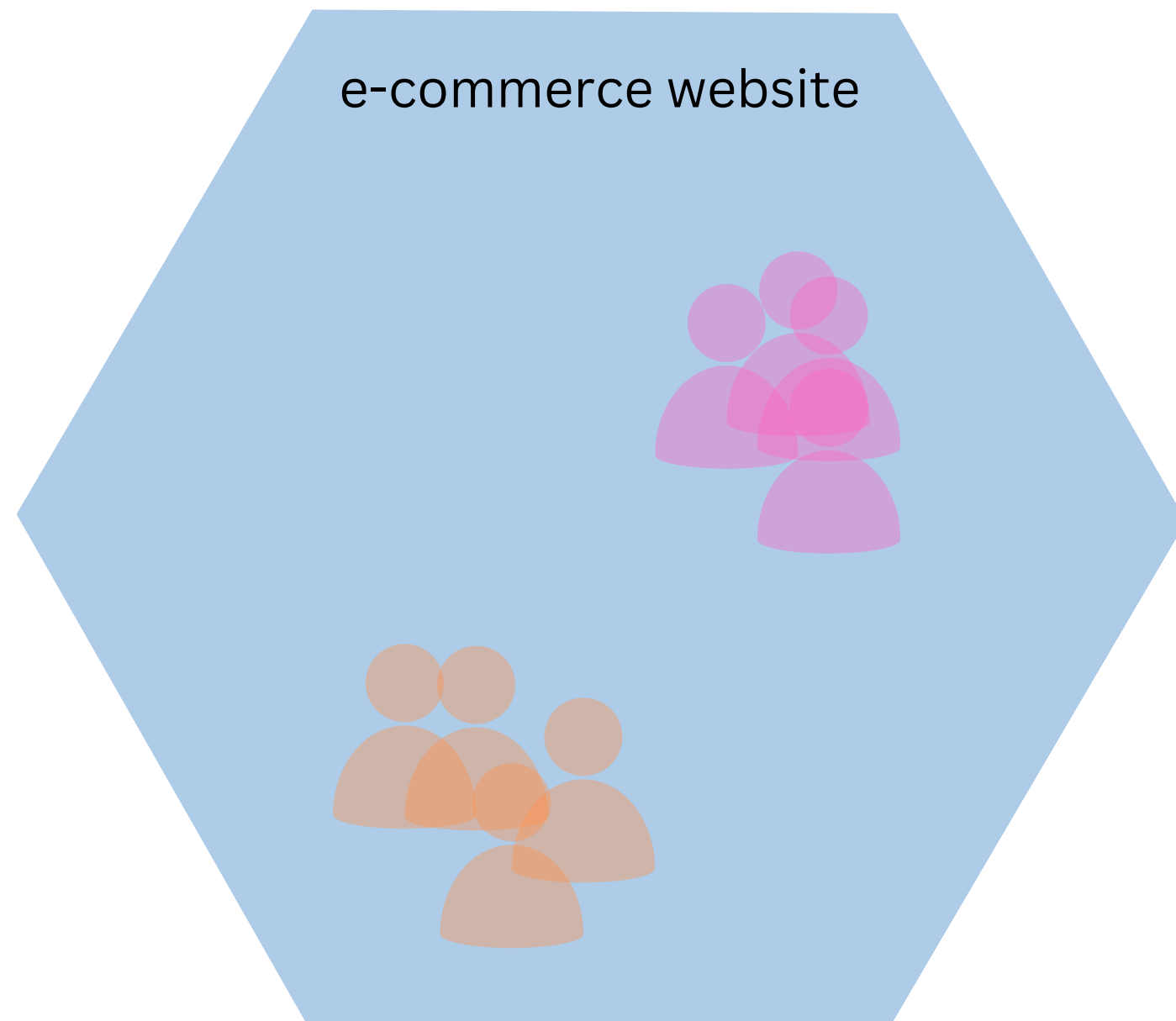


Source

e-commerce
website

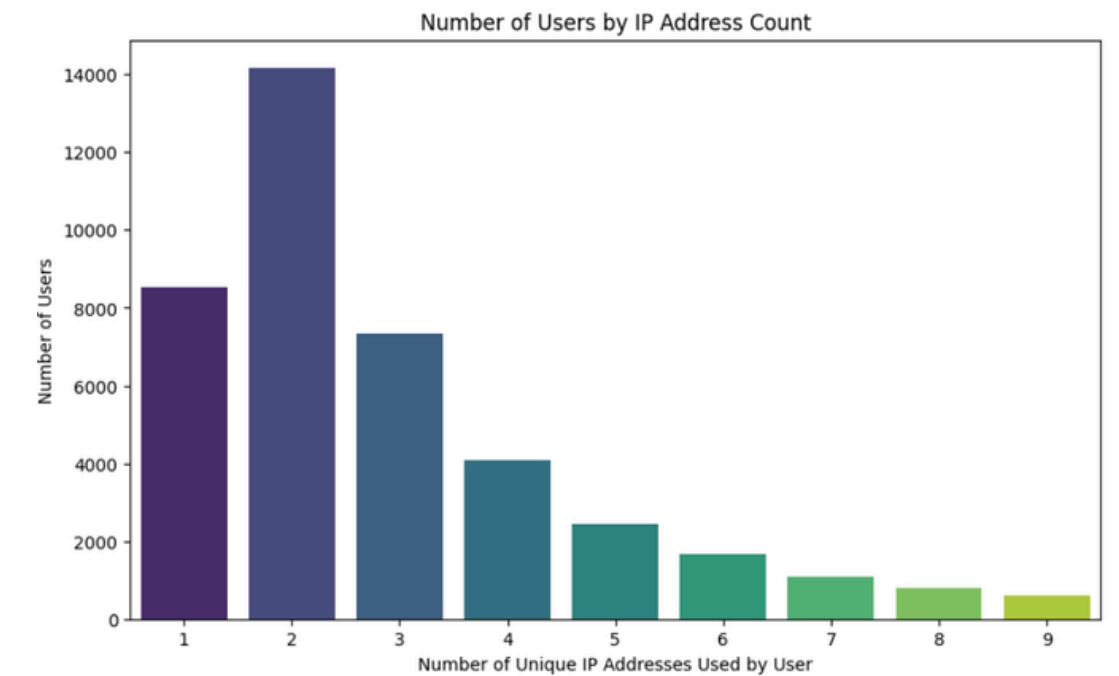
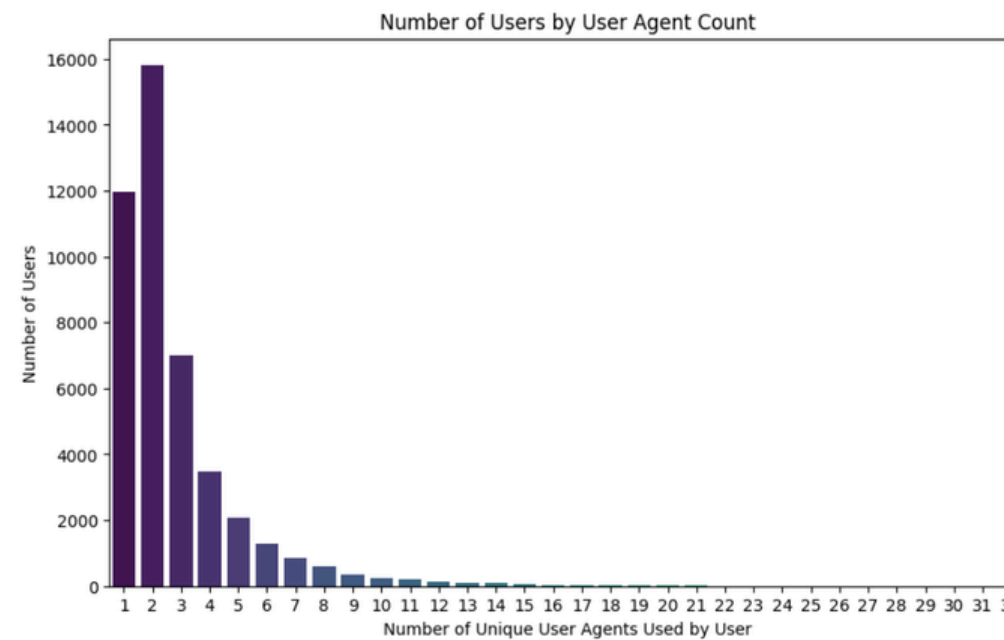
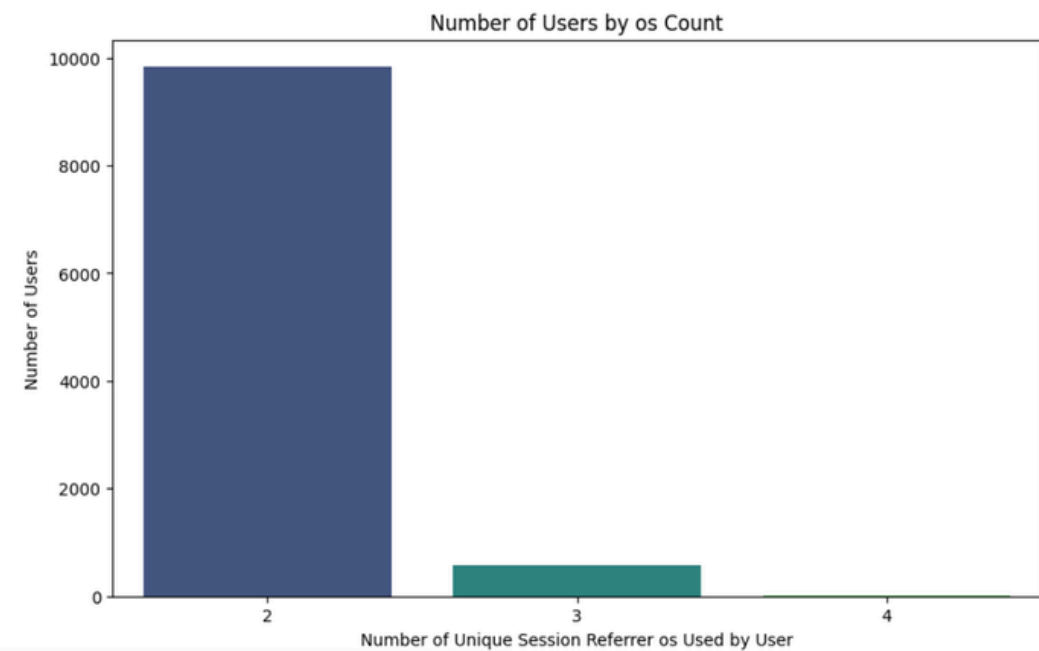
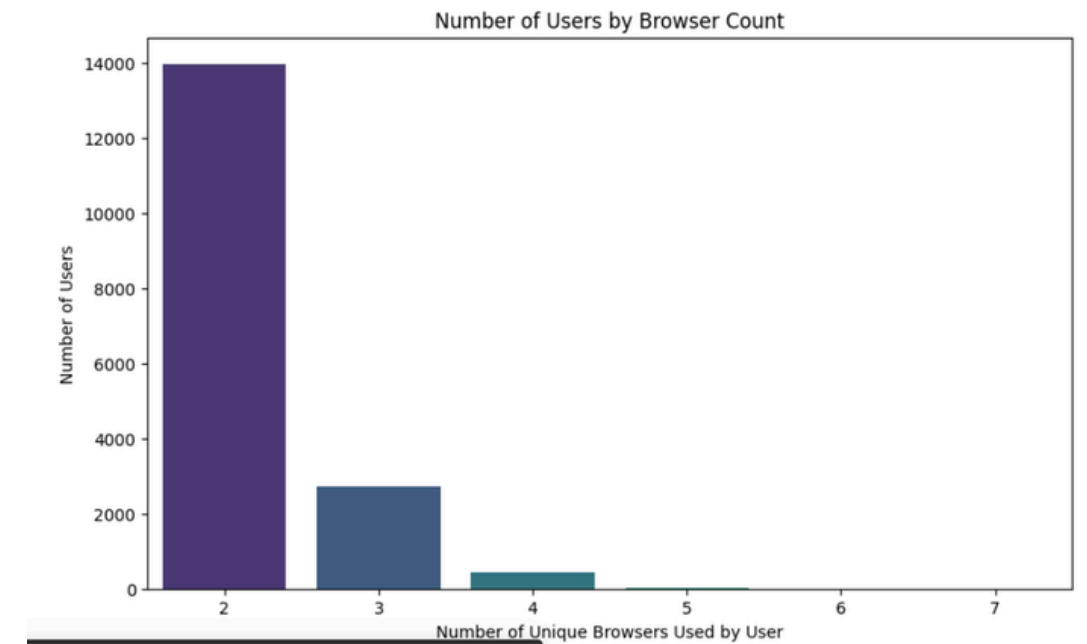
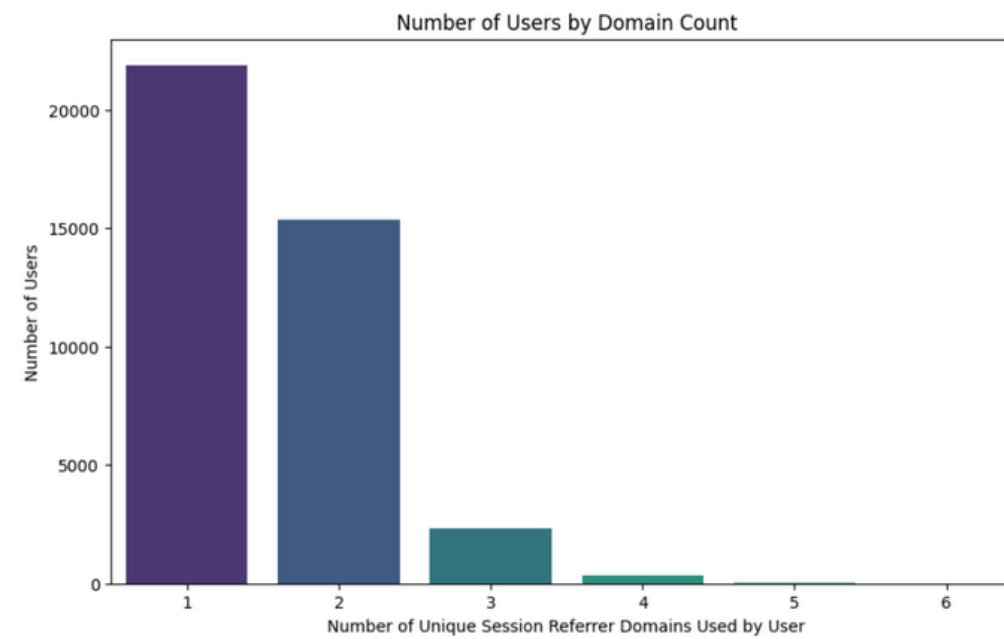
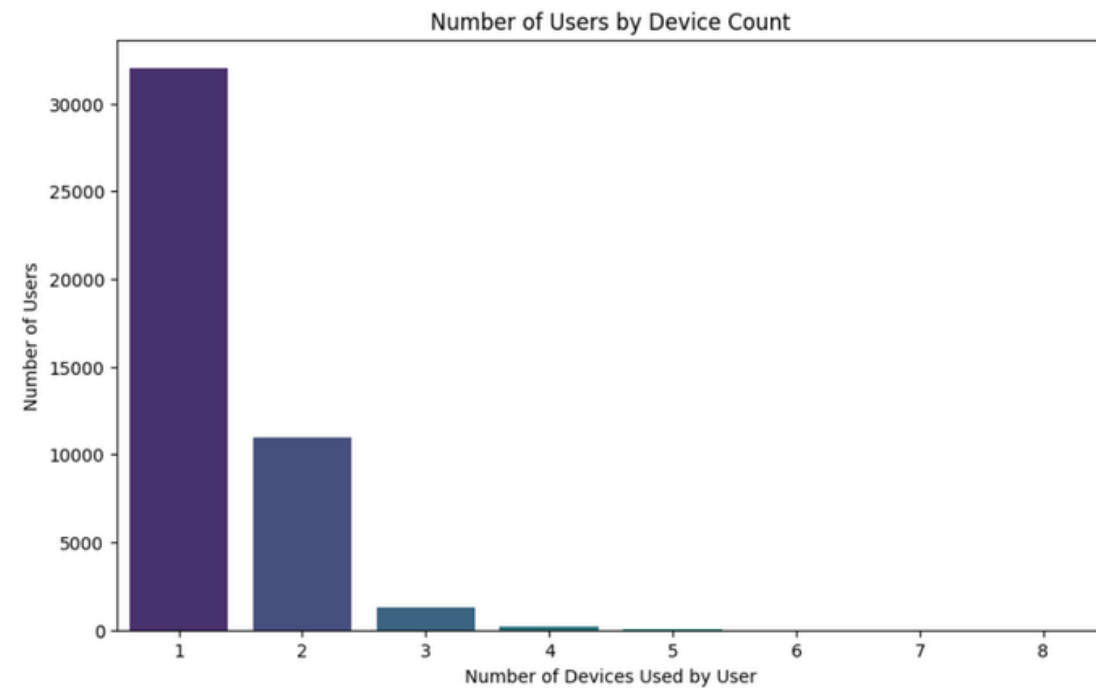
USER TRACES

why we assume that users leave traces?



