

# TEMPORAL FUSION TRANSFORMERS FOR INTERPRETABLE MULTI-HORIZON TIME SERIES FORECASTING

2021.03.28

백지윤

## 목차

1. 연구 의의, 목적 등
2. 용어 정리
3. 모델 구조
4. Loss function
5. 데이터 셋 / 실험 결과
6. 실제 활용 예시
7. 결론
8. 코드

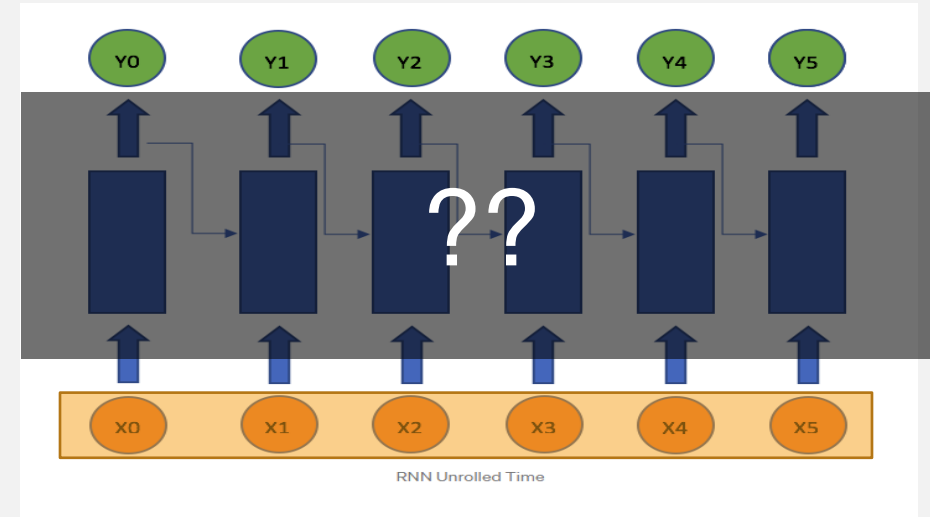
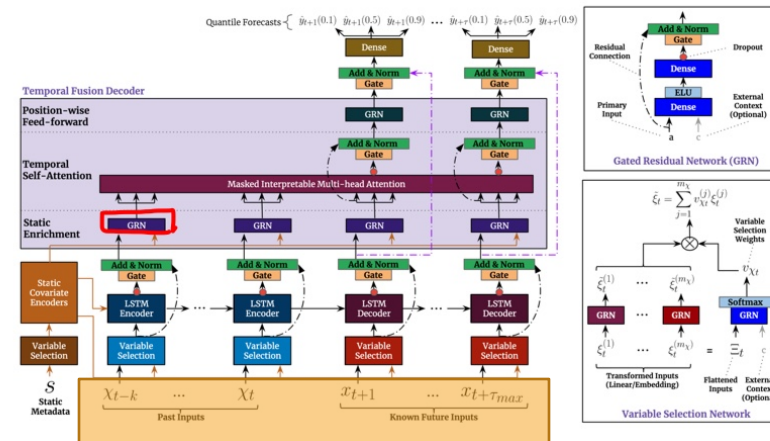
# 1. 연구 의의 및 목적

# TFT VS RNN

**ARCHITECTURE > GATING MECHANISMS**  
**VARIABLE SELECTION NETWORKS**  
**STATIC COVARIATE ENCODERS**

**ARCHITECTURE > GATING MECHANISMS** ✗  
**VARIABLE SELECTION NETWORKS.** ✗  
**STATIC COVARIATE ENCODERS** ✗

## 4. Model Architecture



**input > Static covariates (contexts)**  
**Observed inputs**  
**Known inputs**

**input >**  
**Observed inputs**

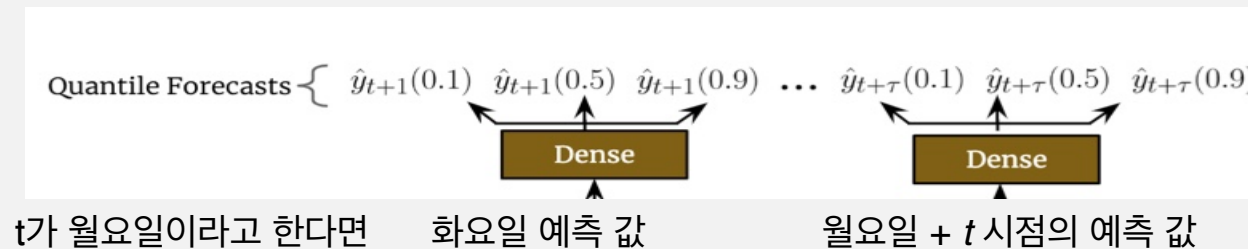
## 연구 목적

- Forecasting 에 영향을 줄 수 있는 보다 유연하고 풍부한 데이터를 모두 활용 할 수 있는 모델을 만들겠다
- 모델 forecasting 도중 해당 시점의 연산에서 필수적인 레이어와 features 만을 필터링하여 사용하겠다 (interpretability)

## 2. 용어 정리

## 용어 정리

- Horizon : 예측 범위 / Multi- Horizon : 여러 개의 예측 범위



- Static (=time invariant) covariates : 독립 변수가 종속 변수에 미치는 효과에 영향을 줄 수 있는 변수 ex. A 수학 문제집과 수학 성적과의 관계에서 학생들의 원 수학 실력 => 메타데이터
- Observed inputs (z), known inputs (x) ex. The way of week at time t

## 용어 정리

Horizon : 예측 범위 / Multi- Horizon : 여러 개의 예측 범위

Static (=time invariant) covariates : 독립 변수가 종속 변수에 미치는 효과에

영향을 줄 수 있는 변수 ex. A 수학 문제집과 수학 성적과의 관계에서 학생들의 원 수학 실력  
=>메타데이터

Observed inputs (z), known inputs (x)

ex. The way of week at time t

$$S_i \in \mathbb{R}^{m_s}$$

$$X_{i,t} = [z_{i,t}^T, x_{i,t}^T]$$

$$y_{i,t} \in \mathbb{R}$$

Static covariates

Inputs ; (observed, known)

Outputs

$$\hat{y}_i(q, t, \tau) = f_q(\tau, y_{i,t-k:t}, z_{i,t-k:t}, x_{i,t-k:t+\tau}, S_i)$$

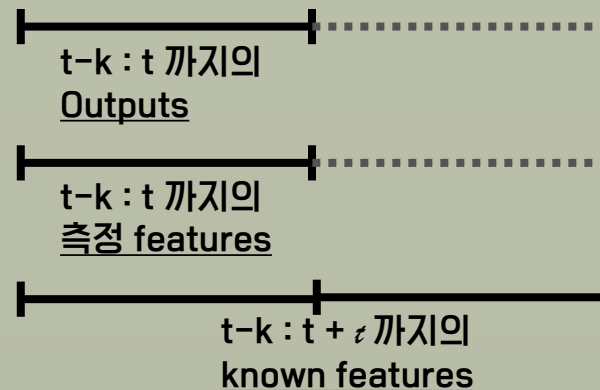


현재 시점 - k시점  
; Window

현재 시점 (t)

t 시점 이후 예측값은?

