# DIN based Look-Alike model

# Workflow
## - Necessary steps

Input →

Kernel computation →

LookAlike model →

| Inactive user elimination |
| --- |

| User feature | User behavior (represented as keywords history) |
| --- | --- |

| DIN model |
| --- |

| User profile generation – user vs keyword correlation |
| --- |

| User profile normalization |
| --- |

| seed_user vs non-seed_user similarity measurements |
| --- |

# Inactive user elimination (user prescreen)

All Users

↓

User's traffic contribution > predefined threshold*

↓
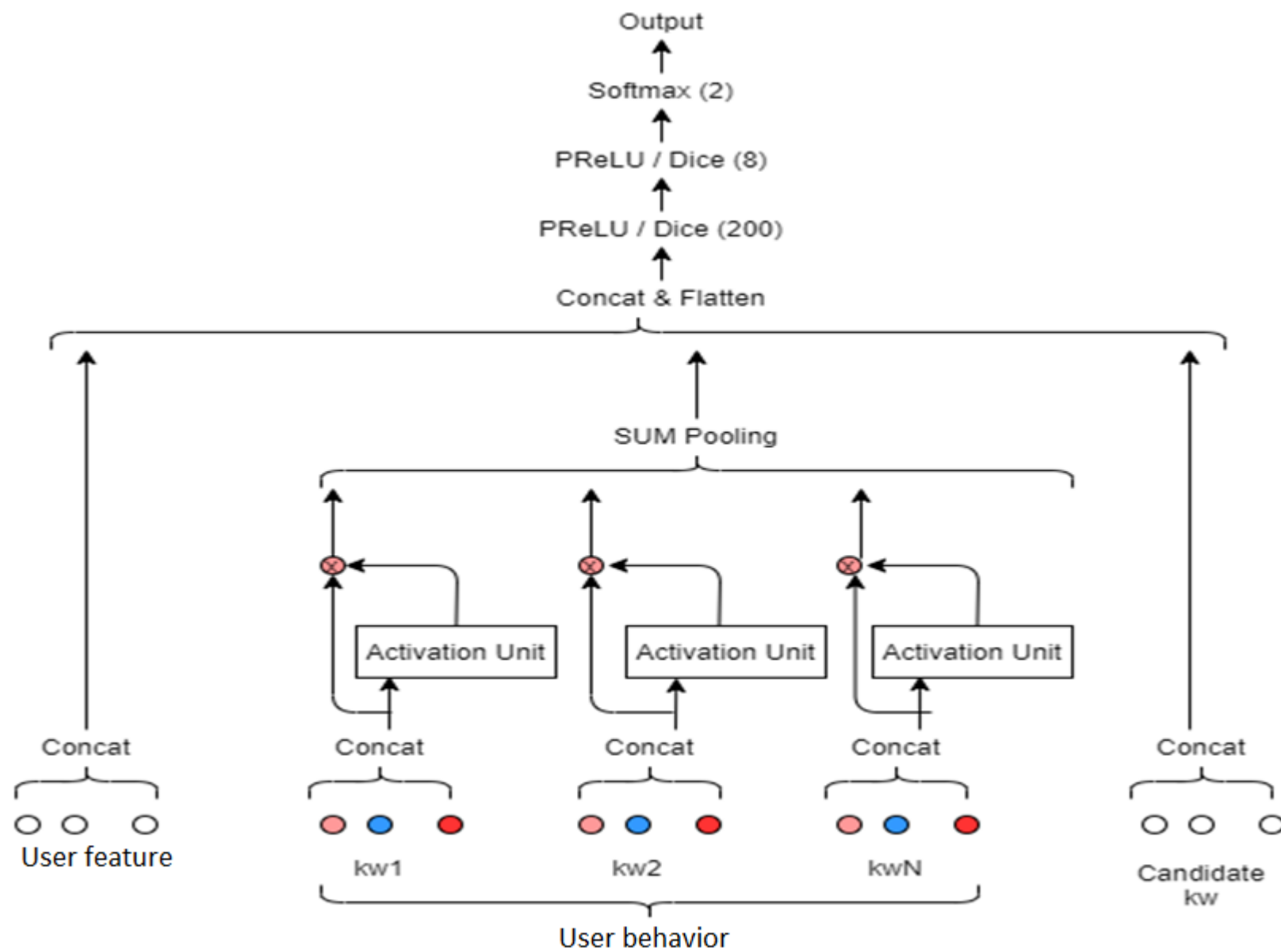
Active Users

* "Prefined threshold" is defined as a range of normal traffic (with low and high bounds) to eliminate:
1. users with consistent low traffic (inactive user, traffic < low bound)
2. users with extremely high traffic for some specific period (robot user, traffic > high bound)

