



A Programming Framework for OpenDP[‡]

Marco Gaboardi[‡]

Michael Hay[§]

Salil Vadhan[¶]

February 11, 2021

Abstract

In this working paper, we propose a programming framework for the library of differentially private algorithms that will be at the core of the OpenDP open-source software project, and recommend programming languages in which to implement the framework.

Note: The current version of this paper is intended as a proposal for discussion at the May 2020 OpenDP Community Meeting, rather than as a scholarly paper that identifies novel contributions. In particular, proper credits to and comparisons with prior work are still missing (but will be added in a future version).

*CC BY 4.0: <https://creativecommons.org/licenses/by/4.0/>

[†]We thank the rest of the OpenDP Design Committee (Mercè Crosas, James Honaker, Gary King, Aleksandra Korolova, and Ilya Mironov) for numerous discussions and feedback that contributed to this paper.

[‡]Department of Computer Sciences, Boston University. Supported by NSF through grants 1565365 and 1845803.

[§]Department of Computer Sciences, Colgate University. Supported by NSF through grant 1409125 and by DARPA and SPAWAR under contract N66001-15-C-4067.

[¶]School of Engineering and Applied Sciences, Harvard University. Supported by a grant from the Sloan Foundation and a Simons Investigator Award.

Contents

1	Goals	2
2	Stable Transformations on Datasets and Pure DP	3
2.1	Measurements	3
2.2	Transformations	4
2.3	Chaining	5
2.4	Composition	6
2.5	Post-processing	7
3	Verifying Privacy Properties	8
4	Varying Types and Distance Measures	9
4.1	Private Data Types	9
4.2	Privacy Relations	12
4.3	Stability Relations	13
4.4	Multiple Distance Measures at Once	14
4.5	Chaining	17
4.6	Assumptions on the Measures	19
4.7	Composition	19
5	Interactive Measurements	20
5.1	Defining Interactive Measurements	20
5.2	Post-processing	24
5.3	Chaining and Composition of Interactive Measurements	26
5.4	Interactive Transformations	28
6	Compute Environment and Interactive Transformations	29
6.1	Working toward Data DAGs	35
7	Putting it Together	37
8	Ensuring Privacy in Implementation	38
9	Discussion of Implementation Language(s)	39
9.1	Desired language support	39
9.2	Proposal: Rust + Python	40
10	Using the Library	41
11	Contributing to the Library	48
12	The Scope of the Framework	49
12.1	Within the Current Framework	49
12.2	Outside the Current Framework	49
13	Comparison with other Programming Frameworks	50
14	Representing Data Domains	51

1 Goals

In this paper, we propose a programming framework for the library of differentially private algorithms that will be at the core of the OpenDP open-source software project, and recommend programming languages in which to implement the framework.

There are a number of goals we seek to achieve with this programming framework and language choice:

Extensibility We would like the OpenDP library to be able to expand and advance together with the rapidly growing differential privacy literature, through external contributions to the codebase. This leads to a number of the other desiderata listed below.

Flexibility The programming framework should be flexible enough to incorporate the vast majority of existing and future algorithmic developments in the differential privacy literature. It should also be able to support many variants of differential privacy. These variants can differ in the type of the underlying sensitive datasets and granularity of privacy (e.g. not only tabular datasets with record level-privacy, but also graph datasets with node-level privacy and datastreams with user-level privacy), as well as in the distance measure between probability distributions on adjacent datasets (e.g. not only pure and approximate differential privacy, but also Rényi and concentrated differential privacy).

Verifiability External code contributions need to be verified to actually provide the differential privacy properties they promise. We seek a programming framework that makes it easier for the OpenDP Editorial Board and OpenDP Committers to verify correctness of contributions. Ideally, most contributions will be written by combining existing components of the library with built-in composition primitives, so that the privacy properties are automatically derived. Human verification should mostly be limited to verifying mathematical proofs for new building blocks, which should usually be small components with clear specifications. At the same time, if an error is later found in a proof for a building block, it should be possible to correct that component locally and have the correction propagate throughout all the code that makes use of that building block.

Programmability The programming framework should make it relatively easy for programmers and researchers to implement their new differentially private algorithms, without having to learn entirely new programming paradigms or having to face excessive code annotation burdens.

Modularity The library should be composed of modular and general-purpose components that can be reused in many differentially private algorithms without having to rewrite essentially the same code. This supports extensibility, verifiability, and programmability.

Usability It should be easy to use the OpenDP Library to build a wide variety of DP Systems that are useful for OpenDP’s target use cases. These may have varying front-end user interfaces (e.g. programming in Python notebooks versus a graphical user interface) and varying back-end storage, compute, and security capabilities.

Efficiency It should be possible to compile algorithms implemented in the programming framework to execute efficiently in the compute and storage environments that will occur in the aforementioned systems. *[[SV: it’d be great if someone can elaborate here, mentioning examples of the range of compute and storage environments]]*

Utility It is important that the algorithms in the library expose their utility or accuracy properties to users, both prior to being executed (so that “privacy loss budget” is not wasted on useless computations) and after being executed (so that analysts do not draw incorrect statistical conclusions). When possible, algorithms should expose uncertainty measures that take into account both the noise due to privacy and the statistical sampling error.

In Sections 2–7, we build up the ideas in our proposed framework by starting with the ideas underlying existing differential privacy systems like PINQ [8], Ektelo [13], and Fuzz [11], and explaining how we generalize

them to achieve the above goals. We present the framework with both mathematical specifications and with Python-like code snippets. The code snippets are meant solely for illustrative purposes, and not meant to be a proposal for actual implementation (which is discussed more in Sections 8 and 9). *[[SV: edited this paragraph a lot]]*

[[SV: it would be nice to include a running example of a concrete dataset (e.g. an education dataset) and queries that one might want to perform on it]]

2 Stable Transformations on Datasets and Pure DP

In the most standard model of differential privacy and in systems like PINQ, a sensitive dataset x is represented as a *multiset* of records, each of which is from some data domain \mathcal{X} , which we'll write as $x \in \text{MultiSets}(\mathcal{X})$. For example, for a dataset where each record has a student's last name, age, and GPA, we might have $\mathcal{X} = \{A, B, \dots, Z\}^* \times \mathbb{N} \times \mathbb{R}$. The *distance* between two datasets $x, x' \in \text{MultiSets}(\mathcal{X})$ is the size of their symmetric difference $d_{\text{Sym}}(x, x') = |x \Delta x'|$.¹ We say that x and x' are *adjacent*, written $x \sim x'$, if $d_{\text{Sym}}(x, x') \leq 1$, i.e. x and x' differ by adding or removing at most one record. The two main kinds of operators in PINQ and Ektelo are *measurements* and *transformations*.

2.1 Measurements

A *measurement* M is a randomized mapping from datasets to outputs of an arbitrary type. That is, $M : \text{MultiSets}(\mathcal{X}) \rightsquigarrow \mathcal{Y}$, where we are using the squiggly arrow \rightsquigarrow to denote a randomized function. A measurement M is ε -DP if for every pair of adjacent datasets $x \sim x' \in \text{MultiSets}(\mathcal{X})$, the random variables $Y = M(x)$ and $Y' = M(x')$ have the property that for every set $T \subseteq \mathcal{Y}$,

$$\Pr[Y \in T] \leq e^\varepsilon \cdot \Pr[Y' \in T].$$

Equivalently, we have $D_\infty(Y||Y') \leq \varepsilon$, where

$$D_\infty(Y||Y') = \sup_{T \subseteq \mathcal{Y}: \Pr[Y' \in T] > 0} \ln \left(\frac{\Pr[Y \in T]}{\Pr[Y' \in T]} \right).$$

Yet another equivalent formulation (known as the “group privacy” property of pure differential privacy) is that for every pair of (not necessarily adjacent) datasets $x, x' \in \text{MultiSets}(\mathcal{X})$,

$$D_\infty(Y||Y') \leq \varepsilon \cdot d_{\text{Sym}}(x, x') \tag{1}$$

An example of a measurement operator is a noisy sum. Suppose $\mathcal{X} = [L, U]$ is an interval of real numbers. Then the following is an ε -DP measurement $\text{NoisySum}_{L,U,\varepsilon} : \text{MultiSets}(\mathcal{X}) \rightsquigarrow \mathbb{R}$:

$$\text{NoisySum}_{L,U,\varepsilon}(x) = \sum_{z \in x} z + \text{Lap}(S/\varepsilon), \text{ where } S = \max\{|L|, |U|\}$$

where $\text{Lap}(S/\varepsilon)$ represents a draw from the Laplace distribution with scale S/ε and the summation $\sum_{z \in x}$ is a summation with multiplicity.²

Here we see that a measurement operator M and its privacy properties are specified by three attributes:

1. The data domain \mathcal{X} , whereby the measurement operator expects inputs from $\text{MultiSets}(\mathcal{X})$.
2. The privacy loss parameter ε .
3. The randomized function from $\text{MultiSets}(\mathcal{X}) \rightsquigarrow \mathcal{Y}$.

¹A multiset $x \in \text{MultiSets}(\mathcal{X})$ can be specified by its *histogram* $h_x : \mathcal{X} \rightarrow \mathbb{N}$, where $h_x(z)$ is the number of occurrences of z in \mathcal{X} . Then $d_{\text{Sym}}(x, x')$ equals the ℓ_1 distance between h_x and $h_{x'}$, i.e. $\sum_z |h_x(z) - h_{x'}(z)|$. *[[SV: new footnotes]]*

²That is, $\sum_{z \in x} z = \sum_{z \in [L,U]: h_z > 0} h_x(z) \cdot z$, where h_x is the histogram of x .

The specification of the input data domain \mathcal{X} is important because the privacy properties of M rely on the assumption that the inputs x come from $\text{MultiSets}(\mathcal{X})$. As above it is often useful to make use of rich data types, like the interval $[L, U]$. (When using the library, one should of course minimize the assumptions about the data that are used for privacy. See Section 8 for discussion on how this can be done.) We don't include the output data domain \mathcal{Y} as an attribute, because without loss of generality we can think of \mathcal{Y} as the countably infinite set of all possible objects representable in our programming language, so it need not be an explicitly specified attribute of M . *[[SV: expanded this paragraph]]*

In code, our noisy sum measurement operator could be implemented as follows:

```

1 class Measurement:
2     input_domain
3     privacy_loss
4     function
5
6 def MakeNoisySum(L: float, U: float, epsilon: float):
7     if L > U or epsilon < 0: raise Exception('Invalid parameters')
8     input_domain = bounded_float(L,U)
9     privacy_loss = epsilon
10    def function(data):
11        data_sum = sum(data)
12        s = max(abs(L), abs(U))
13        z = Laplace(s/epsilon,0)
14        return data_sum + z
15    return Measurement(input_domain,privacy_loss,function)
16
17 # Example
18 l, u, eps = # ... assign to some values
19 NoisySum = MakeNoisySum(l, u, eps)

```

Figure 1: Measurement Example: NoisySum. We define `Measurement` as a class with three attributes and an implicit constructor, but we could have defined it just as some kind of *struct* with three attributes. We use `class` as one such example, but this should not be read as advocating object-oriented languages (discussed further in Section 9). `NoisySum` is a `Measurement` with all parameters fixed (the values are assigned on 18). The `MakeNoisySum` function acts like a constructor that allows a programmer to specify values for the parameters, discussed further in Section 3. Our code examples make some assumptions: we assume the availability of complex data types such as `bounded_float(L,U)`, and of operations over multisets (implicitly defined), like the `sum` operator shown. We also pretend that floating-point arithmetic is a faithful implementation of real-number arithmetic, which it is certainly not and indeed this is known to lead to failures of differential privacy [9]. This is only for pedagogical purposes, to convey the ideas of the programming framework with familiar examples, but in Section 8 we advocate that the OpenDP library entirely avoid floating-point arithmetic.

2.2 Transformations

In PINQ, a *transformation* is a (deterministic) mapping from datasets to datasets. That is, $T : \text{MultiSets}(\mathcal{X}) \rightarrow \text{MultiSets}(\mathcal{Y})$. For $c \in \mathbb{R}$, we say that T is *c-stable* if for all $x, x' \in \text{MultiSets}(\mathcal{X})$,

$$d_{\text{Sym}}(T(x), T(x')) \leq c \cdot d_{\text{Sym}}(x, x'). \quad (2)$$

An example of a stable transformation is *clamping*. Suppose $\mathcal{X} = \mathbb{R}$ and $\mathcal{Y} = [L, U]$ and define $\text{clamp}_{L,U} :$

$\mathcal{X} \rightarrow \mathcal{Y}$ by

$$\text{clamp}_{L,U}(z) = \begin{cases} U & \text{if } z > U \\ z & \text{if } z \in [L, U] \\ L & \text{if } z < L. \end{cases}$$

Then we can define a 1-stable transformation $T : \text{MultiSets}(\mathcal{X}) \rightarrow \text{MultiSets}(\mathcal{Y})$ by:

$$T(x) = \{\text{clamp}_{L,U}(z) : z \in x\} \text{ (with multiplicity).}$$

Here we see that a transformation T and its stability properties are specified by four attributes:

1. The data domain \mathcal{X} for the input datasets.
2. The data domain \mathcal{Y} for the output datasets.
3. The stability parameter c .
4. The function from $\text{MultiSets}(\mathcal{X}) \rightarrow \text{MultiSets}(\mathcal{Y})$.

Here we need to specify the output domain \mathcal{Y} because a transformation should promise that if the input dataset x is in $\text{MultiSets}(\mathcal{X})$, then the output dataset $T(x)$ is in $\text{MultiSets}(\mathcal{Y})$. This is crucial for *chaining* as described later.

In code, the clamping transformation could be implemented as follows:

```
1 class Transformation:
2     input_domain
3     output_domain
4     stability
5     function
6
7 def MakeClamp(L: float, U: float):
8     if L > U: raise Exception('Invalid parameters')
9     input_domain = float
10    output_domain = bounded_float(L,U)
11    stability = 1
12    def function(data):
13        def clamp(x): return max(min(x, U), L)
14        return map(clamp, data)
15    return Transformation(input_domain,output_domain,stability,function)
16
17 # Example
18 Clamping = MakeClamp(1, u) # where l, u are constants defined previously
```

Figure 2: Transformation Example: Clamping. We use the same style as in Figure 1 for measurements. We define `Transformation` as a class with four attributes. `Clamping` is a concrete transformation with fixed parameters, constructed via `MakeClamp`. We use here a function `map` to apply the `clamp` function to every element of the multiset.

2.3 Chaining

From transformations and measurements defined as above, we can build up more complex transformations and measurements through various operators that combine them. One way of combining measurements and transformations is through *chaining*, which is simply function composition (but in the differential privacy literature the term “composition” has a different meaning, described below). There are two forms of chaining:

TT Chaining: If $S : \text{MultiSets}(\mathcal{X}) \rightarrow \text{MultiSets}(\mathcal{Y})$ is a c -stable transformation, and $T : \text{MultiSets}(\mathcal{Y}) \rightarrow \text{MultiSets}(\mathcal{Z})$ is a d -stable transformation, then $T \circ S : \text{MultiSets}(\mathcal{X}) \rightarrow \text{MultiSets}(\mathcal{Z})$ is a cd -stable transformation. Notice that the S and T need to be compatible in the sense that the data domain of the output dataset of S must match the data domain of the input datasets of T .

MT Chaining: Similarly, if $T : \text{MultiSets}(\mathcal{X}) \rightarrow \text{MultiSets}(\mathcal{Y})$ is a c -stable transformation, and $M : \text{MultiSets}(\mathcal{Y}) \rightsquigarrow \mathcal{Z}$ is an ε -DP measurement, then $M \circ T : \text{MultiSets}(\mathcal{X}) \rightsquigarrow \mathcal{Z}$ is a $c\varepsilon$ -DP measurement. Again, this requires compatibility of the intermediate data domain between T and M .

To illustrate the latter (chaining of a transformation and measurement), consider $T = \text{clamp}_{L,U}$ and $M = \text{NoisySum}_{L,U,\varepsilon}$. Since $\text{clamp}_{L,U} : \text{MultiSets}(\mathbb{R}) \rightarrow \text{MultiSets}([L,U])$ is a 1-stable transformation and $\text{NoisySum}_{L,U,\varepsilon} : \text{MultiSets}([L,U]) \rightsquigarrow \mathbb{R}$ is ε -DP, their composition

$$\text{NoisyClampedSum}_{L,U,\varepsilon} = \text{NoisySum}_{L,U,\varepsilon} \circ \text{clamp}_{L,U} : \text{MultiSets}(\mathbb{R}) \rightarrow \mathbb{R}$$

is ε -DP.

In code, we can implement chaining as an operator that takes two transformation objects that are compatible and produces a new transformation object, or similarly for a transformation object and a measurement object. The example above of chaining $\text{clamp}_{L,U}$ and $\text{NoisySum}_{L,U,\varepsilon}$ can be illustrated in code as follows:

[[SV: removed certificate checks from code but added type annotations to inputs]]

```

1 def ChainingTT(trans_2: Transformation, trans_1: Transformation):
2     # makes new transformation: trans_2(trans_1(.))
3     if (trans_1.output_domain!=trans_2.input_domain): raise Exception('Domain mismatch')
4     input_domain = trans_1.input_domain
5     output_domain = trans_2.output_domain
6     stability = trans_1.stability*trans_2.stability
7     def function(data): return trans_2.function(trans_1.function(data))
8     return Transformation(input_domain,output_domain,stability,function)
9
10 def ChainingMT(meas: Measurement, trans: Transformation):
11     # makes new measurement: meas(trans(.))
12     if (trans.output_domain!=meas.input_domain): raise Exception ('Domain mismatch')
13     input_domain = trans.input_domain
14     privacy_loss = trans.stability*meas.privacy_loss
15     def function(data): return meas.function(trans.function(data))
16     return Measurement(input_domain,privacy_loss,function)
17
18 # Example
19 NoisyClampedSum=ChainingMT(NoisySum, Clamping)

```

Figure 3: Chaining. We define two chaining operations, chaining of two transformations, and chaining of a transformation and a measurement. These functions check that the input transformations/measurements are certified to be correct (see Section 3) and that the input and output domains are compatible and they raise an exception otherwise.

2.4 Composition

Informally, the *Basic Composition Theorem* of differential privacy says that if we execute an ε_1 -DP mechanism on a dataset followed by an ε_2 -DP mechanism on the same dataset, the result is $(\varepsilon_1 + \varepsilon_2)$ -DP. This sort of a composition principle has two main applications:

1. It allows for tracking cumulative privacy loss when analysts make many differentially private queries on a dataset. This is what gives rise to the notion of a “privacy loss budget” in differential privacy, where

we can prevent exceeding a desired total privacy loss bound ε by tracking the accumulated privacy loss and refusing to answer any queries that would result in exceeding the overall budget of ε .

2. It allows for building more complex DP mechanisms from multiple executions of simpler mechanisms, deriving an upper bound on the privacy loss of the “outer mechanism” as the sum of the privacy losses of the the “inner mechanisms”.