이용할 features 4가지



Static covariates  
(메타 데이터)

=> features 를 이해할 수 있는 Context 로 같이 넣어줄 예정



### 3. 모델 구조

# 모델 핵심 요소 6가지

Gating Mechanisms

모델 구조

자유도 크게  
자유도 크게  
자유도 작게  
자유도 작게  
자유도 작게  
자유도 작게

=> 정제된 features

Variable Selection Networks

S,X,Y 중 각 시점에 꼭 필요한 features 필터링

Static Covariate Encoders

S 메타 데이터를 features 를 이해할 수 있는 context화

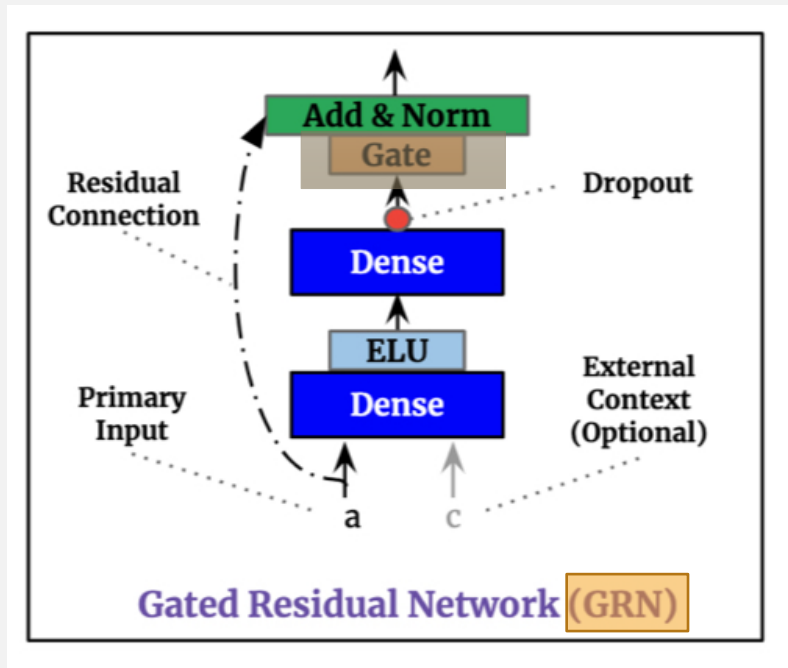
Interpretable Multi-Head Attention

Temporal Fusion Decoder

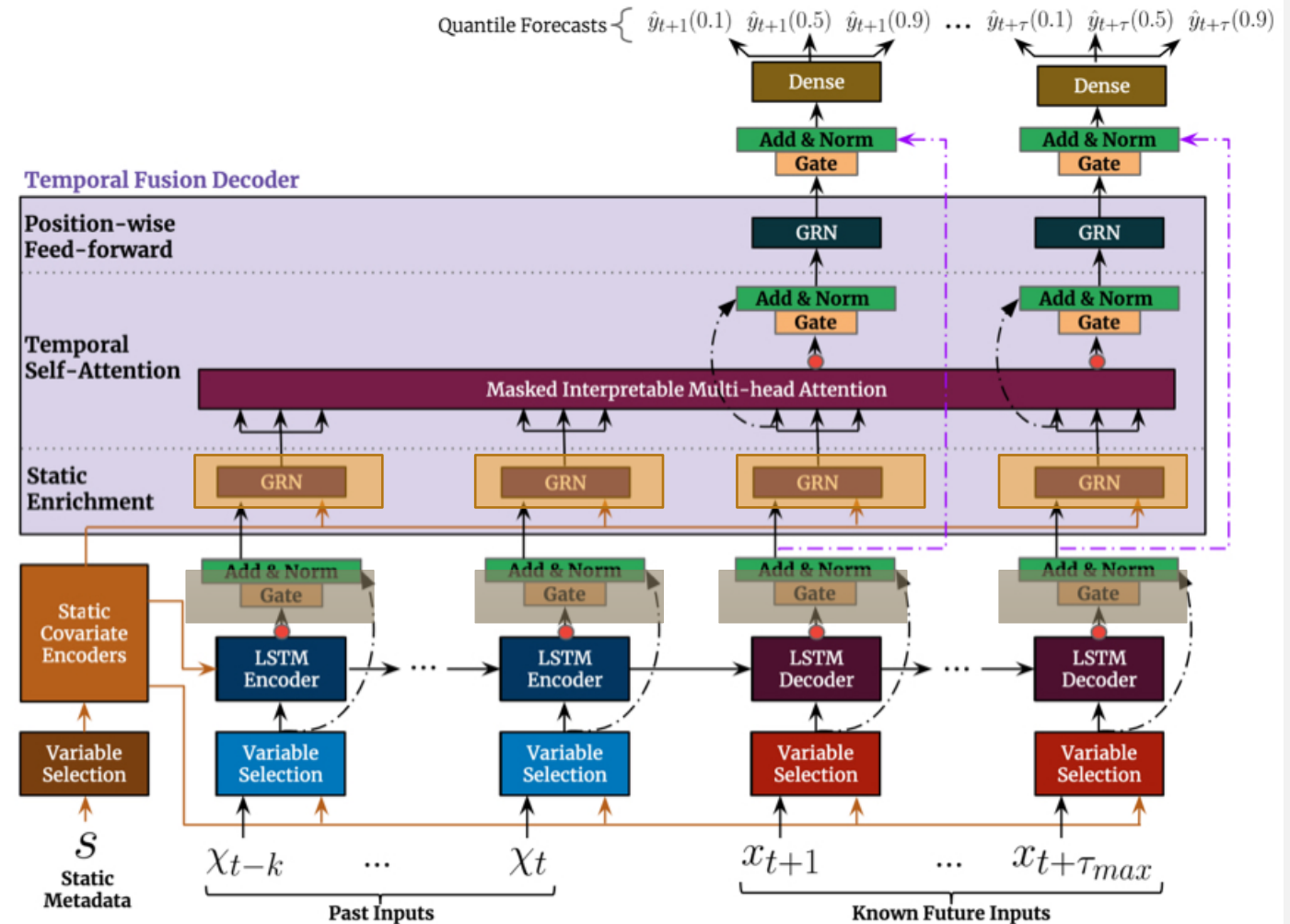
각 time step 의 장기간 상호관계 도출

Quantile Outputs

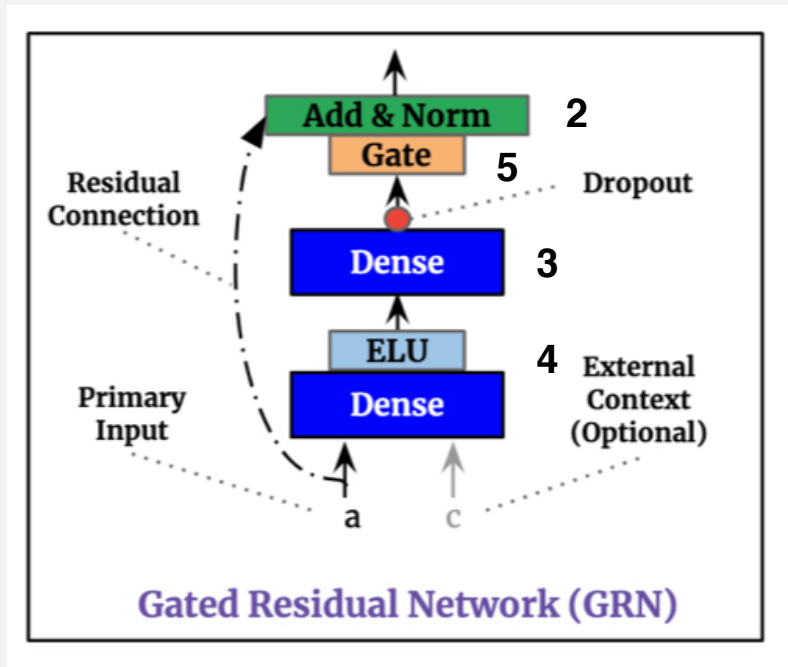
## Gating Mechanisms = GRN layer



Gate 는 TFT 모델의 거의 모든 층에 사용 되는 핵심 테크닉  
GRN layer 에도 gate 가 사용됨 !

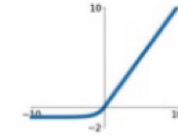


## Gating Mechanisms = GRN layer



**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



$$\eta_2 = \text{ELU} (W_{2,\omega} a + W_{3,\omega} c + b_{2,\omega}), \quad (4)$$

$$\eta_1 = W_{1,\omega} \eta_2 + b_{1,\omega}, \quad (3)$$

Gate

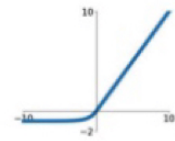
$$\text{GLU}_{\omega}(\gamma) = \sigma(W_{4,\omega} \gamma + b_{4,\omega}) \odot (W_{5,\omega} \gamma + b_{5,\omega}), \quad (5)$$

Dropout (training)

$$\text{GRN}_{\omega}(a, c) = \text{LayerNorm}(a + \text{GLU}_{\omega}(\eta_1)), \quad (2)$$

Gating Mechanisms = GRN layer

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



$$\eta_2 = \text{ELU}(W_{2,\omega} a + W_{3,\omega} c + b_{2,\omega}), \quad (4)$$

$$\eta_1 = W_{1,\omega} \eta_2 + b_{1,\omega}, \quad (3)$$

$$\text{GLU}_\omega(\gamma) = \sigma(W_{4,\omega} \gamma + b_{4,\omega}) \odot (W_{5,\omega} \gamma + b_{5,\omega}), \quad (5) \quad \text{Gate}$$

$$\text{GRN}_\omega(\mathbf{a}, \mathbf{c}) = \text{LayerNorm}(\mathbf{a} + \text{GLU}_\omega(\eta_1)), \quad (2)$$

