# Differential Privacy and the Power of (Formalizing) Negative Thinking
## (Extended Abstract)

Cynthia Dwork

Microsoft Research, Silicon Valley
dwork@microsoft.com

**Abstract.** *Differential privacy* is a promise, made by a data curator to a data subject: you will not be affected, adversely or otherwise, by allowing your data to be used in any study, no matter what other studies, data sets, or information from other sources is, or may become, available. This talk describes the productive role played by negative results in the formulation of differential privacy and the development of techniques for achieving it, concluding with a new negative result having implications related to participation in multiple, independently operated, differentially private databases.

**Keywords:** differential privacy, foundations of private data analysis, lifetime privacy loss, independently operated differentially private databases.

In the digital information realm, loss of privacy is usually associated with failure to control access to information, to control the flow of information, or to control the purposes for which information is employed. *Differential privacy* arose in a context in which ensuring privacy is a challenge even if all these control problems are solved: privacy-preserving statistical analysis of data. Here, even defining the goal is problematic, as the data analyst and the *privacy adversary* are one and the same.

*The Formal Definition.* A database is modeled as a collection of *rows*, with each row containing the data of a different individual. Differential privacy will ensure that the ability of an adversary to inflict harm – of any sort, to any set of people – should be essentially the same, independent of whether any individual opts in to, or opts out of, the dataset. This is done indirectly, by focusing on the probability of any given output of a privacy mechanism and how this probability can change with the addition or deletion of any row. Thus, we concentrate on pairs of databases $(D, D')$ differing only in one row, meaning one is a subset of the other and the larger database contains just one additional row. Finally, to handle worst case pairs of databases, the probabilities will be over the random choices made by the privacy mechanism.

**Definition 1.** [3,4] *A randomized function $\mathcal{K}$ gives $(\varepsilon, \delta)$-differential privacy if for all data sets $D$ and $D'$ differing on at most one row, and all $S \subseteq Range(\mathcal{K})$,*

$$\Pr[\mathcal{K}(D) \in S] \leq \exp(\varepsilon) \times \Pr[\mathcal{K}(D') \in S] + \delta \qquad (1)$$

*where the probability space in each case is over the coin flips of $\mathcal{K}^1$.*

Both the definition and the earliest techniques for achieving it were strongly influenced by negative results [2,5].

Consider a differentially private mechanism answering simple "counting queries" of the form "How many people in the database have property $P$?" Since a differentially private mechanism exhibits a similar probability distribution on answers for neighboring databases $D, D'$, it is clear that the responses given must sometimes be inaccurate; the goal of algorithmic research in this field is to minimize this inaccuracy. The *Laplace method* achieves $(\varepsilon, 0)$-differential privacy for counting queries by adding noise generated according to the Laplace distribution with parameter $1/\varepsilon$ to the true answer, and releasing this "noisy" value [4]. The resulting expected error is on the order of $1/\varepsilon$.

Differential privacy holds regardless of what the adversary knows, now or in the future. In consequence, differential privacy composes obliviously and automatically; the $k$-fold composition of $(\varepsilon, \delta)$-differentially private mechanisms, involing either $k$ operations on a single database or the mutually oblivious operation of $k$ independent databases with arbitrary overlap, is still roughly $(\sqrt{k}\varepsilon, \delta')$-differentially private [7]. It follows that adding independently generated noise with distribution $\text{Lap}(\sqrt{k}/\varepsilon)$ permits $k$ queries to be answered with a total privacy loss of about $\varepsilon$ and expected error per query $\sqrt{k}/\varepsilon$.

A series of results beginning with [1] shows one can do much better – with error depending polylogarithmically on the number of queries – using *coordinated noise*. In fact *coordination is essential* to beat the "$\sqrt{k}$" composition bound [6], so we have reached the end of the line for mutually oblivious, independently operated, differentially private mechanisms running against arbitrarily knowledgeable adversaries. Addressing this newly understood limitation is a fundamental challenge in differentially private data analysis.

# References

1. Blum, A., Ligett, K., Roth, A.: A Learning Theory Approach to Non-Interactive Database Privacy. In: Proc. 40th ACM Symposium on Thoery of Computing (2008)
2. Dinur, I., Nissim, K.: Revealing Information While Preserving Privacy. In: Proc. 22nd ACM Symposium on Principles of Database Systems (2003)
3. Dwork, C.: Differential Privacy. In: Bugliesi, M., Preneel, B., Sassone, V., Wegener, I. (eds.) ICALP 2006. LNCS, vol. 4052, pp. 1–12. Springer, Heidelberg (2006)
4. Dwork, C., McSherry, F., Nissim, K., Smith, A.: Calibrating Noise to Sensitivity in Private Data Analysis. In: Halevi, S., Rabin, T. (eds.) TCC 2006. LNCS, vol. 3876, pp. 265–284. Springer, Heidelberg (2006)
5. Dwork, C., Naor, M.: On the Difficulties of Disclosure Prevention in Statistical Databases or The Case for Differential Privacy. Journal of Privacy and Confidentiality 2(1) (2010)
6. Dwork, C., Naor, M., Vadhan, S.: Coordination is Essential (working title) (manuscript in preparation)
7. Dwork, C., Rothblum, G., Vadhan, S.: Boosting and Differential Privacy. In: Proceedings of the 51st IEEE Symposium on Foundations of Computer Science (2010)

---

[1] Typcially, the literature considers $\delta < 1/\text{poly}(n)$ on a database of size $n$.