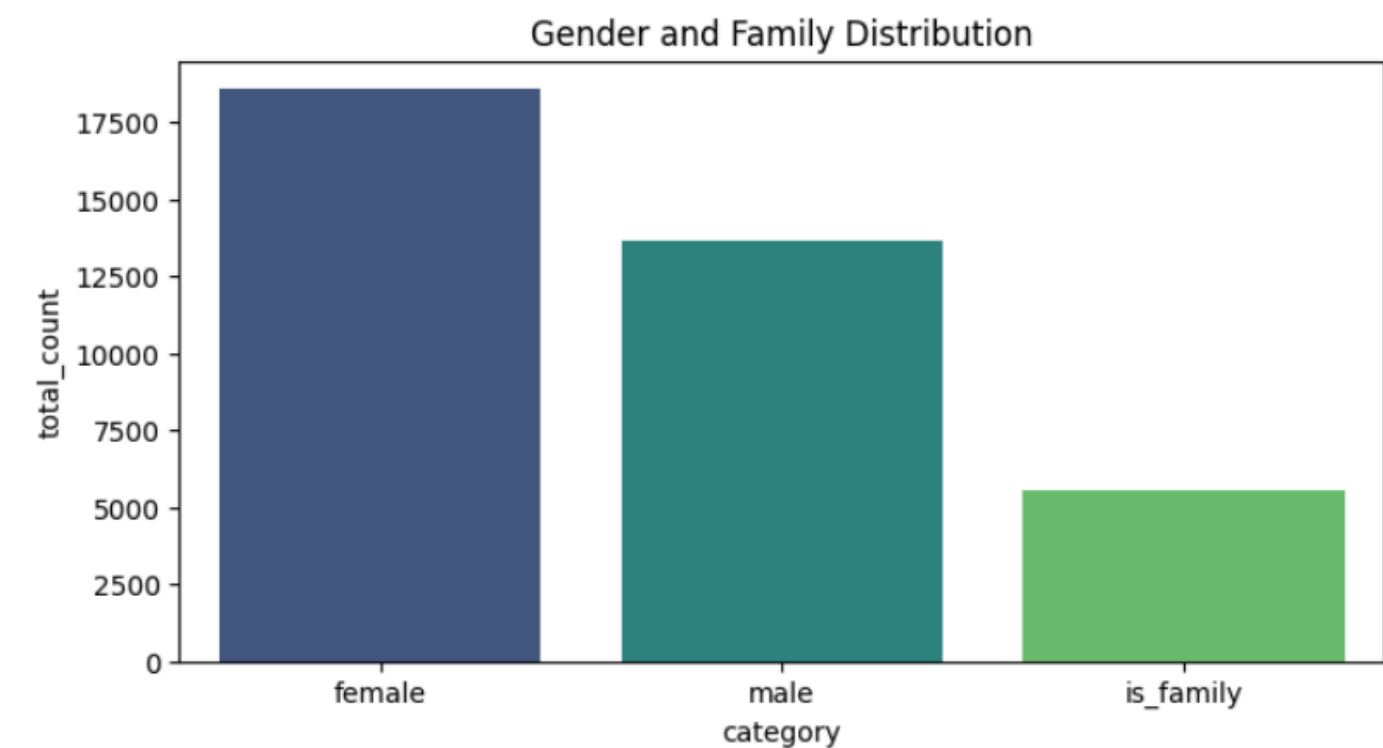
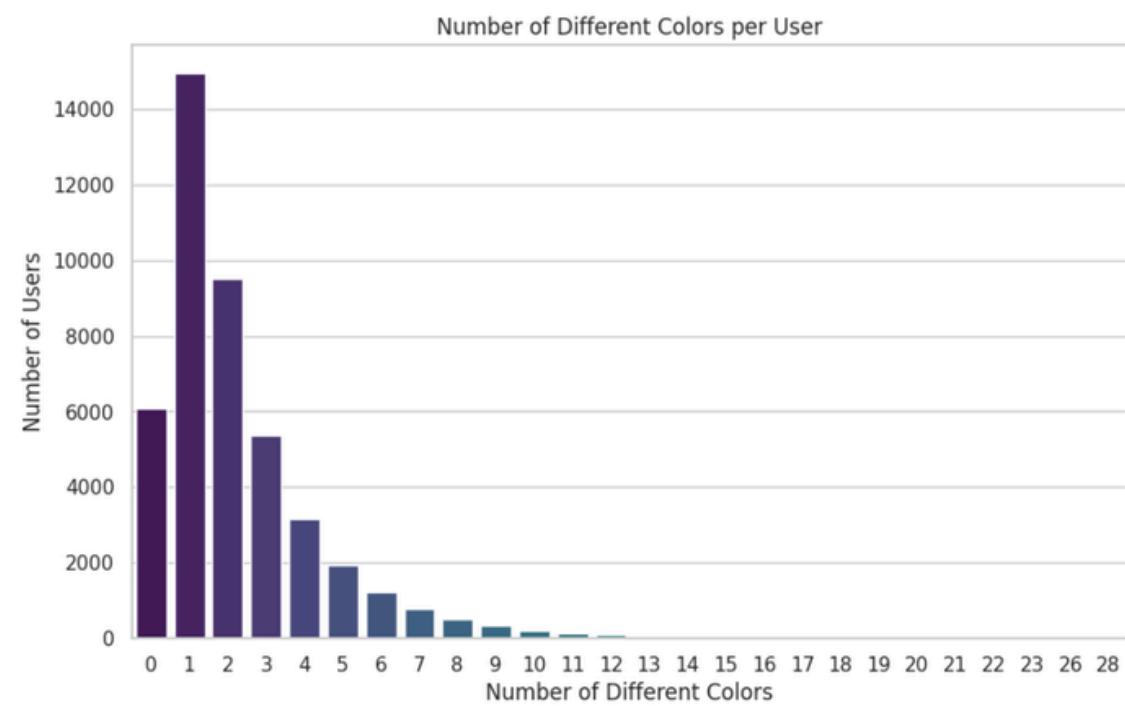
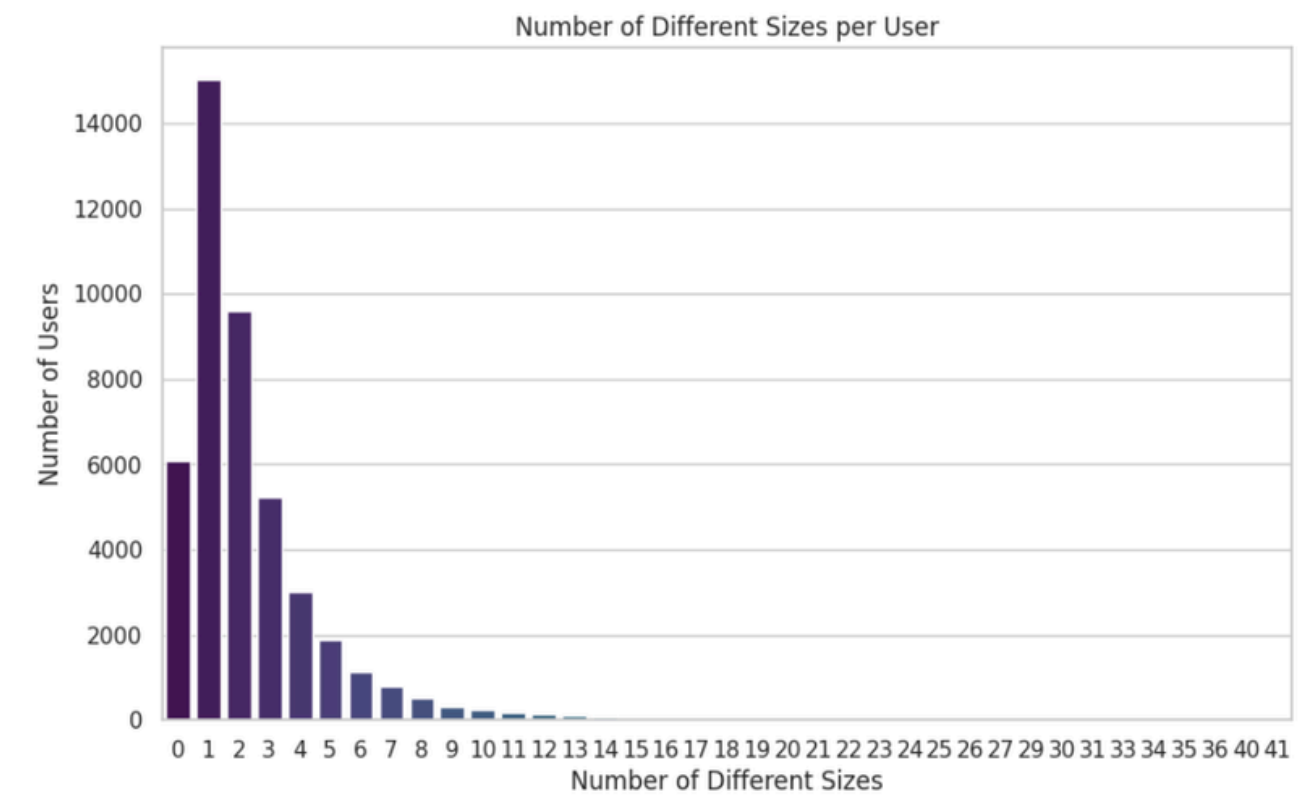
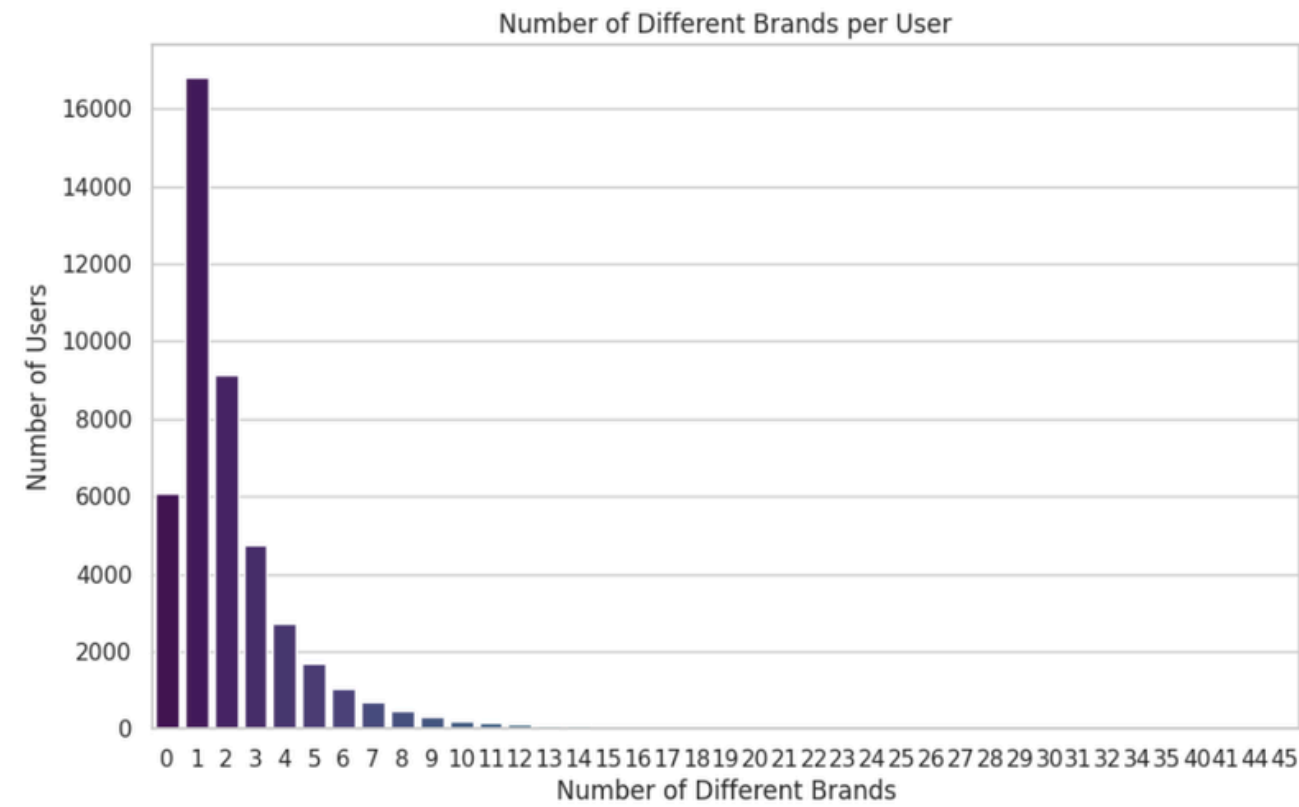
USER TRACES

