

# 1. Embedding

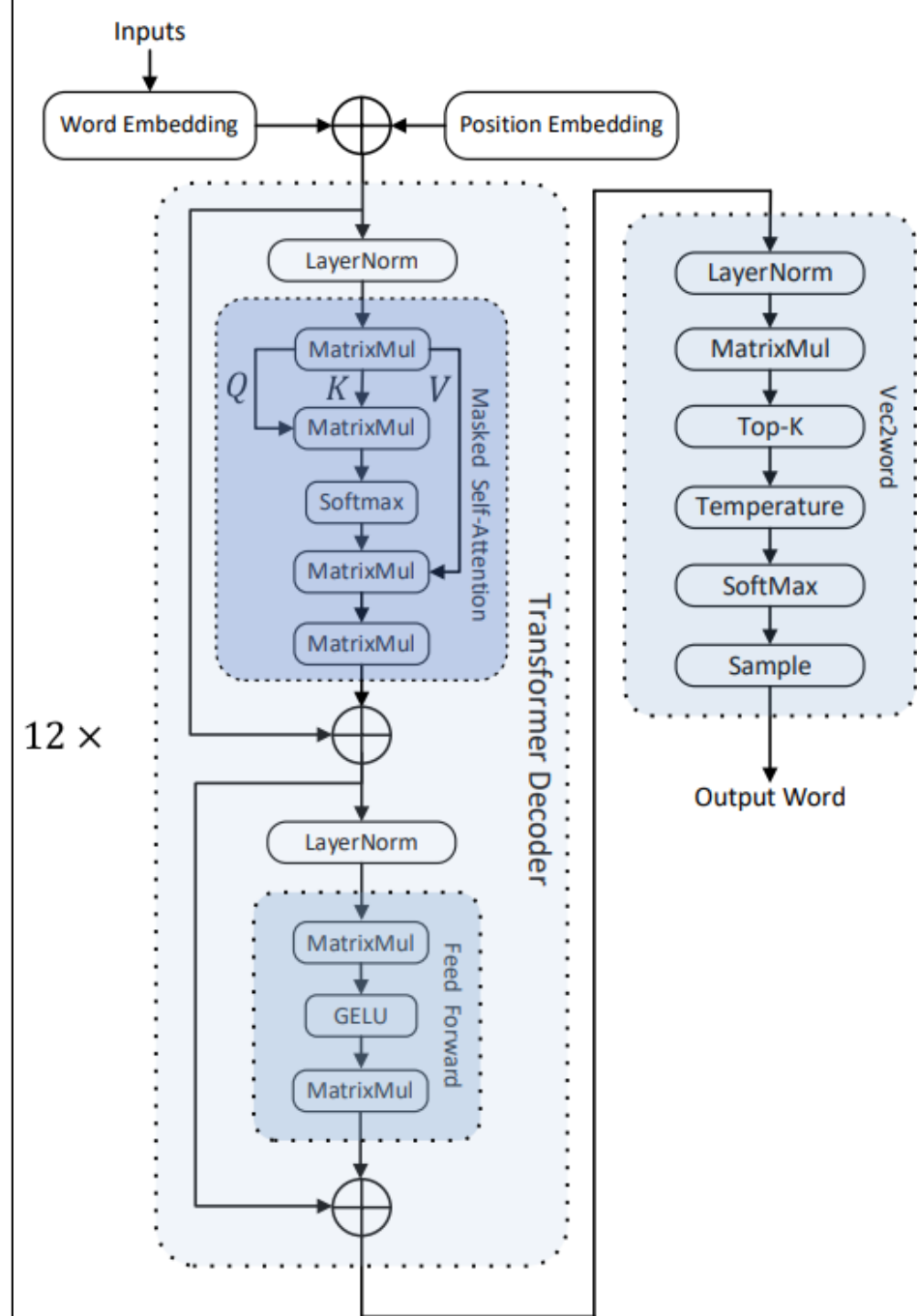
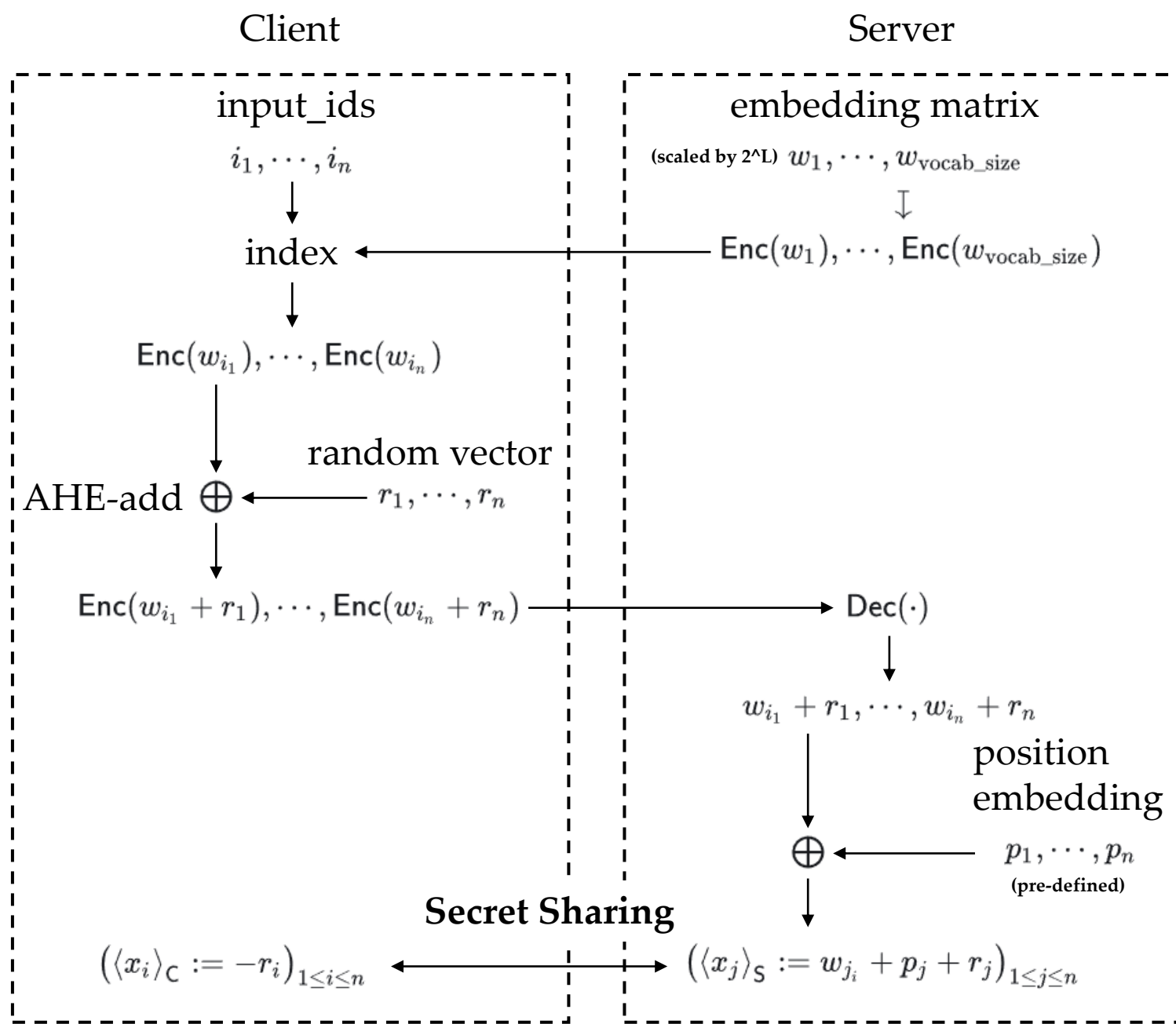


