

Privacy Preserving Machine Learning

SAV Après-Midi, 31 August 2023

Daniel Meier, Swiss Re

Juan Ramón Troncoso-Pastoriza, Tune Insight


Key take-aways

1. Limited data availability for sensitive personal (Life & Health) data in practice (e.g., nFADP, 1 September 2023)
2. For sufficiently large and dense datasets, ML/DL methods outperform traditional models, creating value for policyholders and insurance companies
3. Privacy preserving methods can help to access data

Agenda

- Creation of synthetic health datasets
- Introducing 3 models to create health risk scores: Logistic regression, Cox regression, neural networks
- Homomorphic encryption

Paper and code soon available at actuarialdatascience.org



Actuarial Data Science

An initiative of the Swiss Association of Actuaries

Home	<h3>Actuarial Data Science Tutorials</h3> <p>On this page we present all the tutorials that have been prepared by the working party. We are intensively working on additional ones and we aim to have approx. 10 tutorials, covering a wide range of Data Science topics relevant for actuaries.</p>	<h3>Updates</h3> <p>Below, we provide the most recent changes to the website:</p> <ul style="list-style-type: none">• 15th Mar 23: Publication of our new tutorial: SHAP for Actuaries: Explain any Model• 14th Oct 22: Publication of our new tutorial: Gini Index and Friends
ADS Tutorials		
ADS Strategy		
ADS Lectures / Courses		
ADS Regulatory / Ethics	All tutorials consist of an article and the corresponding code. In the article, we describe the methodology and the statistical model. By providing you with the code you can easily replicate the analysis performed and test it on your own data.	
DS Lectures / Books		
External Courses	Case Study 14: SHAP for Actuaries: Explain any Model	Events <p>Below, we provide upcoming events in Actuarial Data Science:</p> <ul style="list-style-type: none">• 16th May 2023: EAA Data
Newsletter	Article on SSRN	
About Us	Code on GitHub ; Notebook	
	Case Study 13: Gini Index and Friends	
	Article on SSRN	
	Code on GitHub	

(Publicly) available health datasets

- CPRD, <https://cprd.com/data>
- MIMIC, <https://physionet.org/about/database/>
- IPUMS, <https://healthsurveys.ipums.org/>
- NHANES, <https://www.cdc.gov/nchs/nhanes/>
- Nightingale, <https://docs.nightingalescience.org/>
- UK Biobank, <https://www.ukbiobank.ac.uk/>
- IHME, <https://ghdx.healthdata.org/>
- ...
- See also [longitudinal study](#) for other health datasets

- Often, access is restricted to academic institutions and/or limited to a pre-defined research topic
- Data volumes (and density) rather too low for ML
- Access to more data sources (e.g., hospitals, GPs, insurance companies, etc.) – in a privacy preserving manner – is needed

Health risk scores, e.g., QRISK3 providing 10-year risk of a cardio-vascular disease (CVD)

ClinRisk

Welcome to the QRISK[®]-3-2018 risk calculator <https://qrisk.org>

This calculator is only valid if you do not already have a diagnosis of coronary heart disease (including angina or heart attack) or stroke/transient ischaemic attack.

Reset

Information

Publications

About

Copyright

Contact Us

Algorithm

Software

UKCA

About you

Age (25-84): 35

Sex:

☒ Male☐ Female

Ethnicity: White or not stated v

UK postcode: leave blank if unknown

Postcode:

Clinical information

Smoking status: non-smoker v

Diabetes status: none v

Angina or heart attack in a 1st degree relative < 60? ☐

Chronic kidney disease (stage 3, 4 or 5)? ☐

Atrial fibrillation? ☐

On blood pressure treatment? ☐

Do you have migraines? ☐

Rheumatoid arthritis? ☐

Systemic lupus erythematosus (SLE)? ☐

Severe mental illness?
(this includes schizophrenia, bipolar disorder and moderate/severe depression) ☐

On atypical antipsychotic medication? ☐

Are you on regular steroid tablets? ☐

A diagnosis of or treatment for erectile dysfunction? ☐

Leave blank if unknown

Cholesterol/HDL ratio: 3.5

Systolic blood pressure (mmHg): 120

Standard deviation of at least two most recent systolic blood pressure readings (mmHg): 10

Body mass index

Height (cm): 180

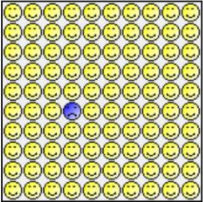
Weight (kg): 75

Your results

Your risk of having a heart attack or stroke within the next 10 years is:

0.6%

In other words, in a crowd of 100 people with the same risk factors as you, 1 are likely to have a heart attack or stroke within the next 10 years.



Risk of
a heart attack or stroke

Your score has been calculated using estimated data, as some information was left blank.

Your body mass index was calculated as 23.15 kg/m².

How does your 10-year score compare?

Your score

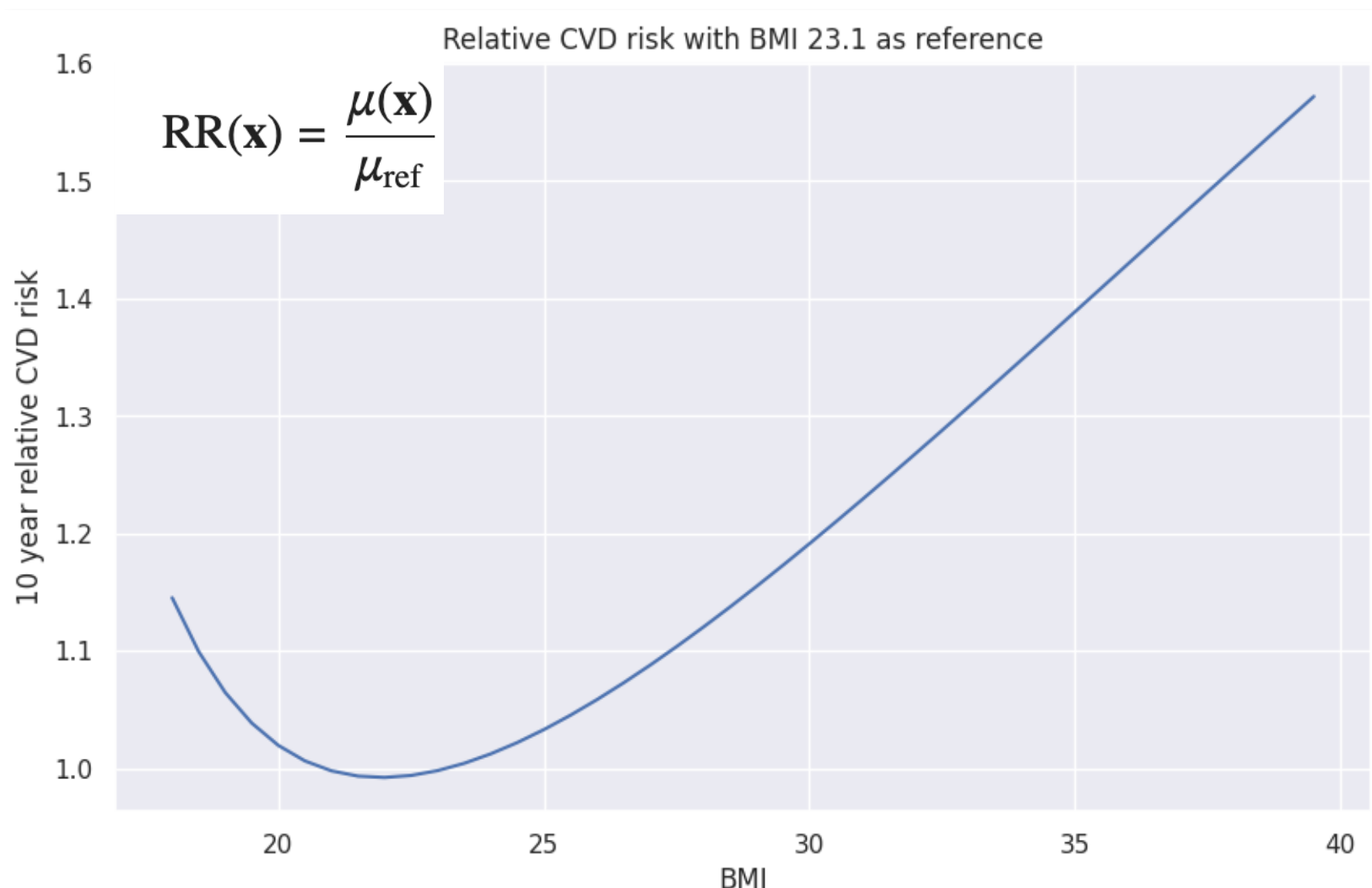
Your 10-year QRISK [®] -3 score	0.6%
The score of a healthy person with the same age, sex, and ethnicity*	0.7%
Relative risk**	0.9
Your QRISK [®] -3 Healthy Heart Age***	35

* This is the score of a healthy person of your age, sex and ethnic group, i.e. with no adverse clinical indicators and a cholesterol ratio of 4.0, a stable systolic blood pressure of 125, and BMI of 25.
** Your relative risk is your risk divided by the healthy person's risk.
*** Your QRISK[®]-3 Healthy Heart Age is the age at which a healthy person of your sex and ethnicity has your 10-year QRISK[®]-3 score.

Calculate risk

Source: qrisk.org

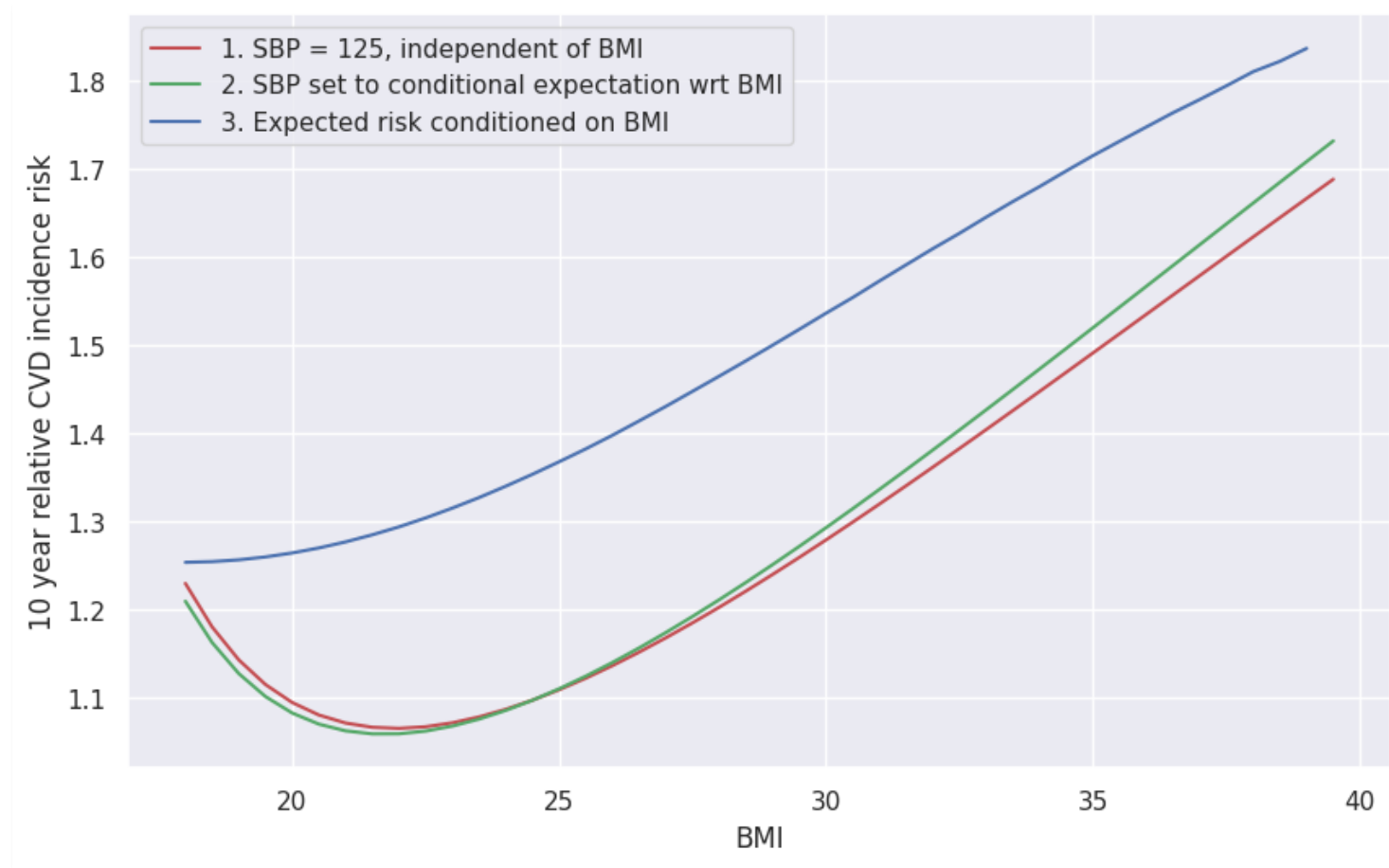
Relative risk with respect to a reference person of same age, gender



Country	Mean BMI females	Mean BMI males
Samoa	33.5	29.9
USA	28.8	28.8
UK	27.1	27.5
Germany	25.6	27.0
Italy	25.2	26.8
France	24.6	26.1
Switzerland	23.8	26.7
Japan	21.7	23.6

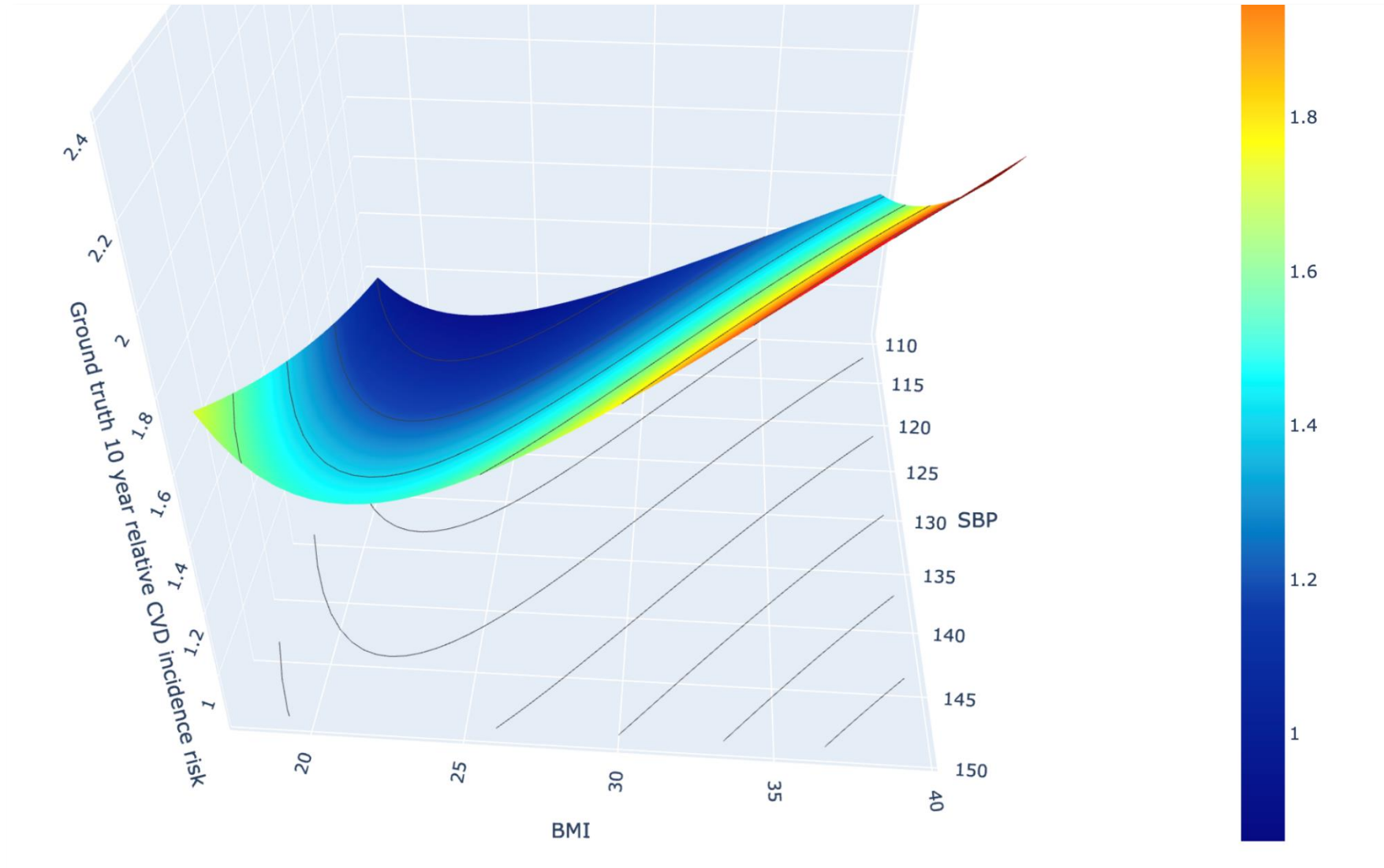
Country	Life exp. females	Life exp. females
Samoa	75.5	71.3
USA	81.5	76.5
UK	83.3	79.6
Germany	83.5	78.8
Italy	85.4	81.1
France	85.6	79.8
Switzerland	85.6	81.9
Japan	87.4	81.4

Various risk factors like BMI, systolic blood pressure (SBP) impact relative risk



1. What is the risk of a person with a given BMI and all other attributes equal to the reference person, $\mu(\text{BMI}, \text{SBP}_{\text{ref}})$?
2. What is the risk of a person with a given BMI, and SBP set to the conditional expectation of SBP given BMI, $\mu(\text{BMI}, \mathbb{E}[\text{SBP}|\text{BMI}])$?
3. What is the expected risk of a person with a given BMI, $\mathbb{E}[\mu(\mathbf{x})|\text{BMI}]$?
4. What is the (causally implied) risk of the reference person when changing BMI, $\mathbb{E}[\mu(\mathbf{x})|do(\text{BMI})]$?

Various risk factors like BMI, systolic blood pressure (SBP) impact relative risk



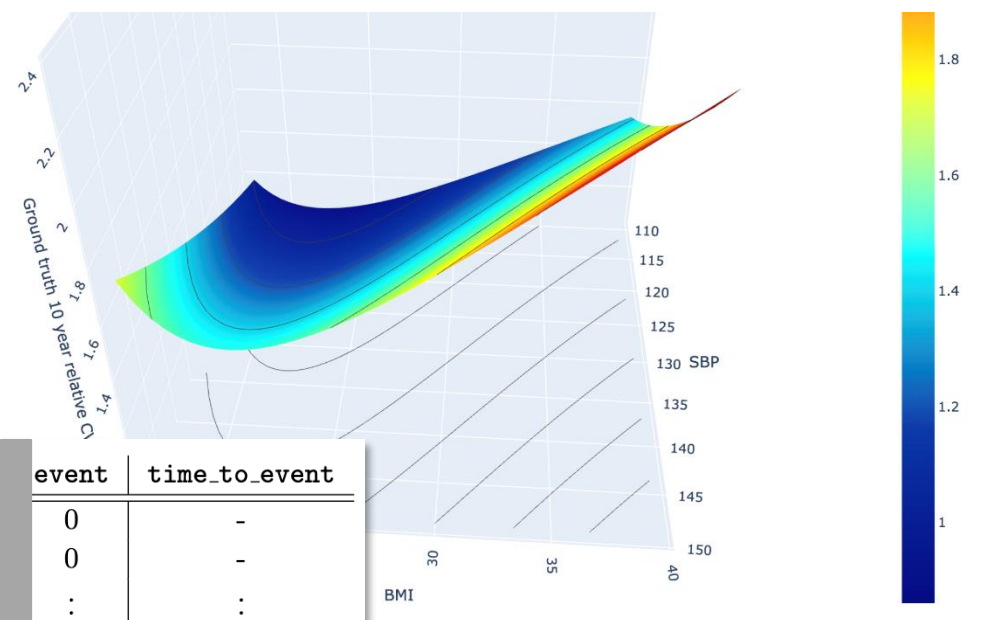
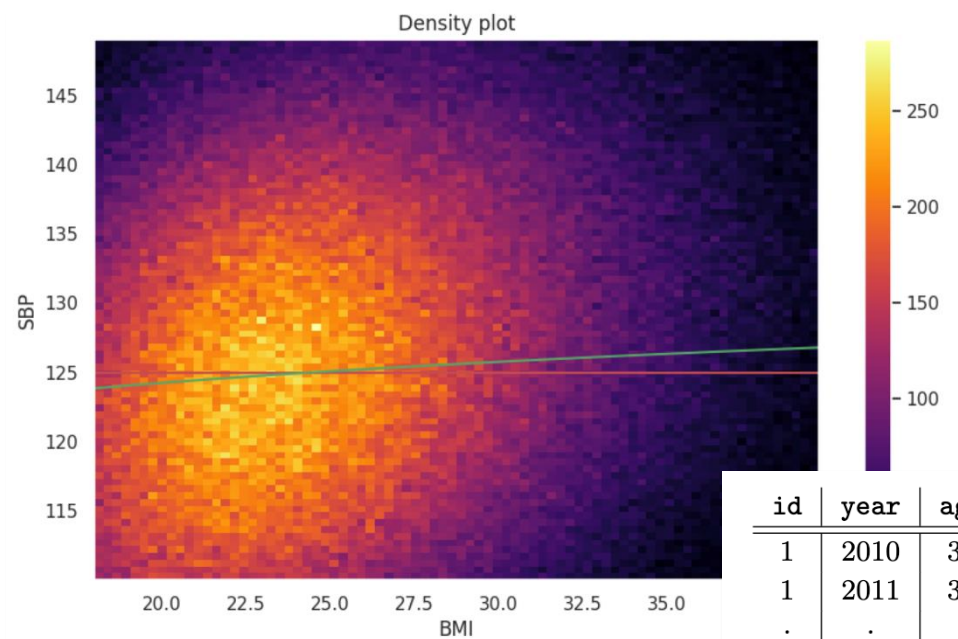
Creation of a synthetic health dataset

