# Privacy–Preserving ML with Fully Homomorphic Encryption

**ZAMA**

# Data is encrypted only during the transport

**Client side**

**Server side**

encrypted *only* during transport

$x$ → $x$

evaluate circuit for $f(\cdot)$

encrypted *only* during transport

$f(x)$ ← $f(x)$

# Fully Homomorphic Encryption

# FHE: encrypting the data also during processing

**Client side**          **Server side**

$x$ → encrypt → $Enc(x)$

secret (enc/dec)ryption key

$sk$

evaluate circuit for $f(\cdot)$ ← $pk_{eval}$    public evaluation key

$f(x)$ ← decrypt ← $Enc(f(x))$

# TFHE: the cryptographic scheme we use

Addition

Multiplication by a scalar

Noise reduction

Table Lookup (TLU)

# Image filtering as a demo



Client Side — Encryption — Server Side — FHE filter execution — Decryption — Client Side
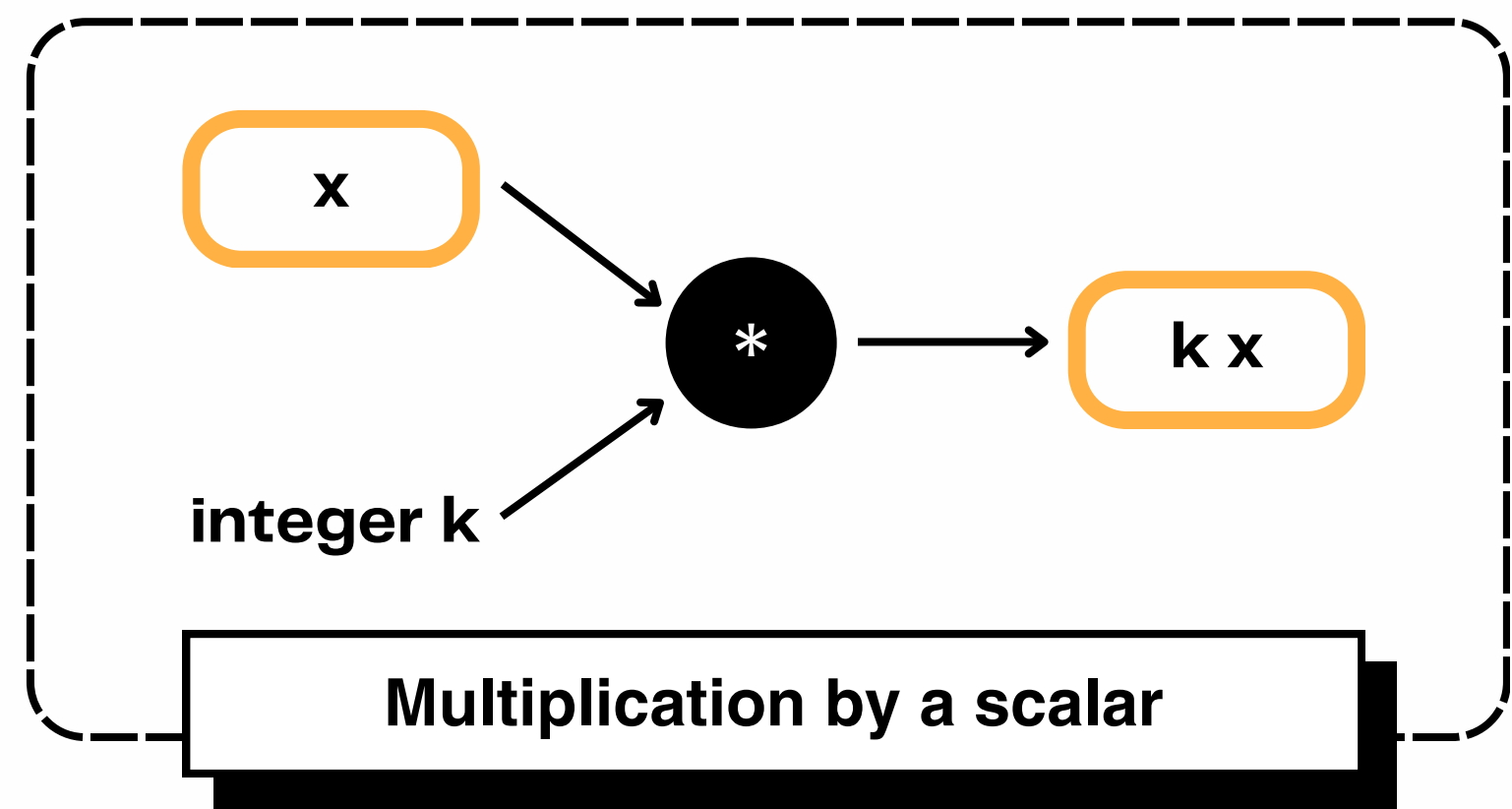
https://huggingface.co/spaces/zama-fhe/encrypted_image_filtering

Privacy-Preserving ML with FHE

# How Does It Work?

# Computation paradigm



circuit of univariate functions

= graph mixing univariate functions and linear combinations

# Linear Models

# Linear models

1. **Train**
   - in clear: linear solver, l-bfgs, quadratic solver
   - on encrypted data: SGD

