

# Solution of iDASH 2024 Track 1

SCB @ Chongqing<sup>\*†</sup>

September 14, 2024

## 1 Workflow

DASHFormer follows a workflow similar to Transformer, consisting of an input layer, encoding layer, attention layer, feedforward layer, pooling layer, and output layer. In this document, all vectors are in row.

### 1.1 Input Layer

- For the Track 1 of iDash 2024, the input contains 1,000 protein sequences, each of length 50, with the alphabet taken from a set of 25 letters  $\{A, B, \dots, Y\}$ . Namely, each sequence is a vector in  $\{A, B, \dots, Y\}^{50}$ .
- But for DASHFormer, each input is a vector  $\mathbf{a} \in \{A, B, \dots, Y\}^{50}$ . However, before starting DASHFormer, we should first tokenize each sequences.
- So the real input for DASHFormer is  $\mathbf{x} \in \mathbb{Z}^{50}$ , which is tokenized by the file named ‘dashformer\_tokenizer.json’.
- Using one-hot encoding, one can get  $\mathbf{X} \in \{0, 1\}^{50 \times 25}$ .

### 1.2 Encodign Layer

#### 1.2.1 Embedding

- Embedding each token into a 128-dimensional vector. Since the alphabet contains only 25 letters, the size of this embedding matrix is  $25 \times 128$ . This matrix can be found in the file named ‘embedding\_Embedding\_weights.txt’.
- After embedding, each  $\mathbf{X} \in \{0, 1\}^{50 \times 25}$  will be converted into a matrix  $\mathbf{X} = \mathbf{X} \cdot \mathbf{W}_e \in \mathbb{R}^{50 \times 128}$ , where  $\mathbf{W}_e \in \mathbb{R}^{25 \times 128}$  is the matrix of embedding weights.

---

<sup>\*</sup>CIGIT, CAS: Jingwei Chen (chenjingwei@cigit.ac.cn), Linhan Yang, Chen Yang, Shuai Wang, Rui Li, Weijie Miao, and Wenyuan Wu (wuwenyuan@cigit.ac.cn).

<sup>†</sup>Sansure Biotech Inc.: Li Yang and Jia Liu.

### 1.2.2 Positional Encoding

- Each token in each sequence has its own position, which can be encoded as a vector of dimension 128. The encoding method is decided by the file named 'positional\_encoding\_Lookup.txt'.
- The positional encoding is  $P \in \mathbb{R}^{50 \times 128}$ .

### 1.2.3 Output of the Encoding Layer

- The output of the encoding layer is  $X := X + P \in \mathbb{R}^{50 \times 128}$ .
- Save  $Y := X$  for future use (Add and LayerNorm).

## 1.3 Attention Layer

### 1.3.1 Attention

- The number of heads:  $H = 4$ . Hence, in what follows,  $h = 0, 1, 2, 3$ .
- Head weights and bias:  $W_Q^{(h)}, W_K^{(h)}, W_V^{(h)} \in \mathbb{R}^{128 \times 32}$  and  $b_Q^{(h)}, b_K^{(h)}, b_V^{(h)} \in \mathbb{R}^{32}$ .
- Queries:  $Q^{(h)} = XW_Q^{(h)} + \mathbf{1}^T \cdot b_Q^{(h)} \in \mathbb{R}^{50 \times 32}$ . Alternatively,  $Q'^{(h)} = X' \cdot W_Q'^{(h)}$  with  $X' = (X, \mathbf{1}^T) \begin{pmatrix} W_Q^{(h)} \\ b_Q^{(h)} \end{pmatrix}$ .
- Keys:  $K^{(h)} = XW_K^{(h)} + \mathbf{1}^T \cdot b_K^{(h)} \in \mathbb{R}^{50 \times 32}$ .
- Values:  $V^{(h)} = XW_V^{(h)} + \mathbf{1}^T \cdot b_V^{(h)} \in \mathbb{R}^{50 \times 32}$ .
- Attention:  $X^{(h)} = \text{softmax}(Q^{(h)} \cdot K^{(h)T} / \sqrt{d}) V^{(h)}$ , where  $d = 32$ , i.e., the number of columns of  $Q$ ,  $K$  and  $V$ .
- Concatenate:  $X = (X^{(1)}, X^{(2)}, X^{(3)}, X^{(4)}) \in \mathbb{R}^{50 \times 128}$ .