In the simplest, “nonadaptive” form of basic composition, we have an ε_1 -DP mechanism $M_1 : \text{MultiSets}(\mathcal{X}) \rightsquigarrow \mathcal{Y}$ and an ε_2 -DP mechanism $M_2 : \text{MultiSets}(\mathcal{X}) \rightsquigarrow \mathcal{Z}$, and obtain an $(\varepsilon_1 + \varepsilon_2)$ -DP mechanism $M : \text{MultiSets}(\mathcal{X}) \rightsquigarrow \mathcal{Y} \times \mathcal{Z}$ by $M(x) = (M_1(x), M_2(x))$, where M_1 and M_2 use independent randomness. For example, if we wanted to compute a differentially private clamped mean, we could take $M_1 = \text{NoisyClampedSum}_{L,U,\varepsilon}$ and $M_2 = \text{NoisyClampedSum}_{1,1,\varepsilon}$ to obtain an (2ε) -DP NoisyPair that provides us with both a noisy clamped sum (from M_1) and a noisy estimate of the size of the dataset (from M_2).

In code, we can implement such basic composition as an operator on measurement objects as follows:

```

1 def Compose(meas_1: Measurement,meas_2: Measurement)
2   if (meas_1.input_domain!=meas_2.input_domain): raise Exception('Domain mismatch')
3   input_domain = meas_1.input_domain
4   privacy_loss = meas_1.privacy_loss+meas_2.privacy_loss
5   def function(data): return (meas_1.function(data),meas_2.function(data))
6   return Measurement(input_domain,privacy_loss,function)
7
8 # Example
9
10 # We first define a new measurement: a noisy count that is built from the
11 # clamping, noisy sum, and chaining operations previously introduced
12 ClampToOne = MakeClamp(1, 1)
13 NoisySumOfOnes = MakeNoisySum(1, 1, eps)
14 NoisyCount=ChainingMT(NoisySumOfOnes, ClampToOne)
15
16
17 # With this we can now define NoisyPair
18 NoisyPair=Compose(NoisyClampedSum, NoisyCount)

```

Figure 4: Compose. We check that the input domains of the two measurements we want to compose are compatible, otherwise we throw an exception. In the example, we define a NoisyCount measurement by reusing previously defined operations and then compose NoisyClampedSum and NoisyCount.

However, it is often useful to use more sophisticated, “adaptive” forms of composition, where the choice of the mechanism M_2 depends on the result of $M_1(x)$, and we also allow even more mechanisms M_3, M_4, \dots to be chosen adaptively. This kind of flexibility is clearly important when allowing an analyst to interactively query a dataset protected by differential privacy. To support this, PINQ and other DP systems often manage the privacy budget and composition at a higher layer that sits above the basic transformations and measurements. (Indeed, PINQ also handles chaining at that higher level, rather than as an operator that produces new transformations and measurements.)

2.5 Post-processing

An important property of differential privacy is that it is closed under post-processing. That is, if $M : \text{MultiSets}(\mathcal{X}) \rightsquigarrow \mathcal{Y}$ is ε -DP and $f : \mathcal{Y} \rightarrow \mathcal{Z}$ is any function, then $f \circ M : \text{MultiSets}(\mathcal{X}) \rightsquigarrow \mathcal{Z}$ is ε -DP.

Now we can actually produce a differentially private mean by composing the NoisyPair mechanism above with a division operation $f(a,b) = a/b$. In code:

```

1 def Postprocess(meas: Measurement, funct):
2     input_domain = meas.input_domain
3     privacy_loss = meas.privacy_loss
4     def function(data): return funct(meas.function(data))
5     return Measurement(input_domain, privacy_loss, function)
6
7 # Example
8 # We can define an auxiliary divide function
9 def divide(x,y): return x/y
10 # We can now redefine NoisyMean using Postprocess.
11 NoisyMean=Postprocess(NoisyPair, divide)

```

Figure 5: Postprocessing. Notice that we don't require the `input_domain` of the function to match the `output_domain` of the measurement. In fact, measurements don't have an `output_domain` attribute. From the privacy perspective this is not required. This can create a mismatch which can generate an exception, but this exception would not create any privacy issue. Using a statically typed language can help to avoid such situations.

3 Verifying Privacy Properties

[[SV: changed title]]

Our goal is to ensure that the only measurements and transformations that can be constructed by the OpenDP Library have mathematically proven privacy properties, based on either:

- A custom proof, which is provided by the contributor and is verified by a human (on the OpenDP editorial board) or by a computer (for components that are amenable to formal verification techniques), or
- An automatically derived proof, if the new component is obtained by combining components that already exist in the library (using combination primitives that exist in the library).

Looking at the examples in Section 2, we see that it is rare that a measurement or transformation would be provided as a single, stand-alone object in the library. Rather they are given as parameterized families of objects, like `clampL,U` is parameterized by L and U and `NoisySumL,U,ε` is parameterized by L , U , and ϵ . So we really need verification procedures for *measurement families* and *transformation families*. In our code examples, we implemented measurement (respectively, transformation) families as constructors that take the parameters and output a measurement (respectively a transformation), or an exception if the parameters are invalid. Thus we propose:

Code should only be accepted to the library if there is a proof that (1) it can only ever construct *valid* measurements or transformations, where valid means that the measurement (resp. transformation) respects the privacy-loss bound (resp., stability bound) promised in its attributes, (2) it does not modify any measurements or transformations that have been constructed, and (3) it does not modify code in the library. *[[MH: we need to define more precisely, particularly around acceptable modes of failure. can a function raise an exception? how do we define what types of exceptions are permissible and what ones are not? is it implicit in our definition that a family can never take the data as input and so exceptions are fine because they cannot depend on the private data?]]*

The proof can assume (by induction) that all measurements and transformations given as inputs or constructed by existing code in the library are valid. In particular, if the new code does not *directly* construct any measurements or transformations on its own (but only using existing code

to do so), does not modify any measurements or transformations that have been constructed, and does not modify code in the library, then it should be possible to verify its validity automatically.

[[SV: items (2) and (3) above are new. in the future, we should probably connect those requirements to the language support section?]]

[[SV: Michael’s counterexample:]]

```

1 # this example only appears in draft=1 mode
2 def MakeNoisyClampedSum(L: float, U: float, epsilon: float):
3     Clamp = MakeClamp(L,U)
4     NoisySum = MakeNoisySum(L,U,eps)
5     NoisySum.epsilon = 0 # --- this is bad
6     return ChainingMT(NoisySum,Clamp) # does not check epsilon

```

[[SV: modified quote per discussion with Marco. is ‘valid’ the best term? or ‘safe’?]] Evaluating the code examples from Section 2 with respect to this principle, note that the various checks they do on their input parameters (e.g. checking that $L \leq U$ and $\epsilon \geq 0$) are crucial for ensuring that is impossible for them to generate invalid measurements or transformations. *[[SV: softened claim that they actually meet the principle, since, for example, we haven’t addressed the floating-point issue yet.]]*

Although the principle above ensures that all measurements and transformations constructed by the library are valid and no further verification is necessary, it might be of interest to decorate measurements and transformations with the “provenance” of their proof of correctness (showing how the proofs of all the library methods used to construct them are stitched together). *[[SV: sentence based on Michael’s comment and code snippets in families.tex. in future, we can discuss more whether this is an important feature, and if so, add the code examples back]]*

[[SV: removed the remaining text in the section — seems unnecessary given the above.]]

4 Varying Types and Distance Measures

4.1 Private Data Types

In the basic PINQ-like framework described above, all sensitive data is represented as a multiset, and the granularity of privacy is captured by the symmetric difference metric d_{Sym} . As realized in Fuzz [11], it is very useful to allow for other sensitive data types and metrics. There are two benefits to such a generalization:

1. Not all sensitive data comes in the form a multiset of records, where the desired granularity of privacy is adding or removing one record. For example, for network data, the sensitive dataset is often a graph (sometimes augmented with labels) and the granularity of privacy may refer to modifications at the level of a node or an edge.
2. Even if the original dataset comes in the form of multiset of records, a differentially private algorithm may produce intermediate data representations (e.g. a histogram). Being able to reason about distances in such intermediate representations allows decomposing differentially private algorithms into smaller modular and reusable components.

Now a transformation $T : \mathcal{X} \rightarrow \mathcal{Y}$ can have arbitrary input and output types \mathcal{X} and \mathcal{Y} , which need not be multisets. For a given metrics $d_{\mathcal{X}}$ and $d_{\mathcal{Y}}$ on \mathcal{X} and \mathcal{Y} , we say that T is a *c-stable transformation from $d_{\mathcal{X}}$ to $d_{\mathcal{Y}}$* if all $x, x' \in \mathcal{X}$, we have $d_{\mathcal{Y}}(T(x), T(x')) \leq c \cdot d_{\mathcal{X}}(x, x')$, generalizing Inequality (2). Similarly, we say that a measurement $M : \mathcal{X} \rightsquigarrow \mathcal{Y}$ is *ϵ -DP with respect to $d_{\mathcal{X}}$* if for all $x, x' \in \mathcal{X}$, $D_{\infty}(M(x)||M(x')) \leq \epsilon \cdot d_{\mathcal{X}}(x, x')$, generalizing Inequality (1).

Now, a transformation should also have the metrics $d_{\mathcal{X}}$ and $d_{\mathcal{Y}}$ associated with it, and a measurement the metric $d_{\mathcal{X}}$. Then chaining works as before, except we must also check compatibility of the intermediate metrics.

As an illustration, we can decompose our $\text{NoisySum}_{L,U,\varepsilon}$ as a chaining of two components $\text{BoundedSum}_{L,U} : \text{MultiSets}([L, U]) \rightarrow \mathbb{R}$ and $\text{BaseLap}_{\max\{|L|, |U|\}/\varepsilon} : \mathbb{R} \rightsquigarrow \mathbb{R}$, where

$$\text{BoundedSum}_{L,U}(x) = \sum_{z \in x} z \text{ (with multiplicity),}$$

and

$$\text{BaseLap}_{\sigma}(y) = y + \text{Lap}(\sigma).$$

The ε -DP property of NoisySum follows from the facts that BoundedSum is $\max\{|L|, |U|\}$ -stable from d_{Sym} to $d_{\mathbb{R}}$, where $d_{\mathbb{R}}(a, b) = |a - b|$, and BaseLap_{σ} is $(1/\sigma)$ -DP with respect to $d_{\mathbb{R}}$.

In code:

```

1 class Measurement:
2     input_metric # --- new ---
3     input_domain
4     privacy_loss
5     function
6
7 class Transformation:
8     input_metric # --- new ---
9     input_domain
10    output_metric # --- new ---
11    output_domain
12    stability
13    function
14
15 def ChainingMT(meas: Measurement, trans: Transformation):
16     if (trans.output_domain!=meas.input_domain
17         or trans.output_metric!=meas.input_metric): raise Exception('Domain/metric mismatch')
18     input_metric = trans.input_metric
19     input_domain = trans.input_domain
20     privacy_loss = trans.stability*meas.privacy_loss
21     def function(data): return meas.function(trans.function(data))
22     return Measurement(input_metric,input_domain,privacy_loss,function)
23
24 def MakeBaseLap(sigma: float):
25     if sigma < 0: raise Exception('Invalid parameter')
26     input_metric = dist_real
27     input_domain = float
28     privacy_loss = 1/sigma
29     def function (data): return data + Laplace(sigma,0)
30     return Measurement(input_metric,input_domain,privacy_loss,function)
31
32 def MakeBoundedSum(L: float, U: float):
33     if L > U: raise Exception('Invalid parameters')
34     input_metric = dist_sym
35     input_domain = multiset(bounded_float(L,U))
36     output_metric = dist_real
37     output_domain = float
38     stability = max(abs(L), abs(U))
39     def function (data): return sum(data)
40     return Transformation(input_metric,input_domain,output_metric,output_domain,
41                           stability,function)
42
43 # Example
44 BaseLaplace = MakeBaseLap(sig) # sig is some constant
45 BoundedSum = MakeBoundedSum(l, u) # l,u defined previously
46 NoisySum=ChainingMT(BaseLaplace, BoundedSum)
47 print(NoisySum.privacy_loss) # prints max(abs(l),abs(u))/sig

```

Figure 6: Private Data Types. We re-define `Measurement` and `Transformation` to include metrics. We also re-define the chaining operations to also check metric compatibility (`ChainingTT` omitted but similar). Notice how the `input_domain` of the `BoundedSum` transformation is now more explicit about the type, which declares that we expect a multiset of bounded floats.

4.2 Privacy Relations

In the definitions of “pure” differential privacy above, we measure the distance between the output distributions $M(x)$ and $M(x')$ using “max-divergence” $D_\infty(M(x)||M(x'))$. In the differential privacy literature, a number of other measures have been proposed, motivated by the greater utility they afford and/or their superior composition properties.

For example, with *approximate differential privacy*, we say that M is (ε, δ) -DP if for all $x \sim x'$, we have $D_\infty^\delta(M(x)||M(x')) \leq \varepsilon$, where D_∞^δ is the *smoothed max-divergence*. This means that for all sets T , we have:

$$\Pr[M(x) \in T] \leq e^\varepsilon \cdot \Pr[M(x') \in T] + \delta.$$

Note that here privacy is not measured by a single number but a pair of numbers (ε, δ) . Moreover, there need not be a single optimal choice of (ε, δ) for a given mechanism, but rather a Pareto curve, where for every $\varepsilon > 0$, we can try to identify the minimum δ for which the mechanism is (ε, δ) -DP. A similar phenomenon occurs with other notions such as concentrated differential privacy and Rényi differential privacy.

Thus, it is no longer sufficient to express privacy as a linear relationship between input distances and output distances, like we did before when we required $D_\infty(M(x)||M(x')) \leq \varepsilon \cdot d_{\mathcal{X}}(x, x')$. Instead, for a measurement $M : \mathcal{X} \rightsquigarrow \mathcal{Y}$ a metric $d_{\mathcal{X}}$ on \mathcal{X} and a similarity measure D on probability distributions, we assert the privacy properties of M by a *relation* $R(d_{\text{in}}, d_{\text{out}})$ that takes an input distance d_{in} and an output distance d_{out} and certifies that “for all $x, x' \in \mathcal{X}$, if x is d_{in} -close to x' under $d_{\mathcal{X}}$, then $M(x)$ is d_{out} -close to $M(x')$ with respect to D .”

Let us illustrate with the example of the Gaussian mechanism $\text{BaseGauss}_\sigma : \mathbb{R} \rightsquigarrow \mathbb{R}$, defined as

$$\text{BaseGauss}_\sigma(y) = y + N(0, \sigma).$$

It is known that if $|y - y'| \leq d$, then for every $\delta > 0$,

$$D_\infty^\delta(\text{BaseGauss}_\sigma(y)||\text{BaseGauss}_\sigma(y')) \leq \frac{d}{\sigma} \cdot \sqrt{2 \ln \left(\frac{1.25}{\delta} \right)},$$

provided that the right-hand side is at most 1.

$$R(d_{\text{in}}, (\varepsilon, \delta)) = \begin{cases} \text{true} & \text{if } \min\{\varepsilon, 1\} \geq (d_{\text{in}}/\sigma)\sqrt{2 \ln(1.25/\delta)} \\ \text{false} & \text{otherwise} \end{cases}$$

To capture this in code, we modify our Measurement class from before in the following ways:

1. In addition to an input metric $d_{\mathcal{X}}$, we also have an attribute that corresponds to the privacy measure D .
2. The pure DP constant is replaced with a privacy relation.

[[SV: replaced ADP with approxDP]]

```

1 class Measurement:
2     input_metric
3     input_domain
4     output_measure          # --- new ---
5     privacy_relation        # --- new (replaces privacy_loss)
6     function
7
8 def MakeBaseGauss(sigma: float):
9     if sigma < 0: raise Exception
10    input_metric = dist_real
11    input_domain = float
12    output_measure = approxDP
13    def privacy_relation (d_in: float,d_out: float*float):
14        return ((d_in/sigma)*sqrt(2*ln(1.25/d_out[2]))) <= min(d_out[1],1)
15    def function (data):
16        z = Gauss(sigma,0)
17        return data + z
18    return Measurement(input_metric,input_domain,output_measure,privacy_relation,function)
19
20 # Example
21 BaseGauss = MakeBaseGauss(sig)

```

Figure 7: Privacy relations. We re-define `Measurement` by changing the `privacy_loss` attribute to a `privacy_relation` attribute. In `BaseGauss`, the `privacy_relation` function returns true if the distance in input and the distance in output satisfy the constraints imposed by the correctness of the Gaussian mechanism for approximate differential privacy – we use a data type `approxDP` as the associated `output_measure`. The relation is intentionally given as *callable* code; as other mechanisms that use the measurement may need to evaluate the relation during execution. The notation `float*float` means the expected type is pair of floats.

[[SV: added type annotations on inputs to relation - does that look right?]] *[[SV: added a comment in the figure about the relation being callable code]]* *[[SV: renamed `relation` to `privacy_relation`]]*

[[SV: add a comment that the privacy relation is intentionally a callable function! in the subsequent subsections we will see that various composition primitives (e.g. chaining) will need to query the privacy relation, possibly at multiple settings of the distance bounds]]

4.3 Stability Relations

Having generalized to arbitrary, multi-parameter privacy measures, it is natural to allow the same flexibility for our distance measures on the sensitive data types. An example that has arisen in the differential privacy literature is that of *joins* [10, 5, 7, 8]. Let $x \in \text{MultiSets}(\mathcal{K} \times \mathcal{A})$ and $y \in \text{MultiSets}(\mathcal{K} \times \mathcal{B})$ be two datasets over records that share a common key from domain \mathcal{K} . The *join* of x and y over \mathcal{K} is a multiset $x \bowtie_{\mathcal{K}} y \in \text{MultiSets}(\mathcal{K} \times \mathcal{A} \times \mathcal{B})$ defined as

$$x \bowtie_{\mathcal{K}} y = \{(k, a, b) : (k, a) \in x, (k, b) \in y\} \text{ (with multiplicity)}$$

A join has unbounded stability with respect to both x and y , so a solution is to first apply a *truncation* operation to each dataset. Define $\text{Trunc}_{\mathcal{K},\ell}(x)$ to be a dataset obtained by discarding records from x so that for every $k \in \mathcal{K}$, the number of records of the form (k, \cdot) in $\text{Trunc}_{\mathcal{K},\ell}(x)$ (counting multiplicity) is at most ℓ . (Specifically, $\text{Trunc}_{\mathcal{K},\ell}(x)$ sorts the elements $(k, a_1), (k, a_2), \dots$ according to a fixed ordering on \mathcal{A} and discards all but the first ℓ records.) Similarly define $\text{Trunc}_{\mathcal{K},m}(y)$. Then we define

$$\text{BoundedJoin}_{\ell,m}(x, y) = \text{Trunc}_{\mathcal{K},\ell}(x) \bowtie_{\mathcal{K}} \text{Trunc}_{\mathcal{K},m}(y).$$

Inputs of BoundedJoin are *pairs* of multisets, so it is useful to measure distances by a pair of numbers. Specifically, we say that (x, y) is (c, d) -close to (x', y') under $d_{\text{Sym}} \times d_{\text{Sym}}$ if $|x \Delta x'| \leq c$ and $|y \Delta y'| \leq d$. This notion arises naturally in DP applications when x and y may each be derived from some underlying sensitive dataset z using transformations of stability c and d , respectively. Then we can express the stability property of BoundedJoin as follows: if (x, y) is (c, d) -close to (x', y') under $d_{\text{Sym}} \times d_{\text{Sym}}$, then $\text{BoundedJoin}_{\ell, m}(x, y)$ is $2(c \cdot m + \ell \cdot d + 2c \cdot d)$ -close to $\text{BoundedJoin}_{\ell, m}(x', y')$ under d_{Sym} [7]. This can be expressed in terms of a stability relation as follows:

$$R((c, d), e) = \begin{cases} \text{true} & \text{if } e \geq 2(c \cdot m + \ell \cdot d + 2c \cdot d) \\ \text{false} & \text{otherwise} \end{cases}$$

[[SV: for future: is bound tight with factor of 2, or can anything be saved?]]