2. **Apply post-training quantization to both inputs and weights**

   Asymmetric quantization

   $$q_a(x) = \text{round}\left(\frac{x}{\Delta}\right) + z \qquad \Delta = \frac{\max(x) - \min(x)}{2^p - 1}$$
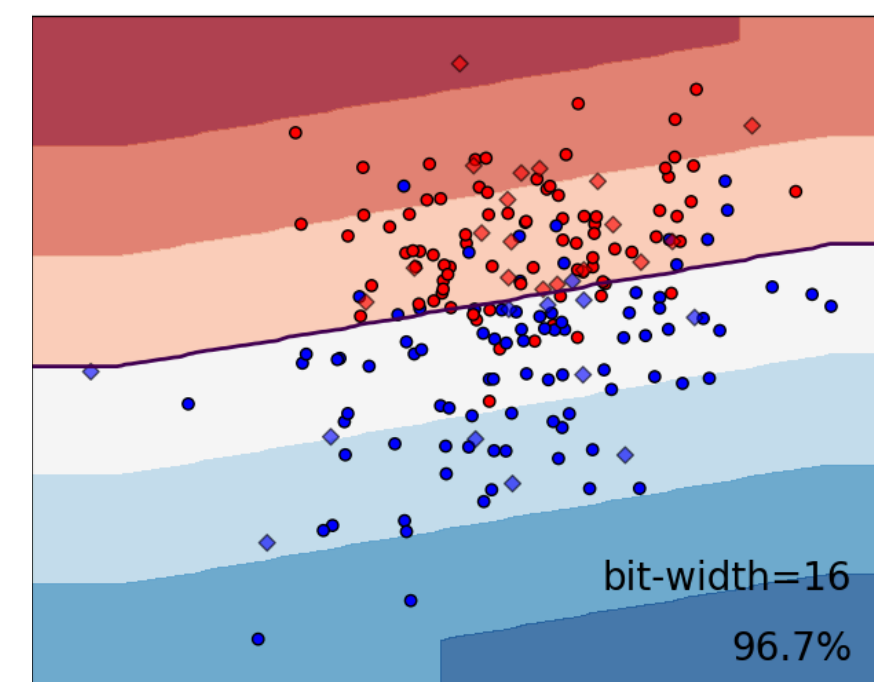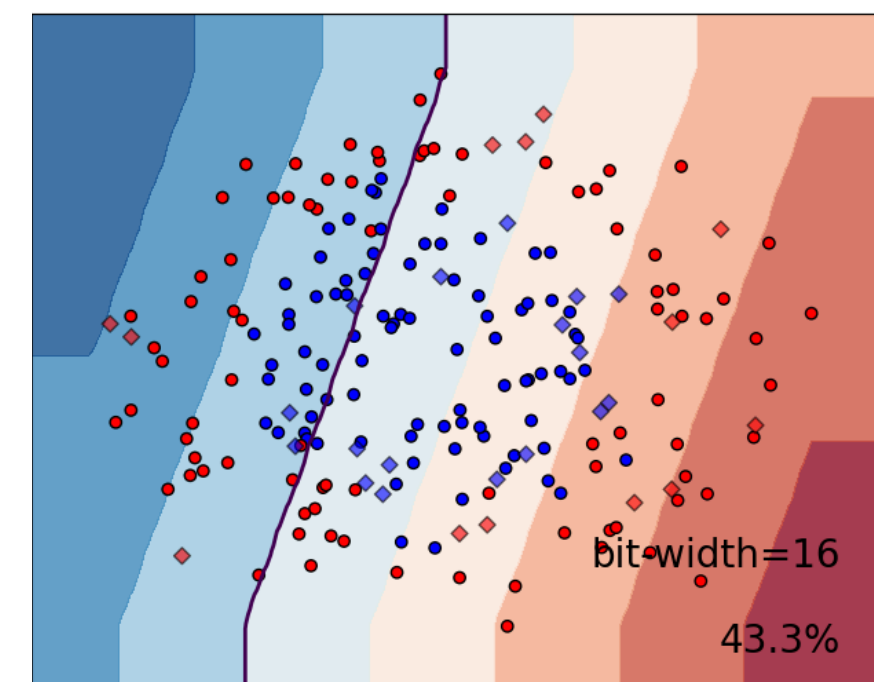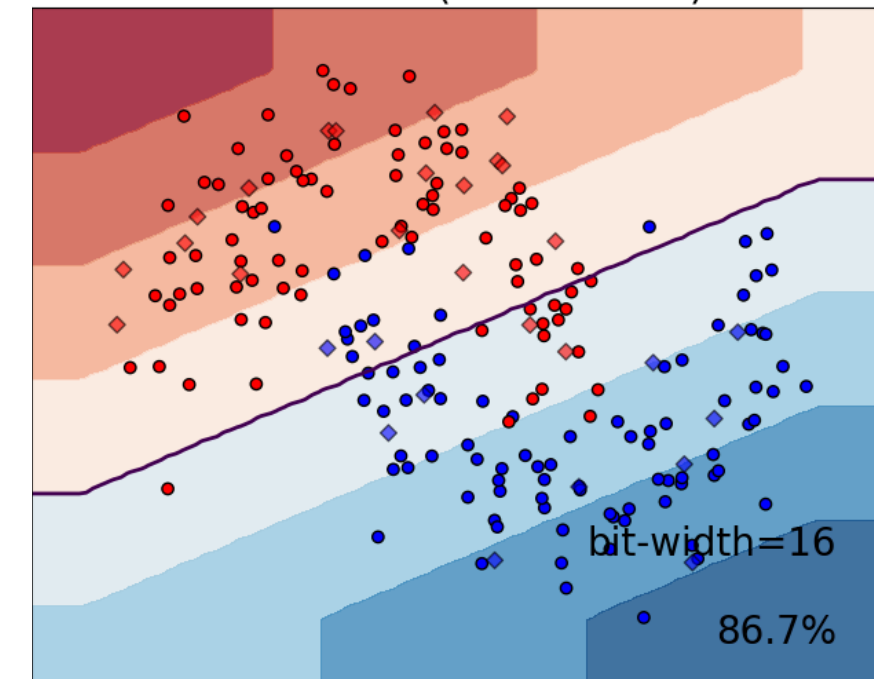   $$z = -\left\lfloor \frac{\min(x)}{\Delta} \right\rfloor$$