- **id**: an ID to uniquely identify a person,
- **year**: observation year of health information,
- **age**: age of the person at time **year**,
- **gender**: male (0)/female (1),
- **bmi**: body-mass-index (BMI), unit kg/m^2 ,
- **sbp**: systolic blood pressure (SBP), unit mmHg,
- **sd_sbp**: standard deviation of systolic blood pressure measurements, unit mmHg,
- **tcl_hdl_ratio**: total cholesterol level (TCL) divided by high-density lipoprotein level (HDL),
- **num1**, **num2**, **num3**: 3 generic numeric health risk factors without specifying their meaning explicitly, e.g., stepcounts, triglycerides, resting heartrate, etc.
- **binary**: a generic binary health risk factor, e.g., smokers, foreign born, etc.,

Creation of a synthetic health dataset

$$\begin{pmatrix} \text{sbp} \\ \log(\text{bmi}) \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 125 \\ 3.2 \end{pmatrix}, \begin{pmatrix} 15^2 & 15 \cdot 0.25\rho \\ 15 \cdot 0.25\rho & 0.25^2 \end{pmatrix} \right)$$

$$\begin{aligned} \mu^*(\mathbf{x}) = & \exp \left(\log(\text{QRISK3}(\mathbf{x})) + r(\text{age}) \right) \\ & + 16(\text{num1} - 0.5)^4 + 4(\text{num2} - 0.5)^2 \text{num3} + \text{num3} + \text{binary} - 1.65 \end{aligned}$$

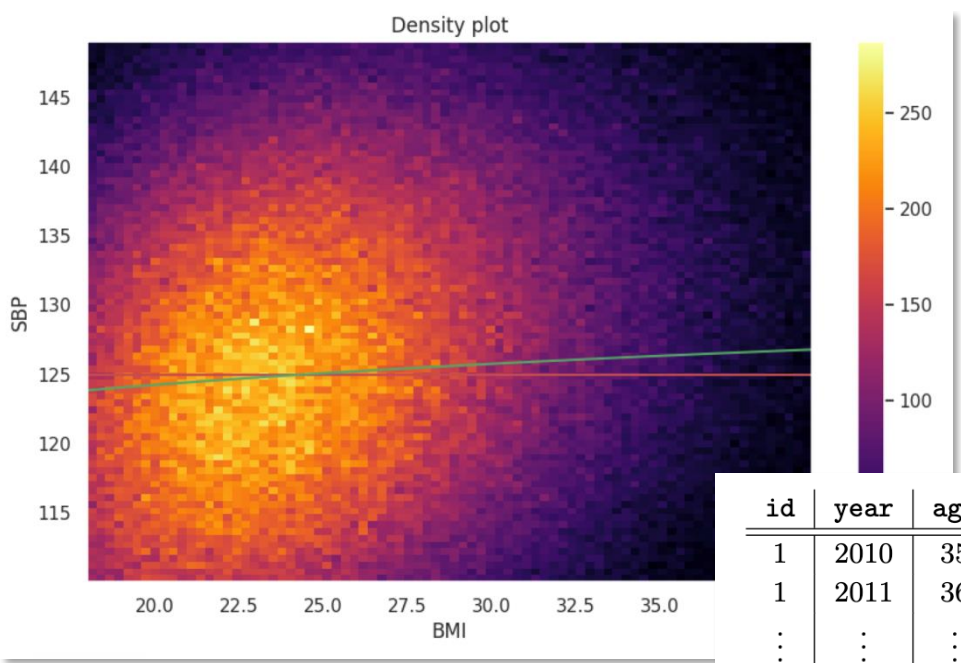


id	year	age	gender	bmi	sbp	event	time_to_event
1	2010	35	m	24	120	0	-
1	2011	36	m	24	120	0	-
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
1	2019	44	m	24	120	0	-
2	2010	35	m	33	145	0	7.5
2	2011	36	m	33	145	0	6.5
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
2	2017	44	m	33	145	1	0.5
3	2010	35	m	26	125	0	-
3	2011	36	m	26	125	0	-
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

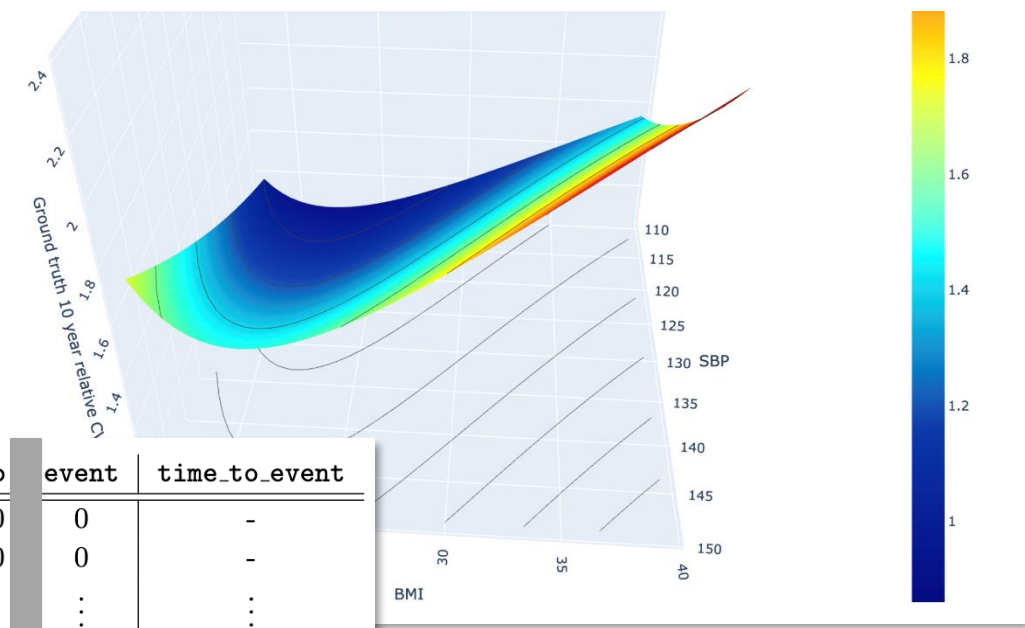


Creation of a synthetic health dataset

$$\begin{pmatrix} \text{sbp} \\ \log(\text{bmi}) \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 125 \\ 3.2 \end{pmatrix}, \begin{pmatrix} 15^2 & 15 \cdot 0.25\rho \\ 15 \cdot 0.25\rho & 0.25^2 \end{pmatrix} \right)$$



$$\begin{aligned} \mu^*(\boldsymbol{x}) = & \exp \left(\log(\text{QRISK3}(\boldsymbol{x})) + r(\text{age}) \right) \\ & + 16(\text{num1} - 0.5)^4 + 4(\text{num2} - 0.5)^2 \text{num3} + \text{num3} + \text{binary} - 1.65 \end{aligned}$$



id	year	age	gender	bmi	sbp	event	time_to_event
1	2010	35	m	24	120	0	-
1	2011	36	m	24	120	0	-
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
1	2019	44	m	24	120	0	-
2	2010	35	m	33	145	0	7.5
2	2011	36	m	33	145	0	6.5
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
2	2017	44	m	33	145	1	0.5
3	2010	35	m	26	125	0	-
3	2011	36	m	26	125	0	-
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

Model 1: Logistic regression/generalized linear model (GLM)

$$\mu_1(\mathbf{x}) = \frac{1}{1 + \exp(-\beta_0 - \beta_1 x_1 - \dots - \beta_k x_k)} , \text{ or equivalently}$$
$$\text{logit}(\mu_1(\mathbf{x})) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k ,$$

```
import statsmodels.formula.api as sm
log_reg = sm.logit(formula="E~SBP+BMI+I(BMI**2)", data=time_to_event_train).fit()
pred = log_reg.predict(time_to_event_test)
```

Odds/log-odds

$$\text{odds}(y = 1 \mid \mathbf{x}) := \frac{P(y = 1 \mid \mathbf{x})}{P(y = 0 \mid \mathbf{x})} = \frac{P(y = 1 \mid \mathbf{x})}{1 - P(y = 1 \mid \mathbf{x})}$$

$$\log(\text{odds}(y = 1 \mid \mathbf{x})) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k .$$

Odds ratios

$$\frac{\text{odds}(y = 1 \mid (x_1, \dots, x_j + 1, \dots, x_k))}{\text{odds}(y = 1 \mid \mathbf{x})} = \frac{\exp(\beta_0 + \beta_1 x_1 + \dots + \beta_j (x_j + 1) + \dots + \beta_k x_k)}{\exp(\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k)}$$
$$= \exp(\beta_j) .$$

id	year	age	gender	bmi	sbp	event	time_to_event
1	2010	35	m	24	120	0	-
1	2011	36	m	24	120	0	-
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
1	2019	44	m	24	120	0	-
2	2010	35	m	33	145	0	7.5
2	2011	36	m	33	145	0	6.5
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
2	2017	44	m	33	145	1	0.5
3	2010	35	m	26	125	0	-
3	2011	36	m	26	125	0	-
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

Model 2: Cox regression

$$h(t \mid \mathbf{x}) = h_0(t) \exp \left(\beta_1 x_1 + \cdots + \beta_k x_k \right)$$

```
import lifelines as ll
cph = ll.CoxPHFitter()
cph.fit(time_to_event, "T", event_col="E", formula="SBP+BMI+I(BMI**2)")
pred = (1 - np.array(cph.predict_survival_function(time_to_event_test))[10,:])
```

From hazard rates to 10-year risk

$$\mu_2(\mathbf{x}) := 1 - \exp \left(- \int_0^{10} h(t \mid \mathbf{x}) dt \right)$$

Hazard ratios

$$\frac{h_0(t) \exp \left(\beta_1 x_1 + \cdots + \beta_j (x_j + 1) + \cdots + \beta_k x_k \right)}{h_0(t) \exp \left(\beta_1 x_1 + \cdots + \beta_k x_k \right)} = \exp(\beta_j)$$

id	year	age	gender	bmi	sbp	event	time_to_event
1	2010	35	m	24	120	0	-
1	2011	36	m	24	120	0	-
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
1	2019	44	m	24	120	0	-
2	2010	35	m	33	145	0	7.5
2	2011	36	m	33	145	0	6.5
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
2	2017	44	m	33	145	1	0.5
3	2010	35	m	26	125	0	-
3	2011	36	m	26	125	0	-
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

Model 3: Neural networks

```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Input, Dense, Activation
model = Sequential([
    Input(shape=(10,)),
    Dense(256),
    Activation(tf.keras.activations.relu),
    Dense(128),
    Activation(tf.keras.activations.relu),
    Dense(64),
    Activation(tf.keras.activations.relu),
    Dense(1, activation = "sigmoid")
])
model.compile(optimizer = opt, loss = 'binary_crossentropy')
model.fit(X_train, y_train, batch_size = 64, epochs = 100, shuffle = True, validation_split = 0.20)
pred = model.predict(X_test).flatten()
```

$$\begin{aligned}\mu_3(\mathbf{x}) &= \mathbf{z}^{(4)} \circ \mathbf{z}^{(3)} \circ \mathbf{z}^{(2)} \circ \mathbf{z}^{(1)}(\mathbf{x}) , \text{ where} \\ z_j^{(\kappa)}(\mathbf{z}) &= \phi^{(\kappa)}\left(\beta_{0,j}^{(\kappa)} + \sum_{\ell=1}^{q_{\kappa}-1} \beta_{\ell,j}^{(\kappa)} z_{\ell}\right) , \text{ for } 1 \leq j \leq q_{\kappa} , \\ \phi^{(\kappa)}(x) &= \begin{cases} 1/(1 + \exp(-x)) & , \kappa = 4, \\ x\mathbb{1}_{\{x \geq 0\}} & , \text{ otherwise,} \end{cases} \\ (q_4, \dots, q_0) &= (1, 64, 128, 256, 10) .\end{aligned}$$

Layer	input size	output size	#parameters	input from layer κ
Fully-conn.	10	256	2'816	0
ReLU ϕ	256	256	0	1
Fully-conn.	256	128	32'896	1
ReLU ϕ	128	128	0	2
Fully-conn.	128	64	8'256	2
ReLU ϕ	64	64	0	3
Fully-conn.	64	1	65	3
Sigmoid output	1	1	0	4

id	year	age	gender	bmi	sbp	event	time_to_event
1	2010	35	m	24	120	0	-
1	2011	36	m	24	120	0	-
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
1	2019	44	m	24	120	0	-
2	2010	35	m	33	145	0	7.5
2	2011	36	m	33	145	0	6.5
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
2	2017	44	m	33	145	1	0.5
3	2010	35	m	26	125	0	-
3	2011	36	m	26	125	0	-
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

Model performance

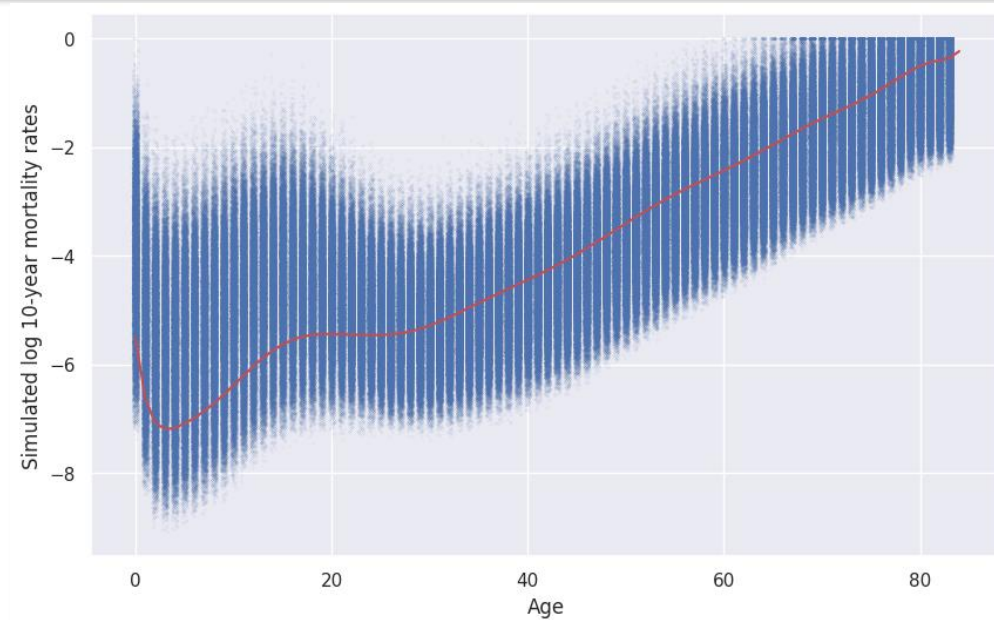
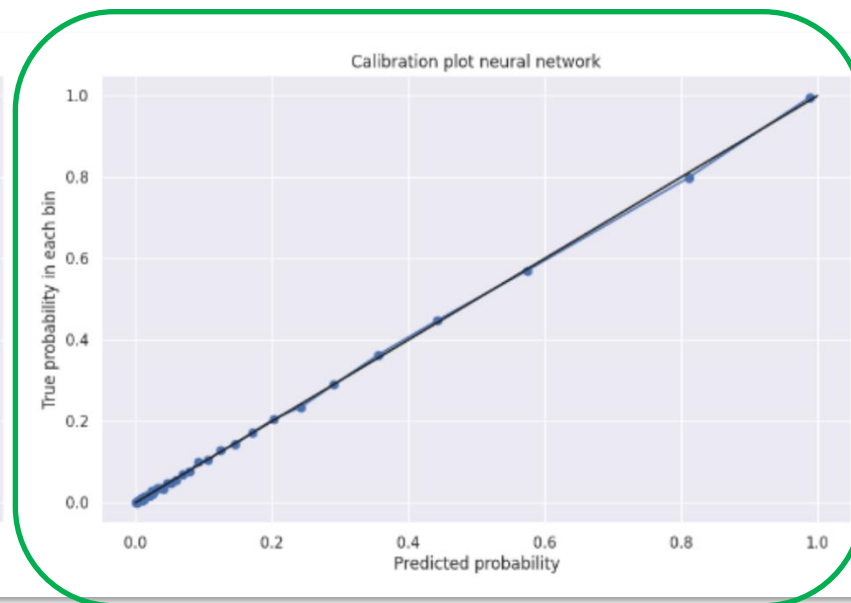
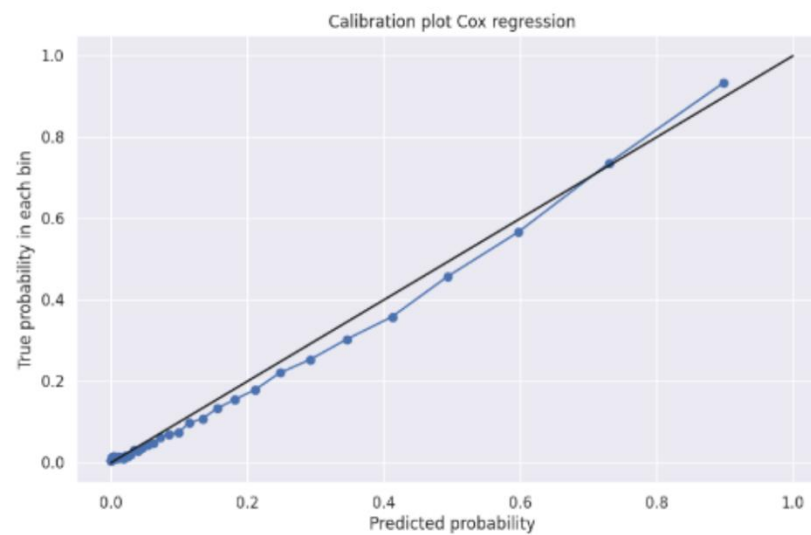
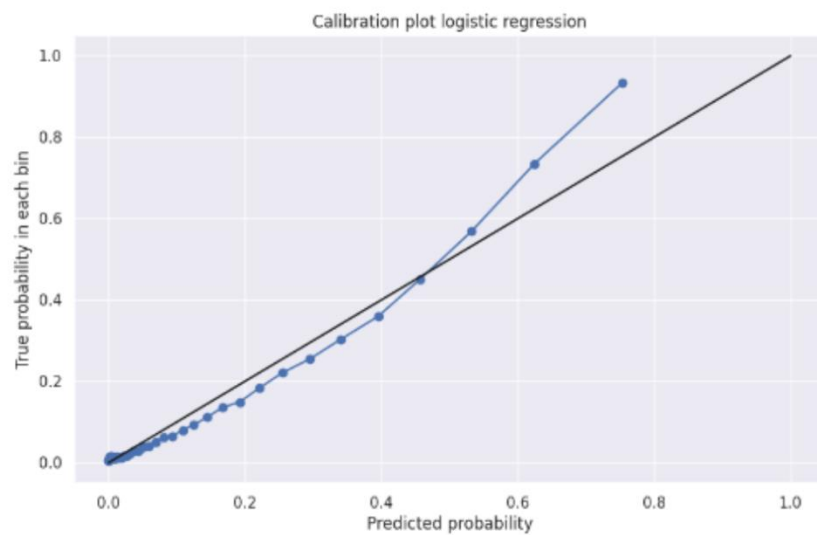
Table 3: Performance metrics on the test data subset of \mathcal{D}_1 .

Performance metric	logistic regression $\mu_1(\mathbf{x})$	Cox regression $\mu_2(\mathbf{x})$	neural net $\mu_3(\mathbf{x})$
ROC AUC	56.17%	56.17%	56.04%
MSE wrt $\log(\mu^*(\mathbf{x}))$	0.0016	0.0016	0.0057
Logistic deviance	9223.88	9223.88	9227.72

Table 4: Performance metrics on the test data subset of \mathcal{D}_2 .

Performance metric	logistic regression $\mu_1(\mathbf{x})$	Cox regression $\mu_2(\mathbf{x})$	neural net $\mu_3(\mathbf{x})$
ROC AUC	90.54%	90.55%	92.05%
MSE wrt $\log(\mu^*(\mathbf{x}))$	1.74	1.75	0.11
Logistic deviance	85383	83994	75732

Model performance



Asymmetric cryptography (public/private key)

- Create a *shared secret* s for symmetric encryption (**stream ciphers**: Salsa20, RC4, ..., **block ciphers**: AES, DES, RC5, ...):

Alice: “Secret message” $\rightarrow m \in (\mathbb{Z}/2\mathbb{Z})^n \rightarrow m + f(s)$ \Rightarrow **Bob**: $m + f(s) + f(s)$ \rightarrow “Secret message”
or $f(m, s)$ or $f(f(m, s), s)$

- Some examples:
 1. RSA (Rivest, Shamir, Adleman, 1977): **Factoring** large integers $n = pq$ (n public, p, q private)
 2. ElGamal (1985): **Discrete logarithm**, (multiplicative) group G , usually $G \subset (\mathbb{Z}/p\mathbb{Z})^* =: \mathbb{F}_p^*$ of order $q = (p-1)/2$ with generator g , solve $x = \log_g h$ (g, G, h public, x private)
 3. Elliptic curves methods (1985): **Discrete logarithm**, where group G is based on elliptic curves
 4. Lattice based methods, e.g., LWE (“learning with errors”, 2005): Solve $Ax + \varepsilon = b \bmod q$ for $x \in (\mathbb{Z}/q\mathbb{Z})^n$, where A is drawn uniformly from $(\mathbb{Z}/q\mathbb{Z})^{m \times n}$, $\varepsilon \in [-q/4, q/4]^m$ is drawn from a “non-trivial” distribution χ , and $b \in (\mathbb{Z}/q\mathbb{Z})^m$ (b, q, A public, x private)
 5. Many more examples from [NIST standardization](#) proposals for post-quantum cryptography (factorization and discrete logarithm can be calculated very efficiently on quantum computers), e.g., CRYSTALS, 2018, while [SIKE](#) had to be removed from the list in August 2022.