[[SV: could also have “pairing operator” that takes transformations $S : \mathcal{X} \rightarrow \mathcal{Y}$ and $T : \mathcal{X} \rightarrow \mathcal{Z}$ and constructs a new transformation $S \times T : \mathcal{X} \rightarrow \mathcal{Y} \times \mathcal{Z}$ whose stability is measured with respect to $d_{\mathcal{Y}} \times d_{\mathcal{Z}}$]]

4.4 Multiple Distance Measures at Once

Additionally, we need not have a fixed pair of input and output distance measures associated with a given transformation or measurement operator, as the same operator can have its stability measured according to multiple metrics and we can exploit relationships between those metrics to automatically translate between them. For example, it is known that any mechanism that is ρ -zCDP (where zCDP stands for “zero-concentrated differential privacy” [2]) is also $(\rho + 2\sqrt{\rho \ln(1/\delta)}, \delta)$ -DP for every $\delta > 0$. Similarly, any two vectors $u, v \in \mathbb{R}^d$ that are α -close in ℓ_1 norm are also $\sqrt{d} \cdot \alpha$ -close in ℓ_2 norm. To support this, we no longer have our input and output metrics as fixed attributes of our transformations and measurements, but we take them as input to the privacy or stability relations. That is, for a measurement operator $M : \mathcal{X} \rightarrow \mathcal{Y}$ we now have a privacy relation $R((d_{\mathcal{X}}, d_{\text{in}}), (D, d_{\text{out}}))$ such that whenever $R((d_{\mathcal{X}}, d_{\text{in}}), (D, d_{\text{out}}))$ is true, the following holds:

1. $d_{\mathcal{X}}$ is a (possibly multidimensional) distance measure on \mathcal{X} , and d_{in} is a valid distance bound under $d_{\mathcal{X}}$.
2. D is a (possibly multidimensional) distance measure on probability distributions on arbitrary sets, and d_{out} is a valid distance bound under D .
3. for all datasets $x, x' \in \mathcal{X}$, if x and x' are d_{in} -close under $d_{\mathcal{X}}$, then the probability distributions $M(x)$ and $M(x')$ are d_{out} -close under D .

For example, when $\mathcal{X} = \text{MultiSets}(\mathcal{Z})$, then $R((d_{\text{Sym}}, 1), (D_{\infty}, \varepsilon))$ being true asserts that M is ε -DP. For the $\text{BaseGauss}_{\sigma} : \mathbb{R} \rightarrow \mathbb{R}$ mechanism defined above, we can assert its privacy properties under zCDP and approximate DP with a relation defined as follows:

$$R((d_{\mathcal{X}}, d_{\text{in}}), (D, d_{\text{out}})) = \begin{cases} \text{true} & \text{if } d_{\mathcal{X}} = d_{\mathbb{R}}, D = \text{zCDP}, d_{\text{in}}, d_{\text{out}} \in \mathbb{R}, \text{ and } d_{\text{out}} \geq d_{\text{in}}^2 / 2\sigma^2 \\ \text{true} & \text{if } d_{\mathcal{X}} = d_{\mathbb{R}}, D = \text{approxDP}, d_{\text{in}} \in \mathbb{R}, d_{\text{out}} = (\varepsilon, \delta) \in \mathbb{R} \times \mathbb{R}, \\ & \text{and } \min\{\varepsilon, 1\} \geq (d_{\text{in}} / \sigma) \sqrt{2 \ln(1.25/\delta)} \\ \text{false} & \text{otherwise} \end{cases}$$

[[SV: comment on the fact that one input (D) determines type of next input (d_{out})... implement using structures or discrete polymorphism or finite data type?]] In code: *[[SV: added new code snippet. notice type checks in relation - syntax OK?]]*

```

1 def MakeBaseGauss(sigma: float):
2     # ... same as before except relation now takes distance measures as input ...
3     def privacy_relation ((m_in,d_in),(m_out,d_out)):
4         if m_in==dist_real and type(d_in)==float:
5             if m_out==approxDP and type(d_out)==float*float:
6                 return ((d_in/sigma)*sqrt(2*ln(1.25/d_out[1]))) <= min(d_out[0],1)
7             elif m_out==zCDP and type(d_out)==float:
8                 return (d_in^2/(2*sigma^2)) <= d_out
9             else return false
10    return Measurement(input_metric,input_domain,output_measure,relation,function)
11
12 # Example
13 BaseGauss = MakeBaseGauss(sig)

```

Figure 8: Privacy relations for multiple distance measures at once. We redefine `MakeBaseGauss` so that it produces a `privacy_relation` that takes distance measures as input in addition to the distances *[[SV: changed indexing of `d_out` to 0 and 1]]*

The `BaseGauss` example illustrates that the privacy relations are only meant to be *sound* but are not necessary *complete*. That is, if the privacy relation says true, it should be considered a certificate that the measurement or transformation has the claimed privacy or stability properties. However, if the privacy relation says false, it only means that no claim is made about the property. For example, using the aforementioned relationship between `zCDP` and approximate DP, it follows that the Gaussian mechanism is (ϵ, δ) -DP whenever

$$\epsilon \geq (d_{\text{in}}^2/2\sigma^2) + (d_{\text{in}}/\sigma)\sqrt{2\ln(1/\delta)},$$

which is satisfied by some values of $d_{\text{out}} = (\epsilon, \delta)$ that do not satisfy the relation above. This example also illustrates the eventual possibility of having some automated translation between different distance measures. That is, we can encode the relationships between different distance measures in the OpenDP library, and have it automatically check whether a privacy or stability property can be derived from other ones.

Working on other approaches:

```

1  # version 1
2  # each measurement has a single privacy relation with fixed in/out metrics
3  # operators that let you "replace" the relation (e.g., zCDP2ApproxDP) by creating a new measurement tha
4
5  class Measurement:
6      input_metric
7      input_domain
8      output_measure          # --- new ---
9      privacy_relation        # --- new (replaces privacy_loss)
10     function
11
12 def zCDP2ApproxDP(meas : Measurement):
13     if not(meas.output_measure == zCDP) RaiseException
14     new_output_measure == approxDP
15     def new_privacy_relation(d_in, d_out):
16         ... use meas.privacy_relation(...)
17     return Measurement(meas.input_metric,meas.input_domain,new_output_measure,new_privacy_relation,function)
18
19 # version 2
20 # each measurement has a set of multiple relations
21
22 class Measurement:
23     input_domain
24     privacy_relation_set
25     # collection of privacy_relations indexed on (input_metric,output_metric) pairs
26     # allow the same pair of metrics to appear multiple times
27     privacy_relation_eval(m_in,d_in,m_out,d_out):
28     # look up all relations in the set that have metrics (m_in,m_out)
29     # return true if any of them output true
30     function
31
32 def zCDP2ApproxDP(meas : Measurement):
33     # search for privacy relations in meas.pr_set that
34     # have output_measure == zCDP
35     # and for each, attach a new privacy relation
36     # with output_measure == approxDP
37
38
39     def new_privacy_relation(d_in, d_out):
40         ... use meas.privacy_relation(...)
41     return Measurement(meas.input_metric,meas.input_domain,new_output_measure,new_privacy_relation,function)
42
43 def MakeBaseGauss(sigma: float):
44     # ... same as before except relation now takes distance measures as input ...
45     def privacy_relation ((m_in,d_in),(m_out,d_out)):
46         if m_in==dist_real and type(d_in)==float:
47             if m_out==approxDP and type(d_out)==float*float:
48                 return ((d_in/sigma)*sqrt(2*ln(1.25/d_out[1]))) <= min(d_out[0],1)
49             elif m_out==zCDP and type(d_out)==float:
50                 return (d_in^2/(2*sigma^2)) <= d_out
51             else return false
52     return Measurement(input_metric,input_domain,output_measure,relation,function)
53
54 # Example
55 BaseGauss = MakeBaseGauss(sig)

```

Figure 9: Privacy relations for multiple distance measures at once: alternatives

```

1 def BasicCompositionApproxDP(ListMeas : Measurement):
2   ...
3   return ComposedMeasurement : Measurement
4
5 def BasicCompositionzCDP(ListMeas : Measurement):
6   ...
7   return ComposedMeasurement : Measurement
8
9 def SelectBetter(queryable_inner):
10  % think of queryable_inner as outer Approx DP composition
11
12
13
14
15
16
17 def SelectBetter(Meas1 : Measurement, Meas2 : Measurement, m_in,d_in,m_out,d_out):
18  Check that Meas1 and Meas2 both use m_in and m_out and both are "monotone" metrics
19  if Meas1.relation(d_in,d_out): return Meas1
20  else if Meas2.relation(d_in,d_out): return Meas2
21
22  def new_relation(d'_in,d'_out):
23    return (d'_in <= d_in and d'_out >= d_out)
24  def new_function(x):
25    if Meas1.relation(d_in,d_out) return Meas1.function(x)
26    else if Meas2.relation(d_in,d_out) return Meas2.function(x)
27    else return UselessQueryable

```

Figure 10: Privacy relations for multiple distance measures at once: alternatives

4.5 Chaining

Let us consider how chaining behaves with these stability and privacy relations. Specifically, how do we derive the relation for a chained transformation or measurement from the relations for its component parts? Let $T : \mathcal{X} \rightarrow \mathcal{Y}$ be a transformation with stability relation R_T and $M : \mathcal{Y} \rightsquigarrow \mathcal{Z}$ be a measurement with privacy relation R_M . Intuitively, the privacy of $M \circ T$ should be derived from the fact that (a) close inputs in \mathcal{X} map to close elements of \mathcal{Y} under T , and (b) close elements of \mathcal{Y} map to close probability distributions under M . Note that we are free to choose any notion of closeness on \mathcal{Y} for which these two statements hold. So the following would be a valid derivation of a privacy relation for the chained measurement $M \circ T : \mathcal{X} \rightsquigarrow \mathcal{Z}$:

$$R_{M \circ T}((d_{\mathcal{X}}, d_{\text{in}}), (D, d_{\text{out}})) = \begin{cases} \text{true} & \text{if } \exists (d_{\mathcal{Y}}, d_{\text{mid}}) R_T((d_{\mathcal{X}}, d_{\text{in}}), (d_{\mathcal{Y}}, d_{\text{mid}})) \wedge R_M((d_{\mathcal{Y}}, d_{\text{mid}}), (D, d_{\text{out}})) \\ \text{false} & \text{otherwise} \end{cases}$$

But in general it is not feasible to implement such a relation, as there may be many choices for the pair $(d_{\mathcal{Y}}, d_{\text{mid}})$ — indeed, infinitely many if d_{mid} is vector of real numbers. So, when chaining, we must provide a hint function that specifies the intermediate notion of closeness to be used: $(d_{\mathcal{Y}}, d_{\text{mid}}) = \text{hint}((d_{\mathcal{X}}, d_{\text{in}}), (D, d_{\text{out}}))$, and we instead use the following relation:

$$R_{M \circ T}((d_{\mathcal{X}}, d_{\text{in}}), (D, d_{\text{out}})) = \begin{cases} \text{true} & \text{if } R_T((d_{\mathcal{X}}, d_{\text{in}}), \text{hint}((d_{\mathcal{X}}, d_{\text{in}}), (D, d_{\text{out}}))) \wedge R_M(\text{hint}((d_{\mathcal{X}}, d_{\text{in}}), (D, d_{\text{out}})), (D, d_{\text{out}})) \\ \text{false} & \text{otherwise} \end{cases}$$

Note that the soundness of this relation does not depend on the choice of the hint function. A poorly chosen hint function may result in a relation that outputs **false** more often than necessary, but will not make the

relation output `true` when the chained measurement does not have the specified privacy property (assuming the relations R_M and R_T are sound).

Let's illustrate this by chaining `BoundedSum $_{L,U}$` and `BaseGauss $_{\sigma}$` to obtain `GaussSum $_{L,U,\sigma}$` . We know that `BoundedSum $_{L,U}$` : `MultiSets([L, U])` \rightarrow \mathbb{R} is a $\max\{|L|, |U|\}$ -stable transformation from d_{Sym} to $d_{\mathbb{R}}$, so here a natural hint function to use would set:

$$\text{hint}((d_{\text{Sym}}, d_{\text{in}}), (D, d_{\text{out}})) = (d_{\mathbb{R}}, \max\{|L|, |U|\} \cdot d_{\text{in}}).$$

In code:

```

1 def ChainingMT(trans: Transformation, meas: Measurement, hint):
2     if trans.output_domain != meas.input_domain: raise Exception('domain mismatch')
3     input_metric = trans.input_metric
4     input_domain = trans.input_domain
5     def function(data): return meas.function(trans.function(data))
6     output_measure = meas.output_measure # --- new ---
7     def privacy_relation((m_in, d_in), (m_out, d_out)): # --- new ---
8         m_mid, d_mid = hint((m_in, d_in), (m_out, d_out))
9         return (trans.stability_relation((m_in, d_in), (m_mid, d_mid))
10                and meas.privacy_relation((m_mid, d_mid), (m_out, d_out)))
11     return Measurement(input_metric, input_domain, output_measure, privacy_relation, function)
12
13 def MakeBoundedSum(L: float, U: float):
14     # ... same as before except relation replaces stability constant ...
15     def stability_relation((m_in, d_in), (m_out, d_out)) =
16         if m_in == dist_sym and type(d_in) == float and m_out == dist_real and type(d_out) == float:
17             return (max(abs(L), abs(U)) <= d_out)
18     return Transformation(input_metric, input_domain, output_metric, output_domain,
19                           stability_relation, function)
20
21 # Example
22 def hintGaussSum((m_in, d_in), (m_out, d_out)):
23     if (m_in == dist_sym and c == approxDP): return (dist_real, max(abs(L), abs(U)) * b)
24 BoundedSum = MakeBoundedSum(l, u)
25 BaseGauss = MakeBaseGauss(sig)
26 GaussSum = ChainingMT(BoundedSum, BaseGauss, hintGaussSum)

```

Figure 11: Chaining revisited. `ChainingMT` has been revised to define a new `privacy_relation` function that simply computes the conjunction of the relation checks for the operations being chained, with a `hint` on the intermediate distance and metric supplied by the caller. As an example, `GaussSum` is made through chaining. *[[MH: chaining arguments should be swapped to be consistent with rest of paper]]*

[[SV: removed basegauss since now it appears in previous code snippet, corrected stability relation in MakeBoundedSum to check the metrics, changed names of inputs to hintGaussSum for consistency, and changed equality checks to ==]]

We anticipate that most cases will be like the above, where the hint is obtained by just using the stability constant of the transformation. Indeed, when a transformation is stable, this information should also be provided by the relation, so that the hint function can specify just the distance measure d_y and need not calculate the number d_{mid} . More generally, if d_y is a univariate, monotone distance measure (i.e. where if y and y' are d -close under d_y and $d' > d$, then y and y' are also d' -close under d_y), then a hint for d_{mid} can also be automatically obtained by doing a binary search on $R_T((d_{\mathcal{X}}, d_{\text{in}}), (d_y, \cdot))$.

4.6 Assumptions on the Measures

A crucial assumption about the privacy measures D is that they satisfy *post-processing*: if X and X' are random variables distributed that are d -close with respect to D and f is any function, then $f(X)$ and $f(X')$ are also d -close with respect to D . (In particular, D should be a well-defined distance measure on probability distributions over all representable objects, and not specific to their domain \mathcal{X} .) This allows us to implement our post-processing operator on measurements just as before.

[[SV: in future- insert code snippet]]

Beyond post-processing, the proposed framework does not need to assume anything else about the measures. In particular, they need not be metrics in the standard mathematical sense. However, it may be useful to annotate the distance measures with additional properties that may be useful to exploit in some privacy analyses:

1. Symmetry: if x and x' are d -close under $d_{\mathcal{X}}$, then x' and x are d -close under $d_{\mathcal{X}}$. (Note that some of the privacy measures like D_{∞} are most natural as non-symmetric measures.)
2. Monotonicity: whenever x and x' are d -close under $d_{\mathcal{X}}$ and $d' > d$, then x and x' are d' -close under $d_{\mathcal{X}}$.
3. Triangle Inequality: if x and x' are d -close under $d_{\mathcal{X}}$ and x' and x'' are d' -close under $d_{\mathcal{X}}$, then x and x'' are $(d + d')$ -close under $d_{\mathcal{X}}$.
4. Nonnegativity: if x and x' are d -close under $d_{\mathcal{X}}$, then $d \geq 0$.

(Some of these only make sense when $d_{\mathcal{X}}$ measures privacy by a single real number, or at least in some ordered domain with an addition operation.) For example, the measure D_{∞} used in defining pure differential privacy does not satisfy symmetry, but it does satisfy monotonicity, the triangle inequality, and nonnegativity, and the triangle inequality in particular underlies the “group privacy” property of pure differential privacy.

[[SV: for future - MH please tell me about (or point me to) the Recursive Mechanism - I'm not familiar with that term]]

4.7 Composition

Let us illustrate how the basic composition of two DP mechanisms can be analyzed with different privacy measures. Composing an $(\varepsilon_1, \delta_1)$ -DP mechanism and an $(\varepsilon_2, \delta_2)$ -DP mechanism yields an $(\varepsilon_1 + \varepsilon_2, \delta_1 + \delta_2)$ -DP mechanism. Composing a ρ_1 -zCDP and a ρ_2 -zCDP mechanism yields a $(\rho_1 + \rho_2)$ -zCDP mechanism. Let $M_1 : \mathcal{X} \rightsquigarrow \mathcal{Y}$ and $M_2 : \mathcal{X} \rightarrow \mathcal{Z}$ be measurements with privacy relations R_1 and R_2 . Then we can construct a privacy relation R for their composition $M : \mathcal{X} \rightsquigarrow \mathcal{Y} \times \mathcal{Z}$ where $M(x) = (M_1(x), M_2(x))$ as follows:

$$R_M((d_{\mathcal{X}}, d_{\text{in}}), (D, d_{\text{out}})) = \begin{cases} \text{true} & \text{if } D = \text{pureDP} \text{ and } \exists \varepsilon_1 R_1((d_{\mathcal{X}}, d_{\text{in}}), (D, \varepsilon_1)) \wedge R_2((d_{\mathcal{X}}, d_{\text{in}}), (D, d_{\text{out}} - \varepsilon_1)) \\ \text{true} & \text{if } D = \text{zCDP} \text{ and } \exists \rho_1 R_1((d_{\mathcal{X}}, d_{\text{in}}), (D, \rho_1)) \wedge R_2((d_{\mathcal{X}}, d_{\text{in}}), (D, d_{\text{out}} - \rho_1)) \\ \text{true} & \text{if } D = \text{approxDP}, d_{\text{out}} = (\varepsilon, \delta) \in \mathbb{R} \times \mathbb{R}, \text{ and} \\ & \exists \varepsilon_1 R_1((d_{\mathcal{X}}, d_{\text{in}}), (D, (\varepsilon_1, \delta/2))) \wedge R_2((d_{\mathcal{X}}, d_{\text{in}}), (D, (\varepsilon - \varepsilon_1, \delta/2))) \\ \text{false} & \text{otherwise} \end{cases}$$

This composition relation is somewhat unsatisfactory because of the existential quantifiers (which can be eliminated by doing a binary search) and because the composition for approximate DP assumes that the δ is split evenly between the two measurements. Later we will describe a much more general “interactive” composition primitive that allows for the privacy parameters of the individual measurements to be precisely specified. *[[SV: for future - maybe add a code snippet here?]]*

5 Interactive Measurements

The framework described before assumes that measurements are one-shot randomized functions $M : \mathcal{X} \rightsquigarrow \mathcal{Y}$. However, many of the useful primitives in the differential privacy literature, such as Adaptive Composition, the Sparse Vector Technique, and Private Multiplicative Weights are actually *interactive* mechanisms, which allow one to ask an adaptive sequence of queries about the dataset. [\[\[SV: future - add citations\]\]](#) Having a library that supports such interactive measurements is useful both for enabling the design of interactive query systems for end users, as well as tools for the design of even noninteractive differentially private algorithms (such as differentially private gradient descent). [\[\[SV: future - add citations for DP gradient descent\]\]](#)

Here we will describe the basic theory and implementation of interactive measurements. For sake of exposition, we will restrict attention to pure differential privacy and stable transformations on multisets as in Section 2, but it can all be combined with the more general forms of privacy relations described in Section 4.

