

# Post quantum secure fair data trading with deterability based on machine learning

Jinhui Liu<sup>1,2</sup>, Yong Yu<sup>3\*</sup>, Hongliang Bi<sup>1</sup>, Yanqi Zhao<sup>3</sup>, Shijia Wang<sup>4</sup> & Huanguo Zhang<sup>5</sup>

<sup>1</sup>*School of Cybersecurity, Northwestern polytechnical university, Xi'an, 710072, China;*

<sup>2</sup>*Research & Development Institute of Northwestern Poly-technical University, Shenzhen 518057, Guangdong, China;*

<sup>3</sup>*School of Cybersecurity, Xi'an University of Posts and Telecommunications, Xi'an, 710121, China;*

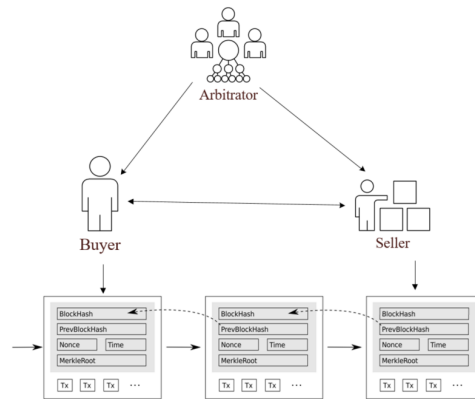
<sup>4</sup>*School of Statistics and Data Science, LPMC & KLMDASR, Nankai University, Tianjin, 300071, China.;*

<sup>5</sup>*School of Cyber Science and Engineering, Wuhan University, Wuhan, 430072, China.*

## Appendix A System architecture and design goals

### Appendix A.1 System architecture

The proposed system architecture consists of four entities: data owner (seller), buyer, blockchain, and arbitrator. Both the data owner and buyer make transactions and maintain the blockchain. The following graph describes the system architecture.



**Figure A1** A framework of our fair data trading

- Data owner: As a data collector, he generates data index list, makes transactions with data payer, performs data exchange and detonates by contact to obtain charges.
- Data payer: He verifies the effectiveness of data, makes transactions with data owner, deploys smart contracts and pays for satisfied data.
- Blockchain: It records transactions, publicises parameters, maintains transaction bills between the data owner and buyer and performs smart contracts.
- Arbitrator: When there exists some disputes, the arbitrator responds to them.

In the system architecture, since the signer uses the public key of the designed verifier seller to make a signature, only the verifier can verify the signature using his own secrete key. There does not need other authentication scheme. Aiming at the arbitrator, which likes a court or a police office, only there exists some disputes it will make a judgment. If the arbitrator is not honest, the arbitrator can also ensures itself to other entities that it is legitimate and is allowed to perform specified (but not malicious) actions by using the method of Quantum2FA [1].

\* Corresponding author (email: yuyongxy@163.com)

## Appendix A.2 Design goals

- **Completeness:** If the data owner and the data payer honestly perform the fair data trading protocol, data owner will obtain the money and data payer will obtain valid data.
- **Fairness:** Data owner obtains money and the data payer obtains valid data, otherwise they achieves nothing.
- **Accountability:** If both data payer and data owner do some mis-behaviors, they lose their deposits.

## Appendix B The proof of Theorem 1-8

**Theorem 1.** The constructed DVDAPS scheme provides **Correctness**.

*Proof.* The property of **Correctness** can be verified by the **Dver** algorithm.

**Theorem 2.** If  $\mathcal{F}$  is a pseudo-random function and  $f$  is a one way function, the DVDAPS has EUF-CMA security with a negligible probability. *Proof.* First, we demonstrate that the non-interactive knowledge proof system in this DVDAPS scheme is simulation sound, extractable under the linear-encryption assumption with soundness error  $\frac{(4d+2)}{p}$ . Therefore, this theorem is proven using a sequence of games.

