# 1 uTraVaS: The Utility-Aware TraVaS

The following supplementary material contains the theoretic foundations and experimental results of *uTraVaS*, an optimization framework for *TraVaS* based on a data utility cost function. We refer to mathematical definitions and the standard *TraVaS* algorithm of the main paper *TraVaS: Differentially Private Trace Variant Selection for Process Mining* (M. Rafiei et al.).

## 1.1 Algorithm Design

As an extension of *TraVaS*, the optimizer framework *uTraVaS* serves as a wrapper around the standard *TraVaS* algorithm and exploits the *composition theorem* and *post-processing immunity* properties of $(\epsilon, \delta)$-DP to spend the privacy budgets, $\epsilon$ and $\delta$, w.r.t. an optimal data utility preservation. Subsequently, we first introduce these properties and then explain the algorithm.

**Theorem 1 ($(\epsilon, \delta)$-DP Mechanism Composition [1])** *Let $L$ be a simple event log, and $\mathcal{M}_i, 1 \leq i \leq n$ be $(\epsilon_i, \delta_i)$-DP mechanisms. The sequential application of these mechanisms on arbitrary sublogs of $L$ leads to an overall worst-case privacy level parameterized by $(\sum_{1 \leq i \leq n} \epsilon_i, \sum_{1 \leq i \leq n} \delta_i)$. If each $\mathcal{M}_i$ operates on strictly disjoint sublogs of $L$, the worst-case privacy level is $(max_{1 \leq i \leq n} \epsilon_i, max_{1 \leq i \leq n} \delta_i)$, so-called parallel composition.*

Post-processing immunity states that $(\epsilon, \delta)$-DP is immune to post-processing activities on mechanism outputs [1]. Hence, parties without additional knowledge about an event log $L$ cannot transform the output in a way that makes the information less differentially private.

According to Theorem 1, a privacy budget of $(\epsilon, \delta)$ may be consumed by either one *travas* query or $n$ subqueries with parameters $(\epsilon/n, \delta/n)$. Algorithm 1 receives the inputs of the *TraVaS* algorithm and the maximum number of subqueries $n$ and returns a merged anonymized event log that preserves more data utility while spending the same privacy budgets. Since the relation between the number of subqueries and the data utility of the resulting anonymized event log is non-linear, *uTraVaS* searches the whole space of possible number of subqueries from 1 to $n$ to find the decomposition structure that makes the best use of the privacy budgets for *TraVaS*.

In line 6, $i$ different anonymized sublogs are compared, and all the unique variants from different sublogs are added to a potential final anonymized log $pL'$. The rationale behind this is that the partition selection mechanism never generates fake variants, and it may remove some infrequent trace variants. Thus, all the variants that appear in different sublogs are among the original trace variants. Note that the privacy degradation due to the knowledge that one can obtain by comparing $i$ different sublogs are already considered by reducing the privacy parameters from $\epsilon$ and $\delta$ to $\epsilon/i$ and $\delta/i$ (see Theorem 1). The frequency of the variants are calculated as we explain in the following. Note that calculating frequency values for the trace variants in $pL'$ is a post-processing step where the noisified frequencies are altered without using original frequency values. Thus, because of the post-processing immunity the provided privacy guarantees are not degraded.

