



**РусКрипто**

**XXVII**

**НАУЧНО-ПРАКТИЧЕСКАЯ  
КОНФЕРЕНЦИЯ**

# Квантово-безопасное федеративное глубокое обучение



РусКрипто

Обзор на основе работы:

Авторы: Kfir Sulimany и др.  
Лаборатория: MIT Research  
Laboratory of Electronics

## Quantum-secure multiparty deep learning

Kfir Sulimany<sup>\*1</sup>, Sri Krishna Vadlamani<sup>1</sup>, Ryan Hamerly<sup>1,2</sup>, Prahlad Iyengar<sup>1</sup>, and Dirk Englund<sup>1</sup>

<sup>1</sup>Research Laboratory of Electronics, Massachusetts Institute of Technology, Cambridge, MA, USA

<sup>2</sup>Physics & Informatics Laboratories, NTT Research, Inc., Sunnyvale, CA, USA

Автор обзора: Воробей Сергей С.

Заместитель начальника отдела лицензирования и сертификации ООО «КуРэйт»

# 1. Предпосылки развития в РФ



Трехуровневая система регулирования позволит создать стимулирующие условия для обмена данными



## Тензорные вычислители

Включен в перечень по 719 и 878 Постановлению Правительства РФ



Слайд из презентации Сорокина Павла Юрьевича, Первый заместитель Министра энергетики Российской Федерации



<a href="#">УБИ. 218</a>	Угроза раскрытия информации о модели машинного обучения	конф.
<a href="#">УБИ. 219</a>	Угроза хищения обучающих данных	конф.
<a href="#">УБИ. 220</a>	Угроза нарушения функционирования («обхода») средств, реализующих технологии искусственного интеллекта	конф.
<a href="#">УБИ. 221</a>	Угроза модификации модели машинного обучения путем искажения («отравления») обучающих данных	цел.
<a href="#">УБИ. 222</a>	Угроза подмены модели машинного обучения	конф. цел.

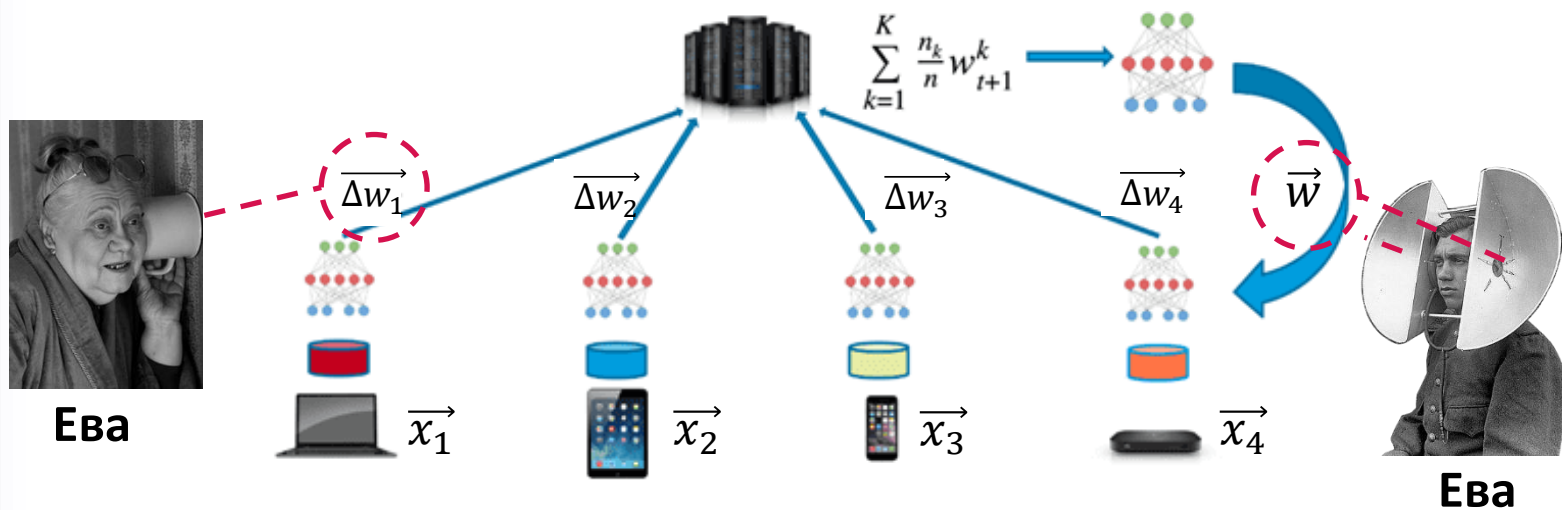
## 2. Конфиденциальность федеративного глубокого машинного обучения