3. **Execute fully-levelled in FHE**
   - execution time: 10-100ms
   - can use up to 8 bits weights and inputs

   ✅ 8 bits quantization - Should be good for all datasets



Linear Classifiers
Linear SVC (Concrete ML)

bit-width=16

86.7%

bit-width=16

43.3%

bit-width=16

96.7%

# Deep Neural Networks

# FHE computation paradigm vs neural networks operation graphs



circuit of univariate functions

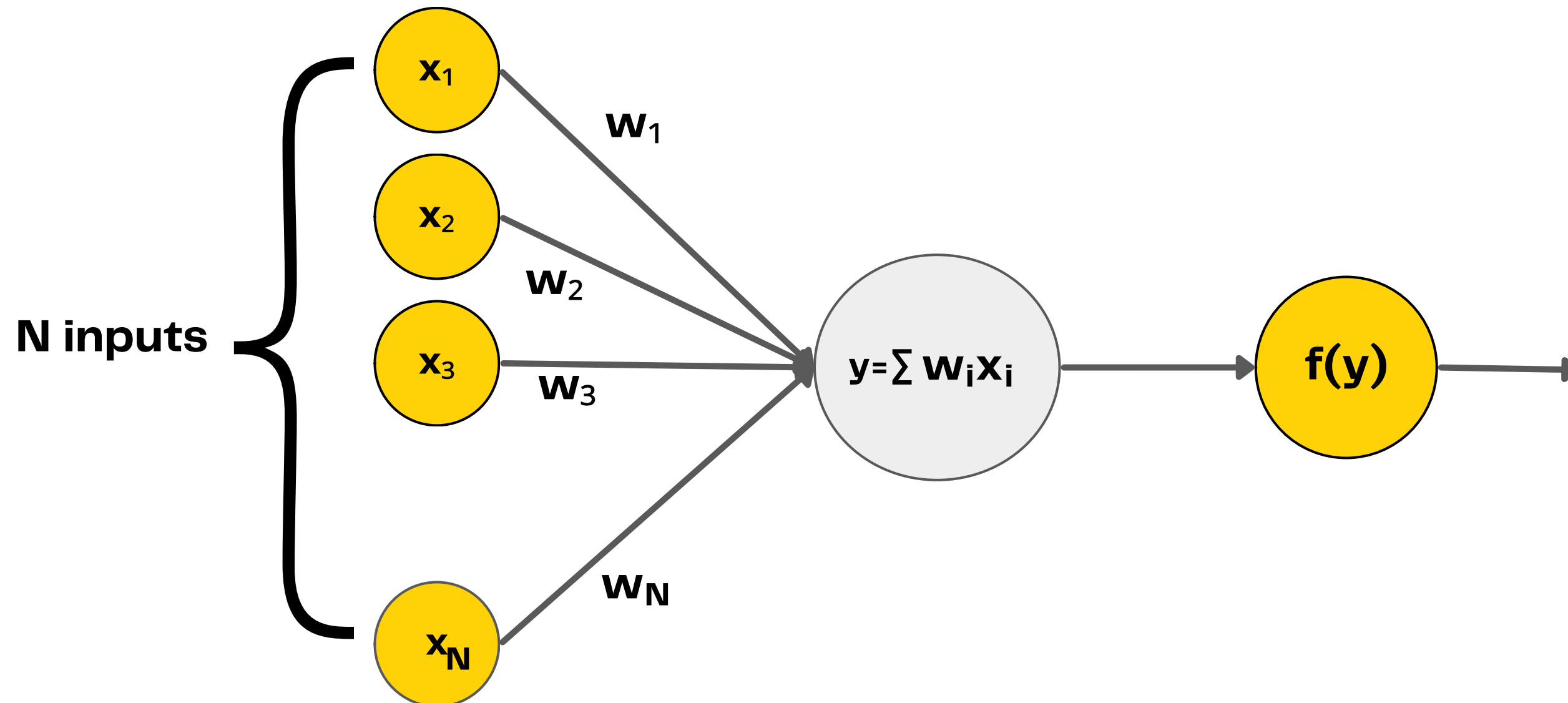= graph mixing univariate functions and linear combinations

**f(...)**
**TLU/PBS**

**Σ**
**LINEAR OPERATION ON INTEGERS**

# Quantization for integer computation

- Inputs, weights and activations are float32
- Need to convert to integer computation
- All values must have at most P bits

N inputs

$x_1$

$w_1$

$x_2$

$w_2$

$x_3$

$w_3$

$x_N$

$w_N$

$y = \sum w_i x_i$

$f(y)$

# Quantization Aware Training (QAT)

**Training:**

find the weights that minimize the objective function when applying the model on the inputs

**Quantization aware training:**

+ the weights are constrained to be representable by integers of a certain bit–width
+ the inputs are represented by integers

**What about gradients under QAT?**

**QAT is (much) more difficult with very low bitwidth**

**Quantizer function**

**Straight through estimator**



Training set loss during training

# Built–in neural networks for toy datasets

- 3–layer MLP
- ReLU Activation

- 2 bits weights
- 4 bits inputs & activations

- Maximum accumulator: 6 bits

# Custom QAT models: results on a tiny dataset

**Fully Connected Model:
3 layers with 192 neurons**



optimized: 4 seconds / image

**CNN:
3 layers with 8–32 3x3 filters, pruned
to 12 max active neurons**



unoptimized: 3 seconds / image

# Optimizations

**Using round PBS / truncate PBS / approximate PBS:**
- replace T[i], where i is on n bits by T'[i'], where:
  - i' is on (n-r) bits
  - i' corresponds to most significant bits of i
- depending on the case:
  - i' = round(i)  or truncate(i)

**Playing with p-error parameters:**
- probability of by-one error can be tuned
- eg, with an error probability of  0.01, it is ~10x faster than a probability of 1E-6