---

**Algorithm 1:** Utility-Aware TraVaS (uTraVaS)

---

    **Input:** Event log $L$, DP-Parameters $(\epsilon, \delta)$, Subquery bound $n$

    **Output:** $(\epsilon, \delta)$-DP log $L'$

**1 function** optimize $(L, \epsilon, \delta, n)$

**2**     **forall** $i \in \{1, 2, \ldots, n\}$ **do**

**3**         create $i$ empty *sublogs*: $\{sL_1, sL_2, \ldots, sL_i\}$         // initialize

**4**         **forall** $j \in \{1, \ldots, i-1, i\}$ **do**

**5**             add travas$(L, \epsilon/i, \delta/i)$ to $sL_j$         // run TraVaS

**6**         add all unique variants $\sigma$ from *sublogs* to a simple log $pL'$

**7**         **forall** $\sigma \in pL'$ **do**

**8**             **forall** $sL \in sublogs$ **do**

**9**                 **if** $\sigma \in sL$ **then**

**10**                     $f_\sigma = sL(\sigma)$         // get frequency from sublog

**11**                 **else**

**12**                     $f_\sigma = $ freqEstimate $(\epsilon, \delta)$         // estimate frequency

**13**                 update $pL'(\sigma)$ to $pL'(\sigma) + f_\sigma$ in $pL'$

**14**         update $pL'(\sigma)$ to $pL'(\sigma)\ /\ i$ in $pL'$         // compute mean

**15**         add $pL'$ to the set *pool*.         // optimization domain

**16**     $L' = \underset{pL' \in pool}{\arg\max}\ \mathrm{utility}(L, pL')$         // optimize data utility

**17**     **return** $L'$

---

For variants contained in all $i$ sublogs, we calculate the mean of the associated noisified frequencies as the best unbiased true frequency estimator due to the symmetric zero-mean of *k-TSGD*. However, there may be some variants that get truncated in some sublogs and do not appear within all sublogs. For such variants, the frequencies must first be estimated before computing the mean. Since we know that the true frequencies lie in the range $[1, 2k]$ due to the k-TSGD interval of $[-k, k]$ and the threshold at $k$, the maximal removed true frequency is $2k$. Such estimations represent conditional expectation values of the form $E[L(\sigma)\ |\ \sigma \in L \wedge L(\sigma) > 0 \wedge L(\sigma) + x_\sigma \leq k]$. Consequently, all possible frequency realizations have to be combined with their likelihoods of producing perturbed outputs under *k-TSGD*. Using the set $T(j) = \{l\ |\ l \in Z \wedge j + l \leq k \wedge |l| \leq k\}$, we formally define the following estimator.

$$\mathrm{freqEstimate}(\epsilon, \delta) = \frac{\sum\limits_{j=1}^{2k} \sum\limits_{l \in T(j)} j \cdot \mathrm{k\text{-}TSGD}[l\ |\ p, k]}{\sum\limits_{j=1}^{2k} \sum\limits_{l \in T(j)} \mathrm{k\text{-}TSGD}[l\ |\ p, k]} \tag{1}$$

After all $n$ anonymized candidate logs are assembled, their respective data utility is measured by means of a suitable utility evaluation function and w.r.t. the original event log. We employ the *earth mover's distance* as a generic data utility function to quantify the similarity between two distributions of trace variants [2].[1] In a final step, the optimizer

---

[1]Note that the utility function can be adjusted w.r.t. any specific utility demand.

then releases the best performing candidate to the public.

We particularly note that in line 16, *uTraVaS* uses the utility score to select the best performing anonymized event log based on the number of *travas* subqueries. This is a selection step that utilizes the private original event log to select one of the possible query decomposition options but does not interfere with the random mechanism itself. The private information is never used during the noise injection process or in the post-processing step. Thus, as long as no internal state information such as the actual number of queries or the utility scores is published the final output of the algorithm is $(\epsilon, \delta)$ differentially private.

The log construction and noise generation routines in Algorithm 1 are parallelizable and only require one data query. As a result, our approach faces a considerably lower complexity compared to traditional DP query methods where iterations need to run sequentially. A more precise computational complexity analysis has been provided as supplementary material in our GitHub repository.[2]

## 1.2 Experiments

For a direct performance comparison between *TraVaS* and *uTraVaS*, we supplement our experimental results of the main paper with a data utility and result utility evaluation on *uTraVaS*. The event data used (Sepsis, BPIC2013) as well as the setup parameters remain unchanged. Figure 1.1 and Figure 1.2 illustrate side-by-side heatmaps of *relative log similarity* and *absolute log difference* on *uTraVaS*, *TraVaS* and our benchmark for BPIC2013 and Sepsis respectively. Both logs demonstrate superior performance of *uTraVaS*, in particular for $\delta > 0.001$. This observation can be explained by the merging process and frequency estimation of Algorithm 1 that supports infrequent variants of largely unstructured event logs.

