

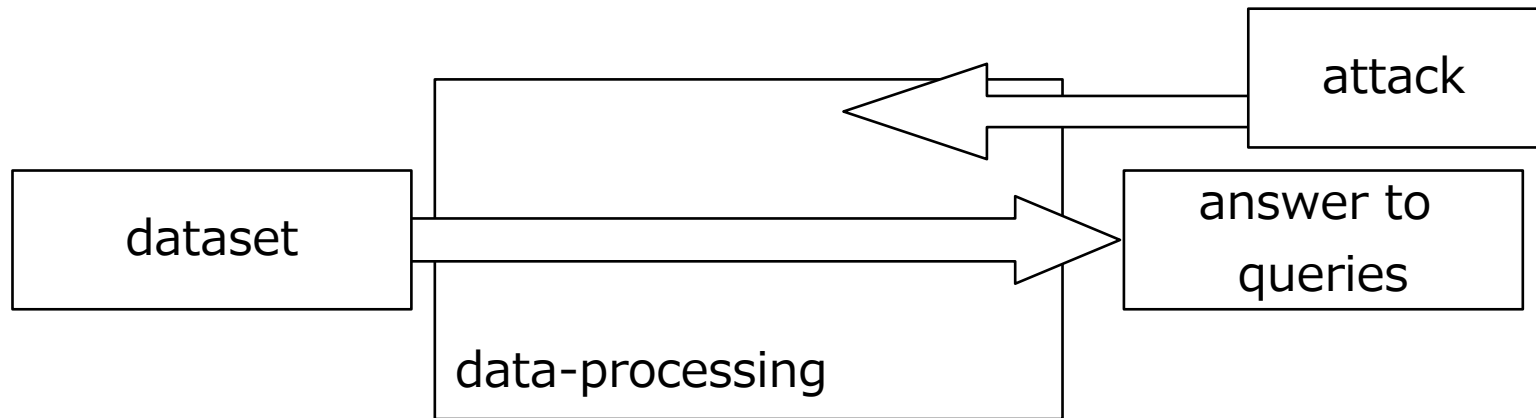
# **Towards Formal Verification of Differential Privacy in Isabelle/HOL**

8th Franco-Japanese Cybersecurity Workshop  
November 29, 2023

Tetsuya Sato(Tokyo Institute of Technology)

(special thanks: Shin-ya Katsumata • Yasuhiko Minamide)

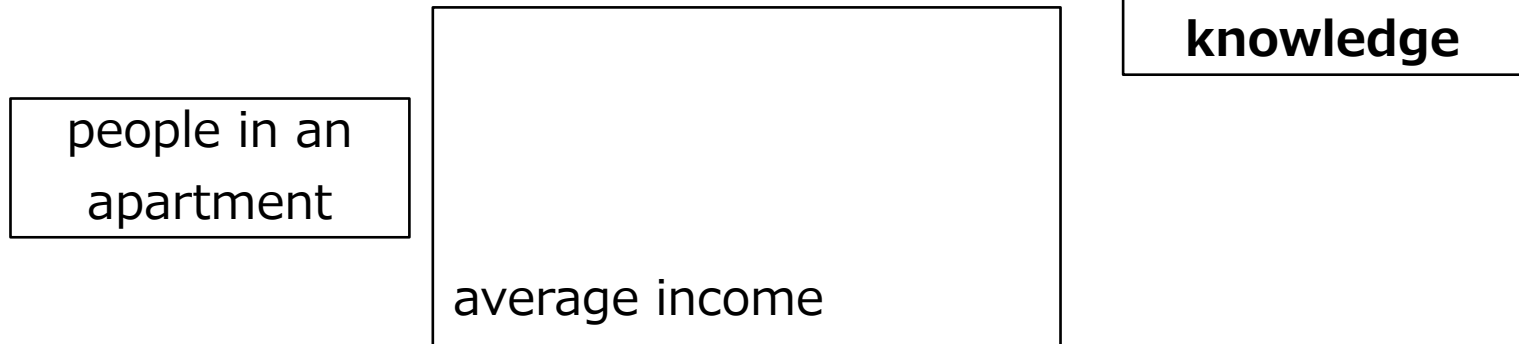
# Background



- Consider a query that processes datasets.  
The dataset contains individuals' private data.
- Even if the query does not show the private data, attacker can steal them from the answers if the attackers have enough background-knowledge.