Environmental Data EDA

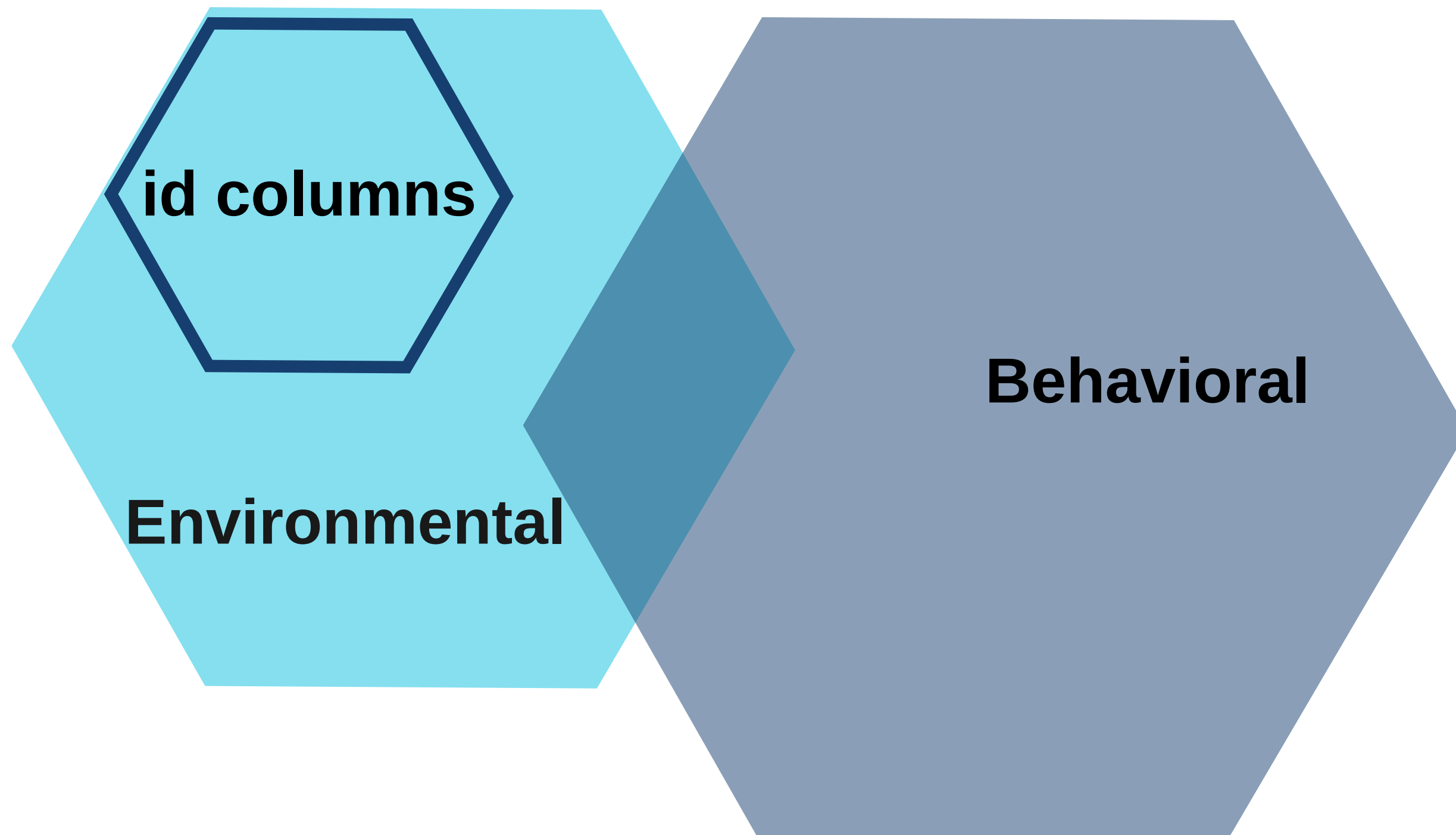


USER TRACES

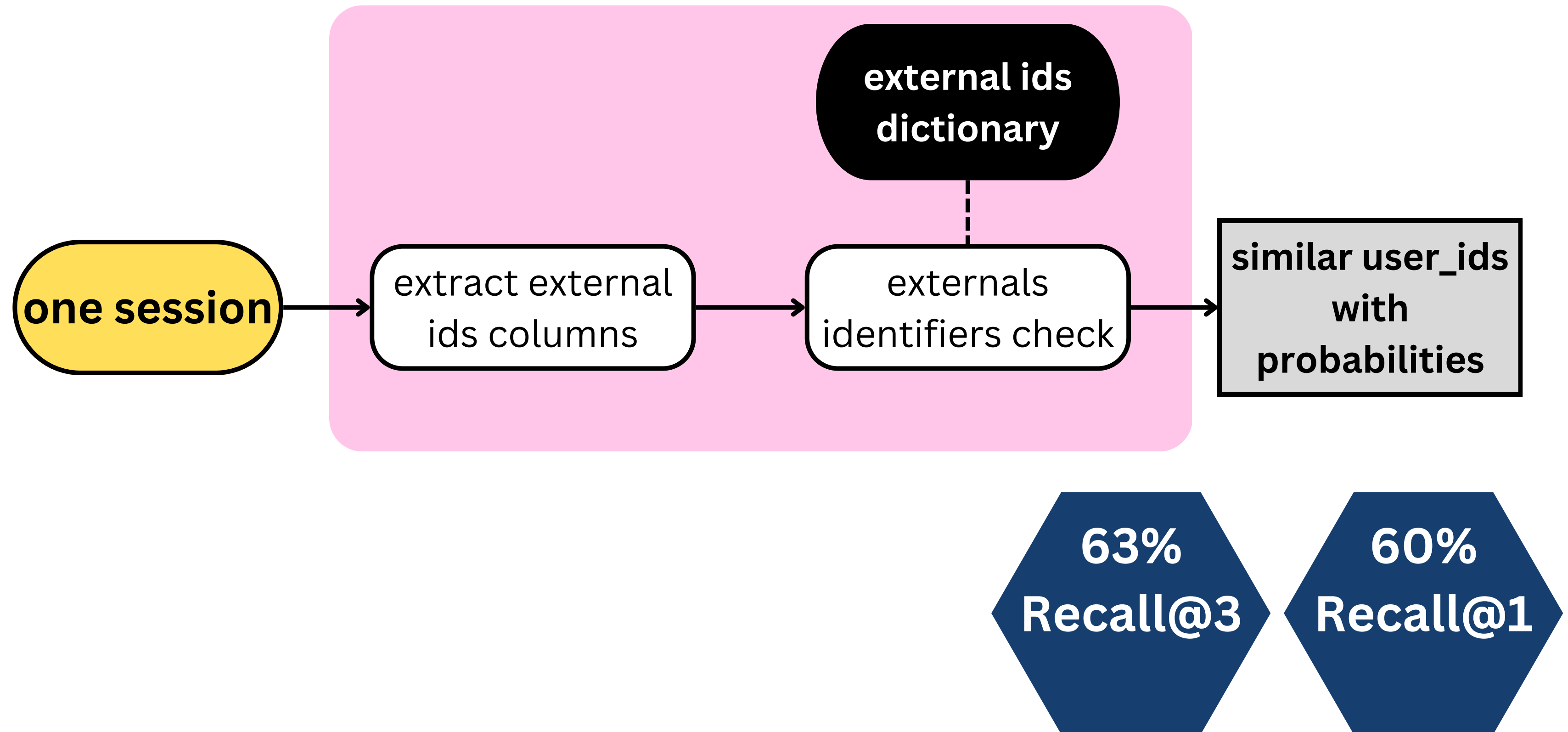
Behavioral Data EDA



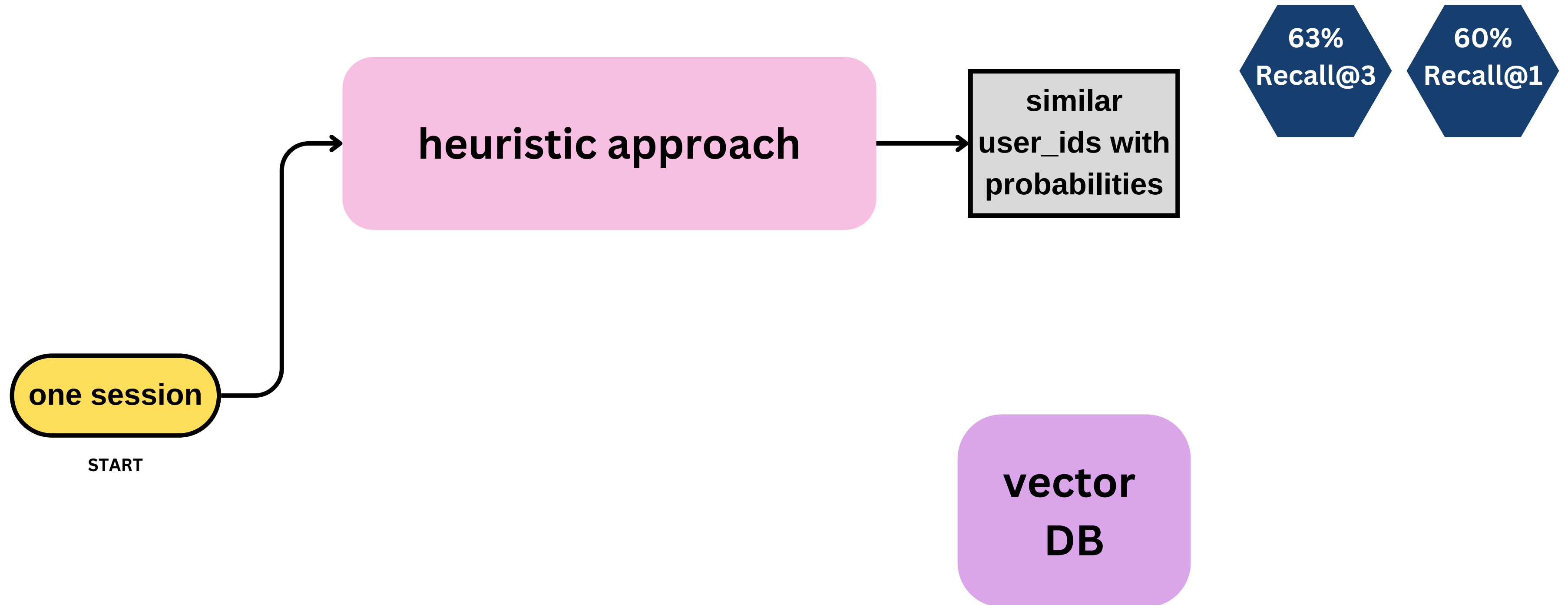
THE DATA



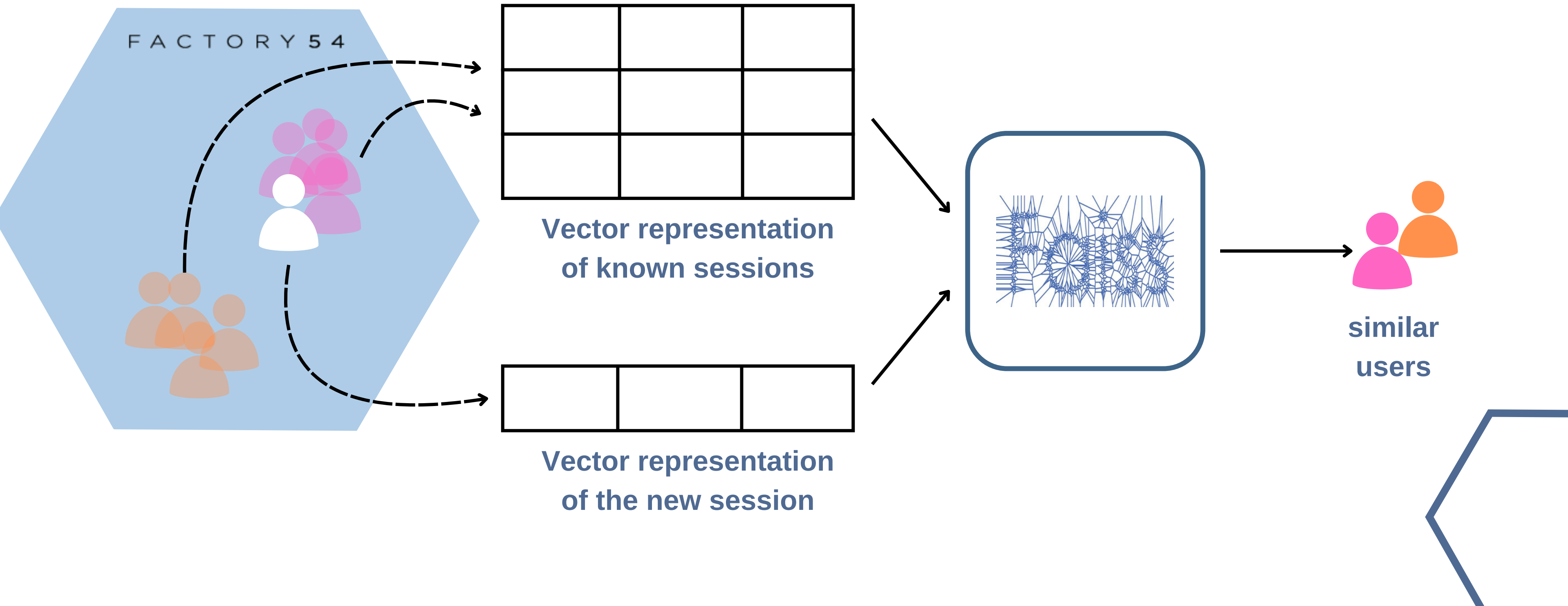
WORKFLOW



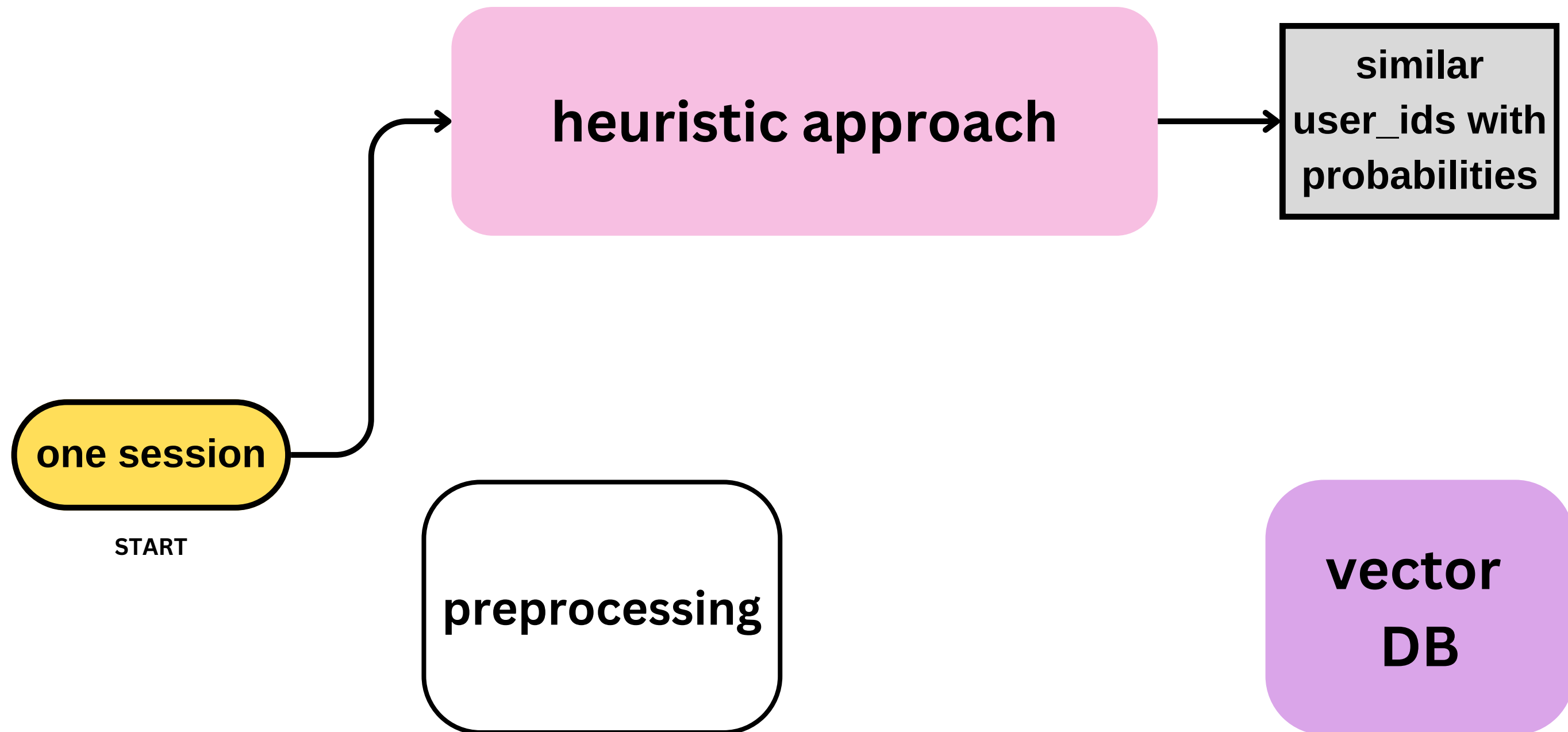
WORKFLOW



FAISS DB



WORKFLOW



DATA PREPROCESSING

business event

```
{{"pageType":"productPage"},
{"events":
  "[
    {"et":"","load","
      \"images\":
[\"/on/demandware.static/-/Sites-master-catalog/default/dw46ad99d3/images/models/893624238_L_2.JPG\",
\"/on/demandware.static/-/Sites-master-catalog/default/dw778a6840/images/large/893624238_P_1.png\",
\"/on/demandware.static/-/Sites-master-catalog/default/dwfcbe4a2c/images/models/893624645 893623745
893624534 893624683 893624238_L T.JPG\",
\"/on/demandware.static/-/Sites-master-catalog/default/dw4f3e7a7d/images/large/893624238_P_2.png\",
\"/on/demandware.static/-/Sites-master-catalog/default/dw73f6e730/images/large/893624238_P_3.png\",