RSA

- [RSA, 1977](#) based on **Euler theorem**: $m^{\varphi(n)} = 1 \bmod n$ for $\gcd(m, n) = 1$, where $\varphi(n) = \#(\mathbb{Z}/n\mathbb{Z})^*$ (Euler's totient function)

- $n = p = 7$, $\varphi(p) = p - 1$

0	1	2	3	4	5	6
---	---	---	---	---	---	---

- $n = pq = 15$, $\varphi(pq) = (p - 1)(q - 1)$

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
---	---	---	---	---	---	---	---	---	---	----	----	----	----	----

- Choose “random” p, q, d with $\gcd(d, \varphi(pq)) = 1$, calculate e with $ed = 1 \bmod \varphi(pq)$ with extended Euclidean algorithm, e, n public key, p, q, d private key
- Encryption: message $m < n, m^e \bmod n$
- Decryption: $m^{ed} = m \bmod n$
- *Proof*: $m^{ed} = m^{k\varphi(n)+1} = m \bmod n$
- Calculating d from e and $n \Leftrightarrow$ calculating $\varphi(n) \Leftrightarrow$ **factoring** $n = pq$
- *Proof idea*: “ \Leftarrow ” 1. $\varphi(pq) = (p - 1)(q - 1)$, 2. extended Euclidean algorithm $ed + b\varphi(pq) = \gcd(d, \varphi(pq)) = 1$
“ \Rightarrow ” 1. $\varphi(pq) = -(p + q) + 1 \bmod n$, 2. $ed - 1 = k\varphi(n)$ sufficient to factor n (see, e.g., [Miller, 1975](#), ERH)
- There are attacks for, e.g., $q < p < 2q, 3d < n^{1/4}$ ([Wiener, 1990](#)) and several others ([Zhang, 1999](#))
- Homomorphic encryption RSA example: $(m_1 m_2)^e = m_1^e m_2^e \bmod n$, in general: $\text{enc}(\text{op}_1(m_1, m_2)) = \text{op}_2(\text{enc}(m_1), \text{enc}(m_2))$

TUNE INSIGHT

**Confidential
collaborative analytics
and machine learning**

Juan R. Troncoso, co-founder & CEO

juan@tuneinsight.com

RSA uses large finite (abelian) groups

$G = (\mathbb{Z}/n\mathbb{Z})^x$ (2048 bits, 4096 bits,...)

To speed things up:

- **Elliptic curve crypto** uses **smaller** groups, whose operations are more expensive.
- **Lattice cryptography** uses **larger** groups, but whose operations are much cheaper.

Lattice-based cryptography is the use of conjectured hard problems on point lattices in \mathbb{R}^n as the foundation for secure cryptographic systems.

Features:

- Apparent resistance to *quantum* attacks (in contrast with most number-theoretic cryptography) [Sho97]
- High asymptotic efficiency and parallelism
- Security under *worst-case* intractability assumptions [Ajt96]
- Versatile and powerful cryptographic objects (FHE [Gen09], ABE [BGG+14], Code obfuscation [GGH+13]...)

Main Milestones in Lattice Cryptography

1982: First use of lattices in cryptanalysis (LLL): knapsack cryptosystems

1996: First crypto schemes based on hard lattice problems: NTRU, Ajtai-Dwork, GGH,...

2009: Fully-Homomorphic Encryption on Lattices

2012: Leveled cryptosystems

What is a Lattice?

A lattice is the set of all integer linear combinations of (linearly independent) basis vectors

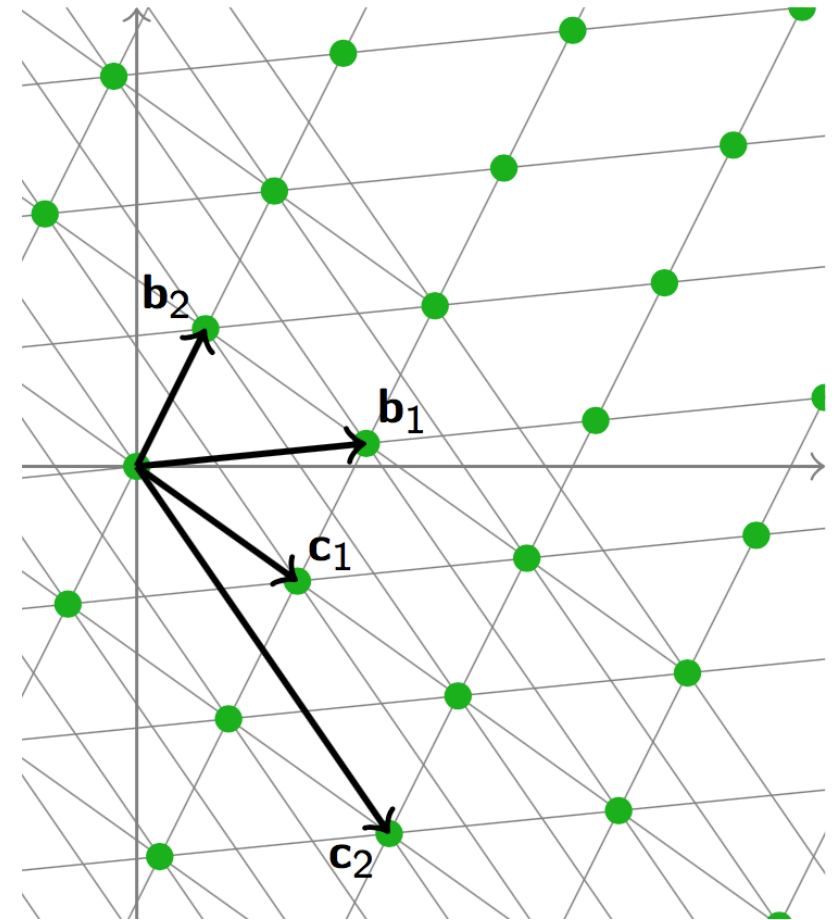
$$\mathbf{B} = \{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n\} \subset \mathbb{R}^n$$

$$\mathcal{L} = \sum_{i=1}^n \mathbf{b}_i \cdot \mathbb{Z} = \{\mathbf{B}\mathbf{x} : \mathbf{x} \in \mathbb{Z}^n\}$$

The same lattice has many different bases:

$$\mathcal{L} = \sum_{i=1}^n \mathbf{c}_i \cdot \mathbb{Z}$$

Lattice: discrete additive subgroup of \mathbb{R}^n



Simple Example (Preliminary Homomorphic Encryption)

Good bases and bad bases: GGH (Goldreich, Goldwasser, Halevi) family

Two lattice bases

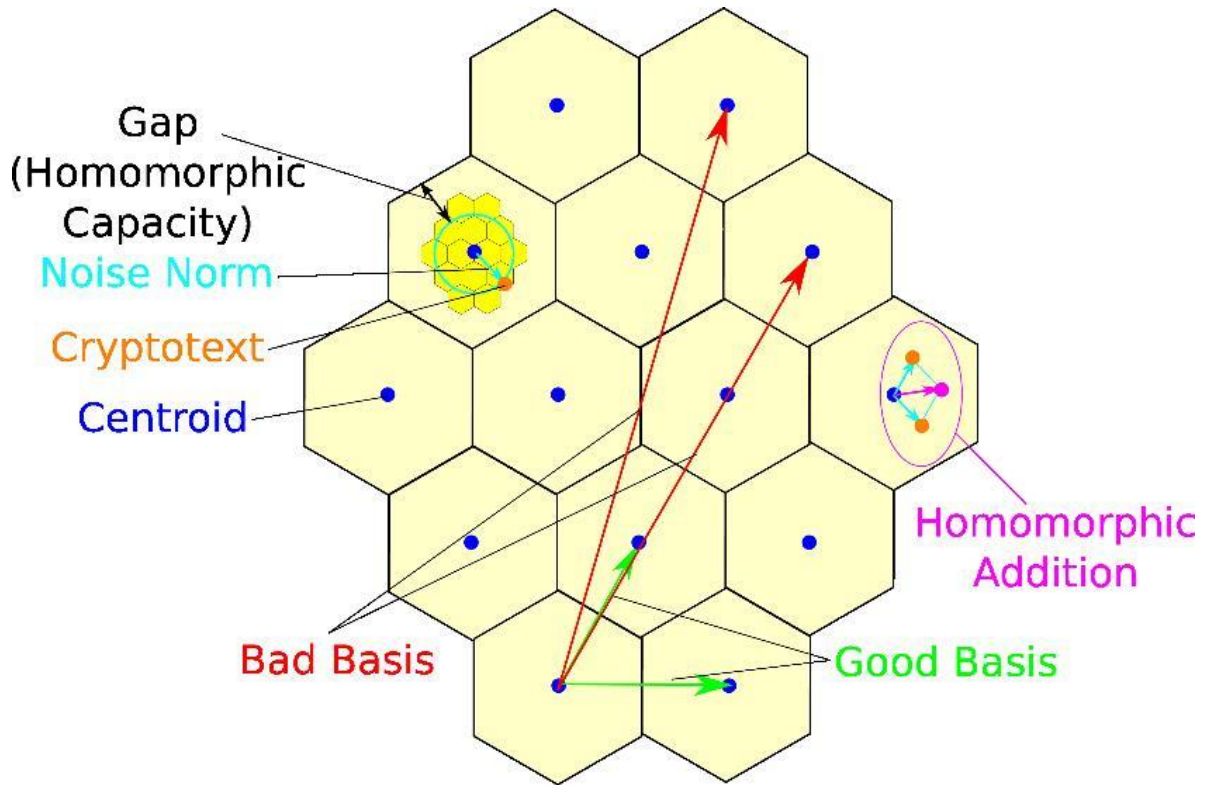
- “Good” basis (\mathbf{B} , private key)
- “Bad” basis (\mathbf{H} , public key, Hermite Normal Form)

Encryption of m : $\mathbf{c} = E(m) = \mathbf{v} + \mathbf{n}[m]$ (lattice point + noise)

Decryption: $D(\mathbf{c}): \hat{\mathbf{v}} = \mathbf{B}[\mathbf{B}^{-1}\mathbf{c}]$

Homomorphism:

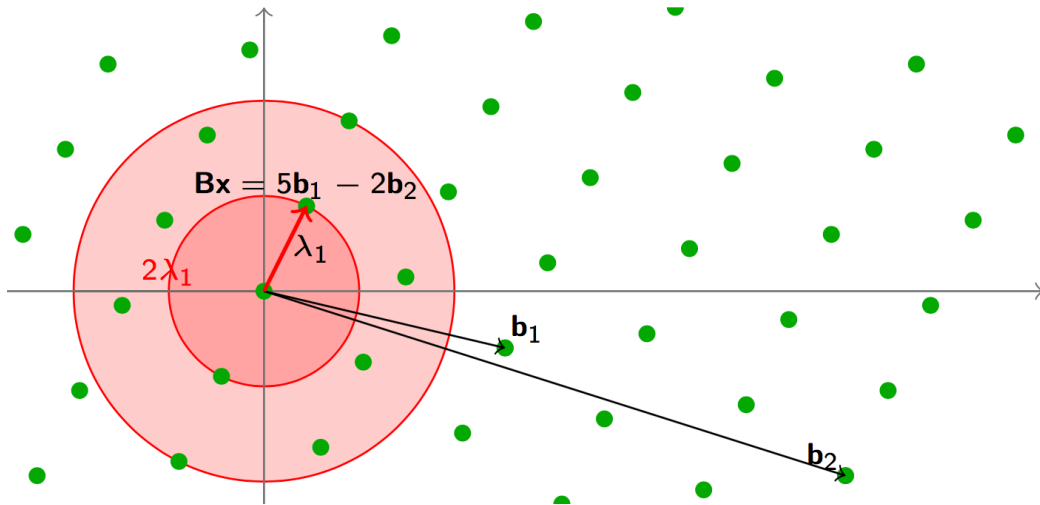
$$\begin{aligned}\mathbf{c}_1 + \mathbf{c}_2 &= (\mathbf{v}_1 + \mathbf{n}[m_1]) + (\mathbf{v}_2 + \mathbf{n}[m_2]) \\ &= \mathbf{v}_3 + \mathbf{n}[m_1 + m_2]\end{aligned}$$



Base Lattice Problems (ex: SVP, CVP)

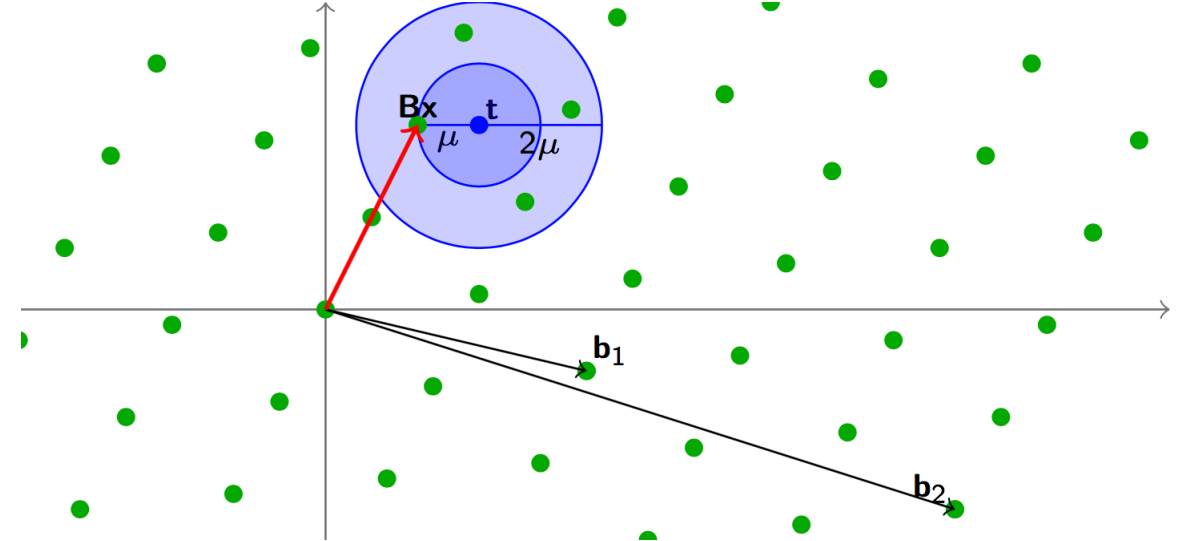
Shortest Vector Problem SVP_γ

Given a lattice $\mathcal{L}(\mathbf{B})$, find a (nonzero) lattice vector \mathbf{Bx} , $\mathbf{x} \in \mathbb{Z}^k$ of length (at most) $\|\mathbf{Bx}\| \leq \gamma \lambda_1$



Closest Vector Problem CVP_γ

Given a lattice $\mathcal{L}(\mathbf{B})$ and a target point \mathbf{t} , find a lattice vector \mathbf{Bx} within distance $\|\mathbf{Bx} - \mathbf{t}\| \leq \gamma \mu$



Ring-LWE distribution: For an $s \in R_q$ (the secret), the ring-LWE distribution $A_{s,\chi}$ over $R_q \times R_q$ is sampled by choosing $a \in R_q$ uniformly at random, $e \leftarrow \chi$, and outputting
$$(a, b = s \cdot a + e \bmod q)$$

Decision-R-LWE: Given m independent samples $(a_i, b_i) \in R_q \times R_q$ where every sample is distributed according to either:

- $A_{s,\chi}$ for a uniformly random $s \in R_q$ (fixed for all samples),
- The uniform distribution

Distinguish which is the case (with non-negligible advantage)

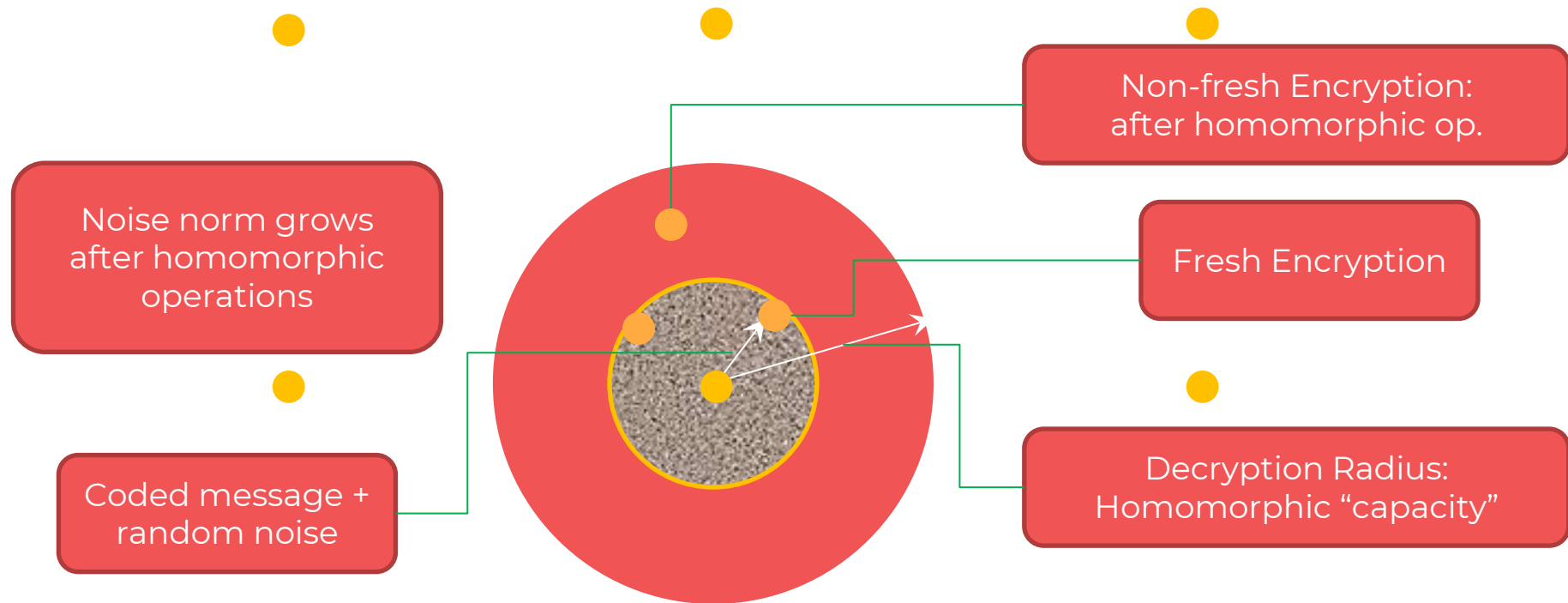
Normal form: secret from $s \leftarrow \chi$

More efficient than LWE (smaller m and FFT-like polynomial products)

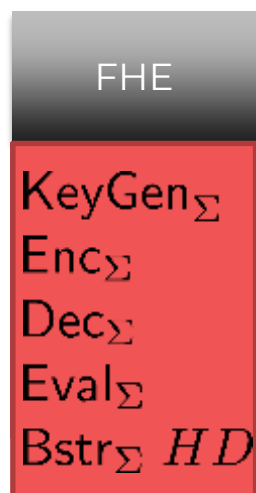
Reduction of $RLWE_{q,\chi,m}$ to quantum SVP_γ [LPR10]

How to build Homomorphic Cryptosystems from RLWE

Noise management is essential in homomorphic cryptosystems

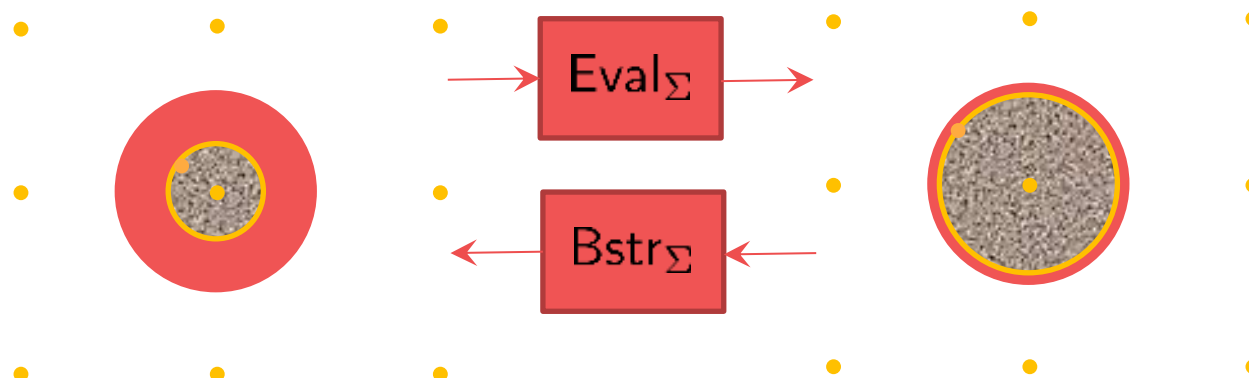


How to build Homomorphic Cryptosystems from R-LWE (Somewhat vs Fully HE)



$$c = \text{Eval}_\Sigma(pk, f, (c_1, \dots, c_n)) \equiv \text{Enc}_\Sigma(pk, f(m_1, \dots, m_n))$$

Only valid when f is of depth $< L$



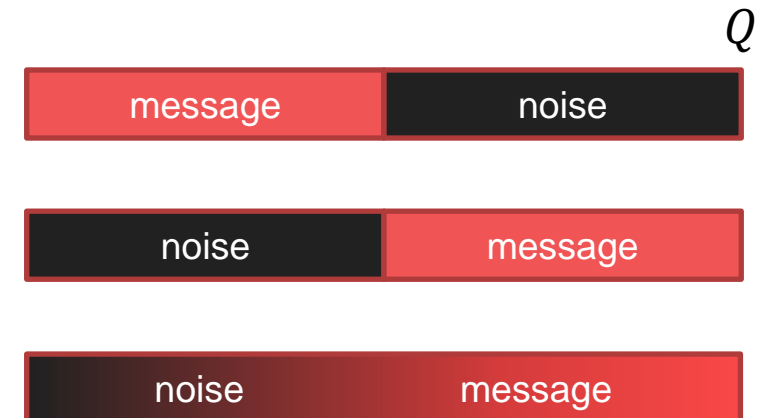
If Dec (squashed) has depth $< L$

$$\text{Bstr}_\Sigma(pk, \llbracket m \rrbracket_{pk}) = \text{Eval}_\Sigma(pk, \text{Dec}_\Sigma, \llbracket m \rrbracket_{pk})$$

