
Feature-level privacy loss modelling in differentially private machine learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We introduce *Tritium*, an automatic differentiation-based sensitivity analysis frame-
2 work for differentially private (DP) machine learning (ML). Optimal noise calibra-
3 tion in this setting requires efficient Jacobian matrix computations and tight bounds
4 on the L2-sensitivity. Our framework achieves these objectives by relying on a func-
5 tional analysis-based method for sensitivity tracking, which we briefly outline. This
6 approach interoperates naturally and seamlessly with static graph-based automatic
7 differentiation, which enables order-of-magnitude improvements in compilation
8 times compared to previous work. Moreover, we demonstrate that optimising the
9 sensitivity of the entire computational graph at once yields substantially tighter
10 estimates of the true sensitivity compared to interval bound propagation techniques.
11 Our work naturally befits recent developments in DP such as individual privacy
12 accounting, aiming to offer improved privacy-utility trade-offs, and represents a
13 step towards the integration of accessible machine learning tooling with advanced
14 privacy accounting systems.

15 1 Introduction

16 Despite the growing availability of high-performance algorithmic tools for advanced statistical
17 modelling and machine learning (ML), solutions to many of the world’s most important problems
18 require access to sensitive or confidential data. Technologies such as differential privacy (DP) can
19 allow drawing insights from such data while objectively allocating and quantifying individual privacy
20 expenditure. Although DP is the gold standard for data protection, its application to everyday ML
21 workflows is –in practice –often constrained. For one, tightly introspecting the privacy attributes
22 of complex models such as deep neural networks can be very challenging. Moreover, substantial
23 expertise is required on the analyst’s behalf to correctly apply DP mechanisms to such models.
24 Software libraries [1, 2, 3, 4] are being developed to alleviate these issues in specific domains such
25 as DP deep learning. They are, however, limited to a small number of programming languages and
26 application programming interfaces (APIs). The democratisation of DP machine learning therefore
27 awaits generic infrastructure, not only compatible with arbitrary workflows, but designed “from first
28 principles” to facilitate the implementation of DP. At its core, contemporary ML is based around
29 the manipulation of multidimensional arrays and the composition of differentiable functions, a
30 programming paradigm referred to as *differentiable programming*. Besides deep learning, some of
31 the most successful ML algorithms [5] and a large number of statistical queries, especially from the
32 domain of *robust statistics*, can be expressed within this paradigm. *Automatic differentiation* (AD)
33 systems are the core of differentiable programming frameworks and are able to track the flow of
34 computation to return precise derivatives with respect to arbitrary computational quantities. Although
35 this functionality may –at first –seem orthogonal to the goals described above, we contend that it is in
36 fact not only highly compatible, but *synonymous* with automatic DP tracking.

37 In the current work, we present *Tritium*, a differentiable programming framework aiming to integrate
 38 the requirements of ML and privacy analysis through the use of AD. We recapitulate the link between
 39 the *sensitivity* of differentiable queries and the *Lipschitz* constant in section 2. We outline our system’s
 40 implementation in section 3 and present the substantial improvements in computational efficiency
 41 and sensitivity bound tightness in section 4. A discussion of prior work can be found in the appendix.

42 2 Theoretical motivation

43 We begin by briefly introducing an interpretation of DP using the language of functional analysis,
 44 which forms the theoretical motivation behind our work. We concentrate on the Gaussian mechanism,
 45 which forms the basis for private data analysis in high dimensions. Differentially private ML can
 46 fundamentally be abstracted as the application of a higher-order function (or *functional*) to private
 47 data. This higher-order function (often termed a *DP mechanism*) receives as its input another function
 48 (termed a *query*) which has been applied to a private dataset, inspects the query to derive its privacy
 49 attributes and modifies it to preserve DP (Definitions 1 and 2).