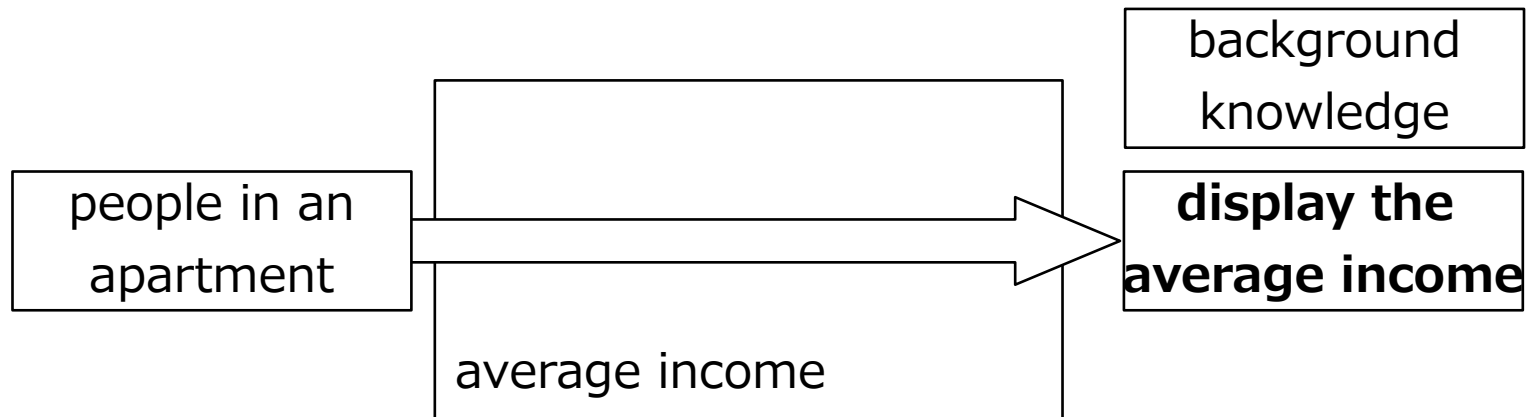
# An Example of Privacy Leakage

- consider a query of average income...



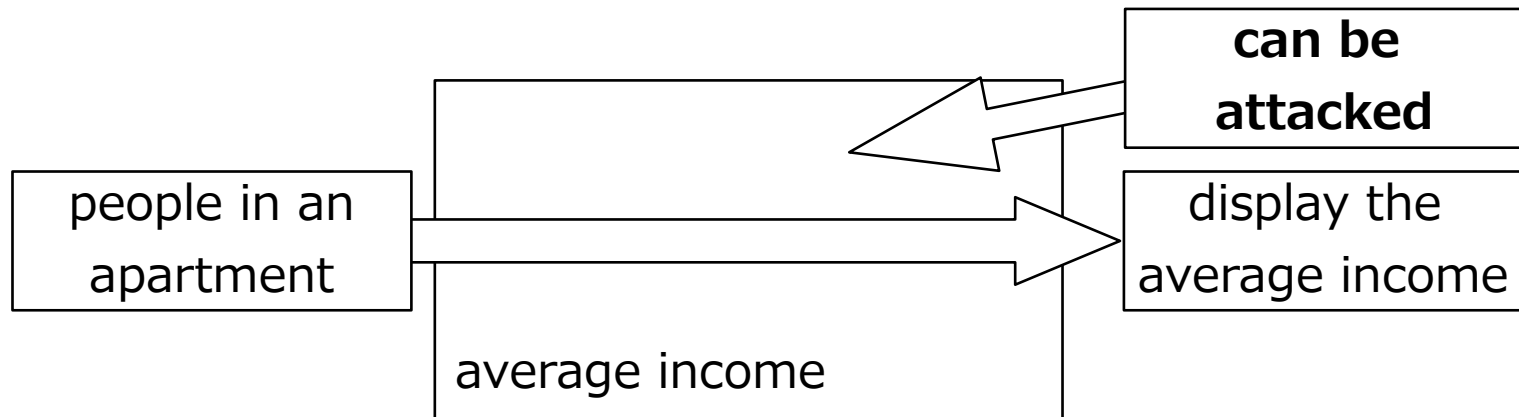
- 3 people lives here. Their average income is \$50,000.

# An Example of Privacy Leakage



- 3 people lives here. Their average income is \$50,000.
- Now, Bill, the 4<sup>th</sup> person joins here.  
Then, the average income changes to \$150,000.