\"/on/demandware.static/-/Sites-master-catalog/default/dwf89e499a/images/large/893624238_P_4.png\"
],
\"productId\":\"893624238\",
\"brand\":\"CALVIN KLEIN\",
\"price\":\"₹ 409.00\",
\"ets\":\"1674580753264\"}],
\"pageType\":\"productPage\"},
{"events":
  "[{et\":\"changeColor\",
\"productId\":\"893624238\",
\"brand\":\"CALVIN KLEIN\",
\"price\":\"₹ 409.00\",
\"previousColor\":\"BEIGE\",
\"pickedColor\":\"BLACK\",
\"ets\":\"1674580756472\"}],
\"pageType\":\"productPage\"},

{"pageType\":\"productPage"}]
```



- size
- top brand
- gender
- language

A decorative graphic in the top-left corner consisting of two overlapping hexagons. The left hexagon is white with a light blue outline, and the right hexagon is solid light blue.

DATA PREPROCESSING

Grouping by session

sum(): document_scroll_height

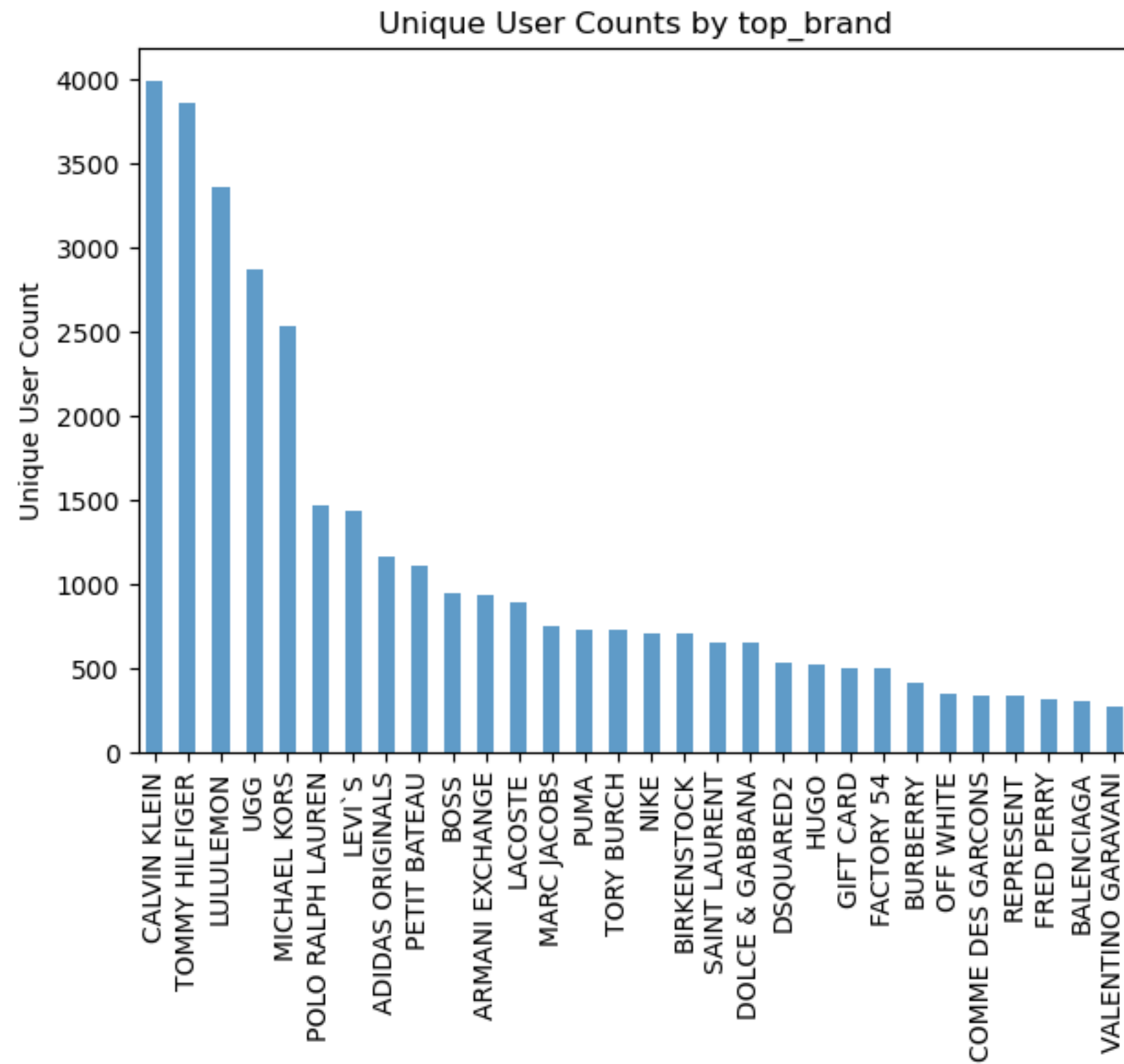
count(): pageview, tabs

mode(): in case of ambiguity during the session

A decorative graphic in the bottom-right corner consisting of a single light blue outlined hexagon.