With Figure 1.3 and Figure 1.4, we present the corresponding *fitness* and *precision* results for BPIC2013 and Sepsis. Compared to Sepsis event data that again show similarly increased scores, BPIC2013 hardly demonstrates any noticeable effect. Since this event log is considerably larger and possesses a lower trace uniqueness, the noise filter of the discovery algorithm suppresses the relatively smaller frequency deviations caused by *uTraVaS*.

We conclude that *uTraVas* provides significant performance increase for most privacy settings, particularly on unstructured, small or highly trace-unique event data.

---

[2] https://github.com/wangelik/TraVaS

Figure 1.1: The *relative log similarity* and *absolute log difference* results of anonymized BPIC2013 event logs generated by *uTraVaS*, *TraVaS*, and our prefix-based benchmark method. Each value represents the mean of 10 algorithm runs.
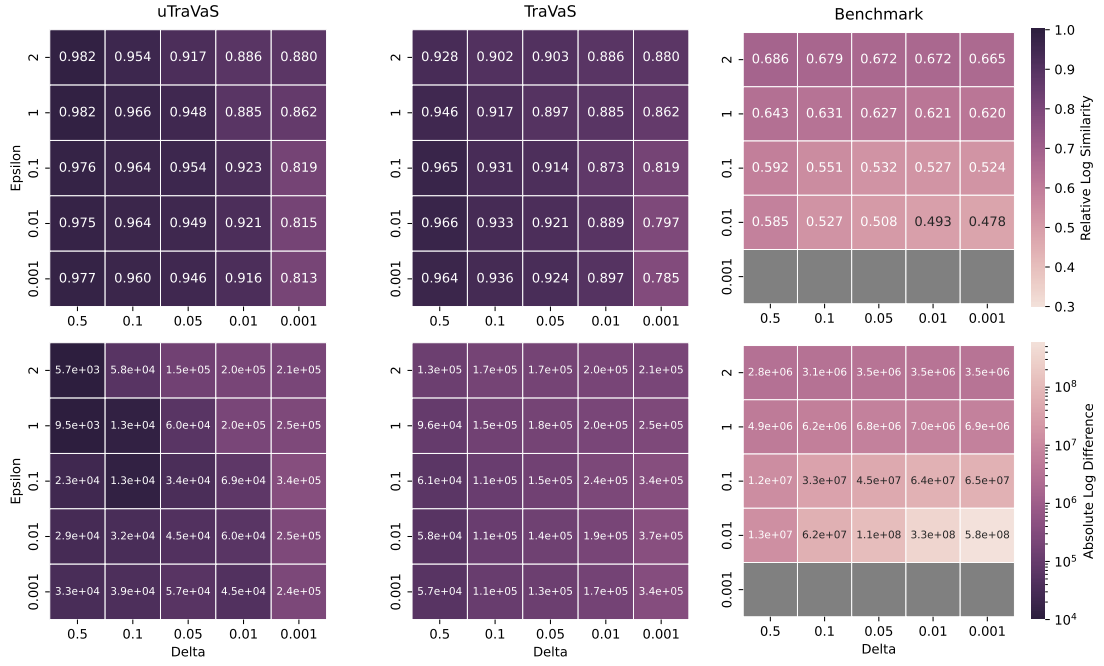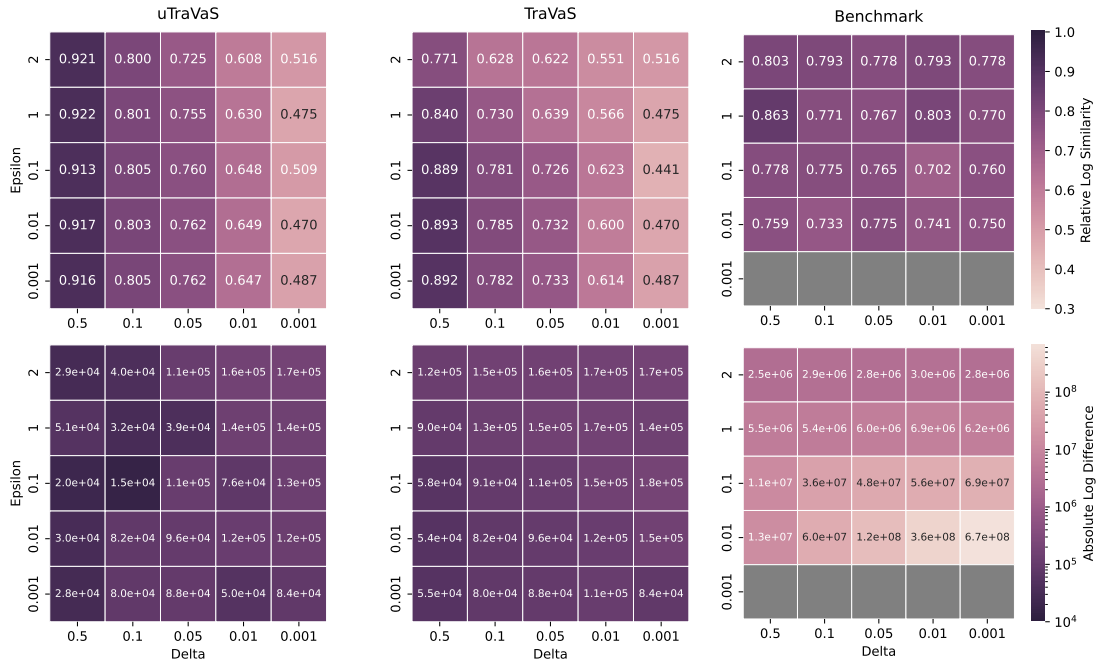