#### (4) 가 모델 복잡도 결정

## (5)-1 가 value scaling 역할

Depends on (5)-1

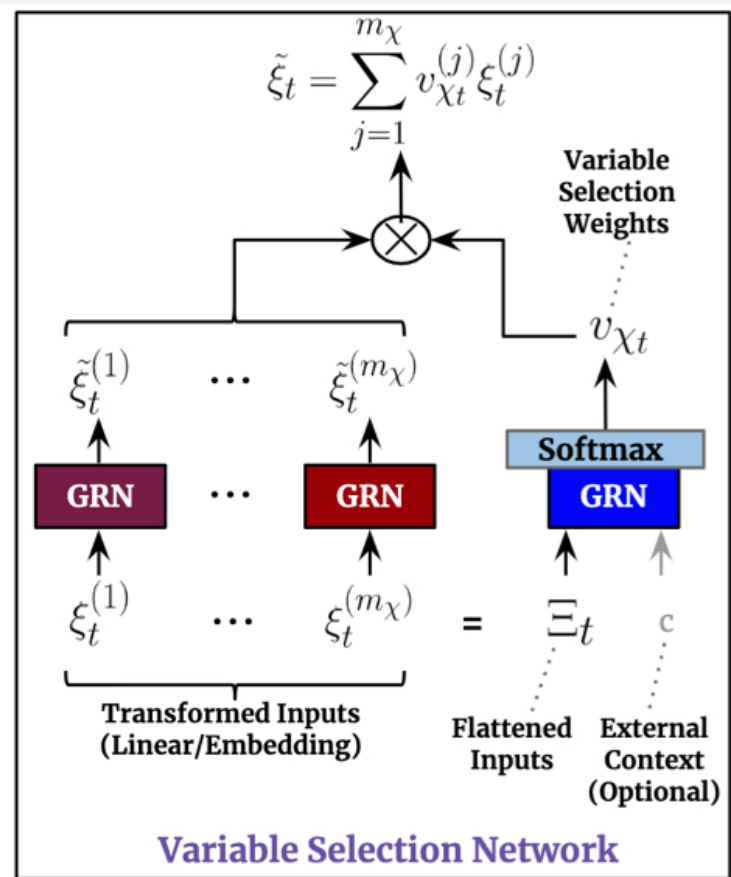
$$\ln(a + \sqrt{a^2 + b^2}) \quad (5)$$
$$\ln(a + \sqrt{a^2 + b^2}) \quad (5)$$
$$\ln(a + \sqrt{a^2 + b^2}) \quad (5)$$

■ ■ ■

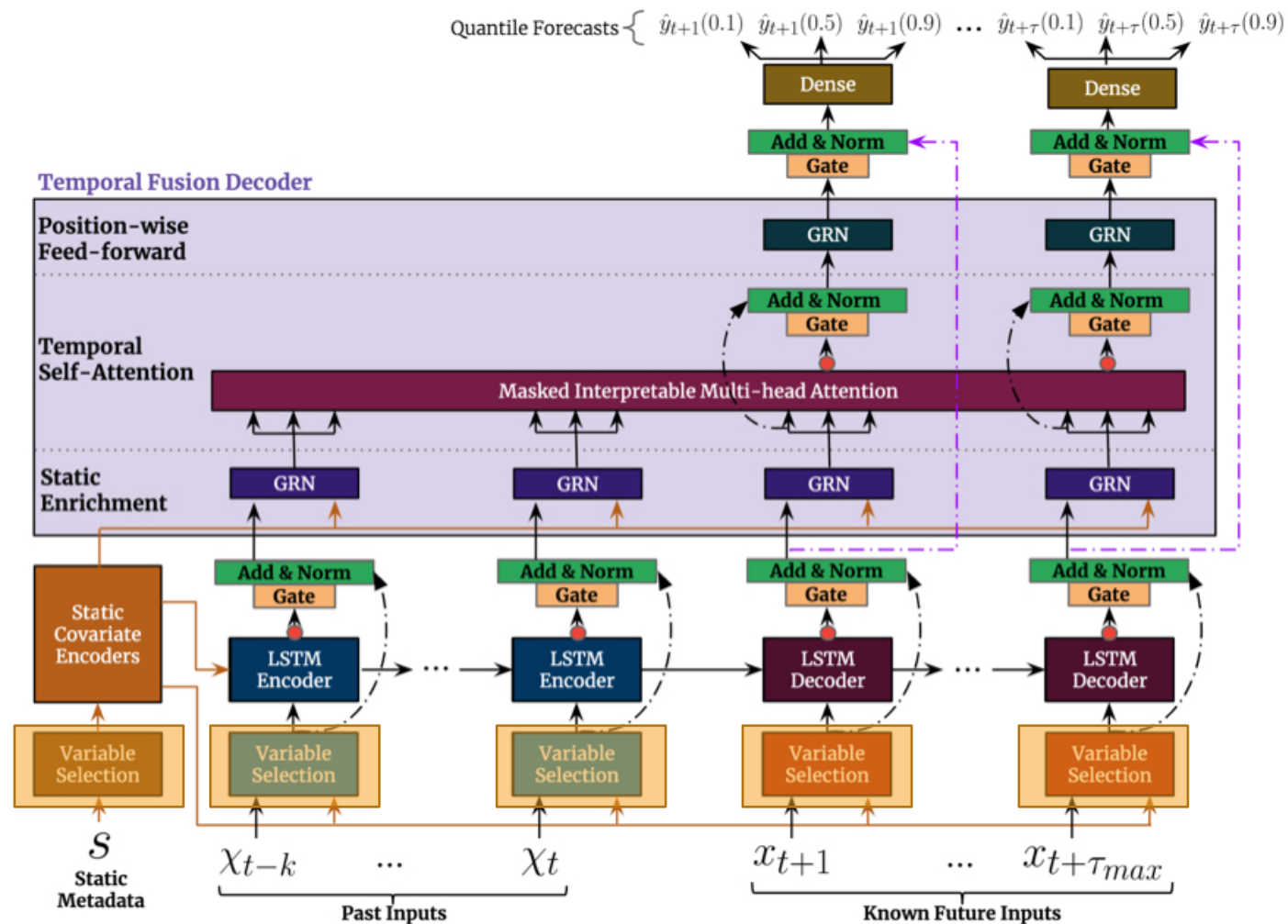
## Final Output

## Variable Selection Networks = VSN

각 시점 input 의 여러 features 중  
예측값에 확실히 관여하는 알맹이들만 남기기



VSN layer 는 GRN layer 을 포함  
모든 inputs 는 VSN layer 을 거침



## Variable Selection Networks = VSN

## 12월 아이스크림의 예측 판매량은 ?

T 일 input 값

공휴일 여부	엄마는 외계인	민트초코
○	200개	100개

# Categorical

Entity embedding ( $D_{Model}$  vector)