5.1 Defining Interactive Measurements

[\[\[SV: revamping text and terminology to match new, state-machine perspective\]\]](#)

An interactive measurement M is a (possibly randomized) function that takes a private dataset $x \in \text{MultiSets}(\mathcal{X})$ and then “spawns” a (possibly randomized) state machine $Q = M(x)$ called a *queryable*. The queryable consists of an initial private state s and an evaluation function Eval . It then receives a query q_1 . Based on q_1 and its current state s_0 , it generates a (possibly randomized) answer a_1 and updates its state to s_1 . That is, $(a_1, s_1) \leftarrow \text{Eval}(q_1, s_0)$. It then receives a new query q_2 , and similarly generates an answer and state update as $(a_2, s_2) \leftarrow \text{Eval}_M(q_2, s_1)$. And so on, arbitrarily long.

To define privacy for interactive measurements, we consider an arbitrary adversarial strategy A interacting with $M(x)$, which selects each query q_i adaptively based on all previous answers (a_1, \dots, a_{i-1}) and any randomness of A . Let $\text{View}(A \leftrightarrow M(x))$ be a random variable denoting the A 's *view* of this entire interaction, namely all of A 's randomness and the answers to all queries. We say that M is an ε -DP *interactive measurement* if for every adversarial strategy A , and every pair $x \sim x'$ of adjacent datasets in $\text{MultiSets}(\mathcal{X})$, the random variables $Y = \text{View}(A \leftrightarrow M(x))$ and $Y' = \text{View}(A \leftrightarrow M(x'))$

$$D_\infty(Y||Y') \leq \varepsilon.$$

[\[\[SV: future - add pointers to literature where this modelling is done\]\]](#)

A noninteractive measurement M can be viewed as an interactive measurement M' where $M'(x)$ returns a queryable whose initial state is $s = M(x)$ and whose transition function always returns answer $a = s$ and does not change the state. However, there are more interesting interactive measurements, such the following *adaptive composition measurement*. Here the queryable constructed keeps the dataset x and the remaining privacy-loss budget ε in its state. It takes (noninteractive) measurements as queries, and evaluates them on the dataset if it can do so within the remaining budget, in which case it decrements the privacy loss of the query from the remaining budget. [\[\[SV: added this text description of adaptive composition\]\]](#)

$\text{AdaptiveComposition}_{\varepsilon, \mathcal{X}}(x)$: return the queryable $Q = (s_0, \text{Eval}_{\text{AC}})$ defined as follows.

- $s_0 = (x, \varepsilon)$.
- Query evaluator $\text{Eval}_{\text{AC}}(q, s)$:
 1. Parse s as $s = (x, \varepsilon)$.
 2. If q is not a syntactically valid noninteractive measurement object for data domain \mathcal{X} , set $a = \text{'improper query'}$ and $\varepsilon' = \varepsilon$.
 3. Else if $\text{privacyloss}(q) > \varepsilon$, set $a = \text{'insufficient budget'}$ and $\varepsilon' = \varepsilon$.
 4. Else set $a \leftarrow q(x)$ and $\varepsilon' = \varepsilon - \text{privacyloss}(q)$.

5. Set $s' = (x, \varepsilon')$ and return (a, s') .

Note that this construction uses the fact that the description of a (noninteractive) measurement includes its privacy-loss parameter (ε) and that this is a publicly exposed attribute.

As we see in the above example, to describe an interactive measurement operator, we need to specify:

1. The domain \mathcal{X} of its data records.
2. The privacy loss ε .
3. The (possibly randomized) function f that generates a queryable from a dataset.

[[SV: replacing Spawn with a function to match with earlier syntax for measurements, now that we can do so]] An example implementation of AdaptiveComposition in code is as follows:

[[SV: added some comments to code below]]

```
1 class InteractiveMeasurement:
2     input_domain
3     privacy_loss
4     function      # --- now the function will output a Queryable ---
5
6 class Queryable:
7     _state        # --- underscore indicates that the state should be private ---
8     eval
9     def query(q):
10        (a, _state)=eval(q,_state)
11        return a
12
13 def AdaptiveComposition(dom,epsilon:float):
14     input_domain = dom
15     privacy_loss = epsilon
16     def function(data):
17         initial_state=(data,epsilon)
18         def eval(query: Measurement, state):
19             (st_data, eps) = state
20             if query.input_domain!=dom: return ('domain mismatch',state)
21             elif query.privacy_loss > eps: return ('insufficient budget',state)
22             else return (query.function(st_data),eps-query.privacy_loss)
23         return Queryable(initial_state,eval)
24
25     return InteractiveMeasurement(input_domain,privacy_loss,function)
26
27 # Example
28 queryable_obj=AdaptiveComposition(float,2).function(dataset)
29 res1=queryable_obj.query(NoisySum)
30 res2=queryable_obj.query(NoisyCount)
```

Figure 12: Adaptive Composition *[[SV: can use some more text here. replaced previous valid_meas check with putting type annotation in input of eval and adding explicit domain match check]]* *[[MH: typo: shouldn't say spawn!]]* *[[MH: typo: eval should (answer, state) and not (answer, eps). here, the state includes not only eps but also data]]*

```

1 def ApproxDPAdaptiveComposition(dom,epsilon:float,delta:float):
2   input_domain = dom
3
4   def privacy_relation(d_in,eps,del):
5     return (eps>=d_in*epsilon AND del>=delta*(exp(d_in*epsilon)-1)/(exp(epsilon)-1))
6
7   def function(data):
8     initial_state=(data,epsilon,delta)
9     def eval((query,eps_query,del_query): Measurement * float * float, state):
10      (st_data, eps,del) = state
11      if query.input_domain!=dom: return ('domain mismatch',state)
12      elif (eps_query > eps OR del_query>del): return ('insufficient budget',state)
13      elif not(query.privacy_relation(1,(eps_query,del_query))): return ('query not private as claimed'
14      else return (query.function(st_data),eps-eps_query,del-del_query)
15      return Queryable(initial_state,eval)
16
17   return InteractiveMeasurement(input_domain,privacy_loss,function)
18 \label{fig:approx_comp}

```

Figure 13: Sketch of Adaptive Composition for Approx DP - in response to Q from Andy

```

1 class InteractiveOdometer:
2     input_domain
3     function      # --- now the function will output a Queryable ---
4
5 class OdometerQueryable:
6     _state        # --- underscore indicates that the state should be private ---
7     eval
8     privacy_loss_eval  # --- NOTE: need to define the semantics of this
9
10    def query(q):
11        (a, _state)=eval(q,_state)
12        return a
13
14    def privacy_loss():      # --- for other privacy measures this might need to be a
15                            # --- relation, but that theory hasn't been worked out
16        return privacy_loss_eval(_state)
17
18 def OdometerToMeasurement(odometer : InteractiveOdometer, epsilon : float)
19     privacy_loss = epsilon
20     input_domain = odometer.input_domain
21
22     def function(data):
23         odometer_queryable = odometer.function(data)
24         initial_state = None
25         def eval(q,state):
26             answer = odometer_queryable.query(q)
27             if odometer_queryable.privacy_loss() <= epsilon: return (answer,None)
28             else return ('exceeded budget',None) # the budget has been spent even though user did not receive
29         return Queryable(initial_state,eval)
30
31     return InteractiveMeasurement(input_domain,privacy_loss,function)
32
33
34 def AdaptiveComposition(dom):
35     input_domain = dom
36
37     def function(data):
38         initial_state=(data,0)
39         def eval(query: Measurement, state):
40             (st_data, eps) = state
41             if query.input_domain!=dom: return ('domain mismatch',eps)
42             else return (query.function(st_data),eps+query.privacy_loss)
43
44         def privacy_loss_eval(state):
45             return state
46
47         return OdometerQueryable(initial_state,eval,privacy_loss_eval)
48
49     return InteractiveOdometer(input_domain,function)
50
51 def ConcurrentCompositionOdometers(dom)
52     input_domain = dom
53
54     def function(data):
55         initial_state = [] # the list of the queryables spawned
56
57         def eval(query, state):
58             query_type = query[0]
59             query_details = query[1]
60
61             if query_type == 'prox odometer' and isinstance(query_details, InteractiveOdometer) and query_details

```

In the example of AdaptiveComposition, the only portion of the state $s = (x, \varepsilon)$ that needs to be kept secret is the dataset x , and we can safely augment Eval with additional queries that allow asking for remaining privacy-loss budget ε . However, other interactive queryables do have additional state that needs to be kept secret for their privacy properties (e.g. the noisy threshold in the Sparse Vector Technique).

5.2 Post-processing

Recall that standard post-processing refers to performing an arbitrary computation with the privacy-protected output of a noninteractive measurement operator. For interactive measurements, we think of the queryable itself as the analogue of the “privacy-protected output” of the measurement operator — no matter how one computes with the queryable (as a black box, without examining its internal state!), privacy is maintained. This gives rise to a post-processing principle for interactive measurements. Specifically we can apply a queryable mapping function P that takes the queryable $Q = M(x)$ produced by M and produces a new queryable $Q' = P(Q)$. Importantly, P does not get to examine the internals of the queryable Q , only interact with it as a (stateful) black box, issuing queries and receiving answers. P can also embed the queryable Q *inside* the queryable Q' , so that whenever Q' receives a query, it can issue some queries to Q to help compute an answer. Note that Q continually updates its state as P and then Q' issues queries to it. $\text{Postprocess}(M, P)$ is the interactive measurement $M'(x) = P(M(x))$ that outputs Q' . *[[SV: replaced long math description of post-processing with above much simpler, albeit slightly less formal, description]]*

The reason privacy is preserved under this interactive form of post-processing is that for every adversary A interacting with $\text{Postprocess}(M, P)$, there is an adversary A^{in} interacting with M and a function f^{out} such that for every dataset x ,

$$\text{View}(A \leftrightarrow \text{Postprocess}(M, P)(x)) = f^{\text{out}}(\text{View}(A^{\text{in}} \leftrightarrow M(x))).$$

That is, views of an adversary interacting with the post-processed mechanism can be obtained by applying a function to the view of an adversary interacting with the original mechanism, and thus differential privacy of the latter implies differential privacy of the former. *[[SV: more detailed proof in old notation commented out below]]*

We illustrate this interactive post-processing with two examples. First, we will show how to obtain an instantiation of a differentially private Statistical Query (SQ) Model as a post-processing of AdaptiveComposition, and then we will see how to obtain an instantiation of noisy gradient descent as a post-processing of the SQ model.

Statistical Query Model. In the SQ Model, we are given a dataset $x \in \text{MultiSets}(\mathcal{X})$, and an analyst can issue up to T queries that can be issued are bounded functions $f : \mathcal{X} \rightarrow [-B, B]$ and obtain noisy estimates of the average $\mathbb{E}_{z \leftarrow x} [f(z)] = (\sum_{z \in x} f(z))/|x|$. We will implement this using $T + 1$ queries to an $\text{AdaptiveComposition}_{\mathcal{X}, \varepsilon}$ queryable, spending half the budget on an initial query to estimate the size of the dataset, and then dividing the remaining budget evenly over the T remaining queries to estimate the sum $\sum_{z \in x} f(x)$. In our implementation, we won't require that the queries f are bounded as specified, but will rather enforce it by evaluating the sums using NoisyClampedSum .

To this end, we will assume that for any function $f : \mathcal{X} \rightarrow \mathcal{Y}$, we can construct the stability 1 transformation $\text{RowTransform}_f : \text{MultiSets}(\mathcal{X}) \rightarrow \text{MultiSets}(\mathcal{Y})$ defined as

$$\text{RowTransform}_{\mathcal{X}, \mathcal{Y}, f}(z) = \{f(w) : w \in z\} \text{ (with multiplicity)}.$$

When $\mathcal{Y} = \mathbb{R}$, we can then use chaining to implement the ε -DP mechanism

$$\text{NoisySumFunction}_{f, L, U, \varepsilon} = \text{NoisyClampedSum}_{L, U, \varepsilon} \circ \text{RowTransform}_{\mathcal{X}, \mathbb{R}, f} : \text{MultiSets}(\mathcal{X}) \rightsquigarrow \mathbb{R}.$$

(We note that allowing arbitrary user-provided code to specify functions f can lead to severe side-channel attacks. We discuss how to address such attacks in Section 8.)

[[SV: can probably use more math description before the code]]

In code, we have:

```

1 def Postprocess(intMeas: InteractiveMeasurement, queryable_map):
2     input_domain = intMeas.input_domain
3     privacy_loss = intMeas.privacy_loss
4     def function(data):
5         queryable_inner=intMeas.function(data)
6         return queryable_map(queryable_inner)
7     return InteractiveMeasurement(input_domain,privacy_loss,function)
8
9 # Example
10 def MakeRowTransform(in_dom, out_dom, f):
11     ...
12
13 def MakeNoisySumFunction(in_dom,f,L,U,epsilon):
14     return(ChainingMT(MakeNoisyClampedSum(L,U,epsilon),
15                       MakeRowTransform(in_dom, float, f)))
16
17
18 def MakeSQmodel(in_dom,T,B,epsilon):
19     def queryable_map(AC_queryable):
20         eps=epsilon/2
21         def sum_query(x): return 1
22         n=AC_queryable.query(MakeNoisySumFunction(in_dom,sum_query,-1,1,eps))
23         initial_state=T
24         def eval(query, state):
25             if state>0:
26                 answer=AC_queryable.query(MakeNoisySumFunction(in_dom,query,-B,B,eps/T))/n
27             else:
28                 answer='no more queries'
29             return (answer,state-1)
30         return Queryable(initial_state,eval)
31     return Postprocess(AdaptiveComposition(in_dom,epsilon),queryable_map)

```

Figure 15: Differentially Private SQ Model *[[SV: can use an elaborated caption]]* *[[MH: this is a corrected version: now RowTransform, MakeNoisySumFunction are appropriately generic; also sum_query just takes a single arg]]*

Now we use our implementation of the SQ Model to obtain an instantiation of noisy gradient descent via post-processing. Suppose we have a dataset $z \in \text{MultiSets}(\mathbb{R} \times \mathbb{R})$ and our goal is to find a differentially private slope θ minimizing $L_z(\theta) = \mathbb{E}_{(x,y) \in z}[\ell_{(x,y)}(\theta)]$, where $\ell_{x,y}(\theta) = (\theta x - y)^2$. In each iteration of the algorithm, we will take our current estimate θ , compute a differentially private estimate a of the derivative $L'_z(\theta) = \mathbb{E}_{(x,y) \in z}[\ell'_{x,y}(\theta)]$, where $\ell'_{x,y}(\theta) = 2\theta \cdot (\theta x - y)$, and update to $a \leftarrow \theta - \eta a$, for a “learning rate” parameter η . After a given number T of iterations, we simply output the current value of θ , embedded in the state of a “noninteractive queryable” that will return θ on any query. Notice that in each iteration, we are simply computing an average over the dataset, which we can estimate using a query to an SQ model queryable. We will also need to specify a bound B to be enforced by the SQ model.

In code:

```

1 def NoninteractiveQueryable(initial_state):
2     def eval(q,s): return (s,s) # always return the state and never change it
3     return Queryable(initial_state,eval)
4
5 # Example
6
7 def MakeNoisyGD(T,B,eta,epsilon):
8     def queryable_map(SQ_queryable):
9         theta=0
10        for i in range(T):
11            def derivative(z : float*float): return (2*theta*(theta*z[0]-z[1]))
12            a=SQ_queryable.query(derivative)
13            theta=theta-eta*a
14        return NoninteractiveQueryable(theta)
15    return Postprocess(MakeSQmodel(float*float,T,B,epsilon),queryable_map)

```

Figure 16: Noisy Gradient Descent *[[SV: can use more commentary in caption]]* *[[MH: shouldn't derivative take a single arg \mathbf{x} , which in this case, would be a 2-tuple of floats?]]* *[[SV: fixed. also changed the type to float*float which I think we used elsewhere in the paper]]*

The examples of AdaptiveComposition and NoisyGD show that it is useful to build interactive measurements from noninteractive ones and vice-versa. Thus, we propose that noninteractive measurements are implemented as a special case of interactive measurements, to avoid duplicating operators (composition, chaining, post-processing, etc.) that can simultaneously deal with both kinds of measurements.

5.3 Chaining and Composition of Interactive Measurements

Chaining generalizes in a straightforward way to interactive measurements. If $T : \text{MultiSets}(\mathcal{X}) \rightarrow \text{MultiSets}(\mathcal{Y})$ is a c -stable transformation and M is an ε -DP interactive measurement for data domain \mathcal{Y} , then their chaining $M \circ T$ is a $c\varepsilon$ -DP interactive measurement for data domain \mathcal{X} . *[[SV: probably later: comment on how we may often want the transformation to be carried out by the database system, rather than bringing x into memory and computing it. will need to expose capabilities of the backend database to the library and have clear specs about what we are assuming about the backend database in terms of functionality, side channels, etc.]]*

As far as composition, it is natural to extend the AdaptiveComposition procedure described above to allow queries that can be *interactive* measurement operators themselves. For example, we should be able to issue a query q_i describing an interactive measurement operator M_{q_i} that spawns an “inner queryable” $M_{q_i}(x)$ within the Adaptive Composition queryable, to which we can issue subsequent queries. We can then choose to allow for either:

1. *Sequential Composition*: all of the queries to the first inner queryable must be completed before another inner queryable is spawned. *[[MH: how does one mark a queryable as complete or finished? it seems like with current formulation, even after budget consumed, one can continue to query indefinitely (even though the answer would be ‘out of budget’).]]*
2. *Concurrent Composition*: multiple inner queryables can be spawned and be simultaneously active, with queries to them arbitrarily interleaved.