# Custom QAT models: VGG9 on CIFAR

**CIFAR10 – 32x32 – VGG–9 with AveragePooling**

**2b weights, 2b activations, 8b inputs**

| Runtime | Rounding | Accuracy |
|---|---|---|
| VGG Torch | None | 88.7 |
| VGG FHE (simulation) | None | 88.7 |
| VGG FHE (simulation) | 8 bits | 88.0 |
| VGG FHE (simulation) | 7 bits | 87.2 |
| VGG FHE (simulation) | 6 bits | 86.0 |
| VGG FHE | 6 bits | 86.0 |

# Custom QAT models: VGG9 on CIFAR



Cifar Run Duration over Time

Legend: CPU, GPU (4), GPU (8)

Data points (CPU):
- 39.0h
- 10.0h
- 6.9h
- 8.9h
- 9.13h
- 2.7h
- 0.51h
- 288s
- 117s
- 83s
- 40s

GPU (4): 33s
GPU (8): 28s

X-axis: Date (2022/12 – 2024/04)
Y-axis: Run Duration (Seconds) - Log Scale

# Tree-Based Models

# Tree-based models and FHE?

- Powerful and common ML models : DecisionTree, RandomForest, XGBoost, …
- Not directly FHE compliant :
  - depend on control-flow operations (if statements)
  - work with floating points by default

# Hummingbird method

$n_1$   $n_2$   $n_3$   $n_4$

A

| 0 | 0 | 1 | 0 | ← $x_1$ |
| 0 | 1 | 0 | 0 | ← $x_2$ |
| 1 | 0 | 0 | 0 | ← $x_3$ |
| 0 | 0 | 0 | 1 | ← $x_4$ |
| 0 | 0 | 0 | 0 | ← $x_5$ |

X

| 1 | 2 | 7 | 5 | 4 |
| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |

×

P

| 7 | 2 | 1 | 5 |

>

B

| 3 | 5 | 1 | 3 |

Conditions

Q

| 1 | 0 | 0 | 1 |

$e_1$   $e_2$   $e_3$   $e_4$   $e_5$

C

| 1 | 1 | 1 | -1 | -1 | ← $n_1$ |
| 1 | -1 | -1 | 0 | 0 | ← $n_2$ |
| 0 | 0 | 0 | 1 | -1 | ← $n_3$ |
| 0 | 1 | -1 | 0 | 0 | ← $n_4$ |

×

R

| 1 | 2 | 0 | -1 | -1 |

==

D

| 2 | 2 | 1 | 1 | 0 |

Max sum per path

E

| 10 | $e_1$ |
| 40 | $e_2$ |
| 50 | $e_3$ |
| 30 | $e_4$ |
| 20 | $e_5$ |

×

S

| 0 | 1 | 0 | 0 | 0 |

T

| 40 |

# Hummingbird method



$$P \leftarrow X \times A \qquad Q \leftarrow P < B \qquad R \leftarrow Q \times C \qquad S \leftarrow R == D \qquad T \leftarrow \sum S * L_Q$$

max bits : p       p       log2(d)       log2(d)       p

                                                       + 1          + 1

Max bit-width reached : $\mathbf{max(p, log_2(d) + 1)}$

- p : number of bits of quantization used for inputs

  and outputs

- d : tree depth

# Experimental results : XGBoost

XGBoost classifier model (max_depth=3, n_estimators=50) on data set "Spambase"

# Experimental results : spambase

|  |  | accuracy | f1 | AP | #nodes | Time (s) | FHE/Clear |
|---|---|---|---|---|---|---|---|
| spambase | FHE-DT | 91.0% | 88.0% | 84.3% | 23 | 1.313 | 825x |
|  | FP32-DT | 90.3% | 87.4% | 82.4% | - | 0.002 | - |
|  | FHE-XGB | 93.1% | 90.9% | 87.7% | 350 | 7.020 | 4617x |
|  | FP32-XGB | 93.6% | 91.7% | 88.3% | - | 0.002 | - |
|  | FHE-RF | 90.9% | 87.5% | 84.6% | 750 | 16.248 | 8520x |
|  | FP32-RF | 91.8% | 89.0% | 86.0% | - | 0.002 | - |

n_bits=6, max_depth={4,5}, n_estimators=50

8 cores

# Concrete / Concrete ML

# Concrete stack



**Concrete ML**  |  **3rd party applications**  |  **3rd party frameworks**  —  **Applications**

**Python frontend**  |  **3rd party frontends**  —  **Frontends**

**MLIR based compiler**  |  **3rd party compilers**  —  **Compilers**

**CPU**  |  **GPU**  |  **3rd party backends**  —  **Backends**

# Scikit–learn API's for ML

```python
from concrete.ml.sklearn import LogisticRegression

model = LogisticRegression(n_bits=12)
model.fit(X_train, y_train)
model.predict(X_test)
model.compile(X_train)
model.predict(X_test, fhe="simulate")
model.predict(X_test, fhe="execute")
```

```python
from concrete.ml.sklearn import XGBClassifier

model = XGBClassifier(n_bits=8)
model.fit(X_train, y_train)
model.predict(X_test)
model.compile(X_train)
model.predict(X_test, fhe="simulate")
model.predict(X_test, fhe="execute")
```

**No need to know cryptography!**

# Torch API's for DL

also support for: Tensorflow

ONNX

```python
from transformers import AutoModel
from concrete.ml.torch.compile import compile_torch_model

# Load model from Hugging Face Hub
model = AutoModel.from_pretrained("dacorvo/mnist-mlp")

# Convert to FHE
q_module = compile_torch_model(
    model,
    torch_inputset=data,
)

fhe_outputs = q_module.forward(test_data, fhe="execute")
```

