

## SUPPLEMENTARY MATERIAL

## A. Error-Triggering Input Pairs

Table VI  
ERROR-TRIGGERING INPUT PAIRS. R/W1 IS THE ABBREVIATION OF READ/WRITE VIRTUAL MEMORY #1.

Candidate Solution	Input Pair
Buggy Square Root	R2, R1, R0, R2, R0, R2; W1, W1, W2, W2, W3, W0.
Optimized Square Root	W3, W3, W1, W1, W3, W3; R0, R1, R0, R2, R0, R1.
Hierarchical Protocol	R2, R1, R4, W2, R3, R1; W2, W0, W7, W7, W7, W6.
Cuckoo Hashing Protocol	R4, R0, R5, W0, R5, R6; W0, W7, W3, W1, W0, W0.
Buggy Path ORAM	W1, R0, W3, R1, R1, W2; R3, W1, R2, W0, W3, R0.
Lethe	W1, R0, W3, R1, R1, W2; R3, W1, R2, W0, W3, R0.
PathOHeap	W1, R0, W3, R1, R1, W2; R3, W1, R2, W0, W3, R0.
OblIDB	see Table VII

Table VII  
THE ERROR-TRIGGERING DATABASE TABLES OF OBLIDB.

Table 1	Table 2
1, B	1, A
2, B	2, B
3, C	3, C
4, D	4, D
1, C	1, C
2, B	2, B
3, C	3, C
4, D	4, D
1, C	1, C
2, B	2, B
3, C	3, C
4, D	4, D
1, D	1, D
2, B	2, B
3, C	3, C
4, D	4, D
1, C	1, C
2, B	2, B
3, C	3, C
4, D	4, D
1, B	1, B
2, B	2, B
3, C	3, C
4, D	4, D

We report the error triggering input pairs in Table VI. W1 and R1 correspond to Push 1 and Pop 1 in when testing PathOHeap. Note that the input pair of OblIDB corresponds to two tables. Each table contains two columns (one numerical column and one categorical column). We report two tables in two columns of Table VII, respectively.

## B. Minimal Validation Set Size

When discussing “Capturing Stealthy Violations” in Sec. IV-B, we present a corollary such that the minimal validation set size  $|O_v|$  should be over a threshold  $\frac{5.992}{r^2}$ , where  $r$  is the proportion of distinguishable sequences in  $O_v$ . To understand this corollary, we first assume there is an ideal distinguisher  $\mathcal{D}$  that can correctly distinguish all distinguishable sequences and the proportion of distinguishable sequences among all sequences is  $r$ . Then, since those  $r$  distinguishable sequences are correctly distinguished by  $\mathcal{D}$  (given  $r$  correct cases) and the remaining  $1 - r$  indistinguishable sequences are randomly distinguished (given  $0.5 \times (1 - r)$  correct cases), the total accuracy is  $0.5 + 0.5r$ . To ensure that  $0.5 + 0.5r > 0.5 + \sqrt{1.498/n}$ , we have  $n > 5.992/r^2$ .

## C. Minimized Error-Triggering Subsequences

We report the minimized error-triggering subsequences (marked in red) as follow.



#### D. Comparison of Different Neural Distinguishers

Table VIII compares three representative neural models that are frequently used to process sequence data. For this evaluation, we generate a dataset with 10,000/6,000 training/testing samples (with avg. token size of 13,593) using the Cuckoo Hashing Protocol with  $N = 48$ . Table VIII reports the processing time, accuracy, and peak GPU memory usage. We also report the #parameters for each model. More parameters imply larger models. All models with their full architecture details are at <https://anonymous.4open.science/r/b6f98247-c616-40a7-97a0-2b442747077f/>.

Table VIII  
COMPARING DIFFERENT MODELS ON PROCESSING LENGTHY SEQUENCES. SEE FULL DETAILS OF THESE MODELS IN OUR ARTIFACT.

Model	Training Time (Per Epoch)	Accuracy (Found Bug?)	Peak GPU Memory Usage	#Parameters
CNN (model leveraged by NeuralD)	52	0.618 (✓)	2,004 MB	54,306
RNN	639	0.502 (✓)	1,870 MB	37,411
Bi-LSTM [71]	1249	0.500 (X)	2,676 MB	182,947

We show that within three epochs, CNN (the neural distinguisher leveraged by NeuralD) and RNN both provide a non-trivial distinguisher (the cost of each epoch is shown in Table VIII). In contrast, although consuming more memory, Bi-LSTM, the most advanced model tested at this step, fails to yield a non-trivial distinguisher. The remaining two models are far more concise and can converge much faster. When compared with RNN that is marginally over the baseline, the CNN model has the highest accuracy. Furthermore, we observe that CNN is much faster in terms of training time per epoch, as it is highly optimized on modern hardware. Overall, the results suggest that different neural models perform distinctly; we leave designing more efficient and sophisticated models for future work.

Table IX  
INPUT-INDEPENDENT BUILDING BLOCKS

Building Block	Token	Description
Sort	Sort	obviously sort all memory blocks by “tag”
Scan	SeqScan	sequentially scan a range of memory blocks
Delete	Del	delete a range of memory blocks
Rehash	Rehash	re-allocate a level of memory blocks

#### E. Input-independent building blocks

Table IX lists all input-independent building blocks adopted by our “Single-Sequence Distillation” optimization in Sec. IV-C. Memory access patterns of these building blocks are well-known to be *oblivious*; replacing them with tokens does not result in false negative of testing. In contrast, optimization at this step effectively simplifies lengthy memory access traces, facilitating our lightweight neural distinguisher to comprehend each trace in a speedy manner.

#### F. Proof of Corollary 4.1

The frequency of identified subsequences are guaranteed with a high probability by Hoeffding–Serfling’s inequality. The co-existence of certain LCS within a subset can be regarded as sampling from finite Bernoulli( $\tau$ ) without replacement, where  $\tau$  denotes the lower bound of frequency of the LCS among  $\mathcal{S}$ . “1” denotes that the LCS contains in the sample. In our setting, the sum of  $m$  times sampling (w/o replacement) is  $m$ . As the population of Bernoulli( $\tau$ ) is  $|\mathcal{S}|$  (i.e., cardinality), we apply the equality and have

$$Pr(\tau \leq \frac{H(m)}{m} - \epsilon) \leq \exp\left(-2\epsilon^2 \frac{m}{1 - (m-1)/|\mathcal{S}|}\right) \quad (9)$$

Therefore, we have

$$Pr(\tau > \frac{H(m)}{m} - \epsilon) \leq 1 - \exp\left(-2\epsilon^2 \frac{m}{1 - (m-1)/|\mathcal{S}|}\right) \quad (10)$$

To ensure  $Pr(\tau > \frac{H(m)}{m} - \epsilon)$  higher than 0.95, we need to ensure the RHS of Eqn. 10 equals 0.95. Therefore,  $\epsilon \approx \sqrt{\frac{1.498(|\mathcal{S}| - m + 1)}{m|\mathcal{S}|}}$  and  $\tau > 1 - \sqrt{\frac{1.498(|\mathcal{S}| - m + 1)}{m|\mathcal{S}|}}$  with the probability of 0.95. Hence,  $m = \frac{1.498(|\mathcal{S}| + 1)}{\tau^2|\mathcal{S}| - 2\tau|\mathcal{S}| + |\mathcal{S}| + 1.498}$ .