Figure 2: The architecture and workflow of GPT.

## 2. LayerNorm $x_i^{(j)}$ denotes the j-th component of vector $x_i$

Client

Server

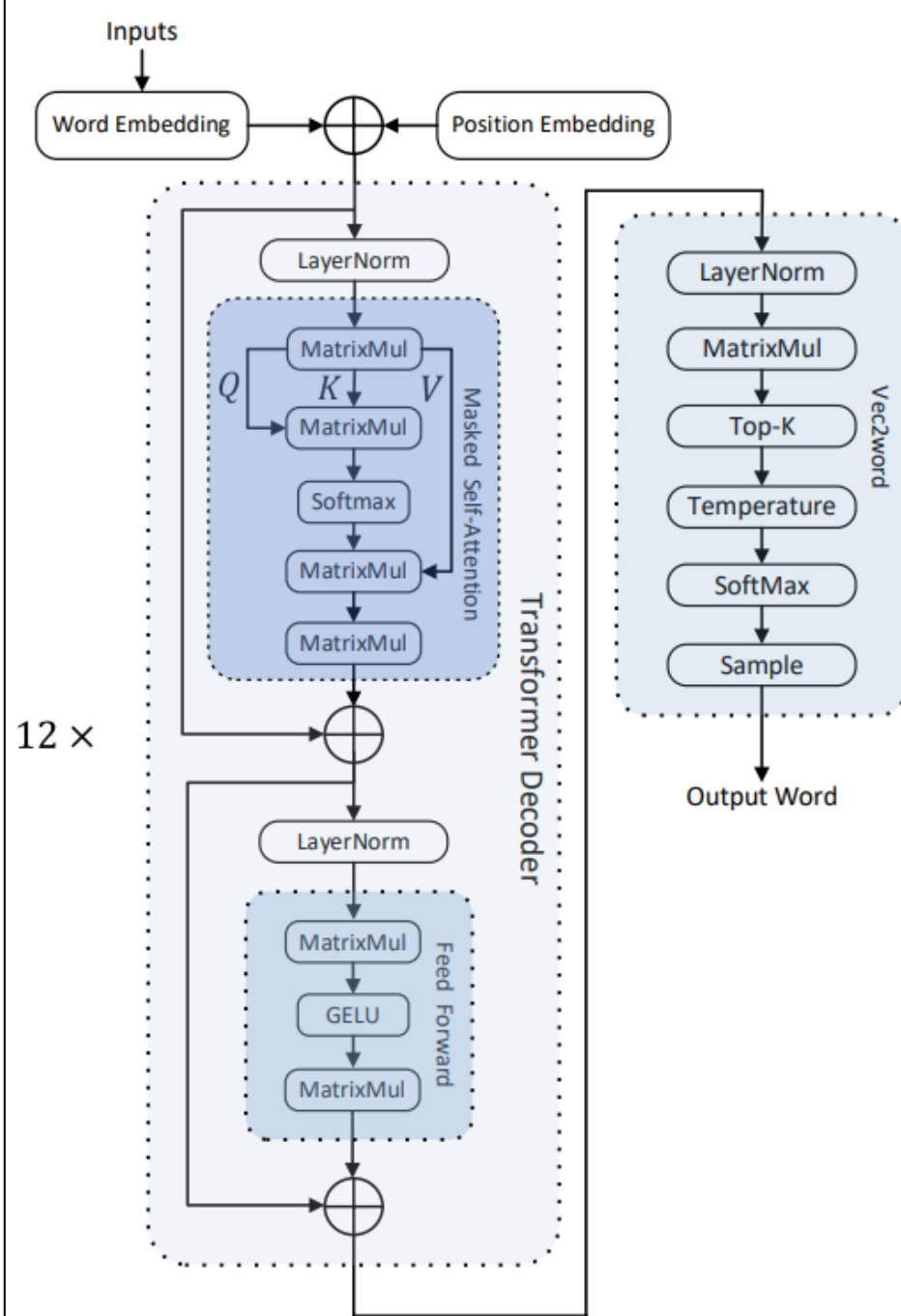
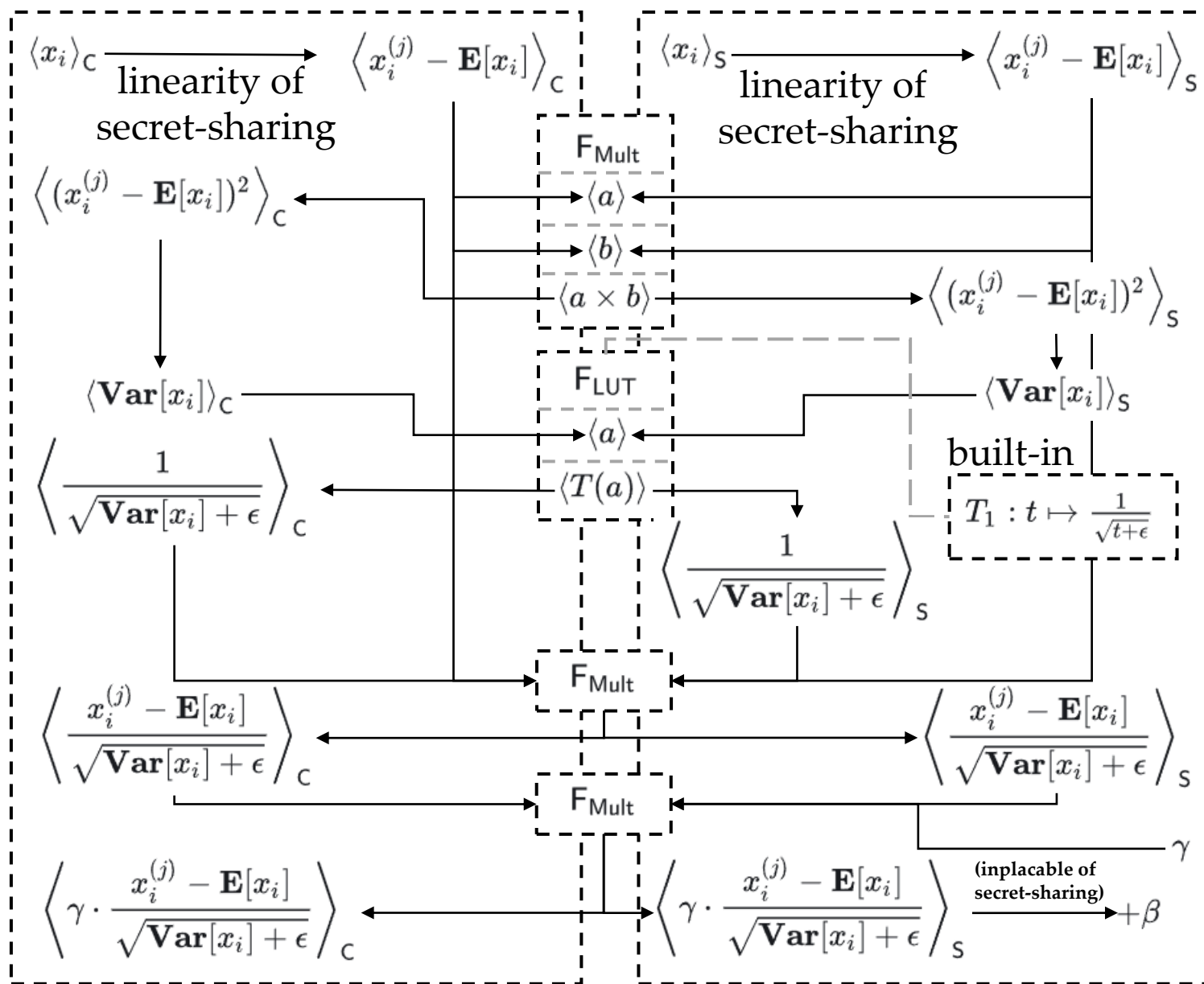
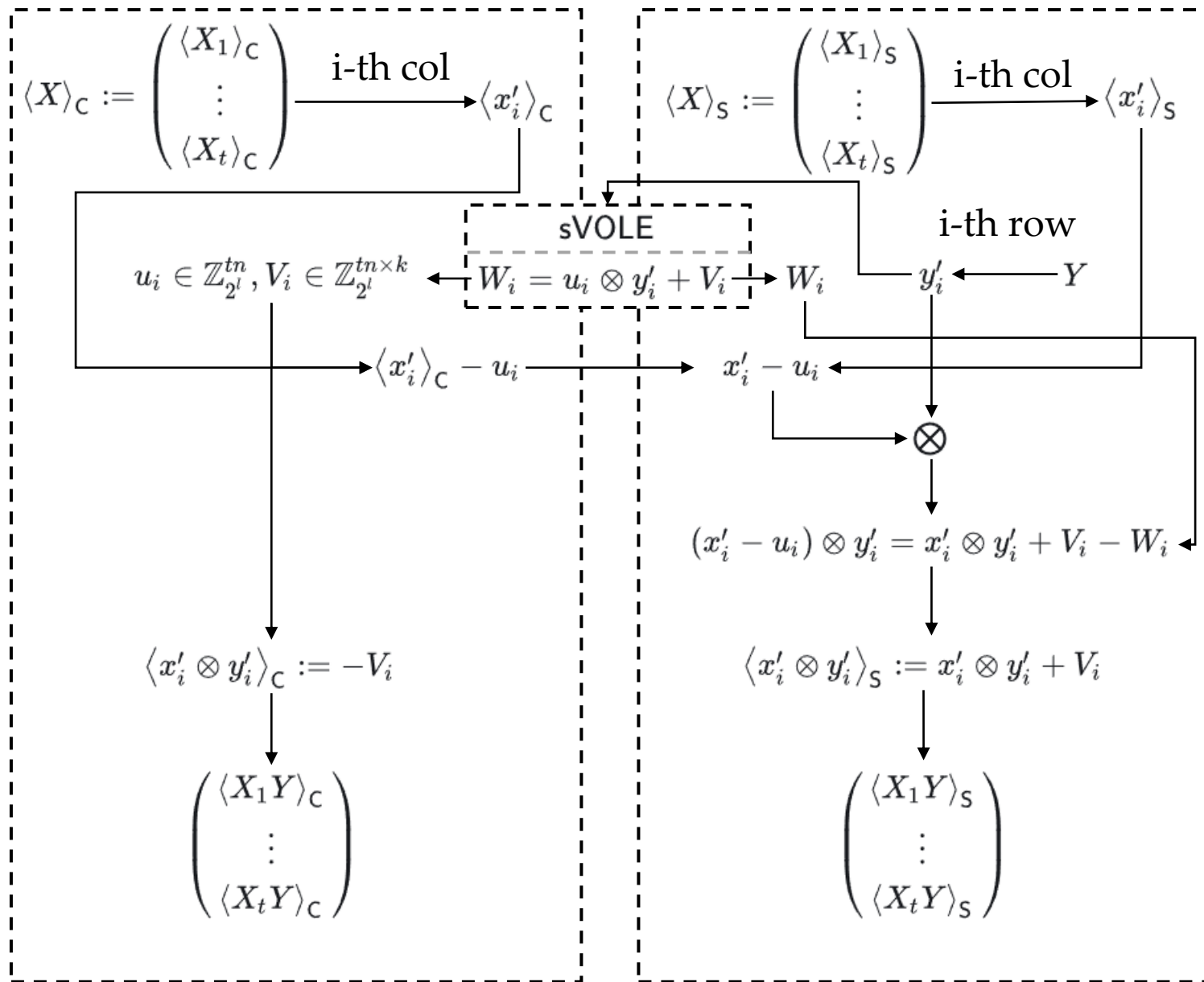