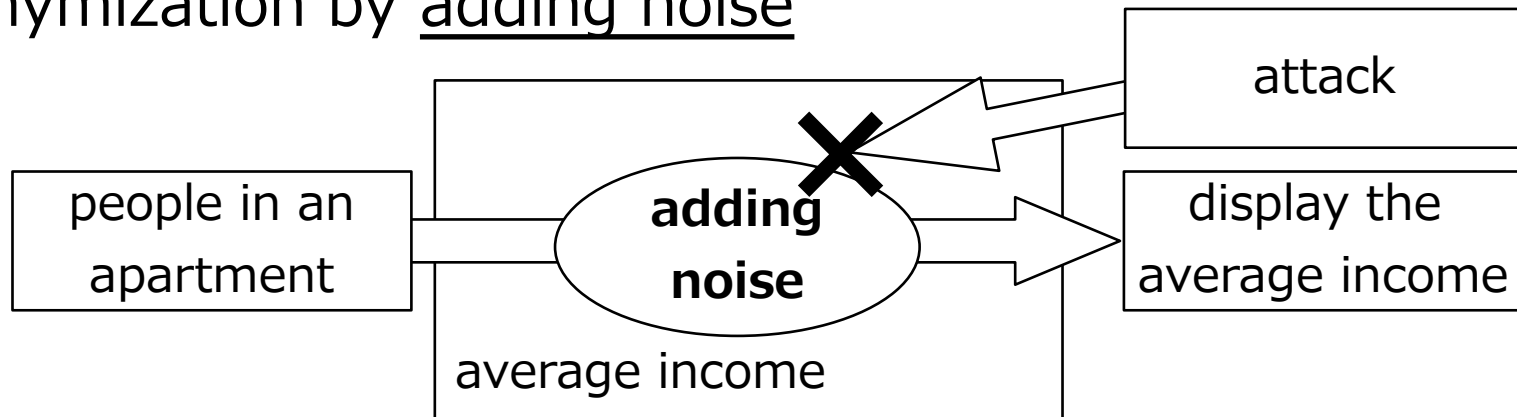
# An Example of Privacy Leakage



- 3 people lives here. Their average income is \$50,000.
- Now, Bill, the 4<sup>th</sup> person joins here.  
Then, the average income changes to \$150,000.
- While the “average income” query shows only the average, but it **leaks Bill’s income**, \$450,000.

# Differential Privacy

## 1. Anonymization by adding noise



- By adding noise, we make query hard to leak private data.
- The method is robust against background knowledge attack.

## 2. Standards of privacy in such randomized queries

# Definition of Differential Privacy

[Dwork+, TCC 2006]

- A randomized mechanism  $\mathbf{M}: \mathbf{X} \rightarrow \mathbf{Prob}(\mathbf{Y})$  is  **$(\epsilon, \delta)$ -differentially private(DP)** if for “adjacent” datasets  $D_1 \sim D_2$ , the following inequality holds for any  $S \subseteq \mathbf{Y}$  :

$$\Pr[M(D_1) \in S] \leq \exp(\epsilon) \Pr[M(D_2) \in S] + \delta$$

– intuition:

The ratio of probability is bounded by  $\epsilon$  except in probability  $\delta$ .

(If  $(\epsilon, \delta) = (0, 0)$  then, the distributions are equal. )

# Reformulating DP via Divergences

[Barthe & Olmedo, ICALP 2013]

- $M: X \rightarrow \text{Prob}(Y)$  satisfied  $(\epsilon, \delta)$ -DP
- $\Leftrightarrow$  for adjacent datasets  $D_1 \sim D_2$ ,

$$\Pr[M(D_1) \in S] \leq \exp(\epsilon) \Pr[M(D_2) \in S] + \delta$$

$\Leftrightarrow$  for adjacent datasets  $D_1 \sim D_2$ ,

$$\sup_{S \in \Sigma_Y} (\Pr[M(D_1) \in S] - \exp(\epsilon) \Pr[M(D_2) \in S]) \leq \delta$$