# Confidential training

## Federated Learning

- Federated learning protects training data
- Works for big models

Models trained with FL can be **deployed with Concrete ML**

```python
from concrete.ml.sklearn import LogisticRegression

with open("federated_trained_model.pkl", "rb") as f:
    federated_model = pickle.load(f)

fhe_model = LogisticRegression.from_sklearn_model(model)
fhe_model.compile()

model.predict(X_test, fhe="execute")
```

## Training on encrypted data

- FHE protects training data
- Encrypted models are trained on encrypted data
- Works for small models like LogisticRegression, MLP, and later for larger models

```python
from concrete.ml.sklearn import SGDClassifier

sgd_clf_encrypted = SGDClassifier(fit_encrypted=True)
sgd_clf_encrypted.fit(X_binary, y_binary, fhe="execute")

y_pred = sgd_clf_encrypted.predict(X_binary, fhe="execute")
```
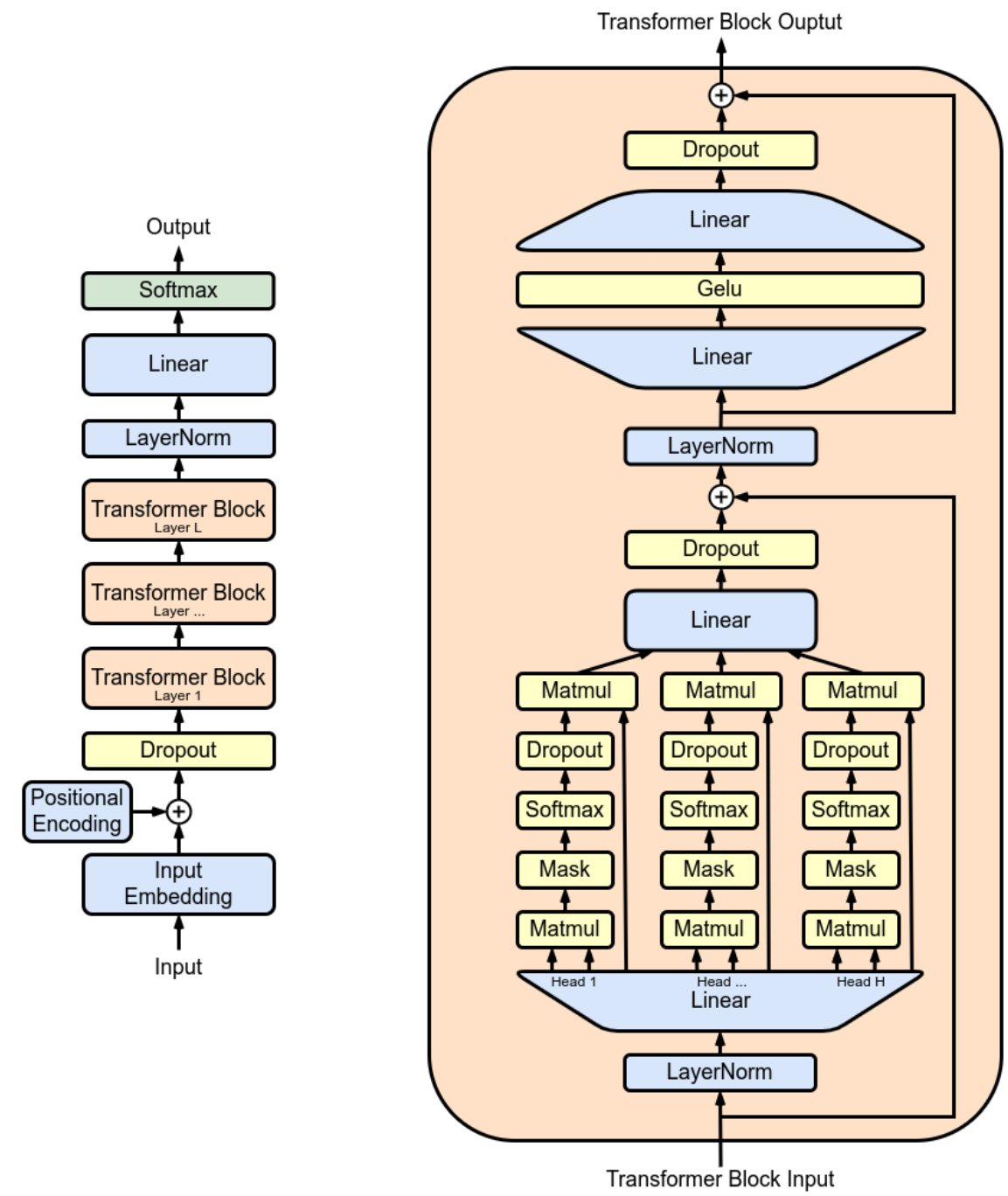
# Large Language Models

# LLM in FHE

## Operations

- Linear
- Embedding
- Non-Linear Activation (e.g. Gelu)
- Attention
  - Softmax
  - Encrypted matrix multiplication

# LLM in FHE

## Operations

## FHE Complexity

- Linear
- Embedding
- Non-Linear Activation (e.g. Gelu)
- Attention
  - Softmax
  - Encrypted matrix multiplication