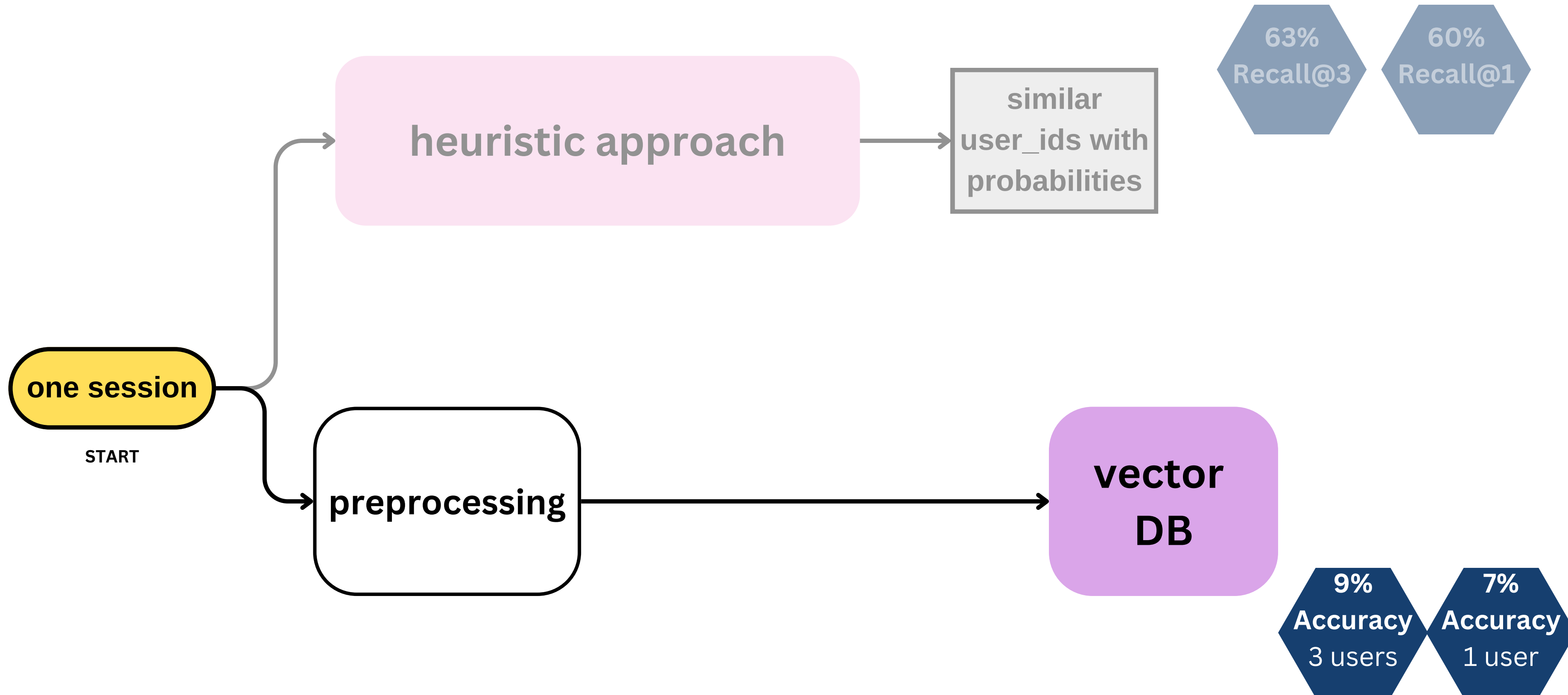
DATA PREPROCESSING

Handling categorical features

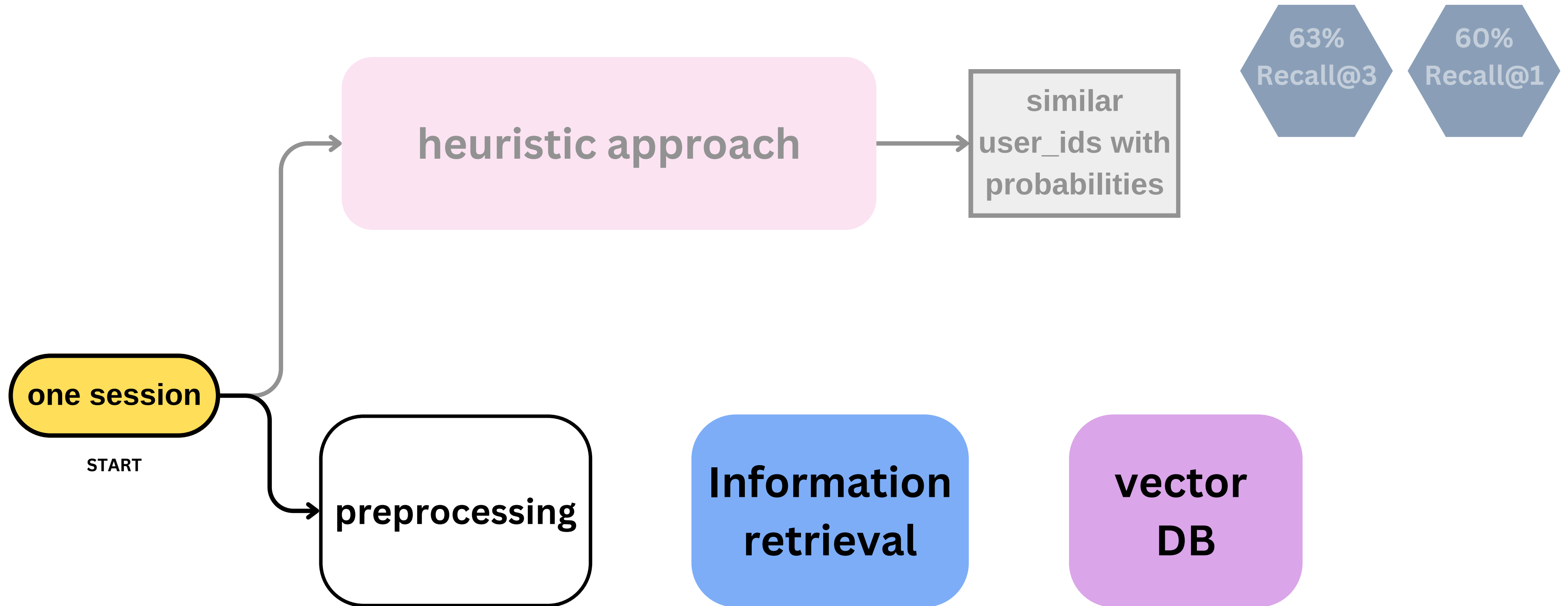


**For nominal categories:
find the top ~10 values that
distinguish users**

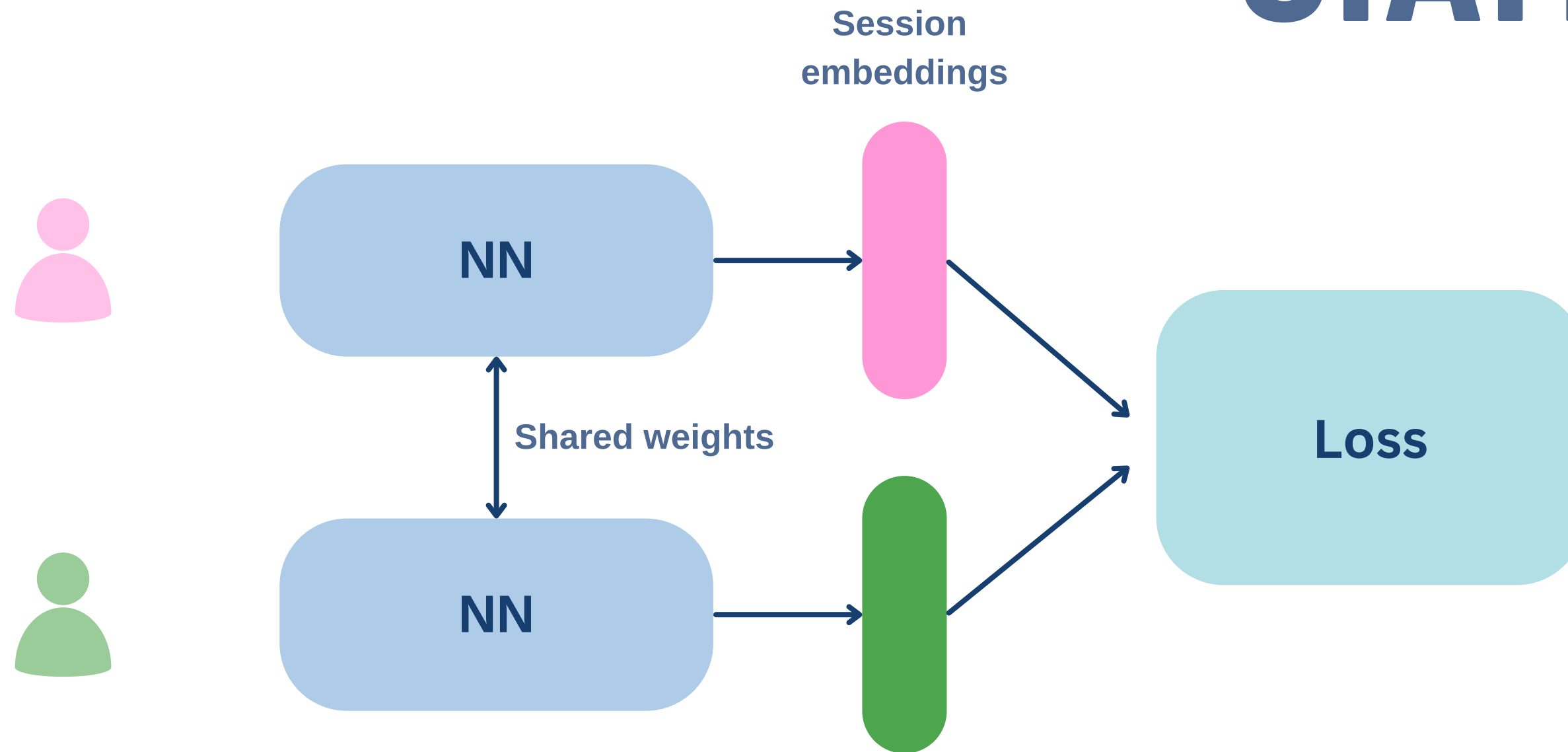
WORKFLOW



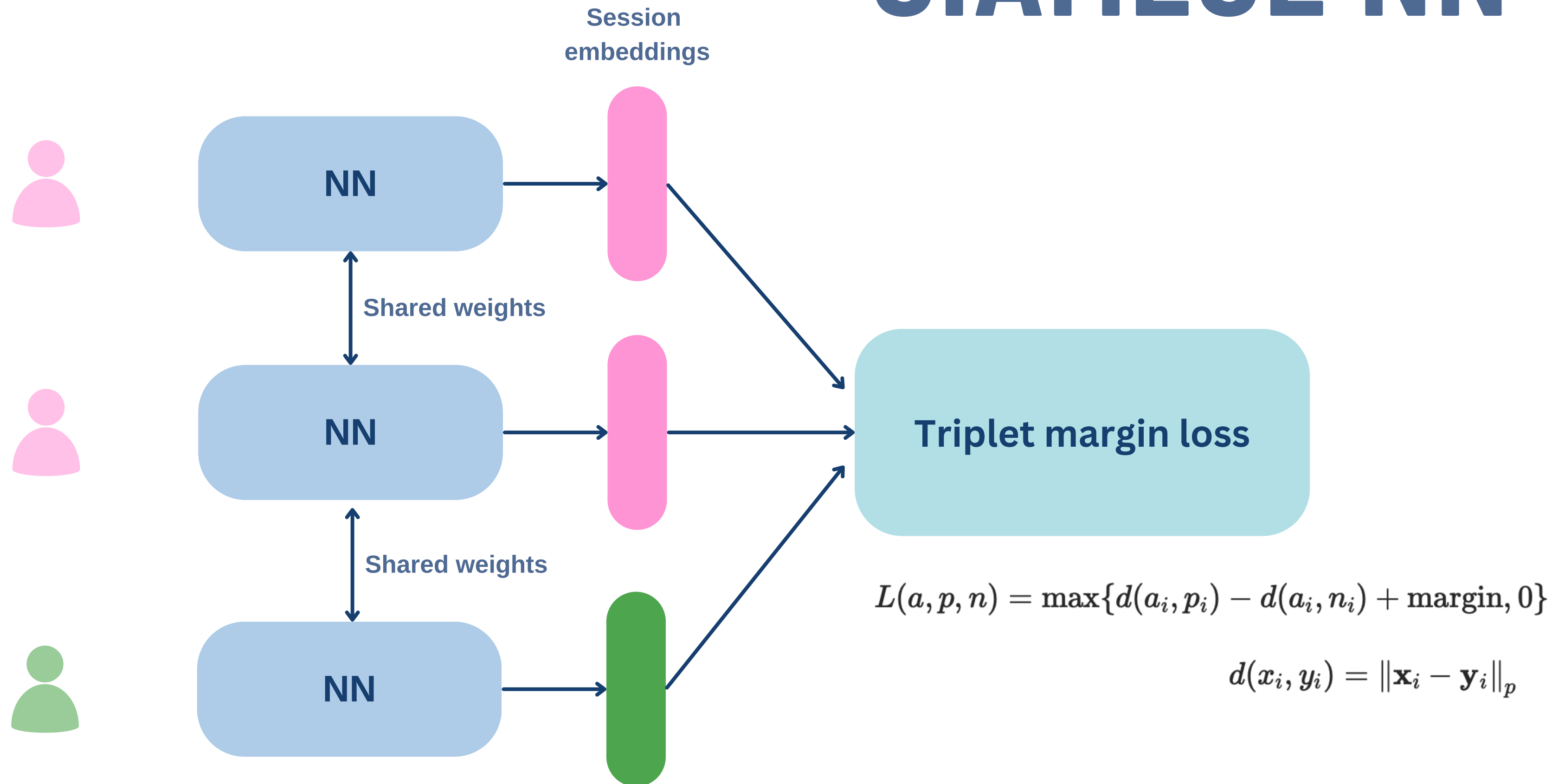
WORKFLOW



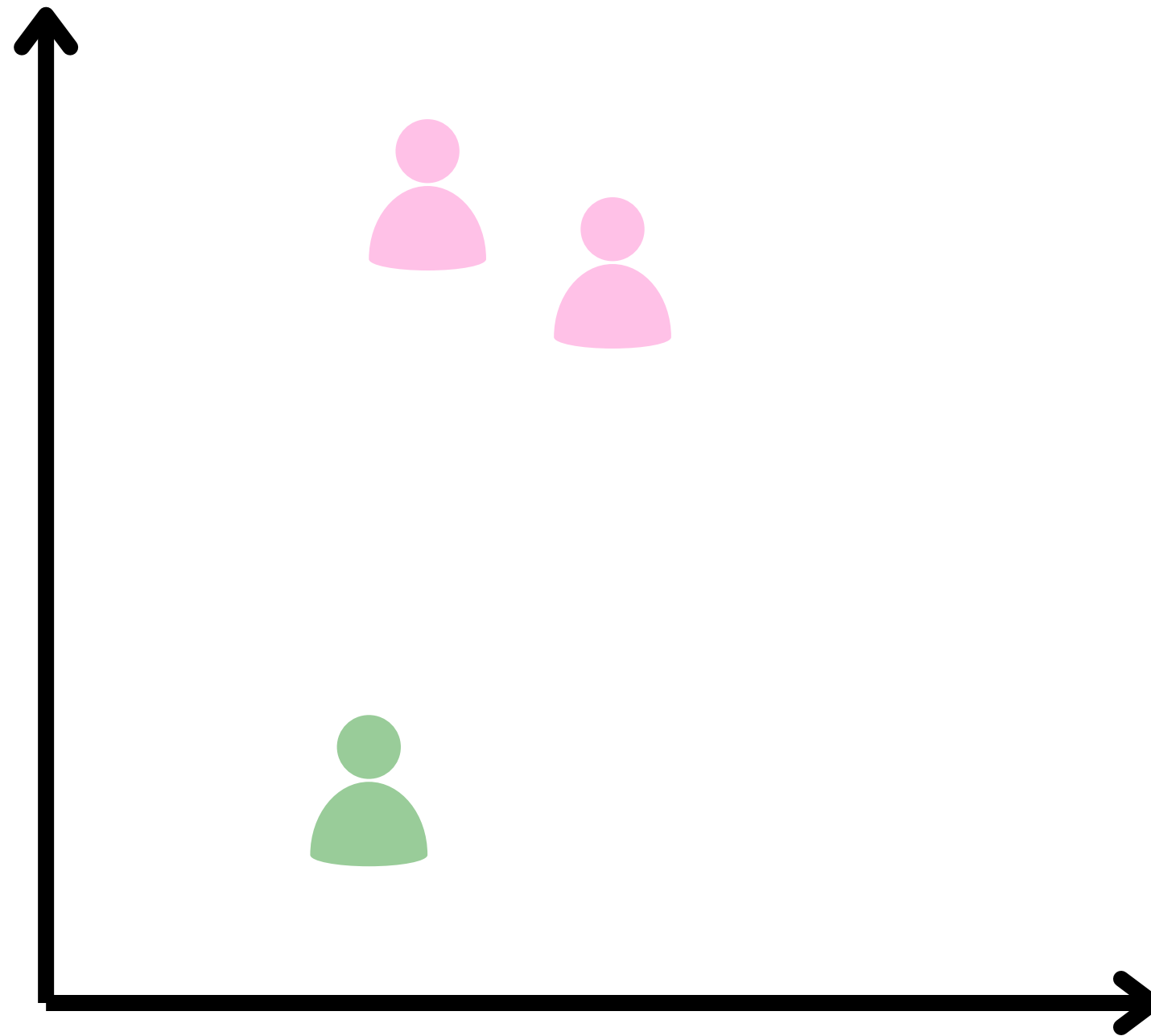
SIAMESE NN



SIAMESE NN



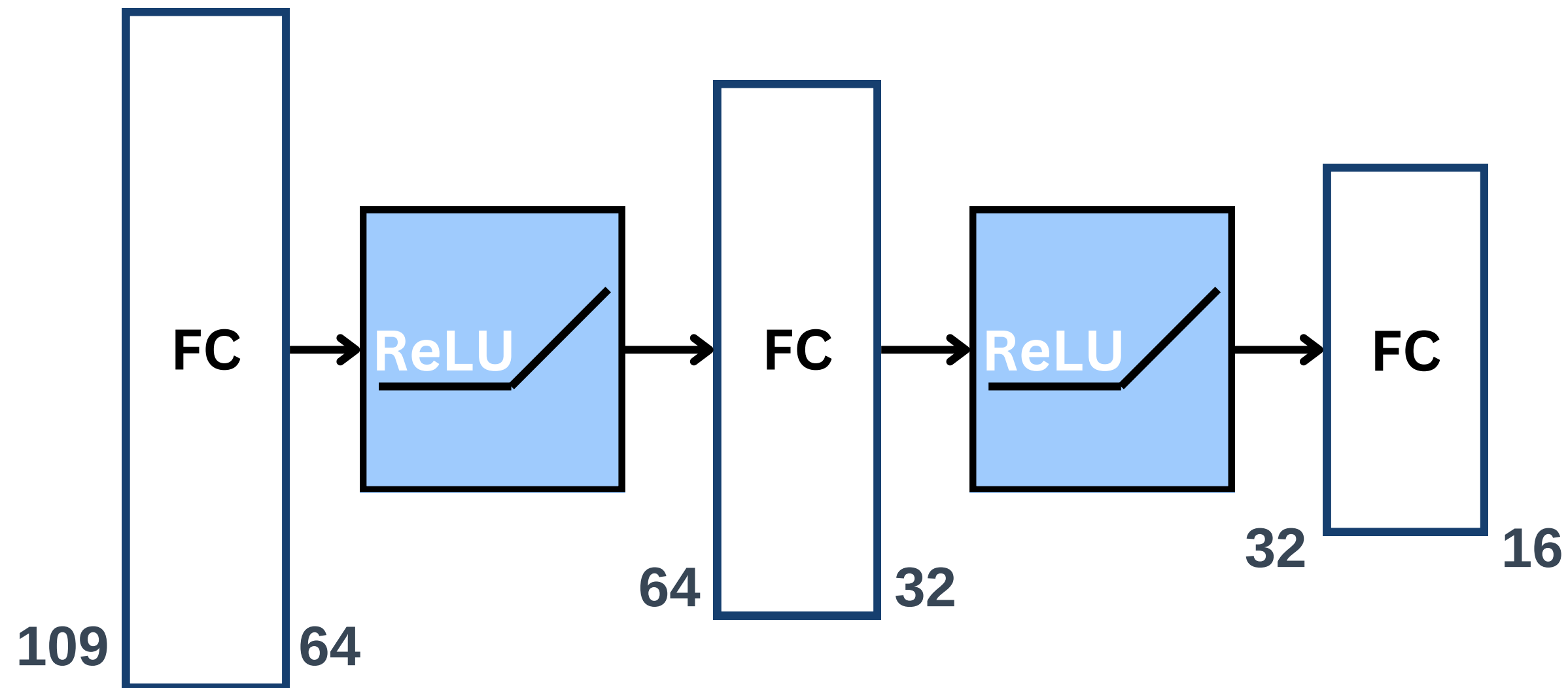
SIAMESE NN