$$\Delta^\epsilon(M(D_1) || M(D_2))$$

$\Delta^\epsilon$  ... the divergence for  $(\epsilon, \delta)$ -DP

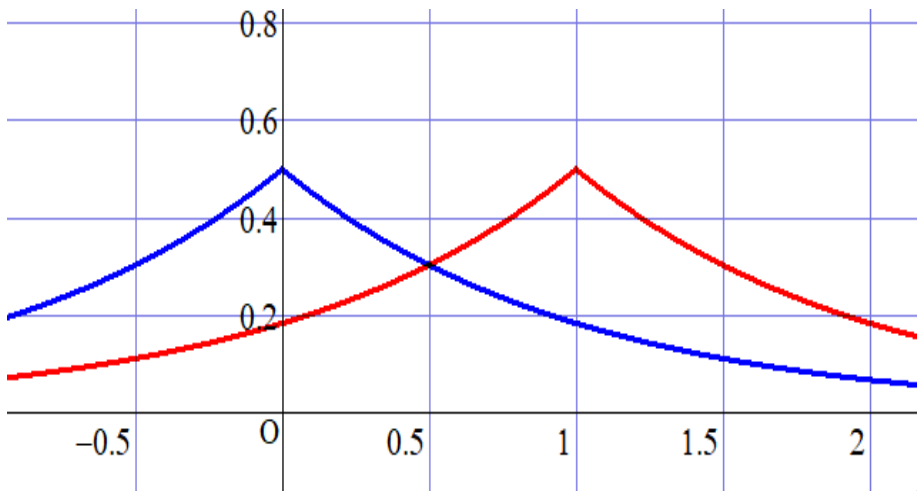


# Example: Laplace Mechanism

- The mechanism adding the noise sampled from Laplacian distribution. It is  $(\epsilon, 0)$ -DP if the adjacency is  $|x-y| \leq 1$ .

$$\text{Lap}_\epsilon : \mathbb{R} \rightarrow \text{Prob}(\mathbb{R})$$

$\text{Lap}_\epsilon(x)$  is the Laplacian distribution(avg.  $x$ , var.  $2\epsilon^2$ )

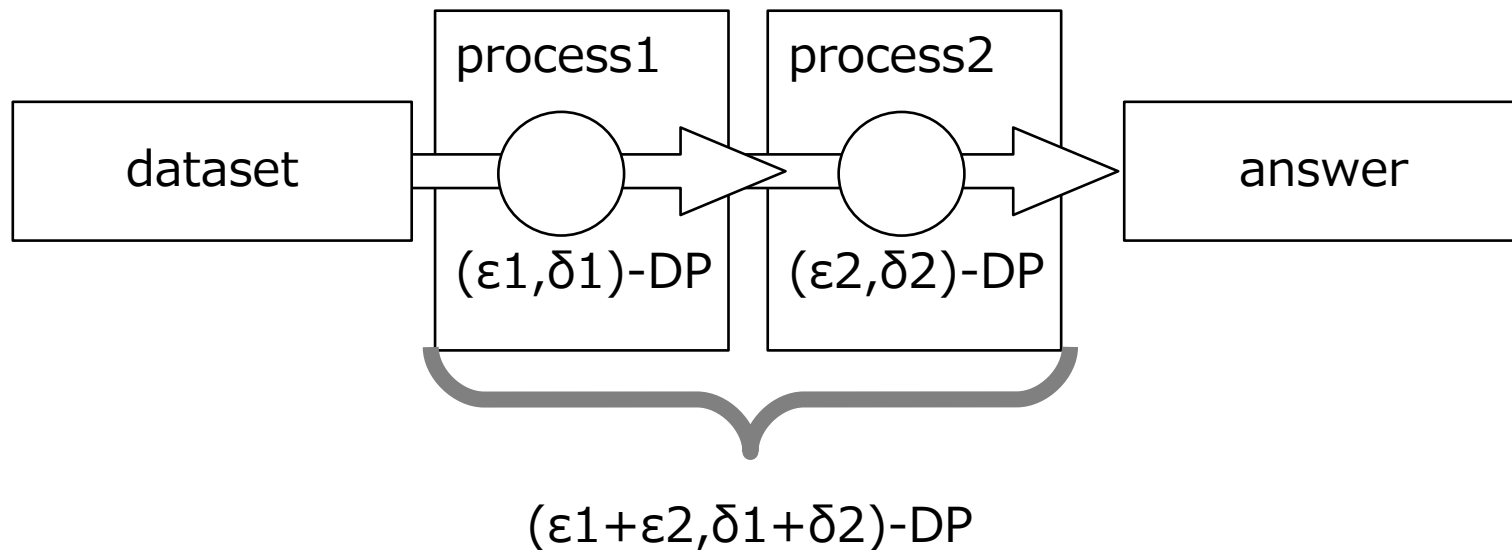


$$\begin{aligned} \Pr[\text{Lap}_\epsilon(x) = z] &= \frac{1}{2\epsilon} \exp\left(-\frac{|x-z|}{\epsilon}\right) \\ &\leq \frac{1}{2\epsilon} \exp\left(-\frac{|y-z| - |x-y|}{\epsilon}\right) \\ &\leq \exp(\epsilon) \frac{1}{2\epsilon} \exp\left(-\frac{|y-z|}{\epsilon}\right) \\ &= \exp(\epsilon) \Pr[\text{Lap}_\epsilon(y) = z] \end{aligned}$$

The ratio is bounded by  $\exp(\epsilon)$  everywhere.