## Continuous

## Linear transformation `nn.linear(1,D Model vector)`

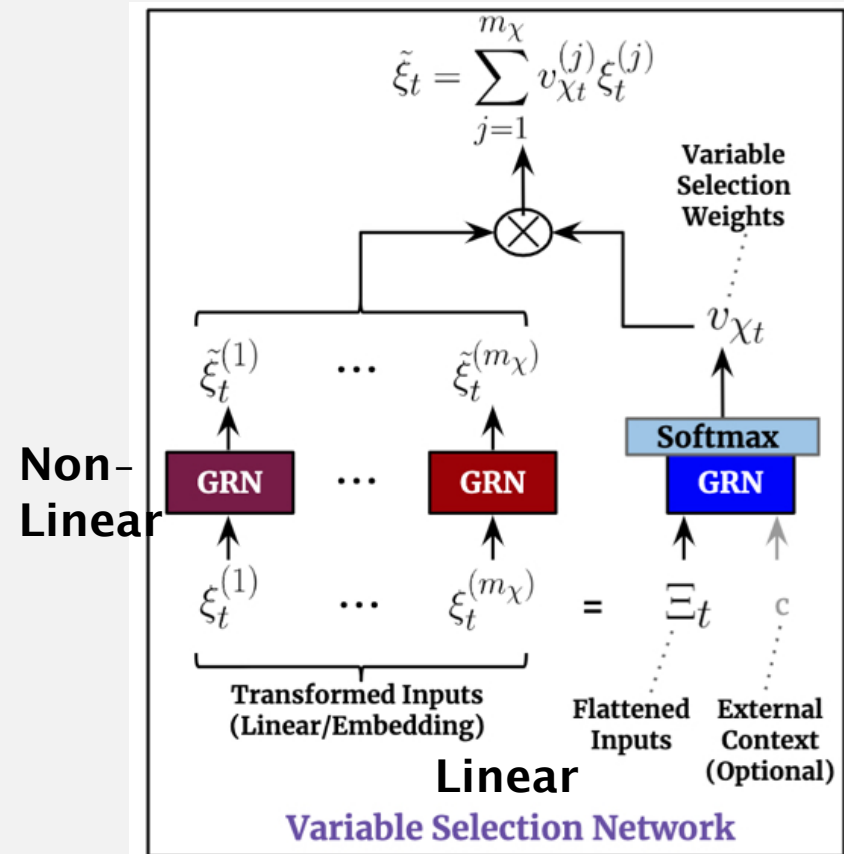


## Flattened Inputs



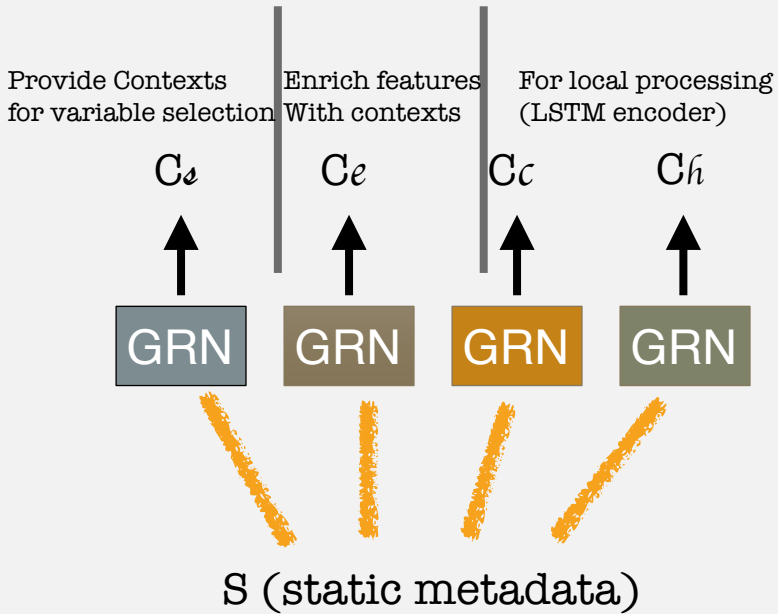
## Variable Selection Weights

Feature 개수만큼 (j개)

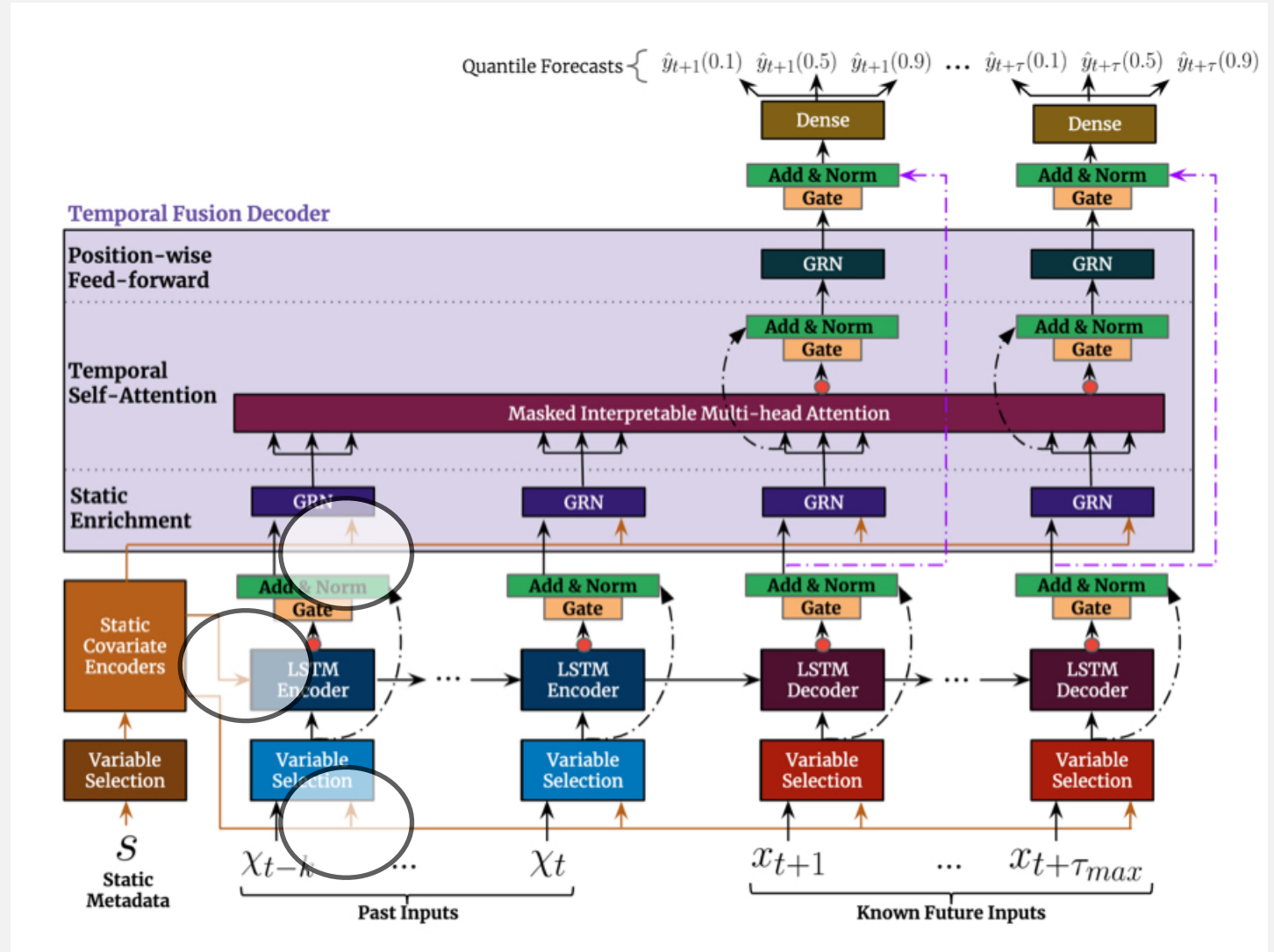


## Static Covariate Encoders

S 메타 데이터를 features 을 이해할 수 있는 context 로 사용



각기 다른 4개의 GRN 을 사용하여서  
쓰임이 다른 4개의 문맥 생성

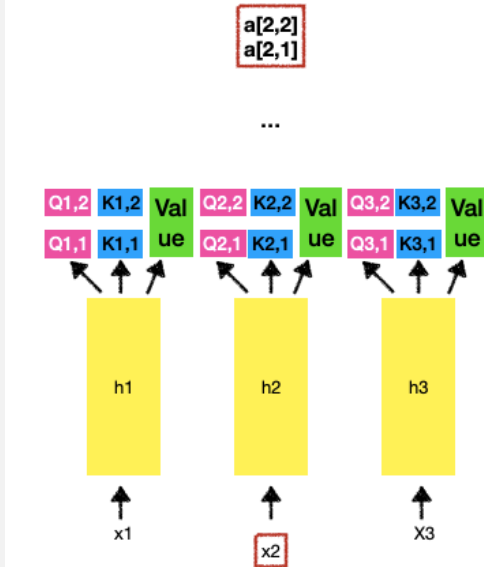
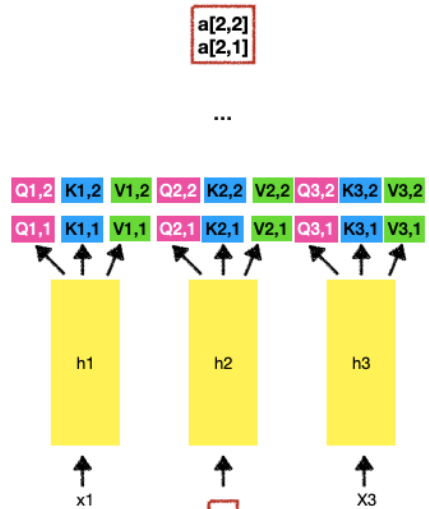