50 **Definition 1** (Query). *A query is a function $q : \mathbb{R}^{m \times \mathcal{D}} \rightarrow \mathbb{R}^{n \times \mathcal{D}}$, where \mathcal{D} represents arbitrary
 51 (possibly unused) dimensions and $n, m \geq 1$ which receives as input some private dataset \mathbf{x} , $|\mathbf{x}| \geq 1$
 52 and outputs a result \mathbf{y} representing the result of a computation over \mathbf{x} (e.g. a mean calculation or the
 53 output of a neural network).*

54 **Definition 2** (DP mechanism). *A DP mechanism is a higher-order function M which receives
 55 as its input one or more query functions q_1, q_2, \dots, q_n and outputs $M(q_1 \circ q_2 \circ \dots \circ q_n) =$
 56 $q_n(\dots q_2(q_1(\mathbf{x}))) + \xi$, where $\xi \sim \mathcal{N}(0, C)$ and C is selected based on the privacy properties
 57 of $q_1 \circ q_2 \circ \dots \circ q_n$.*

58 The tight characterisation of these privacy properties is central to enabling privacy expenditure
 59 tracking. The effect on inputs on the output of the query functions is reflected in query *sensitivity*.
 60 We use the *Lipschitz* constant to reason about sensitivity.

61 **Definition 3** (Lipschitz constant and sensitivity). *Let $q : X \rightarrow Y$ be a function between metric
 62 spaces X and Y with distance metrics d_X and d_Y , respectively. Then q is Lipschitz continuous with
 63 constant K_q (equivalently, “ K -Lipschitz”) if*

$$d_Y(f(x), f(x')) \leq K d_X(x, x') \forall x, x' \in X \quad (1)$$

64 *The smallest value of K corresponds to the sensitivity of q , $\Delta_2(q)$ [6]:*

$$\Delta_2(q) = \max_{x, x'} \|q(x) - q(x')\|_2 \quad (2)$$

65 *where $\|\cdot\|_2$ is the L_2 distance. Recall that in this case, $d(x, x') = 1$ as x, x' are adjacent. Thus for
 66 differentiable query functions, $K_q \equiv \Delta_2(q)$ when X and Y are Euclidean spaces endowed with the
 67 L_2 -norm. Then,*

$$K_q = \sup \|\mathcal{J}(q)\|_2 \quad (3)$$

68 *where \mathcal{J} is the Jacobian matrix (the differential operator).*

69 This equivalence between Lipschitz constant and query sensitivity allows, in principle, to reason over
 70 the privacy attributes of individual query functions and calibrate noise appropriately. Typical functions
 71 with globally bounded sensitivity are affine queries or linear functions of the form $q(x) = \alpha x(+\beta)$,
 72 with $K_q = \alpha$. However, queries exist for which the Lipschitz constant is not defined over the
 73 entire input domain. One example of such a function is $q(x) = x^2 = x \cdot x$ with $K_q = x$, whose
 74 sensitivity is *unbounded*, as it depends on the value of x . The sensitivity analysis of such queries
 75 sometimes therefore depends on (private) properties of the dataset. We term such a case as *data-*
 76 *dependent sensitivity*. Reasoning over sensitivity in such cases is complicated by a requirement to
 77 propagate this data dependency effect through function composition. Previous works on Lipschitz
 78 analysis of machine learning algorithms [7, 8] achieve this through techniques such as *interval bound*
 79 *propagation* [9], that carry the bounds on input variables (which, for DP, should be defined in a
 80 data-agnostic manner) through the computation flow. This technique can easily be made compatible
 81 with *tracing-type* AD systems which are widely used in contemporary machine learning. However,
 82 due to well-known limitations of interval arithmetic (such as interval dependency [10]) and due to

83 the fact that the Lipschitz constant is defined by inequality, the resulting sensitivity terms may be
 84 valid, but too loose to be of any practical utility (e.g. 10^5). The last challenge relates to the fact that
 85 the *actual* effect of an individual’s data on the query’s output may, in fact, be much smaller than the
 86 worst case assumed by the definition, resulting in more noise being added by the mechanism than
 87 would be required for the guarantee to hold. Consequently, although the worst-case sensitivity value
 88 has typically been used for privacy accounting, newer techniques perform accounting based not only
 89 the worst case but combine it with the actual output L_2 -norm [11].