Figure 1.2: The *relative log similarity* and *absolute log difference* results of anonymized Sepsis event logs generated by *uTraVaS*, *TraVaS*, and our prefix-based benchmark method. Each value represents the mean of 10 algorithm runs.

**uTraVaS — Fitness**

| Epsilon \ Delta | 0.5 | 0.1 | 0.05 | 0.01 | 0.001 |
|---|---|---|---|---|---|
| 2 | 0.995 | 0.995 | 0.995 | 0.997 | 0.995 |
| 1 | 0.995 | 0.995 | 0.995 | 0.996 | 0.995 |
| 0.1 | 0.995 | 0.995 | 0.995 | 0.997 | 0.997 |
| 0.01 | 0.995 | 0.995 | 0.995 | 0.997 | 0.994 |
| 0.001 | 0.995 | 0.995 | 0.995 | 0.996 | 0.988 |

**TraVaS — Fitness**

| Epsilon \ Delta | 0.5 | 0.1 | 0.05 | 0.01 | 0.001 |
|---|---|---|---|---|---|
| 2 | 0.995 | 0.995 | 0.995 | 0.997 | 0.995 |
| 1 | 0.995 | 0.995 | 0.995 | 0.996 | 0.995 |
| 0.1 | 0.995 | 0.995 | 0.995 | 0.996 | 0.997 |
| 0.01 | 0.995 | 0.995 | 0.995 | 0.997 | 0.994 |
| 0.001 | 0.995 | 0.995 | 0.995 | 0.996 | 0.988 |

**Benchmark — Fitness**

| Epsilon \ Delta | 0.5 | 0.1 | 0.05 | 0.01 | 0.001 |
|---|---|---|---|---|---|
| 2 | 0.967 | 0.955 | 0.962 | 0.980 | 0.970 |
| 1 | 0.963 | 0.970 | 0.968 | 0.955 | 0.965 |
| 0.1 | 0.951 | 0.954 | 0.965 | 0.935 | 0.966 |
| 0.01 | 0.982 | 0.930 | 0.974 | 0.965 | 0.963 |
| 0.001 | | | | | |