### 1.3.2 Combining Head (A Linear Layer)

- Head combining weight and bias:  $W_c \in \mathbb{R}^{128 \times 128}$  and  $b_c \in \mathbb{R}^{128}$ .
- Combining:  $X = XW_c + \mathbf{1}^T \cdot b_c \in \mathbb{R}^{50 \times 128}$ .

### 1.3.3 Add and LayerNorm-1

- Add:  $X := X + Y$ .
- LayerNorm-1:  $X := \text{LayerNorm}_1(X)$ , where the parameters can be found in the file named 'transformer\_block\_LayerNorm1\_weights.txt'.
- Update  $Y := X$ .

## 1.4 Feed Forward with ReLU

- $\mathbf{X} := \mathbf{X}\mathbf{W}_1 + \mathbf{1}^T \cdot \mathbf{b}_1$  with  $\mathbf{W}_1 \in \mathbb{R}^{128 \times 256}$  and  $\mathbf{b}_1 \in \mathbb{R}^{256}$ .
- $\mathbf{X} := \text{ReLU}(\mathbf{X})$ .
- $\mathbf{X} := \mathbf{X}\mathbf{W}_2 + \mathbf{1}^T \cdot \mathbf{b}_2$  with  $\mathbf{W}_2 \in \mathbb{R}^{256 \times 128}$  and  $\mathbf{b}_2 \in \mathbb{R}^{128}$ , where the parameters can be found in the file named ‘transformer\_block\_FFN\_weights.txt’.

### 1.4.1 Add and LayerNorm-2

- Add:  $\mathbf{X} := \mathbf{X} + \mathbf{Y}$ .
- LayerNorm-2:  $\mathbf{X} := \text{LayerNorm}_2(\mathbf{X})$ , where the parameters can be found in the file named ‘transformer\_block\_LayerNorm2\_weights.txt’.

## 1.5 Average Pooling

- Given  $\mathbf{X} \in \mathbb{R}^{50 \times 128}$ , output  $\mathbf{x} = \mathbf{1} \cdot \mathbf{X} \in \mathbb{R}^{128}$ .

## 1.6 Dense Classification

- Taking  $\mathbf{x} \in \mathbb{R}^{128}$  as input, output  $\mathbf{y} = \mathbf{x}\mathbf{W}_d + \mathbf{b}_d$ , where  $\mathbf{W}_d \in \mathbb{R}^{128 \times 25}$  and  $\mathbf{b}_d \in \mathbb{R}^{25}$ .

# 2 Implementation

We strictly follow the workflow introduced in the previous section for cipher-text computation. Specifically, non-polynomial functions are handled using approximation methods. At the same time, we fully leverage the SIMD support provided by the CKKS scheme [2] to adopt flexible data packing methods. This allows us to use tensor-matrix computations to complete the prediction of 100 samples simultaneously, and successfully applies the Baby-Step Giant-Step (BSGS) method [3] to accelerate the matrix-matrix multiplications.

## 2.1 Approximate Non-polynomial Functions

### 2.1.1 Softmax

The Softmax of a vector  $\mathbf{x} = (x_i)_i$  is defined as

$$\text{softmax}(\mathbf{x}) = \left( \frac{\exp(x_i)}{\sum_i \exp(x_i)} \right)_i.$$

Instead of high-degree polynomial approximations of the exponential function in softmax, we use the square softmax (sqmax) in the following form:

$$\text{sqmax}(\mathbf{x}) = \left( \frac{(x_i + c)^2}{\sum_i (x_i + c)^2} \right)_i.$$

Furthermore, we use a constant as the denominator to avoid the ciphertext division, i.e.,

$$\text{sqmax}(\mathbf{x}) = \left( \frac{(x_i + c)^2}{\delta} \right)_i.$$

Since there are four heads in DASHformer, we choose different  $c$  and  $\delta$  for each head. These coefficients are stored in “`config/config.go`”.