ms
ms
ms

seconds
minutes

# LLM in FHE

## Operations

## Weight distribution

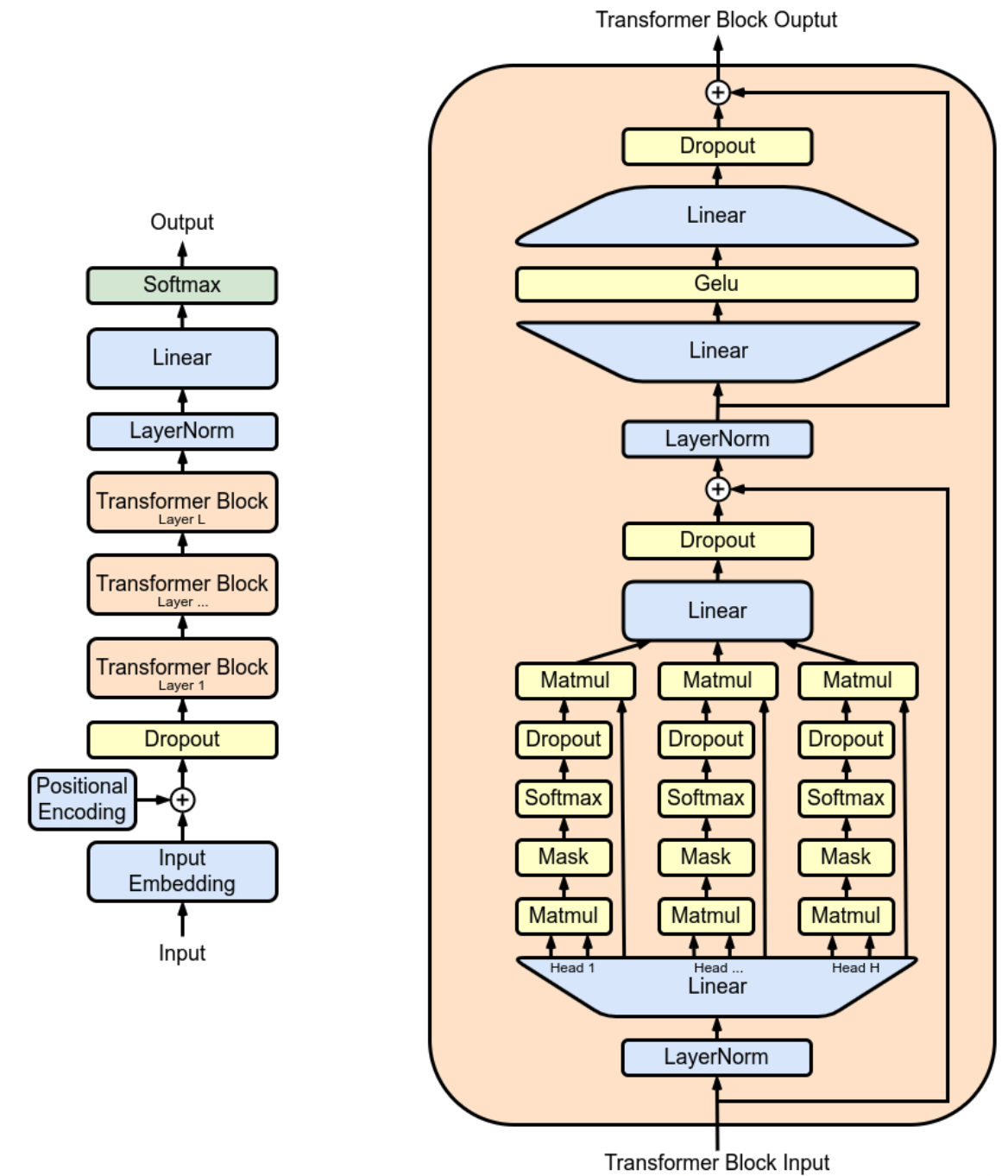- Linear                                    **99%**
- Embedding                                 1%