**uTraVaS — Precision**

| Epsilon \ Delta | 0.5 | 0.1 | 0.05 | 0.01 | 0.001 |
|---|---|---|---|---|---|
| 2 | 0.877 | 0.877 | 0.877 | 0.877 | 0.877 |
| 1 | 0.877 | 0.877 | 0.877 | 0.876 | 0.877 |
| 0.1 | 0.877 | 0.877 | 0.877 | 0.901 | 0.931 |
| 0.01 | 0.877 | 0.877 | 0.877 | 0.894 | 0.931 |
| 0.001 | 0.877 | 0.877 | 0.877 | 0.898 | 0.940 |

**TraVaS — Precision**

| Epsilon \ Delta | 0.5 | 0.1 | 0.05 | 0.01 | 0.001 |
|---|---|---|---|---|---|
| 2 | 0.877 | 0.877 | 0.877 | 0.869 | 0.876 |
| 1 | 0.877 | 0.877 | 0.877 | 0.873 | 0.877 |
| 0.1 | 0.877 | 0.877 | 0.877 | 0.878 | 0.908 |
| 0.01 | 0.877 | 0.877 | 0.877 | 0.894 | 0.931 |
| 0.001 | 0.877 | 0.877 | 0.877 | 0.898 | 0.940 |

**Benchmark — Precision**

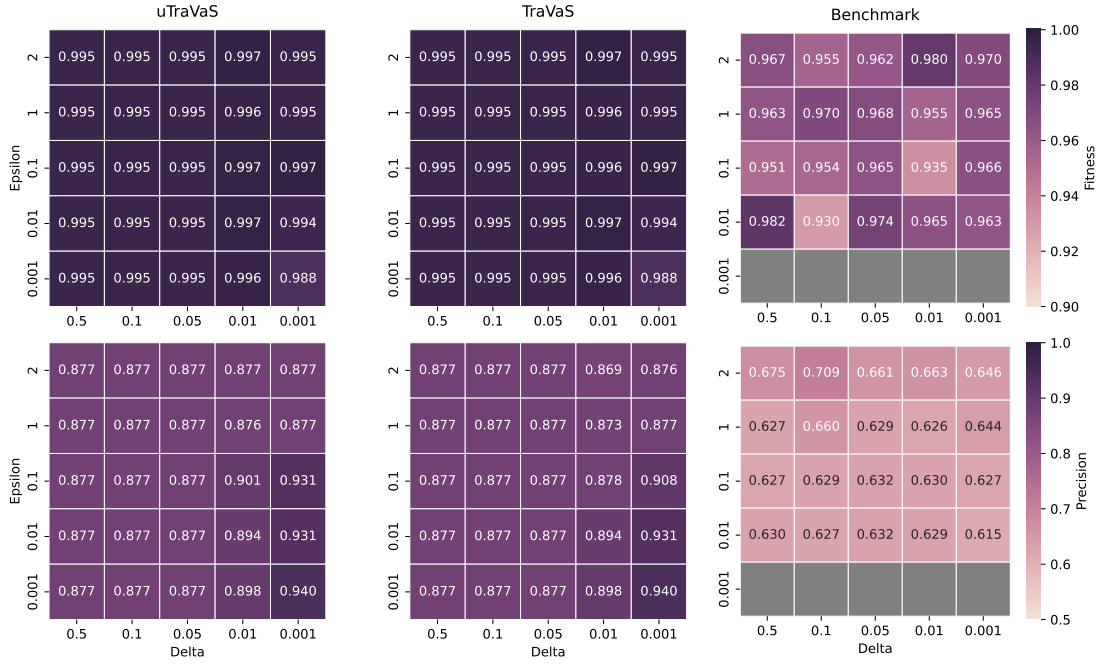| Epsilon \ Delta | 0.5 | 0.1 | 0.05 | 0.01 | 0.001 |
|---|---|---|---|---|---|
| 2 | 0.675 | 0.709 | 0.661 | 0.663 | 0.646 |
| 1 | 0.627 | 0.660 | 0.629 | 0.626 | 0.644 |
| 0.1 | 0.627 | 0.629 | 0.632 | 0.630 | 0.627 |
| 0.01 | 0.630 | 0.627 | 0.632 | 0.629 | 0.615 |
| 0.001 | | | | | |