# DIN Model

# DIN Model Output
## – user vs keyword correlation (user profile generation)

| | Keyword$_1$ | Keyword$_2$ | Keyword$_3$ | Keyword$_4$ | ... | Keyword$_m$ |
|---|---|---|---|---|---|---|
| User$_1$ | score$_{11}$ | score$_{12}$ | score$_{13}$ | score$_{14}$ | ... | score$_{1m}$ |
| User$_2$ | score$_{21}$ | score$_{22}$ | score$_{23}$ | score$_{24}$ | ... | score$_{2m}$ |
| ... | ... | ... | ... | ... | ... | ... |
| User$_n$ | score$_{n1}$ | score$_{n2}$ | score$_{n3}$ | score$_{n4}$ | ... | score$_{nm}$ |

← DIN

# DIN Model Output
## – user profile normalization

DIN

| | Keyword$_1$ | Keyword$_2$ | Keyword$_3$ | Keyword$_4$ | ... | Keyword$_m$ |
|---|---|---|---|---|---|---|
| User$_1$ | score$_{11}$ | score$_{12}$ | score$_{13}$ | score$_{14}$ | ... | score$_{1m}$ |
| User$_2$ | score$_{21}$ | score$_{22}$ | score$_{23}$ | score$_{24}$ | ... | score$_{2m}$ |
| ... | ... | ... | ... | ... | ... | ... |
| User$_n$ | score$_{n1}$ | score$_{n2}$ | score$_{n3}$ | score$_{n4}$ | ... | score$_{nm}$ |

| Normalization constant |
|---|
| C$_1$ |
| C$_2$ |
| ... |
| C$_n$ |

Score normalization

| | Keyword$_1$ | Keyword$_2$ | Keyword$_3$ | Keyword$_4$ | ... | Keyword$_m$ |
|---|---|---|---|---|---|---|
| User$_1$ | norm_score$_{11}$ | norm_score$_{12}$ | norm_score$_{13}$ | norm_score$_{14}$ | ... | norm_score$_{1m}$ |
| User$_2$ | norm_score$_{21}$ | norm_score$_{22}$ | norm_score$_{23}$ | norm_score$_{24}$ | ... | norm_score$_{2m}$ |
| ... | ... | ... | ... | ... | ... | ... |
| User$_n$ | norm__Score$_{n1}$ | norm_score$_{n2}$ | norm_score$_{n3}$ | norm_score$_{n4}$ | ... | norm_score$_{nm}$ |

$$norm\_score_{ij} = \frac{score_{ij}}{C_i}$$

$$C_i = \sqrt{\sum_{j=1}^{m} score_{ij}^2}$$

# DIN Model Output
## – user similarity measurement

$User_i's\ normalized\ profile$:

$$S_i = \{norm\_score_{i1},\ \ norm\_score_{i2},\ \ \ldots\ \ norm\_score_{im}\}$$

$Cross\ user\ similarity$:

$$Similarity(S_i,\ \ S_j) = S_i\ \cdot\ S_j = \sum_{k=1}^{m} norm\_score_{ik} \times norm\_score_{jk}$$

# DIN based Look-Alike model
## – seed_user vs non-seed_user similarity measure

| | Seed_user$_1$ | Seed_user$_2$ | …… | Seed_user$_m$ |
|---|---|---|---|---|
| Nonseed_user$_1$ | Similary$_{11}$ | Similary$_{12}$ | …… | Similary$_{1m}$ |
| Nonseed_user$_2$ | Similary$_{21}$ | Similary$_{22}$ | …… | Similary$_{2m}$ |
| Nonseed_user$_3$ | Similary$_{31}$ | Similary$_{32}$ | …… | Similary$_{3m}$ |
| Nonseed_user$_4$ | Similary$_{41}$ | Similary$_{42}$ | …… | Similary$_{4m}$ |
| …… | …… | …… | …… | …… |
| Nonseed_user$_n$ | Similary$_{n1}$ | Similary$_{n2}$ | …… | Similary$_{nm}$ |