### 2.1.2 ReLU

For  $x \in \mathbb{R}$ ,  $\text{ReLU}(x) = \max(0, x)$ . Currently, we use a polynomial of degree 6 to approximate the ReLU function:

$$f(x) = a_6x^6 + a_5x^5 + a_4x^4 + a_3x^3 + a_2x^2 + a_1x + a_0,$$

where the coefficients are stored in “`config/config.go`”.

### 2.1.3 LayerNorm

Given a vector  $\mathbf{x} = (x_i)_i$ ,

$$\text{LayerNorm}(\mathbf{x}) = \left( \gamma \cdot \frac{x_i - \mu}{\sqrt{\sigma^2}} + \beta \right)_i$$

with  $\mu$  and  $\sigma^2$  the mean and variance of  $\mathbf{x}$ . The LayerNorm function applies to all the rows of the matrix  $\mathbf{X} \in \mathbb{R}^{50 \times 128}$ .

Instead of computing  $\frac{1}{\sigma^2}$  in ciphertext, we compute it in plaintext with all 1000 samples to get a 50-dimensional vector for the inverse of the standard deviation so that

$$\text{LN}(\mathbf{x}) = (\gamma \cdot \tilde{\sigma} \cdot (x_i - \mu) + \beta)_i,$$

where  $\tilde{\sigma}$  is the learned standard deviation from the given 1000 samples. If we use the matrix form, we have

$$\begin{aligned} \text{LN}(\mathbf{X}) &= \text{diag}\left(\frac{1}{\sqrt{\sigma_i^2}}\right)(\mathbf{X} - \mathbf{U})\text{diag}(\gamma_i) + \mathbf{1}^T \cdot \boldsymbol{\beta} \\ &= \text{diag}\left(\frac{1}{d\sqrt{\sigma_i^2}}\right)(d\mathbf{X} - \mathbf{X} \cdot \mathbf{E})\text{diag}(\gamma_i) + \mathbf{1}^T \cdot \boldsymbol{\beta} \\ &=: \boldsymbol{\Sigma}^{-1} \cdot \mathbf{X} \cdot \boldsymbol{\Gamma} + \mathbf{1}^T \cdot \boldsymbol{\beta}. \end{aligned}$$

Here  $d$  is the number of columns of  $\mathbf{X}$ ,  $\boldsymbol{\Sigma}^{-1} = \text{diag}\left(\frac{1}{d\sqrt{\sigma_i^2}}\right)$  and

$$\boldsymbol{\Gamma} = (d\mathbf{I} - \mathbf{E}) \cdot \text{diag}(\gamma_i),$$

where  $\mathbf{I}$  is the identity matrix and all entries of  $\mathbf{E}$  are 1. The parameters  $\boldsymbol{\Sigma}^{-1}$ ,  $\boldsymbol{\Gamma}$  and  $\boldsymbol{\beta}$  are stored in “`data/layerNorm1_Reciprocal_SqrtVariance.txt`” and “`data/layerNorm2_Reciprocal_SqrtVariance.txt`”.

## 2.2 Multiplicative Depth

Let  $\mathbf{X}_0 \in \{0, 1\}^{50 \times 25}$  be the one-hot encoding for an input  $\mathbf{x} \in \mathbb{Z}^{50}$ . Then

$$\begin{aligned}
& \text{DASHFormer}(\mathbf{X}) \\
&= \mathbf{1} \cdot \mathbf{X} \cdot \mathbf{W}_d + \mathbf{b}_d \\
&= \mathbf{1} \cdot \text{LN}_2(\mathbf{X} + \mathbf{Y}_2) \cdot \mathbf{W}_d + \mathbf{b}_d \\
&= \mathbf{1} \cdot (\Sigma_2^{-1}(\mathbf{X} + \mathbf{Y}_2)\Gamma_2 + \mathbf{1}^T \beta_2) \cdot \mathbf{W}_d + \mathbf{b}_d \\
&= \mathbf{1} \cdot \Sigma_2^{-1} \mathbf{X} \Gamma_2 \cdot \mathbf{W}_d + \mathbf{1} \cdot \Sigma_2^{-1} \mathbf{Y}_2 \Gamma_2 \cdot \mathbf{W}_d + \mathbf{1} \cdot \mathbf{1}^T \beta_2 \cdot \mathbf{W}_d + \mathbf{b}_d. \tag{1}
\end{aligned}$$