Figure 1.3: The *fitness* and *precision* results of anonymized BPIC2013 event logs generated by *uTraVaS*, *TraVaS*, and our prefix-based benchmark method. Each value represents the mean of 10 algorithm runs.



**uTraVaS — Fitness**

| Epsilon \ Delta | 0.5 | 0.1 | 0.05 | 0.01 | 0.001 |
|---|---|---|---|---|---|
| 2 | 0.982 | 0.973 | 0.970 | 0.972 | 0.893 |
| 1 | 0.979 | 0.969 | 0.970 | 0.980 | 0.805 |
| 0.1 | 0.970 | 0.960 | 0.961 | 0.955 | 0.906 |
| 0.01 | 0.971 | 0.966 | 0.963 | 0.964 | 0.800 |
| 0.001 | 0.972 | 0.963 | 0.947 | 0.957 | 0.849 |

**TraVaS — Fitness**

| Epsilon \ Delta | 0.5 | 0.1 | 0.05 | 0.01 | 0.001 |
|---|---|---|---|---|---|
| 2 | 0.982 | 0.973 | 0.970 | 0.971 | 0.893 |
| 1 | 0.979 | 0.969 | 0.970 | 0.980 | 0.747 |
| 0.1 | 0.970 | 0.960 | 0.952 | 0.945 | 0.774 |
| 0.01 | 0.971 | 0.966 | 0.963 | 0.957 | 0.800 |
| 0.001 | 0.972 | 0.959 | 0.947 | 0.957 | 0.771 |

**Benchmark — Fitness**

| Epsilon \ Delta | 0.5 | 0.1 | 0.05 | 0.01 | 0.001 |
|---|---|---|---|---|---|
| 2 | 0.990 | 0.994 | 0.988 | 0.992 | 0.986 |
| 1 | 0.995 | 0.984 | 0.996 | 0.993 | 0.994 |
| 0.1 | 0.995 | 0.984 | 0.997 | 0.998 | 0.987 |
| 0.01 | 0.987 | 0.969 | 0.991 | 0.945 | 0.919 |
| 0.001 | | | | | |

**uTraVaS — Precision**

| Epsilon \ Delta | 0.5 | 0.1 | 0.05 | 0.01 | 0.001 |
|---|---|---|---|---|---|
| 2 | 0.481 | 0.533 | 0.514 | 0.474 | 0.710 |
| 1 | 0.553 | 0.486 | 0.460 | 0.468 | 0.823 |
| 0.1 | 0.498 | 0.568 | 0.529 | 0.548 | 0.846 |
| 0.01 | 0.480 | 0.506 | 0.496 | 0.525 | 0.903 |
| 0.001 | 0.440 | 0.512 | 0.569 | 0.528 | 0.850 |

**TraVaS — Precision**

| Epsilon \ Delta | 0.5 | 0.1 | 0.05 | 0.01 | 0.001 |
|---|---|---|---|---|---|
| 2 | 0.398 | 0.430 | 0.485 | 0.449 | 0.517 |
| 1 | 0.380 | 0.486 | 0.460 | 0.393 | 0.823 |
| 0.1 | 0.417 | 0.519 | 0.529 | 0.548 | 0.846 |
| 0.01 | 0.454 | 0.445 | 0.465 | 0.525 | 0.835 |
| 0.001 | 0.438 | 0.512 | 0.555 | 0.514 | 0.850 |