Figure 2: The architecture and workflow of GPT.

**0. MatrixMul**  $\otimes : \mathbb{Z}_{2^l}^n \times \mathbb{Z}_{2^l}^k \ni (x, y) \mapsto x \otimes y := (x^{(j)} \cdot y^{(r)})_{1 \leq j \leq n, 1 \leq r \leq k} \in \mathbb{Z}_{2^l}^{n \times k}$

Client

Server



C has  $X_1, \dots, X_t$  and S has  $Y$ .  
They want to get the secret-sharing of  $X_1 * Y, \dots, X_t * Y$ .

**Lemma :**  $X \in \mathbb{Z}_{2^l}^{n \times m}, Y \in \mathbb{Z}_{2^l}^{m \times k}$ , then :

$$X \cdot Y = \sum_{j=1}^m \underbrace{X_{\cdot, j}}_{j\text{-th col}} \otimes \underbrace{Y_{j, \cdot}}_{j\text{-th row}}$$

### 3. Self-Attention $X := (x_i^{(j)})_{1 \leq i \leq n, 1 \leq j \leq m}$

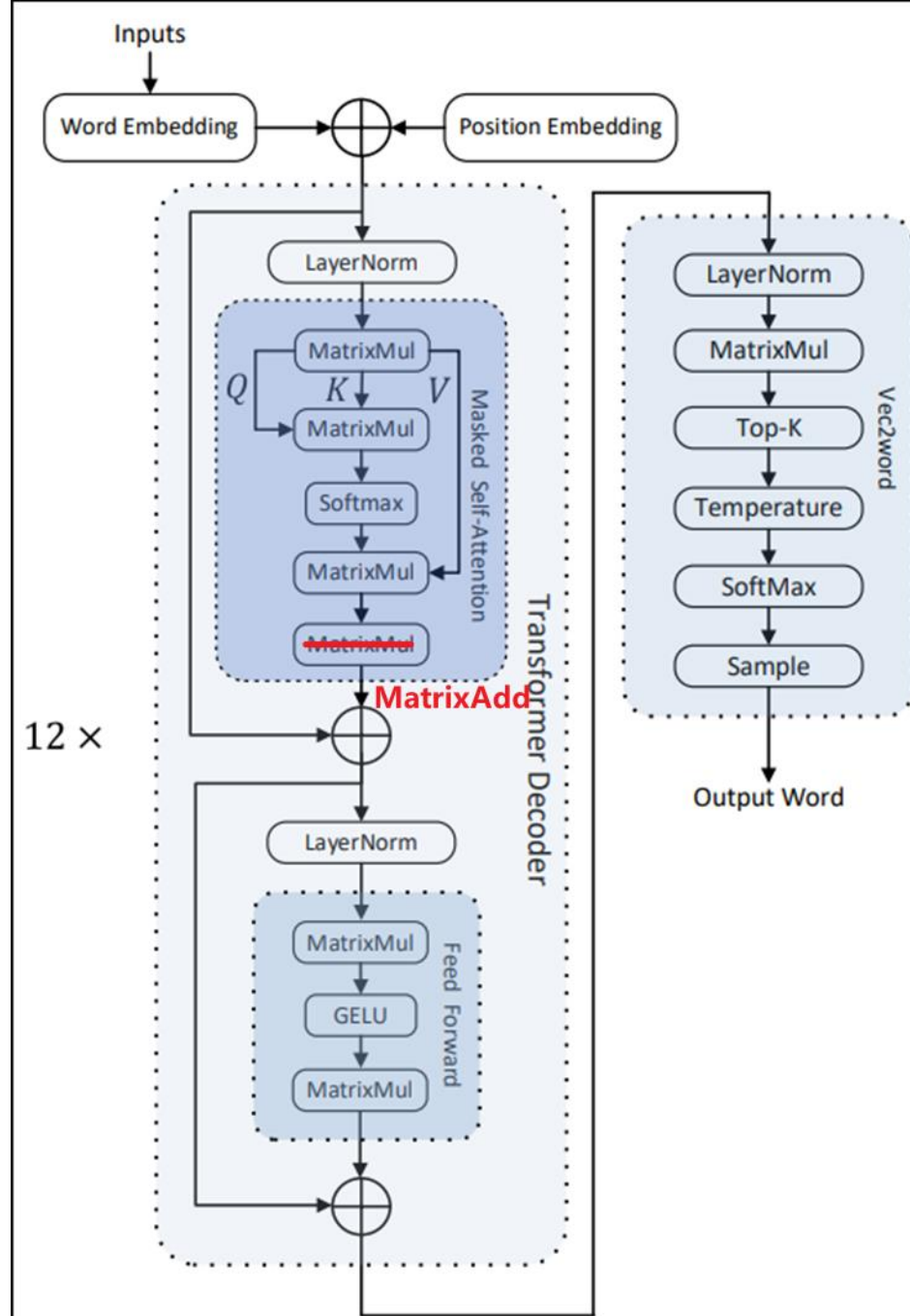
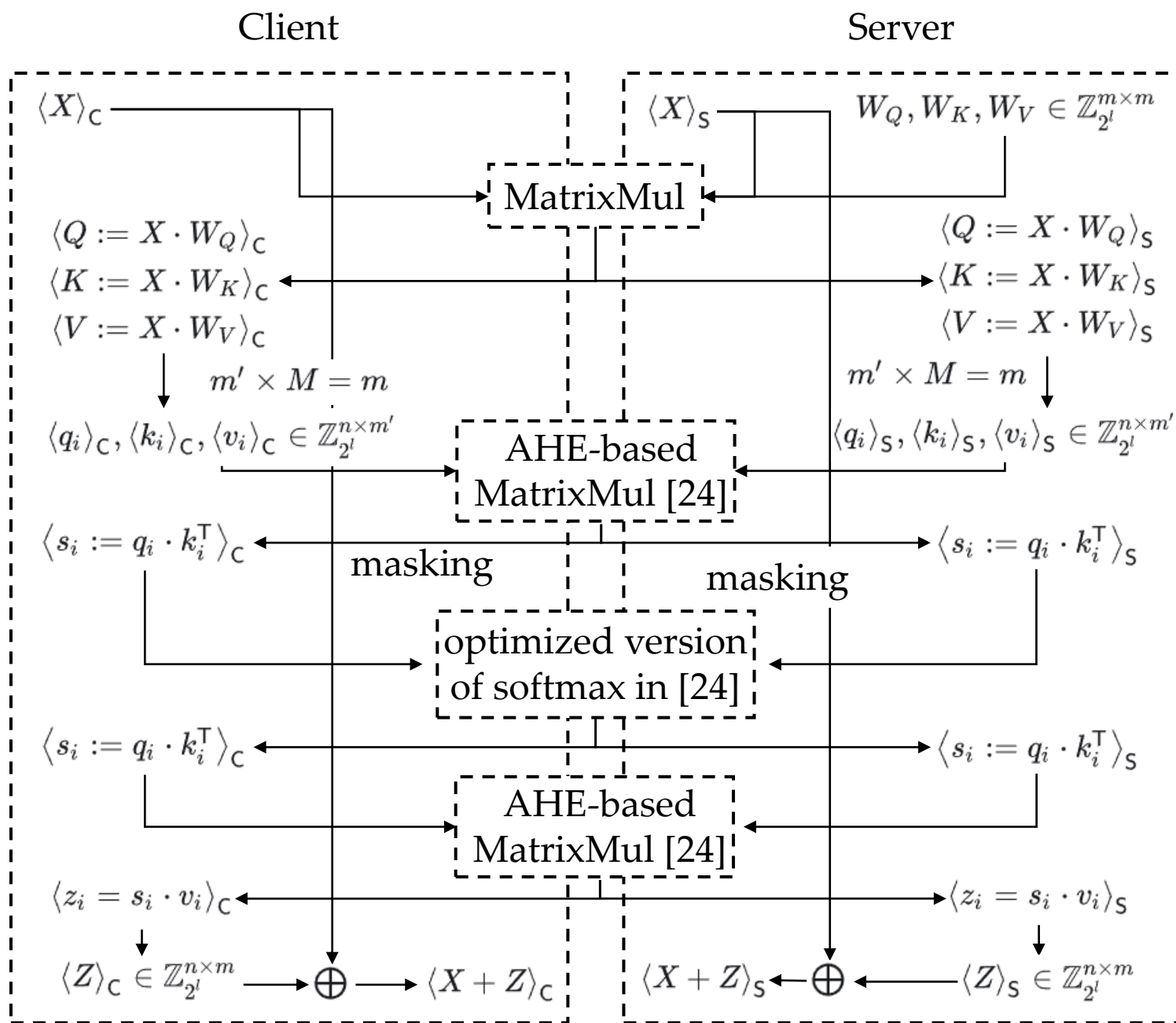