# (Sequential) Composability of DP

- Differential privacy of a sequential composition of processes can be estimated by ones of its components.



- DP of a fixed number of loop of private mechanism can be estimated by the DP of the loop body.

# Naive Report Noisy-Max (RNM)

- Consider the following simple mechanism:

$\text{RNM}_\varepsilon : \text{list}(\mathbb{R}) \rightarrow \text{Prob}(\mathbb{R})$   
input :  $D = [x_1, \dots, x_n] \in \text{list}(\mathbb{R})$   
1. sample  $y_k \leftarrow \text{Lap}_\varepsilon(x_k)$  ( $1 \leq k \leq n$ )  
2. return  $\max\{y_k | 1 \leq k \leq n\}$

- Using the composability, we have  $(\varepsilon, 0)$ -DP
  - when the adjacency is defined by

$$D_1 \sim D_2 \iff \underbrace{|D_1| = |D_2|}_{\text{same length}} \wedge \underbrace{\|D_1 - D_2\|_1 \leq 1}_{\text{L1-norm}}$$

**Today's topic: we formalize this fact.**

# Proof Sketch

$\text{RNM}_\varepsilon : \text{list}(\mathbb{R}) \rightarrow \text{Prob}(\mathbb{R})$

input :  $D = [x_1, \dots, x_n] \in \text{list}(\mathbb{R})$

1. sample  $y_k \leftarrow \text{Lap}_\varepsilon(x_k)$  ( $1 \leq k \leq n$ )
2. return  $\max\{y_k | 1 \leq k \leq n\}$

- Show a bit stronger statement, by induction on length  $n$ :

$$|D_1| = |D_2| = n \wedge \|D_1 - D_2\|_1 \leq r$$

$$\implies \Delta^{r\varepsilon}(\text{RNM}_\varepsilon(D_1) || \text{RNM}_\varepsilon(D_2)) \leq 0$$

- (case:  $n = 1$ ) Using the DP of Laplacian mechanism:

$$|x_1 - x_2| \leq r \implies \Delta^{r\varepsilon}(\text{Lap}_\varepsilon(x_1) || \text{Lap}_\varepsilon(x_2)) \leq 0$$

- (case:  $n = k + 1$ ) Use I.H. and the below equation:

$$\text{RNM}_\varepsilon(x :: xs)$$

$$= (\text{Lap}_\varepsilon(x) \otimes \text{RNM}_\varepsilon(xs)) \succcurlyeq (\lambda(x, y). \text{return max}(x, y))$$

# Proof Assistant

- **Proof Assistant:**

- a tool that assists with writing formal proofs.
- We can program definitions, theorems and proofs, and certificate their validity.

(Isabelle/HOL example)

```
fun func1::"nat list  $\Rightarrow$  nat" where  
  "func1 [] = 0" | "func1(x # xs) = x + func1 (xs)"
```

```
lemma  
  fixes a xs  
  shows "func1 (xs @ [a]) = func1 (a # xs)"  
  by (induction xs, auto)
```

# Probability Theory in Isabelle/HOL

- Isabelle/HOL's standard library contains:

- Measure type `'a measure`  $(X, \Sigma_X, \mu)$ 
  - underlying set `"space"`  
`:: "'a measure  $\Rightarrow$  'a set"`  $X$
  - $\sigma$ -algebra `"sets"`  
`:: "'a measure  $\Rightarrow$  'a set set"`  $\Sigma_X$
  - evaluation `"Sigma_Algebra.measure"`  
`:: "'a measure  $\Rightarrow$  'a set  $\Rightarrow$  real"`  $\mu(A) \quad A \in \Sigma_X$
  - measurable functions `"( $\rightarrow_M$ )"`  
`:: "'a measure  $\Rightarrow$  'b measure  $\Rightarrow$  ('a  $\Rightarrow$  'b) set"`
- Monad for probability `"prob_algebra M"`  $\text{Prob}(X, \Sigma_X)$   
`:: "'a measure measure"`
  - bind `"( $\gg$ )"`  
`:: "'a measure  $\Rightarrow$  ('a  $\Rightarrow$  'b measure)  $\Rightarrow$  'b measure"`
  - return `"return"`  
`:: "'a measure  $\Rightarrow$  'a  $\Rightarrow$  'a measure"`
- Radon-Nikodym derivative `"RN_deriv M N"`  
`:: "'a  $\Rightarrow$  ennreal"`  $\frac{d\mu_M}{d\mu_N}$
- Lebesgue measure `"lborel"`  
`:: "real measure"`