Common characteristics of modern RLWE cryptosystems:

- Cyclotomic polynomial $f(x) = x^n + 1$, n power of two
- Ciphertext modulus $Q = \prod q_i$
- Ring $R_Q = \mathbb{Z}_Q / (f(x))$
- Error distribution χ with power $\|\chi\| < B$
- Plaintext modulus $\ll Q$, scale factor Δ
- Key generation:
- Secret key: $s \leftarrow \chi$
- Public key is an RLWE sample: E.g., $(a_0 = -(a_1 s + e), a_1)$, with $a_1 \leftarrow R_Q$, $e \leftarrow \chi$
- Encryptions are vectors of polynomials in R_Q , with the encoded message
- The decryption function is of the form

$$\sum_{i=0}^v c_i s^i$$



Efficiently using lattice cryptosystems: packing in the coefficient domain

All the encryptions over RLWE work with polynomials of degree d

Each coefficient is a plaintext slot in \mathbb{Z}_Q : $\mathbf{a} \equiv a \equiv \sum_{i=0}^{d-1} a_i x^i$

SIMD homomorphic operations:

Polynomial addition

$$E(\mathbf{a}) + E(\mathbf{b}) = E\left(\sum_{i=0}^{d-1} a_i x^i\right) + E\left(\sum_{i=0}^{d-1} b_i x^i\right) = E\left(\sum_{i=0}^{d-1} (a_i + b_i) x^i\right)$$

Polynomial multiplication (modular polynomial $f(x) = x^d + 1$)

$$E(\mathbf{a}) \cdot E(\mathbf{b}) = E\left(\sum_{i=0}^{d-1} a_i x^i\right) \cdot E\left(\sum_{i=0}^{d-1} b_i x^i\right) = E\left(\sum_{i=0}^{d-1} \left(\sum_{j=0}^i (a_j \cdot b_{i-j})\right) x^i - \sum_{i=0}^{d-2} \left(\sum_{j=i+1}^{d-1} (a_j \cdot b_{d+i-j})\right) x^i\right)$$

Nega-cyclic homomorphic convolution: $E(c') = E(a') \cdot E(b')$

After decryption: $c = \sum_{i=0}^{d-1} c'_i (-1)^{-\frac{i}{d}} x^i \Rightarrow c = \sum_{i=0}^{d-1} \left(\sum_{j=0}^i (a_j \cdot b_{i-j})\right) x^i + \sum_{i=0}^{d-2} \left(\sum_{j=i+1}^{d-1} (a_j \cdot b_{d+i-j})\right) x^i$

Efficiently using lattice cryptosystems: packing in the slot domain

Use an automorphism as message coding that switches domain

Equivalent to an NTT (Number Theoretic Transform) or DFT (Discrete Fourier Transform)

Ex. DFT (for inputs $x, X \in \mathbb{C}^d$)

$$DFT[x] = X_k = \sum_{n=0}^{d-1} x_n \cdot e^{-\frac{j2\pi kn}{d}}$$

$$DFT^{-1}[X] = x_n = \frac{1}{d} \sum_{k=0}^{d-1} X_k \cdot e^{\frac{j2\pi kn}{d}}$$

Important properties:

- Circular convolution: $DFT^{-1}[X \cdot Y]_n = \sum_{i=0}^n x_i y_{n-i} + \sum_{i=n+1}^{d-1} x_i y_{d+n-i} = x_n \circledast y_n$
- Duality: $DFT[x \cdot y]_k = \frac{1}{d} \left(\sum_{i=0}^n X_i Y_{n-i} + \sum_{i=n+1}^{d-1} X_i Y_{d+n-i} \right) = \frac{1}{d} X_k \circledast Y_k$
- Parseval's theorem: $\sum_{n=0}^{d-1} x_n y_n^* = \sum_{k=0}^{d-1} X_k Y_k^*$

Homomorphic operations become component-wise when the message is in the slot domain



Function evaluation: polynomial approximations

Ring operations are additions and products

Non-polynomial functions have to be:

- approximated by a polynomial
- run on universal gates nand / xor with binary arithmetic

Let $f(x): [a, b] \subset R \rightarrow R$, with $c \in [a, b]$

- **Taylor approximation:** Error bounded, but not uniform in $[a, b]$
Preferred when input distribution is denser around c (e.g., Gaussian)
- **Least-squares approximation:** minimizes average square error in $[a, b]$
Preferred when input is uniformly distributed and high homomorphic capacity
- **Chebyshev approximation:** Bounded maximum error, converges with d to the minimax polynomial that minimizes this maximum approximation error $[a, b]$.
Preferred to avoid overflows and for numerical stability

Practical example: evaluating a Logistic regression under encryption with Tune Insight's Python cryptolib

Evaluation of a Logistic regression prediction

$$y_i = \mu \left(\beta_0 + \sum_{j=1}^k a_{i,j} \beta_j \right)$$

For a dataset with l records and k features

Practical example: evaluating a Logistic regression under encryption with Tune Insight's Python cryptolib

1. Parameterization and cryptosystem instantiation

```
from tuneinsight.cryptolib.hefloat import hefloat

# Parameterization: scale/precision and circuit depth
log_scale = 45                                # Fixed-point arithmetic floating point scaling factor in bits (log2(Delta))
levels = 7                                    # Circuit depth
log_qi = [log_scale+5] + levels*[log_scale]    # 5 additional bits for the lowest level, to account for plaintext growth
log_pi = [log_scale+5]                        # Auxiliary module used for relinearization (usually, at least of the same size as the lowest level q0)

# In order to generate an instance of the cryptosystem, the RLWE ring degree is automatically chosen to ensure at least 128-bit of security
# A context stores the scheme cryptographic parameters and a key generator
context = hefloat.new_context(log_qi = log_qi, log_default_scale= log_scale, log_pi = log_pi)

#Print some information about the cryptographic parameters
print(f'Log2 N: {context.parameters.log_n()}')
print(f'Log2 Moduli Chain: Q{log_qi} + P{log_pi}')
print(f'Log2 QP: {context.parameters.log_q() + context.parameters.log_p()}')
print(f'Log2 Slots: {context.parameters.log_slots()}')
print(f'Available Depth: {levels}')
```

```
Log2 N: 14
Log2 Moduli Chain: Q[50, 45, 45, 45, 45, 45, 45, 45] + P[50]
Log2 QP: 414.9999999903431
Log2 Slots: 13
Available Depth: 7
```

Practical example: evaluating a Logistic regression under encryption with Tune Insight's Python cryptolib

2. Key generation

```
# Generate a fresh secret key
sk = context.new_secret_key()

# Instantiate an evaluator with a relinearization key
# The relinearization key is at public-evaluation key required to ensure ciphertext x ciphertext compactness
# The resulting evaluator object contains only public information and can be freely shared
evaluator = context.new_evaluator(context.new_relinearization_key(sk))
```

3. Polynomial approximation of the activation function

```
import numpy.polynomial.chebyshev as chebyshev
import numpy as np

# Expected interval of the encrypted values after the scalar product
a = -12
b = 12

# Interpolates the Sigmoid in the interval [-12, 12] and returns the coefficients
# for the Chebyshev approximation polynomial in the Chebyshev basis
coeffs = chebyshev.chebinterpolate(lambda x: 1/(1+np.exp(-((b-a)/2 * x + (b+a)/2))), 63)
```

Practical example: evaluating a Logistic regression under encryption with Tune Insight's Python cryptolib

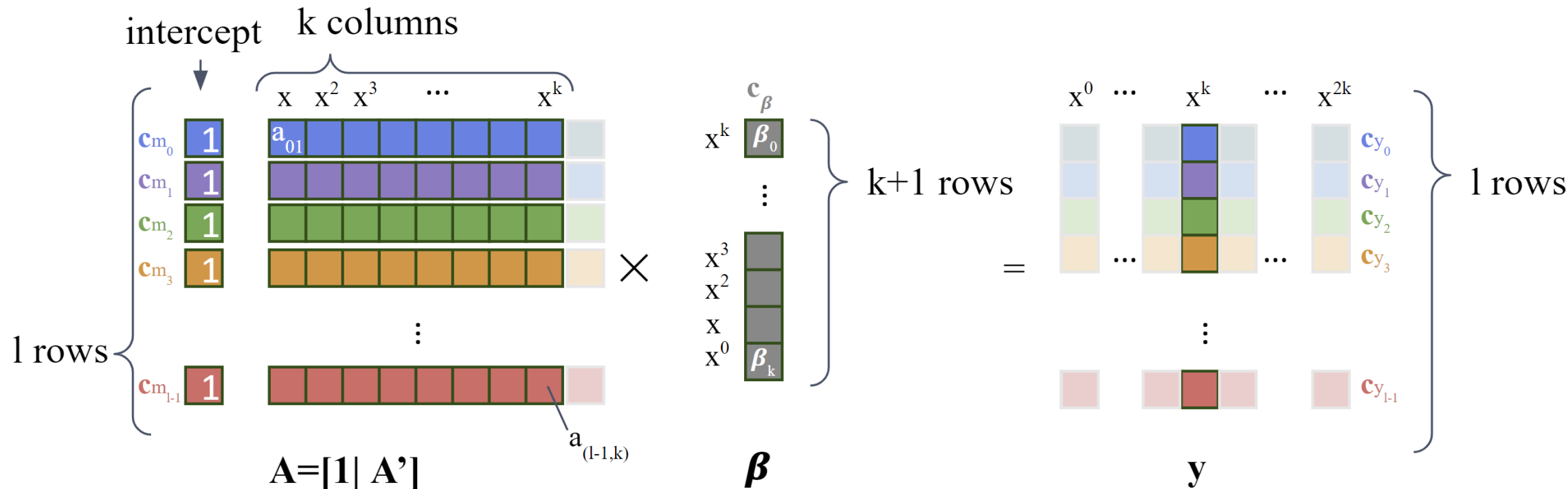
4. Synthetic data generation

```
## Synthetic data generation:  
# Number of samples to process in parallel (available plaintext slots that one encryption can hold)  
batch_size = context.slots()  
  
# Number of features (k=200)  
features = 200  
  
# Generate random data in [-0.5, 0.5]. This is the matrix A'  
data = np.random.rand(batch_size, features)-0.5  
  
# Generate random regression weights in [0, 1]. These represent beta_i, i=1,...,k  
weights = np.random.rand(features, 1)  
  
# Generate random bias (intercept coefficient) in [0, 1]. This represents beta_0  
bias = np.random.rand(1)
```

Practical example: evaluating a Logistic regression under encryption with Tune Insight's Python cryptolib

5. Packed encryption of all inputs

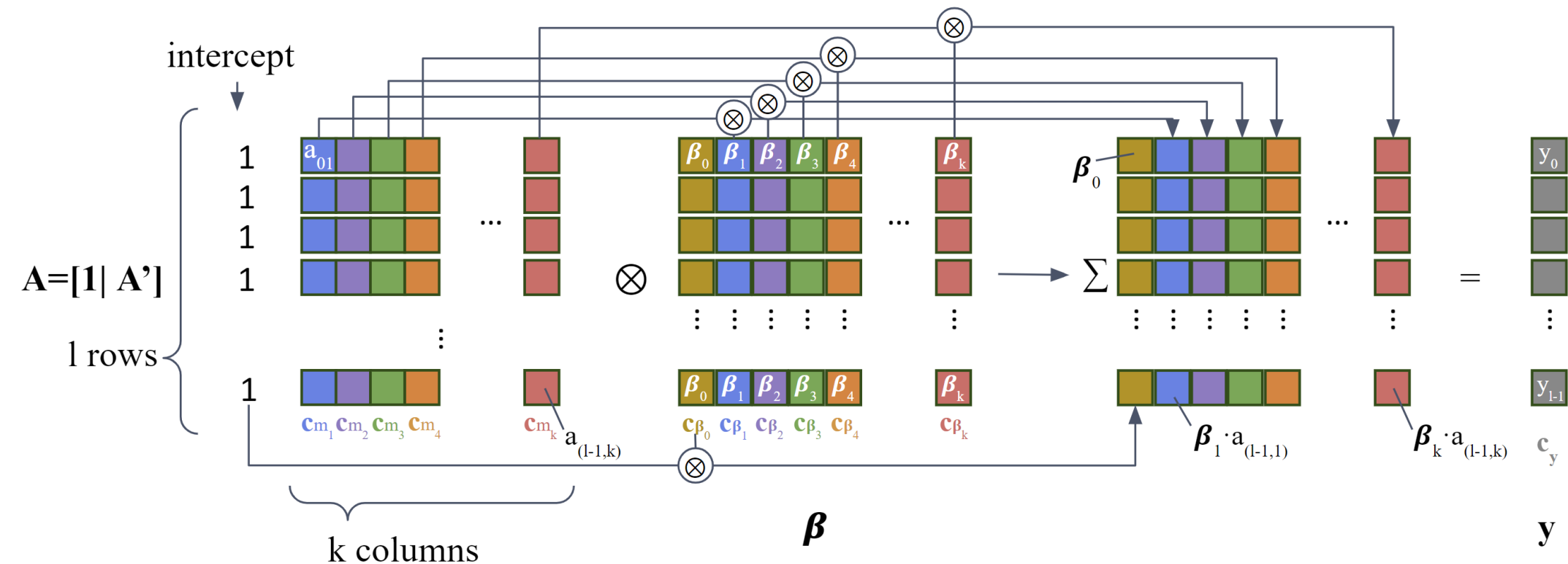
Option 1: horizontal packing



Practical example: evaluating a Logistic regression under encryption with Tune Insight's Python cryptolib

5. Packed encryption of all inputs

Option 2: vertical packing



Practical example: evaluating a Logistic regression under encryption with Tune Insight's Python cryptolib

5. Packed (batched) encryption of all inputs

```
# This optional parameter defines whether the input vectors will be encoded in the coefficients domain (if batched=False)
# or in the slots domain (if batched=True). The latter is the default behavior, and it enables component-wise homomorphic operations
# (additions and products)
batched = True

# The encrypt function can receive a two-dimensional matrix as input, in which case it encrypts each row of the input matrix in one ciphertext.
# Therefore, we transpose the input A', in order to encrypt each column of A' in one ciphertext.
# We need to explicitly make a copy to ensure a correct memory
# alignment when passing C pointers of arrays to the Go wrapper.
# The function returns an object that stores a vector of ciphertexts.
encrypted_data = context.encrypt(data.transpose().copy(), sk, batched)

# As for the regression coefficients, we encrypt each of the weights replicated in all slots of the corresponding ciphertext.
# For this, we apply repetition coding (with tile) and pass the resulting matrix as input to the encrypt function, so that each row is encrypted in a separate ciphertext.
# The result is an object that stores a vector of ciphertexts, each containing one regression coefficient replicated in all its slots.
encrypted_weights = context.encrypt(np.tile(weights, (1, batch_size))* 2/(b-a), sk, batched)

# The intercept coefficient or bias is also encrypted in its own ciphertext, with the same repetition coding as all the other regression coefficients
encrypted_bias = context.encrypt(np.tile(bias, (1, batch_size))* 2/(b-a) + (-a-b)/(b-a), sk, batched)
```

Practical example: evaluating a Logistic regression under encryption with Tune Insight's Python cryptolib

6. Homomorphic evaluation of the model prediction under encryption

```
# Encrypted evaluation of data @ weights computed as np.sum(data.transpose() * np.tile(bias, (1, batch_size)), axis=0)
# This is faster, but equivalent, to doing evaluator.sum(evaluator.mul(encrypted_data, encrypted_weights), axis=0)
encrypted_scalar_product = evaluator.scalar_product(encrypted_data, encrypted_weights)

# Encrypted evaluation of data @ weights + bias
encrypted_scalar_product_plus_bias = evaluator.add(encrypted_bias, encrypted_scalar_product)

# Encrypted evaluation of sigmoid(data @ weights + bias)
encrypted_prediction = evaluator.polynomial(encrypted_scalar_product_plus_bias, coeffs=coeffs, basis="Chebyshev")
```


Practical example: evaluating a Logistic regression under encryption with Tune Insight's Python cryptolib

7. Decryption of results

```
# Decrypts the values
prediction = context.decrypt(encrypted_prediction, sk)[: , :batch_size]
```

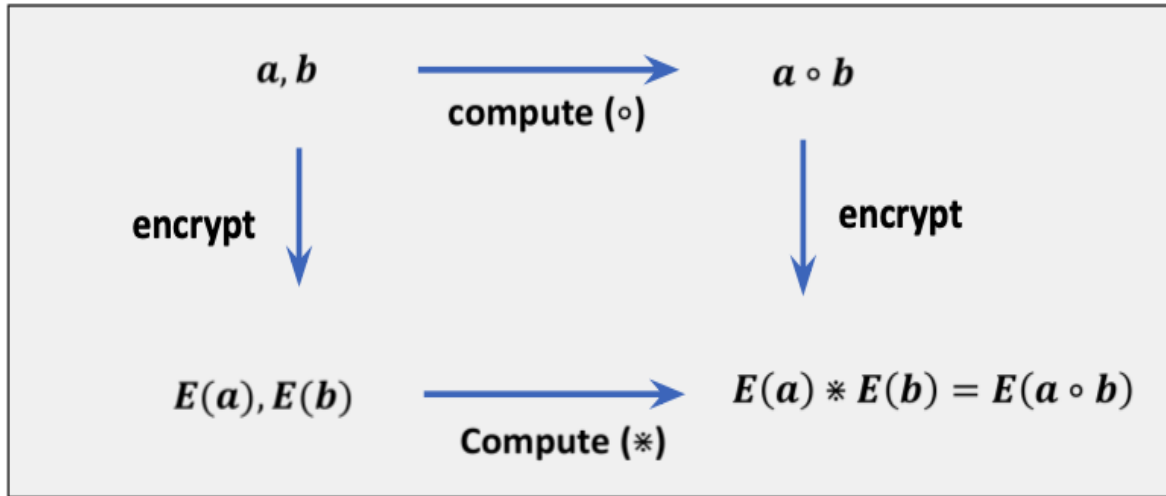
8. Accuracy comparison with the clear-text process

```
from math import log
# Finally, we evaluate the plaintext circuit
clear_target = 1/(np.exp(-(data @ weights + bias))+1)

# And compare with the decrypted result
print(f'Obtained: {prediction}')
print(f'Clear_tg: {clear_target.transpose()}')
print(f'Precision as -log2(avg_l2(obtained-clear_tg))): {-log(np.sqrt(np.sum((prediction-clear_target.transpose())**2))/batch_size, 2)}')
```

```
Obtained: [[0.87919843 0.17177785 0.22382661 ... 0.99488169 0.06266606 0.45500863]]
Clear_tg: [[0.87919843 0.17177785 0.22382656 ... 0.99488167 0.06266607 0.45500862]] Precision as -log2(avg_l2(obtained-clear_tg))):
31.818855691640856
```

Recap on Homomorphic Encryption



1. Put your gold in a locked box.
2. Keep the key.
3. Let your jeweler work on it through a glove box.
4. Unlock the box when the jeweler is done!

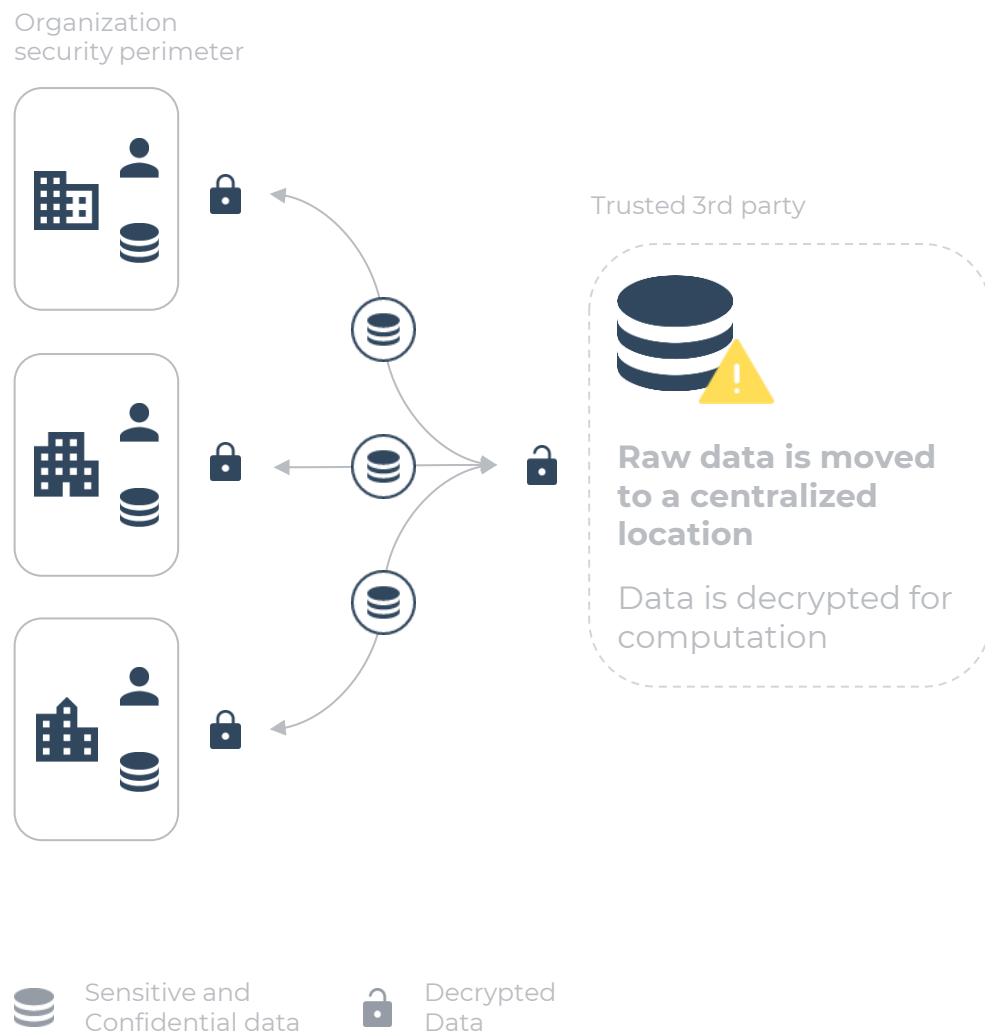


Figure from Prof. Kristin Lauter ("Private AI: Machine Learning on Encrypted Data", 2021)

**Homomorphic encryption enables computations directly on encrypted data:
"compute on the data without seeing the data"**

...but what happens if the raw data cannot be moved or centralized?