90 3 Implementation details

91 Our work presents *Tritium*, an automatic-differentiation-based machine learning and sensitivity
 92 analysis system engineered to address the above-mentioned challenges. It consists of the following
 93 components:

- 94 1. A user-facing front-end to specify a query $\mathbf{q} = q_1 \circ \dots \circ q_n$ *abstractly*, i.e. without directly
 95 utilising private data during model creation. This is achieved through the utilisation of
 96 abstract tensors with pre-defined dimensions. The system creates an optimised computational
 97 graph \mathcal{G} based on this specification.
- 98 2. During model specification, the user can impose *bounds* on the quantities (e.g. inputs,
 99 weights) used in the model.
- 100 3. The user selects the desired *privacy parameters*, e.g. ϵ and δ values or a maximum allowed
 101 sensitivity.
- 102 4. A *compiler* then emits a program which receives a private dataset and outputs an appropri-
 103 ately privatised result.

104 Internally, *Tritium* undertakes the following steps:

- 105 1. The computational graph \mathcal{G} is compiled into a program which outputs $\mathcal{J}(\mathbf{q})$ with respect to
 106 the inputs.
- 107 2. $K_q = \sup \|\mathcal{J}(\mathbf{q})\|_2$ is computed given the input bounds.
- 108 3. Finally, \mathcal{G} is compiled into the program described in step (4) above which receives a private
 109 dataset \mathbf{x} , computes $\mathbf{q}(\mathbf{x})$, potentially *clips* out-of-bound values to preserve the required
 110 K_q , adds noise $\xi \sim \mathcal{N}(0, C)$ with C proportional to K_q to satisfy the required ϵ value for a
 111 given δ and outputs $M(\mathbf{q})$.

112 This system architecture has several benefits: It avoids utilising private data until the moment the final
 113 computation is executed (*data minimisation*). Moreover, it provides a tight sensitivity calculation by
 114 optimising the entire query function at once [12] instead of the above-mentioned forward-propagation,
 115 which can lead to vacuous sensitivity values. Furthermore, it utilises the pre-specified bounds on
 116 the input variables to not only enable the calculation of *data-dependent* sensitivity, but also greatly
 117 accelerate the process. Moreover, it is *agnostic* to the method used to actually *obtain* the desired
 118 sensitivity. For example, *Lipschitz neural network layers* [13, 14] or activation functions with
 119 bounded outputs and gradients [15] can be used for model building, but bounded sensitivity can also
 120 be enforced by clipping, as is common in DP-SGD [16]. In addition, the system is able to compute
 121 the full Jacobian matrix (which is required in DP-SGD) as well as arbitrary higher-order derivative
 122 matrices (which can be used to accelerate the sensitivity computation). Additionally, as the system
 123 outputs both the Lipschitz constant and the norm of the outputs, it can be leveraged to provide tighter
 124 privacy guarantees through the use of e.g. *individual privacy accounting*, as shown below. Finally, the
 125 system is designed to output privatised values by default instead of outputting non-private values and
 126 relying on the user to perform an appropriate privatisation step. This can reduce both user workload
 127 and the probability of failure due to incorrect application of DP mechanisms on the user’s behalf.

128 4 Experimental evaluation

129 4.1 Exact sensitivity calculations through *a posteriori* optimisation

130 To assess the benefits of computing query sensitivity by assessing the entire computational graph
 131 at once instead of forward-propagating interval bounds, we constructed a small neural network

132 comprising 4 *linear* layers with *logistic sigmoid* activations and the *binary cross-entropy* cost function.
133 We set the bounds for the neural network weights and the input bounds to the $[0, 1]$ interval. We
134 then calculated the sensitivity using two techniques: *Interval Bound Propagation* (IBP [9]) and our
135 proposed method optimising the entire computational graph at once. The estimate of the sensitivity
136 was returned as $[0.0, 22175.37]$ by IBP (which is a valid, but vacuous bound) and as 0.99929 by
137 *Tritium*. The IBP bounds are also similar to previous work (compare e.g. [17]). However, IBP
138 was faster, requiring 103 ms to return a result, compared with 905 ms for our technique (excluding
139 compilation time of ca. 3 s). These results are summarised in Table 1 in the appendix.