Using the definition of zero knowledge, we constructed two PPT simulators  $Sim_1$  and  $Sim_2$ . Here  $(crs, vrs, trap) \leftarrow Sim_1(1^\lambda, \mathcal{R})$  and  $(\hat{A}, \hat{B}, \hat{C}) \leftarrow Sim_2(crs, u, trap)$ . Suppose that  $trap = (\alpha, \beta, \delta, s)$ ,  $Sim_2$  picks random values  $a, b \in \mathbf{Z}_p$  and calculates

$$C(s) = \frac{ab - \alpha\beta - \sum_{i \in I_{public}} d_i(\beta a_i(s) + \alpha b_i(s) + c_i)}{\delta}.$$

It generates the LWE-based signature  $\sigma = (z, \hat{A}, \hat{B}, \hat{C})$  satisfying  $\hat{A} = \text{Enc}(a)$ ,  $\hat{B} = \text{Enc}(b)$ ,  $\hat{C} = \text{Enc}(C(s))$ . Thus  $ct$  satisfies the **Dver** algorithm and its distribution is statistically distinguishable from the output of **Dsign**.

If, for any non-negligible function  $negl(\lambda)$ , a PPT-knowledge-soundness adversary  $\mathcal{A}_{snd}$  exists, an extractor  $\mathcal{E}_{\mathcal{A}_{snd}}$  with witness  $\omega$  can be constructed to show that the non-interactive proof system in DVDAPS can extract knowledge soundness. Therefore,, we can compute  $\alpha, \beta, \delta, s$  with probability  $1 - \frac{4d+2}{p}$  using the Schwartz-Zippel lemma.

$$A(\alpha, \beta, \delta, s) = A_\alpha \alpha + A_\beta \beta + A_\delta \delta + A(s) + \sum_{i \in I_{mid}} A_i \frac{\beta a_i(s) + \alpha b_i(s) + c_i(s)}{\delta} + A_h(s) \frac{t(s)}{\delta},$$

we can also construct  $B(\alpha, \beta, \delta, s)$  and  $C(\alpha, \beta, \delta, s)$  in a similar way, where  $A_\alpha, A_\beta, A_\delta, B_\alpha, B_\beta, B_\delta$  and  $C_\alpha, C_\beta, C_\delta$  are over  $\mathbf{Z}_p$ ,  $A(s), A_h(s), B(s), B_h(s), C(s), C_h(s)$  are the polynomials of degree  $d$ .

Let  $S_i$  be the winning event of the game  $G_i$ .

**Game  $G_0$ :** The original game.

**Game  $G_1$ :** As  $G_0$ , we modify **DkeyGen**( $1^\lambda$ ), let  $(crs, \tau) \leftarrow S_{1, \Pi}(1^\lambda)$  and store  $\tau$ .

Both  $G_1$  and  $G_0$  are indistinguishable under zero knowledge of  $\Pi$  with probability distance less than  $\text{Adv}_{\mathcal{A}, \Pi}^{\text{Sim}}(\lambda)$ .

**Game  $G_2$ :** Because  $G_1$ , we modify **Dsign** and let  $\pi \leftarrow S_{2, \Pi}(crs, \tau, pk)$ .

Both  $G_2$  and  $G_1$  are also indistinguishable under zero knowledge of  $\Pi$  with a probability distance less than  $\text{Adv}_{\mathcal{A}, \Pi}^{\text{ZK}}(\lambda)$ .

**Game  $G_3$ :** Because  $G_2$ , we modify **DkeyGen** and **Dsign**, but let  $\hat{c} \leftarrow \mathbf{Z}_q$  in **DkeyGen** and  $\rho \leftarrow \mathbf{Z}_q$  in **DVsign**.

We construct a PRF challenger  $\mathcal{C}$  against the pseudo-random function  $\mathcal{F}$  which means  $\hat{c} \leftarrow \mathcal{C}(\hat{\beta})$  in **DkeyGen** and  $\rho \leftarrow \mathcal{C}(m_0)$  in **Dsign**. Thus, an adversary distinguishes  $G_3$  and  $G_2$  means that  $\mathcal{A}$  distinguishes the PRF from a random function with a probability distance less than  $\text{Adv}_{\mathcal{F}}(\lambda)$ .

**Game  $G_4$ :** Because  $G_3$ , we trace all  $(m_0, \rho)$  pairs in  $\mathcal{Q}$ . If  $(m'_0, \rho), (m''_0, \rho) \in \mathcal{Q}$  exists such that  $m'_0 \neq m''_0$ , we abort.