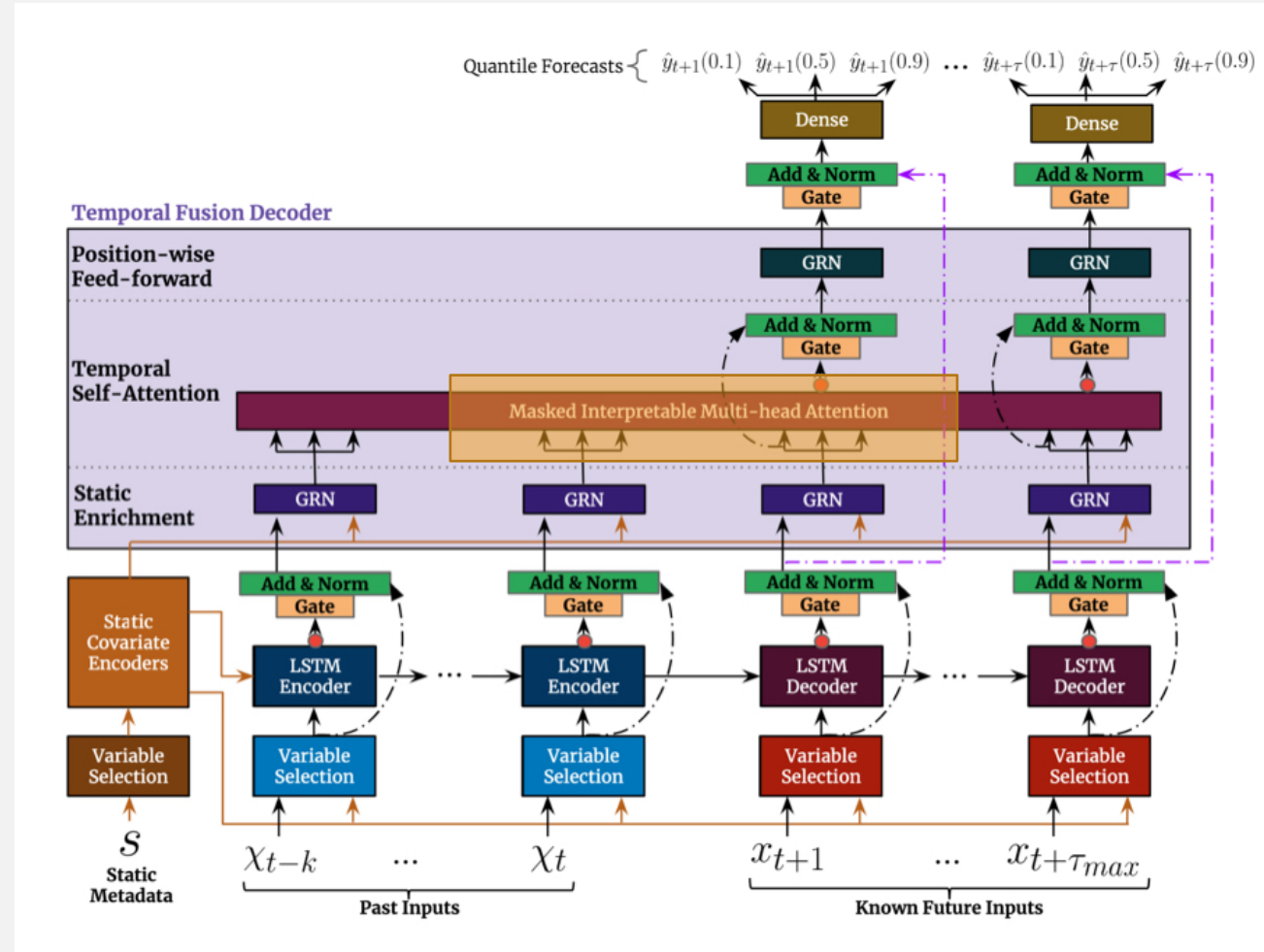
## Interpretable Multi-Head Attention

각 time step 의 장기간 상호관계 도출

원 Multi-head attention TFT Multi-head attention



Multi-attention 아키텍처 그대로 갖고가되,  
query,key,value 중 value 는 모든 head 에서 동일



## Interpretable Multi-Head Attention

각 time step 의 장기간 상호관계 도출

TFT Multi-head attention

$$\text{MultiHead}(Q, K, V) = [H_1, \dots, H_{m_H}] W_H, \quad (11)$$

$$H_h = \text{Attention}(Q W_Q^{(h)}, K W_K^{(h)}, V W_V^{(h)}), \quad (12)$$

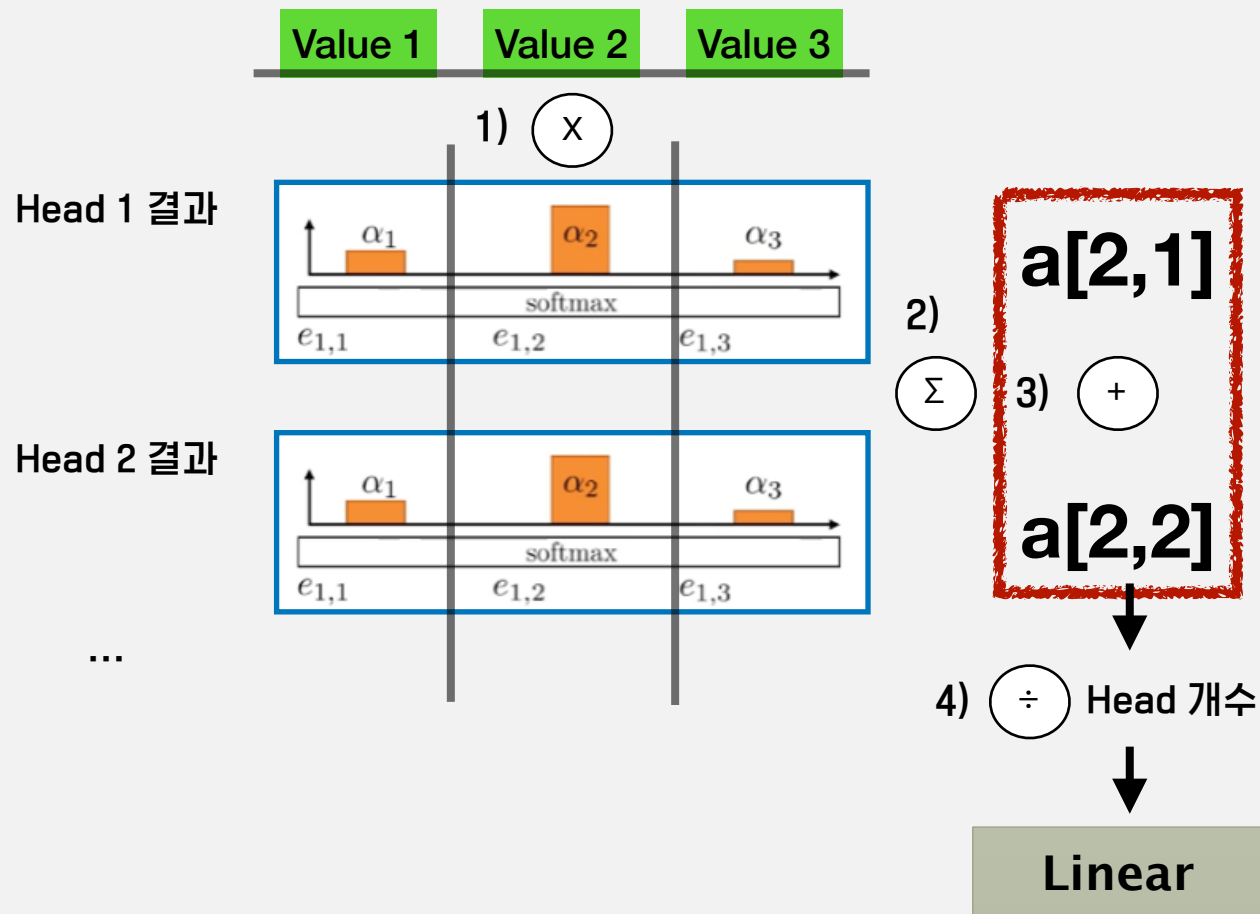
$$\text{InterpretableMultiHead}(Q, K, V) = \tilde{H} W_H, \quad (13)$$

$$\tilde{H} = \tilde{A}(Q, K) V W_V, \quad (14)$$

$$= \left\{ 1/H \sum_{h=1}^{m_H} A(Q W_Q^{(h)}, K W_K^{(h)}) \right\} V W_V, \quad (15)$$

$$= 1/H \sum_{h=1}^{m_H} \text{Attention}(Q W_Q^{(h)}, K W_K^{(h)}, V W_V), \quad (16)$$