Data collaborations: Centralized approach



- **Single point of failure at the central database**
- **Individual sites lose control over their data**
- **Not always feasible across jurisdictions**

Organization
security perimeter



- **Requires trust on the aggregation server**
- **Vulnerable to re-identification and reconstruction attacks**
 - B. Hitaj, G. Ateniese, and F. Perez-Cruz. Deep models under the GAN: Information leakage from collaborative deep learning. In ACM CCS, 2017.
 - Z. Wang, M. Song, Z. Zhang, Y. Song, Q. Wang, and H. Qi. Beyond inferring class representatives: User-level privacy leakage from federated learning. In IEEE INFOCOM, 2019.
 - L. Zhu, Z. Liu, and S. Han. Deep leakage from gradients. In NIPS. 2019.
 - L. Melis, C. Song, E. De Cristofaro, and V. Shmatikov. Exploiting unintended feature leakage in collaborative learning. In IEEE S&P, 2019.
 - M. Nasr, R. Shokri, and A. Houmansadr. Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning. In IEEE S&P, 2019.



Sensitive and
Confidential data



Decrypted
Data



Local Partial
Results

Secure Multiparty Computation (SMC)

Problem statement:

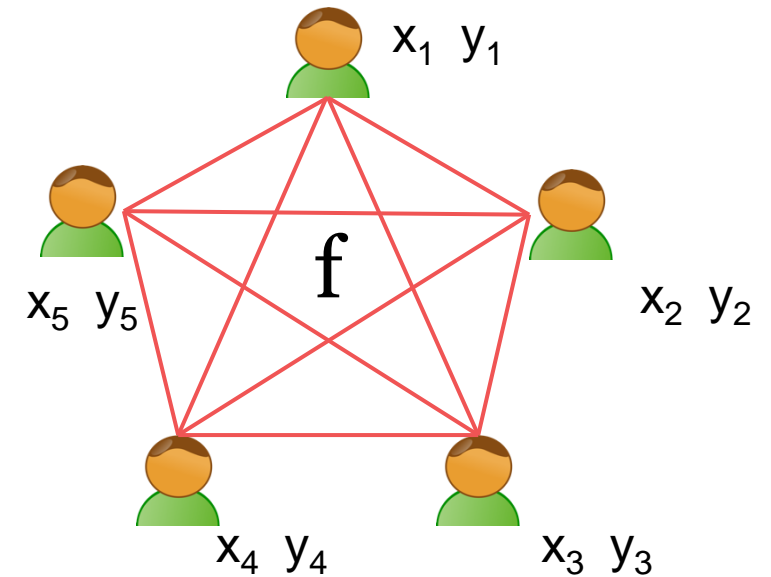
A set of players $\mathcal{P} = \{P_1, P_2, \dots, P_N\}$ would like to compute a function $f(x_1, x_2, \dots, x_N) = (y_1, y_2, \dots, y_N)$ of their joint inputs.

Requirements:

1. *Privacy*
No party should learn anything more than its prescribed output
2. *Correctness*
Each party is guaranteed that the output that it receives is correct

Realization:

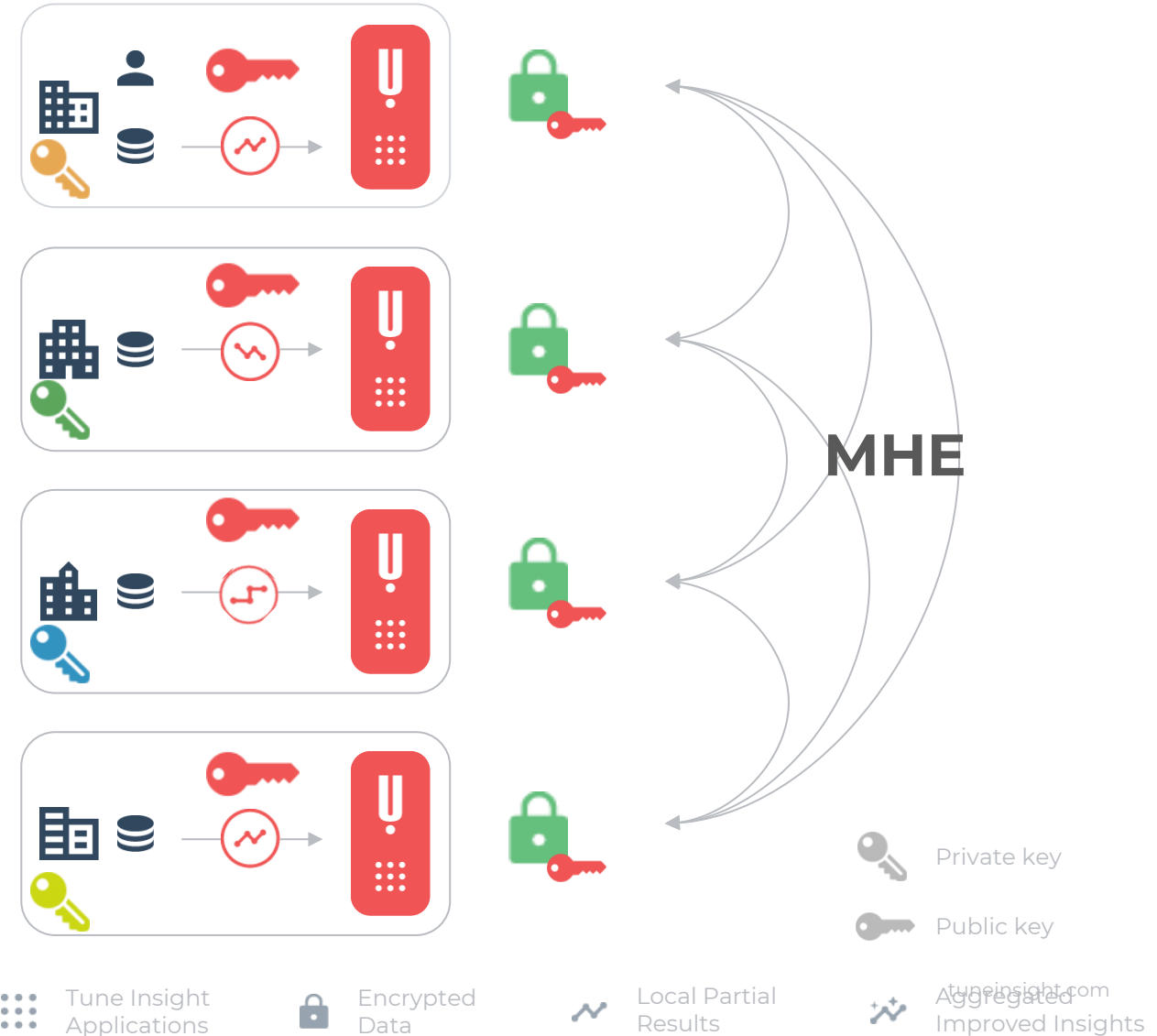
An (interactive) multiparty cryptographic protocol



MHE (Multiparty Homomorphic Encryption)

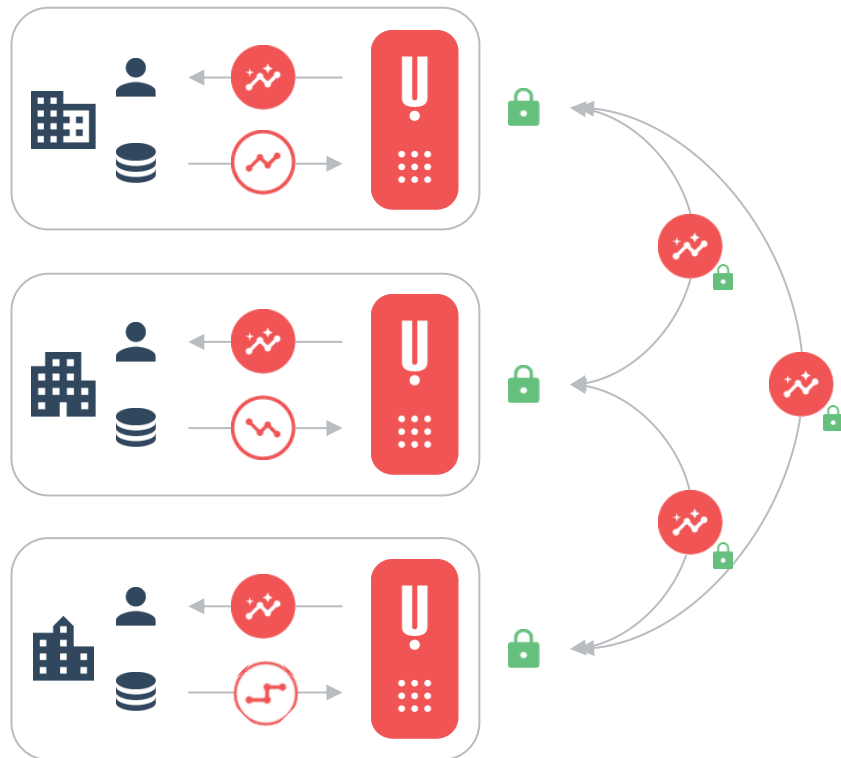
Combination of:

- HE**
encrypted computation → **Security and Privacy**
 - SMC**
collective key → **Trust Distribution**
 - FL**
Iterative local & collective training → **Data Minimization**
- ✓ **Policy enforcement embedded in the protocol**
 - ✓ **Raw data does not move**
 - ✓ **Computation is encrypted end-to-end**



Data collaborations: Secure and distributed approach

Organization
security perimeter



-  Minimization of transfers
- Always aggregated & encrypted data
- Computation over encrypted data
- Controlled computation

 Tune Insight
Applications

 Encrypted
Data

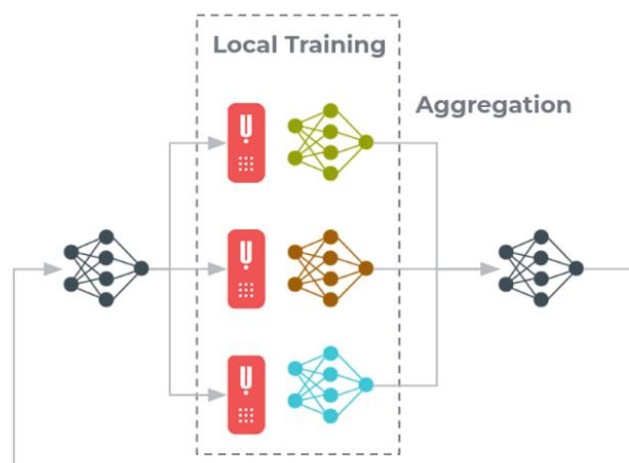
 Local Partial
Results

 Aggregated
Collective Insights

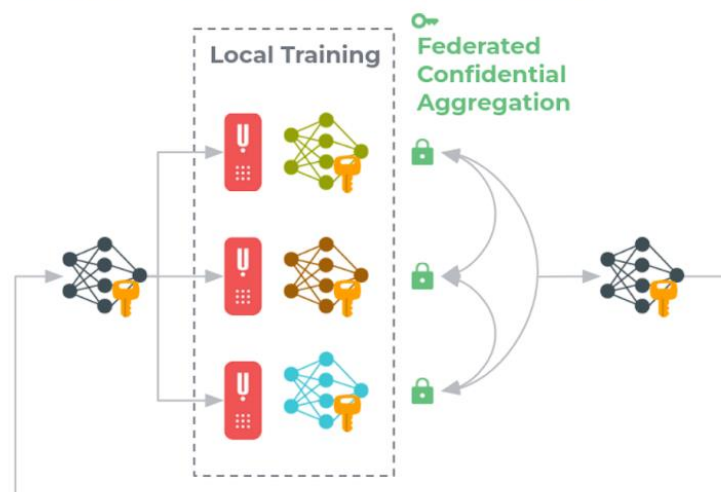
Practical example: training mortality models on D_2 dataset with federated data using Tune Insight's platform and Python SDK

Comparison of three scenarios

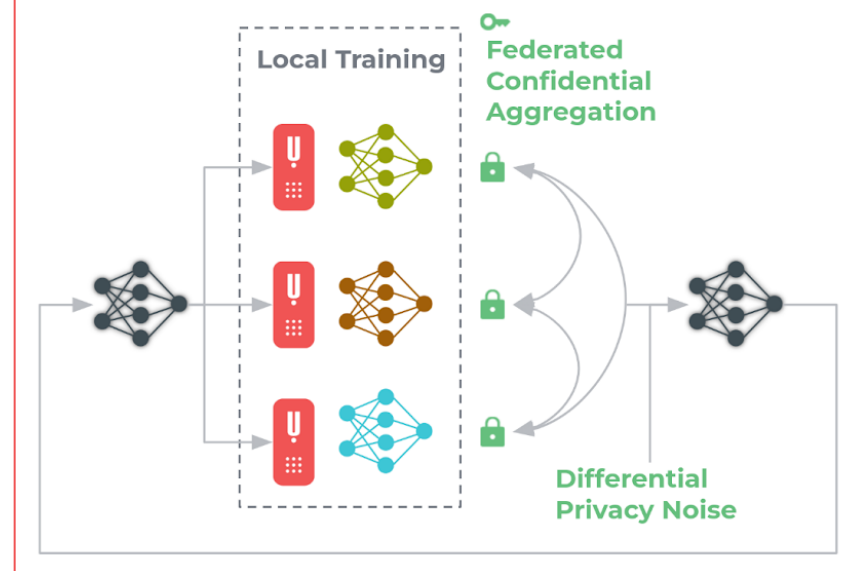
Federated Learning



Encrypted Federated Learning



Hybrid Federated Learning



Practical example: training mortality models on D_2 dataset with federated data using Tune Insight's platform and Python SDK

1. Model parameter definition (Cox and Logistic Regression)

```
task_id_cox = 'mortality_cox'
task_def_cox = {
    "n_inputs": 11,
    "n_classes": 2,
    "model": {
        "type": "cox",
        "pretrained": False,
    },
    "shuffle_data": True,
    "balance_train_classes": True,
    "batch_size": 2048,
    "drop_last_batch": True,
    "loss_criteria": "cross_entropy_loss"
}
```

```
task_id_logreg = 'mortality_logreg'
task_def_logreg = {
    "n_inputs": 11,
    "n_classes": 2,
    "model": {
        "type": "logreg",
        "pretrained": False,
    },
    "shuffle_data": True,
    "balance_train_classes": True,
    "batch_size": 2048,
    "drop_last_batch": True,
    "loss_criteria": "cross_entropy_loss"
}
```

Practical example: training mortality models on D_2 dataset with federated data using Tune Insight's platform and Python SDK

Secure Federated Learning workflow: Training parameters

```
learning_params_federated_secure = models.HybridFLLearningParams(  
    fl_rounds = 2,  
    local_epochs = 1,  
    num_workers = 8,  
    batch_size = 2048,  
    learning_rate = 0.01,  
    momentum = 0.9,  
    gradient_clipping = 0.1,  
    epsilon = 1,  
    delta = 0.001,  
)
```

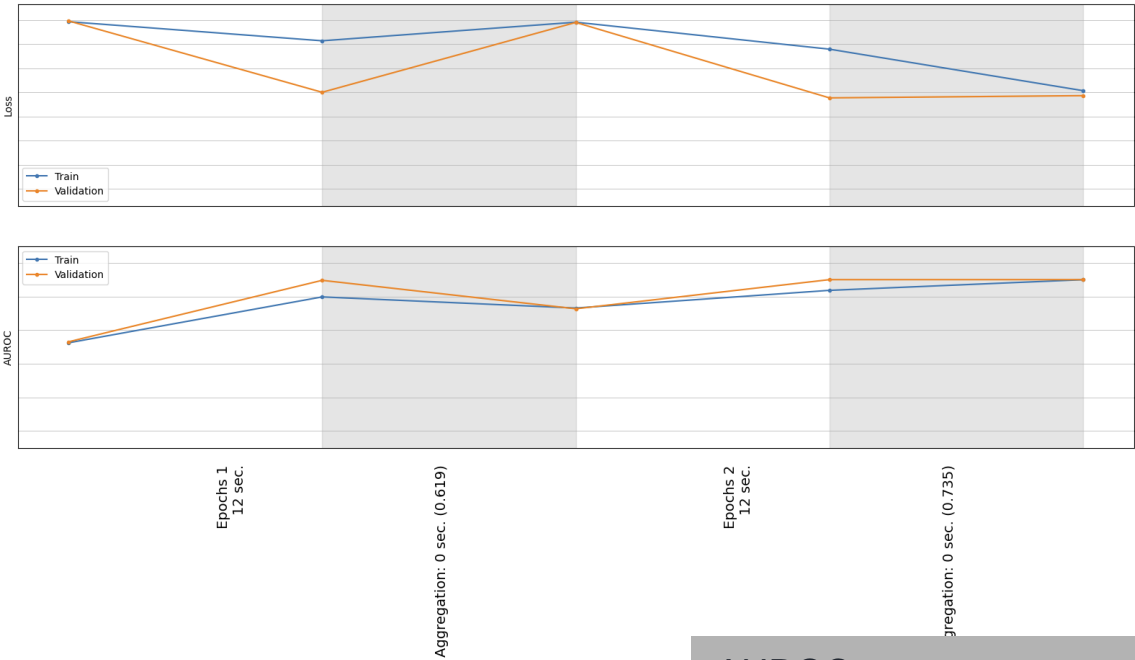
Running the computation

```
hybrid_fl = project.new_hybrid_fl()  
hybrid_fl.max_timeout = 300 * 60 * time.second  
_ = hybrid_fl.create_from_params(task_id=task_id_cox, learning_params=learning_params_federated_secure, task_def=task_def_cox)
```

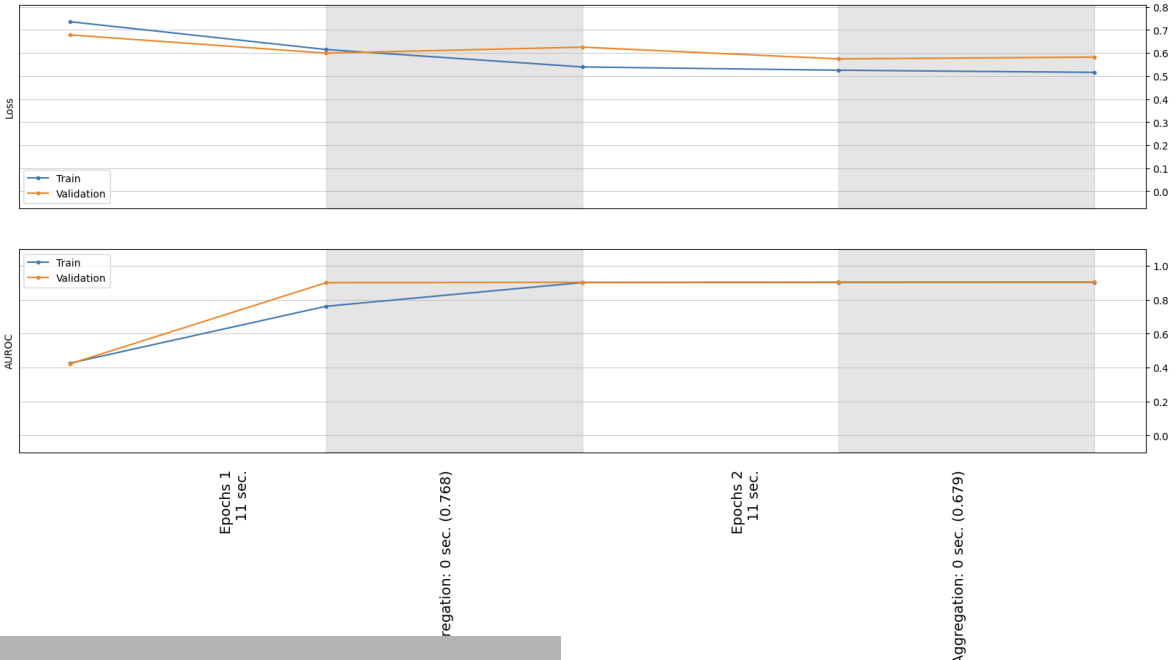
Practical example: training mortality models on D_2 dataset with federated data using Tune Insight's platform and Python SDK

Training performance

Cox Regression



Logistic Regression



AUROC		
	Cox regression	Logistic regression
Local	0.916300	0.905600
Federated	0.916500	0.904500
Secure Federated	0.900600	0.904500

“Technical solutions such as multiparty homomorphic encryption (MHE) that combine these three technical measures while still allowing for the possibility to query and analyse encrypted data without decrypting it have significant potential to **provide effective security measures that facilitate cross-borders transfers of personal data in high-risk settings.**”

Compagnucci et al., Supplementary Measures and Appropriate Safeguards for International Transfers of Personal Data after Schrems II (February 23, 2022). <https://ssrn.com/abstract=4042000>

Contact us for a full analysis of the platform benefits and risk minimization, addressing the relevant GDPR recitals.