The basic, additive composition of privacy loss for pure differential privacy applies for both sequential composition and concurrent composition; indeed, PINQ allows for concurrent composition and a formal proof of its correctness is given in [3]. *[[SV: check if this is correct]]*

Another restriction of AdaptiveComposition as described above is that when specifying the i 'th query $q_i = M_i$, the privacy-loss $\varepsilon_i = \text{privacyloss}(q_i)$ is also specified, and is effectively “consumed” immediately

when $\text{Spawn}_{q_i}(x)$ is executed. In contrast, in PINQ, each “inner queryable” does not have a preallocated budget but separately tracks its accumulated privacy loss, and these are summed in order to track the overall privacy loss. To support this in our proposed framework, we would need to implement a theory of *privacy odometers* [12], which are variants of interactive measurements that track accumulated privacy loss rather than specify a total privacy loss at the beginning. It should be relatively straightforward to modify the programming framework to implement composition, chaining, and post-processing for pure-DP odometers, as well as conversion of odometers to standard interactive measurements (by not answering a query if it will make the accumulated privacy loss exceed the budget).

However, we remark that the aforementioned generalizations of composition appear to be more complex for variants of differential privacy (such as approximate differential privacy). Indeed, the Advanced and Optimal Composition Theorems for approximate differential privacy [\[\[SV: future-cites\]\]](#) do not hold as is for approximate-DP odometers, and in fact they require the entire sequence of privacy-loss parameters $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_T$ to be specified before any queries are made. (Specifying the sequence of privacy-loss parameters in advance is often a good thing in any case, as it allows for verifying that the budget won’t inadvertently be consumed before an analysis is complete.) As far as we know, there has also been no rigorous analysis of concurrent composition for approximate differential privacy or other variants of differential privacy; this also seems to be an important direction for future work. (The cryptography literature indicates that concurrent composition is often much more subtle than sequential composition. [\[\[SV: future-cites\]\]](#))

[\[\[SV: later: data access and where transformations are executed\]\]](#)

5.4 Interactive Transformations

```
1 def MakeLazyCollection(im,source_name):
2   def queryable_map(im_queryable):
3     initial_state={source_name: NoOp()}
4     def eval(query, state):
5       if query.type == "Transform":
6         tfm = state[query.source] # tfm is what "produced" the source dataset
7         new_tfm = query.op # the transformation user wants to apply to source
8         try:
9           state[query.dest] = chainTT(new_tfm, tfm)
10        except InvalidChaining:
11          answer = "Cannot apply {new_tfm} to {query.source}!"
12        else:
13          answer = "Applied transformation {new_tfm} to {query.source} to produce {query.dest}"
14      elif query.type == "Measure":
15        tfm = state[query.source] # look up the tfm that produced source
16        meas = query.op # the measurement user wants to apply to source
17        try:
18          measurement = chainMT(meas, tfm)
19        except InvalidChaining:
20          answer = "Cannot apply {meas} to {query.source}!"
21        else:
22          answer=im_queryable.query(measurement)
23      else:
24        answer='unrecognized query type'
25      return (answer,state)
26    return Queryable(initial_state,eval)
27  return Postprocess(im,queryable_map)
```

Figure 17: Interactive Transformation done lazily via post processing. The user can now produce new transformed datasets and query them, but these datasets are produced lazily when a measurement request arrives. This can be applied to any InteractiveMeasurement operator; though it will only produce useful results if the IM operator accepts Measurements as queries.

```

1 def AdaptiveCompositionWithInteractiveTransformations(dom,epsilon:float):
2     input_domain = dom
3     privacy_loss = epsilon
4     def function(data):
5         initial_state=({'root': (data, NoOp()), '_epsilon': epsilon})
6         def eval(query: QueryObject, state):
7             if query.type == "Transform":
8                 (d, tfm) = state[query.source] # previously computed private data d and the tfm that made it
9                 new_tfm = query.op # the transformation user wants to apply to source
10                try:
11                    chained_tfm = chainTT(new_tfm, tfm)
12                except InvalidChaining:
13                    answer = "Cannot apply {new_tfm} to {query.source}!"
14                else:
15                    new_d = chained_tfm.function(d)
16                    state[query.dest] = (new_d, chained_tfm)
17                    answer = "Applied transformation {new_tfm} to {query.source} to produce {query.dest}"
18            elif query.type == "Measure":
19                (d, tfm) = state[query.source] # look up the tfm that produced source
20                eps = state['_epsilon']
21                meas = query.op # the measurement user wants to apply to source
22                try:
23                    new_meas = chainMT(meas, tfm) # this new_meas operator will be made but never run
24                except InvalidChaining:
25                    answer = "Cannot apply {meas} to {query.source}!"
26                else:
27                    if new_meas.input_domain!=dom: answer = 'domain mismatch'
28                    elif new_meas.privacy_loss > eps: answer = 'insufficient budget'
29                else:
30                    answer = meas.function(d) # note: we apply meas, not new_meas!
31                    state['_epsilon'] = eps-query.privacy_loss
32            else:
33                answer='unrecognized query type'
34            return (answer,state)
35        return Queryable(initial_state,eval)
36
37    return InteractiveMeasurement(input_domain,privacy_loss,function)

```

Figure 18: Eager evaluation

6 Compute Environment and Interactive Transformations

The basic data structure in the proposed compute environment is a `DataTree`, which represents a tree of derived datasets. The root represents the original data(set). The children of a node represent data derived by applying (deterministic, pure) functions to the parent data. The pointers between a parent and child are labelled with the function used to derive the child. Only the root needs to actually store the data; the data at other nodes is implicitly determined by the sequence of functions from the root to the node. The data at non-root nodes is only computed when actually needed using `read-data`. At each node we also store a “hint” which can be a request to either cache the data at the node or erase the data at the node after the next `read-data` command.

```

1 class DataTree:
2     _data      # --- a dataset
3     children   # --- a list of pairs [f,T] where T is the child DataTree labelled by f
4     _parent    # --- a pair [f,T] where T is the parent DataTree and self is the child of T labelled by f
5     hint       # --- takes values from {None,cache,erase}
6
7     def __init__():
8         _data=None
9         children=[]
10        _parent=None
11        hint=None
12
13    def fill_data(data): # --- fill in the data at the root of the tree
14        if _parent==None:
15            _data=data
16            hint='cache'
17
18    # --- apply_functions returns the DataTree obtained by applying function_sequence
19    # --- function sequence is a list of functions alternating with hints
20    # --- it does not actually compute the derived data, but just returns the
21    # --- corresponding node of the datatree (creating it if necessary)
22
23    def apply_functions(function_sequence):
24        # --- base case: no functions to apply
25        if function_sequence==[]: return self
26
27        # --- recursive case: pull off first function to apply and corresponding hint
28        f = function_sequence[0]
29        hint = function_sequence[1]
30        remaining_functions = function_sequence.pop(0).pop(0)
31
32        # --- if already have a child labelled by f
33        # --- then recursively apply remaining_functions to that child
34        for fTpair in self.children:
35            if fTpair[0]==f:
36                child = fTpair[1]
37                child.hint = hint
38                return child.apply_functions(remaining_functions)
39
40        # --- if no child is labelled by f
41        # --- create a child labelled by f and recursively apply remaining_functions to it
42        child = DataTree()
43        child._parent = [f,self]
44        child.hint = hint
45        self.children.append([f,child])
46        return child.apply_functions(remaining_functions)
47
48    # --- actually compute and return the derived data corresponding to current node
49    # --- derived recursively from data at parent
50    def read_data():
51        if _data is None:
52            data=self._parent[0](self._parent[1].read_data())
53            if hint='cache': _data=data
54        else
55            data=_data
56            if hint='erase': _data=None
57        return data

```

```

1 class InteractiveMeasurement:
2     input_domain
3     privacy_loss
4     function      # --- now the function will take a DataTree as input ---
5
6 class Transformation:
7     input_domain
8     output_domain
9     stability
10    function_sequence # --- now an odd-length list of alternating functions and hints
11
12 def MakeClamp(L: float, U: float):
13     if L > U: raise Exception('Invalid parameters')
14     input_domain = float
15     output_domain = bounded_float(L,U)
16     stability = 1
17     def function(data):
18         def clamp(x): return max(min(x, U), L)
19         return map(clamp, data)
20     return Transformation(input_domain,output_domain,stability,[function])
21     # --- now we put the function in a list
22
23 # add hint for whether to cache or erase intermediate value
24 def ChainingTT(trans_2: Transformation, trans_1: Transformation, hint):
25     if (trans_1.output_domain!=trans_1.input_domain): raise Exception('Domain mismatch')
26     input_domain = trans_1.input_domain
27     output_domain = trans_2.output_domain
28     stability = trans_1.stability*trans_2.stability
29     function_sequence = trans_1.function_sequence + [hint] + trans_2.function_sequence # --- concatenatio
30     return Transformation(input_domain,output_domain,stability,function_sequence)
31
32
33 def ChainingMT(meas: InteractiveMeasurement, trans: Transformation, hint):
34     # makes new interactive measurement
35     if (trans.output_domain!=meas.input_domain): raise Exception ('Domain mismatch')
36     input_domain = trans.input_domain
37     privacy_loss = trans.stability*meas.privacy_loss
38     # --- note that the function does not actually compute or read the transformed data.
39     # --- that will only happen inside meas.function
40     def function(datatree : DataTree):
41         return meas.function(datatree.apply_functions(trans.function_sequence+[hint]))
42     return Measurement(input_domain,privacy_loss,function)
43
44
45 def MakeNoisySum(L: float, U: float, epsilon: float):
46     if L > U or epsilon < 0: raise Exception('Invalid parameters')
47     input_domain = bounded_float(L,U)
48     privacy_loss = epsilon
49     def function(datatree : DataTree):
50         data_sum = datatree.apply_functions([sum,None])
51         s = max(abs(L), abs(U))
52         z = Laplace(s/epsilon,0)
53         return NoninteractiveQueryable(data_sum.read_data() + z) # --- transformations only evaluated a
54     return InteractiveMeasurement(input_domain,privacy_loss,function)

```

Figure 20: Transformations, Interactive Measurements, and Chaining with DataTrees

```

1  # first set up an interactive measurement overall_meas in any data independent way (like in PSI exampl
2  ...
3  # check that overall_meas has the desired privacy properties
4  assert (overall_meas.input_domain == MultiSets(strings,n) &&
5         overall_meas.input_metric == change-one &&
6         overall_meas.output_measure == approxDP &&
7         overall_meas.privacy_relation(1,(epsilon,delta))
8
9  # initialize a datatree for your dataset data
10 datatree=DataTree()
11 datatree.fill_data(data)
12
13 # create a queryable by applying overall_meas to the datatree
14 queryable = overall_meas.function(datatree)
15
16 # now access the data only through queryable
17 ...

```

Figure 21: Setting up and Using a DataTree

```

1  def MakeLazyCollection(im,source_name):
2      def queryable_map(im_queryable):
3          initial_state={source_name: NoOp()}
4          def eval(query, state):
5              if query.type == "Transform":
6                  tfm = state[query.source] # tfm is what "produced" the source dataset
7                  new_tfm = query.op # the transformation user wants to apply to source
8                  try:
9                      state[query.dest] = ChainingTT(new_tfm, tfm, 'cache')
10                 except InvalidChaining:
11                     answer = "Cannot apply {new_tfm} to {query.source}!"
12                 else:
13                     answer = "Applied transformation {new_tfm} to {query.source} to produce {query.dest}"
14             elif query.type == "Measure":
15                 tfm = state[query.source] # look up the tfm that produced source
16                 meas = query.op # the measurement user wants to apply to source
17                 try:
18                     measurement = ChainingMT(meas, tfm, 'cache')
19                 except InvalidChaining:
20                     answer = "Cannot apply {meas} to {query.source}!"
21                 else:
22                     answer=im_queryable.query(measurement)
23             else:
24                 answer='unrecognized query type'
25             return (answer,state)
26         return Queryable(initial_state,eval)
27     return Postprocess(im,queryable_map)

```

Figure 22: Interactive Transformations with Caching

To-do on the compute environment:

1. Better memory management. Right now nodes of the tree are never erased, only the data fields are. How can we track whether we can erase a part of the tree, e.g. when none of the data fields in that part of the tree are supposed to be cached and when no code still has pointers to that part of the tree?
2. Applying randomized functions. Need to distinguish between applications of a function with the same randomness vs. fresh randomness. Maybe the functions fed to the datatree should always be deterministic, so we need to hardcode random seeds into randomized functions before applying them.
3. Creating a registry of functions that can be applied to allow for using external compute resources.
4. Partitions and joins: will this move us from trees to dags? Partitions don't seem to require dags.
5. Spark already has some form dag. Spark literature - maybe some useful tricks. Spark has transformations and actions. Compile the plan.

Comments: reason to only have execution happen within measurements is because measurements are responsible for all privacy protection, including timing channels.

6.1 Working toward Data DAGs

```
1 class StraightLineProgram:
2     num_lines    # the number of lines in the program
3     num_inputs  # the number of inputs
4     num_outputs # the number of outputs
5     # --- the first num_inputs lines are the inputs, and the last num_outputs lines are the outputs
6     # --- so num_lines >= num_inputs+num_outputs
7
8     parents      # a list of length num_lines such that for all i>=num_inputs,
9                 # parents[i] is the list of parents of line i
10
11    functions    # a list of length num_lines such that for all i>=num_inputs,
12                # functions[i] is a function of arity parents[i]
13
14    requests     # a list of length num_lines such that for all i>=num_inputs,
15                # request[i] is the cache/erase request (or hint)
16
17 class DerivedData:
18     _program     # a straight-line program with num_inputs=num_outputs=1
19     _data        # a dataset obtained by applying _program
20
21 class DataStore:
22     # --- a data structure that maintains an initial_data
23     # --- and collection of triples (datakey,program,derived_data,request) where
24     # --- program is a SLP with num_inputs=num_outputs=1
25     # --- and derived_data = program(initial_data)
26     # --- and request is cache or erase
27     # --- each triple represents a single node in the dynamic compute DAG
28     # --- which comes from overlaying all the individual programs
29     # --- on top of each other according to shared functions.
30     # --- and is also stored in this dat structure
31
32     def __init__(initial_data):
33         # --- creates triple (datakey,identity,initial_data)
34         return datakey
35
36     def lookup(program,datakeys):
37         # --- program.num_outputs should equal length(datakeys)
38         # --- this looks for a node in the DAG that is
39         # --- syntactically equivalent (up to topological sort) to
40         # --- applying program to the nodes indexed by the list of datakeys
41         # --- if it doesn't find one, it recursively looks for
42         # --- or constructs the parents of the output line,
43         # --- and then creates this node as a child
44         return datakey
45
46     def read_data(datakey):
47         # --- recursive evaluation, storing or erasing
48         # --- intermediate results according to cache requests
49
50     def erase(datakey): ...
```

Figure 23: The Compute Environment with DAGs

```

1 class InteractiveMeasurement:
2     input_domain # --- now a tuple of data domains (for multi-arity functions like join)
3     privacy_loss
4     function      # --- now the function will take a datastore and list of datakeys ---
5
6 class Transformation:
7     input_domain # --- now a tuple of data domains
8     output_domain # --- now a tuple of data domains
9     stability
10    program # --- now a straight-line program with matching arities
11
12 def MakeClamp(L: float, U: float):
13     if L > U: raise Exception('Invalid parameters')
14     input_domain = float
15     output_domain = bounded_float(L,U)
16     stability = 1
17     def function(data):
18         def clamp(x): return max(min(x, U), L)
19         return map(clamp, data)
20     return Transformation([input_domain],[output_domain],stability,SLP(function))
21     # --- SLP(function) returns the corresponding two-line straight-line program
22     # with unspecified cache request on the output
23
24 def ChainingTT(trans_2: Transformation, trans_1: Transformation, request):
25     if (trans_1.output_domain!=trans_1.input_domain): raise Exception('Domain mismatch') # actually need
26     input_domain = trans_1.input_domain
27     output_domain = trans_2.output_domain
28     stability = trans_1.stability*trans_2.stability
29     program = Concatenate(trans_1.program,trans_2.program,request)
30     # --- concatenate the straight-line programs, with appropriate
31     # --- requests on the intermediate nodes
32     # --- note that 'request' itself can be a tuple of requests
33     return Transformation(input_domain,output_domain,stability,program)
34
35
36 def ChainingMT(meas: InteractiveMeasurement, trans: Transformation, requests):
37     # makes new interactive measurement
38     if (trans.output_domain!=meas.input_domain): raise Exception ('Domain mismatch') # actually check equ
39     input_domain = trans.input_domain
40     privacy_loss = trans.stability*meas.privacy_loss
41     # --- note that the function does not actually compute or read the transformed data.
42     # --- that will only happen inside meas.function
43     def function(datastore,input_datakeys):
44         for i=0 to trans.program.num_outputs-1:
45             subprogram = i'th output of trans.program with request[i] attached
46             output_datakeys[i] = datastore.lookup(subprogram,input_datakeys)
47         return meas.function(datastore,output_datakeys)
48     return Measurement(input_domain,privacy_loss,function)

```

Figure 24: Transformations, Interactive Measurements, and Chaining with DataDAGs

```

1 def MakeNoisySum(L: float, U: float, epsilon: float):
2     if L > U or epsilon < 0: raise Exception('Invalid parameters')
3     input_domain = [bounded_float(L,U)]
4     privacy_loss = epsilon
5     def function(datastore,datakeys):
6         data_sum = datastore.lookup(SLP(sum),datakeys)
7         s = max(abs(L), abs(U))
8         z = Laplace(s/epsilon,0)
9         return NoninteractiveQueryable(datastore.read_data(data_sum) + z) # --- transformations only ev
10    return InteractiveMeasurement(input_domain,privacy_loss,function)
11
12 def MakePartition(rowdomain, predicate):
13     input_domain = [rowdomain]
14     output_domain = [rowdomain,rowdomain] # output is a pair of datasets
15     stability = 1
16     program = SLP(filter(predicate),filter(not(predicate)))
17     # --- the straight-line program whose first output is the
18     # --- dataset filtered by the predicate, and whose second
19     # --- output is the one filtered by its negation
20    return Transformation(input_domain,output_domain,stability,program)
21
22 def ParallelComposition(meas0 : InteractiveMeasurement, meas1 : InteractiveMeasurement):
23     input_domain = [meas0.input_domain,meas1.input_domain]
24     privacy_loss = max(meas0.privacy_loss,meas1.privacy_loss)
25
26     def function(datastore,datakeys):
27         queryable0=meas0.function(datastore,[datakeys[0]])
28         queryable1=meas1.function(datastore,[datakeys[1]])
29         return Compose(queryable0,queryable1)
30         # --- should specify sequential or concurrent composition
31
32    return InteractiveMeasurement(input_domain,privacy_loss,function)

```

Figure 25: MakeNoisySum, Partitioning, and Parallel Composition with DataDAGs

Notes on Compute DAGs:

- For simplicity
 - all data domains are MultiSets
 - input and output metrics are always Sum of the AddRemove metric
 - PrivacyMetric is always PureDP
- Should we construct intermediate nodes of DAG during lookup or read-data?

7 Putting it Together

In Section 2, we described a basic framework of pure DP measurements and stable transformations on multisets. In Sections 3, 4, and 5, we described three orthogonal extensions to the basic framework:

1. A framework for ensuring that measurements and transformations have their claimed privacy or stability properties, by only allowing procedures that have been vetted (e.g. through the OpenDP editorial board reviewing a written proof) the right to directly construct measurements and transformations.
2. Allowing for many sensitive datatypes (not just multisets), many measures of distance between sensitive data items (not symmetric difference), and many measures of closeness between probability distributions (not just the multiplicative notion used to define pure differential privacy). With this generalization, stability and privacy are no longer measured by single parameters, but are given by *relations* that tell us whether a given kind of “closeness” of inputs implies a given kind of “closeness” of outputs.
3. Generalizing measurements to spawn interactive query procedures (randomized state machines called *queryables*). This allows for modelling adaptive composition theorems as measurements, as well as more sophisticated interactive DP algorithms. It also provides for in-library (and thus vetted) mechanisms for managing privacy budgets, rather than delegating this to external systems.

