

Privacy-preserving datasets of eye-tracking samples with applications in XR: Supplementary Material

Brendan David-John, *Member, IEEE*, Kevin Butler *Senior Member, IEEE* and Eakta Jain, *Member, IEEE*

1 THREAT SCENARIO k -ANONYMITY DETAILS

Age and gender demographics are generalized by grouping values into ranges to achieve k -anonymity. The number of data rows for each unique combination of age and gender ranges must be k or greater to maintain the privacy guarantee. The combined dataset of ET-DK2 and 360_em consists of 24 individuals with age and gender values listed in Table 1.

Table 1: Age and Gender demographics for ET-DK2 and 360_em datasets. Note that Subject ID 1 from both datasets were excluded from analysis due to data loss and subject sickness during data collection, respectively.

Dataset	Subject ID	Age	Gender
ET-DK2	2	M	43
ET-DK2	3	F	27
ET-DK2	4	M	29
ET-DK2	5	M	32
ET-DK2	6	F	28
ET-DK2	8	M	26
ET-DK2	9	F	23
ET-DK2	10	M	30
ET-DK2	11	F	28
ET-DK2	12	M	26
ET-DK2	13	M	52
ET-DK2	14	M	26
ET-DK2	15	M	35
ET-DK2	16	M	50
ET-DK2	17	M	33
ET-DK2	18	M	31
ET-DK2	19	M	32
ET-DK2	20	M	36
360_em	2	M	38
360_em	3	M	29
360_em	4	F	23
360_em	5	F	31
360_em	6	M	27
360_em	7	M	31
360_em	8	F	23
360_em	9	M	24
360_em	10	M	23
360_em	11	M	27
360_em	12	M	23
360_em	13	M	23
360_em	14	M	32

Ranges were selected for each value of k that maximized the total number of groups while ensuring each group had at least k rows matching the ranges of age and gender. The ranges of age and gender used to establish k -anonymity are listed in Table 2.

Table 2: Gender and age ranges used to generalize the ET-DK2 and 360_em demographics for k -anonymity. For each value of k the data rows are mapped into the listed ranges based on actual values. For example, (Male, 23-31) would be assigned to all Males between the age of 23 and 31. Male/Female refers to the data rows not specifying either value for Gender.

k	Gender & Age Generalization
4	(Female, 23-31), (Male, 23-27), (Male, 29-31), (Male, 32-33), (Male, 35-52)
6	(Female, 23-31), (Male, 23-27), (Male, 29-33), (Male, 35-52)
8	(Male/Female, 23-27), (Male/Female, 28-31), (Male/Female, 32-52)
15	(Male/Female, 23-28), (Male/Female, 29-52)

2 PRIVACY MECHANISM PSEUDOCODE

2.1 k -same-synth

```

1: procedure  $k$ -SAME-SYNTH( $k$ ,  $sample\_data$ ,  $fix\_event\_params$ ,  $sacc\_event\_params$ )
2: Parameters:  $k$  -  $k$ -anonymity parameter
3:    $sample\_data$  - Time series of gaze samples, indexed by stimulus  $m$ , identity  $i$ , and fixation/saccade events  $e$ 
4:    $fix\_event\_params$  - Fixation Gaussian parameters, indexed by stimulus  $m$ , identity  $i$ , and event  $e$ 
5:    $sacc\_event\_params$  - Velocity profile parameters, indexed by stimulus  $m$ , identity  $i$ , and event  $e$ 
6:    $fix\_event\_params \leftarrow k$ -same-select sequence( $k, fix\_event\_params$ ) ▷ Make fixation params  $k$ -anonymous
7:    $sacc\_event\_params \leftarrow k$ -same-select-sequence( $k, sacc\_event\_params$ ) ▷ Make saccade params  $k$ -anonymous
8:   for  $m = 1$  to  $num\_stimuli$  do ▷ Process events from each stimulus independently
9:     for  $i = 1$  to  $num\_identities$  do ▷ Process samples for each identity
10:       $fix\_data\_params \leftarrow fix\_event\_params[m, i, :]$  ▷ List of fixation event parameters
11:      for  $e = 1$  to  $num\_fixations$  do
12:         $\mu_x, \mu_y, \sigma_x, \sigma_y, t \leftarrow fix\_data\_params[e]$ 
13:         $sample\_data[m, i, e] \leftarrow SynthFixation(\mu_x, \mu_y, \sigma_x, \sigma_y, t)$  ▷ Synthesize samples for fixation  $e$  by sampling 2D Normal distribution
14:         $sacc\_data\_params \leftarrow sacc\_event\_params[m, i, :]$  ▷ List of saccade event parameters
15:        for  $e = 1$  to  $num\_saccades$  do
16:           $a, b, c, t \leftarrow sacc\_data\_params[e]$ 
17:           $sample\_data[m, i, e] \leftarrow SynthSaccade(a, b, c, t)$  ▷ Synthesize samples for saccade  $e$  using velocity profile from Gaussian model
return  $sample\_data$ 