Figure 2: The architecture and workflow of GPT.

4. Feed Forward

Client

Server

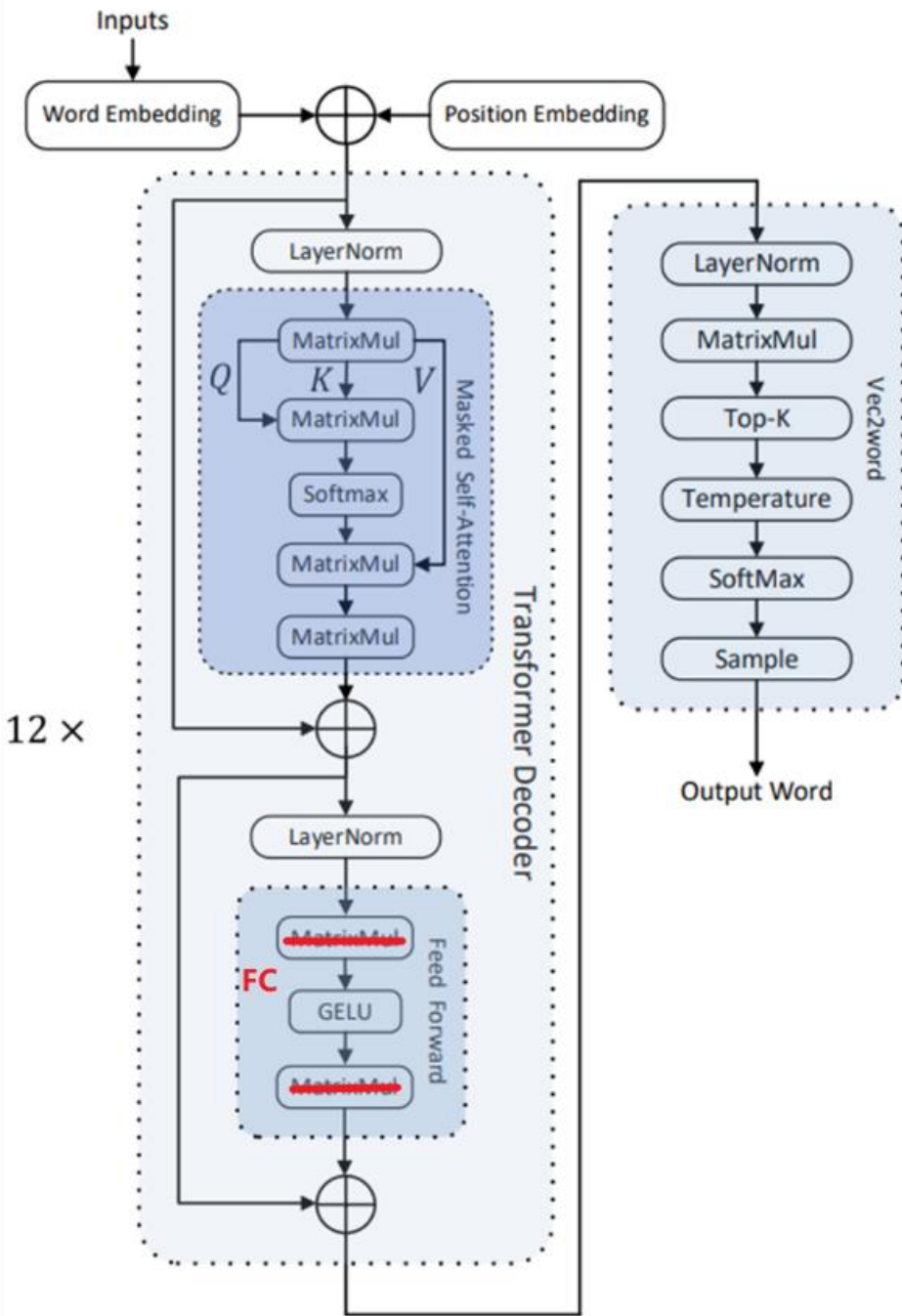


Figure 2: The architecture and workflow of GPT.