**Benchmark — Precision**

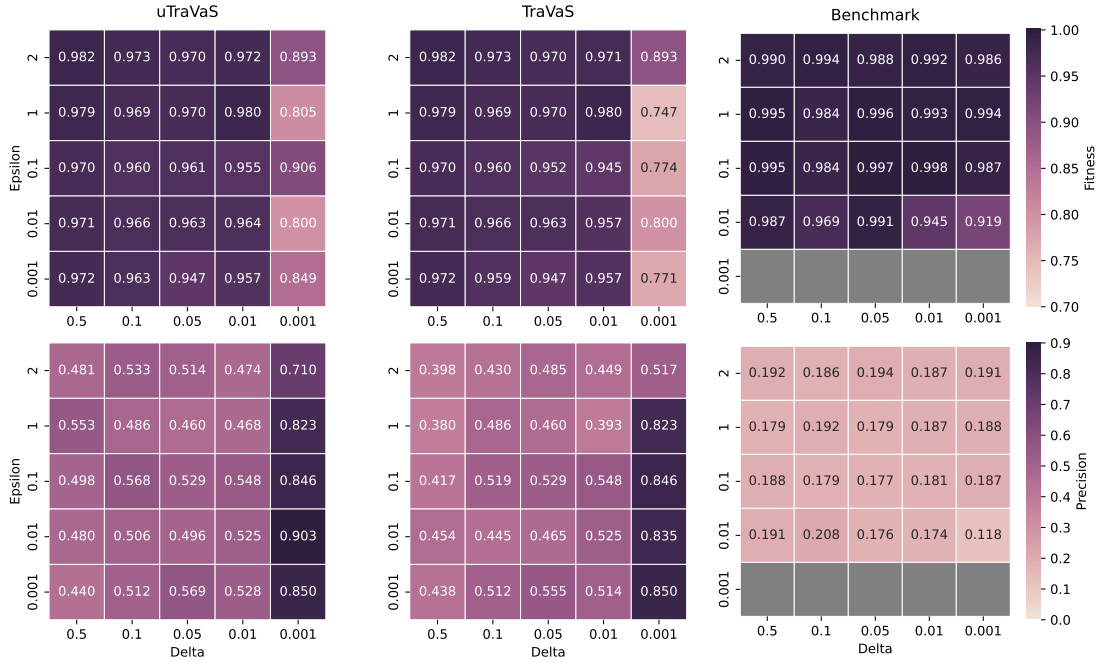| Epsilon \ Delta | 0.5 | 0.1 | 0.05 | 0.01 | 0.001 |
|---|---|---|---|---|---|
| 2 | 0.192 | 0.186 | 0.194 | 0.187 | 0.191 |
| 1 | 0.179 | 0.192 | 0.179 | 0.187 | 0.188 |
| 0.1 | 0.188 | 0.179 | 0.177 | 0.181 | 0.187 |
| 0.01 | 0.191 | 0.208 | 0.176 | 0.174 | 0.118 |
| 0.001 | | | | | |

Figure 1.4: The *fitness* and *precision* results of anonymized Sepsis event logs generated by *uTraVaS*, *TraVaS*, and our prefix-based benchmark method. Each value represents the mean of 10 algorithm runs.

## 1.3 Note of Authorship

This document is part of the supplementary material created for the paper "TraVaS: Differentially Private Trace Variant Selection for Process Mining", written by Majid Rafiei, Frederik Wangelik and Wil M.P. Van der Aalst. Please contact *frederik.wangelik@rwth-aachen.de* for further information.

# References

[1] C. Dwork, "Differential privacy: A survey of results," in *Theory and Applications of Models of Computation, 5th International Conference*, M. Agrawal, D. Du, Z. Duan, and A. Li, Eds. Springer, 2008.

[2] M. Rafiei and W. M. P. van der Aalst, "Towards quantifying privacy in process mining," in *Process Mining Workshops - ICPM 2020 International Workshops*, ser. Lecture Notes in Business Information Processing. Springer, 2020.