```

2.2 event-synth-PD

```

1: procedure EVENT-SYNTH-PD( $k$ ,  $\gamma$ ,  $sample\_data$ ,  $fix\_event\_params$ ,  $sacc\_vel\_profiles$ ,  $CVAE_{enc}$ ,  $CVAE_{dec}$ )
2: Parameters:  $k, \gamma$  - plausible deniability parameters
3:    $sample\_data$  - Time series of gaze sample, indexed by stimulus  $m$ , identity  $i$ , and fixation/saccade events  $e$ 
4:    $fix\_event\_params$  - Fixation Gaussian parameters, indexed by stimulus  $m$ , identity  $i$ , and event  $e$ 
5:    $sacc\_vel\_profiles$  - Saccade velocities and conditions, indexed by stimulus  $m$ , identity  $i$ , and event  $e$ 
6:    $CVAE_{enc}$  - Encoder network of C-VAE, maps input to latent space distributions defined by  $\mu$  and  $\sigma$ 
7:    $CVAE_{dec}$  - Decoder network of C-VAE, maps input random samples  $z \oplus c$  to synthetic velocities
8:   for  $m = 1$  to  $num\_stimuli$  do ▷ Process events from each stimulus independently
9:     for  $i = 1$  to  $num\_identities$  do ▷ Process samples for each identity
10:       $fix\_data\_params \leftarrow fix\_event\_params[m, i, :]$  ▷ List of fixation event parameters
11:      for  $e = 1$  to  $num\_fixations$  do ▷ Synthesize fixation samples until PD criterion is met
12:         $d = (\mu_x, \mu_y, \sigma_x, \sigma_y, t) \leftarrow fix\_data\_params[e]$  ▷ Params for fixation  $e$ 
13:         $\mathbf{M}_{fix} \leftarrow N(x, y)$  ▷ 2D Normal distribution that returns  $t$  values
14:         $result \leftarrow False$ 
15:        while  $result == False$  do
16:           $y \leftarrow \mathbf{M}_{fix}(d)$  ▷ Generate  $t$  samples from distribution with curr params
17:           $Pr_d \leftarrow Pr\{y \leftarrow \mathbf{M}_{fix}(d)\}$  ▷ Probability real seed generated synthetic samples  $y$ 
18:           $result \leftarrow PD\ Event\ Privacy\ Test(k, \gamma, Pr_d, \mathbf{M}_{fix}, fix\_event\_params[m, \neq i, :])$  ▷  $\neq i$  indicates all individual data besides  $i$ 
19:           $sample\_data[m, i, e] \leftarrow y$ 
20:       $sacc\_data \leftarrow sacc\_vel\_profiles[m, i, :]$  ▷ List of real data saccade profiles
21:      for  $e = 1$  to  $num\_saccades$  do ▷ Synthesize fixation samples until PD criterion is met
22:         $d = (\mu_1, \sigma_1, \dots, \mu_L, \sigma_L) \leftarrow C - VAE_{enc}(sacc\_data[e])$ 
23:         $\mathbf{M}_{sacc} \leftarrow N_1, \dots, N_L$  ▷ Define  $\mathbf{M}$  as  $L$  independent Normal distributions
24:         $result \leftarrow False$ 
25:        while  $result == False$  do
26:           $y = (z_1, \dots, z_L) \leftarrow \mathbf{M}_{sacc}(d)$ 
27:           $Pr_d \leftarrow Pr\{y \leftarrow \mathbf{M}_{sacc}(d)\}$  ▷ Probability real seed generated synthetic samples  $y$ 
28:           $result \leftarrow PD\ Event\ Privacy\ Test(k, \gamma, Pr_d, \mathbf{M}_{sacc}, sacc\_vel\_profiles[m, \neq i, e])$  ▷  $\neq i$  indicates all individual data besides  $i$ 
29:           $sample\_data[m, i, e] \leftarrow y$ 
return  $sample\_data$ 