140 4.2 Compilation improvements

141 In this section, we compare the compilation and execution time improvements of *Tritium* to the previ-
142 ously described framework by [18] on neural network architectures with the architecture described
143 above, but with increasing width of linear layers. We recall that the system proposed by the authors
144 of this work relies on a scalar AD implementation and authors report compilation times quadratic
145 in the number of model parameters (that is, around 60 hours for a 2.5 million parameter model). In
146 contrast, the here-presented implementation utilises the *Aesara* library (a fork of the now-defunct
147 *Theano* framework [19]) as a computational back-end. This allows both a memory-efficient vectorised
148 execution of tensor operations and a more mature and faster compilation back-end. For all but the
149 smallest architectures, this back-end achieved substantially faster compilation times which were
150 independent of the number of parameters and able to leverage caching to accelerate re-compilation.
151 A similar effect was observed in execution times, where our system achieved considerably higher
152 performance. These results are visualised in Figure 1 in the appendix.

153 5 Discussion and conclusion

154 We propose *Tritium*, an automatic differentiation-based system for differentially private machine
155 learning. Our framework relies on an interpretation of DP queries and mechanisms through the
156 language of functional analysis, linking them by the definition of Lipschitz continuity. We found
157 the combination of an efficient computational and compilation back-end with the consideration of
158 the entire query function at once to yield both improved performance and tighter sensitivity bounds
159 compared to previous work. Our proposed framework relies on *static graph-based* AD, which
160 can apply specific compiler optimisations to the entire computational graph and is responsible for
161 *Tritium*'s high performance. However, such systems have noteworthy limitations. For instance,
162 the definition of control flow statements is cumbersome and such systems are not well-suited for
163 utilisation with *just-in-time* compilers. Most currently used machine learning frameworks utilise
164 *tracing/leager execution* AD back-ends, which, while more user-friendly, cannot always leverage the
165 same optimisations. An alternative AD implementation, *source-to-source translation* can combine
166 the benefits of dynamic graph specification with the high performance and optimisations of static
167 compilation. It forms the basis of a recent paradigm in programming language design (e.g. [20]),
168 attempting to merge a general-purpose programming language with differentiable programming
169 primitives. A limitation of our method stems from the computational hardness of exactly calculating
170 the Lipschitz constant [21]. Our system will output the true bound using *Simplicial Homology*
171 *Global Optimisation* [22] if the bound exists and can be found, and the application of constraints can
172 substantially accelerate this process. However, it will output a warning and switch to an approximate
173 algorithm without a guaranteed bound otherwise. If such a bound is undesirable, an alternative
174 technique is the utilisation of model components with known (or manually adjustable) Lipschitz
175 constants, which can allow one to avoid the utility penalty imposed by clipping-based approaches,
176 both enabling the design of algorithms with milder privacy-utility penalties and additionally reaping
177 the benefits of well-defined model sensitivity, such as (certifiable) robustness to perturbations by
178 adversarial samples. In conclusion, our work serves as a first proof-of-concept for the the design of
179 generic infrastructure exposing familiar APIs to data scientists while automatically tracking privacy
180 loss through the computation flow. We view the further development of such systems as an accelerator
181 for the wide-spread adoption of privacy-preserving machine learning algorithms across data-driven
182 research disciplines.