## 5. Vec2Word

$k := \text{vocab\_size}$

Client

Server

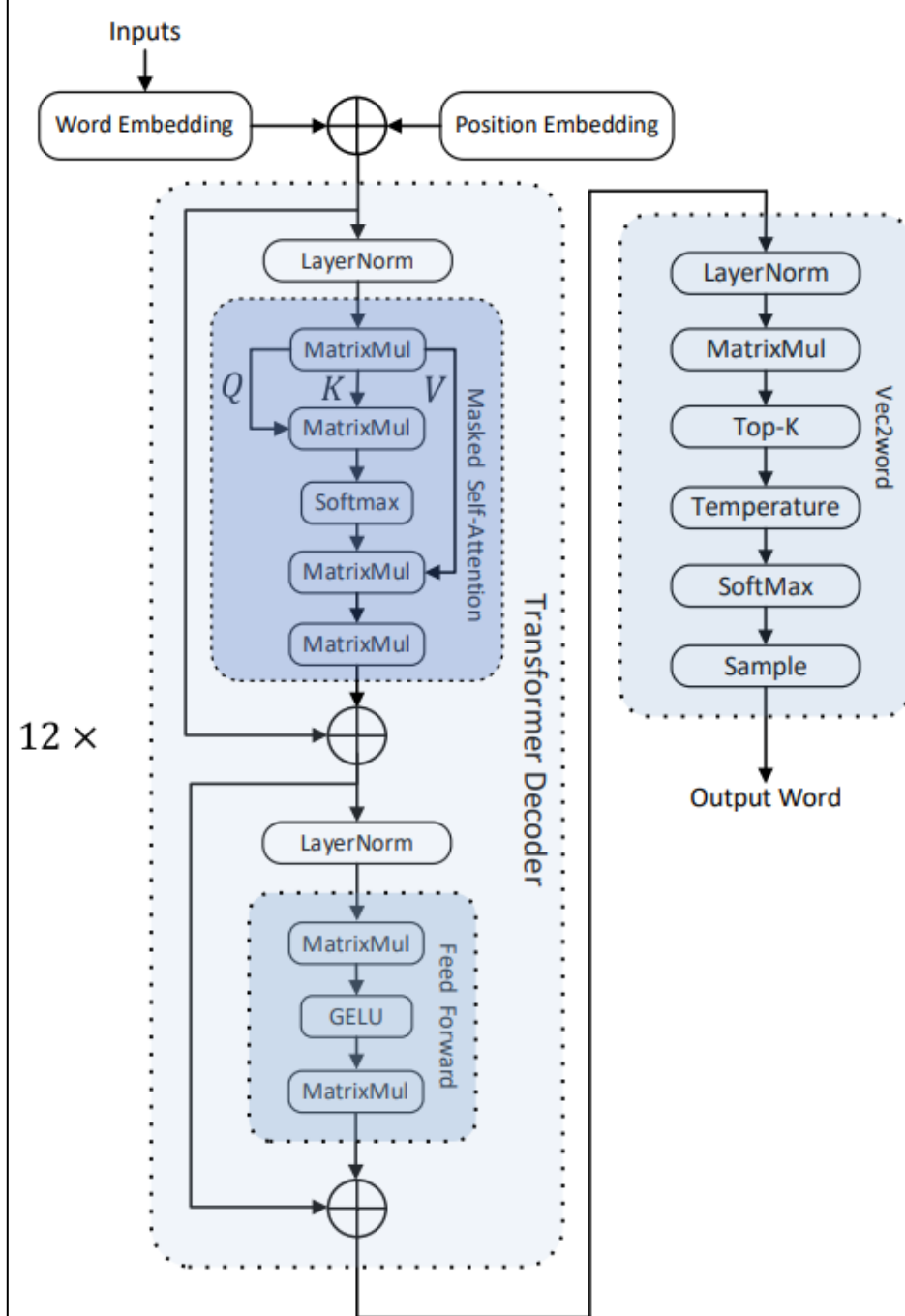