$G_3$  and  $G_4$  proceed identically unless an abort event exists. Therefore, the probability distance between them is bounded by  $\frac{Q_{\text{sign}}}{q}$ , where  $Q_{\text{sign}}$  denotes the number of signature queries.

**Game  $G_5$ :** Because  $G_4$ , we modify **DVsign** and let  $\hat{z} \leftarrow \mathbf{Z}_q$ .

Because  $\rho$  is uniformly random without being revealed, and  $\hat{z}$  is also uniformly random,  $G_5$  and  $G_4$  proceed identically.

**Game  $G_6$ :** Similar to Game  $G_5$ , we modify **DkeyGen**, let  $(crs, \tau, \xi) \leftarrow \mathcal{E}_\pi(1^\lambda)$  and store  $(\tau, \xi)$ .

$G_6$  and  $G_5$  are indistinguishable under the simulation sound extractability property of  $\Pi$ .

**Game  $G_7$ :** Because  $G_6$ , we use an extractor to obtain  $sk_\Sigma^* \leftarrow \mathcal{E}(crs, \xi, (\hat{\beta}, \hat{c}, m_0, \hat{z}_1, \hat{\mathbf{b}}, \hat{A}, \hat{B}, \hat{C}), \pi)$  and abort if the extraction fails.

The probability distance between Games  $G_7$  and Game  $G_6$  is bounded by  $\text{Adv}_{\mathcal{A}, \mathcal{E}, \Pi}^{\text{Ext}}(\lambda)$ . Finally, Game  $G_7$  presents a reduction when engaging with an OWF challenger. Thus, the success probability of a successful game is bounded by  $\text{Adv}_{\mathcal{A}, \mathcal{E}, \Pi}^{\text{Ext}}(\lambda)$ .

Hence the probability  $\text{Pr}[S_7]$  is bounded by a negligible probability and the proposed construct scheme provides EUF-CMA security.

**Theorem 3.** The proposed DVDAPS scheme provides double signature extractability if the non-interactive proof system is simulation-sound extractable and the pseudo-random function  $\mathcal{F}$  is computational fixed value key binding, and the DVDAPS provides double signature extractability.

*Proof.* We also use games to prove this theorem. Let  $S_i$  be the winning event in game  $G_i$ . Suppose that  $m = (m_0, m_1), m' = (m_0, m_2), \sigma_1 = (\cdot, z_1, \pi_1), \sigma_2 = (\cdot, z_2, \pi_2)$ .

**Game  $G_0$ :** The original double signature extractability game.

**Game  $G_1$ :** Because  $G_0$ , we modify **DkeyGen** and let  $(crs, \tau) \leftarrow S_{1, \Pi}(1^\lambda)$  and store  $\tau$ .

**Game  $G_2$ :** Because  $G_1$ , we modify **DkeyGen** and let  $(crs, \tau, \xi) \leftarrow \mathcal{E}_{1, \Pi}(1^\lambda)$  and store  $\xi$ .

**Game  $G_3$ :** As  $G_2$ , but we use the extractor to obtain  $sk_{\text{PRF}}^* \leftarrow \mathcal{E}_{1,\Pi}(crs, \xi, pk, m, m', \pi)$  and aborts if the extraction fails.

**Game  $G_4$ :** Because  $G_3$ , if  $sk_{\text{PRF}}^* \neq sk_{\text{PRF}}$  we abort.

From the  $G_0 \Rightarrow$  transition  $G_1 \Rightarrow$  transition  $G_2 \Rightarrow$  transition  $G_3 \Rightarrow$  and  $G_4$  transitions, the games are indistinguishable under adaptive zero knowledge, simulation-sound extractability of the proof system and the assumption of fixed-value-key-binding.

**Theorem 4.** The proposed DVDAPS scheme provides the signer’s privacy property .

*Proof.* Because the LWE-based encryption scheme  $\Sigma$  is IND-CCA2-secure with the security parameter  $\lambda$  [?], we construct an algorithm adversary  $\mathcal{A}'$  as follows:

$\mathcal{A}'$  selects a pair of verification keys  $(sk_B, pk_B)$  and two pairs of signing keys  $(sk_{A_0}, pk_{A_0})$  and  $(sk_{A_1}, pk_{A_1})$ .