In Eq. (1), we have

$$\begin{aligned}
\mathbf{Y}_2 &= \text{LN}_1(\mathbf{X} + \mathbf{Y}_1) \\
&= \Sigma_1^{-1} \mathbf{X} \Gamma_1 + \Sigma_1^{-1} \mathbf{Y}_1 \Gamma_1 + \mathbf{1}^T \cdot \beta_1 \\
&= \Sigma_1^{-1} \mathbf{X} \Gamma_1 + \Sigma_1^{-1} \mathbf{X}_0 \mathbf{W}_e \Gamma_1 + \Sigma_1^{-1} \mathbf{P} \Gamma_1 + \mathbf{1}^T \cdot \beta_1 \tag{2}
\end{aligned}$$

where  $\mathbf{Y}_1 = \mathbf{X}_0 \mathbf{W}_e + \mathbf{P}$ . In Eq. (2),

$$\begin{aligned}
\Sigma_1^{-1} \mathbf{X} \Gamma_1 &= \Sigma_1^{-1} (\mathbf{X} \mathbf{W}_c + \mathbf{1}^T \cdot \mathbf{b}_c) \Gamma_1 \\
&= \Sigma_1^{-1} \mathbf{X} \mathbf{W}_c \Gamma_1 + \Sigma_1^{-1} \cdot \mathbf{1}^T \cdot \mathbf{b}_c \Gamma_1 \\
&= \Sigma_1^{-1} \cdot \text{sqmax}(\mathbf{Q} \mathbf{K}^T / \sqrt{d}) \mathbf{V} \mathbf{W}_c \Gamma_1 + \Sigma_1^{-1} \cdot \mathbf{1}^T \cdot \mathbf{b}_c \Gamma_1,
\end{aligned}$$

where

$$\mathbf{Q} = \mathbf{X} \mathbf{W}_Q + \mathbf{1}^T \cdot \mathbf{b}_Q = \mathbf{X}_0 \mathbf{W}_e \mathbf{W}_Q + \mathbf{P} \mathbf{W}_Q + \mathbf{1}^T \cdot \mathbf{b}_Q,$$

and similarly for  $\mathbf{K}$  and  $\mathbf{V}$ .

In Eq. (1),

$$\begin{aligned}
& \mathbf{1} \cdot \Sigma_2^{-1} \mathbf{X} \Gamma_2 \cdot \mathbf{W}_d \\
&= \mathbf{1} \cdot \Sigma_2^{-1} (\mathbf{X} \mathbf{W}_2 + \mathbf{1}^T \cdot \mathbf{b}_2) \Gamma_2 \cdot \mathbf{W}_d \\
&= \mathbf{1} \cdot \Sigma_2^{-1} \mathbf{X} \mathbf{W}_2 \cdot \Gamma_2 \cdot \mathbf{W}_d + \mathbf{1} \cdot \Sigma_2^{-1} \mathbf{1}^T \cdot \mathbf{b}_2 \cdot \Gamma_2 \cdot \mathbf{W}_d \tag{3}
\end{aligned}$$

In Eq. (3),

$$\mathbf{1} \cdot \Sigma_2^{-1} \mathbf{X} \mathbf{W}_2 \cdot \Gamma_2 \cdot \mathbf{W}_d = \mathbf{1} \cdot \Sigma_2^{-1} \text{ReLU}(\mathbf{X} \mathbf{W}_1 + \mathbf{1}^T \cdot \mathbf{b}_1) \mathbf{W}_2 \cdot \Gamma_2 \cdot \mathbf{W}_d \tag{4}$$