Although we described these extensions separately for sake of exposition, we propose that all three be combined in the OpenDP Programming Framework. They combine with each other in natural ways:

- Our definition of privacy for interactive measurements readily extends to other privacy notions, in particular by simply asking that the views of the adversary on two input data items (that are close according to some distance measure) are close according to any desired notion of closeness between probability distributions (capturing, for example, approximate differential privacy or concentrated differential privacy). We can replace the Adaptive Composition interactive measurement for pure DP with other interactive composition procedures that can be analyzed using other privacy notions, and the resulting bounds would be encoded by the privacy relation.
- The verification principle from Section 3 can and should be applied to more general measurements and transformations from Sections 4 and 5. In particular, now a measurement or transformation is considered *valid* if its privacy or stability relation is sound: it should never claim that a false privacy or stability bound holds.

8 Ensuring Privacy in Implementation

Like with any other implementation of security or privacy critical code, we need to beware of places where the mathematical abstraction used in the privacy analysis may deviate from the implementation.

For example, in the framework above we abstract transformations as functions $T : \mathcal{X} \rightarrow \mathcal{Y}$ and measurements as randomized functions $M : \mathcal{X} \rightsquigarrow \mathcal{Y}$. Implicit in this modelling is that when executing $T(x)$, the only private information it has access to is x , and that it does not produce any observable output other than $T(x)$. That is, we need to prevent leakage of information through *side channels*.

One known side channel of concern is the *timing channel*, where the amount of time to execute a function reveals information about its input. Standard mitigations against the timing channel include introducing delays or truncating long executions to make a function essentially constant time. This problem and side channels in general are exacerbated when a transformation or measurement is built using a user-defined function f (which may come from an interactive query issued by an adversary). For example, `RowTransformf` can cause privacy leakage if the function f can write to an I/O device observable to the adversary or if the function f intentionally changes its execution time based on the data record it is fed. Side-channel attacks using user-defined functions can be mitigated by some programming language support. (See Section 9.)

Another concern comes from the implementation of arithmetic. Like in much of statistics, in the mathematical analysis of differential privacy, we often model data using real numbers, use continuous noise distributions, and allow probabilities that are arbitrary real numbers in $[0, 1]$. Indeed, we followed this tradition above, starting from the beginning of Section 2 when we presented `NoisySum` as taking input datasets whose records are elements of the real interval $[L, U]$ and as adding continuous Laplace noise to the output.

In the code examples, however, we used floating-point numbers as a proxy for real numbers. Unfortunately, as shown by Mironov [9] the discrepancy between floating-point numbers and true real numbers leads to severe violations of differential privacy. For this reason, following [1], we advocate the use of *fixed-point* numbers and integer arithmetic as much as possible in the OpenDP library. For code that is in the fully certified library, any deviations from standard arithmetic (e.g. overflow of floating-point arithmetic) should be explicitly modelled in the privacy analysis. (There may be versions of the library that are subject to less rigorous constraints for the purpose of prototyping and experimentation.) *[[SV: have different levels of ‘certification’ depending on which floating-point and side-channel issues are addressed. should it be library code or the measurements and transformations themselves that have the level of certification? e.g. if we produce a grade A adaptive composition measurement, then it should also only accept grade A queries (and also be careful with any use of floating point arithmetic it uses for budget management)]]*

Differentially private algorithms also rely crucially on high-quality randomness. To avoid this being a source of vulnerability, the OpenDP library should use a cryptographic-strength pseudorandom generator for its random bits.

9 Discussion of Implementation Language(s)

The choice of programming language, or languages, for the OpenDP library implementation is an important one. It has implications on performance, reliability and also impacts the degree to which the library will achieve each of the goals identified in the introduction. This section identifies some features of the programming language that might be particularly beneficial and offers a concrete proposal for OpenDP. It is likely that no single language meets all criteria and that different components of the library will be implemented using different languages. Our initial prototyping was done using Python with a mix of object-oriented and functional paradigms, as reflected in the code examples above, and many of the features identified below are a reflection of that experience.

Section 1 outlines high-level goals including *usability* and *programmability*. To support the goal of usability, it will be important that the library offer APIs in languages that are commonly used in OpenDP’s target use cases. The most common (open-source) languages for statistical/data analysis include Python and R. By providing APIs in one or both of these languages, programmers will be able to integrate the library into the conventional data workflows. There will also be a need to integrate with a variety of back-end data access layers, which may necessitate re-implementing some core functionality in a back-end compatible language such as SQL. It is anticipated that this re-implementation would be limited to data-intensive transformations, though prior work has also explored implementing measurements in the data layer [6]. *[[SV: future-cite published version]]* To support the goal of programmability, it may be preferable to select a language in the more familiar imperative paradigm rather than adopt a purely functional language, though as noted below, choosing a language that has support for some functional features may be desirable.

9.1 Desired language support

Below are some specific programming language features that are desired.

Memory safety, encapsulation and immutability. Given the library’s principal goal of mediating access to private data, securing the data against unauthorized access is a paramount concern. While systems hardening will likely happen at many layers, some features at the programming language level can help avoid inadvertent data leakage and reduce the verification effort for new contributions. These features include memory safety and state encapsulation (in the object-oriented sense). For example, it should be impossible for the user to access the internal state of a queryable. In addition, the ability to control mutability may be desirable. For instance, if a query is implemented as a mutable object, then the programmer could potentially manipulate the state of the query object *while it is being evaluated by the queryable*. Immutable structures may help avoid programmer mistakes and make it easier to automatically analyze a program for its safety.

Abstract (parameterized) data types. Many of the code examples presented earlier make use of specialized data types. In many cases these types are similar to conventional types but with additional

constraints on allowable values (e.g., a float that is bounded). Being able to precisely constrain the data type has tangible benefits in the context of a DP library: it often simplifies privacy semantics for individual functions and invites more modular design (upstream transformations often refine the type to make downstream processing easier). It also reduces the burden for verifying contributed functions: since the privacy promise is conditional on type, narrower types reduce the scope of the claim to be verified. But to do all of this, the library needs support from the programming language to enable the creation of specialized data types. Further, these types are often parameterized (e.g., the bounds of a bounded float may not be specified until runtime).

Type safety. Strictly speaking, the library provides its own layer of type enforcement (e.g., chaining verifies type compatibility between chained functions). However, it may be beneficial to use a language that helps programmers write type-safe code, so languages that are strongly and statically typed may be preferable.

Functional features and “pure” functions. A number of functional paradigm features proved quite useful in prototyping, including: higher-order functions, function composition (chaining), and partial function evaluation. Some transformations, such as RowTransform, take an arbitrary, user-defined function and apply it to the data. To ensure privacy, the library must enforce that this function is free of side effects, as discussed in Section 8. This can be achieved by purifying user-defined functions (as is done in PINQ [8]) or requiring users to express computations in a restricted domain-specific language (as is done in Psi [4]).

Generics. We would like to allow the programmer to construct general-purpose transformations (or measurements) that are agnostic to the data type to which they are applied. An example where this would be useful is RowTransform: ideally one implements a single generic routine that applies the user-defined function to map a multiset of one type to a multiset of another type, where the types are generic and can be specified at compile/run time. This kind of behavior can be supported even in statically typed languages through the use of generics (parametric polymorphism). *[[SV: maybe can use elaboration or updating given our new treatment of families?]]* *[[MH: updated]]*

Structures or object-orientation. In our implementation, we took advantage of the object-oriented support of Python to define classes that allowed us to link together related functionality (e.g., a function and its privacy relation) into a single object that can be treated as a unit. But similar functionality is present in languages that are not considered object-oriented (e.g., Rust).

9.2 Proposal: Rust + Python

[[MH: should title change: Rust + bindings for Python/R?]] *[[SV: changed some of the pronouns and verb tenses to distinguish the work of the OpenDP Project more broadly from that of the authors of this paper]]*

We and the OpenDP Design Committee propose to implement the core library using Rust. Rust is a fast-growing, imperative functional language. It is both high performance and reliable with good memory safety and type safety. It also meets many of the desiderata identified previously: structures, generics, higher-order functions, and values are immutable by default. It has a strong ownership model that gives memory safety without the necessity of garbage collection. Collectively, this makes Rust code much easier to reason about concretely. For example, there are no dangling pointers, there are no memory races as there are in C++.³

Where Rust may fall short is in the programmability criteria, as it is not a widely used language (though interest in it is growing, and it has been ranked the “most loved” language on Stackoverflow, edging out Python for the last four years) and it can be challenging to program (at least in part because of the built-in safety features).

The OpenDP project has already been developing an end-to-end differential privacy system as a joint project with Microsoft, and for this project required a library to stand in while the core OpenDP library is not yet ready. The development team built a prototype library in Rust, using ideas from the PSI library [4] and some advances in the direction of those proposed in this paper. All of the developers were new to Rust,

³The Microsoft Security Response Center (MSRC) released reports stating that 70% of security vulnerabilities addressed through Microsoft security patches are memory safety issues that would have been avoided in Rust, <https://msrc-blog.microsoft.com/2019/07/18/we-need-a-safer-systems-programming-language/>, and explicitly suggests moving away from C++ to Rust <https://msrc-blog.microsoft.com/2019/07/22/why-rust-for-safe-systems-programming/>.

and all picked it up straightforwardly. Also, one graduate student who wished to contribute code from a new research paper was able to quickly port their code to Rust with some small direction.

Our proposal is that Rust should be the language for the robust, polished, production code that enters the maximally trusted OpenDP library. However, to encourage contributions from researchers and other developers, as well as to encourage experimentation and sandboxing of performance, the library should be able to work with Python and R code components. We propose to enable this through Python and R bindings (through a Protocol Buffer API described next). Initial code could be contributed in Python (and R), with the goal that final, vetted and performant code be ported to Rust, either by contributors when possible, or by the core OpenDP development team, or in collaboration. *[[SV: given the simplification of our vetting process described in Section 3, it seems that only code that directly constructs measurements or transformations actually needs to be in Rust for the sake of privacy]]*

APIs. Rust uses the same Foreign Function Interfaces as C++ for bindings so the bindings compiled from either Rust or C++ are identical. However, in the OpenDP project with Microsoft, the code to write these bindings has been found to be clearer than in C++. That project has created Python bindings to the prototype differential privacy library that have been heavily tested and work well. R bindings are in construction and have been found to be straightforward too.

In that project, an API to the library was built using Protocol Buffers⁴, a serialized data language, like JSON or XML, but heavily typed and very compact. This allowed for defining any sequence of calls to the library as a plan, in the style of ϵ ktelo or more specifically, a computational graph of library components. The protobuf API calls Rust components for the runtime. The Python and R bindings also call the API which calls the Rust library. In the case of Python, these bindings can be distributed on PyPI and pip installed, and executed in a Python notebook seemingly natively, although the actual execution is in Rust. The OpenDP Design Committee proposes to release a Protocol Buffer API for the core library as part of the OpenDP Commons, by which it will also offer Python and R bindings. This API would facilitate other possible bindings, as well a way to call the library in more complicated system architectures.

10 Using the Library

The programming framework proposed in this paper is meant to capture all the mathematical reasoning about privacy needed for use of differential privacy. However, we do not expect it to be directly used by end users; instead it is meant as a substrate to build larger, user-friendly end-to-end differential privacy systems. Here we discuss some aspects of how such a system should use the library envisioned here.

Calling the Library from a DP System.

1. The system should first determine, or elicit from its user, the following:
 - (a) The sensitive dataset on which differential privacy is going to be applied.
 - (b) The type of the dataset (e.g. is it a multiset of records, or social network graph, or a collection of tables). Do the records consist of a known set of attributes?
 - (c) The granularity of privacy protection. For example, if we want to protect privacy at the level of individual human beings, we need to know whether every individual human corresponds to at most one record in the multiset, at most one node in a social network graph, or at most one key in a collection of tables.
 - (d) The privacy measure to be used (e.g. pure, approximate, or concentrated DP), and the desired privacy-loss parameters.

⁴<https://developers.google.com/protocol-buffers/>

Regarding the type of the dataset, it is best for the system to make as few assumptions as possible about the data type, as long as the granularity of privacy can be well-specified, to reduce the chance that an incorrect assumption leads to a violation of privacy. For example, it is safer to assume that the dataset is a multiset of records of arbitrary type or a multiset of strings, rather than assuming it is a multiset of records that each consist of 10 floating-point numbers in the interval $[0, 1]$. Additional public, non-sensitive information about the dataset (e.g. a full schema) can typically be exploited later for greater utility without being used for the privacy calculus. Note that specifying the granularity of privacy requires keeping information in the datasets that one might be tempted to remove beforehand — in order to enforce privacy at the level of individuals in a database where an individual might contribute to multiple records, we need the data records to have a key that corresponds to distinct individuals.

2. Then the system should select an interactive measurement family from the library (typically some form of adaptive composition) that is compatible with the data type, granularity of privacy, and privacy measure, and use it to spawn a queryable on the sensitive dataset that will mediate *all* future access to the dataset.
3. If additional type information about the dataset is known (e.g. about the individual records, or even a more complex property like being a social network of bounded degree), then before applying any queries to the dataset, the library can be used to apply a stable “casting” transformation that ensures that the dataset actually has the specified properties. For example, for multisets we can apply RowTransform_f for a function f that maps any record that does not have the specified type to a default value with that type.
4. All queries to the dataset should be made through the initial queryable. The queryable automatically ensures that we stay within the initially specified privacy-loss budget. Deciding how to apportion the budget among current and future queries is for the system to decide, but we expect that statistical algorithms in the library will expose their utility in ways that assist a system in making such decisions (e.g. through a priori utility estimates provided analytically or through benchmarking). *[[SV: add pointer to statistics section of OpenDP whitepaper]]*

[[SV: do we need to talk again about linking multiple datasets? we have talked in one of the technical sections about how joins would be supported.]]

Exposing System Features to the Library. To make efficient and safe use of the library, it is also important to expose some system features to the library:

- Exposing the data-access model to the library. What means do the library functions have to read the sensitive data? If, for example, the data is only available in a streaming format, that will limit the choice of methods in the library that can be applied.
- Exposing the capabilities of the back-end data system to the library. Global data transformations, aggregates, sorting, and joins can be performed much more efficiently if they are done by an optimized database engine than if the data needs to be extracted from the database and computed on externally. Thus, there should be a way of identifying which transformations in the library can be done by the database system. On the other hand, this means that the resulting privacy properties also depend on the database engine faithfully and securely implementing the specified transformation, without observable side effects (including timing) that might leak sensitive data. More complex capabilities that the database may have include randomized procedures like subsampling.

[[SV: the function attribute of a transformation or measurement can include implementations in SQL, Spark, etc. need to have common datatypes are meaningful across the different environments, as well as ways of moving data between environments]]

- Defining the partition between secure and insecure storage and compute according to the threat model (to what might the privacy adversary have access), so that private data does not cross the “privacy barrier” until after the protections of differential privacy have been applied. Recall that in addition to the original sensitive dataset, the results of transformations, and the state of queryables need to be kept secret from the privacy adversary.

User interfaces. A system should also provide a user-friendly way for data custodians and data analysts without expertise in differential privacy to make effective use of the library. In particular, we envision that many existing DP systems can be built as layers on top of the library, using the certified code in the library for all the privacy analysis. Such interfaces could include:

- A SQL-like interface as in PINQ and ϵ ktelo, where queries are specified as plans that are compiled into DP measurements.
- A GUI for selecting from a predefined set of statistics and transformations, as in PSI.
- A Python-like notebook interface for issuing queries to the data, like in the OpenDP System being built in collaboration with Microsoft. *[[SV: future - add citation]]*

[[SV: add cites]] Another important aspect of the user interface is how accuracy and statistical confidence are exposed to users (so that analysts don’t draw incorrect conclusions due to the noise for privacy), and it is important that the algorithms in the library provide support for this. We defer discussion of this to the Statistics section of the OpenDP Whitepaper.

[[SV: below is sketch of PSI implementation using the library, in draft=1 mode only]] In Figures 27, 28, and 29 is pseudocode for how the PSI system could be implemented using library as described.

- The only privacy-critical commands are (meant to be) at the end of Figure 27, namely the `assert` command and the one where `queryable` is created, and at the beginning of Figure 29, where the depositor’s queryable is cloned to allow for per-analyst budgets.
- The code assumes all of the generalizations of the previous sections, *except* allowing multiple metrics or privacy measures to be associated with a single transformation or measurement (so the input and output metrics are attributes of the transformation). It also uses the fact that for transformations that are stable in the usual sense, we do not need to provide an `hint` function when chaining.
- `var_names` and `var_ranges` are not explicitly used in the code. However, they are used by the PSI interface to help the depositor and analysts and provide default ranges (which can be changed) to use for measurements that require them.
- To help the depositor or analyst select their batches of non-interactive measurements, the PSI GUI would be making extensive use of the library to convert between privacy-loss parameters and a priori measures of accuracy (fixing other auxiliary parameters of a measurement family), and also to specify in advance how much of the budget a batch would consume. None of this requires touching the dataset, so has no privacy implications, and is functionality offered for the sake of utility. To avoid having to actually construct the actual measurement and transformation objects during this query selection/-tuning phase, these could be helper functions associated with the measurement families (which again do not need to be inspected for privacy validation).
- PSI also allows for privacy amplification if the dataset is a random sample of a larger population, which is omitted from the pseudocode below. It seems like we should handle that similarly to the way we decide to handle trusting the database infrastructure to perform certain transformations. (Now the subsampling is trusted to have been done by the external data collection process.)
- The code also illustrates the value of being able to have transformations persist. Currently the new variables previously defined by the depositor and analyst are recomputed on every batch of queries

```

1  # version 1
2  def set_policy(input_domain,input_metric,output_measure,input_distance,output_distance,int_meas,data)
3      if isinstance(int_meas, InteractiveMeasurement):
4          if ((int_meas.input_domain == input_domain) and (int_meas.input_metric == input_metric) and (
5              int_meas.privacy_relation(input_distance,output_distance)):
6              return int_meas.function(data)
7          elif return('interactive measurement doesn't satisfy policy, try again')
8          elif return('need to specify an interactive measurement, try again')
9
10 # using v1 in the PSI example
11 queryable = set_policy(MultiSets(strings,n),change-one,approxDP,1,(epsilon,delta),AC_meas,x) # --- wi
12 queryable = set_policy(MultiSets(strings,n),change-one,approxDP,1,(epsilon,delta),overall_meas,x) # -
13
14 # what I don't like in the above is that there are multiple lines that have the dataset and that spec
15
16 # version 2: enforces that there is only one line where the policy is specified and the data is touch
17 def set_policy(input_domain,input_metric,output_measure,input_distance,output_distance,data):
18     def privacy_relation(d_in,d_out): return (d_in <= input_distance and d_out>=output_distance)
19     def function(x):
20         initial_state = TRUE
21         def eval(query, state):
22             if state:
23                 if isinstance(query, InteractiveMeasurement):
24                     if ((query.input_domain == input_domain) and (query.input_metric == input_metric) and (query
25                         query.privacy_relation(input_distance,output_distance)):
26                         return(function(x),FALSE)
27                     elif return('interactive measurement doesn't satisfy policy',state)
28                     elif return('need to specify an interactive measurement',state)
29                     elif return (state.query(query),state)
30 return InteractiveMeasurement(input_domain,input_metric,output_measure,privacy_relation,function).fun
31
32 # using v2 in the PSI example
33 # the next line should be the only time the dataset x is ever directly referenced. after that, only
34 global_policy_enforcer = set_policy(MultiSets(strings,n),change-one,approxDP,1,(epsilon,delta),x)
35 queryable=global_policy_enforcer.query(AC_meas) # --- will return 'interactive measurement doesn't sa
36 queryable=global_policy_enforcer.query(overall_meas) # --- will not return the error, and now can que