```

```

1: procedure PD EVENT PRIVACY TEST( $k, \gamma, Pr_d, \mathbf{M}, D$ )
2: Parameters:  $k, \gamma$  - plausible deniability parameters,  $Pr_d$  - Probability of real seed for  $y$ ,  $Pr\{y \leftarrow \mathbf{M}(d)\}$ 
3:    $\mathbf{M}$  - generative model that synthesized  $y$ ,  $D$  - data records from identities other than input
4:    $i' \leftarrow$  unique integer  $i'$ ,  $s.t. \gamma^{-i'-1} < Pr_d \leq \gamma^{-i'}$ 
5:    $k' \leftarrow 0$ 
6:   for  $i = 1$  to  $num\_identities$  do
7:      $D_i \leftarrow D[i]$ 
8:     for  $d_a \in D_i$  do
9:       if  $\gamma^{-i'-1} < Pr\{y = \mathbf{M}(d_a)\} \leq \gamma^{-i'}$  then
10:         $k' \leftarrow k' + 1$ 
11:        Break ▷ Move for loop for  $i$  onto the next identity
12:   if  $k' \geq k - 1$  then return Pass
13:   else return Fail

```

2.3 Kaleido

The pseudocode below details the kaleido approach for a stream of n_{raw} gaze samples g_1, \dots, n_{raw} , window size w , privacy parameter ϵ , sample distance threshold l_{thresh} , sample skipping parameter t_{skip} , spatial parameter r , and ratio of testing to publishing privacy budget h .

The adaptive algorithm includes several parameters that allow for privacy budget savings while processing the gaze sample at each timestamp. First, a fixed time duration $t_{skip} = 50ms$ is used to skip gaze samples that arrive within t_{skip} of the last published gaze position. Next, after t_{skip} has passed since the last published gaze point, the algorithm moves on to the testing phase. If the current gaze position is within the fixation threshold determined by l_{thresh} and ϵ^{test} , then the previously published position is re-used, and only ϵ^{test} of the budget for the current time window is consumed. The algorithm enters the publishing phase if the new gaze position is farther than the threshold. A noisy gaze position is generated using the ϵ^{pub} budget with a Planar Laplacian mechanism [1]. The amount of the ϵ^{pub} budget used decreases adaptively to preserve as much utility as possible while maintaining ϵ -DP guarantee within each time window. This process is repeated for each time window, and any leftover ϵ^{pub} budget is recycled into the next window. A complete description of the proof that each window consumes at most ϵ of the privacy budget is available in the original paper [2].

```

1: procedure KALEIDO DP( $g_1, \dots, n_{raw}, w, \epsilon, l_{thresh}, t_{skip}, r, h$ )
2: Parameters:  $g_1, \dots, n_{raw}$  - Stream of gaze positions,  $w$  - Window size (# samples),  $\epsilon$  - DP privacy level
3:            $l_{thresh}$  - Distance threshold for testing,  $t_{skip}$  - # of samples to skip over during testing
4:            $r$  - Privacy radius for DP,  $h$  - Ratio of privacy budget used for testing
5:            $n_{test} \leftarrow \lceil w/t_{skip} \rceil$  ▷ Number of points to test for each window
6:            $\epsilon_{test} \leftarrow \epsilon/(h \cdot n_{test})$  ▷ Privacy budget allocated to test each sample
7:            $i_{test} \leftarrow null$  ▷ Index of the last tested gaze position.
8:            $i_{pub} \leftarrow null$  ▷ Index of the last published gaze position.
9:            $g'_i \leftarrow zeros(n_{raw})$  ▷ Published gaze position for sample  $i$ , initialized to zeros.
10:           $\epsilon_i^{pub} \leftarrow zeros(n_{raw})$  ▷ List of privacy budget consumed for sample  $i$ , initialized to zeros.
11:          for  $i = 1$  to  $num_{raw}$  do ▷ Process each window of raw gaze samples
12:            if  $i_{test} \neq null$  AND  $t(i) - t(i_{test}) < t_{skip}$  then ▷ Check if sample should be skipped based on  $t_{skip}$  parameter
13:               $g'_i \leftarrow g'_{i_{pub}}$ 
14:               $\epsilon_i^{pub} \leftarrow 0$ 
15:              Continue
16:             $i_{test} = i$ 
17:             $l_{dis} = d(g_i, g'_{i_{pub}})$  ▷ Distance between gaze sample  $i$  and last published
18:             $\eta \sim Lap(1/\epsilon_{test})$  ▷ Sample from Laplace distribution, small values of  $\epsilon_{test}$  introduce more noise
19:            if  $l_{dis} \neq null$  AND  $l_{dis} \leq l_{thresh} + \eta$  then ▷ Test if current gaze is close enough to last published to repeat
20:               $g'_i \leftarrow g'_{i_{pub}}$ 
21:               $\epsilon_i^{pub} \leftarrow 0$ 
22:              Continue
23:             $i_{pub} \leftarrow i$  ▷ Publish a new gaze sample, update index of last published
24:             $\epsilon_{rem} \leftarrow \epsilon - \epsilon/h - \sum_{k=i-n_{raw}+1}^{i-1} \epsilon_k^{pub}$  ▷ Compute remaining privacy budget for this window
25:             $\epsilon_i^{pub} \leftarrow \epsilon_{rem}/2$ 
26:             $g'_i \leftarrow PlanarLap(g_i, \epsilon_i^{pub}/r)$ 
return  $g'$ 