In Eq. (4),

$$\text{ReLU}(\mathbf{X} \mathbf{W}_1 + \mathbf{1}^T \cdot \mathbf{b}_1) = \text{ReLU}(\mathbf{Y}_2 \mathbf{W}_1 + \mathbf{1}^T \cdot \mathbf{b}_1)$$

with  $\mathbf{Y}_2$  as computed in Eq. (2). So

$$\mathbf{Y}_2 \mathbf{W}_1 = \Sigma_1^{-1} \mathbf{X} \Gamma_1 \mathbf{W}_1 + \Sigma_1^{-1} \mathbf{X}_0 \mathbf{W}_e \Gamma_1 \mathbf{W}_1 + \Sigma_1^{-1} \mathbf{P} \Gamma_1 \mathbf{W}_1 + \mathbf{1}^T \cdot \beta_1 \mathbf{W}_1.$$

We now can give the multiplicative depth cost of our homomorphic evaluation of the DASHFormer as in Table 1. It shows that one can evaluate the whole circuit with at most 10 multiplicative depths (including both ciphertext-ciphertext and plaintext-ciphertext multiplication), which allows us to choose  $\log N = 14$  to achieve 128-bit security [1].

Table 1: Multiplicative depth for the whole DASHFormer

Task	Before ReLU	ReLU	After ReLU	Total
Depth	6	3	1	10

### 2.3 Numeric Error-Control

To ensure that the ciphertext computation results (at the lowest level) still can be decrypted correctly, we scale the coefficients used during the last step in Section 1.6. Namely, we multiply a constant  $c = 1/2649.372705$  with  $\mathbf{W}_d$  and  $\mathbf{b}_d$ , respectively. Correspondingly, to make the decrypted results consistent with the plaintext results, we also multiplied the decrypted results by  $1/c$ .

### 2.4 Encrypted Tensor Product

For each sample  $\mathbf{x}$  to be classified, the inference of DASHFormer is as stated in the previous subsection. However, we have to deal with 100 samples during the evaluation phase of Track 1. Instead of a  $50 \times 25$  matrix  $\mathbf{X}_0$  as input, we have to compute with a tensor  $\mathbf{X}_0 \in \mathbb{R}^{100 \times 50 \times 25}$ . When packing and encrypting data for the CKKS scheme, we treat each  $100 \times 50$  slice of the tensor as a ciphertext. This makes us choose  $\log N = 14$  for CKKS without using the conjugate invariant, i.e., the number of slots is  $8192 > 5000$ . In this setting, our solution can evaluate at most 163 samples.

### 2.5 Encrypted Matrix Multiplication

The most time-consuming part of the entire process is to compute  $\text{sqmax}(\mathbf{Q}\mathbf{K}^T)\mathbf{V}$ . Since  $\text{sqmax}$  is a second-order approximation of softmax, this step essentially involves the sequential multiplication of three ciphertext matrices. We encode the matrices column by column. Thus, after  $\mathbf{C} = \mathbf{Q}\mathbf{K}^T$  is computed, the diagonal encoding of  $\mathbf{C}$  is obtained. Then, when computing  $\mathbf{C} \cdot \mathbf{V}$ , the computation is performed by  $\mathbf{C}\mathbf{v}_i$  using the Halevi-Shoup algorithm [3], where  $\mathbf{v}_i$  is the  $i$ th column of  $\mathbf{V}$ . Ultimately, the ciphertexts of the columns of the resulting matrix  $\mathbf{R} = \text{sqmax}(\mathbf{Q}\mathbf{K}^T)\mathbf{V}$  are obtained, allowing for subsequent computations.

**Double-BSGS.** Through some non-trivial packing techniques, we can fully utilize the so-called baby-step-giant-step (BSGS) technique *twice* during computation of  $\mathbf{R}$  to reduce the number of required rotations. We call this technique *double-BSGS*. Assuming  $\mathbf{Q}$ ,  $\mathbf{K}$ , and  $\mathbf{V}$  are all  $n \times m$  matrices, our method reduces the required number of ciphertext rotations from  $m^2 \cdot n$  to  $m(m+1)(\ell+k)$  with  $n \leq \ell \cdot k$ . In our implementation of DASHFormer,  $(n, m, p, k, \ell) = (50, 128, 25, 8, 7)$ . Particularly, we reduce the required number of rotations from 819,200 to 247,680, saving 571,520 (about 70%) rotations.

