

Anonymization in Intelligent Surveillance Systems

Hauke Vagts, Christoph Bier and Jürgen Beyerer

Fraunhofer Institute of Optronics, System Technologies and Image Exploitation

Fraunhoferstr. 1 D-76131 Karlsruhe

Karlsruhe Institute of Technology—Vision and Fusion Laboratory

Adenauerring 4, D-76131 Karlsruhe

Abstract—Modern surveillance systems collect a massive amount of data. In contrast to conventional systems that store raw sensor material, modern systems take advantage of smart sensors and improvements in image processing. They extract relevant information about the observed objects of interest, which is then stored and processed during the surveillance process. Such high-level information is, e.g., used for situation analysis and can be processed in different surveillance tasks. Modern systems have become powerful, can potentially collect all kind of user information and make it available to any surveillance task. Hence, direct access to the collected high-level data must be prevented. Multiple approaches for anonymization exist, but they do not consider the special requirements of surveillance tasks. This work examines and evaluates existing metrics for anonymization and approaches for anonymization. Even though all kinds of data can be collected, position data is still the one with the highest demand. Hence, this work focuses on the anonymization of position data and proposes an algorithm that fulfills the requirements for anonymization in surveillance.

I. INTRODUCTION

Data protection and (video-) surveillance is an up-to-date topic. Since September 11th, 2001, public space is observed in almost every country. The United Kingdom is still the most observed country and current estimations about installed cameras differ between one million and 4.2 millions, thereof 500,000 in London. These impressive numbers are just estimated but they make the need for privacy-aware surveillance obvious.

Conventional video systems collect all available information, store it and perform situation analysis on the raw material. With the growing number of cameras and other sensors it is essential to extract required information as soon as possible to reduce the amount of data. In addition, data collected by other sensors (RFID, acoustics, etc.) must be combined with information gained out of the video material. Hence, data must be represented and processed on a high level of abstraction. The abstraction leads to new opportunities for privacy enforcement and a framework for surveillance systems that follows the fair information practice principles [1]. If personal data is only stored in conjunction with the corresponding object, anonymization strategies can be applied to the records. It is possible to maximize privacy efficiently fulfill surveillance tasks at the same time.

Metrics and strategies for anonymization already exist. This work examines to what extent they can be used in the context of surveillance. After a short introduction of the most important metrics for anonymization, it is discussed which metrics should be used in surveillance. In the following, the

requirements for anonymization in surveillance are pointed out and an algorithm for anonymization of position data in surveillance is presented.

II. RELATED WORK

Multiple architectures for intelligent surveillance systems exist, an overview of them can, e.g., be found in [2].

A. Intelligent Surveillance and Privacy

One of the architectures, which is following a task-oriented approach, is NEST [3]. Basically, an architecture of a modern and intelligent surveillance system consists of three parts. These are (smart) sensors that collect all relevant information for a surveillance task, a central storage for data and intelligent modules that process the stored data to fulfill the surveillance task. In NEST data is stored in an object-oriented world model [4], which is a virtual representation of a part of the real world, but other solutions are possible as well. However, anonymization requires that data is not stored in its raw format (e.g. video), but methods must rather be applied to sets of abstract data, e.g. position data of the observed objects, that contain fused information of all sensors. The anonymization itself is then independent from the data sources.

Existing approaches for privacy in surveillance aim at adding privacy to the video source itself. This is not sufficient, if different types of sensors are used and systems are working with data on a high level of abstraction. Schiff et al. [5] propose a system that identifies employees by marks that are applied to the observed objects. However, if the recognition of an objects fails, privacy cannot be enforced. A similar solution is proposed by Senior et al. [6]. They make use of a “privacy-preserving console” that manipulates a video stream and hides sensitive details. Fleck [7] proposes a more extreme approach that makes use of smart cameras. These cameras do not transmit video data, but rather high-level information, e.g., the position of a human combined with the information, whether he is standing or has fallen to the ground. Fidaleo et al. [8] propose a framework for video surveillance that uses a privacy buffer. According to privacy policies, the information is filtered and presented to the user.