Parallel computed and only maximum value for each row need to be stored

| All Seed Users |
|---|
| $\underset{i}{mean}(top10\ similarity_{1i})$ |
| $\underset{i}{mean}(top10\ similarity_{2i})$ |
| $\underset{i}{mean}(top10\ similarity_{3i})$ |
| $\underset{i}{mean}(top10\ similarity_{4i})$ |
| …… |
| $\underset{i}{mean}(top10\ similarity_{ni})$ |
| sort |

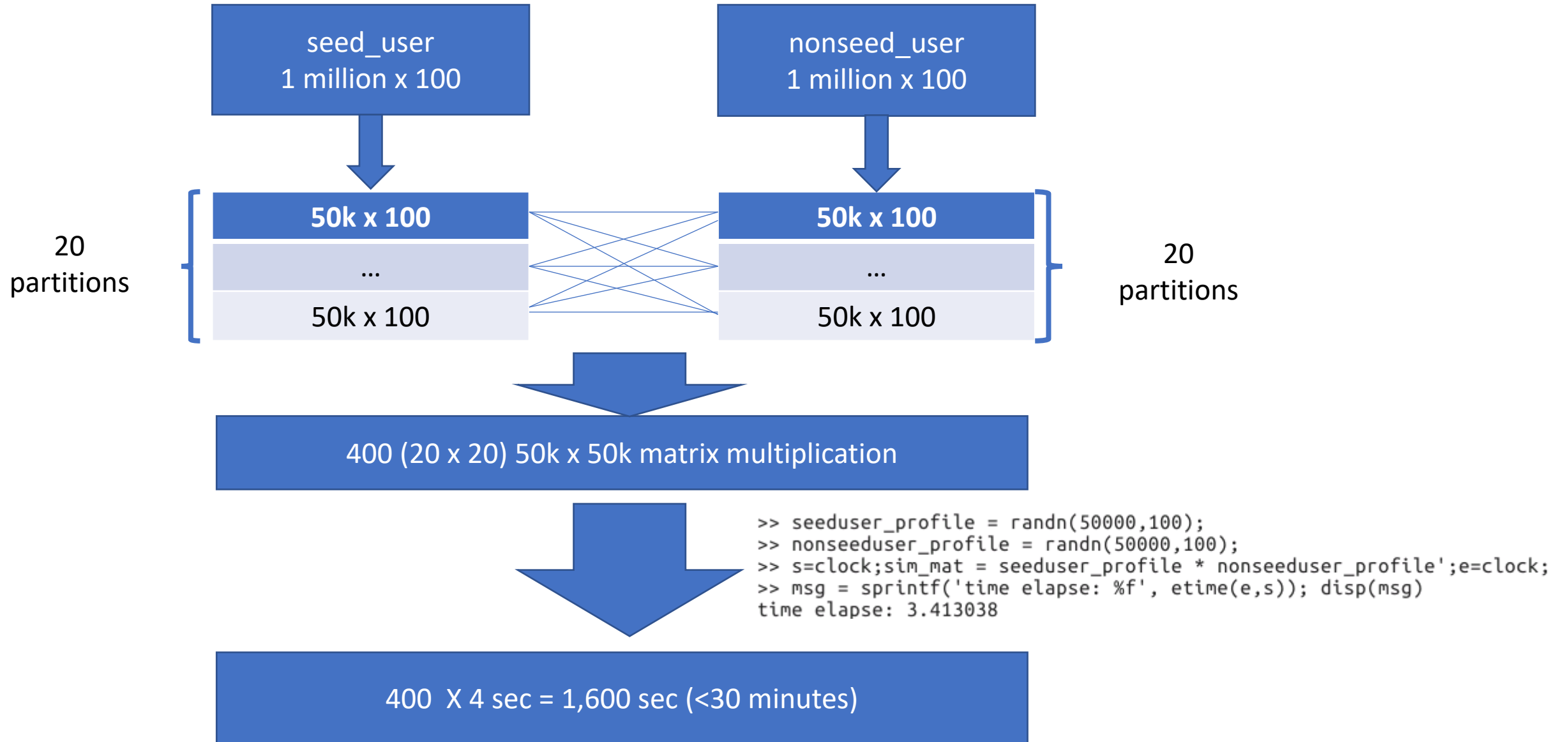| Rank$_1$ nonseed_user |
|---|
| Rank$_2$ nonseed_user |
| Rank$_3$ nonseed_user |
| Rank$_4$ nonseed_user |
| … |
| Rank$_n$ nonseed_user |