### 3 Performance

We use the CKKS scheme in Lattigo (v5.0.2) [4] to implement our solution. The parameters are chosen as the following:

- $\log N = 14$ , where  $N$  is the polynomial degree of the ring,
- $\log Q = 38 + 33 \cdot 10 = 368$ , where  $Q$  is the largest ciphertext modulus,
- $\log P = 36 \cdot 2 = 72$ , where  $P$  is the special prime for key-switching,
- $\log \Delta = 33$ , where  $\Delta$  is the scaling factor.

According to the latest security guidelines for homomorphic encryption [1], the above setup can achieve 128-bit security.

In addition, we use *four* threads to parallelize our implementation. On a node (Intel Xeon Gold 6248R 3.00 GHz and 128 GB memory) of a cluster with Centos 8 system (Platform 1) and a PC (i9-12900K 3.2 GHz and 32 GB memory) with Ubuntu 22.04 (Platform 2), we tested our implementation. The maximum memory used during our execution is about 29 GB. The timings are obtained with limited 4 cores and listed in Table 2.

Table 2: Timings in second

Task	Before ReLU	ReLU	After ReLU	Total (with i/o)
Platform 1	314.4	10.5	5.8	333.9
Platform 2	209.4	6.2	3.9	222.9

We execute our implementation several times for randomly chosen 100 samples; the micro-AUC ranges from 0.91 to 0.97, as shown in Figure 1.

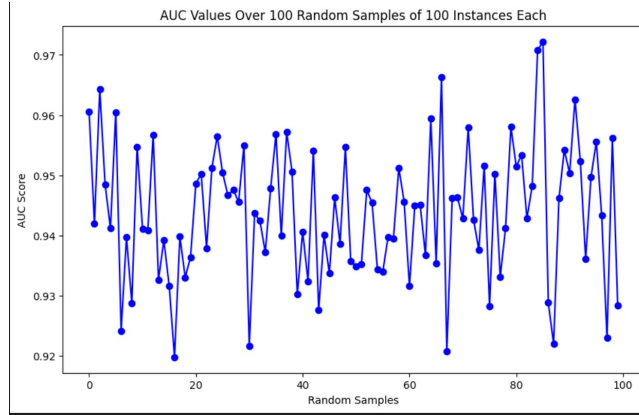


Figure 1: Micro-AUC for randomly chosen Samples

## References

- [1] M. Albrecht, M. Chase, H. Chen, J. Ding, S. Goldwasser, S. Gorbunov, S. Halevi, J. Hoffstein, K. Laine, K. Lauter, S. Lokam, D. Micciancio, D. Moody, T. Morrison, A. Sahai, and V. Vaikuntanathan. Homomorphic encryption security standard. Technical report, HomomorphicEncryption.org, Toronto, Canada, November 2018. <https://eprint.iacr.org/2019/939.pdf>.
- [2] J. H. Cheon, A. Kim, M. Kim, and Y. Song. Homomorphic encryption for arithmetic of approximate numbers. In T. Takagi and T. Peyrin, editors, *Proceedings of ASIACRYPT 2017 – 23rd International Conference on the Theory and Applications of Cryptology and Information Security (December 3-7, 2017, Hong Kong, China), Part I*, volume 10624 of *Lecture Notes in Computer Science*, pages 409–437. Springer, Heidelberg, 2017. [https://doi.org/10.1007/978-3-319-70694-8\\_15](https://doi.org/10.1007/978-3-319-70694-8_15).
- [3] S. Halevi and V. Shoup. Algorithms in HELib. In J. A. Garay and R. Genaro, editors, *Advances in Cryptology – CRYPTO 2014 (Santa Barbara, USA, August 17-21, 2014)*, volume 8616 of *Lecture Notes in Computer Science*, pages 554–571. Springer, Heidelberg, 2014. [https://doi.org/10.1007/978-3-662-44371-2\\_31](https://doi.org/10.1007/978-3-662-44371-2_31).
- [4] Tune-Insight. Lattigo (v5.0.2), Accessed in July, 2024. <https://github.com/tuneinsight/lattigo/releases/tag/v5.0.2>.