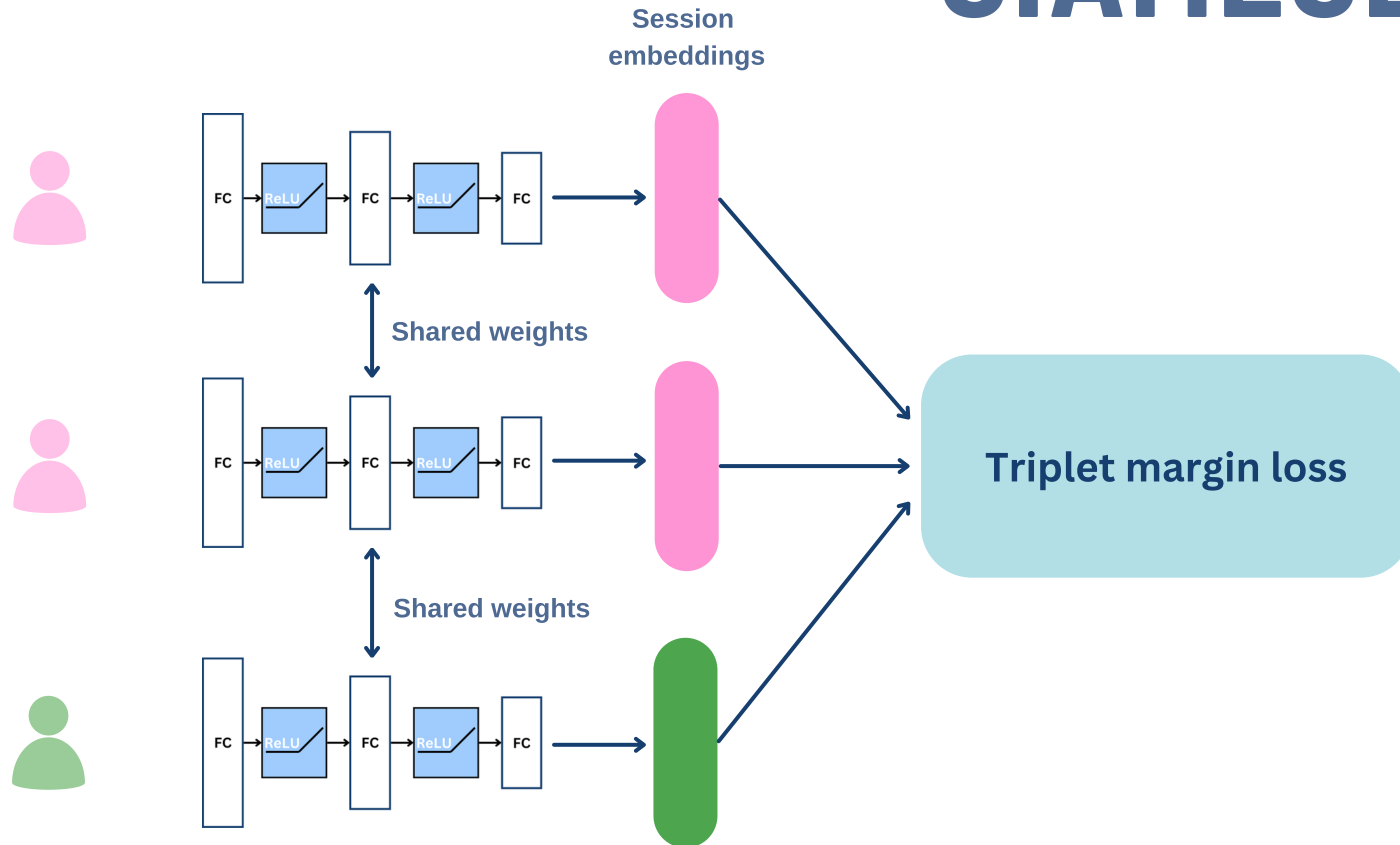
MODEL ARCHITECTURE

Custom architecture

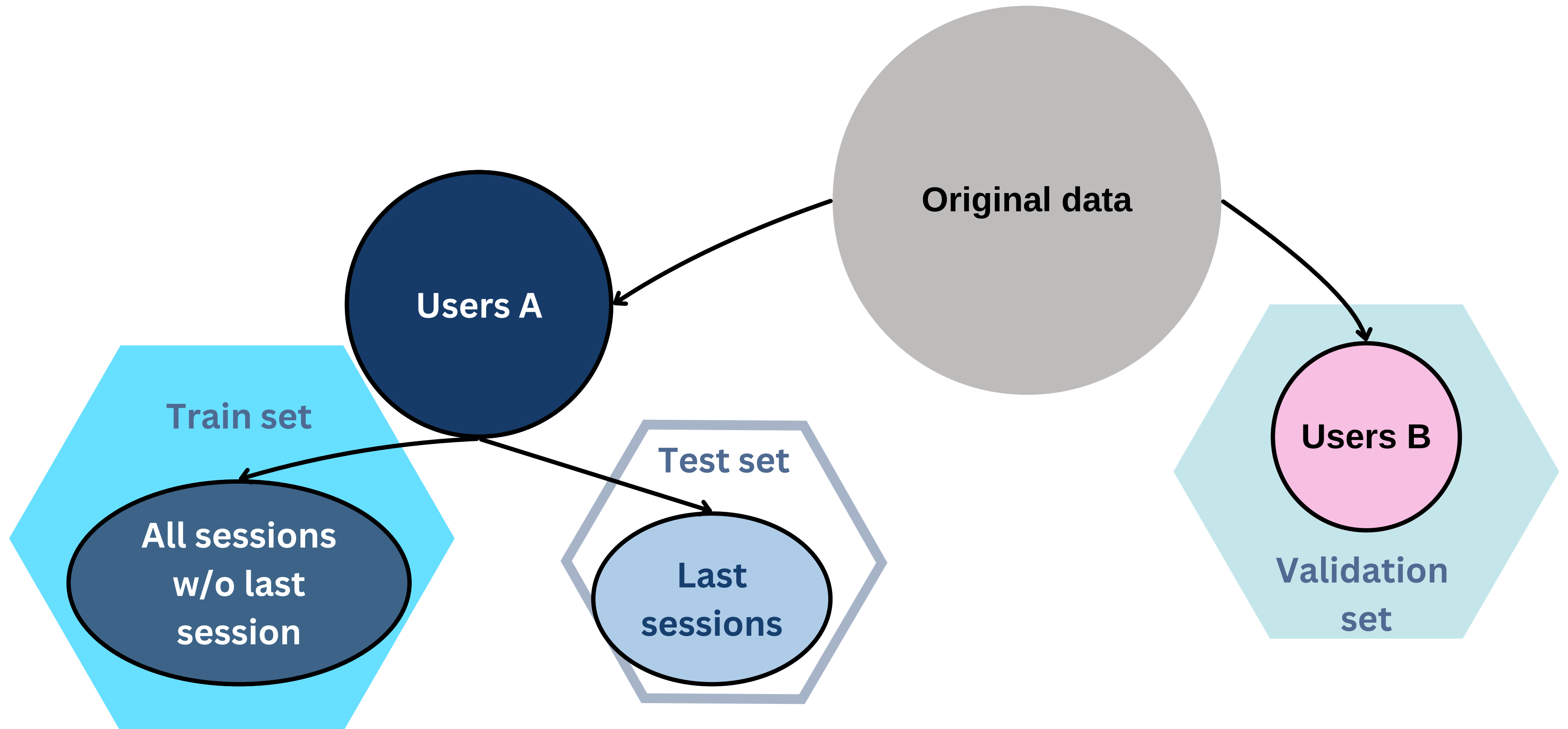
Total params: 8,432



SIAMESE NN

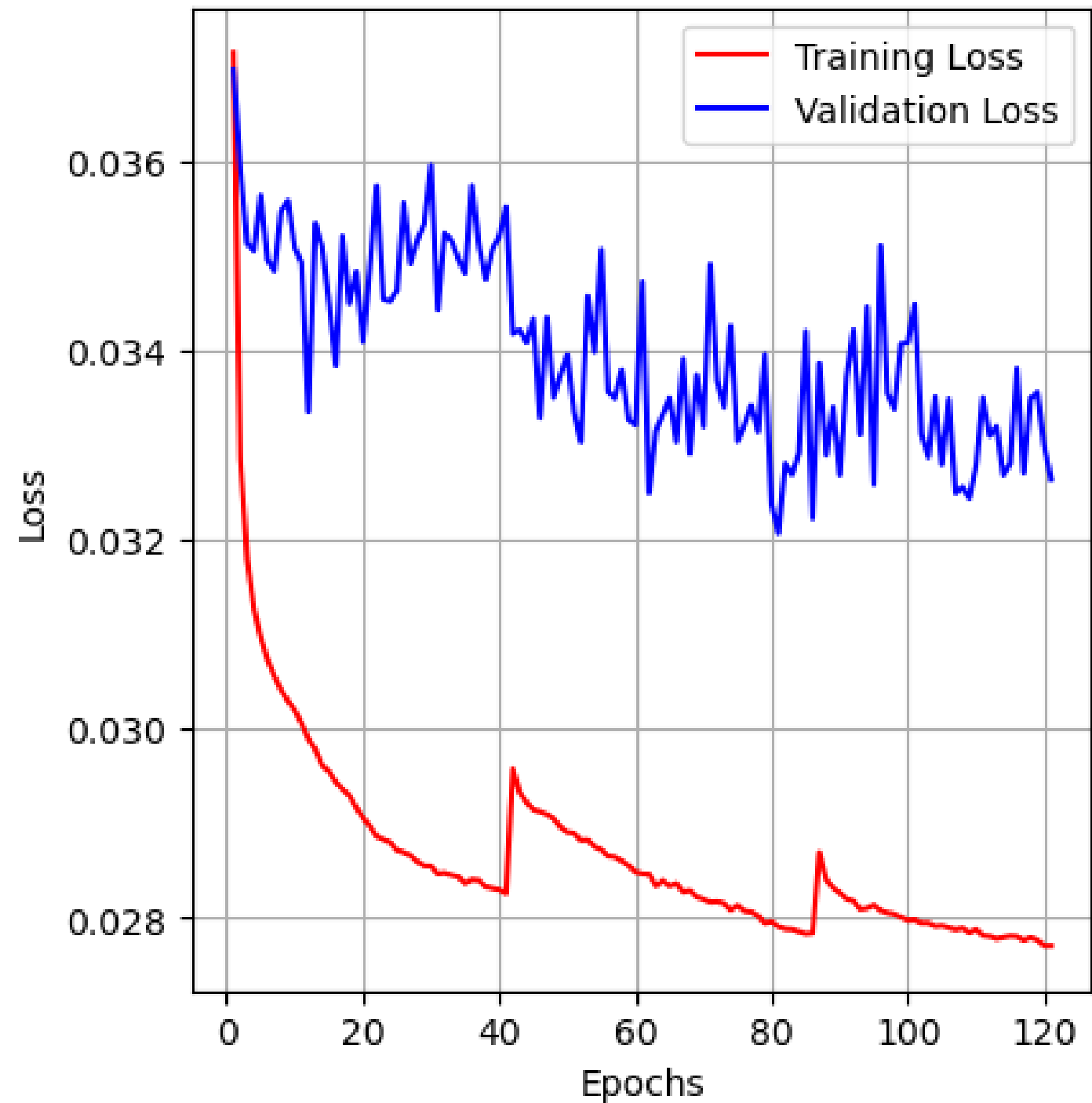


DATASETS

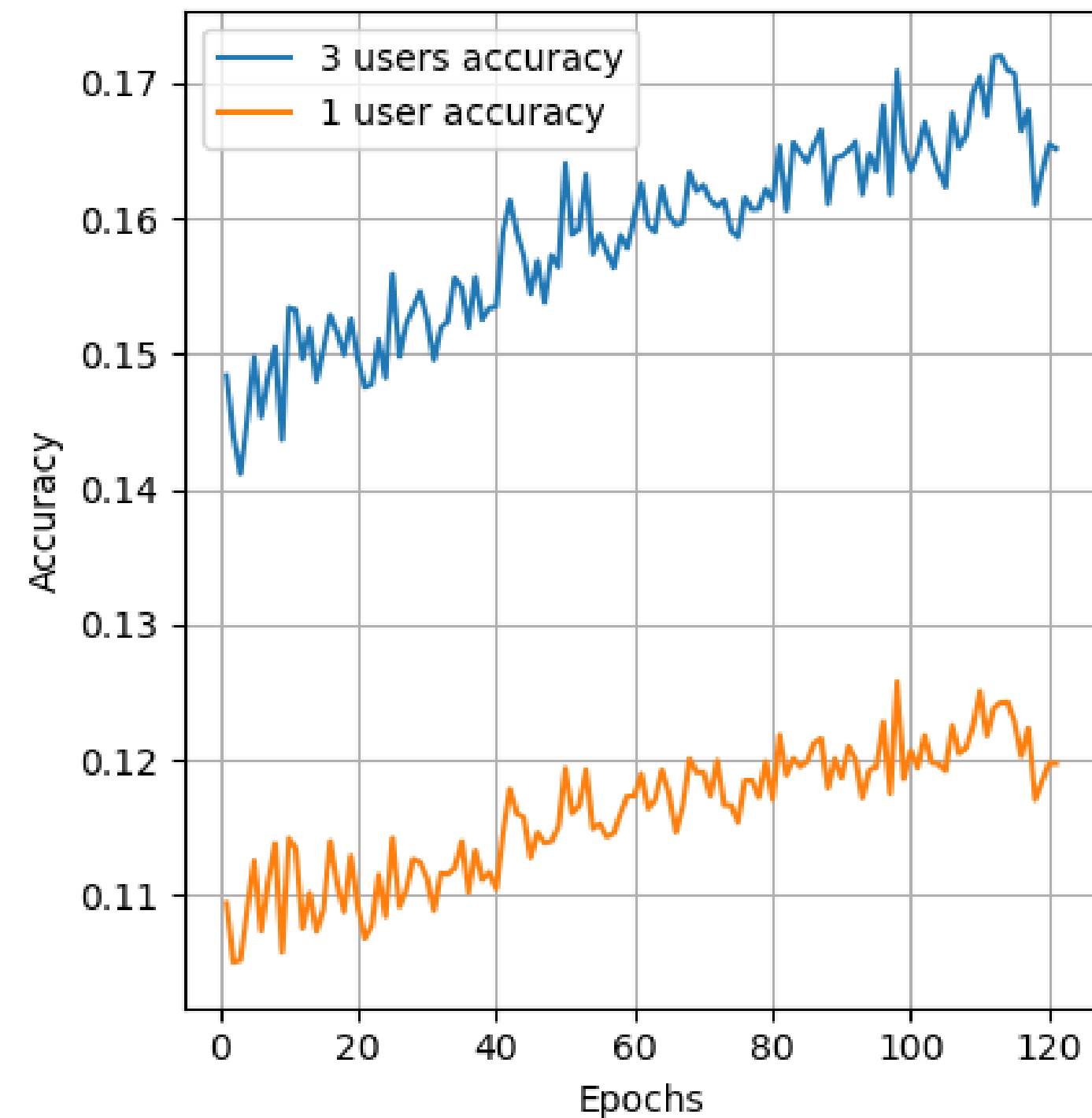


MODEL TRAINING. PROCESS

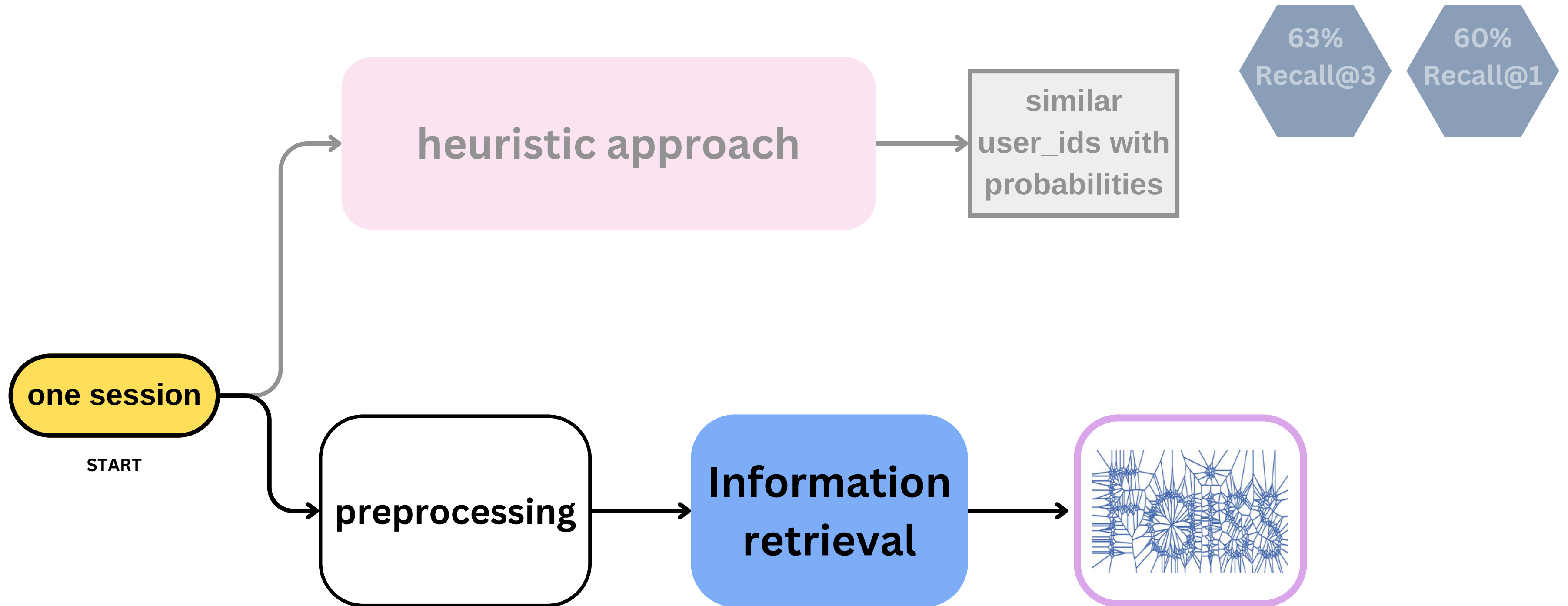
Training and Validation Loss



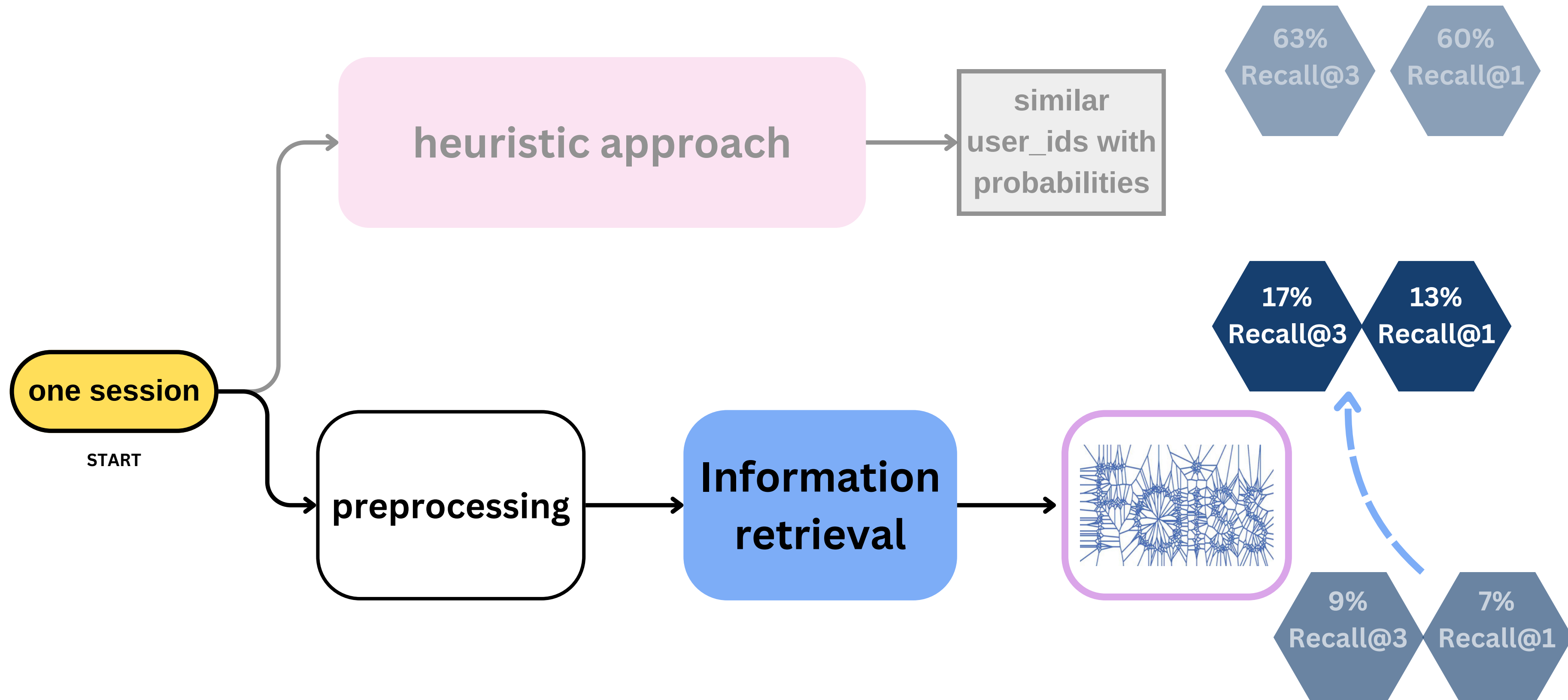
Accuracy Trends During Model Training for the Test set