# Formalizing the Divergence for DP in Isabelle/HOL

$\text{ereal} = [-\infty, \infty]$

(definition)

```
definition DP_divergence:: "'a measure  $\Rightarrow$  'a measure  $\Rightarrow$  real  $\Rightarrow$  ereal " where  
  "DP_divergence M N  $\epsilon$  = ( $\bigcup$  A  $\in$  (sets M). ereal( measure M A - (exp  $\epsilon$ ) * measure N A))"
```

(non-negativity)

```
lemma DP_divergence_nonnegativity:  
  assumes M: "M  $\in$  space (prob_algebra L)" and N: "N  $\in$  space (prob_algebra L)"  
  shows "0  $\leq$  DP_divergence M N  $\epsilon$  "
```

(basic properties (for detail, [Olmedo, Phd thesis 2014]))

```
lemma DP_divergence_monotonicity:  
  assumes M: "M  $\in$  space (prob_algebra L)" and N: "N  $\in$  space (prob_algebra L)"  
  and " $\epsilon_1 \leq \epsilon_2$  "  
  shows "DP_divergence M N  $\epsilon_2 \leq$  DP_divergence M N  $\epsilon_1$  "
```

```
lemma DP_reflexivity:  
  shows " DP_divergence M M 0 = 0 "
```

"locale" structure providing  
the assumption  $M, N \in \text{Prob}(N)$

```
theorem (in comparable_probability_measures) DP_composability:  
  assumes f: "f  $\in$  measurable L (prob_algebra K)"  
  and g: "g  $\in$  measurable L (prob_algebra K)"  
  and div1: "DP_divergence M N  $\epsilon_1 \leq (\delta_1::\text{real})"$   
  and div2: " $\forall x \in$  (space L). DP_divergence (f x) (g x)  $\epsilon_2 \leq (\delta_2::\text{real})"$   
  and "0  $\leq \epsilon_1$ " "0  $\leq \epsilon_2$ "  
  shows "DP_divergence (bind M f) (bind N g) ( $\epsilon_1 + \epsilon_2$ )  $\leq \delta_1 + \delta_2$ "
```

# Proof Sketch of Composability of DP

$$\begin{aligned}
 & \Pr[\mu \ggg f \in S] - \exp(\varepsilon_1 + \varepsilon_2) \Pr[\nu \ggg g \in S] && \text{expanding the bind with densities} \\
 &= \int f(-)(S) \cdot \frac{d\mu}{d\pi} d\pi - \exp(\varepsilon_1 + \varepsilon_2) \int g(-)(S) \cdot \frac{d\nu}{d\pi} d\pi && \text{linearity of integrals} \\
 &= \int f(-)(S) \cdot \frac{d\mu}{d\pi} - \exp(\varepsilon_1 + \varepsilon_2) g(-)(S) \cdot \frac{d\nu}{d\pi} d\pi \\
 &\leq \int (\max(0, f(-)(S) - \delta_2) + \delta_2) \cdot \frac{d\mu}{d\pi} - \exp(\varepsilon_1) \min(1, \exp(\varepsilon_2) g(-)(S)) \cdot \frac{d\nu}{d\pi} d\pi \\
 &= \int \max(0, f(-)(S) - \delta_2) \cdot \frac{d\mu}{d\pi} - \exp(\varepsilon_1) \min(1, \exp(\varepsilon_2) g(-)(S)) \cdot \frac{d\nu}{d\pi} d\pi + \int \delta_2 \frac{d\mu}{d\pi} d\pi \\
 &\leq \int_B \left( \frac{d\mu}{d\pi} - \exp(\varepsilon_1) \frac{d\nu}{d\pi} \right) \cdot \min(1, \exp(\varepsilon_2) \cdot g(-)(S)) d\pi + \int \delta_2 \frac{d\mu}{d\pi} d\pi \\
 &\leq \delta_1 + \delta_2 && \text{taking the positive part} \\
 &= \Delta^{\varepsilon_1}(\mu || \nu) + \sup_x \Delta^{\varepsilon_2}(f(x) || g(x)) && \text{assumptions}
 \end{aligned}$$