References

- 183
- 184 [1] Opacus PyTorch library. Available from opacus.ai, 2021.
- 185 [2] TensorFlow Privacy. Available from TensorFlow Privacy, 2021.
- 186 [3] Naoise Holohan, Stefano Braghin, Pól Mac Aonghusa, and Killian Levacher. Diffprivlib: The
187 ibm differential privacy library, 2019.
- 188 [4] Differential Privacy. Available from google/differential-privacy, 2021.
- 189 [5] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of*
190 *the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages
191 785–794, 2016.
- 192 [6] Sofya Raskhodnikova and Adam Smith. Lipschitz extensions for node-private graph statis-
193 tics and the generalized exponential mechanism. In *2016 IEEE 57th Annual Symposium on*
194 *Foundations of Computer Science (FOCS)*, pages 495–504. IEEE, 2016.
- 195 [7] Aritra Bhowmick, Meenakshi D’Souza, and G Srinivasa Raghavan. Lipbab: Computing exact
196 lipschitz constant of relu networks. *arXiv preprint arXiv:2105.05495*, 2021.
- 197 [8] Sungyoon Lee, Jaewook Lee, and Saerom Park. Lipschitz-certifiable training with a tight outer
198 bound. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances*
199 *in Neural Information Processing Systems*, volume 33, pages 16891–16902. Curran Associates,
200 Inc., 2020.
- 201 [9] Sven Gowal, Krishnamurthy Dvijotham, Robert Stanforth, Rudy Bunel, Chongli Qin, Jonathan
202 Uesato, Relja Arandjelovic, Timothy Mann, and Pushmeet Kohli. On the effectiveness of
203 interval bound propagation for training verifiably robust models, 2019.
- 204 [10] Walter Krämer. Generalized intervals and the dependency problem. *PAMM*, 6(1):683–684,
205 December 2006.
- 206 [11] Vitaly Feldman and Tijana Zrnic. Individual privacy accounting via a renyi filter. *arXiv preprint*
207 *arXiv:2008.11193*, 2020.
- 208 [12] Henry Gouk, Eibe Frank, Bernhard Pfahringer, and Michael J Cree. Regularisation of neural
209 networks by enforcing lipschitz continuity. *Machine Learning*, 110(2):393–416, 2021.
- 210 [13] Yonadav Shavit and Boriana Gjura. Exploring the use of lipschitz neural networks for automating
211 the design of differentially private mechanisms. 2019.
- 212 [14] Cem Anil, James Lucas, and Roger Grosse. Sorting out lipschitz function approximation. In
213 *International Conference on Machine Learning*, pages 291–301. PMLR, 2019.
- 214 [15] Nicolas Papernot, Abhradeep Thakurta, Shuang Song, Steve Chien, and Úlfar Erlingsson.
215 Tempered sigmoid activations for deep learning with differential privacy. *arXiv preprint*
216 *arXiv:2007.14191*, 2020.
- 217 [16] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar,
218 and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC*
219 *conference on computer and communications security*, pages 308–318, 2016.
- 220 [17] Huan Zhang, Pengchuan Zhang, and Cho-Jui Hsieh. Recurjac: An efficient recursive algorithm
221 for bounding jacobian matrix of neural networks and its applications. In *Proceedings of the*
222 *AAAI Conference on Artificial Intelligence*, volume 33, pages 5757–5764, 2019.
- 223 [18] Alexander Ziller, Dmitrii Usynin, Moritz Knolle, Kritika Prakash, Andrew Trask, Rickmer
224 Braren, Marcus Makowski, Daniel Rueckert, and Georgios Kaissis. Sensitivity analysis in
225 differentially private machine learning using hybrid automatic differentiation. 2021.
- 226 [19] Theano Development Team. Theano: A Python framework for fast computation of mathematical
227 expressions. *arXiv e-prints*, abs/1605.02688, May 2016.
- 228 [20] Brennan Saeta, Denys Shabalin, Marc Rasi, Brad Larson, Xihui Wu, Parker Schuh, Michelle
229 Casbon, Daniel Zheng, Saleem Abdulrasool, Aleksandr Efremov, Dave Abrahams, Chris Lattner,
230 and Richard Wei. Swift for tensorflow: A portable, flexible platform for deep learning, 2021.
- 231 [21] Kevin Scaman and Aladin Virmaux. Lipschitz regularity of deep neural networks: analysis and
232 efficient estimation. *arXiv preprint arXiv:1805.10965*, 2018.

- 233 [22] Stefan C Endres, Carl Sandrock, and Walter W Focke. A simplicial homology algorithm for
234 lipschitz optimisation. *Journal of Global Optimization*, 72(2):181–217, 2018.
- 235 [23] Yuichi Yoshida and Takeru Miyato. Spectral norm regularization for improving the generaliz-
236 ability of deep learning. *arXiv preprint arXiv:1705.10941*, 2017.
- 237 [24] Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. Certified
238 robustness to adversarial examples with differential privacy. In *2019 IEEE Symposium on*
239 *Security and Privacy (SP)*, pages 656–672. IEEE, 2019.
- 240 [25] Varun Gupta, Christopher Jung, Seth Neel, Aaron Roth, Saeed Sharifi-Malvajerdi, and Chris
241 Waites. Adaptive machine unlearning. *arXiv preprint arXiv:2106.04378*, 2021.
- 242 [26] Chike Abuah, Alex Silence, David Darais, and Joe Near. Dduo: General-purpose dynamic
243 analysis for differential privacy. *arXiv preprint arXiv:2103.08805*, 2021.
- 244 [27] Paolo Pistone. From identity to difference: A quantitative interpretation of the identity type.
245 *arXiv preprint arXiv:2107.06150*, 2021.