같은 timestep 은 다른 head 에서도  
동일한 value 를 갖게 함으로써 앙상블하는 방식으로 작용

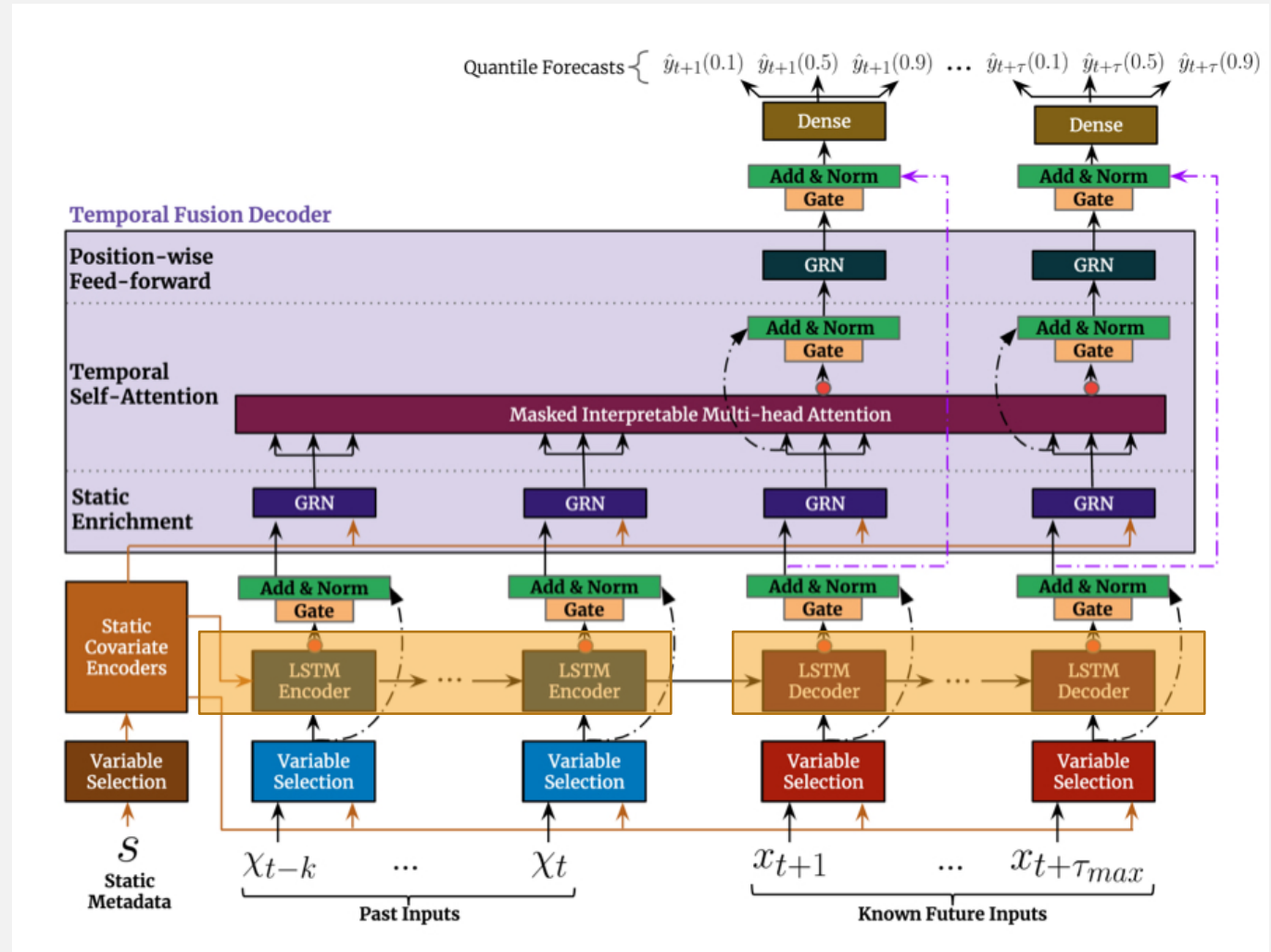


## Temporal Fusion Decoder - (1) Seq2Seq layer

각 시점의 특징 추출, 각 시점 정보 추가

$$\tilde{\phi}(t, n) = \text{LayerNorm} \left( \tilde{\xi}_{t+n} + \text{GLU}_{\tilde{\phi}}(\phi(t, n)) \right), \quad (17)$$

Temporal Fusion Decoder 에 들어가는 final inputs

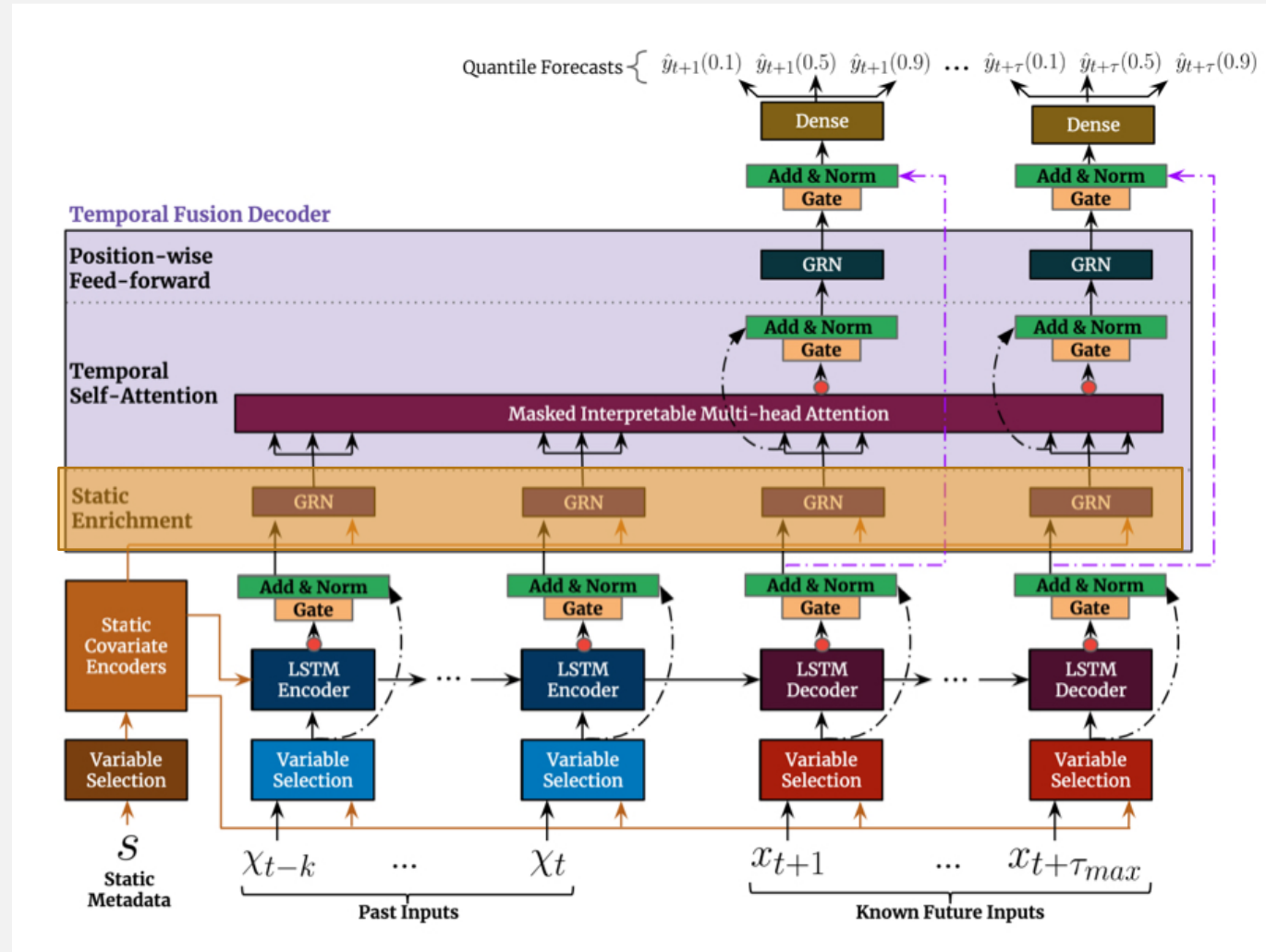


## Temporal Fusion Decoder - (2) Static Enrichment Layer

temporal features 에 메타데이터를 활용해  
풍부한 문맥 추가

$$\theta(t, n) = \text{GRN}_{\theta} \left( \tilde{\phi}(t, n), c_e \right), \quad (18)$$