```

Figure 26: Functions to set and enforce a policy with one line of code

```

1  # Data Depositor set-up
2  #   use the PSI interface to set the following values
3
4  n = number of records in dataset # PSI assumes n is public and correct
5  x = pointer and credentials to access dataset (CSV format) in Dataverse
6  d = the number of variables supposedly in each record
7  var_names = a d-tuple of the names of the variables
8  var_types = a d-tuple of the types of the variables (e.g. int, float, bool, categorical)
9  var_ranges = a d-tuple of default ranges for each variable
10         (e.g. upper and lower bounds for numerical, bins for categorical)
11  epsilon, delta = approxDP privacy loss parameters
12  analysts = a description of the set of Dataverse accounts that will be
13         given their own per-analyst privacy budgets
14
15  # define initial row_transformation
16
17  row_domain = d-tuples of values matching var_types
18  input_metric = change-one
19
20  def parse(z: string)
21      return z parsed/cast into row_domain
22
23  parse_trans = MakeRowTransform(string,row_domain,parse)
24
25  # create an adaptive composition measurement where changing at most one row
26  # of input implies (epsilon,delta)-closeness of outputs
27  AC_meas = MakeAdaptiveComposition(MultiSets(row_domain,n),change-one,1,approxDP,(epsilon,delta))
28
29  # hint can be inferred because parse_trans is 1-stable
30  overall_meas = ChainingMT(AC_meas,parse_trans)
31
32  # check privacy properties and start mechanism
33  # only this block of code needs to be trusted!
34
35  assert (overall_meas.input_domain == MultiSets(strings,n) &&
36         overall_meas.input_metric == change-one &&
37         overall_meas.output_measure == approxDP &&
38         overall_meas.privacy_relation(1,(epsilon,delta)))
39
40  queryable = overall_meas.function(x)

```

Figure 27: PSI: set up by data depositor *[[MH: the “Make” function names are obsolete; I think MakeAdaptiveComposition is just AdaptiveComposition now.]]* *[[MH: here we’re using the simple definition of MakeRowTransform but now we’ve started including input domains/metrics and relations, so MakeRowTransform constructor probably needs more inputs.]]*

```

1 [continuing from Figure~\ref{fig:PSI-setup}]
2 # variables to keep track of the current set of variables (since depositor may create new ones)
3 # and mapping from original variables to current set
4 d_curr = d
5 var_names_curr = var_names
6 var_types_curr = var_types
7 var_ranges_curr = var_ranges
8 row_domain_curr = row_domain
9 var_mapping_curr = identity function on row_domain_curr
10
11 Repeat until depositor decides to leave rest of budget for future analysts:
12     print(queryable.query('remaining budget'))
13
14     # create any new variables
15     Repeat as long as depositor wants:
16         depositor defines a quadruple (f,vname,vtype,vrange)
17         where f is a DSL-defined (Transformer) function from row_domain_curr to vtype
18         row_domain_curr = row_domain_curr * vtype # * is concatenation/cross product
19         var_mapping_curr = var_mapping_curr * compose(f,var_mapping_curr)
20         var_types_curr = var_types_curr * vtype
21         var_ranges_curr = var_ranges_curr * vrange
22         d_curr = d_curr+1
23
24     var_transform = MakeRowTransform(row_domain,row_domain_curr,var_mapping_curr)
25
26     # select a batch of non-interactive measurements on subsets of variables
27     Depositor uses GUI to select
28     ((meas_1,S_1),(meas_2,S_2),...,(meas_m,S_m)) such that
29     each S_j is a sequence (i_1,...,i_k) of elements of
30     {1,...,d_curr} (selecting variables) and
31     meas_j.input_domain == MultiSets(var_types_curr[i_1]*var_types_curr[i_2]*...*var_types_curr[i_k]),
32     meas_j.noninteractive == true, and
33     meas_j.output_measure == approxDP
34
35     # convert these to measurements on all the variables in row_domain_curr
36     for j=1 to m:
37         def proj(z): return z[S_j]
38         proj_transform = MakeRowTransform(row_domain_curr,meas_j.input_domain,proj)
39         actualmeas_j = ChainingMT(meas_j,proj_transform)
40
41     # batch the measurements together using optimal composition of approx-DP
42     # and convert to measurements on the original variables
43     batch_measurement_curr=MakeOptimalComposition(actualmeas_1,...,actualmeas_j)
44     batch_measurement = ChainingMT(batch_measurement_curr,var_transform)
45
46     publish_for_world(queryable.query(batch_measurement))
47
48 depositor_state = [queryable,row_domain,d_curr,row_domain_curr,var_mapping_curr,
49                   var_types_curr,var_ranges_curr]
50
51 store depositor_state

```

Figure 28: PSI: releases by data depositor

```

1  authenticate user as member of the set analysts and as not having
2  made queries on this dataset before
3
4  # clone the state where the depositor finished off --- each analyst can
5  # consume the rest of the the depositor's budget.
6  # note that we are cloning a queryable, which is a dangerous operation!
7  # we need to do this because the current framework does not support
8  # non-colluding adversary models
9
10 [queryable,row_domain,d_curr,row_domain_curr,var_mapping_curr,
11 var_types_curr,var_ranges_curr]=clone(depositor_state)
12
13 # rest is the same as code for depositor releases,
14 # except results are stored for this analyst only
15
16 Repeat:
17     print(queryable.query('remaining budget'))
18
19     # create any new variables
20     Repeat as long as analyst wants:
21         analyst defines a quadruple (f,vname,vtype,vrange)
22         where f is a DSL-defined (Transformer) function from row_domain_curr to vtype
23         row_domain_curr = row_domain_curr * vtype # * is concatenation/cross product
24         var_mapping_curr = var_mapping_curr * compose(f,var_mapping_curr)
25         var_types_curr = var_types_curr * vtype
26         var_ranges_curr = var_ranges_curr * vrange
27         d_curr = d_curr+1
28
29     var_transform = MakeRowTransform(row_domain,row_domain_curr,var_mapping_curr)
30
31     # select a batch of non-interactive measurements on subsets of variables
32     Analyst uses GUI to select
33     ((meas_1,S_1),(meas_2,S_2),...,(meas_m,S_m)) such that
34     each S_j is a sequence (i_1,...,i_k) of elements of
35     {1,...,d_curr} (selecting variables) and
36     meas_j.input_domain == MultiSets(var_types_curr[i_1]*var_types_curr[i_2]*...*var_types_curr[i_k]),
37     meas_j.noninteractive == true, and
38     meas_j.output_measure == approxDP
39
40     # convert these to measurements on all the variables in row_domain_curr
41     for j=1 to m:
42         def proj(z): return z[S_j]
43         proj_transform = MakeRowTransform(row_domain_curr,meas_j.input_domain,proj)
44         actualmeas_j = ChainingMT(meas_j,proj_transform)
45
46     # batch the measurements together using optimal composition of approx-DP
47     # and convert to measurements on the original variables
48     batch_measurement_curr=MakeOptimalComposition(actualmeas_1,...,actualmeas_j)
49     batch_measurement = ChainingMT(batch_measurement_curr,var_transform)
50
51     publish_for_analyst_only(queryable.query(batch_measurement))

```

Figure 29: PSI: interactive queries with per-analyst budgets

11 Contributing to the Library

Building on the discussion on Section 3, we elaborate here on the different types of code contributions we envision to the library.

1. A new measurement or transformation family constructed by combining existing measurement and transformation families in the library using operators that already exist in the library (like composition, chaining, and post-processing).

This is the easiest kind of contribution, as the certificates of privacy or stability are automatically generated by the built-in components, and we also expect it to be the most frequent contribution once the library has a large enough base of building-block components. Although no proof of privacy or stability is needed for such contributions, they will still need to be reviewed for having demonstrated utility (see the Statistics section of the OpenDP Whitepaper), code quality, and documentation.

2. A new “atomic” measurement or transformation family, where the privacy or stability properties are proven “from scratch.”

For the sake of modularity and ease of verification, such contributions should be broken down into as small and simple sub-measurements or sub-transformations as possible, each of which has a proof from scratch. Together with each such a contribution, one must also provide a mathematical proof that it actually achieves its claimed privacy or stability properties. One day, we hope that some of these proofs can be formally verified using the rapidly advancing set of tools for formal verification of differential privacy. [\[\[SV: future - add citations\]\]](#) However, in many cases, we expect that what will be provided is a written proof together with a pseudocode description of the algorithm and its claimed privacy properties. The DP researchers on the OpenDP editorial board will be responsible for verifying that the proof is correct for the pseudocode. The OpenDP Committers, with expertise on the programming framework and its implementation, will be responsible for verifying that the actual code is a faithful implementation of the pseudocode.

Sometimes the privacy or stability proof of a new contribution will be based on some property of an existing component in the library that is not captured by the privacy or stability properties of that component. For example, “stability methods” for approximate differential privacy (like “propose–test–release”) often rely on the *accuracy* of other DP mechanism (e.g. that with probability at least $1 - \delta$, the Laplace mechanism produces an estimate that is correct to within $c \log(1/\delta)/\epsilon$). Whenever this occurs, there should also be a proof provided that the said component has the given property, this proof should be attached as a certified assertion to the implementation of that component, and the privacy or stability relation for the new contribution should check the certificate. This way, if the component later changes and the assertion has to be removed, the soundness of the contribution relying on it is preserved.

3. A new combining primitive for measurements or transformations.

This can be a new form or analysis of composition or chaining, or something more sophisticated like privacy amplification by subsampling (viewed as an operator on measurements). In the proposed framework, these are also measurement or transformation families and the process for contributing them and having them vetted and accepted is the same as above.

4. Adding a new private data type, distance measure between private data types, or privacy notion (distance measure between distributions).

There should be a standardized format for adding such elements to the library, and what information needs to be provided and checked (again by a combination of the editorial board and committers). For example, for a new distance measure, one may provide a statement of the data types to which it applies, specify whether it satisfies properties like monotonicity and symmetry, give conversions between it and other distance measures, etc. For a new privacy notion, it needs to be verified that it satisfies the post-processing property.

5. Adding a new type of privacy or stability calculus that is not supported by the existing framework.

Examples include introducing frameworks for privacy odometers, local sensitivity and variants (like smooth sensitivity), or per-analyst privacy budgets. These will require a more holistic architectural review before being incorporated.

12 The Scope of the Framework

[[SV: future- lots of citations missing in this section]]

The programming framework presented in Sections 2–7 is meant to be very general and flexible, and we believe it can support a great deal of the privacy calculus in the differential privacy literature. However, it does not support all important ways of reasoning about differential privacy. Thus, below we list some concepts that we believe can be easily supported and others that would require an expansion of the framework.

12.1 Within the Current Framework

Many different private dataset types and granularities of privacy. This includes unbounded DP (multisets, adjacency by adding or removing a record), bounded DP (datasets consisting of a known number of records, adjacency by changing a record), networks with edge-level or node-level privacy, a data stream with person-level or event-level privacy, multirelational databases with person-level privacy, DP under continual observation.

Many different measures of privacy loss. This includes pure DP, approximate DP, zero-concentrated DP, Rényi, and Gaussian DP. (We require the privacy-loss measure to be closed under postprocessing, so the original formulation of concentrated DP would not be supported.) Group privacy is also captured by the proposed privacy relations, as it allows for interrogating about the privacy loss at inputs of different distances.

Combining primitives to build complex mechanisms from simpler ones. This includes many of the sophisticated composition theorems for differential privacy (such as the advanced and optimal composition theorems for approximate differential privacy, the composition of zCDP, and more), “parallel composition” via partitioning, privacy amplification by subsampling, the sample-and-aggregate framework.

Global-sensitivity-based mechanism families. The exponential mechanism, the geometric mechanism, the (discrete) Gaussian mechanism, the Matrix Mechanism, and the Sparse Vector Technique.

Common database transformations. Subsetting according to a user-defined predicate, partitioning according to a user-defined binning, clamping into a given range/ball, vectorizing into a histogram, join (or variant) on a multi-relational dataset, user-defined row-by-row transformations, and summing over a dataset.

Restricted sensitivity and Lipschitz projections. Both of these refer to (global) stability properties of transformations with respect to rich data types, which the framework supports.

12.2 Outside the Current Framework

Below are some general concepts in differential privacy that the proposed framework does not currently support reasoning about. When these concepts are used only as an intermediate step to proving a standard privacy property, then the entire mechanism can be added to the library as a measurement. But the framework would not support reasoning directly these concepts without further augmentation.

Local sensitivity and mechanisms based on it. We envision that the local sensitivity of a transformation could be captured by adding a “local stability” relation to the attributes of a transformation, which takes a dataset as an additional argument. Such a relation would also need to be verified similarly to the already-proposed stability relations. By appropriate augmentations like this to the framework, we hope it will be possible to support reasoning about differentially private mechanisms based on variants of local sensitivity, such as smooth sensitivity, stability-based mechanisms, and propose-test release.

Privacy odometers. These are variants of interactive measurements that track accumulated privacy loss rather than specify a total privacy loss at the beginning [12]. As discussed in Section 5, it should be relatively straightforward to modify the programming framework to incorporate at least pure-DP odometers.

Complex adversary models. These include pan-privacy with a limited number of intrusions, computational differential privacy, adversaries with limited collusion (e.g. when we give different analysts budgets of their own)

Randomized transformations. It is possible and sometimes useful to reason about stability properties of randomized transformations, often via couplings. For example, we can call a randomized transformation $T : \text{MultiSets}(\mathcal{X}) \rightsquigarrow \text{MultiSets}(\mathcal{Y})$ *c-stable* if for every two adjacent datasets $x \sim x'$, there is a coupling (Y, Y') of the random variables $Y = T(x)$ and $Y' = T(x')$ such that it always holds that $d_{\text{Sym}}(Y, Y') \leq c$. The framework should be easily extendable to this way of capturing the stability properties of randomized transformations, but more complex coupling properties may be more challenging to capture.

Interactive transformations. In some DP systems, like PINQ, transformations can also be performed in an interactive and adaptive manner, creating a multitude of derived datasets (behind the “privacy barrier”) on which queries can subsequently be issued. The framework proposed so far only allows for measurements to be interactive, which may mean that the same derived dataset is recomputed multiple times if it is an intermediate dataset in multiple measurements formed by chaining. This can be potentially remedied by modifying the execution engine to retain previously derived datasets and recognize when they are being recreated. However, it may eventually be useful to extend the programming framework to directly work with some form of interactive transformations, just as we allow interactive measurements. (But a general theory of interactive transformations and how they may be combined with measurements may need to be developed first.)

13 Comparison with other Programming Frameworks

- Comparison with existing programming frameworks: Michael or Marco?

Maybe Marco’s notes from the design cmtc are useful? <https://docs.google.com/document/d/1hATwn0DQPss5Hyw6ipE-v-91GQUX1Dw20x8n13pftuU/edit?usp=sharing>

[[SV: Commented out old material from Fall discussions]]

14 Representing Data Domains

```
1 class DataDomain:
2     name          # e.g. Float, BoundedFloat, MultiSet
3     parameters    # bounds, the type of elements it contains, etc.
4     valid_item    # tests whether a given object is in this data domain
5     default_item  # a default element of the data domain
6     samedomain    # certifies that another DataDomain object represents the same domain
7
8 def MakeFloat():
9     name = 'Float'
10    parameters = None
11    def valid_item(x): return isinstance(x,float)
12    default_item = 0.0
13    def samedomain(x : DataDomain):
14        return (x.name == 'Float')
15    return DataDomain(name,parameters,valid_item,default_item,samedomain)
16
17 def MakeBoundedFloat(lower : float,upper : float):
18     if lower > upper: raise Exception
19     name = 'BoundedFloat'
20     parameters = [lower,upper]
21     def valid_item(x):
22         if not(isinstance(x,float)):
23             return False
24         elif (x>upper) or (x<lower):
25             return False
26     default_item = (lower+upper)/2
27     def samedomain(x : DataDomain):
28         return (x.name == 'BoundedFloat' and x.parameters == self.parameters)
29     return DataDomain(name,parameters,valid_item,default_item,samedomain)
30
31 def MakeMultiSet(rowdomain : DataDomain):
32     name = 'MultiSet'
33     parameters = rowdomain
34     def valid_item(x):
35         if isinstance(x,list):
36             for y in x:
37                 if not(rowdomain.valid_item(y)): return False
38             return True
39         elif return False
40     default_item = []
41     def samedomain(x : DataDomain):
42         return (x.name == 'MultiSet' and
43             self.parameters.samedomain(x.parameters)) # note nontrivial equality check!
44     return DataDomain(name,parameters,valid_item,default_item,samedomain)
```

Figure 30: DataDomains

```

1 class Measurement:
2     input_metric
3     input_domain : DataDomain
4     privacy_loss
5     function
6
7 class Transformation:
8     input_metric
9     input_domain : DataDomain
10    output_metric
11    output_domain : DataDomain
12    stability
13    function
14
15 def ChainingMT(meas: Measurement, trans: Transformation):
16     if (not(trans.output_domain.samedomain(meas.input_domain))
17         or trans.output_metric!=meas.input_metric): raise Exception('Domain/metric mismatch')
18     input_metric = trans.input_metric
19     input_domain = trans.input_domain
20     privacy_loss = trans.stability*meas.privacy_loss
21     def function(data): return meas.function(trans.function(data))
22     return Measurement(input_metric,input_domain,privacy_loss,function)
23
24 def MakeBaseLap(sigma: float):
25     if sigma < 0: raise Exception('Invalid parameter')
26     input_metric = dist_real
27     input_domain = MakeFloat()
28     privacy_loss = 1/sigma
29     def function (data): return data + Laplace(sigma,0)
30     return Measurement(input_metric,input_domain,privacy_loss,function)
31
32 def MakeBoundedSum(L: float, U: float):
33     if L > U: raise Exception('Invalid parameters')
34     input_metric = dist_sym
35     input_domain = MakeMultiSet(MakeBoundedFloat(L,U))
36     output_metric = dist_real
37     output_domain = MakeFloat()
38     stability = max(abs(L), abs(U))
39     def function (data): return sum(data)
40     return Transformation(input_metric,input_domain,output_metric,output_domain,stability,function)
41
42 # Example
43 BaseLaplace = MakeBaseLap(sig) # sig is some constant
44 BoundedSum = MakeBoundedSum(l, u) # l,u defined previously
45 NoisySum=ChainingMT(BaseLaplace, BoundedSum)
46 print(NoisySum.privacy_loss) # prints max(abs(l),abs(u))/sig