# Similarity computation estimation

$$M_{seed} = \begin{bmatrix} norm\_score_{1,1} & \cdots & norm\_score_{1,m} \\ \cdots & \cdots & \cdots \\ norm\_score_{n_{seed},1} & \cdots & norm\_score_{n_{seed},m} \end{bmatrix}$$

$$M_{nonseed} = \begin{bmatrix} norm\_score_{1,1} & \cdots & norm\_score_{1,m} \\ \cdots & \cdots & \cdots \\ norm\_score_{n_{nonseed},1} & \cdots & norm\_score_{n_{nonseed},m} \end{bmatrix}$$

$$M_{similarity} = M_{seed} \times M_{nonseed}^{T}$$
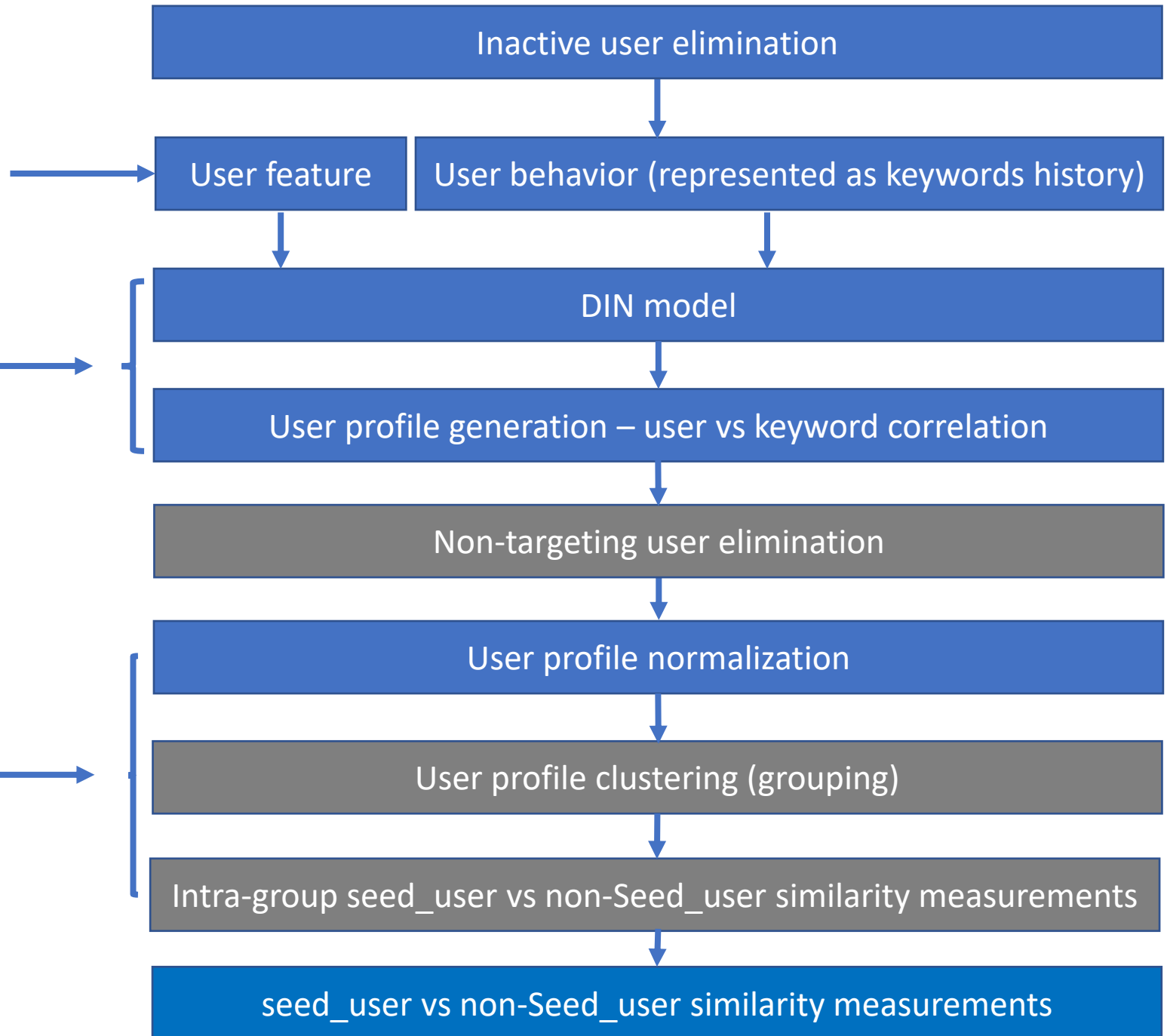
# Similarity computation estimation