- Non-Linear Activation (e.g. Gelu)         0%
- Attention
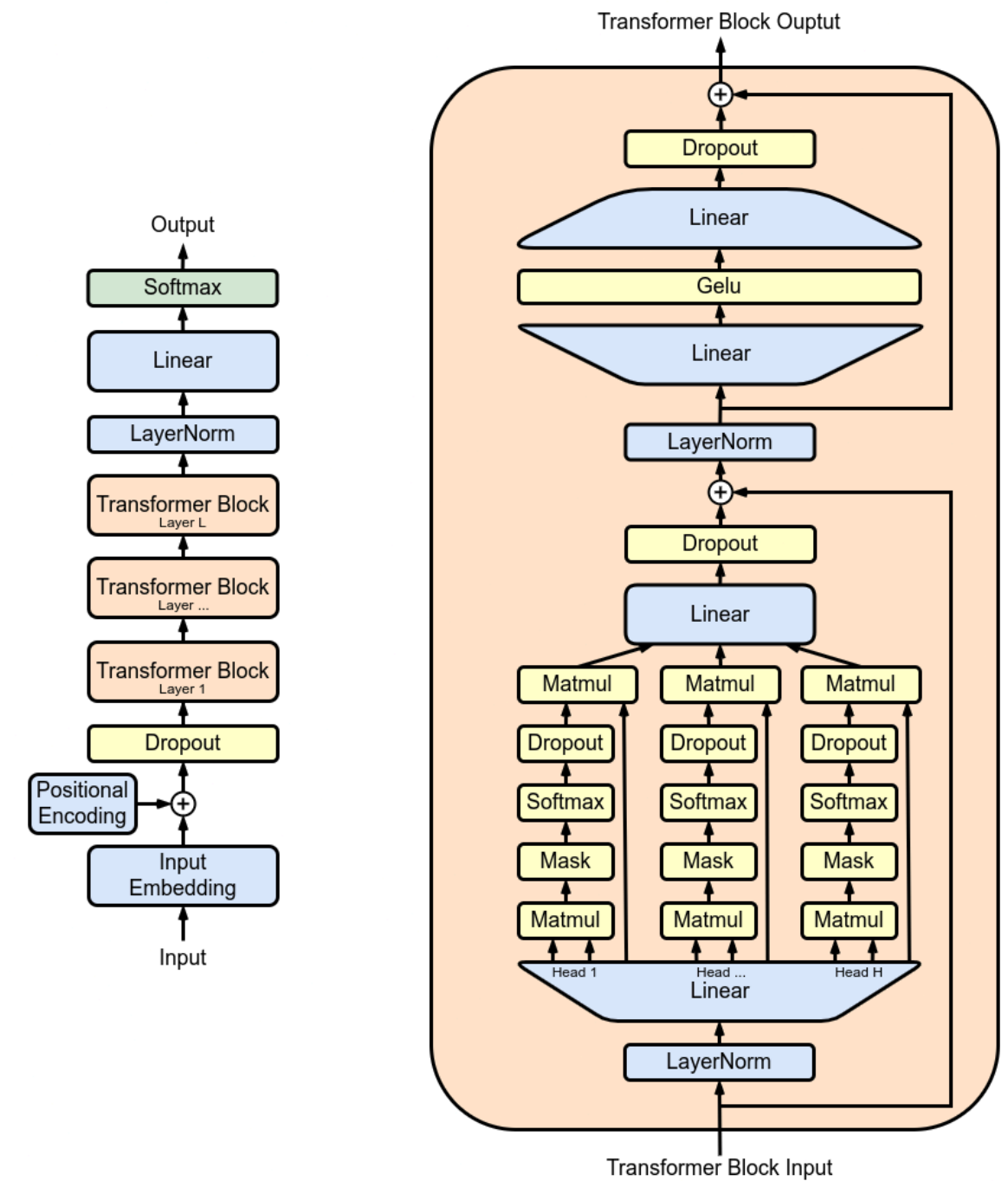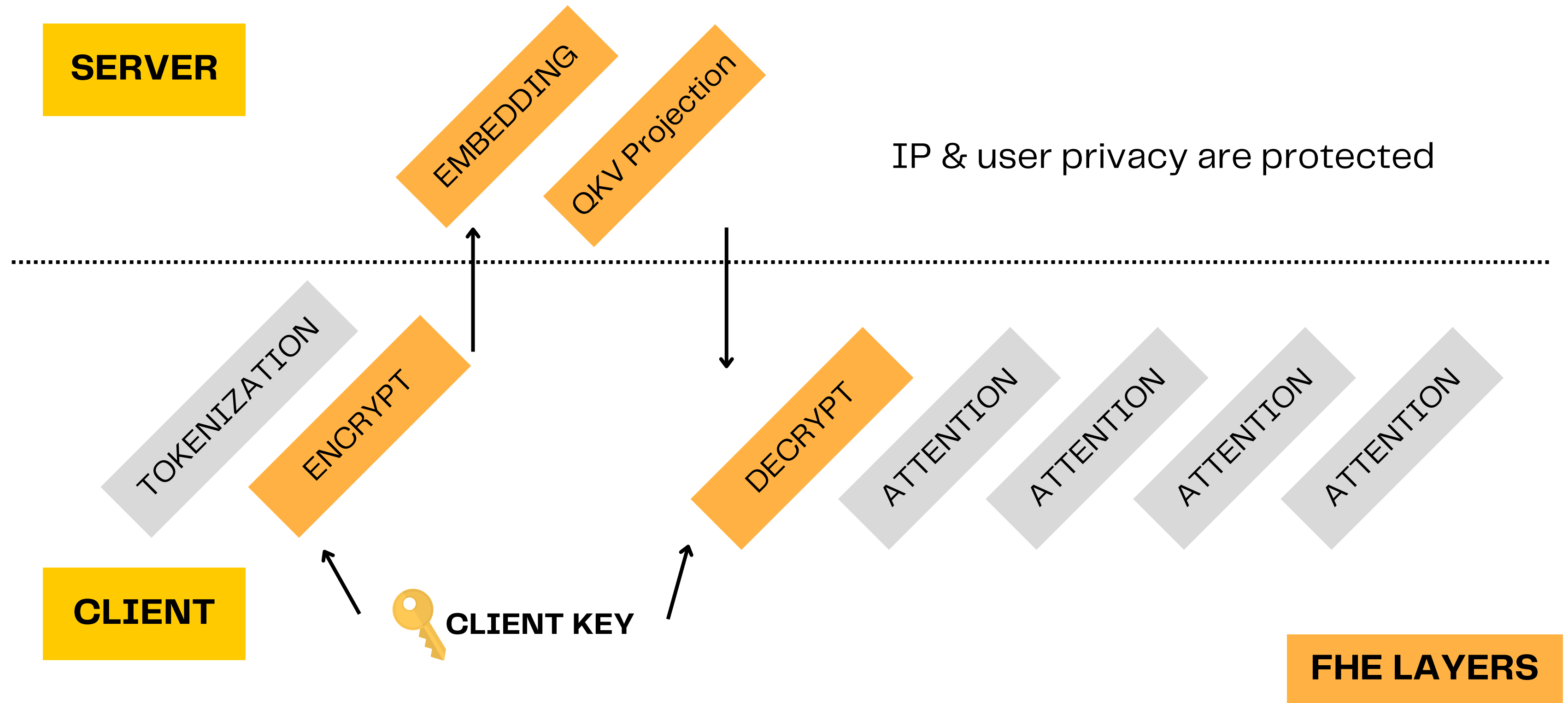  - Softmax                                 0%
  - Encrypted matrix multiplication         0%