## Temporal Fusion Decoder 에 들어가는 final inputs



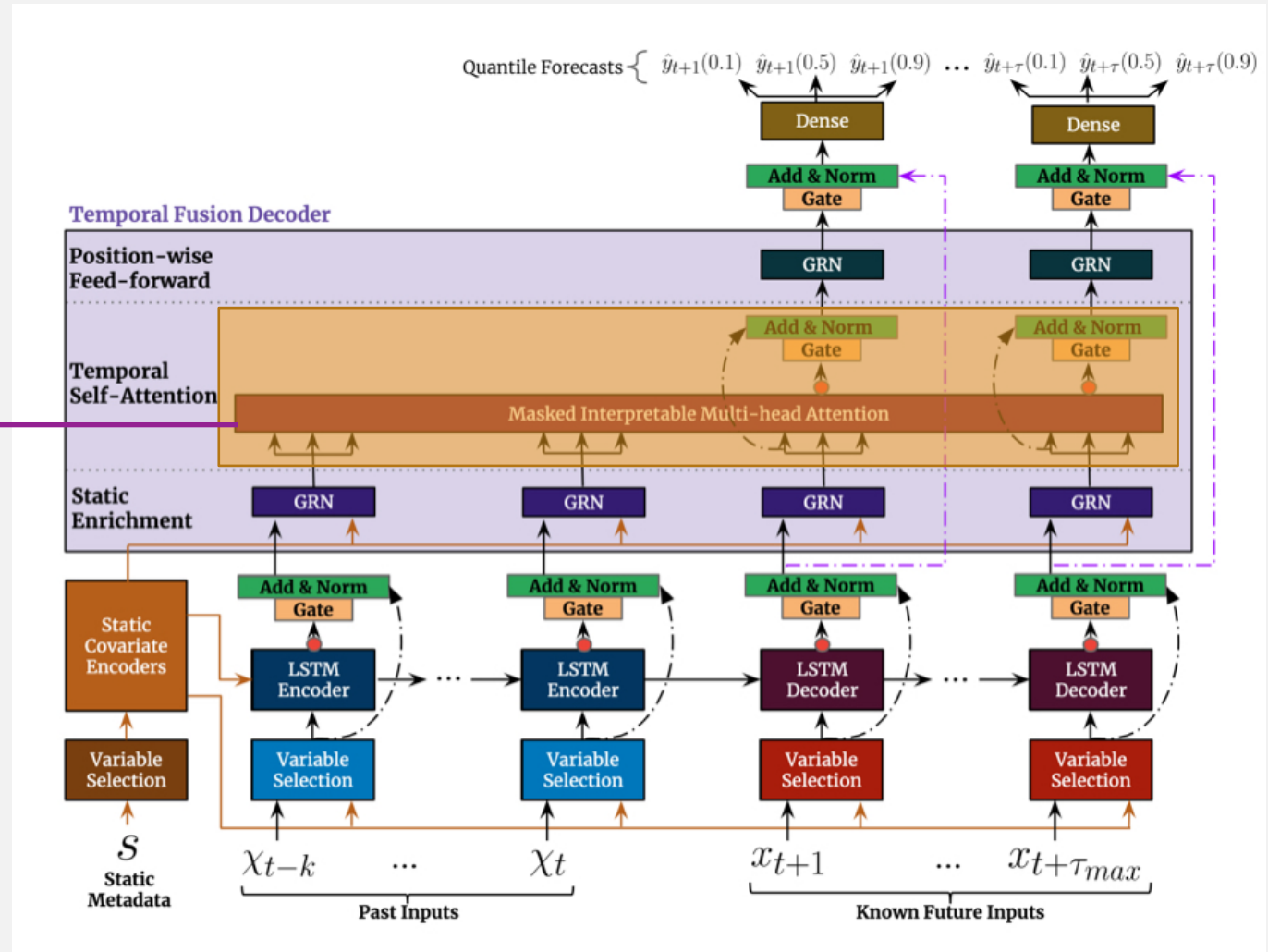
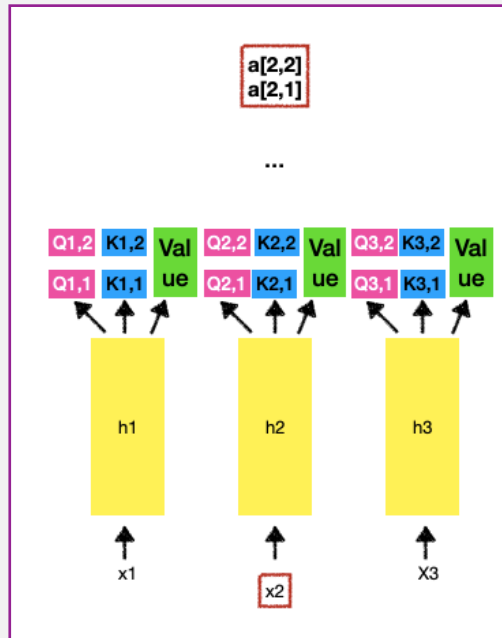
## Temporal Fusion Decoder - (3) Temporal Self-Attention Layer

## Interpretable Multi-Head Attention

이 투입된 layer

각 time step 의 장기간 상호관계 도출

TFT Multi-head attention



## Temporal Fusion Decoder - (3) Temporal Self-Attention Layer

## Interpretable Multi-Head Attention

이 투입된 layer

각 time step 의 장기간 상호관계 도출

$$\theta(t) = [\theta(t, \kappa), \dots, \theta(t, \tau)]^T$$

Static enrichment layer 에서 나온 값들을 하나의 단일 벡터로 묶이기



TFT masked Multi-head attention : 이전 시점들과의 관계(attention) 만을 이용하도록 하기 위해서



$$\delta(t, n) = \text{LayerNorm}(\theta(t, n) + \text{GLU}_{\delta}(\beta(t, n))). \quad (20)$$