For any signing query,  $\mathcal{A}$  and  $\mathcal{A}'$  answer either  $A_0$  or  $A_1$  using secret keys. For any verification query,  $\mathcal{A}$  and  $\mathcal{A}'$  answer  $B$  and the decryption oracle uses secret keys.

$\mathcal{A}$  outputs a message  $m^*$  and  $\mathcal{A}'$  computes its two signatures  $\sigma_0, \sigma_1$  using the  $\text{Dsign}$  algorithm. Subsequently,  $\mathcal{A}'$  queries signatures and sends these queries to the IND-CCA2 challenger  $\mathcal{C}$ .  $\mathcal{C}$  encrypts  $\sigma_b, b \in \{0, 1\}$ .

$\mathcal{A}'$  sends the challenge to  $\mathcal{A}$  which is addressed to challenge  $\mathcal{C}$ .

$\mathcal{A}$  outputs a bit  $b'$ .

Owing to the definition of  $\mathcal{A}, b' = b$  with the advantages  $\text{Adv}_{\text{DVDAPS}, \mathcal{A}}^{\text{PSI-CMA}}, \mathcal{A}'$  distinguishes two signatures  $\sigma_0$  and  $\sigma_1$  with the advantages  $\text{Adv}_{\text{DVDAPS}, \mathcal{A}'}^{\text{IND-CCA}} = \text{Adv}_{\text{DVDAPS}, \mathcal{A}}^{\text{PSI-CMA}}$ .

Thus, the proposed DVDAPS scheme provides signer privacy. This theorem is proven.

**Theorem 5.** The proposed data-trading protocol has **completeness**.

*Proof.* Buyer  $B$  and seller  $S_i$  perform successfully in Phase 1. This means that  $B$  makes a signature  $\sigma$  without a double signature,  $S_i$  receives  $d$  bitcoins as their reward, and  $B$  will obtain his/her deposit  $d'$  bitcoins. Therefore, the proposed data-trading protocol has the property of **completeness** .

**Theorem 6.** The proposed data-trading protocol has **fairness**.

*Proof.* If there is an honest buyer  $B$  and a dishonest seller  $S_i, B$  receives its deposit; however,  $B$  does not receive data before deadline time  $t$ . In this case,  $S_i$  cannot achieve reward  $d$  bitcoins. If  $S_i$  wants to obtain the deposit  $d'$  bitcoins, he must run **phase 2** with  $A$  successfully. Because  $B$  can receive the signature from  $A$ , a contradiction exists. Thus, the probability of success for the dishonest seller  $S_i$  is negligible. If there is a dishonest buyer  $B$  and an honest seller  $S_i, B$  obtains a true verification result and returns his deposit if  $B$  cheats successfully and  $S_i$  cannot obtain anything. Suppose that  $S_i$  receives nothing, the buyer and seller must come to an agreement phase. If the entire computation cannot be completed before the deadline,  $S_i$  cannot receive the signature or obtain double signatures on the colliding messages. Otherwise,  $S_i$  runs **Phase 3** and obtains a deposit with the help of  $A$ . This contradicts the assumption. Thus, the probability of success that  $S_i$  is negligible after completing the entire computation is negligible.

**Theorem 7.** The proposed data-trading protocol has **accountability**.

*Proof.* Assume that **Phase 2** is executed, and before the deadline time,  $S_i$  does not obtain the signature from  $A$ . Therefore,  $A$  conspires with  $B$ , or the computation task is not completed before the deadline. In another case,  $B$  sends a signature to  $A$ , whereas  $A$  does not send it to  $S_i$ . Thus, we must perform **Phase 3**. Finally,  $S_i$  obtains the signature of  $A$  on transaction  $T_D$  and can obtain the deposit of  $B$ , which cannot be obtained.

It is assumed that **Phase 3** is executed. If  $A$  conspires with  $S_i$ , it implies that  $B$  has completed the entire work. In this case, it is expedient to perform Phase 2. Subsequently,  $S_i$  obtains the signature of  $A$  on transaction  $T_D$ , and obtains an *abort* token.  $S_i$  transfers transactions  $T_{get}$  to the blockchain network and they cannot obtain the deposit; however, they can find the cheating actions of  $A$ . This is because  $A$  performs both **Phase 2** and **Phase 3**.