Article 25 Data protection by design and by default	Article 32 Security of processing	Article 33 Breach notification to supervisory authority
Article 34 Breach communication to the data subject	Article 35 Data protection impact assessment	Article 46 Transfers subject to appropriate safeguards

Other applications of secure federated analytics

Hospitals & Pharma

Collective survival analysis in oncology

Lab reference data

Train image classifiers in dermatology

Insurance & Re-Insurance

Train collective risk models

Cross-vertical collaboration (Value-Based Healthcare)

Cyber Security

Cross-organization alert enrichment

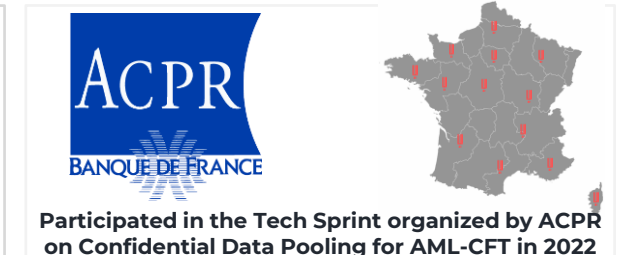
Collective threat intelligence models

Private search of IoCs/alerts

Financial Services

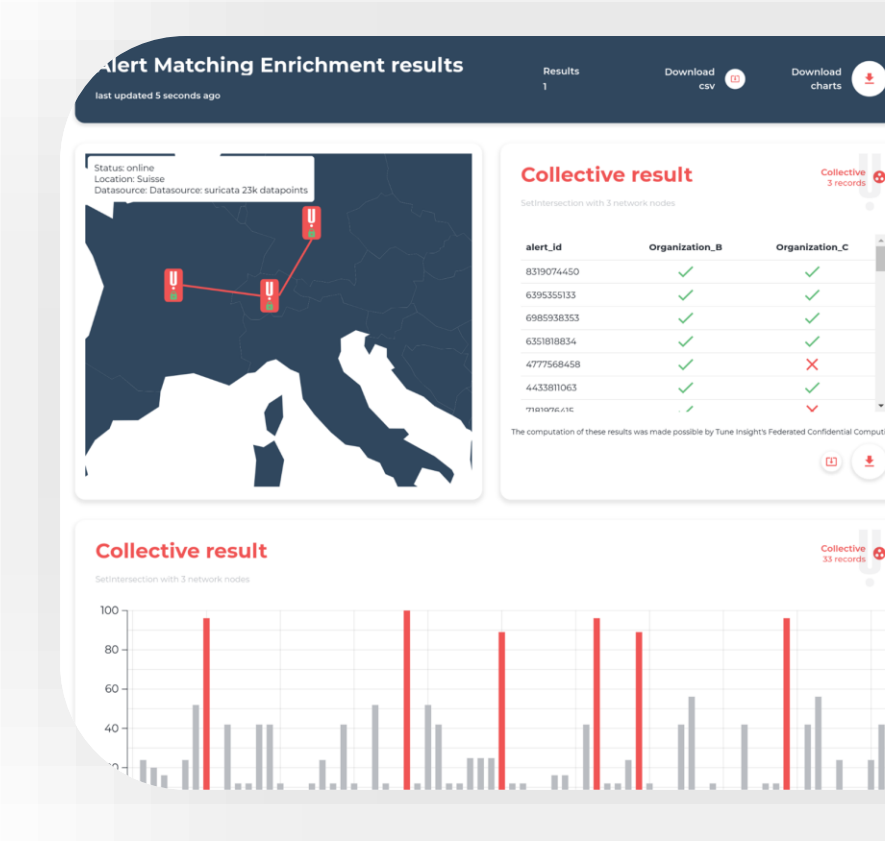
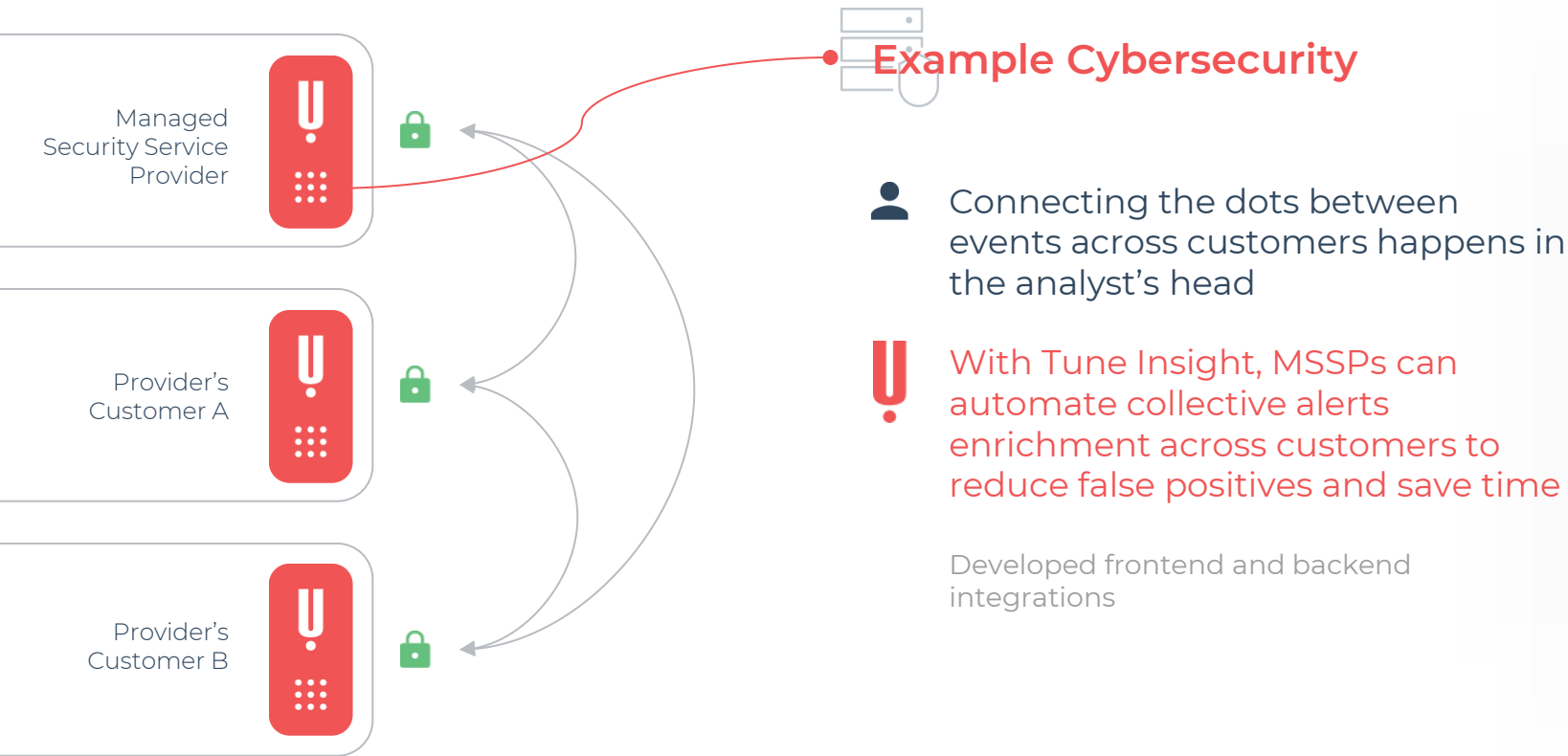
Collaborative analytics

Sensitive data pooling, AML-CFT



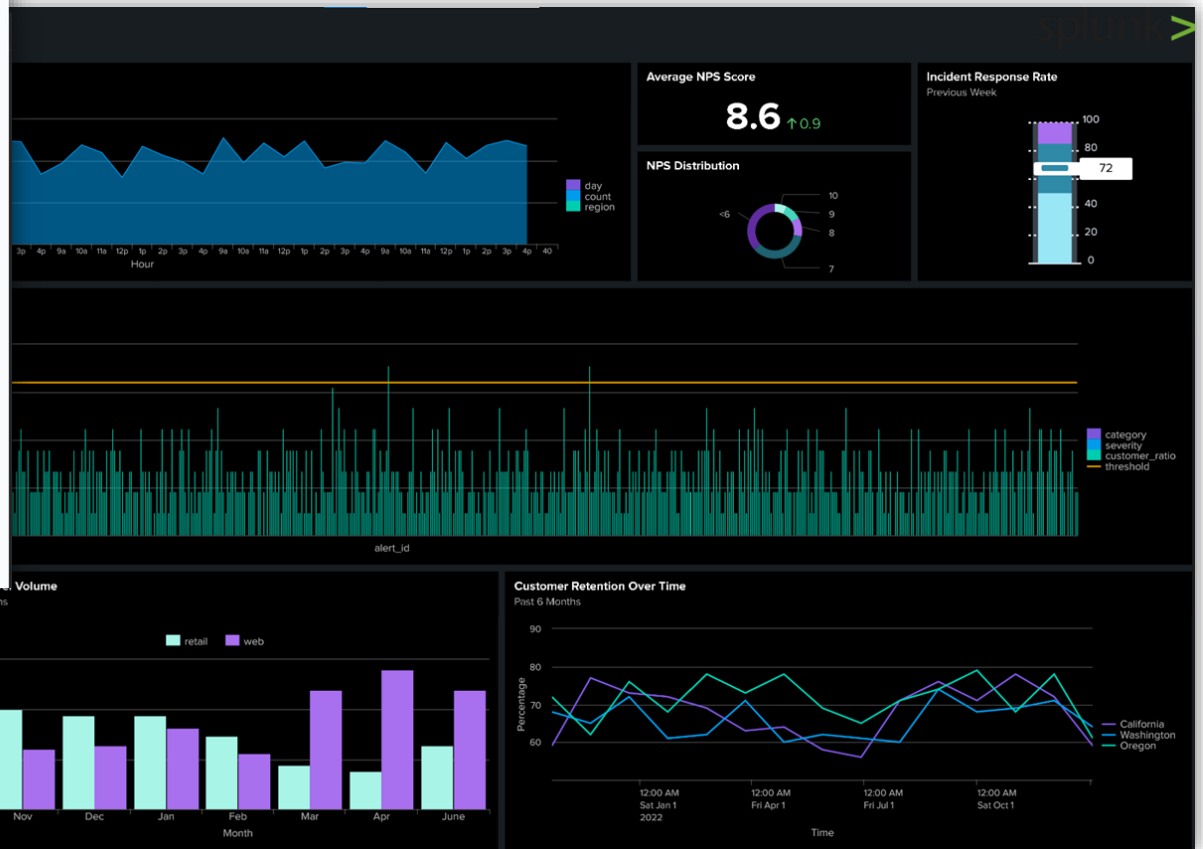
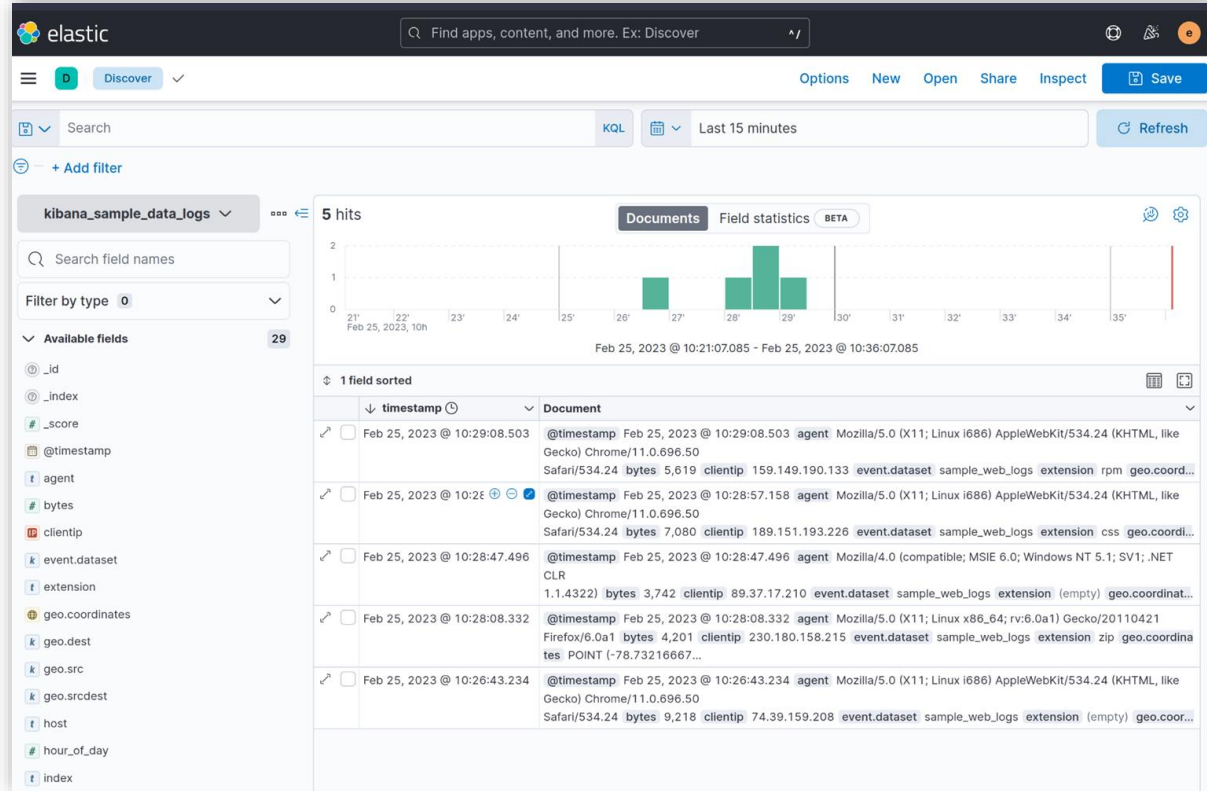
Confidential Collaborative Analytics and Machine Learning

With Tune Insight, organizations can collaborate on their most sensitive cybersecurity data to collectively better defend against cyber attacks

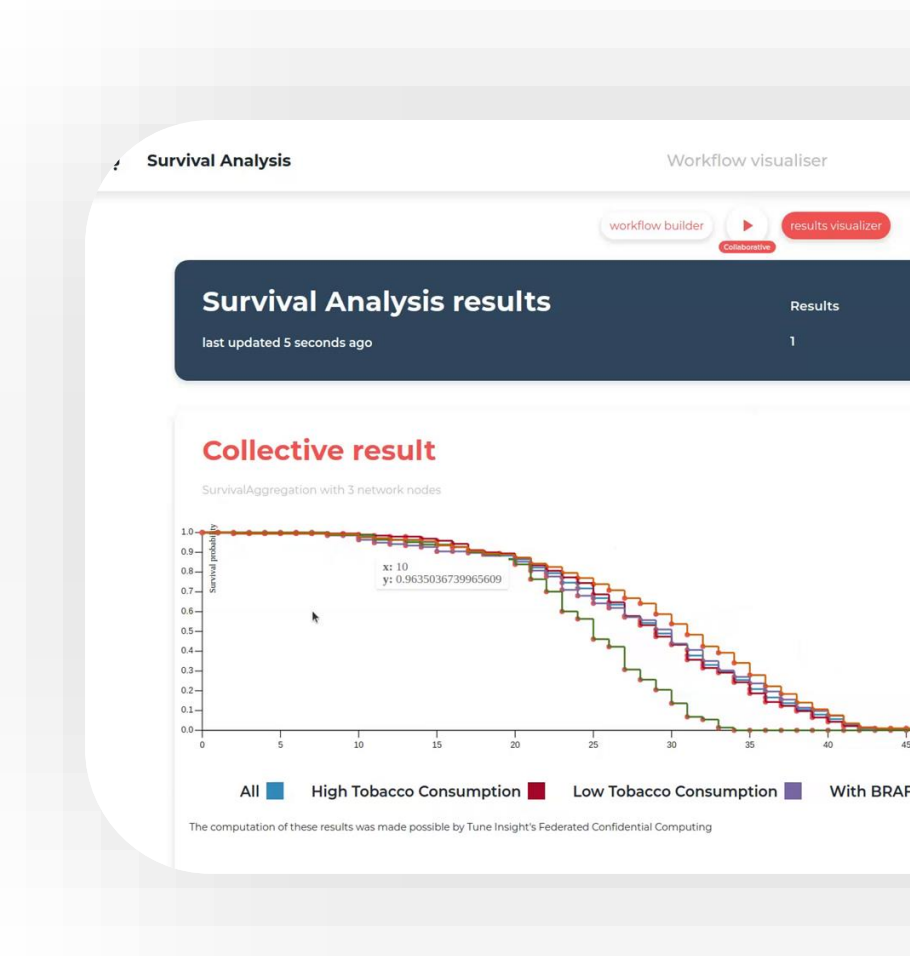
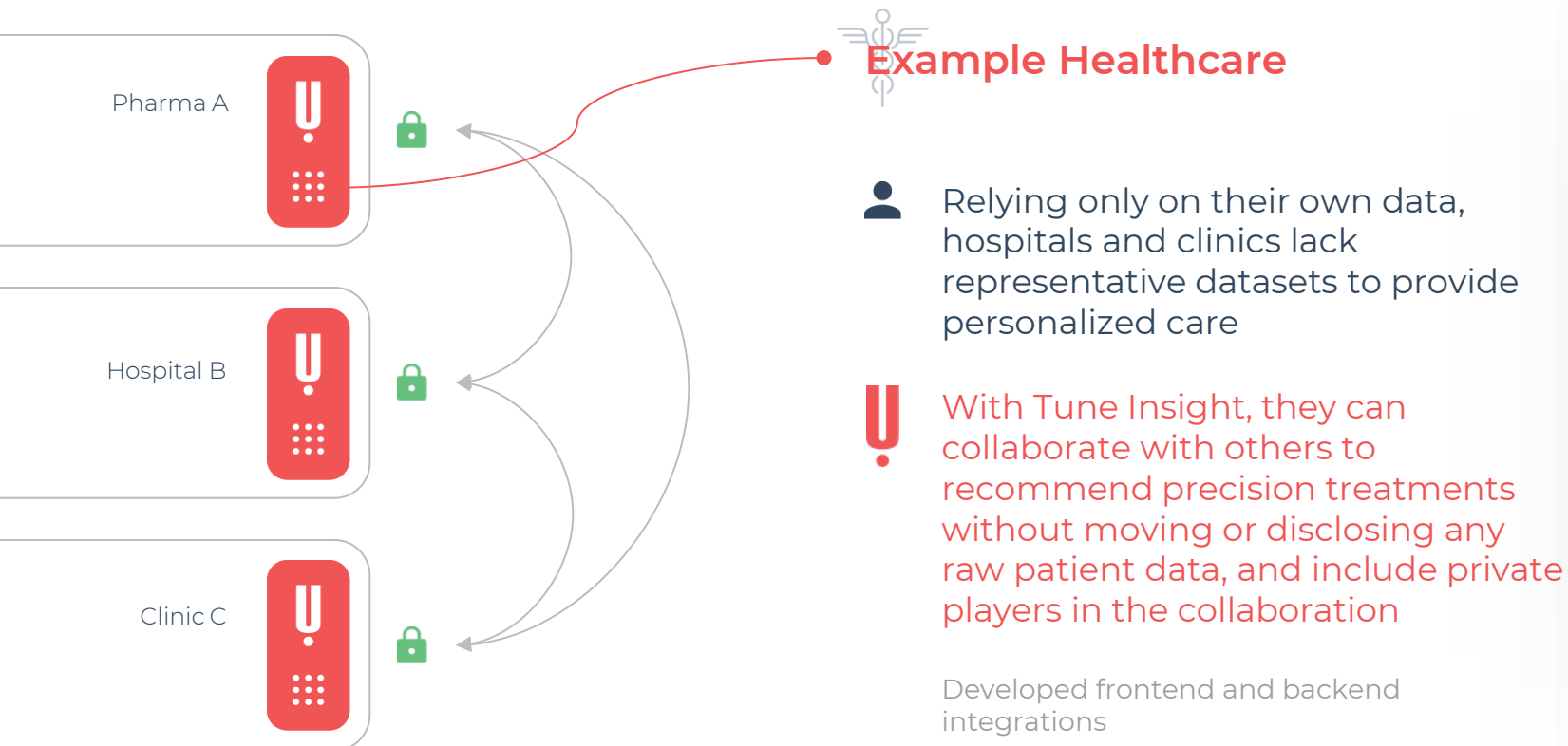


Cyber: Integration with existing platforms and dashboards

Use case: enriching alerts with data from multiple parties,
integrated in the organization's existing tools and workflows



Based on the same core technology, we address similar problems in other verticals

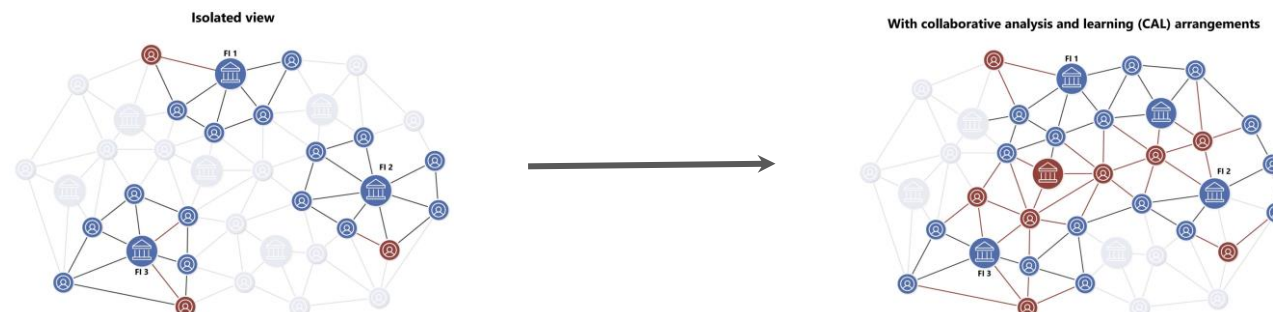


Challenges in Fraud Detection and AML²

- *Isolated view*
- Data interoperability
- Data protection and privacy

WORLD ECONOMIC FORUM “PETs can fundamentally change the nature of data sharing in financial services, creating new value for customers and addressing institutions' most pressing problems in a way that is acceptable to customers, regulators, and society at large.”¹
WEF. September 2019