# Workflow
## - Necessary + optional steps

Input

Kernel
computation

LookAlike model

**Inactive user elimination**

**User feature** | **User behavior (represented as keywords history)**

**DIN model**

**User profile generation – user vs keyword correlation**

**Non-targeting user elimination**

**User profile normalization**

**User profile clustering (grouping)**

**Intra-group seed_user vs non-Seed_user similarity measurements**

**seed_user vs non-Seed_user similarity measurements**

# DIN Model Output
## – non-targeting user elimination

| | Keyword$_1$ | Keyword$_2$ | Keyword$_3$ | Keyword$_4$ | ... | Keyword$_m$ |
|---|---|---|---|---|---|---|
| User$_1$ | score$_{11}$ | score$_{12}$ | score$_{13}$ | score$_{14}$ | ... | score$_{1m}$ |
| User$_2$ | score$_{21}$ | score$_{22}$ | score$_{23}$ | score$_{24}$ | ... | score$_{2m}$ |
| ... | ... | ... | ... | ... | ... | ... |
| User$_n$ | score$_{n1}$ | score$_{n2}$ | score$_{n3}$ | score$_{n4}$ | ... | score$_{nm}$ |

DIN

$$User_i's\ profile: S_i = \{score_{i1},\quad score_{i2},\quad ...\quad score_{im}\}$$

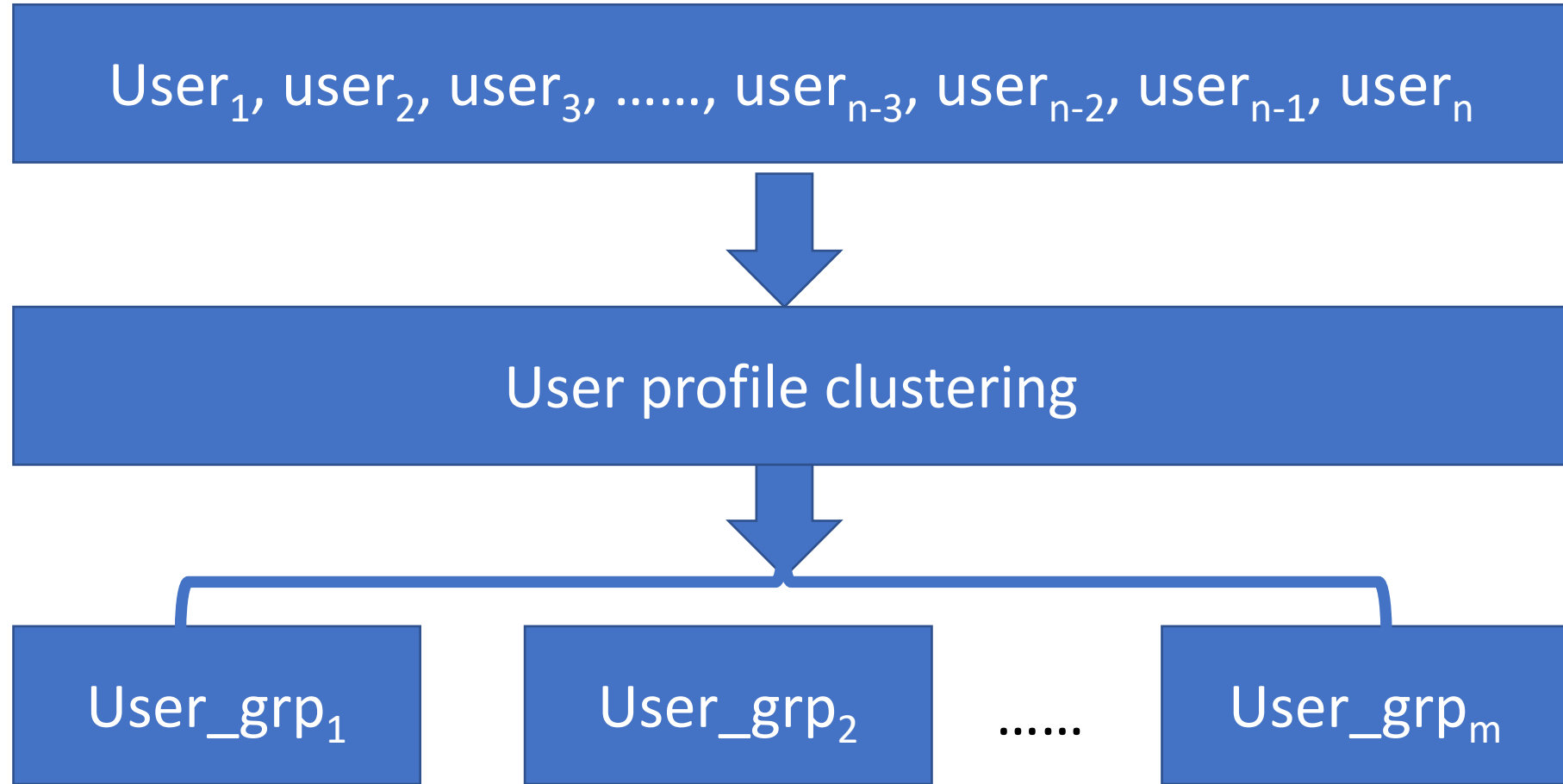$$Eliminate\ max(S_i) = \max_j score_{ij} < prefined\ threshold$$