РусКрипто

Вычисление  $f(\vec{x}, \vec{w}) = \Delta \vec{w}$ , не раскрывая аргументов:

- $\vec{x}$  – приватные данные клиента,
- $\vec{w}$  – веса модели сервера,
- $\Delta \vec{w}$  – промежуточные градиенты



### 3. Недостатки гомоморфного шифрования (homomorphic encryption - HE)



РусКрипто

$$Enc(f(x)) = f(Enc(x))$$

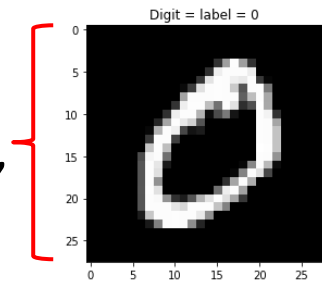
$O(n^k)$  – сложность, где

$n$  – размер входных данных:

- $\vec{x} \in \mathbb{R}^N$  – локальный вектор данных клиента
- $\vec{w} = vec(\mathbf{W}), \mathbf{W} \in \mathbb{R}^{N \times M}$  – вектор когерентных состояний весов. MNIST:  $N = 784$  (28×28 пикселей),  
 $M$  – число нейронов в слое
- $k$  – сложность ( $P(2 \leq k \leq 3) \approx 1$ )

Недостатки гомоморфного шифрования:

- Вычислительная сложность
- Уязвимо к кв.вычислениям



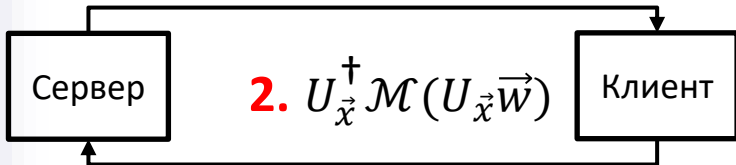
$$\vec{x} = vec(\mathbf{Fig})$$



# 4. Архитектура решения

$$\underbrace{\vec{W}}_{\text{сервер}} \rightarrow \underbrace{U_{\vec{x}} \rightarrow \mathcal{M}(U_{\vec{x}}\vec{W}) \rightarrow U_{\vec{x}}^\dagger}_{\text{клиент}} \rightarrow \rho_v$$