¹ <https://www.weforum.org/whitepapers/the-next-generation-of-data-sharing-in-financial-services-using-privacy-enhancing-techniques-to-unlock-new-value/>



Traditional Siloed Rule-Based

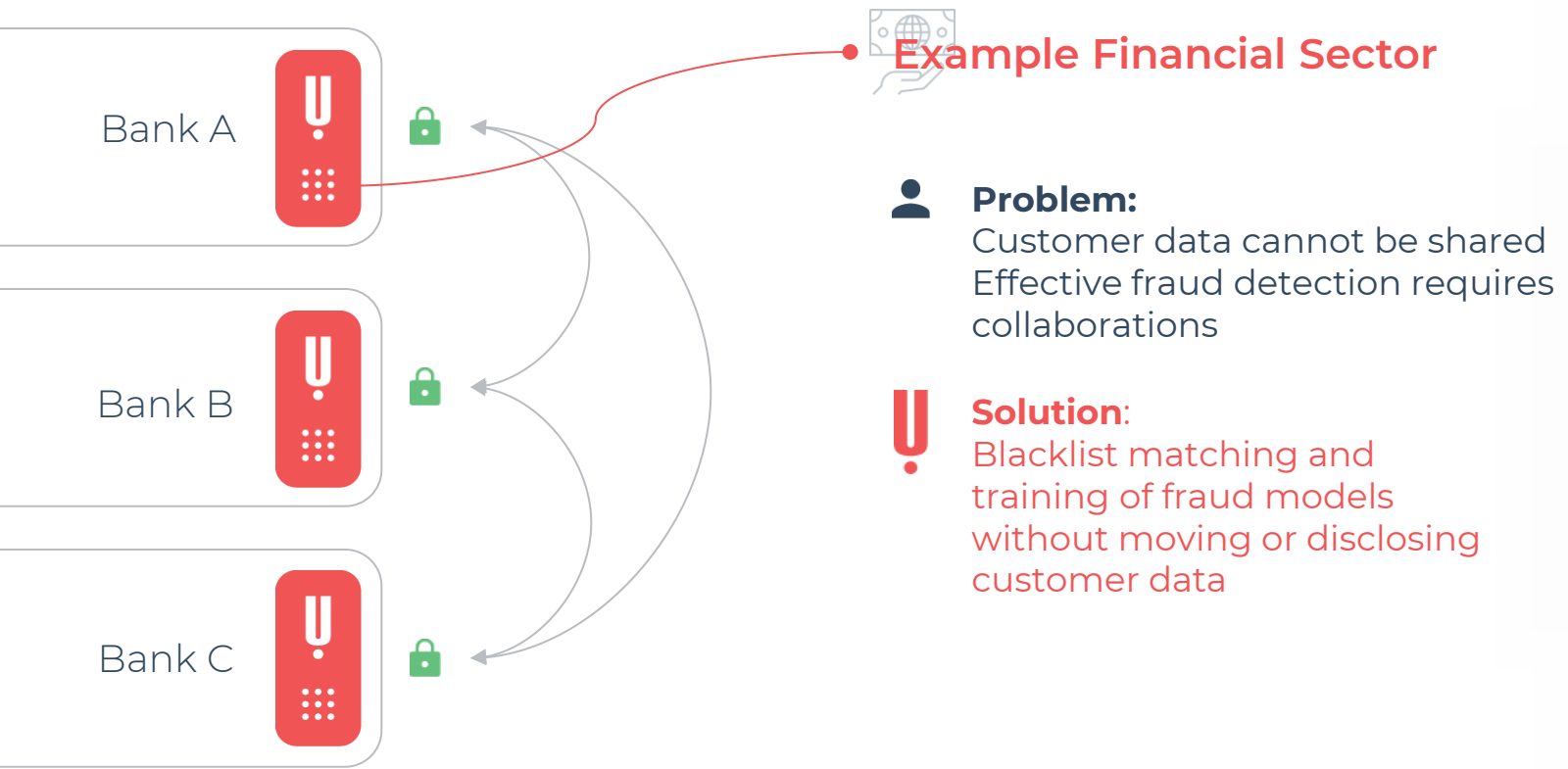
- Traditional rule-based systems can result in high false-positives and false-negatives (90%-95% of the generated alerts are FP)²
- ML not fully effective when data from multiple sources is not available (siloed views)

² <https://www.bis.org/about/bisih/topics/fmis/aurora.htm>

Collaborative Analysis and Learning

- Cross-border ML monitoring can reduce FPs 75% vs rule-based siloed.²
- PET-enabled CAL with machine learning-based network analysis appears to reduce the number of FPs by up to 80% compared with the siloed rule-based method.

Based on the same core technology, we address similar problems in other verticals



Collective Blacklist Search

Check if an IBAN has been blacklisted by a bank in this network

IBAN
BE05416941631736997173182998

Verify

- The IBAN you search for is not disclosed to the other participating
- If there is a match, its source is hidden
- Your blacklist never leaves your security perimeter
- Your blacklist is never revealed to other participating banks

This IBAN was **found** in the collective blacklist.

Collective statistics and time-series information about suspicious account activity

Collective IBAN Usage Monitoring using Private Search

In this notebook, we showcase *Tune Insight's Collective Private Search* through a simple use case: A group of three independent financial institutions hold databases of transactions and analysts want to query the usage of specific IBANs over a period of time across all financial institutions.

In practice, responding to such query would require the institutions to share or centralize their transactional data and view the query from the user, affecting confidentiality and privacy across stakeholders.

Using our *Collective Private Search* distributed computation, the analyst's query can effectively be treated while ensuring:

- The transactions databases are not shared across institutions.
- The analyst's query is not visible to any of the institutions.

These guarantees are made possible through the use of a Private Information Retrieval (PIR) protocol, a cryptographic primitive that allows a user to retrieve information from a server without revealing which information they are retrieving. The PIR protocol is implemented securely using homomorphic encryption.

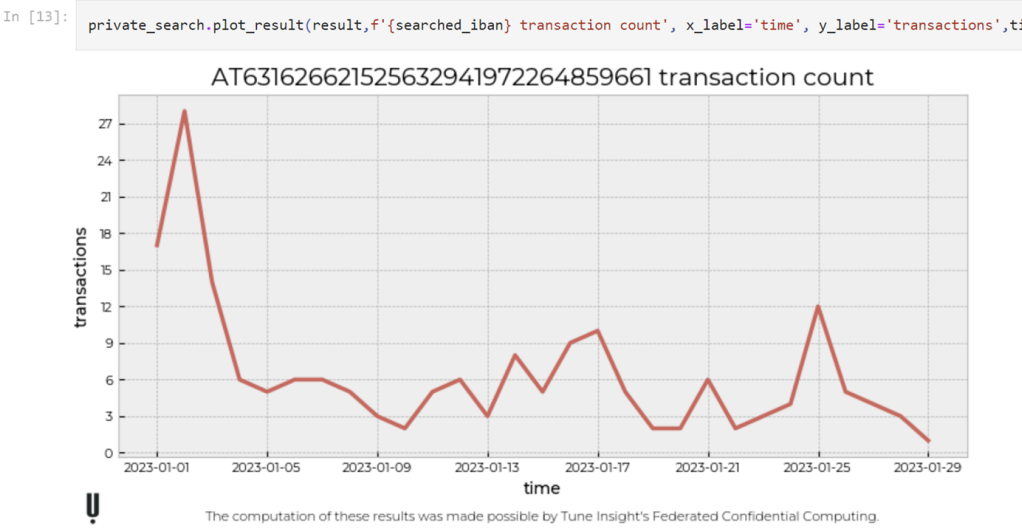
In [12]:

```
# Valid Query
searched_iban = "AT631626621525632941972264859661"
result = private_search.query(searched_iban)
display(result)
```

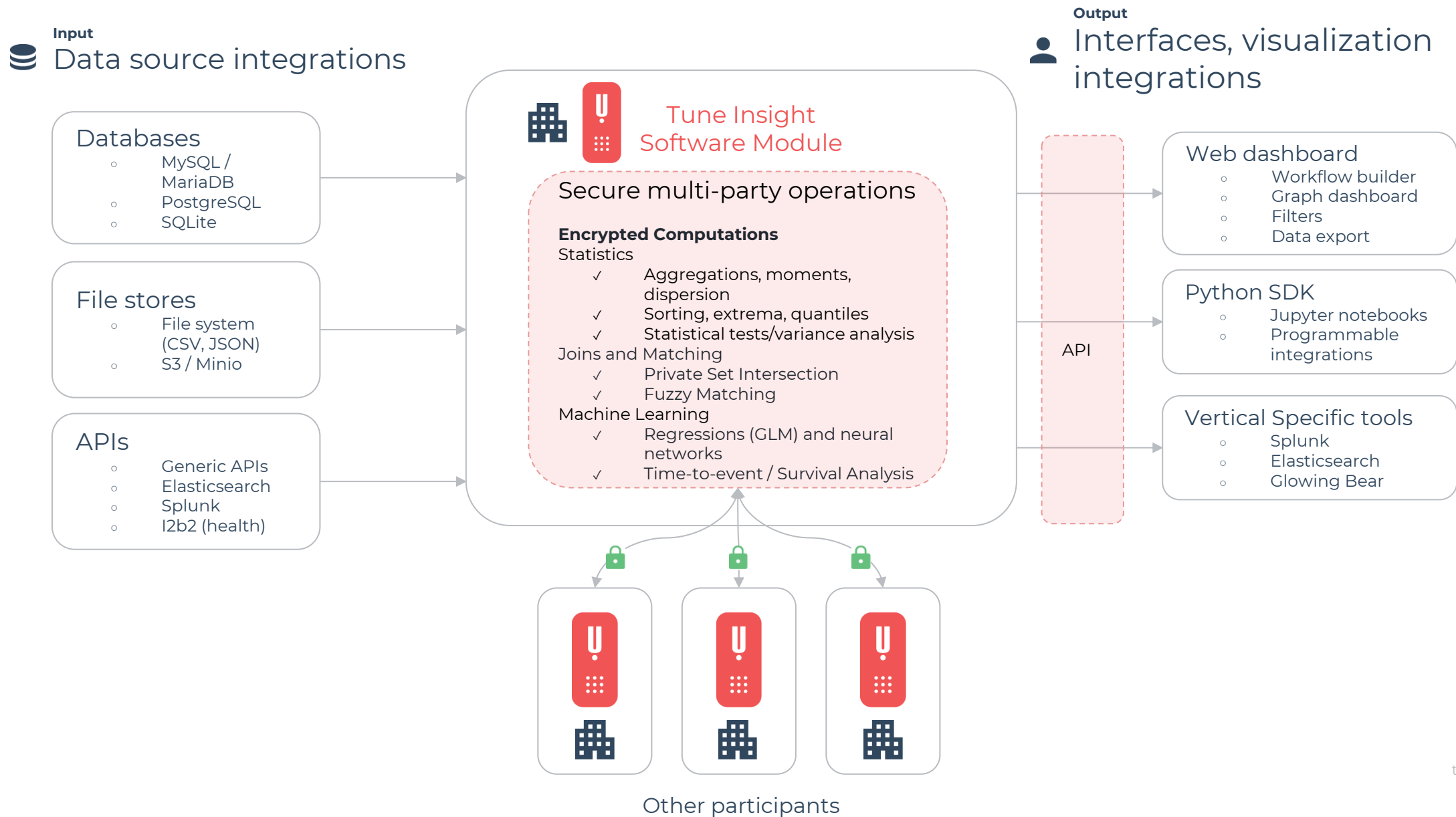
	2023-01-01	2023-01-02	2023-01-03	2023-01-04	2023-01-05	2023-01-06	2023-01-07	2023-01-08	2023-01-09	2023-01-10	...	2023-01-20	2023-01-21	2023-01-22	2023-01-23	2023-01-24	2023-01-25	2023-01-26	2023-01-27	2023-01-28	2023-01-29
0	17	28	14	6	5	6	6	5	3	2	...	2	6	2	3	4	12	0	0	0	0

1 rows × 29 columns

The results can then be plotted using the library or used in further processing.



Product Overview: Tune Insight Software Module



IT Security Assessment

- International security frameworks
 - OWASP Application Security Verification Standard
 - NIST Cybersecurity Framework
- Swiss and domain specific frameworks
 - ICT minimum standard
 - Hospitals and H+ guidelines
- State-of-the-art security technologies
 - OpenID Connect
 - Attribute-based Access Control
 - Key Management Service integration
- Static Application Security Testing
 - Snyk, Trivy, GitHub Dependabot
- Dynamic Application Security Testing
 - Penetration testing by ImmuniWeb



Detected Vulnerabilities Statistics



Diagram 1: Number of vulnerabilities in your web grouped by risk levels

Use case for SPO: Federated analytics platform for research and molecular tumor board

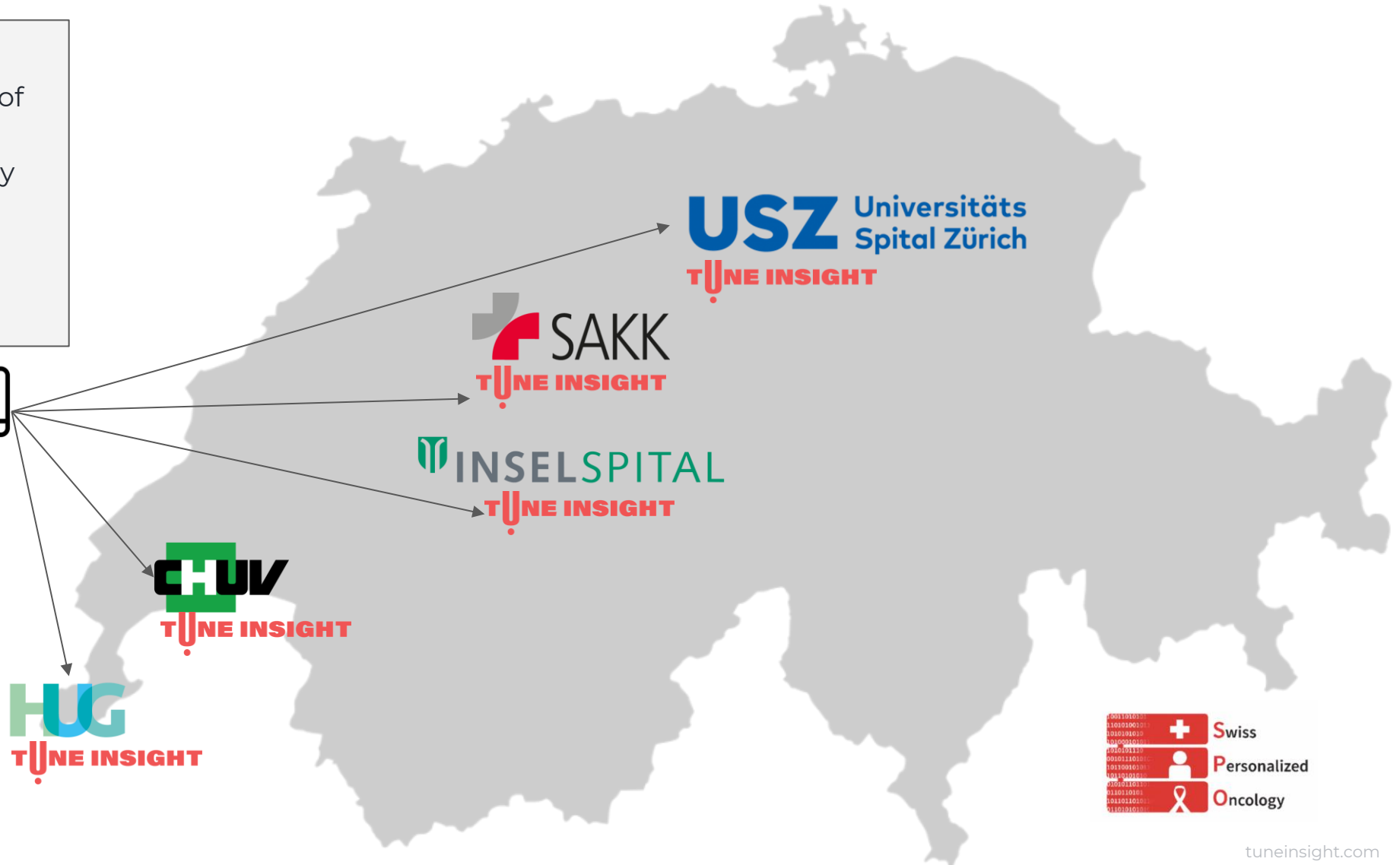
Q1: How many adult cancer patients consenting on reuse of routine data for research with diagnosis of a malignancy on or after 1st January 2015, mutations in BRAF gene and under anti-PD-1 are there?

Explore



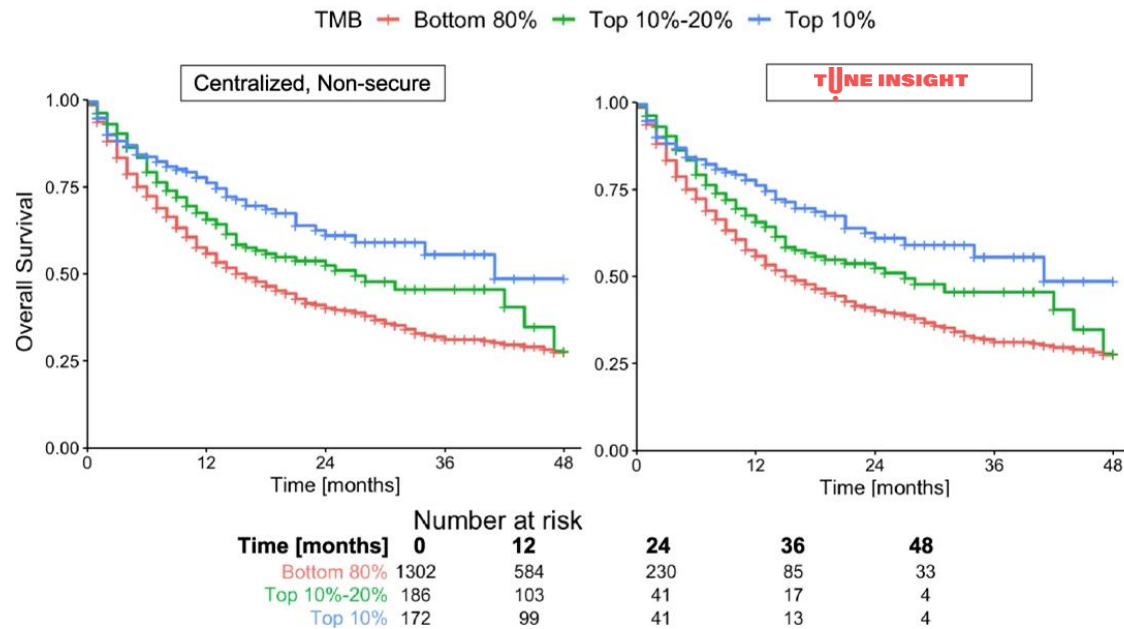
Analysis

Q2: Among these patients, what is the overall survival for patients with and without a mutation on position 600 of the BRAF gene?