```

3 C-VAE MODEL TRAINING PROCEDURE

The C-VAE model for generating synthetic saccade profiles was trained using tensorflow version 1.13.1. Models were trained independently for each dataset using data from all individuals and stimuli. Training was performed using 75% of the available data with the remaining 25% used as a validation set.

All models were trained with an ADAM optimizer using tensorflow's Model compile and fit functions. The loss function was defined as

$$L(x, \mathbf{D}(z)) = \|x - \mathbf{D}(z)\|_2 - \mathbf{KL}(\mathbf{N}(\mu, \sigma), \mathbf{N}(0, 1)),$$

where the first term is Mean Squared Error for the reconstructed synthetic profile and the second terms employs KL Divergence to enforce latent space sampling that follows a normal distribution with zero mean.

4 C-VAE MODEL HYPER-PARAMETER OPTIMIZATION

Hyper-parameters were tuned using the EHTask dataset as it contained a longer duration of data compared to the DGaze dataset. Grid search optimization was performed over the following sets of values, with optimal parameters in bold:

- Learning Rate: 0.001, **0.01**
- Batch Size: 20, 60, **100**
- Number of Epochs: 10, **20**, 30
- Encoder Hidden Layer with ReLU activation function: **32**, 64, 96 Nodes
- Latent Space Dimension: 32, **64**, 96
- Decoder Hidden Layer with linear activation function: 32, 64, **96** Nodes

The optimal parameters produced an average loss of 0.33 on the validation set.

ACKNOWLEDGMENTS

Authors acknowledge funding from the National Science Foundation (Awards FWHTF-2026540, CNS-1815883, and CNS-1562485), the National Science Foundation GRFP (Awards DGE-1315138 and DGE-1842473), and the Air Force Office of Scientific Research (Award FA9550-19-1-0169).

REFERENCES

- [1] M. E. Andrés, N. E. Bordenabe, K. Chatzikokolakis, and C. Palamidessi. Geo-indistinguishability: Differential privacy for location-based systems. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, pp. 901–914, 2013.
- [2] J. Li, A. R. Chowdhury, K. Fawaz, and Y. Kim. KalEido: Real-time privacy control for eye-tracking systems. In *29th USENIX Security Symposium (USENIX Security 20)*, 2020.