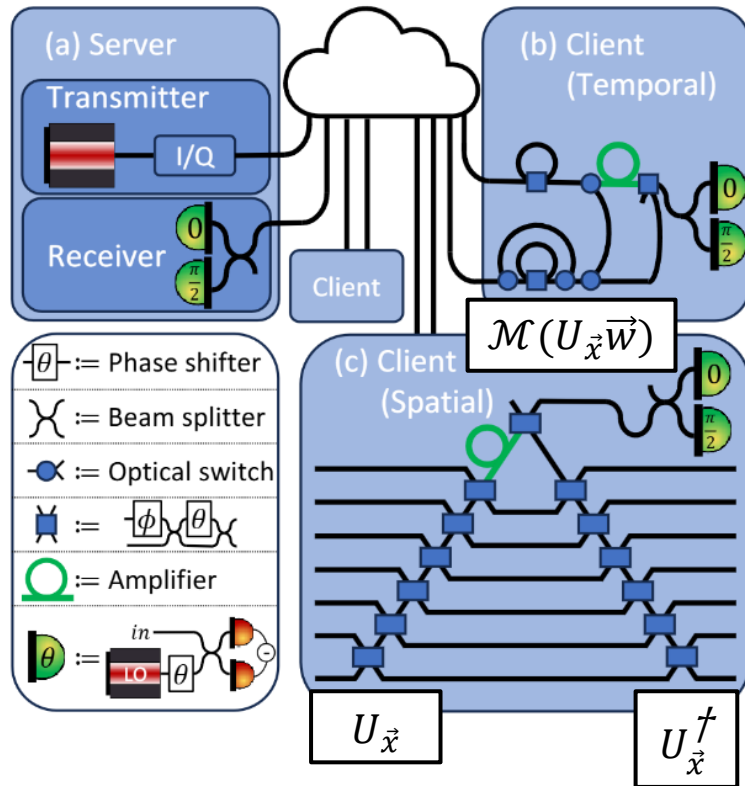
1.  $W_i \rightarrow \vec{w}_i = \{|\alpha_0\rangle, \dots\}$  – Когерентные СОСТОЯНИЯ ВЕСОВ



3.  $\rho_v$  – Верификационное СОСТОЯНИЕ

$\alpha_j = \frac{\sqrt{\mu} W_{ij}}{\|W\|_{RMS}}$ , где  $\mu$  – ср. число фотонов

$$\vec{x} \rightarrow \hat{x} = \frac{\vec{x}}{\|\vec{x}\|}$$





# 5. Операции на клиенте

$$\vec{w} \rightarrow \underbrace{U_{\hat{x}} \rightarrow \mathcal{M}(U_{\hat{x}} \vec{w})}_{\text{Client}} \rightarrow U_{\hat{x}}^\dagger \rightarrow \rho_v$$

## 1 Измерение

$$U_{\hat{x}} \vec{w} = (\vec{w} \cdot \hat{x}, v_2, \dots, v_N) = (\alpha, v_2, \dots, v_N)$$

Клиент измеряет обе квадратуры ( $\hat{X}$ ,  $\hat{P}$ ) когерентного состояния  $|\alpha\rangle$  с помощью гомодинного детектирования.  $|\alpha\rangle$  имеет дисперсию в 1 SNU (Shot Noise Unit)

$$\langle (\Delta \hat{X})^2 \rangle = \langle (\Delta \hat{P})^2 \rangle = 1 \text{ SNU}$$

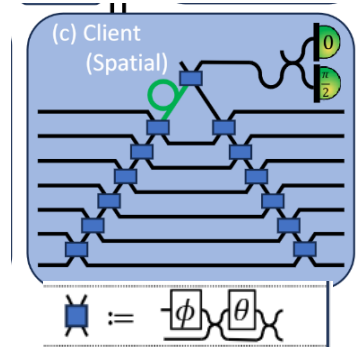
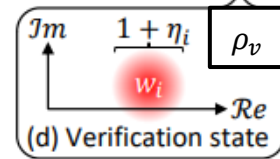
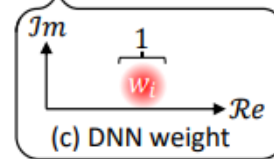
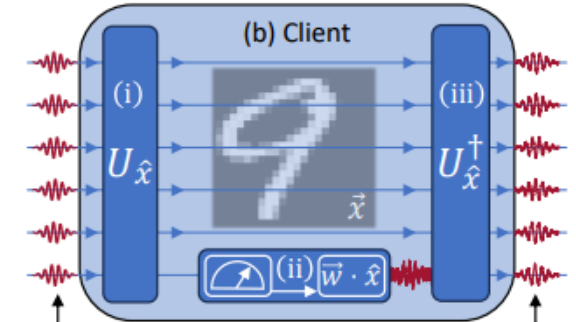
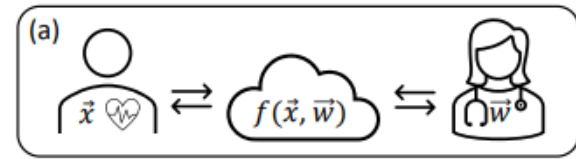
## 2 Прямое распространение

$$\mathcal{M}(U_{\hat{x}} \vec{w}) = (\tilde{\alpha}, v_2, \dots, v_N)$$

Операция  $\mathcal{M}$  обобщается через  $\mathcal{G}$ , которая включает:

- Усиление с коэффициентом  $G$  (phase-insensitive amplification);
- Разделение луча с соотношением  $1 - \frac{1}{G} : \frac{1}{G}$ ;

Где  $\tilde{\alpha}$  – новое состояние в результирующей моде



## 6. Операции на клиенте

$$\vec{w} \rightarrow U_{\vec{x}} \rightarrow \mathcal{M}(U_{\vec{x}} \vec{w}) \rightarrow \underbrace{U_{\vec{x}}^\dagger}_{\rho_v}$$

### Прямое распространение

После измерения  $U_{\vec{x}}^\dagger$  добавленный шум распределяется по всем модам верификационного состояния  $\rho_v$ , для  $i$ -й моды шум:

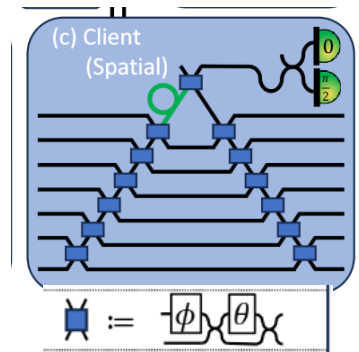
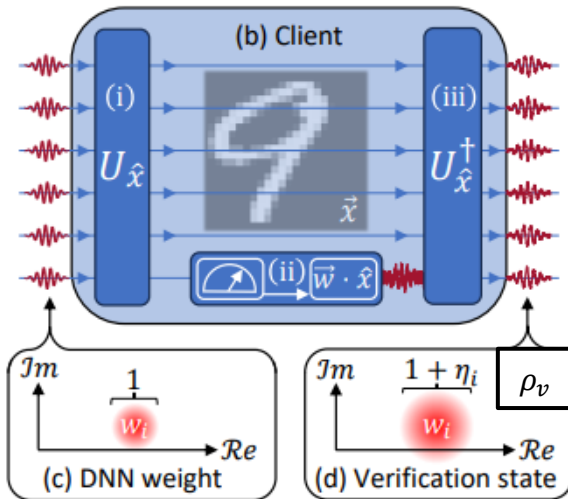
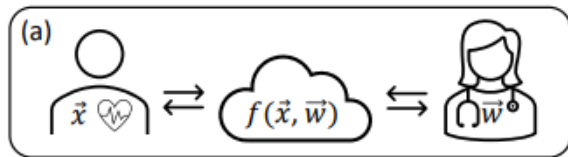
$$\eta_i = \left(2 - \frac{2}{G}\right) |\hat{x}_i|^2$$

$$\rho_v = U_{\vec{x}}^\dagger \mathcal{M}(U_{\vec{x}} \vec{w})$$

Среднее значение  $\rho_v$  совпадает с  $\vec{w}$ :  $\langle \tilde{\alpha} \rangle = \vec{w} \cdot \hat{x}$ , но дисперсия каждой моды увеличивается на  $\eta_i$ :

$$\sigma^2 = 1 + \left(2 - \frac{2}{G}\right) |\hat{x}_i|^2 = 1 + \eta,$$

где  $\eta$  – избыточный шум, что маскирует информацию о  $\vec{x}$ .