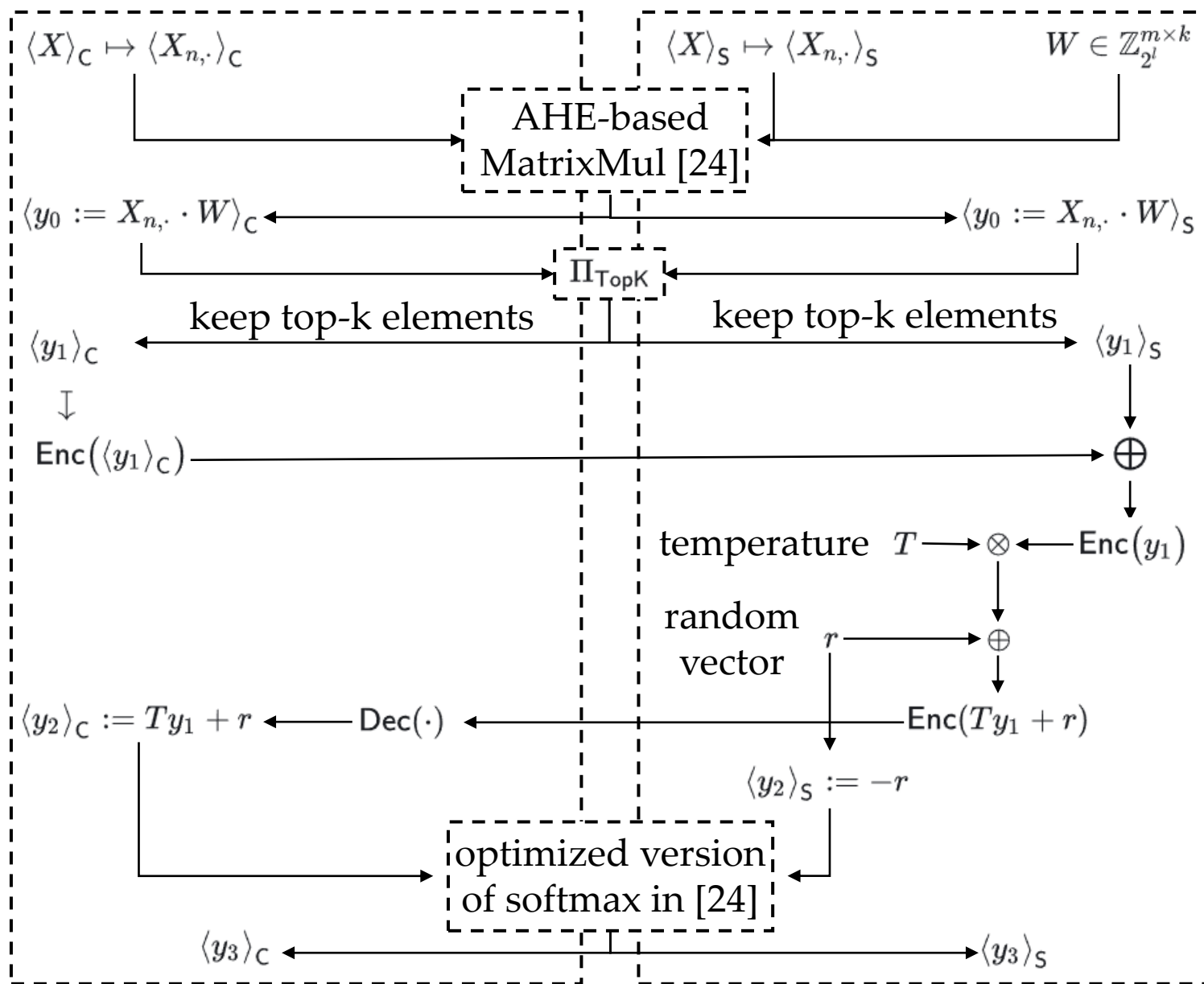


Figure 2: The architecture and workflow of GPT.