**Theorem 8.** The proposed data-trading protocol has a signer **privacy**.

*Proof.* Theorem 8 can be proven using Theorem 4 because the protocol is based on the construction of DVDAPS.

Then we compare the safety goals with some representative data trading schemes. The comparison results are depicted in the following table B1.

**Table B1** The comparison results of schemes

Scheme	Auditable	Post quantum secure	No the third party	Fairness	Privacy of data
Zhao et al. [4]	✓	×	✓	×	×
Gao et al. [5]	✓	×	✓	✓	×
Delgado et al. [6]	✓	×	✓	×	×
Karame et al. [7]	×	×	✓	×	×
Our construction	✓	✓	✓	✓	✓

## Appendix C Computational cost and performance evaluation

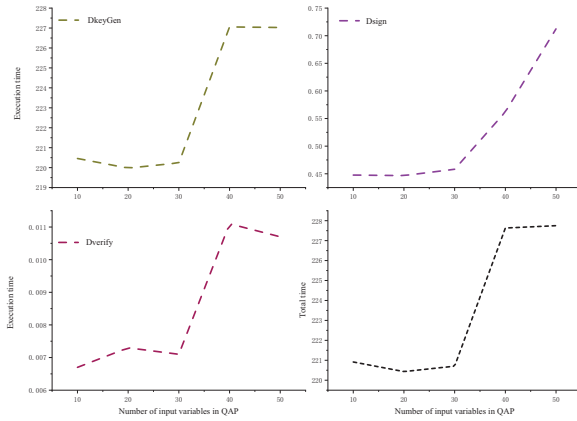
Suppose that  $ct$  is the length of the ciphertext of the LWE-based encryption scheme with the security parameter  $\lambda$ . The length of signature  $\sigma$  is  $3|ct| + q$ . For real applications, we select the size of the plaintext  $q = 2^{32} - 5$ , statistical distance

parameter  $\kappa = 32$ , and dimension of the lattice  $n = 1470$  [2]. The bit length of the proposed construction was evaluated. The common reference string comprises  $m + 2d + 5$  LWE ciphertexts, public key  $pk$ , and  $q$ . The bit length of the common reference string (CRS) is  $\lambda + (2^{22} + \#pk + 1) \times \log q$  and the proof size is  $(3n + 4)\log q$ . We provide the bit length in Table C1 using the security parameters recommended in [3], where  $|CRS|$  represents the bit size of the common reference string.

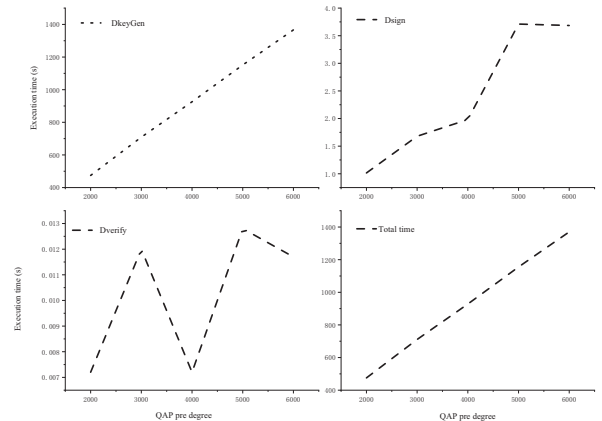
**Table C1** Parameter comparison

Scheme	Circuit	Post quantum	Proof size	$ CRS $	Assumption
Our construction	QAP	✓	405.72KB	15.5M	Linear-only

We used an Intel i7-7700 CPU @ 3.60GHz and an Ubuntu16 Linux operation system. The extraction algorithm consists only of a point-multiplication operation in which the running time can be omitted. Figure 2 presents the experimental results for the construction with different QAP numbers for the input variables. The number of input variables in the QAP is 10, 20, 30, 40, and 50, and the corresponding degrees of the QAP are 1024, 1024, 1032, 1032, and 1056. Figure 3 shows the results obtained from the preliminary analysis of the construction for different degrees in QAP, where the QAP degrees are 2048, 3072, 4096, 5120, and 6140, and the corresponding QAP numbers of variables are 2002, 3002, 4002, 5002, and 6002. Figures 2 and 3 show that, the total running time of DVDAPS was approximately 228 s and 1400 s, respectively. If we precompute the DkeyGen algorithm, the running time of the proposed algorithm is less than 1 s (Figure 2) and less than 4 s (Figure 3). Some errors may occur because the execution time of the third graph in Figure C1-C2 is considerably small. In the Dextract algorithm, only a point-multiplication operation exists, and we omit its execution time.