	Upper bound	Computation time (ms)
Ours	0.99929	905
IBP	22175.37	103

Table 1: Comparison of the sensitivity bounds and computation times for our proposed framework (*Ours*) vs. Interval Bound Propagation (*IBP*).

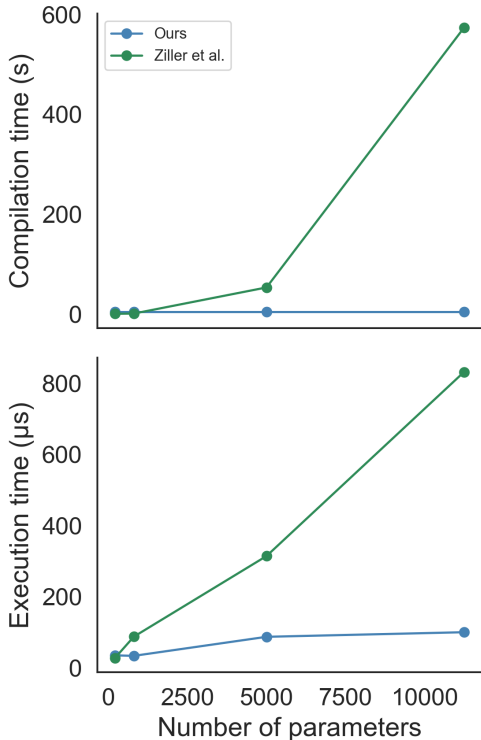


Figure 1: Performance comparison between our proposed framework (*Ours*, blue) vs. [18] (green). Compilation times in s and execution times in μs are shown for neural network architectures of increasing sizes.

247 **B Related works**

248 Our work can be seen as a natural evolution of the previous study by [18] in the context of AD-
 249 based sensitivity analysis for DP machine learning. In comparison to this work, *Tritium* relies on
 250 a vectorised, GPU-compatible execution engine and a mature graph compiler which drastically
 251 improves performance, as shown in the experimental section. The properties of Lipschitz continuous
 252 functions have been leveraged in several domains beyond DP. Works such as [13, 23, 14, 17] attempt
 253 to constrain the Lipschitz constant to reason over and control the properties of neural networks. The
 254 utilisation of this approach has been proposed for network certification against adversarial samples
 255 [24], whereby a network that is ϵ -certified is provably robust to input perturbations within a norm
 256 ball of radius ϵ . Additionally, constraining the Lipschitz constant has been proposed for DP model
 257 training, as this allows to calibrate the noise addition based on the bounded Lipschitz constant
 258 [13]. Moreover, certain works [25] have addressed the problem of *machine unlearning*, providing
 259 methods for a reliable removal of contributions associated with an individual in the context of neural
 260 network training. We note that the approach to sensitivity analysis employed in these studies is

261 orthogonal to our work, as our work is compatible with manual sensitivity constraints (such as directly
262 adjusting the Lipschitz constant of neural network layers through appropriate layers as described
263 above, which however may impair their expressivity) but also with sensitivity tracking for privacy
264 loss calculation. Several recent works concentrate on the computation of accurate estimates of the
265 Lipschitz constant mostly focused on *ReLU* networks such as [21, 14], but most of these obtain
266 the upper bounds rather than the exact values of the Lipschitz constant, often resulting in valid, but
267 extremely loose approximations that are not intended to be applied in DP training. A line of work
268 centred on languages for differentially private programming also exists. Among these, the recently
269 proposed *DDuo* framework [26] performs dynamic sensitivity analysis in the context of DP algorithm
270 specification. As shown above however, this approach does not attempt to derive a tight bound on
271 sensitivity in the setting of unbounded queries, declaring sensitivity as *infinite* and relying exclusively
272 on clipping. More general approaches, including *category-theoretical* views on the intersection of
273 differentiable programming and differential privacy such as [27] have also recently been proposed.