B. Metrics for Anonymization

Metrics for anonymity have been compared [9], [10], but both works do not consider the requirements of surveillance systems. A detailed description of them would go beyond this

work. Hence, only k-Anonymity and l-diversity are roughly introduced. The definition of a table and attributes should be intuitive, for more details see the referenced work.

An explicit identifier is an attribute that can identify an object solely, without other attributes. A quasi-identifier (QI) is a set of attributes that can identify an object in combination with each other.[11]

It is important to distinguish between sensitive and non-sensitive attributes: sensitive attributes of a specific object must be hidden from attackers. Non-sensitive attributes are common knowledge or not related with any privacy concerns. It can be estimated that most attributes are QIs. In most work [12], [13], [14], QIs and sensitive attributes are considered to be not disjunct. This assumption cannot be made in surveillance. In particular, position data is a sensitive attribute and a QI.

The idea of k-Anonymity is motivated by the observation that QIs can destroy anonymity of a data set. An attacker might use background knowledge to identify objects [12]. K-anonymity has weaknesses for numeric values, a solution (k,e)-Anonymity is proposed in [14].

Definition 1 (k-Anonymity): Let $T(A_1, \dots, A_m)$ be a table and QI_T be the Quasi-identifier associated with it. T is said to satisfy k-Anonymity in relation to QI_T , if and only if each tuple $t \in T$ k-1 other tuples $t_{i_1}, t_{i_2}, \dots, t_{i_{k-1}} \in T$ exist, such that $\forall QI \in QI_T : t[QI] = t_{i_1}[QI] = t_{i_2}[QI] = \dots = t_{i_{k-1}}[QI]$.

L-Diversity extends k-Anonymity to prevent a distribution of the values of sensitive attributes in a manner an attacker can link a specific sensitive value to a specific object.

Definition 2 (l-Diversity principle): An equivalence class $QI-EC$ of a table T is l-diverse, if it contains at least l values for a sensitive attribute SA . A table T is l-diverse, if all of its $QI-EC$ are l-diverse.

Definition 3 (Entropy l-Diversity): A table T is Entropy l-diverse, if for all $QI-EC$ of T :

$$-\sum_{s \in SA} p_s \log_2 p_s \geq \log_2(l)$$

with p_s as part of the tuple in the equivalence class with $t[SA] = s$. SA denotes a sensitive attribute.

It follows, because $-x \log_2(x)$ is concave, for Entropy l-Diversity that the entropy of the entire table is at least $\log_2(l)$ [13]. This is a very high requirement, hence Entropy l-Diversity is sometimes too restrictive.

III. METRICS FOR THE ANONYMIZATION OF SURVEILLANCE DATA

Intelligent surveillance sets a standard for anonymization and the requirements differ depending on the surveillance task. This work focuses on surveillance of a single object. Position data is still of major importance and has a special characteristic that must be used during anonymization.

As mentioned above, position data is a sensitive attribute and a QI. In addition, it has its own semantic, which does not allow the use of regular methods for anonymization of numeric attributes.

In general, k-Anonymity and l-Diversity can be used for the anonymization of position data. In the context of surveillance k-Anonymity means that an anonymized position may be just as accurate that it is appropriate for at least k individuals. But K-Anonymity is not sufficient as a measure for the anonymity of position data. If there are many people together in a confined space the semantic meaning of their positions is almost identical. Hence, position data must be treated in a special manner.

The solution to this problem is the separation of the quasi-identifier and the sensitivity aspect of the position data. Therefore we introduce an auxiliary attribute area and an unique assignment f_h of a position (x,y) to an area r : $f_h : (x,y) \rightarrow r$. Following the externalization the position data is simply a quasi-identifier and area is the sensitive attribute. Afterwards you can use l-Diversity to handle the sensitive aspect.

IV. AN APPROACH FOR ANONYMIZATION IN SURVEILLANCE

After determining satisfying metrics it must be specified, which methods for anonymization are sufficient to achieve a specific level of anonymity.

A. Time

In existing work, the factor time has not been considered. However, time has an extensive impact on anonymization in surveillance. It can be used in three dimensions.

Firstly, as a *temporal variance* Δz . If the date of the observation is published with reduced accuracy, the position can be published with higher precision in return, while the anonymity is still at the same level.