**Figure C1** Simulation results of DkeyGen, Dsign and Dver in DVDAPS.



**Figure C2** Simulation results of DkeyGen, Dsign and Dver in DVDAPS.

We measured the running time of steps 3 and 4 in the data-trading protocol using MATLAB R2019a software on a desktop with 8 core Intel (R) Core (TM) i7-9700 CPU and 16G RAM running on Windows 10. The machine-learning algorithm used in this study was a backpropagation (BP) neural network. We normalised the features of the data to the interval [0,1] before network construction. The BP network was then constructed using the newff function. We constructed a four-layer network (an input, output, and two hidden layers). The number of nodes in the input layer is the same as the number of features in each label. We used tansigoids as the transfer function for the hidden layer. All neurons in the hidden layers were set to 30. The output layer adopts a linear-transfer function with 0 or 1 nodes. The heart disease dataset was used to train the BP network. The dataset included 13 features, such as maximum heart rate, resting ECG results, fasting blood glucose, and serum cholesterol. We adopted a ten-fold cross-validation method. After training, we verified the norm of the difference between the prediction and test results. The smaller the norm, the more accurate the prediction is. If the norm is zero, the prediction is accurate. Figure C3 presents an overview of the running times of Steps 3 and 4. From the figure, it is evident that the testing time is in training model which is only 0.07 s.

Generally, these results indicate that the proposed data-trading protocol is highly efficient, specifically when we precompute the DkeyGen algorithm.

**References**

1 Wang Q, Wang D, Cheng C, et al. Quantum2FA: Efficient Quantum-Resistant Two-Factor Authentication Scheme for Mobile Devices. *IEEE Transactions on Dependable and Secure Computing*, 2021.

```

Time taken to build model: 0.07 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances      258           85.1485 %
Incorrectly Classified Instances    45           14.8515 %
Kappa statistic                    0.6997
Mean absolute error                 0.2108
Root mean squared error             0.3574
Relative absolute error             42.4975 %
Root relative squared error        71.7609 %
Total Number of Instances          303

=== Detailed Accuracy By Class ===

                TP Rate  FP Rate  Precision  Recall  F-Measure  MCC      ROC Area  PRC Area  Class
                0.879   0.181   0.853     0.879   0.866     0.700   0.889   0.892   A1
                0.819   0.121   0.850     0.819   0.834     0.700   0.889   0.876   A0
Weighted Avg.   0.851   0.154   0.851     0.851   0.851     0.700   0.889   0.885

=== Confusion Matrix ===

  a  b  <-- classified as
145 20 |  a = A1
 25 113 | b = A0

```

**Figure C3** Simulation results of machine learning in the data trading protocol

- 2 Chillotti I, Gama N, Georgieva M, et al. Faster fully homomorphic encryption: Bootstrapping in less than 0.1 seconds. In: *Advances in theory and application of cryptology and information security*, Berlin: Springer, 2016. 3-33.
- 3 Gennaro R, Minelli M, Nitulescu A, et al. Lattice-based zk-SNARKs from square span programs. In: *Advances in CCS*. ACM, 2018. 556-573.
- 4 Zhao Y Q, Yu Y , Li Y N, Han G, Du X J. Machine learning based privacy-preserving fair data trading in big data market. *Information Science*, 2019, 478: 449-460.
- 5 Gao J, Wu T, Li X. Secure, fair and instant data trading scheme based on bitcoin. *Journal of Information Security and Applications*, 2020, 53: 102511-102516.
- 6 Delgado Segura S , Pérez SolàC , Navarro Arribas G , Herrera JoancomartíJ . A fair protocol for data trading based on bitcoin transactions. *Future Generat Comput Syst*, 2017, 34(7):1.
- 7 Karame G O, Androulaki E, Capkun S. Double-spending fast payments in bitcoin. In: *Advances in CCS*, ACM, 2012: 906-917.