# Work in Progress

# Secure anonymization for private ChatGPT

https://huggingface.co/spaces/zama-fhe/encrypted-anonymization

# Collaborative computation

- Several companies wanting to co-compute a common function, without sharing their data
- Companies manage the collaborative decryption of results thanks to treshold decryption

```python
def f(enc_a, enc_b, enc_c):

    encrypted_concatenated_features = np.concatenate([enc_a, enc_b, enc_c])

    encrypted_predictions = fhe_model.run(encrypted_concatenated_features)

    return encrypted_predictions
```

# Data-frames

- Persist encrypted tabular data while allowing FHE processing
- Encapsulate pre-processing before encryption to streamline FHE processing

## Client Side

- Encrypt pandas data-frames

```
client = ClientEngine(keys_path=client_1_keys_path)
df_left = pandas.read_csv("df_left.csv")
df_left_enc = client.encrypt_from_pandas(df_left)
```

| | index | day | time | size |
|---|---|---|---|---|
| **0** | 2 | Thur | Lunch | 2 |
| **1** | 5 | Sat | Dinner | 3 |
| **2** | 9 | Sun | Dinner | 2 |

## Server Side

- Join encrypted data-frames

```
df_left_enc = load_encrypted_dataframe(df_left_enc_path)
df_right_enc = load_encrypted_dataframe(df_right_enc_path)


df_joined_enc_server = df_left_enc.merge(df_right_enc, how=HOW, on=ON)
```
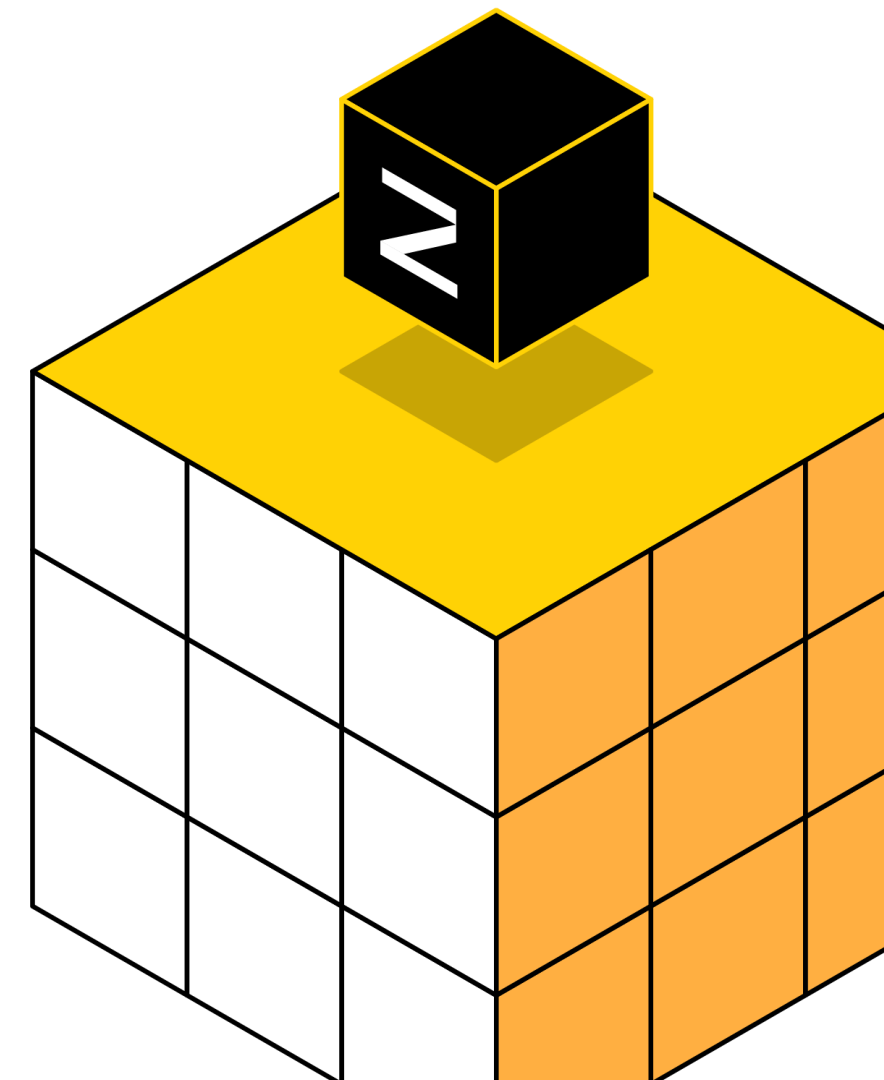
| index | day | time | size |
|---|---|---|---|
| ..48d4814937.. | ..dd6b288e52.. | ..497a80e2dd.. | ..41f496fe3a.. |
| ..0a19fbfc58.. | ..047a92f5bc.. | ..7f7a6f1167.. | ..5ca8e5edfc.. |
| ..79c726effe.. | ..6835b68ece.. | ..4ae3bca370.. | ..f4eb2bde07.. |

# Starting Building

# We are open-source

- We're open source: have a look to
    - **repo**: https://github.com/zama-ai/concrete-ml
    - **docs**: http://docs.zama.ai/concrete-ml
    - **live demos on Hugging Face**: https://huggingface.co/zama-fhe

- And if you're interested by the scientific side:
    - Neural Network Training on Encrypted Data with TFHE
    - Deep Neural Networks for Encrypted Inference with TFHE (CSCML 2023)
    - Privacy-Preserving Tree-Based Inference with TFHE (MSPN 2023)

# How to learn and make yours

- Reproduce examples from the documentation
- Then:
    - create some task in ML (linear or tree-based models)
    - once it works, compile it with Concrete ML
    - once familiar, continue with an easy DL example, have a look to QAT

- We'll be happy to support you on discord.fhe.org/

```
import concrete.numpy as hnp

def add(x, y):
    return x + y

print(f"Compiling...")
circuit = compiler.compile_on_inputset(inputset)

inputset = [(2, 3), (0, 0), (1, 6), (7, 7), (7, 1),
compiler = hnp.NPFHECompiler(add

examples = [(3, 4), (1, 2), (7, 7),

for example in examples:

    result = circuit.run("example
    print(f"Evaluation of {' + '.join(

                                orphically = (resu
```

# Contact and Links

benoit.chevalliermames@zama.ai

jordan.frery@zama.ai

zama.ai

github.com/zama-ai

community.zama.ai/

discord.fhe.org

**ZAMA**