Secondly, as a latency λ , i. e., the time a surveillance system can wait until it provides an answer to a query. More future data can be used for anonymization, when allowing a latency.

Thirdly, as the frequency for requests ϕ , which has an influence on the traceability of objects.

B. Grid versus Graph

Two approaches exist for anonymization of position data. Either the observed area can be split into a grid, or a graph can be used. In case of the latter, nodes represent the objects and the (weighted) edges represent the distance between them. When using a grid, it is the objective to find the smallest set of neighboring fields that fulfills the anonymity requirements. When using graphs a problem similar to the clique problem must be solved, which may result in a bad performance. However, even if a grid approach is not optimal, it has performance advantages and it can also be chosen in what direction the selected region is extended on the grid, which makes it more difficult to find an object in a raster field.

C. Finding a Suitable Algorithm

When following the grid-based approach, an algorithm can either work top-down or bottom-up. In [15] Grutser and Grundwald propose a hierarchical top-down algorithm that is based on a Quad Tree Algorithm. In each step the segment,

which contains the object, is picked and is split in four squares of the same size. The algorithm stops, if the number of objects in a segment is $< k$. The position is then replaced by the segment that was split. No matter, whether top-down or bottom-up, when using a hierarchical algorithm, the accuracy is drastically reduced with each step. Through the hierarchical structure the way of anonymization is predetermined.

A compromise between speed (grid-based method) and accuracy is the approach from Bamba and Lui [16]. Furthermore, it is the only approach that considers l-Diversity and k-Anonymity. The area that is released instead of a position must contain $k - 1$ other objects and must span over l regions. The algorithm is also based on a grid, but in each step only one segment of the grid is added to the released region (north, east, west, south). The algorithm offers the option to separate k-Anonymity, which is related to the objects and l-Diversity, which is related to the segments of the released region. Hence, our attribute region is not allocated to positions, instead it is related to segments of the grid. This results in many advantages:

- If k-Anonymity and l-Diversity are both related to objects, k is high on the one side and l is low on the other side. This effect is prevented.
- When observing a single object, l-Diversity ensures that the region is only related to the observed object. No prediction for other objects can be made by considering the borders.
- As k-anonymity and l-Diversity are decoupled, both parameters can be changed independently and according to the scenario. Both metrics can be weighted differently.
- The sensitivity of regions does not depend on the number of people that are comprised.
- The approach can be extended to consider multiple levels of sensitivity for different regions.

The named separation can only be performed, because a grid-based approach has been chosen. A trade-off between performance and accuracy can be achieved by setting the density of the grid and choosing the algorithm for extension of the region. A top-down approach is faster, if the size of the final segment is close to the size of the observed area. Such a region only contains a minimum of information and is not useful in practice. Hence, a bottom-up approach should be chosen.

As shown above, time is an important factor. Thus, an algorithm should be extended with another dimension. The grid then consists of cubes instead of squares. It must also be considered that time has only an influence on k-Anonymity and not on l-Diversity. This leads to three QoS parameters that should be used. The maximal temporal variance Δz_{max} , the maximal latency λ_{max} and the maximal frequency for requests ϕ_{max} .

V. AN ALGORITHM FOR POSITION DATA

As shown above, an algorithm that fulfills the requirements for anonymization of position data in surveillance should follow the approach from Bamba and Liu [16]. The existing

approach must be extended with the temporal dimension and the parameters k and l are to be handled separately. The algorithm determines a k-anonymous and l-diverse space-time cuboid.

Algorithm 1 Position Privacy

Require: $\{ID, z\}, \{\delta_{max}, \Delta z_{max}, \lambda\}, \{k, l\}$

Ensure: $\{[x_1, x_2], [y_1, y_2], [z_1, z_2]\}$

- 1: $z_2 \leftarrow \text{random}(z - \Delta z_{max}, z + \min \Delta z_{max}, \lambda)$
 - 2: $x_2 \leftarrow \text{xPosOf}(ID, z_2)$
 - 3: $y_2 \leftarrow \text{yPosOf}(ID, z_2)$
 - 4: $C \leftarrow \text{getCuboidOf}(x_2, y_2, z_2)$
 - 5: $C \leftarrow \text{FUNCTION_FIND_K_CUBOID}(C, z, \{\delta_{max}, \Delta z_{max}, \lambda\}, k)$
 - 6: $C \leftarrow \text{FUNCTION_FIND_L_CUBOID}(C, \delta_{max}, l)$
 - 7: **return** C.XYZ
-

In lines 1 to 4 the Position Privacy Algorithm determines the initial cuboid C for the anonymization. To extend the cuboid in the temporal dimension, a starting point must be chosen randomly out of the valid time interval. The anonymization after k (line 5) and l (line 6) itself takes place in two separated functions.