* The purpose is to eliminate users that have no interest of any keywords (ineffective traffic)

$$targeting\ users$$

# DIN Model Output
## – user clustering

# DIN based Look-Alike model
## – group-wise seed_user vs non-seed_user similarity measure (active user only)

| | Seed user in $grp_1$ | Seed user in $grp_2$ | …… | Seed user in $grp_m$ |
|---|---|---|---|---|
| Nonseed user in $grp_1$ | similarity matrix$_{11}$ | 0 | …… | 0 |
| Nonseed user in $grp_2$ | 0 | similarity matrix$_{22}$ | …… | 0 |
| …… | …… | …… | …… | 0 |
| Nonseed user in $grp_m$ | 0 | 0 | …… | similarity matrix$_{mm}$ |

# DIN based Look-Alike model
## – within group seed_user vs non-seed_user similarity measure

*Similarity matrix$_{ii}$*

| | Seed_user$_{grpi,1}$ | Seed_user$_{grpi,2}$ | …… | Seed_user$_{grpi,m}$ |
|---|---|---|---|---|
| Nonseed_user$_{grpi,1}$ | Similary$_{11}$ | Similary$_{12}$ | …… | Similary$_{1m}$ |
| Nonseed_userg$_{rpi,2}$ | Similary$_{21}$ | Similary$_{22}$ | …… | Similary$_{2m}$ |
| Nonseed_user$_{grpi,3}$ | Similary$_{31}$ | Similary$_{32}$ | …… | Similary$_{3m}$ |
| Nonseed_user$_{grpi,4}$ | Similary$_{41}$ | Similary$_{42}$ | …… | Similary$_{4m}$ |
| …… | …… | …… | …… | …… |
| Nonseed_user$_{grpi,n}$ | Similary$_{n1}$ | Similary$_{n2}$ | …… | Similary$_{nm}$ |

Parallel computed and top 10 values for each row need to be stored

| All Seed Users in grpi |
|---|
| $\operatorname{mean}_{i}(\text{top10 } similarity_{1i})$ |
| $\operatorname{mean}_{i}(\text{top10 } similarity_{2i})$ |
| $\operatorname{mean}_{i}(\text{top10 } similarity_{3i})$ |
| $\operatorname{mean}_{i}(\text{top10 } similarity_{4i})$ |
| …… |
| $\operatorname{mean}_{i}(\text{top10 } similarity_{ni})$ |
| sort |

| Rank$_1$ nonseed_user |
|---|
| Rank$_2$ nonseed_user |
| Rank$_3$ nonseed_user |
| Rank$_4$ nonseed_user |
| … |
| Rank$_n$ nonseed_user |