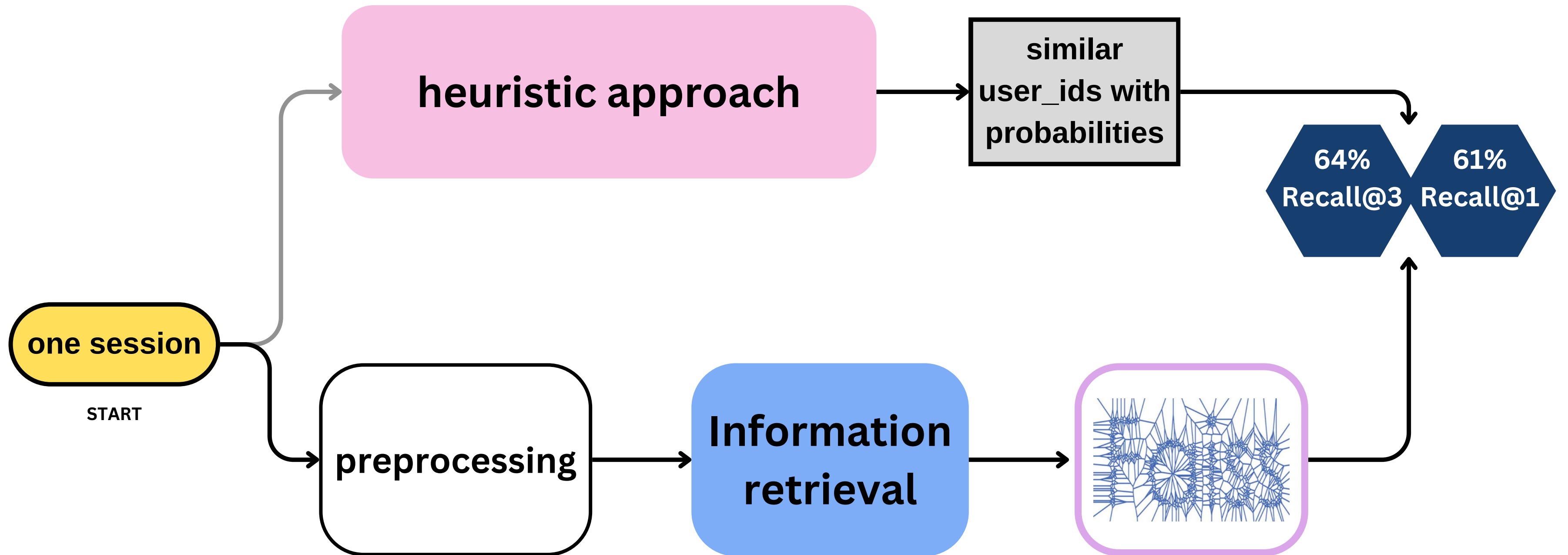
WORKFLOW



WORKFLOW



WORKFLOW



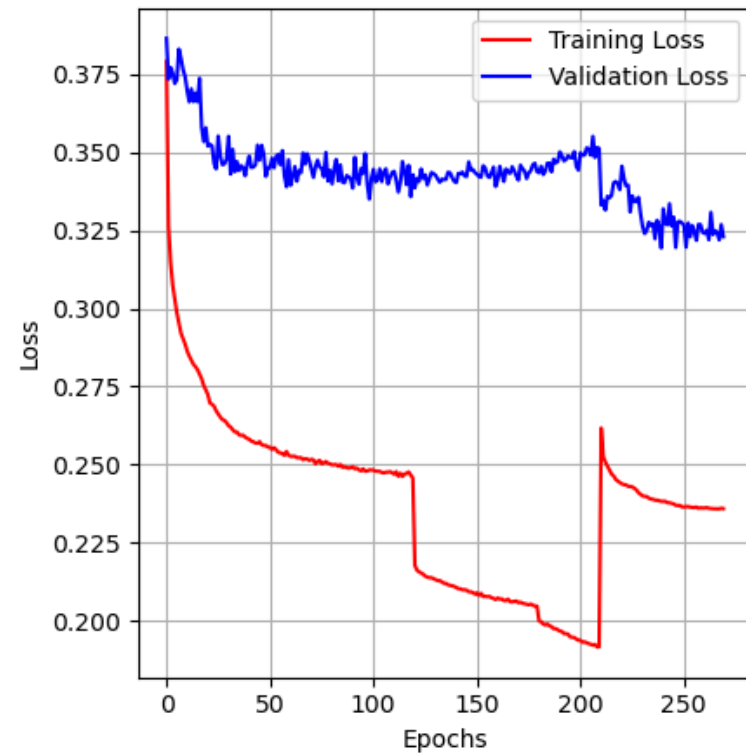


CHALLENGES

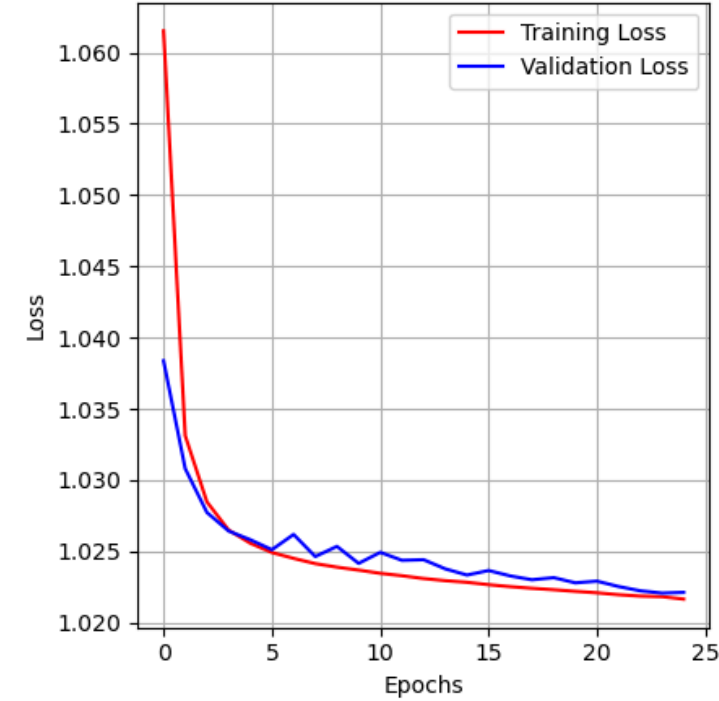
- **Data preprocessing**
- **Handling categorical features**
- **Training NN models from scratch**
 - different model architectures
 - high variation in parameters
 - time consuming

APP. 1. EXPERIMENTS

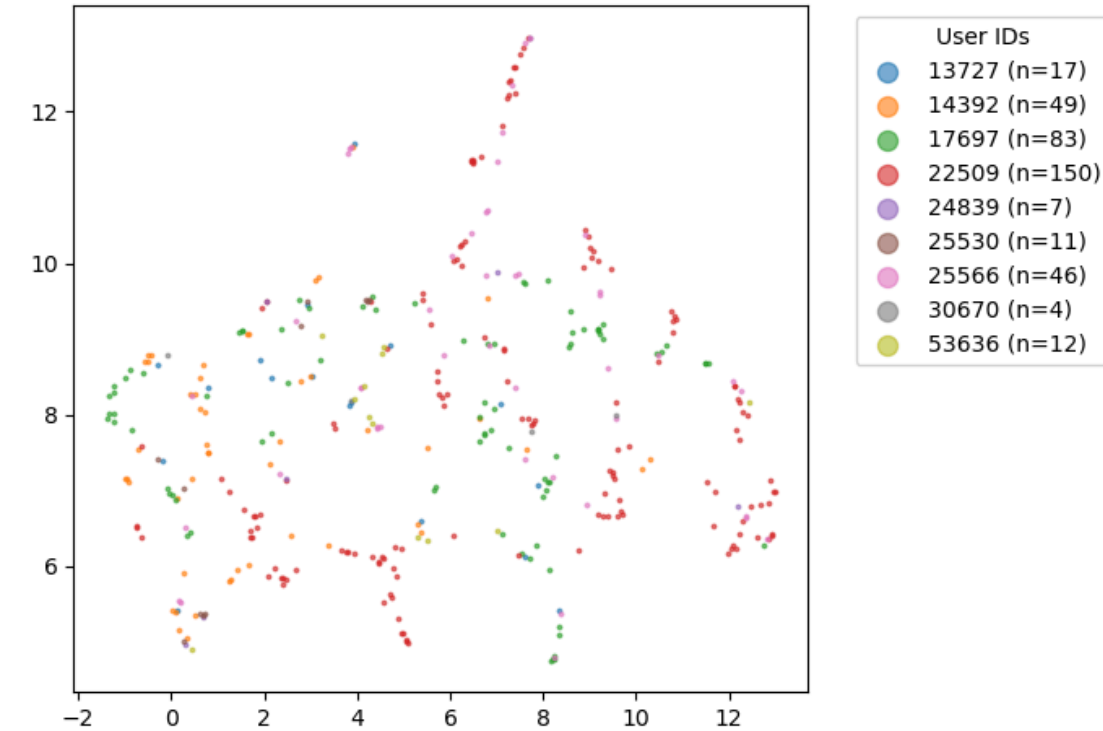
Training and Validation Loss. Custom architecture