At first, the cuboid is extended to contain at least k objects (Algorithm 2, line 1). This is done within the restrictions given by the variables δ_{max} , Δz_{max} and λ (line 2). In the lines 5 to 18, the increment of k for the extension of the cuboid is determined in the different directions and dimensions (if an extension is possible). In the last step, the extension that leads to the highest increment of k is chosen (lines 19 to 23). This is repeated until the cuboid complies with the k value.

To achieve l-Diversity for the location, the cuboid must contain fields of the grid in a suitable diversity. The approach is similar to the k value (Algorithm 3). Time is not considered, as the room layout is static. Each field of the grid is assigned to the region ID of the region, which it covers for the most part.

VI. CONCLUSION AND FUTURE WORK

A lot of research has been done in the area of privacy and anonymization. Intelligent surveillance systems can imperil privacy. Hence, this work has analyzed the suitability of existing metrics and approaches for anonymization. Each observable attribute can result in privacy issues, but the position is the most important one. Thus, an algorithm for anonymization of position data in intelligent surveillance has been proposed.

The anonymization of other attributes is currently being analyzed and an approach is being developed. Currently, the presented approach for position data is implemented in a demonstration system of the NEST architecture. The system must then be tested under real time conditions.

Algorithm 2 FUNCTION_FIND_K_CUBOID

Require: Starting cuboid C , z , $\{\delta_{max}, \Delta z_{max}, \lambda\}$, k **Ensure:** Resulting cuboid C

```
1: while C.kValue < k do
2:   if C.sizeX + gridElementSizeX ≤  $\delta_{max}$  = false
     and C.sizeY + gridElementSizeY ≤  $\delta_{max}$  = false
     and C.sizeZ + gridElementSizeZ ≤  $\Delta z_{max}$  = false
   then
3:     return PrivacyNotPossibleError
4:   end if
5:   if C.sizeX + gridElementSizeX ≤  $\delta_{max}$  then
6:     cuboidExtension(S) ← C ∪ southern 3D row
7:     cuboidExtension(N) ← C ∪ northern 3D row
8:   end if
9:   if C.sizeY + gridElementSizeY ≤  $\delta_{max}$  then
10:    cuboidExtension(E) ← C ∪ eastern 3D row
11:    cuboidExtension(W) ← C ∪ western 3D row
12:  end if
13:  if C.sizeZ + gridElementSizeZ ≤  $\Delta z_{max}$  then
14:    if C.upperZ + gridElementSizeZ ≤  $z + \lambda$  then
15:      cuboidExtension(F) ← C ∪ 3D row in future
16:    end if
17:    cuboidExtension(P) ← C ∪ 3D row in the past
18:  end if
19:  for all  $d \in \{S, N, E, W, F, P\}$  do
20:    if C.kValue < cuboidExtension(d).kValue then
21:      C ← cuboidExtension(d)
22:    end if
23:  end for
24: end while
```