## 7. Утечка $w_i$ , оценка границы Холево



РусКрипто

Операция  $\mathcal{M}$  добавляя квантовый шум:

- Защищает  $\vec{w}$  сервера, так как клиент измеряет только проекцию  $\vec{w} \cdot \hat{x}$ , а остальной шум маскирует исходные веса.

$$I_{w_i} \leq \chi = S(\rho_{AB}) - S(\rho_{B|A}) \rightarrow I_{w_i} \leq g(v_1) + g(v_2) - g(v_3)$$

$$g(v) = \left(\frac{v+1}{2}\right) \log_2 \left(\frac{v+1}{2}\right) - \left(\frac{v-1}{2}\right) \log_2 \left(\frac{v-1}{2}\right)$$

$$a = 2\mu + 1, b = 2\mu + 1 + \eta_i, c = \sqrt{4\mu^2 + 2\mu + 1}, z = \sqrt{(a+b)^2 - 4c^2}$$
$$v_{1,2} = \frac{1}{2}(z \pm [b-a]), v_3 = a - \frac{c^2}{b+1}$$

## 8. Утечка $x_i$ , оценка границы Крамера-Рао



РусКрипто

Операция  $\mathcal{M}$  добавляя квантовый шум:

- Защищает  $\vec{x}$  клиента, т.к. обратно высылается верификационное состояние и только в дисперсии содержится информация о данных. Сервер видит только шум  $\eta_i$ , зависящий от  $|\hat{x}_i|^2$ :

Квантовая граница Крамера-Рао (QCRB) – оценка границы дисперсии случайной квантовой величины.

$$\text{Var}(\hat{x}_i) \geq \frac{1}{M \mathcal{F}_R[\rho_x]}, \text{ где } M \text{ – число измерений, } \mathcal{F}_R[\rho_x] = \frac{4|\hat{x}_i|^2 \left(2 - \frac{2}{G}\right)^2}{\sigma_i^4}$$

$$I_{x_i} = \frac{1}{2} \log_2 \left( 1 + \frac{|\hat{x}_i|^2}{\text{Var}(\hat{x}_i)} \right)$$

$$I_{x_i} \leq \frac{1}{2} \log_2 \left( 1 + k \cdot \frac{8M(G-1)^2 |\hat{x}_i|^4}{G^2 \sigma_i^4} \right), \quad k = 2$$

(следствие из границы Крамера-Рао)

## 9. Компромисс точности и безопасности



РусКрипто

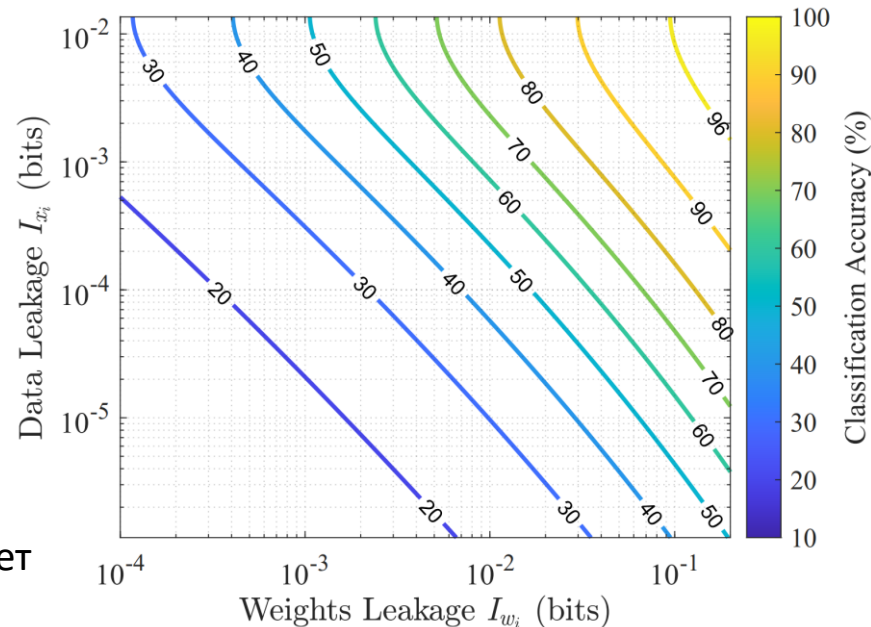
Естественный компромисс: уменьшение  $\mu$  снижает  $I_{w_i}$ , но требует роста  $G$ , увеличивая  $I_{x_i}$ .