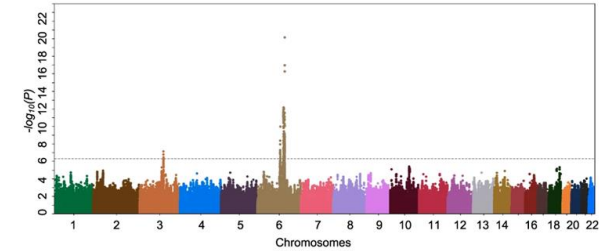
Privacy-Preserving Federated Analytics for Precision Medicine

Survival analysis

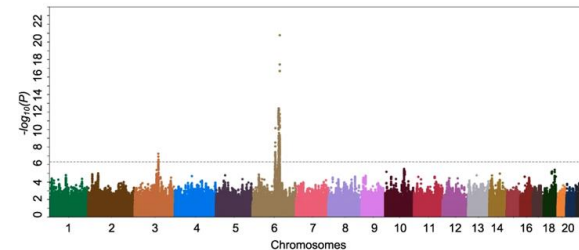


GWAS

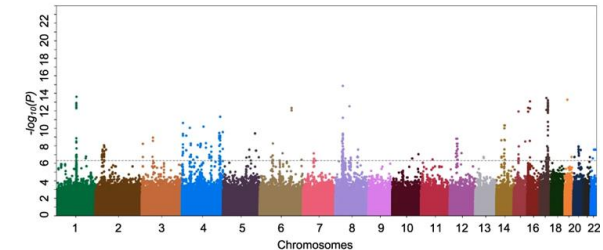
(a) Original Approach



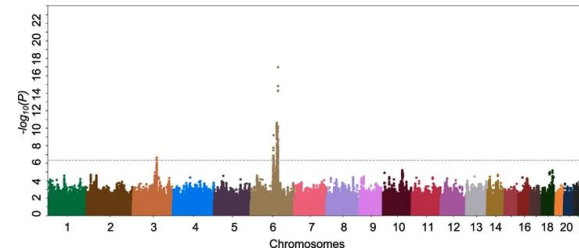
(b) TUNE INSIGHT-GWAS



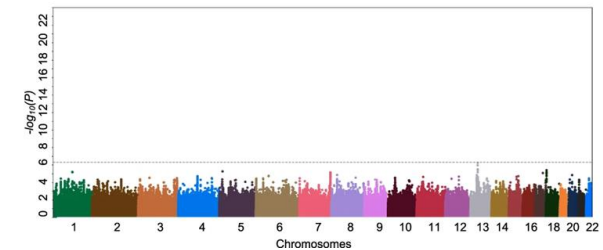
(c) Meta-analysis Approach



(d) TUNE INSIGHT-FastGWAS



(e) Independent Approach

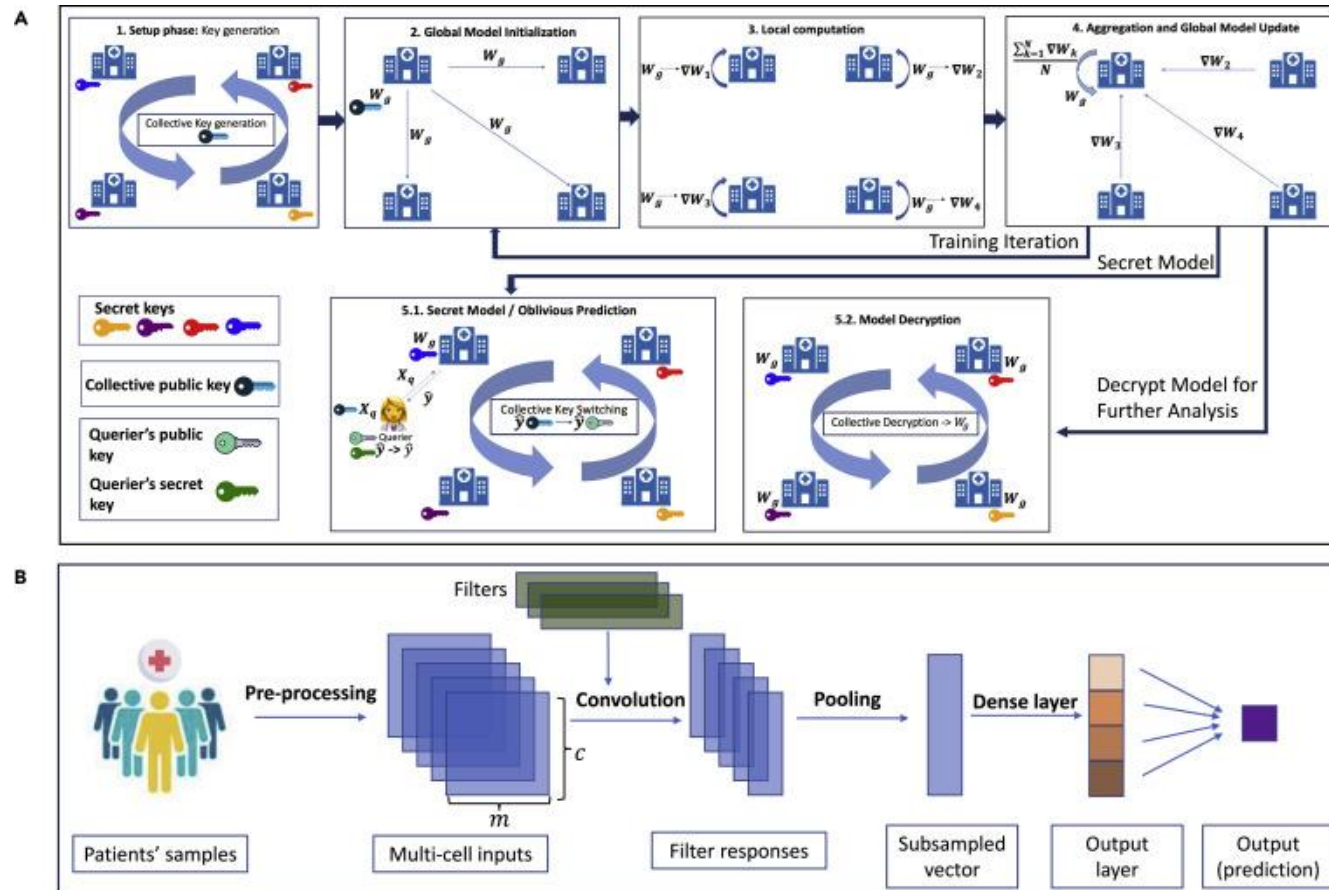


D. Froelicher, J.R. Troncoso-Pastoriza, et al. "Truly privacy-preserving federated analytics for precision medicine with multiparty homomorphic encryption", Nat Commun 12, 5910 (2021). <https://doi.org/10.1038/s41467-021-25972-y>

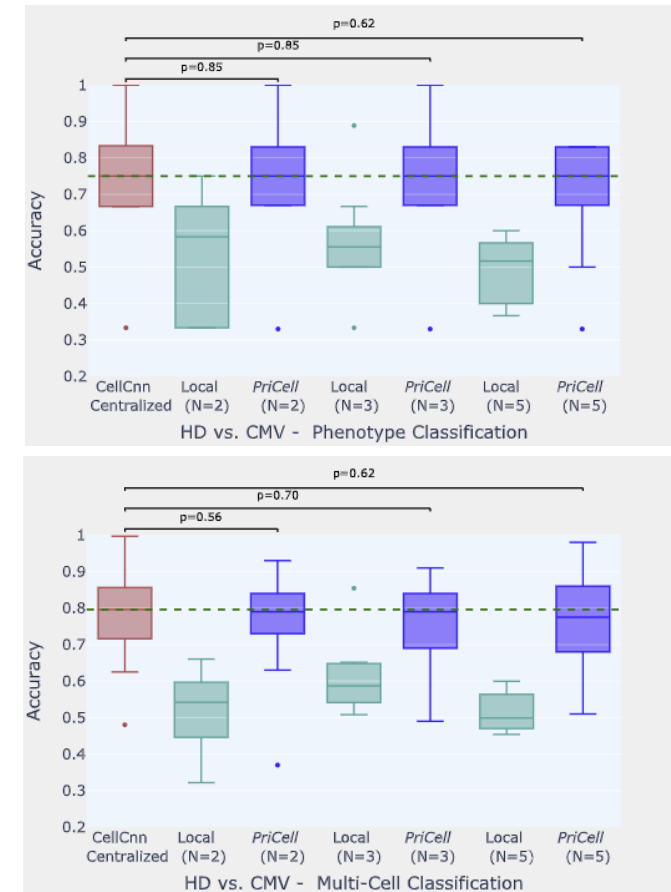
H. Cho, D. Froelicher, J. Chen, M. Edupalli, A. Pyrgelis, J.R. Troncoso-Pastoriza, J.-P. Hubaux, B. Berger. "Secure and Federated Genome-Wide Association Studies for Biobank-Scale Datasets", bioRxiv, 2022 <https://www.biorxiv.org/content/10.1101/2022.11.30.518537v1>

Privacy-Preserving Single-Cell Analysis

System Model



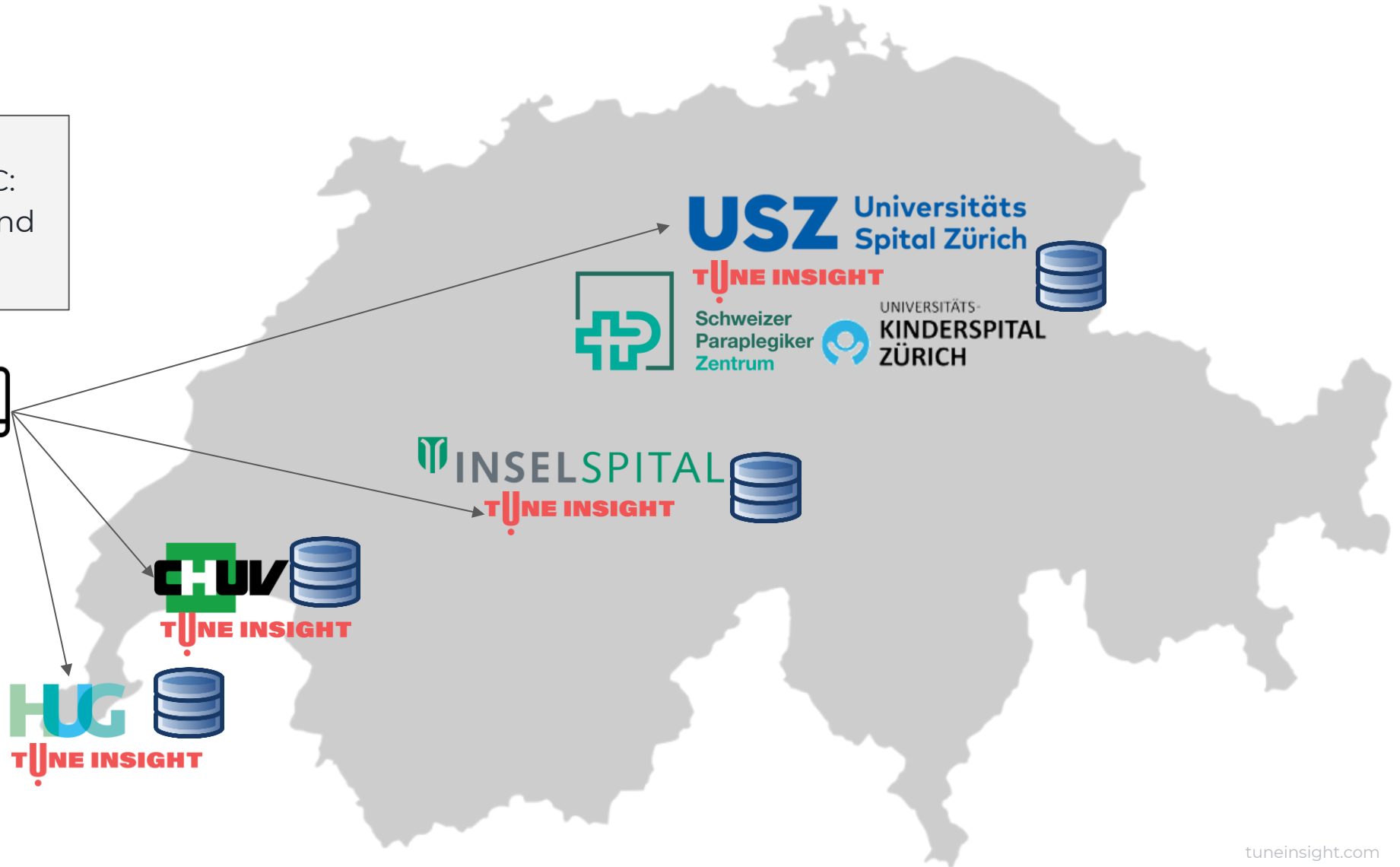
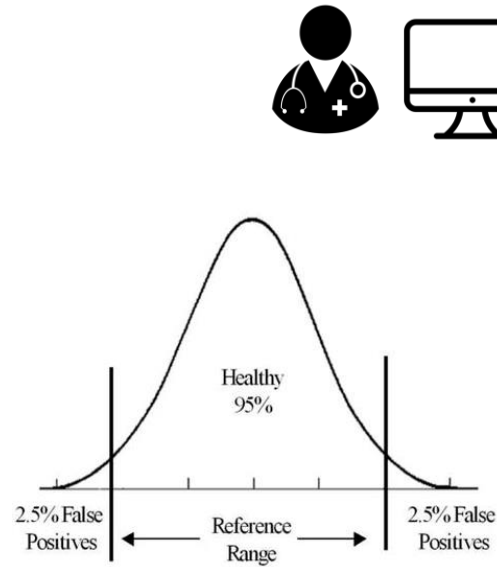
Results



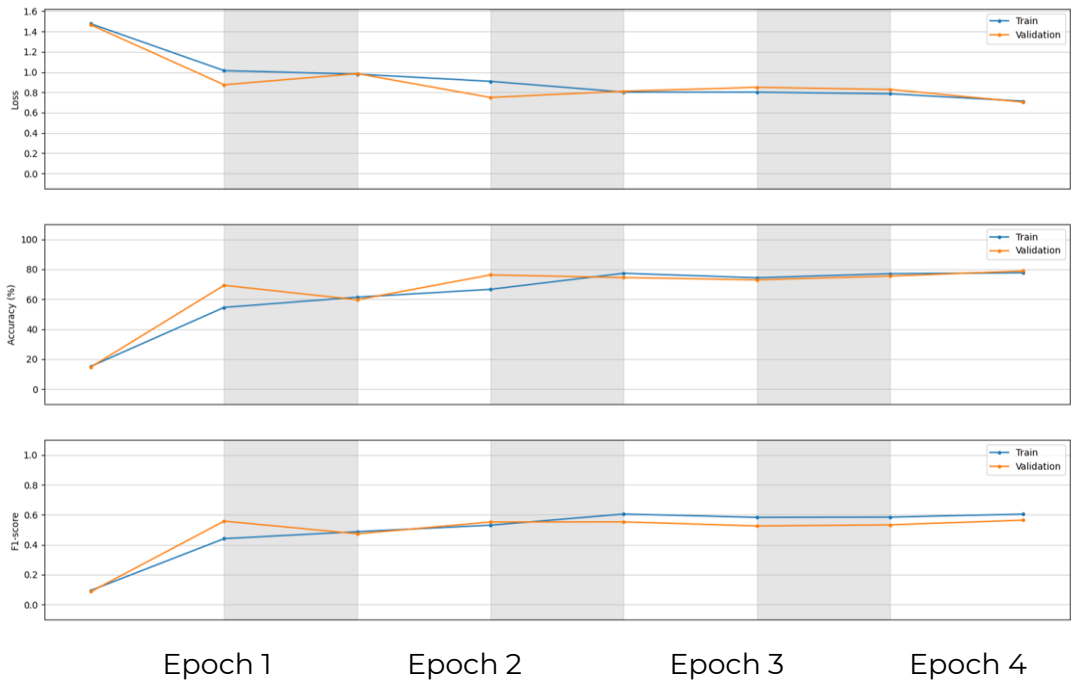
S. Sav, J.P. Bosuat, J.R. Troncoso-Pastoriza, et al. "Privacy-preserving federated neural network learning for disease-associated cell classification." Cell Patterns 2022. <https://doi.org/10.1016/j.patter.2022.100487>

Use case for Swiss BioRef: real-time personalized lab reference ranges

Q1: What is the reference range for Creatinine (LOINC: 14682-9) for 50y-old male and heart failure?



Secure Federated Training of Deep Neural Networks on Dermatology Images with combination of HE, MPC, FL, and DP



Dataset: Fitzpatrick17k, ~30k images
(<https://github.com/mattgroh/fitzpatrick17k>)

Model:
Type: ViT with 4-layers embedding
Size: 5,528,259 parameters, 44.3MB

4 epochs	Local training baseline	Secure federated training
Nodes	1 node with 10909 samples	3 nodes (~3635 samples each)
Training accuracy	72.16%	77.65%
Training F1-score	0.279431	0.604438
Validation accuracy	72.13%	78.88%
Validation F1-score	0.279364	0.564171
Privacy params	N/A	$\epsilon = 1.0, \delta = 0.0001$
Time overhead	0	~10% (w.r.t. vanilla FL)

100 seconds/epoch on a g4dn.2xlarge AWS EC2 instance with a Nvidia T4 GPU (16GB memory)

Selected recent publications from Tune Insight

D. Froelicher et al. “**Scalable and Privacy-Preserving Federated Principal Component Analysis.**” IEEE Symposium on Security and Privacy (2023).
<https://doi.ieeecomputersociety.org/10.1109/SP46215.2023.00051>

H. Cho et al. “**Secure and Federated Genome-Wide Association Studies for Biobank-Scale Datasets.**” Biorxiv, 2022.
<https://www.biorxiv.org/content/10.1101/2022.11.30.518537v1>

S. Sav, J.P. Bossuat, J.R. Troncoso-Pastoriza, et al. “**Privacy-preserving federated neural network learning for disease-associated cell classification.**” Cell Patterns 2022.
<https://doi.org/10.1016/j.patter.2022.100487>

J.P. Bossuat, C. Mouchet, J.R. Troncoso-Pastoriza, J.P. Hubaux. “**Efficient Bootstrapping for Approximate Homomorphic Encryption with Non-Sparse Keys.**” EUROCRYPT 2022

D. Froelicher, J.R. Troncoso-Pastoriza, et al. “**Truly privacy-preserving federated analytics for precision medicine with multiparty homomorphic encryption.**” Nature Communications 12, 5910 (2021).
<https://doi.org/10.1038/s41467-021-25972-y>

J. Scheibner, J.L. Raisaro, J.R. Troncoso-Pastoriza, M. Ienca, J. Fellay, E. Vayena, J.-P. Hubaux. “**Revolutionizing Medical Data Sharing Using Advanced Privacy Enhancing Technologies: Technical, Legal and Ethical Synthesis.**” Journal of Medical Internet Research, 2021.
<https://pubmed.ncbi.nlm.nih.gov/33629963/>

C. Mouchet, J. R. Troncoso-Pastoriza, J.P. Bossuat and J.P. Hubaux. “**Multiparty Homomorphic Encryption from Ring-Learning-with-Errors.**” PETS 2021. IACR ePrint Archive: Report 2020/304, 2020.

D. Froelicher, J. R. Troncoso-Pastoriza, A. Pyrgelis, S. Sav, J. S. Sousa, J.P. Bossuat, and J.P. Hubaux. “**Scalable Privacy-Preserving Distributed Learning.**” PETS 2021. CoRR abs/2005.09532, 2020.

S. Sav, A. Pyrgelis, J.R. Troncoso-Pastoriza, J.-P. Bossuat, J.S. Sousa, J.-P. Hubaux, “**POSEIDON: Privacy-Preserving Federated Neural Network Learning.**” CoRR abs 2009.00349, 2020. NDSS 2021

Lattigo library. <https://github.com/tuneinsight/lattigo>



<https://tuneinsight.com>



<https://linkedin.com/company/tuneinsight>



<https://twitter.com/tuneinsight>

- [Ajt96] M. Ajtai. Generating hard instances of lattice problems. Quaderni di Matematica, 13:1-32, 2004. Preliminary version in STOC 1996
- [BDOP16] C. Baum, I. Damgård, S. Oechsner and C. Peikert. Efficient Commitments and Zero-Knowledge Protocols from Ring-SIS with Applications to Lattice-based Threshold Cryptosystems. Cryptology ePrint Archive, report 2016/997. 2016
- [BGG+14] D. Boneh, C. Gentry, S. Gorbunov, S. Halevi, V. Nikolaenko, G. Segev, V. Vaikuntanathan, D. Vinayagamurthy. Fully key-homomorphic encryption, arithmetic circuit ABE and compact garbled circuits. In EUROCRYPT, pp 533-556. 2014
- [BGV12] Z. Brakerski, C. Gentry, V. Vaikuntanathan. (Leveled) fully homomorphic encryption without bootstrapping. TOCT 6(3):13, 2014. Preliminary version in ITCS 2012
- [BLLN13] J. W. Bos, K. Lauter, J. Loftus, M. Naehrig. Improved security for a ring-based fully homomorphic encryption scheme. In IMACC 2013, Oxford, UK, Dec 17-19. LNCS 8308, pp. 45–64. 2013.

- [CS15] A. Costache and N.P. Smart. Which Ring Based Somewhat Homomorphic Encryption Scheme is Best? Cryptology ePrint Archive, report 2015/889. 2015
- [FV12] J. Fan and F. Vercauteren. Somewhat Practical Fully Homomorphic Encryption. Cryptology ePrint Archive, report 2012/144. 2012
- [Gen09] C. Gentry. Fully Homomorphic Encryption using ideal lattices. In STOC, pp. 169-178. 2009
- [GGH+13] S. Garg, C. Gentry, S. Halevi, M. Raykova, A. Sahai, B. Waters. Candidate indistinguishability obfuscation and functional encryption for all circuits. In FOCS, pp. 40-49, 2013
- [HPS96] J. Hoffstein, J. Pipher, J.H. Silverman. NTRU: A Public Key Cryptosystem. NTRU Cryptosystems, Inc. (now Security Innovation)

- [LLL82] A.K. Lenstra, H.W. Lenstra, Jr., L. Lovász. Factoring polynomials with rational coefficients. *Mathematische Annalen*, 261(4):515-534, Dec. 1982
- [LNV10] K. Lauter, M. Naehrig, V. Vaikuntanathan. Can homomorphic encryption be practical?. In *ACM CCSW '11*. ACM, New York, NY, USA, pp. 113-124. 2011
- [LP11] R. Lindner and C. Peikert. Better key sizes (and attacks) for LWE-based encryption. In *CT-RSA*, pp. 319-339. 2011
- [LPR10] V. Lyubashevsky, C. Peikert, O. Regev. On Ideal lattices and learning with errors over rings. *Journal of the ACM*, 60(6):43:1-43:35, Nov 2013. Preliminary version in *Eurocrypt 2010*
- [Pei16] C. Peikert. A Decade of Lattice Cryptography. *Cryptology ePrint Archive*, report 2015/939. 2015. Last revised 17 Feb 2016

- [PTP17] A. Pedrouzo-Ulloa, J. R. Troncoso-Pastoriza, F. Pérez-González. Number Theoretic Transforms for Secure Signal Processing, in IEEE Transactions on Information Forensics and Security, vol. 12, no. 5, pp. 1125-1140, May 2017.
- [PR07] C. Peikert and A. Rosen. Lattices that admit logarithmic worst-case to average-case connection factors. In STOC, pp. 478-487. 2007
- [RS10] M. Rückert and M. Schneider. Estimating the security of lattice-based cryptosystems. Cryptology ePrint Archive, report 2010/137. 2010
- [Sch87] C.-P. Schnorr. A hierarchy of polynomial time lattice basis reduction algorithms. Theor. Comput. Sci., 53:201-224, 1987
- [Sho97] P.W. Shor. Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer
- [SS13] D. Stehlé and R. Steinfeld. Making NTRUEncrypt and NTRUSign as Secure as Standard Worst-Case Problems over Ideal Lattices. Cryptology ePrint Archive, report 2013/004. 2013

Key take-aways

1. Limited data availability for sensitive personal (Life & Health) data in practice (e.g., nFADP, 1 September 2023)
2. For sufficiently large and dense datasets, ML/DL methods outperform traditional models, creating value for policyholders and insurance companies
3. Privacy preserving methods can help to access data