REFERENCES

- [1] H. Vagts and J. Beyerer, "Security and privacy challenges in modern surveillance systems," in *Future Security: 4th Security Research Conference*, P. Elsner, Ed. Fraunhofer Verlag, Oct. 2009, p. 94.
- [2] M. Valera and S. A. Velastin, "Intelligent distributed surveillance systems: a review," *IEEE Proceedings - Vision, Image and Signal Processing*, vol. 152, no. 2, p. 192, 2005.
- [3] J. Moßgraber, F. Reinert, and H. Vagts, "An architecture for a task-oriented surveillance system – a service and event based approach," in *Proc. Fifth International Conference on Systems ICONS*, 11–16 April 2010.
- [4] A. Bauer, T. Emter, H. Vagts, and J. Beyerer, "Object oriented world model for surveillance systems," in *Future Security: 4th Security Research Conference*, P. Elsner, Ed. Fraunhofer Verlag, 2009, p. 339.
- [5] J. Schiff, M. Meingast, D. K. Mulligan, S. Sastry, and K. Goldberg, "Respectful cameras: Detecting visual markers in real-time to address privacy concerns," in *Protecting Privacy in Video Surveillance*. Springer, 2009, p. 65.
- [6] A. W. Senior, S. Pankanti, A. Hampapur, L. M. G. Brown, Y. li Tian, A. Ekin, J. H. Connell, C.-F. Shu, and M. Lu, "Enabling video privacy through computer vision," *IEEE Security & Privacy*, vol. 3, no. 3, p. 50, 2005.
- [7] S. Fleck and W. Strasser, "Smart camera based monitoring system and its application to assisted living," *Proceedings of the IEEE*, vol. 96, no. 10, p. 1698, 2008.
- [8] D. A. Fidaleo, H.-A. Nguyen, and M. Trivedi, "The networked sensor tapestry (nest): a privacy enhanced software architecture for interactive analysis of data in video-sensor networks," in *VSSN '04: Proceedings of the ACM 2nd international workshop on Video surveillance & sensor networks*. New York, NY, USA: ACM, 2004, p. 46.

Algorithm 3 FUNCTION_FIND_L_CUBOID

Require: Starting cuboid C , δ_{max} , l **Ensure:** Resulting cuboid C

```
1: while C.LValue < l do
2:   if C.sizeX + gridElementSizeX ≤  $\delta_{max}$  = false
     and C.sizeY + gridElementSizeY ≤  $\delta_{max}$  = false
   then
3:     return PrivacyNotPossibleError
4:   end if
5:   if C.sizeX + gridElementSizeX ≤  $\delta_{max}$  then
6:     cuboidExtensionS ← C ∪ southern 3D row
7:     cuboidExtensionN ← C ∪ northern 3D row
8:   end if
9:   if C.sizeY + gridElementSizeY ≤  $\delta_{max}$  then
10:    cuboidExtensionE ← C ∪ eastern 3D row
11:    cuboidExtensionW ← C ∪ western 3D row
12:  end if
13:  for all  $d \in \{S, N, E, W\}$  do
14:    if C.LValue < cuboidExtension(d).LValue then
15:      C ← cuboidExtension(d)
16:    end if
17:  end for
18: end while
```

- [9] C. Andersson and R. Lundin, "On the fundamentals of anonymity metrics," in *The Future of Identity in the Information Society*, vol. 262. Springer, 2008, p. 325.
- [10] D. J. Kelly, R. A. Raines, M. R. Grimaila, R. O. Baldwin, and B. E. Mullins, "A survey of state-of-the-art in anonymity metrics," in *NDA '08: Proceedings of the 1st ACM workshop on Network data anonymization*. New York, NY, USA: ACM, 2008, p. 31.
- [11] L. Sweeney, "Achieving k-anonymity privacy protection using generalization and suppression," *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 10, 2002.
- [12] P. Samarati and L. Sweeney, "Protecting privacy when disclosing information: k-anonymity and its enforcement through generalization and suppression," Computer Science Laboratory, SRI International, Tech. Rep., 1998.
- [13] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkatasubramanian, "L-diversity: Privacy beyond k-anonymity," *ACM Transactions on Knowledge Discovery from Data*, vol. 1, no. 1, p. 3, 2007.
- [14] Q. Zhang, N. Koudas, D. Srivastava, and T. Yu, "Aggregate query answering on anonymized tables," in *Proceedings of the 23rd International Conference on Data Engineering*. Los Alamitos, CA, USA: IEEE Computer Society, 2007, p. 116.
- [15] M. Gruteser and D. Grunwald, "Anonymous usage of location-based services through spatial and temporal cloaking," in *MobiSys '03: Proceedings of the 1st international conference on Mobile systems, applications and services*. New York, NY, USA: ACM, 2003, p. 31.
- [16] B. Bamba and L. Liu, "Privacygrid: Supporting anonymous location queries in mobile environments," GIT-CERCS, Tech. Rep., 2007.