$$I_{w_i} < 0.1 \text{ (для } \mu = 4\text{)}$$

$$I_{x_i} < 0.01 \text{ (для } G = 3\text{)}$$

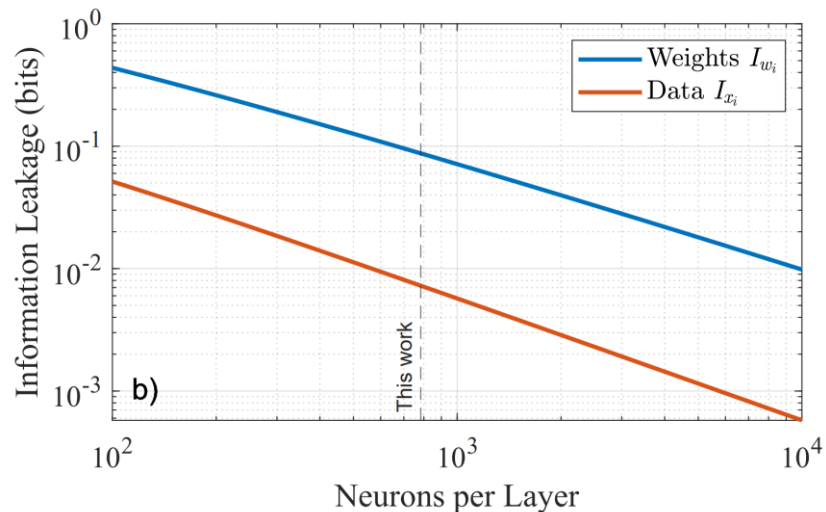
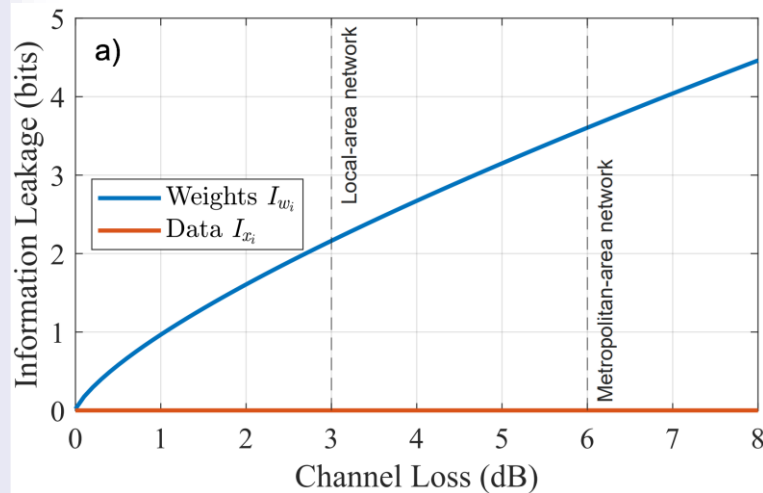
Современные стандарты квантования, необходимые для обеспечения минимальной точности для DNN, соответствуют 1 биту\*.

Даже зная утечку, атакующий не сможет восстановить модель без доступа к обучающим данным.



\* – Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference. In Low-Power Computer Vision, pages 291–326. Chapman and Hall/CRC, 2022.

# 10. Зависимость утечки от параметров



Потери в каналах. При 6 дБ – типичных для локальных сетей –  $I_{w_i}$  растет до 4 бит, но  $I_{x_i}$  не зависит от потерь. С увеличением числа нейронов утечка падает: шум распределяется по модам. Для больших DNN, с тысячами нейронов, безопасность только возрастает.



РусКрипто

СПАСИБО  
ЗА ВНИМАНИЕ

QRATE