```

Figure 31: Using DataDomains

```

1  # version 1: requires data domains to have a "valid" and "default_item" methods
2
3  def MakeUnstableTransformation(in_dom,out_dom,f : carrier(in_dom)->carrier(out_dom),T):
4      if in_dom.valid == None or in_dom.default_item == None then: RaiseException
5          input_domain = in_dom
6          output_domain = out_dom
7          stability = infinity
8          def g(x):
9              y=SafeEval(f,T,x) # SafeEval(f,T) evaluates f(x) for exactly T time steps, returns a NaN (missing d
10             if out_dom.valid(y): return y
11             else: return out_dom.default_item
12         return Transformation(in_dom,out_dom,stability,g)
13
14 def MakeRowTransform(T : Transformation):
15     input_domain = MultiSets(T.input_domain)
16     output_domain = MultiSets(T.output_domain)
17     stability = 1
18     def function(data):
19         [T.function(x) | for x in data]
20     return Transformation(input_domain,output_domain,stability,function)
21
22 # version 2: requires data domains to have an "unrestricted" attribute - meaning any element of the car
23
24 def MakeUnstableTransformation(in_dom,out_dom,f : carrier(in_dom)->carrier(out_dom),T):
25     if not(in_dom.unrestricted and out_dom.unrestricted) then: RaiseException('only for unrestricted doma
26     input_domain = in_dom
27     output_domain = out_dom
28     stability = infinity
29     def g(x):
30         return SafeEval(f,T,x) # SafeEval(f,T) evaluates f(x) for exactly T time steps, returns a NaN (miss
31     return Transformation(in_dom,out_dom,stability,g)
32
33 def Clamp(input_domain,output_domain):
34     if input_domain==Float and output_domain==BoundedFloat:
35         def function(x): return min(max(x,output_domain.lower_bound),output_domain.upper_bound)
36     else if input_domain==Int and output_domain==BoundedInt:
37         def ...
38     else: RaiseException('Clamp does not support these domains')
39     stability = infinity
40     return Transformation(input_domain,output_domain,function,stability)
41
42 def Impute(input_domain,output_domain):
43     if input_domain==FloatWithNaN and output_domain==FloatWithoutNaN:
44         def function(x):
45             if x==NaN: return 0.0 # this is the 'default item'
46             else return x
47     else if input_domain==IntWithNaN...
48     else: RaiseException('Impute does not support these domains')
49     stability = infinity
50     return Transformation(input_domain,output_domain,function,stability)

```

Figure 32: Ideas for MakeRowTransform

```

1 class DataMetric:
2     name          # e.g. dist-float, L1, L2, add-remove
3     domain        # a *list* [dom_1,dom_2,...,dom_k] s.t.
4                   # the metric applies to dom_1*dom_2*...*dom_k
5     parameters    # e.g. weights, or metrics on dom_i's
6     same_metric   # tests whether another DataMetric is the same
7
8     # the following are optional parameters
9     single_real   # is the distance given by a single nonnegative real parameter that is
10                  # monotonic (if x,y d-close and d'>d, then x,y are d'-close)?
11                  # and if x,y are 0-close then x==y
12     integer       # if single_real, are the distances integers?
13     upper_bound   # if single_real, upper bound on distances
14     triangle      # if single_real, is the triangle inequality satisfied?
15     path          # if integer and triangle, can distances
16                  # always be realized by a path of distance 1 steps?
17     symmetric     # x,y d-close implies y,x d-close
18
19 # ordinary distance between floats:
20 # x and y are d-close for a float d if |x-y|<=d
21 def MakeDistFloat():
22     name = 'Dist-Float'
23     domain = [MakeFloat()]
24     parameters = None
25     def same_metric(dist): return (dist.name=='Dist-Float')
26     single_real = True
27     integer = False
28     upper_bound = None
29     triangle = True
30     path = False
31     symmetric = True
32     return DataMetric(name,domain,parameters,same_metric,
33                       single_real,integer,upper_bound,triangle,path,symmetric)

```

Figure 33: DataMetrics

```

1  # ordinary distance between bounded floats:
2  # x and y are d-close for a float d if |x-y|<=d
3  def MakeDistBoundedFloat(L,U):
4      name = 'Dist-BoundedFloat'
5      domain = [MakeBoundedFloat(L,U)]
6      parameters = [L,U]
7      def same_metric(dist : DataMetric):
8          return (dist.name=='Dist-BoundedFloat' and dist.parameters==[L,U])
9      single_real = true
10     integer = false
11     upper_bound = |U-L|
12     triangle = true
13     path = false
14     symmetric = true
15     return DataMetric(name,domain,parameters,same_metric,
16         single_real,integer,upper_bound,triangle,path,symmetric)
17
18 # symmetric difference between multisets:
19 # x and y are d-close if the symmetric difference between x and y is at most d
20
21 def MakeDistAddRemove(rowdomain):
22     name = 'Dist-AddRemove'
23     domain = [MakeMultiSet(rowdomain)]
24     parameters = rowdomain
25     def same_metric(dist : DataMetric):
26         return (dist.name=='Dist-AddRemove' and self.parameters.samedomain(dist.parameters))
27     single_real = true
28     integer = true
29     upper_bound = None
30     triangle = true
31     path = true
32     symmetric = true
33     return DataMetric(name,domain,parameters,same_metric,
34         single_real,integer,upper_bound,triangle,path,symmetric)

```

Figure 34: More DataMetrics

```

1  # combines two single_real metrics to get sum metric on ordered pairs
2  # i.e. the distance between (x,y) and (x',y') is the sum
3  # of the distances between x and x and between y and y'
4  # this is the metric on the output of a Partition
5
6  def MakeSumMetric(metric1 : DataMetric, metric2: DataMetric):
7      if not(metric1.single_real and metric2.single_real): raise Exception
8      name = 'Dist-Sum'
9      domain = metric1.domain+metric2.domain
10     parameters = [metric1,metric2]
11     def same_metric(dist : DataMetric):
12         return (dist.name == 'Dist-Sum' and self.parameters[0].same_metric(dist.parameters[0])
13             and self.parameters[1].same_metric(dist.parameters[1]))
14     single_real = true
15     integer = (metric1.integer and metric2.integer)
16     upper_bound = metric1.upper_bound+metric2.upperbound # if neither is None
17     triangle = (metric1.triangle and metric2.triangle)
18     path = (metric1.path and metric2.path)
19     symmetric = (metric1.symmetric and metric2.symmetric)
20     return DataMetric(name,domain,parameters,same_metric,
21         single_real,integer,upper_bound,triangle,path,symmetric)
22
23     # combines two metrics to get a metric on ordered pairs
24     # where (x,y) is (c,d)-close to (x',y') if
25     # x is c-close to x' AND y is d-close to y'
26     # this is what we get if we pair the result of two transformations
27     # and is what's needed as the input to joins
28
29     def MakeAndMetric(metric1 : DataMetric, metric2: DataMetric):
30         name = 'Dist-And'
31         domain = metric1.domain+metric2.domain
32         parameters = [metric1,metric2]
33         def same_metric(dist : DataMetric):
34             return (dist.name == 'Dist-And' and self.parameters[0].same_metric(dist.parameters[0])
35                 and self.parameters[1].same_metric(dist.parameters[1]))
36         single_real = false
37         integer = false
38         upper_bound = None
39         triangle = false
40         path = false
41         symmetric = (metric1.symmetric and metric2.symmetric)
42         return DataMetric(name,domain,parameters,same_metric,
43             single_real,integer,upper_bound,triangle,path,symmetric)

```

Figure 35: DataMetrics on Pairs

```

1 class PrivacyMetric:
2     name          # e.g. pure-DP, approx-DP, zCDP
3     parameters    # e.g. threshold for truncated zCDP
4     same_metric   # tests whether another PrivacyMetric is the same
5
6     # the following are optional parameters
7     single_real   # is the distance given by a single nonnegative real parameter that is
8                   # monotonic (if x,y d-close and d'>d, then x,y are d'-close)?
9                   # and if x,y are 0-close then x==y
10    upper_bound   # if single_real, upper bound on distances
11    triangle      # if single_real, is the triangle inequality satisfied?
12    symmetric     # x,y d-close implies y,x d-close
13
14 def MakePureDP():
15     name = 'PureDP'
16     parameters = None
17     def same_metric(dist): return (dist.name=='PureDP')
18     single_real = True
19     upper_bound = None
20     triangle = True
21     symmetric = False
22     return PrivacyMetric(name,parameters,same_metric,
23                           single_real,upper_bound,triangle,symmetric)
24
25 def MakeApproxDP():
26     name = 'ApproxDP'
27     parameters = None
28     def same_metric(dist): return (dist.name=='ApproxDP')
29     single_real = False
30     upper_bound = None
31     triangle = False
32     symmetric = False
33     return PrivacyMetric(name,parameters,same_metric,
34                           single_real,upper_bound,triangle,symmetric)

```

Figure 36: PrivacyMetrics

```

1 class Measurement:
2     input_domain # *list* of DataDomains
3     input_metric : DataMetric # should apply to input_domain
4     privacy_metric : PrivacyMetric # e.g. pure-dp, approx-dp, etc.
5     privacy_constant # only allowed if both metrics are single_real
6     privacy_relation # automatically constructed if privacy_constant given
7     function
8
9 class Transformation:
10    input_domain # *list* of DataDomains
11    input_metric : DataMetric on input_domain
12    output_domain # *list* of DataDomains
13    output_metric : DataMetric on output_domain
14    stability_constant # only allowed if both metrics are single_real
15    stability_relation # automatically constructed if stability_constant given
16    function
17
18 def MakeBaseLap(sigma: float):
19     if sigma < 0: raise Exception('Invalid parameter')
20     input_metric = MakeDistFloat()
21     input_domain = [MakeFloat()]
22     privacy_metric = MakePureDP()
23     def privacy_constant = 1/sigma
24     def function (data): return [data + Laplace(sigma,0)]
25     return Measurement(input_domain,input_metric,privacy_metric,privacy_constant,None,function)
26
27 def MakeBoundedSum(L: float, U: float):
28     if L > U: raise Exception('Invalid parameters')
29     rowdomain=MakeBoundedFloat(L,U)
30     input_metric = MakeDistAddRemove(rowdomain)
31     input_domain = [MakeMultiSet(rowdomain)]
32     output_metric = MakeDistFloat()
33     output_domain = [MakeFloat()]
34     stability_constant = max(abs(L),abs(U))
35     def function (data): return [sum(data)]
36     return Transformation(input_domain,input_metric,output_domain,output_metric,stability_constant,None)

```

Figure 37: Using DataMetrics and PrivacyMetrics

```

1 def ChainingMT(meas: Measurement, trans: Transformation, hint):
2     if (not(equalequaldomainlists(trans.output_domain,meas.input_domain)
3         or trans.output_metric.same_metric(meas.input_metric))): raise Exception('Domain/metric mismatch')
4     input_domain = trans.input_domain
5     input_metric = trans.input_metric
6     privacy_metric = meas.output_metric
7     mid_metric = trans.output_metric
8     if (hint==None and mid_metric.single_real==false): raise Exception('need hint')
9
10    def function(data): return meas.function(trans.function(data))
11
12    if (trans.stability_constant!=None and meas.privacy_constant!=None):
13        privacy_constant = trans.stability_constant*meas.privacy_constant
14        return Measurement(input_domain,input_metric,privacy_constant,None,function)
15
16    if (hint==None and trans.stability_constant!=None):
17        def hint(d_in,d_out) = d_in*trans.stability_constant
18
19    if (hint==None and meas.privacy_constant!=None):
20        def hint(d_in,d_out) = d_out/meas.privacy_constant
21
22    if (hint!=None):
23        def privacy_relation(d_in,d_out):
24            d_mid = hint(d_in,d_out)
25            return (trans.stability_relation(d_in,d_mid)
26                and meas.privacy_relation(d_mid,d_out))
27    else:
28        def privacy_relation(d_in,d_out): # binary search over d_mid in the interval [0,U],
29            # where U is either upper_bound or the largest float
30            # return true if find d_mid s.t. trans.stability_relation(d_in,d_mid) and meas.privacy_relation(d
31
32    return Measurement(input_domain,input_metric,privacy_metric,None,privacy_relation,function)
33
34 # Example
35 BaseLaplace = MakeBaseLap(sig) # sig is some constant
36 BoundedSum = MakeBoundedSum(l, u) # l,u defined previously
37 NoisySum=ChainingMT(BaseLaplace, BoundedSum)
38 print(NoisySum.privacy_loss) # prints max(abs(l),abs(u))/sig

```

Figure 38: Chaining with DataMetrics

```

1 def MakePartition(rowdomain : DataDomain, predicate):
2     input_domain = [MakeMultiSet(rowdomain)]
3     input_metric = MakeDistAddRemove(rowdomain)
4     output_domain = [input_domain,input_domain] # output is a pair of datasets
5     output_metric = MakeSumMetric(input_metric,input_metric)
6     stability_constant = 1
7     def function(x):
8         z=[[ ],[ ]]
9         for y in x:
10            if predicate(y): z[0]=z[0]+[y]
11            else: z[1]=z[1]+[y]
12        return z
13    return Transformation(input_domain,input_metric,output_domain, output_metric,stability_constant,none,)
14
15 def ParallelComposition(meas0 : Measurement, meas1 : Measurement):
16     if not(meas0.input_metric.single_real and meas1.input_metric.single_real): raise Exception('inapplicabl
17     if not(meas0.privacy_metric.same_metric(meas1.privacy_metric)): raise Exception('mismatched metrics')
18
19     input_domain = [meas0.input_domain,meas1.input_domain]
20     input_metric = MakeSumMetric(meas0.privacy_metric,meas1.privacy_metric)
21     privacy_metric = meas0.privacy_metric
22
23     def function(x): return [meas0.function(x[0]),meas1.function(x[1])]
24     # --- function will be more involved with interactive measurements.
25     # --- need to specify sequential vs. concurrent composition
26
27     # --- for pure DP, we can just take max of privacy constants
28     if (meas0.privacy_constant!=None and meas1.privacy_constant!=None):
29         privacy_constant=max(meas0.privacy_constant,meas1.privacy_constant)
30         return Measurement(input_domain,input_metric,privacy_metric,privacy_constant,None,function)
31
32     # --- but for approx DP, we restrict to path metrics on input
33     # --- and apply group privacy
34     if (meas0.input_metric.path and meas1.input_metric.path):
35         def privacy_relation(d_in,d_out): return (d_in<1 or (d_in>=1 and meas0.privacy_relation(1,d_out) and
36         return GroupPrivacy(Measurement(input_domain,input_metric,privacy_metric,None,privacy_relation))
37
38     raise Exception('inapplicable input metric/privacy metric combo')
39
40 def GroupPrivacy(meas : Measurement):
41     if not(meas.input_metric.path): raise Exception('only path metrics')
42
43     input_domain = meas.input_domain
44     input_metric = meas.input_metric
45     privacy_metric = meas.privacy_metric
46     function = meas.function
47
48     if meas.privacy_metric.name='PureDP':
49         privacy_constant = smallest d s.t. meas.privacy_relation(1,d) # binary search
50         return Measurement(input_domain,input_metric,privacy_metric,privacy_constant,None,function)
51
52     if meas.privacy_metric.name='ApproxDP':
53         def privacy_relation(d_in,d_out):
54             d_in = floor(d_in)
55             eps = d_out[0]/d_in
56             delta = d_out[1]*(exp(eps)-1)/(exp(d_out[0])-1)
57             return meas.privacy_relation(1, [eps,delta])
58         return Measurement(input_domain,input_metric,privacy_metric,None,privacy_relation,None,function)
59
60     raise Exception('only for Pure and Approx DP')

```

Notes on DataDomains:

- Commented out: a subset testing function for data domains (and allowed chaining when the output domain of the first function is a subset of the input domain of the second). Dropped because I think it's better to keep the description of data domains as simple as possible. Increasing the domain can be achieved by transformations, for example the the identity/inclusion map from BoundedFloat([L,U]) to Float is a 1-stable transformation with respect to dist-real. Applying this transformation allows us to cast a BoundedFloat as a Float.
- important correction 8/4: nontrivial samedomain check for MultiSet DataDomain
- Potentially omit: valid-item and default-item for DataDomains
- Potentially add: equality check for elements of a DataDomain, comparator for ordered DataDomains (to support generic sorting), relation to actually check distance for a DataMetric
- I've added binary search for metrics that are single-real.
- Changed order of domains and metrics in last block of code

References

- [1] V. Balcer and S. Vadhan. Differential Privacy on Finite Computers. *Journal of Privacy and Confidentiality*, 9(2), September 2019. Special Issue on TDP 2017. Preliminary versions in ITCS 2018 and posted as arXiv:1709.05396 [cs.DS].
- [2] M. Bun and T. Steinke. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In M. Hirt and A. D. Smith, editors, *Theory of Cryptography - 14th International Conference, TCC 2016-B, Beijing, China, October 31 - November 3, 2016, Proceedings, Part I*, volume 9985 of *Lecture Notes in Computer Science*, pages 635–658, 2016.
- [3] H. Ebadi and D. Sands. Featherweight pinq. *Journal of Privacy and Confidentiality*, 7(2), 2016.
- [4] M. Gaboardi, J. Honaker, G. King, J. Murtagh, K. Nissim, J. Ullman, and S. Vadhan. Psi ($\{\Psi\}$): a private data sharing interface. *arXiv preprint arXiv:1609.04340*, 2016.
- [5] N. Johnson, J. P. Near, and D. Song. Towards practical differential privacy for sql queries. *Proceedings of the VLDB Endowment*, 11(5):526–539, 2018.
- [6] N. M. Johnson, J. P. Near, J. M. Hellerstein, and D. Song. Chorus: Differential privacy via query rewriting. *CoRR*, abs/1809.07750, 2018.
- [7] I. Kotsogiannis, Y. Tao, X. He, M. Fanaeepour, A. Machanavajjhala, M. Hay, and G. Miklau. Privatesql: a differentially private sql query engine. *Proceedings of the VLDB Endowment*, 12(11):1371–1384, 2019.
- [8] F. D. McSherry. Privacy integrated queries: an extensible platform for privacy-preserving data analysis. In *Proceedings of the 2009 ACM SIGMOD International Conference on Management of data*, pages 19–30. ACM, 2009.
- [9] I. Mironov. On significance of the least significant bits for differential privacy. In *Proceedings of the 2012 ACM conference on Computer and communications security*, pages 650–661, 2012.
- [10] C. Palamidessi and M. Stronati. Differential privacy for relational algebra: Improving the sensitivity bounds via constraint systems. *Electronic Proceedings in Theoretical Computer Science*, 85:92–105, Jul 2012.

- [11] J. Reed and B. C. Pierce. Distance makes the types grow stronger: a calculus for differential privacy. In P. Hudak and S. Weirich, editors, *Proceeding of the 15th ACM SIGPLAN international conference on Functional programming, ICFP 2010, Baltimore, Maryland, USA, September 27-29, 2010*, pages 157–168. ACM, 2010.
- [12] R. Rogers, A. Roth, J. Ullman, and S. Vadhan. Privacy odometers and filters: Pay-as-you-go composition. In *Advances in Neural Information Processing Systems 29 (NIPS '16)*, pages 1921–1929, December 2016. Full version posted as arXiv:1605.08294 [cs.CR].
- [13] D. Zhang, R. McKenna, I. Kotsogiannis, M. Hay, A. Machanavajjhala, and G. Miklau. Ektelo: A framework for defining differentially-private computations. In *Proceedings of the 2018 International Conference on Management of Data*, pages 115–130. ACM, 2018.