$$n > t$$

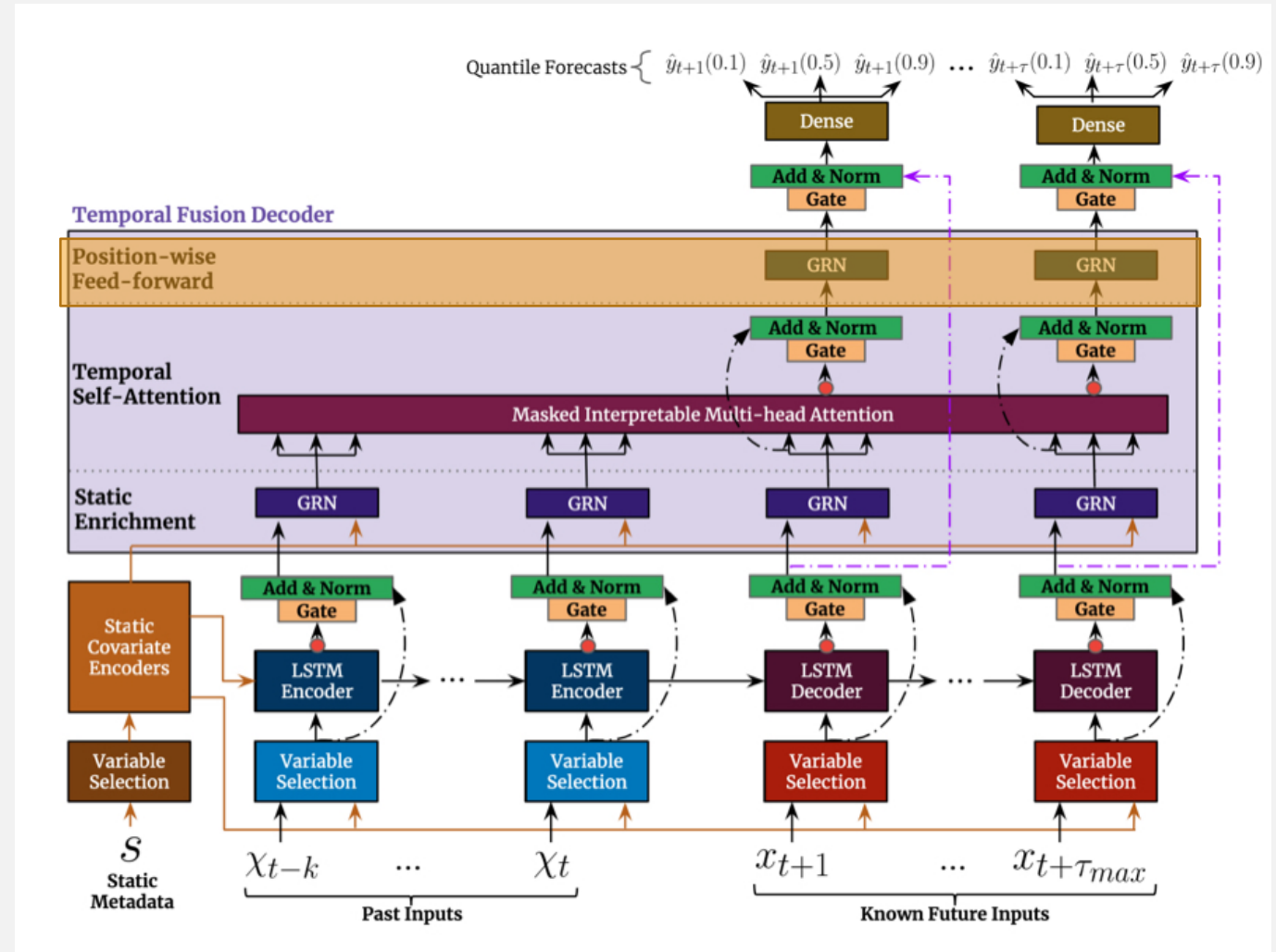
## Temporal Fusion Decoder - (4) Position-wise Feed-forward layer

non-linear 층 (GRN) 추가

$$\psi(t, n) = \text{GRN}_{\psi}(\delta(t, n)), \quad (21)$$

$$\bar{\psi}(t, n) = \text{LayerNorm}(\bar{\phi}(t, n) + \text{GLU}_{\bar{\psi}}(\psi(t, n))), \quad (22)$$

$n > t$



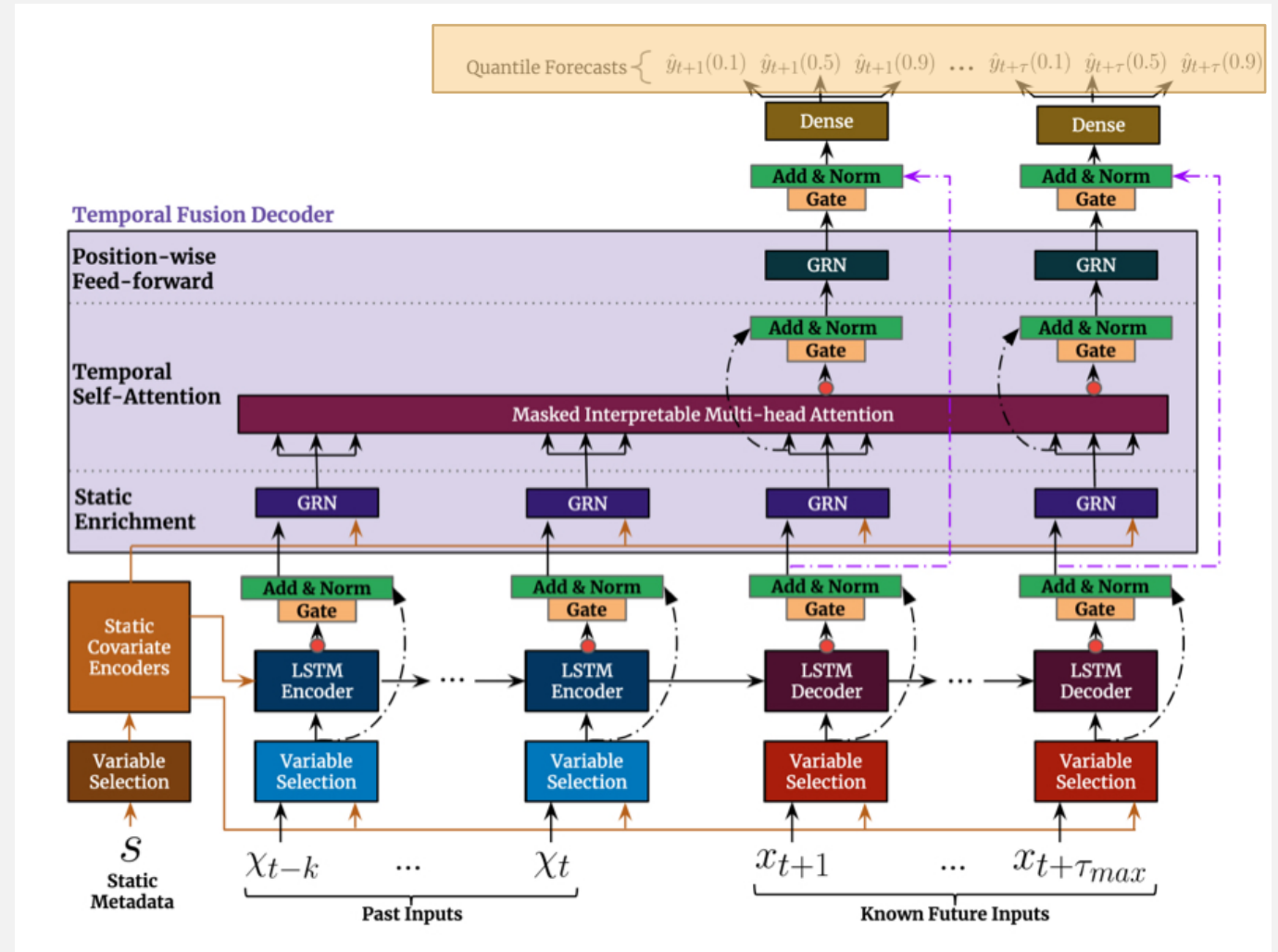


## Quantile Outputs

최종 quantile 확률별 outputs 출력

각 quantile 마다 각자의 linear 층으로 output 값 계산  
Quantile : 해당 예측값이 나올 확률이 quantile %

$$\hat{y}(q, t, \tau) = W_q \bar{\psi}(t, \tau) + b_q, \quad (23)$$





## 4. LOSS FUNCTION

# LOSS FUNCTION

Quantile Forecasts  $\{ \hat{y}_{t+1}(0.1) \ \hat{y}_{t+1}(0.5) \ \hat{y}_{t+1}(0.9) \ \dots \ \hat{y}_{t+\tau}(0.1) \ \hat{y}_{t+\tau}(0.5) \ \hat{y}_{t+\tau}(0.9) \}$

$$\mathcal{L}(\Omega, W) = \sum_{y_t \in \Omega} \sum_{q \in Q} \sum_{\tau=1}^{\tau_{max}} \frac{QL(y_t, \hat{y}(q, t - \tau, \tau), q)}{M \tau_{max}} \quad (24)$$

해당 데이터셋 内      예측할 개수

$$QL(y, \hat{y}, q) = q(y - \hat{y})_+ + (1 - q)(\hat{y} - y)_+, \quad (25)$$

$\max(0, x)$

- Quantile loss function 출처 : <https://arxiv.org/pdf/1711.11053.pdf>
- 조금 변형된 loss function 으로 out of sample test 도 같이 진행

$$q\text{-Risk} = \frac{2 \sum_{y_t \in \tilde{\Omega}} \sum_{\tau=1}^{\tau_{max}} QL(y_t, \hat{y}(q, t - \tau, \tau), q)}{\sum_{y_t \in \tilde{\Omega}} \sum_{\tau=1}^{\tau_{max}} |y_t|}, \quad (26)$$

where  $\tilde{\Omega}$  is the domain of test samples. Full details on hyperparameter optimization and training can be found in Appendix A.