# Formal Proof of Composability of DP in Isabelle/HOL

- Most of the formal proof can be done according to the sketch.

```

have "(measure (M  $\gg$  f) A) - exp ( $\epsilon_1 + \epsilon_2$ ) * (measure (N  $\gg$  g) A) [3 lines]
also have "... = ( $\int$  x. (dM x) * (measure (f x) A)  $\partial$ (sum_measure M N)) - ( $\int$  x. (exp ( $\epsilon_1 + \epsilon_2$ )) * (dN x) * (measure (g x) A)  $\partial$ (sum_measure M N))" [1 lines]
also have "... = ( $\int$  x. (dM x) * (measure (f x) A) - (exp ( $\epsilon_1 + \epsilon_2$ )) * (dN x) * (measure (g x) A)  $\partial$ (sum_measure M N))" [1 lines]
also have "... = ( $\int$  x. (dM x) * (measure (f x) A) - (exp  $\epsilon_1$ ) * (exp  $\epsilon_2$ ) * (dN x) * (measure (g x) A)  $\partial$ (sum_measure M N))" [1 lines]
also have "...  $\leq$  ( $\int$  x. (dM x) * (max 0 (measure (f x) A -  $\delta_2$ ) +  $\delta_2$ ) - (exp  $\epsilon_1$ ) * (dN x) * min 1 ((exp  $\epsilon_2$ ) * (measure (g x) A))  $\partial$ (sum_measure M N))" [15 lines]
also have "... = ( $\int$  x. (dM x) * (max 0 (measure (f x) A -  $\delta_2$ )) - (exp  $\epsilon_1$ ) * (dN x) * min 1 ((exp  $\epsilon_2$ ) * (measure (g x) A) + (dM x) *  $\delta_2$   $\partial$ (sum_measure M N))" [1 lines]
also have "...  $\leq$  ( $\int$  x. (dM x) * (min 1 ((exp  $\epsilon_2$ ) * (measure (g x) A))) - (exp  $\epsilon_1$ ) * (dN x) * min 1 ((exp  $\epsilon_2$ ) * (measure (g x) A) + dM x *  $\delta_2$   $\partial$ (sum_measure M N))" [12 lines]
also have "... = ( $\int$  x. ((dM x) - (exp  $\epsilon_1$ ) * (dN x)) * min 1 ((exp  $\epsilon_2$ ) * (measure (g x) A)) + dM x *  $\delta_2$   $\partial$ (sum_measure M N))" [1 lines]
also have "... = ( $\int$  x. ((dM x) - (exp  $\epsilon_1$ ) * (dN x)) * min 1 ((exp  $\epsilon_2$ ) * (measure (g x) A))  $\partial$ (sum_measure M N)) + ( $\int$  x. dM x *  $\delta_2$   $\partial$ (sum_measure M N))" [1 lines]
finally have *: "(measure (M  $\gg$  f) A) - exp ( $\epsilon_1 + \epsilon_2$ ) * (measure (N  $\gg$  g) A)  $\leq$  ( $\int$  x. ((dM x) - (exp  $\epsilon_1$ ) * (dN x)) * min 1 ((exp  $\epsilon_2$ ) * (measure (g x) A))  $\partial$ (sum_measure M N)) + ( $\int$  x. dM x *  $\delta_2$   $\partial$ (sum_measure M N))".

have "( $\int$  x. dM x *  $\delta_2$   $\partial$ (sum_measure M N)) = ( $\int$  x.  $\delta_2$   $\partial$ (density (sum_measure M N) dM))" [1 lines]
also have "... = ( $\int$  x.  $\delta_2$   $\partial$ (density (sum_measure M N) (ennreal 0 dM)))" [1 lines]
also have "... = ( $\int$  x.  $\delta_2$  dM)" [1 lines]
also have "... =  $\delta_2$  * measure M (space M)" [1 lines]
also have "...  $\leq \delta_2$ " [1 lines]
finally have **: "( $\int$  x. dM x *  $\delta_2$   $\partial$ (sum_measure M N))  $\leq \delta_2$ ".

let ?B = "{x  $\in$  space (sum_measure M N). 0  $\leq$  ((dM x) - (exp  $\epsilon_1$ ) * (dN x)) }"

have mble10: "?B  $\in$  sets (sum_measure M N)" [1 lines]

have "( $\int$  x. ((dM x) - (exp  $\epsilon_1$ ) * (dN x)) * min 1 ((exp  $\epsilon_2$ ) * (measure (g x) A))  $\partial$ (sum_measure M N))  $\leq$  ( $\int$  x  $\in$  ?B. ((dM x) - (exp  $\epsilon_1$ ) * (dN x)) * min 1 ((exp  $\epsilon_2$ ) * (measure (g x) A))  $\partial$ (sum_measure M N))"
proof(rule integral_drop_negative_part2) [17 lines]
qed
also have "...  $\leq$  ( $\int$  x  $\in$  ?B. ((dM x) - (exp  $\epsilon_1$ ) * (dN x))  $\partial$ (sum_measure M N))" [11 lines]
also have "... = ( $\int$  x  $\in$  ?B. (dM x)  $\partial$ (sum_measure M N)) - ( $\int$  x  $\in$  ?B. ((exp  $\epsilon_1$ ) * (dN x))  $\partial$ (sum_measure M N))" [8 lines]
also have "... = ( $\int$  x  $\in$  ?B. (dM x)  $\partial$ (sum_measure M N)) - (exp  $\epsilon_1$ ) * ( $\int$  x  $\in$  ?B. (dN x)  $\partial$ (sum_measure M N))" [1 lines]
also have "... = measure M ?B - (exp  $\epsilon_1$ ) * (measure N ?B)" [42 lines]
also have "...  $\leq \delta_1$ " [1 lines]
finally have ***: "( $\int$  x. ((dM x) - (exp  $\epsilon_1$ ) * (dN x)) * min 1 ((exp  $\epsilon_2$ ) * (measure (g x) A))  $\partial$ (sum_measure M N))  $\leq \delta_1$ ".

show "measure (M  $\gg$  f) A - exp ( $\epsilon_1 + \epsilon_2$ ) * measure (N  $\gg$  g) A  $\leq \delta_1 + \delta_2$ "
using * ** *** by auto
qed