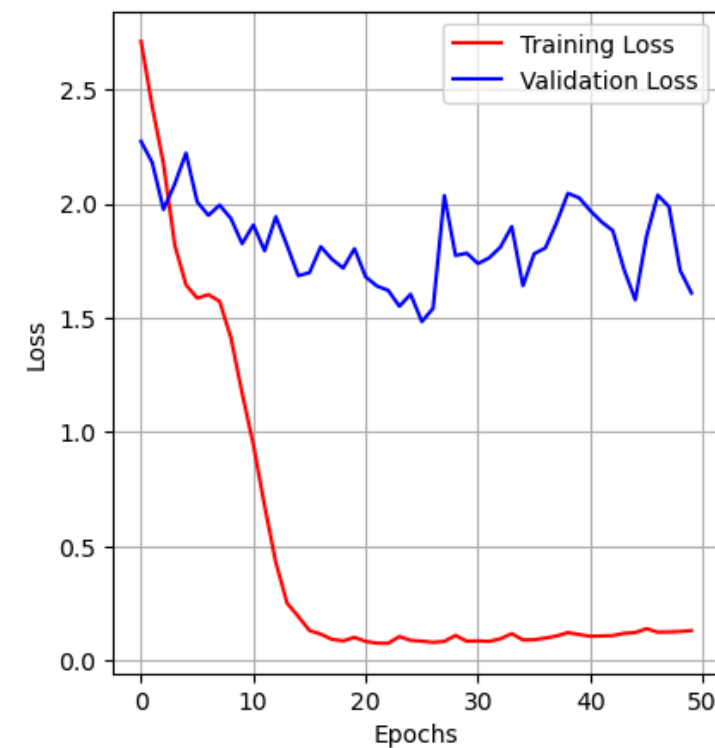
Training and Validation Loss. Custom architecture



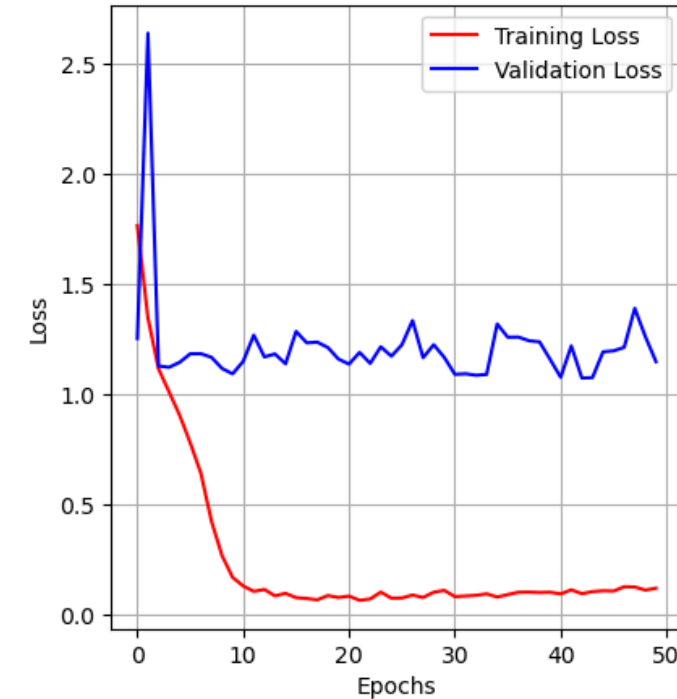
UMAP Dimensionality Reduction



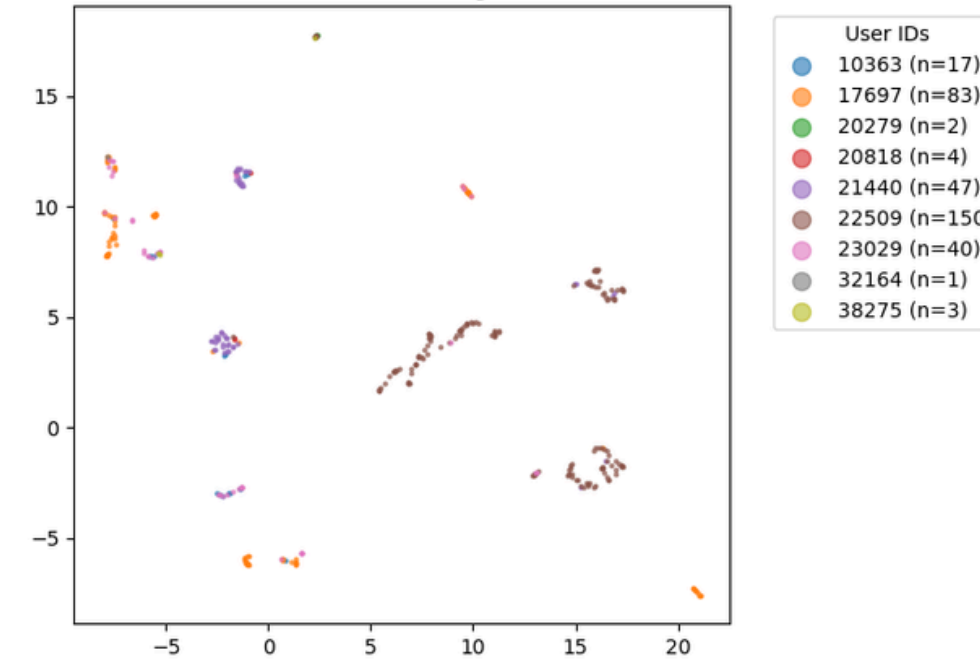
Training and Validation Loss. TabNetEmbeddingModel triplets



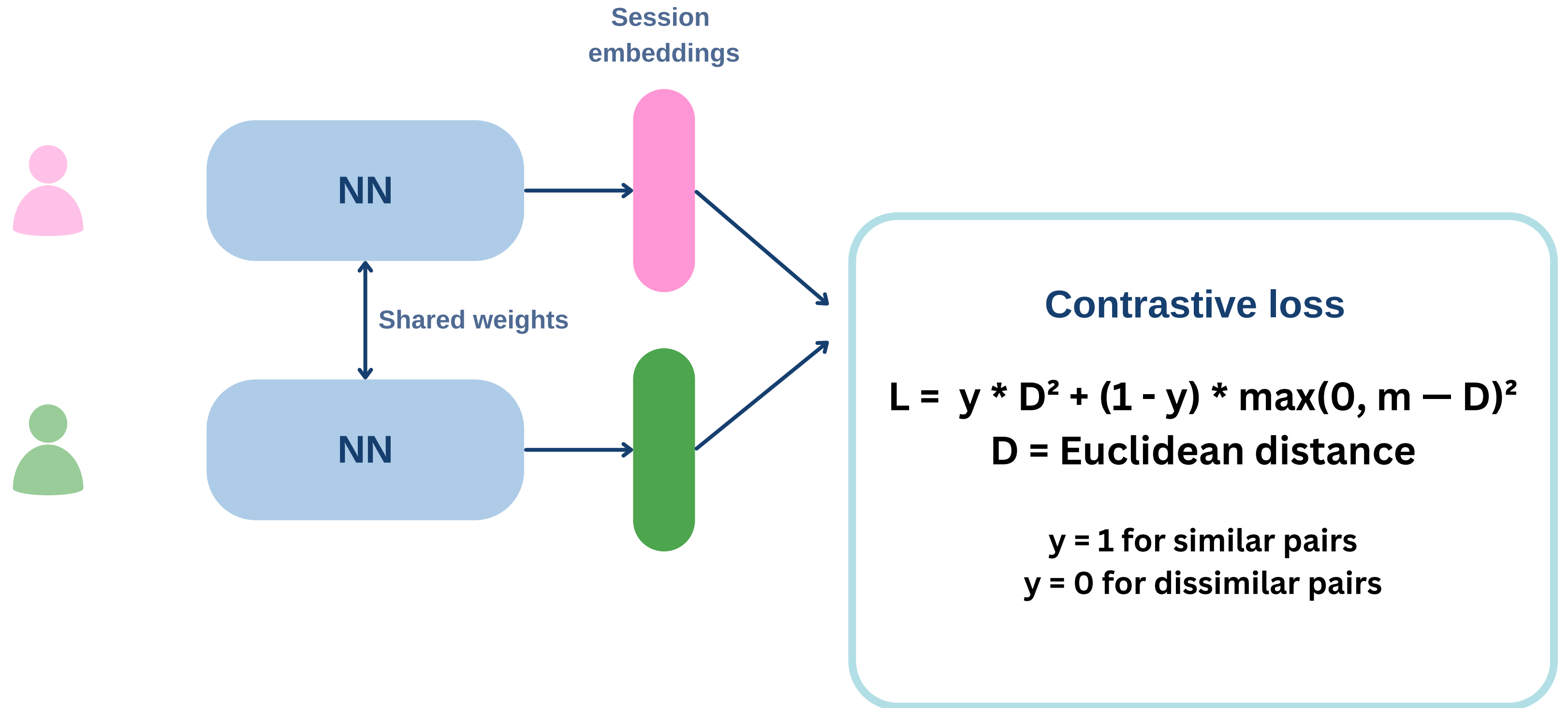
Training and Validation Loss. TabNetEmbeddingModel triplets



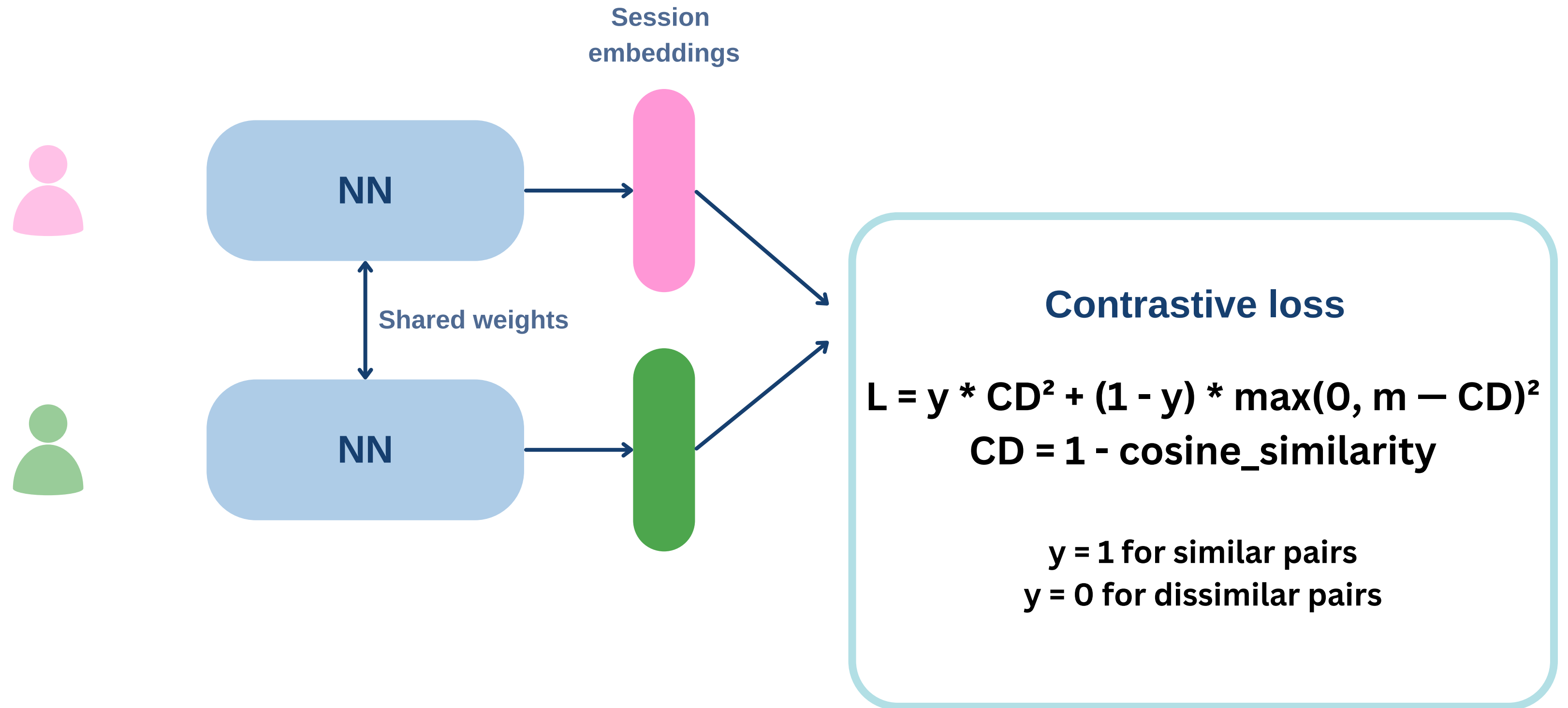
UMAP Dimensionality Reduction



APP. 2. SIAMESE NN



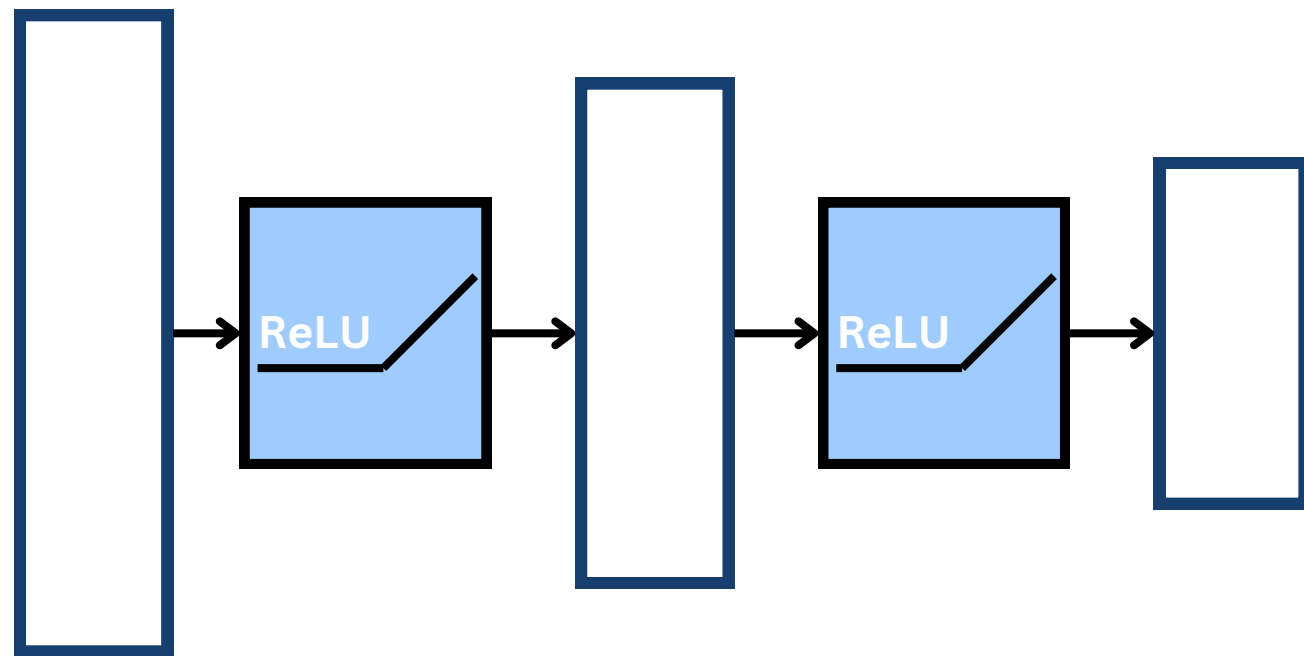
APP. 2. SIAMESE NN



APP. 3. MODEL ARCHITECTURES

Custom architecture

Total params: 8,432



TabNet architecture

Total params: 15,092