```

# Formalizing the Laplace Mechanism in Isabelle/HOL

- (definition)

```
definition laplace_density :: "real  $\Rightarrow$  real  $\Rightarrow$  real  $\Rightarrow$  real" where
  "laplace_density l m x = (if l > 0 then (exp(-| x - m | / l) / (2* l)) else 0)"

definition Lap_mechanism :: "real  $\Rightarrow$  real  $\Rightarrow$  real measure"
  where "Lap_mechanism  $\epsilon$  x = (if  $\epsilon \leq 0$  then return lborel x
    else (density lborel (laplace_density (1/ $\epsilon$ ) x)))"
```

- (measurability)

```
lemma measurable_Lap_mechanism[measurable]:
  shows "Lap_mechanism  $\epsilon \in$  lborel  $\rightarrow_M$  prob_algebra lborel"
```

- (differential privacy(formalized via divergence DP))

```
proposition DP_Lap_mechanism':
  fixes x y  $\epsilon$  :: real
  assumes " $\epsilon > 0$ " and " $| x - y | \leq r$ "
  shows "DP_divergence (Lap_mechanism  $\epsilon$  x) (Lap_mechanism  $\epsilon$  y) (r *  $\epsilon$ )  $\leq$  (0::real)"
```

# Formalizing the Naive RNM in Isabelle/HOL

- (definition)

```
fun RNM :: "real  $\Rightarrow$  real list  $\Rightarrow$  real measure "  
  where  
    "RNM  $\epsilon$  [] = (return lborel 0)" | (* empty case, it is dummy *)  
    "RNM  $\epsilon$  [x] = (Lap_mechanism  $\epsilon$  x)" |  
    "RNM  $\epsilon$  (x # xs) = do { x1  $\leftarrow$  (Lap_mechanism  $\epsilon$  x);  
      x2  $\leftarrow$  (RNM  $\epsilon$  xs); (return lborel (max x1 x2)) }"
```

- (measurability)

```
lemma measurable_RNM:  
  shows "(RNM  $\epsilon$ )  $\in$  (listM lborel)  $\rightarrow_M$  (prob_algebra lborel)"
```

- (differential privacy)

```
theorem DP_RNM:  
  fixes xs ys :: "real list" and  $\epsilon$  :: real and n :: nat and r :: real  
  assumes pose [arith]: " $\epsilon > (0 :: real)$ "  
    and adj: "length xs = n  $\wedge$  length ys = n  
     $\wedge$  (  $\sum_{i \in \{1..n\}} | \text{nth xs } (i-1) - \text{nth ys } (i-1) |$ )  $\leq r$ "  
    and posr [arith]: " $r \geq 0$ "  
  shows "DP_divergence (RNM  $\epsilon$  xs) (RNM  $\epsilon$  ys) (r *  $\epsilon$ )  $\leq 0$ "
```

# Concluding Remark

- Formal verification of DP in the discrete setting is already implemented in Coq[Barthe+, TOPLAS2013].
- We now aim to develop an Isabelle/HOL library for formal verification of DP in the continuous setting.
  - Today, we have formalized DP of naive report noisy-max in the continuous setting. It is the first formalization example.
  - Next,
    - We are formalizing DP of (true) report noisy-max.
    - We need to optimize what we have implemented.
    - We want to formalize relaxation of DP, such as RDP[Mironov, CSF2017], zCDP[Bun+, TCC2016].

**Thank you!**